COMPLEMENTARY SCIENCE SERIES



# Crystallography Made Crystal Clear

A GUIDE FOR USERS OF MACROMOLECULAR MODELS



# Gale Rhodes



Crystallography Made Crystal Clear

#### WHAT IS THE COMPLEMENTARY SCIENCE SERIES?

We hope you enjoy this book. If you would like to read other quality science books with a similar orientation see the order form and reproductions of the front and back covers of other books in the series at the end of this book.

The **Complementary Science Series** is an introductory, interdisplinary, and relatively inexpensive series of paperbacks for science enthusiasts. The series covers core subjects in chemistry, physics, and biological sciences but often from an interdisciplinary perspective. They are deliberately unburdened by excessive pedagogy, which is distracting to many readers, and avoid the often plodding treatment in many textbooks.

These titles cover topics that are particularly appropriate for self-study although they are often used as complementary texts to supplement standard discussion in textbooks. Many are available as examination copies to professors teaching appropriate courses.

The series was conceived to fill the gaps in the literature between conventional textbooks and monographs by providing real science at an accessible level, with minimal prerequisites so that students at all stages can have expert insight into important and foundational aspects of current scientific thinking.

Many of these titles have strong interdisciplinary appeal and all have a place on the bookshelves of literate laypersons.

Potential authors are invited to contact our editorial office: j.hayhurst@elsevier.com. Feedback on the titles is welcome.

Titles in the *Complementary Science Series* are detailed at the end of these pages. A 15% discount is available (to owners of this edition) on other books in this series—see order form at the back of this book.

#### Physics

Physics in Biology and Medicine, 2nd Edition; Paul Davidovits, 0122048407

Introduction to Relativity; John B. Kogut, 0124175619

Fusion: The Energy of the Universe; Gary McCracken and Peter Stott, 012481851X

#### Chemistry

The Physical Basis of Chemistry, 2nd Edition; Warren S. Warren, 0127358552

*Chemistry Connections: The Chemical Basis of Everyday Phenomena*, 2nd Edition; Kerry K. Karukstis and Gerald R. Van Hecke, 0124001513

Fundamentals of Quantum Chemistry; James E. House, 0123567718

Introduction to Quantum Mechanics; S. M. Blinder, 0121060519

#### Geology

Earth Magnetism; Wallace Hall Campbell, 0121581640

www.books.elsevier.com

This Page Intentionally Left Blank

## Crystallography Made Crystal Clear

A Guide for Users of Macromolecular Models

**Third Edition** 

### Gale Rhodes

Chemistry Department University of Southern Maine Portland, Maine CMCC Home Page: www.usm.maine.edu/~rhodes/CMCC



AMSTERDAM • BOSTON • HEIDELBERG • LONDON NEW YORK • OXFORD • PARIS • SAN DIEGO SAN FRANCISCO • SINGAPORE • SYDNEY • TOKYO



Academic Press is an imprint of Elsevier

Acquisitions Editor: Jeremy Hayhurst Project Manager: Jeff Freeland Marketing Manager: Linda Beattie Cover Design: Eric DeCicco Composition: Cepha Imaging Private Limited Cover Printer: Transcontinental Printing Book Group Interior Printer: Transcontinental Printing Book Group

Academic Press is an imprint of Elsevier 30 Corporate Drive, Suite 400, Burlington, MA 01803, USA 525 B Street, Suite 1900, San Diego, California 92101-4495, USA 84 Theobald's Road, London WC1X 8RR, UK

This book is printed on acid-free paper.  $\bigotimes$ 

Copyright © 2006, Elsevier Inc. All rights reserved.

No part of this publication may be reproduced or transmitted in any form or by any means, electronic or mechanical, including photocopy, recording, or any information storage and retrieval system, without permission in writing from the publisher.

Permissions may be sought directly from Elsevier's Science & Technology Rights Department in Oxford, UK: phone: (+44) 1865 843830, fax: (+44) 1865 853333, E-mail: permissions@elsevier.com. You may also complete your request on-line via the Elsevier homepage (http://elsevier.com), by selecting "Support & Contact" then "Copyright and Permission" and then "Obtaining Permissions."

#### Library of Congress Cataloging-in-Publication Data

Rhodes, Gale.
Crystallography made crystal clear / Gale Rhodes.– 3rd ed.
p. cm. – (Complementary science series)
Includes bibliographical references and index.
ISBN 0-12-587073-6 (alk. paper)
1. X-ray crystallography. 2. Macromolecules–Structure. 3. Proteins–Structure. I.
Title. II. Series.

QP519.9.X72R48 2006 547'.7-dc22

2005057239

#### British Library Cataloguing-in-Publication Data

A catalogue record for this book is available from the British Library.

ISBN 13: 978-0-12-587073-3 ISBN 10: 0-12-587073-6

For information on all Academic Press Publications visit our Web site at www.books.elsevier.com

Printed in Canada 06 07 08 09 10 8 7 6 5 4 3 2 1

### Working together to grow libraries in developing countries

www.elsevier.com | www.bookaid.org | www.sabre.org

ELSEVIER BOOK AID Sabre Foundation

Like everything, for Pam.

This Page Intentionally Left Blank

| Pro<br>Pro<br>Pro | eface<br>eface<br>eface | to the Th<br>to the Sec<br>to the Fir | ird Edition   | xv<br>xix<br>xxiii |  |
|-------------------|-------------------------|---------------------------------------|---|--------------------|--|
| 1                 | Mod                     | lel and Molecule                      |   |                    |  |
| 2                 | An (                    | Overview of Protein Crystallography   |   |                    |  |
|                   | 2.1                     | Introduction                          |   |                    |  |
|                   |                         | 2.1.1                                 | Obtaining an image of a microscopic object                | 8                  |  |
|                   |                         | 2.1.2                                 | Obtaining images of molecules                             | 9                  |  |
|                   |                         | 2.1.3                                 | A thumbnail sketch of protein crystallography             | 9                  |  |
|                   | 2.2                     | Crystals                              | ;   | 10                 |  |
|                   |                         | 2.2.1                                 | The nature of crystals                                    | 10                 |  |
|                   |                         | 2.2.2                                 | Growing crystals  | 11                 |  |
|                   | 2.3                     | Collecti                              | ng X-ray data   | 13                 |  |
|                   | 2.4                     | Diffract                              | tion  |                    |  |
|                   |                         | 2.4.1 Simple objects                  |   | 15                 |  |
|                   |                         | 2.4.2                                 | Arrays of simple objects: Real and reciprocal lattices    | 16                 |  |
|                   |                         | 2.4.3                                 | Intensities of reflections                                | 16                 |  |
|                   |                         | 2.4.4                                 | Arrays of complex objects                                 | 17                 |  |
|                   |                         | 2.4.5                                 | Three-dimensional arrays                                  | 18                 |  |
|                   | 2.5                     | Coordinate systems in crystallography |   |                    |  |
|                   | 2.6                     | The mat                               | thematics of crystallography: A brief description         | 20                 |  |
|                   |                         | 2.6.1                                 | Wave equations: Periodic functions                        | 21                 |  |
|                   |                         | 2.6.2                                 | Complicated periodic functions: Fourier series and sums   | 23                 |  |
|                   |                         | 2.6.3                                 | Structure factors: Wave descriptions of X-ray reflections | 24                 |  |
|                   |                         | 2.6.4                                 | Electron-density maps                                     | 26                 |  |
|                   |                         | 2.6.5                                 | Electron density from structure factors                   | 27                 |  |
|                   |                         | 2.6.6                                 | Electron density from measured reflections                | 28                 |  |
|                   |                         | 2.6.7                                 | Obtaining a model   | 30                 |  |

| 3                | Prot                | Protein Crystals 3                     |  |          |  |  |
|------------------|---------------------|--|--|----------|--|--|
|                  | 3.1                 | Propert                                | verties of protein crystals                                    |          |  |  |
| 3.1.1            |                     |  | Introduction   | 31       |  |  |
|                  |                     | 3.1.2                                  | Size, structural integrity, and mosaicity                      | 31       |  |  |
|                  |                     | 3.1.3                                  | Multiple crystalline forms                                     | 33       |  |  |
|                  |                     | 3.1.4                                  | Water content  | 34       |  |  |
|                  | 3.2                 | Eviden                                 | ce that solution and crystal structures are similar            | 35       |  |  |
|                  |                     | 3.2.1                                  | Proteins retain their function in the crystal                  | 35       |  |  |
|                  |                     | 3.2.2                                  | X-ray structures are compatible with other structural evidence | 36       |  |  |
|                  |                     | 3.2.3                                  | Other evidence   | 37       |  |  |
|                  | 3.3                 | Growing protein crystals               |  |          |  |  |
|                  | 3.3.1 Introduction  |  | Introduction   | 37       |  |  |
|                  |                     | 3.3.2                                  | Growing crystals: Basic procedure                              | 38       |  |  |
|                  |                     | 3.3.3                                  | Growing derivative crystals                                    | 40       |  |  |
|                  |                     | 3.3.4                                  | Finding optimal conditions for crystal growth                  | 41       |  |  |
|                  | 3.4                 | Judging crystal quality                |  |          |  |  |
|                  | 3.5                 | Mounting crystals for data collection  |  |          |  |  |
|                  |                     |  |  |          |  |  |
| 4                | Coll                | ecting D                               | iffraction Data  | 49       |  |  |
| •                | 4.1                 | Introdu                                | ction  | 49       |  |  |
|                  | 4.2                 | Geometric principles of diffraction    |  |          |  |  |
|                  | 1.2                 | 4 2 1                                  | The generalized unit cell                                      | 49       |  |  |
|                  |                     | 422                                    | Indices of the atomic planes in a crystal                      | 50       |  |  |
|                  |                     | 423                                    | Conditions that produce diffraction: Bragg's law               | 55       |  |  |
|                  |                     | 424                                    | The reciprocal lattice   | 57       |  |  |
|                  |                     | 425                                    | Bragg's law in reciprocal space                                | 60       |  |  |
|                  |                     | 426                                    | Number of measurable reflections                               | 64       |  |  |
|                  |                     | 427                                    | Unit-cell dimensions   | 65       |  |  |
|                  |                     | 428                                    | Unit-cell symmetry   | 65       |  |  |
|                  | 43                  | Collect                                | ing X-ray diffraction data                                     | 73       |  |  |
|                  | ч.5                 | 431                                    | Introduction   | 73       |  |  |
|                  |                     | 432                                    | Y ray sources  | 73       |  |  |
|                  |                     | 433                                    | Detectors  | 73<br>77 |  |  |
|                  |                     | 4.3.3                                  | Comeros  | 80       |  |  |
|                  |                     | 4.3.4                                  | Calling and postrafinament of intensity data                   | 85       |  |  |
|                  |                     | 4.3.3                                  | Determining unit call dimensions                               | 0J<br>86 |  |  |
|                  |                     | 4.5.0                                  | Summetry and the strategy of collecting date                   | 00       |  |  |
|                  | 11                  | 4.3.7<br>Summo                         |  | 00<br>80 |  |  |
|                  | 4.4                 | Summa                                  | шу   | 09       |  |  |
| 5                | Free                | n Diffra                               | ction Data to Electron Density                                 | 01       |  |  |
| 3                | <b>F FOI</b><br>5 1 | Un Din action Data to Electron Density |  |          |  |  |
| 5.1 Introduction |                     | Equation                               | cuoil  | 91       |  |  |
|                  | 3.2                 | rourier                                | One dimensional waves  | 92       |  |  |
|                  |                     | 5.2.1                                  |  | 92       |  |  |
|                  |                     | 5.2.2                                  | I nree-aimensional waves                                       | 94       |  |  |

x

|   |     | 5.2.3     | The Fourier transform: General features                    | 96  |
|---|-----|-----------|--|-----|
|   |     | 5.2.4     | Fourier this and Fourier that: Review                      | 97  |
|   | 5.3 | Fourier   | mathematics and diffraction                                | 98  |
|   |     | 5.3.1     | Structure factor as a Fourier sum                          | 98  |
|   |     | 5.3.2     | Electron density as a Fourier sum                          | 99  |
|   |     | 5.3.3     | Computing electron density from data                       | 100 |
|   |     | 5.3.4     | The phase problem  | 101 |
|   | 5.4 | Meanin    | g of the Fourier equations                                 | 101 |
|   |     | 5.4.1     | Reflections as terms in a Fourier sum: Eq. (5.18)          | 101 |
|   |     | 5.4.2     | Computing structure factors from a model: Eq. (5.15)       |     |
|   |     |           | and Eq. (5.16)   | 104 |
|   |     | 5.4.3     | Systematic absences in the diffraction pattern: Eq. (5.15) | 105 |
|   | 5.5 | Summa     | ry: From data to density                                   | 107 |
| 6 | Obt | aining P  | hases  | 109 |
|   | 6.1 | Introdu   | ction  | 109 |
|   | 6.2 | Two-di    | mensional representation of structure factors              | 112 |
|   |     | 6.2.1     | Complex numbers in two dimensions                          | 112 |
|   |     | 6.2.2     | Structure factors as complex vectors                       | 112 |
|   |     | 6.2.3     | Electron density as a function of intensities and phases   | 115 |
|   | 6.3 | Isomor    | phous replacement  | 117 |
|   |     | 6.3.1     | Preparing heavy-atom derivatives                           | 117 |
|   |     | 6.3.2     | Obtaining phases from heavy-atom data                      | 119 |
|   |     | 6.3.3     | Locating heavy atoms in the unit cell                      | 124 |
|   | 6.4 | Anoma     | lous scattering  | 128 |
|   |     | 6.4.1     | Introduction   | 128 |
|   |     | 6.4.2     | Measurable effects of anomalous scattering                 | 128 |
|   |     | 6.4.3     | Extracting phases from anomalous scattering data           | 130 |
|   |     | 6.4.4     | Summary  | 132 |
|   |     | 6.4.5     | Multiwavelength anomalous diffraction phasing              | 133 |
|   |     | 6.4.6     | Anomalous scattering and the hand problem                  | 135 |
|   |     | 6.4.7     | Direct phasing: Application of methods from small-molecule |     |
|   |     |           | crystallography  | 135 |
|   | 6.5 | Molecu    | lar replacement: Related proteins as phasing models        | 136 |
|   |     | 6.5.1     | Introduction   | 136 |
|   |     | 6.5.2     | Isomorphous phasing models                                 | 137 |
|   |     | 6.5.3     | Nonisomorphous phasing models                              | 139 |
|   |     | 6.5.4     | Separate searches for orientation and location             | 139 |
|   |     | 6.5.5     | Monitoring the search                                      | 141 |
|   |     | 6.5.6     | Summary of molecular replacement                           | 143 |
|   | 6.6 | Iterativo | e improvement of phases (preview of Chapter 7)             | 143 |
| 7 | Obt | aining a  | nd Judging the Molecular Model                             | 145 |
|   | 7.1 | Introdu   | ction  | 145 |
|   | 7.2 | lterative | e improvement of maps and models-overview                  | 146 |

|   | 7.3 | First m   | aps  | 149 |
|---|-----|-----------|--|-----|
|   |     | 7.3.1     | Resources for the first map                              | 149 |
|   |     | 7.3.2     | Displaying and examining the map                         | 150 |
|   |     | 7.3.3     | Improving the map  | 151 |
|   | 7.4 | The Mo    | odel becomes molecular                                   | 153 |
|   |     | 7.4.1     | New phases from the molecular model                      | 153 |
|   |     | 7.4.2     | Minimizing bias from the model                           | 154 |
|   |     | 7.4.3     | Map fitting  | 156 |
|   | 7.5 | Structu   | re refinement  | 159 |
|   |     | 7.5.1     | Least-squares methods                                    | 159 |
|   |     | 7.5.2     | Crystallographic refinement by least squares             | 160 |
|   |     | 7.5.3     | Additional refinement parameters                         | 161 |
|   |     | 7.5.4     | Local minima and radius of convergence                   | 162 |
|   |     | 7.5.5     | Molecular energy and motion in refinement                | 163 |
|   |     | 7.5.6     | Bayesian methods: Ensembles of models                    | 164 |
|   | 7.6 | Conver    | gence to a final model                                   | 168 |
|   |     | 7.6.1     | Producing the final map and model                        | 168 |
|   |     | 7.6.2     | Guides to convergence                                    | 171 |
|   | 7.7 | Sharing   | g the model  | 173 |
|   |     |           |  |     |
| 8 | A U | ser's Gu  | iide to Crystallographic Models                          | 179 |
|   | 8.1 | Introdu   | ction  | 179 |
|   | 8.2 | Judging   | g the quality and usefulness of the refined model        | 181 |
|   |     | 8.2.1     | Structural parameters                                    | 181 |
|   |     | 8.2.2     | Resolution and precision of atomic positions             | 183 |
|   |     | 8.2.3     | Vibration and disorder                                   | 185 |
|   |     | 8.2.4     | Other limitations of crystallographic models             | 187 |
|   |     | 8.2.5     | Online validation tools: Do it yourself!                 | 189 |
|   |     | 8.2.6     | Summary  | 192 |
|   | 8.3 | Readin    | g a crystallography paper                                | 192 |
|   |     | 8.3.1     | Introduction   | 192 |
|   |     | 8.3.2     | Annotated excerpts of the preliminary (8/91) paper       | 193 |
|   |     | 8.3.3     | Annotated excerpts from the full structure-determination |     |
|   |     |           | (4/92) paper   | 198 |
|   | 8.4 | Summa     | ary  | 209 |
|   |     |           |  |     |
| 9 | Oth | er Diffra | action Methods   | 211 |
|   | 9.1 | Introdu   | ction  | 211 |
|   | 9.2 | Fiber d   | iffraction   | 211 |
|   | 9.3 | Diffrac   | tion by amorphous materials (scattering)                 | 219 |
|   | 9.4 | Neutro    | n diffraction  | 222 |
|   | 9.5 | Electro   | n diffraction and cryo-electron microscopy               | 227 |
|   | 9.6 | Laue d    | iffraction and time-resolved crystallography             | 231 |
|   | 9.7 | Summa     | ary  | 235 |

xii

| 10 | Othe                            | er Kinds  | of Macromolecular Models                                  | 237 |  |  |
|----|---------------------------------|-----------|---|-----|--|--|
|    | 10.1                            | Introduc  | ction   | 237 |  |  |
|    | 10.2                            | NMR m     | odels   | 238 |  |  |
|    |                                 | 10.2.1    | Introduction  | 238 |  |  |
|    |                                 | 10.2.2    | Principles  | 239 |  |  |
|    |                                 | 10.2.3    | Assigning resonances                                      | 251 |  |  |
|    |                                 | 10.2.4    | Determining conformation                                  | 252 |  |  |
|    |                                 | 10.2.5    | PDB files for NMR models                                  | 257 |  |  |
|    |                                 | 10.2.6    | Judging model quality                                     | 257 |  |  |
|    | 10.3                            | Homolo    | gy models   | 259 |  |  |
|    |                                 | 10.3.1    | Introduction  | 259 |  |  |
|    |                                 | 10.3.2    | Principles  | 260 |  |  |
|    |                                 | 10.3.3    | Databases of homology models                              | 263 |  |  |
|    |                                 | 10.3.4    | Judging model quality                                     | 265 |  |  |
|    | 10.4                            | Other th  | eoretical models  | 267 |  |  |
| 11 | <b>T</b> I                      | e. C(     | 1 to Manager I.c.   | 2(0 |  |  |
| 11 | 1001                            | S IOP Stu | dying Macromolecules                                      | 269 |  |  |
|    | 11.1                            | Commut    | armodele of molecules                                     | 209 |  |  |
|    | 11.2                            |           | Two dimensional images from apardinates                   | 209 |  |  |
|    |                                 | 11.2.1    | Into three dimensions: Basic modeling operations          | 209 |  |  |
|    |                                 | 11.2.2    | Three dimensional display and perception                  | 270 |  |  |
|    |                                 | 11.2.3    | Types of graphical models                                 | 272 |  |  |
|    | 11.3                            | Touring   | a molecular modeling program                              | 275 |  |  |
|    | 11.5                            | 11 3 1    | Importing and exporting coordinate files                  | 275 |  |  |
|    |                                 | 11.3.1    | Loading and saving models                                 | 270 |  |  |
|    |                                 | 11.3.2    | Viewing models  | 278 |  |  |
|    |                                 | 11.3.5    | Editing and labeling the display                          | 280 |  |  |
|    |                                 | 11.3.1    | Coloring  | 281 |  |  |
|    |                                 | 11.3.5    | Measuring   | 281 |  |  |
|    |                                 | 11.3.7    | Exploring structural change                               | 282 |  |  |
|    |                                 | 11.3.8    | Exploring the molecular surface                           | 282 |  |  |
|    |                                 | 11.3.9    | Exploring intermolecular interactions: Multiple models    | 286 |  |  |
|    |                                 | 11.3.10   | Displaying crystal packing                                | 287 |  |  |
|    |                                 | 11.3.11   | Building models from scratch.                             | 287 |  |  |
|    |                                 | 11.3.12   | Scripts and macros: Automating routine structure analysis | 287 |  |  |
|    | 11.4                            | Other to  | bols for studying structure                               | 288 |  |  |
|    |                                 | 11.4.1    | Tools for structure analysis and validation               | 288 |  |  |
|    |                                 | 11.4.2    | Tools for modeling protein action                         | 290 |  |  |
|    | 11.5                            | Final no  | te  | 291 |  |  |
|    |                                 |           |   |     |  |  |
| Ap | AppendixViewing Stereo Images29 |           |   |     |  |  |

#### Index

295

This Page Intentionally Left Blank

### Preface to the Third Edition

►

Three prefaces make quite a moat to dig around this little castle of crystallography, so if you are tempted to get inside more quickly by skipping the introductions, at least take a quick look at the Preface to the First Edition, which still stands as the best guide to my aims in writing this book, and to your most efficient use of it. The second and this third preface are sort of like release notes for new versions of software. They are mostly about changes from previous editions. In brief, in the first edition, I taught myself the basics of crystallography by writing about it, drawing on a year or two of sabbatical experience in the field, preceded by quite a few years of enthusiastic sideline observation. In the second edition, I added material on other diffraction methods (neutron and electron diffraction, for instance) and other kinds of models (NMR and homology), and updated the crystallography only superficially. This time, the main subject, macromolecular crystallography, got almost all of my attention, and I hope the result is clearer, more accurate, and more up-to-date. One thing for sure, it's more colorful. Modern publishing methods have made color more affordable, and I found it very liberating to use color wherever I thought it would make illustrations easier to understand.

Just before writing this edition, I took a course, "X-Ray Methods in Structural Biology," at Cold Spring Harbor Laboratory. Professor David Richardson of Duke University, one of many accomplished crystallographers who contributed to this excellent course, seemed surprised to find me among the students there. When I told him I was looking for help in updating my book, he quickly offered this advice: "After taking this course, you will be tempted to complicate your book. Don't." I tried to keep David's words in mind as I worked on this edition. My main goal was to make the crystallography chapters more timely, accurate, and clear, by weeding out withered ideas and methods, replanting with descriptions of important new developments, culling out errors that readers of previous editions kindly took the trouble to point out, and adding only those new ideas and details

#### Preface to the Third Edition

that will truly help you to get a feeling for how crystallography produces models of macromolecules.

So what is new in crystallography since the last edition? First of all, it is faster than ever. Three developments—more powerful multiwavelength X-ray sources, low-temperature crystallography, and fast molecular biology methods for producing just about any protein and abundant variants of it-have set off an explosion of new structures. Automation has reached into every nook and cranny of the field, to the point that "high-throughput" crystallography is giving us models of proteins faster than we can figure out their functions. Just now I searched the Protein Data Bank (PDB), the world's primary repository for macromolecular models, for entries in which the protein function is listed as "unknown." I found almost 800 entries, many also marked "structural genomics." This means that the structures were determined as part of sweeping efforts, a prominent current one called the Protein Structure Initiative, to determine the structure of every protein in sight. Well, not quite; a research group participating in this effort usually focuses on a specific organism, like the tuberculosis bacterium, and works to determine the structures of proteins from every open reading frame (ORF) in its genome. For the first five years of this initiative, participating groups emphasized developing the technology to automate every step of structure determination: expressing and purifying the proteins, crystallizing them, collecting X-ray data, solving the structures, and refining the models. As I write these words, they are just beginning to turn their attention to cranking out new structures, although there are debates about whether the technology is ready for mass production. The goal of the Protein Structure Initiative is 10,000 structures by 2010, and even if the initiative falls a few thousand short, high-throughput crystallography is here to stay. You might not need to determine the structure of that protein whose action you just detected for the first time. It may already be in the Protein Data Bank, marked "structural genomics, unknown function."

Second, if the structure of that new protein of yours is not already lurking in the PDB, you might be able to determine it yourself. Methods of crystallization, data collection, and structure determination are more transparent and user-friendly than ever. Some of my fellow students at Cold Spring Harbor had already determined protein structures before they arrived, guided by modern crystallization screens, automated data collection, fast new software, and usually, a post-doctoral colleague with some crystallography experience. Now they wanted to know what goes on under the hood, in case a future venture stalls and requires an expert mechanic to make a few adjustments. If the next steps in your research would profit from your knowing the structure of a new protein, consider adding crystallography to your research skills. It's no longer necessary to make it your whole career.

Third, even if you never do crystallography, you are in a better position than ever to use models wisely. Powerful software and online tools allow you to make sound decisions about whether a model will support the conclusions you would like to draw from it, and with greater ease and clarity than ever. Today's validation

#### Preface to the Third Edition

tools can tell you a great deal about model quality, even if the original model publication is very sketchy on experimental methods and results.

Although the pace of crystallography is quickening toward mass production, I still wrote this edition about crystallography the old-fashioned way, one model at a time, with attention to the details of every step. Why? Because these details are essential to understanding crystallography and to assessing the strengths and limitations of each model. If you know the whole story, from purified protein to refined model, then you have a better understanding of the model and all that it might tell you. And if you try crystallography yourself, you will know something about the decisions the software is making for you, and when to ask if there are alternative routes to a model, perhaps better ones.

Many people helped me with this edition. At the top of the list are the instructors at the Cold Spring Harbor course who, for sixteen years, have offered what many crystallographers tout as the best classroom and hands-on diffraction training session on the planet—2.5 weeks, 9 AM to 9 PM, packed with labs, lectures, and computer tutorials, with homework for your spare time. The course gave me great confidence in choosing what to keep, what to revise, and what to throw out. The four organizers and main instructors, Bill Furey, Gary Gilliland, Alex McPherson, and Jim Pflugrath, were patient, helpful, and brimming with good ideas about how to do and teach crystallography. My fifteen CSH classmates (the oldest among them about half my age) were friendly, helpful, and inspirational. It was sad to realize that my presence in the course had displaced another one like them.

Thanks also to readers who pointed out errors in the first two editions, and to reviewers for their careful readings and helpful suggestions. Thanks to USM colleagues for granting me a sabbatical leave for this project (again!). Thanks to my wife, Pam, for proofreading, editing, and helpful suggestions on text and figures. I had to pay her, but she finally read my book. Thanks to the staff at Elsevier/Academic Press for guiding my words and pictures through the international maze of operations needed to get a book to market, and especially to Jeremy Hayhurst for talking me into doing this again, and to Jeff Freeland for overseeing production. Finally, thanks to all my students for constantly reminding me that teachers, whether they teach by lecturing, writing books, or building web pages, have more fun than people.

> Gale Rhodes Portland, Maine May 2005

This Page Intentionally Left Blank

### Preface to the Second Edition

The first edition of this book was hardly off the press before I was kicking myself for missing some good bets on how to make the book more helpful to more people. I am thankful that heartening acceptance and wide use of the first edition gave me another crack at it, even before much of the material started to show its age. In this new edition, I have updated the first eight chapters in a few spots and cleaned up a few mistakes, but otherwise those chapters, the soul of this book's argument, are little changed. I have expanded and modernized the last chapter, on viewing and studying models with computers, bringing it up-to-date (but only fleetingly, I am sure) with the cyberworld to which most users of macromolecular models now turn to pursue their interests, and with today's desktop computers—sleek, friendly, cheap, and eminently worthy successors to the five-figure workstations of the eighties.

My main goal, as outlined in the Preface to the First Edition, which appears herein, is the same as before: to help you see the logical thread that connects those mysterious diffraction patterns to the lovely molecular models you can display and play with on your personal computer. An equally important aim is to inform you that not all crystallographic models are perfect and that cartoon models do not exhaust the usefulness of crystallographic analysis. Often there is both less and more than meets the eye in a crystallographic model.

So what is new here? Two chapters are entirely new. The first one is "Other Diffraction Methods." In this chapter (the one I should have thought of the first time), I use your new-found understanding of X-ray crystallography to build an overview of other techniques in which diffraction gives structural clues. These methods include scattering of light, X-rays, and neutrons by powders and solutions; diffraction by fibers; crystallography using neutrons and electrons; and time-resolved crystallography using many X-ray wavelengths at the same time. These methods sound forbidding, but their underlying principles are precisely the same as those that make the foundation of single-crystal X-ray crystallography.

#### Preface to the Second Edition

The need for the second new chapter, "Other Types of Models," was much less obvious in 1992, when crystallography still produced most of the new macromolecular models. This chapter acknowledges the proliferation of such models from methods other than diffraction, particularly NMR spectroscopy and homology modeling. Databases of homology models now dwarf the Protein Data Bank, where all publicly available crystallographic and NMR models are housed. Nuclear magnetic resonance has been applied to larger molecules each year, with further expansion just a matter of time. Users must judge the quality of *all* macromolecular models, and that task is very different for different kinds of models. By analogies with similar aids for crystallographic models, I provide guidance in quality control, with the hope of making you a prudent user of models from all sources.

Neither of the new chapters contains full or rigorous treatments of these "other" methods. My aim is simply to give you a useful feeling for these methods, for the relationship between data and structures, and for the pitfalls inherent in taking any model too literally.

By the way, some crystallographers and NMR spectroscopists have argued for using the term *structure* to refer to the results of experimental methods, such as X-ray crystallography and NMR, and the term *model* for theoretical models such as homology models. To me, molecular structure is a book forever closed to our direct view, and thus never completely knowable. Consequently, I am much more comfortable with the term *model* for all of the results of attempts to know molecular structure. I sometimes refer loosely to a model as a *structure* and to the process of constructing and refining models as *structure determination*, but in the end, no matter what the method, we are trying to construct models that agree with, and explain, what we know from experiments that are quite different from actually looking at structure. So in my view, models, experimental or theoretical (an imprecise distinction itself), represent the best we can do in our diverse efforts to know molecular structure.

Many thanks to Nicolas Guex for giving to me and to the world a glorious free tool for studying proteins—Swiss-PdbViewer, since renamed DeepView— along with plenty of support and encouragement for bringing macromolecular modeling to my undergraduate biochemistry students; for his efforts to educate me about homology modeling; for thoughtfully reviewing the sections on homology modeling; and for the occasional box of liqueur-loaded Swiss chocolates (whoa!). Thanks to Kevin Cowtan, who allowed me to adapt some of the clever ideas from his *Book of Fourier* to my own uses, and who patiently computed image after image as I slowly iterated toward the final product. Thanks to Angela Gronenbom, Duncan McRee, and John Ricci for thorough, thoughtful, and helpful reviews of the manuscript. Thanks to Jonathan Cooper and Martha Teeter, who found and reported subtle and interesting errors lurking within figures in the first edition. Thanks to all those who provided figures—you are acknowledged alongside the fruits of your labors. Thanks to Emelyn Eldredge at Academic Press for inducing me to tiptoe once more through the minefields of Microsoft Word to update this

#### Preface to the Second Edition

little volume, and to Joanna Dinsmore for a smooth trip through production. Last and most, thanks to Pam for generous support, unflagging encouragement, and amused tolerance for over a third of a century. Time certainly does fly when we're having fun.

> Gale Rhodes Portland, Maine March 1999

This Page Intentionally Left Blank

### Preface to the First Edition

►

Most texts that treat biochemistry or proteins contain a brief section or chapter on protein crystallography. Even the best of such sections are usually mystifying—far too abbreviated to give any real understanding. In a few pages, the writer can accomplish little more than telling you to have faith in the method. At the other extreme are many useful treatises for the would-be, novice, or experienced crystallographer. Such accounts contain all the theoretical and experimental details that practitioners must master, and for this reason, they are quite intimidating to the noncrystallographer. This book lies in the vast and heretofore empty region between brief textbook sections on crystallographer. I hope there is just enough here to help the noncrystallographer understand where crystallographic models come from, how to judge their quality, and how to glean additional information that is not depicted in the model but is available from the crystallographic study that produced the model.

This book should be useful to protein researchers in all areas; to students of biochemistry in general and of macromolecules in particular; to teachers as an auxiliary text for courses in biochemistry, biophysical methods, and macromolecules; and to anyone who wants an intellectually satisfying understanding of how crystallographers obtain models of protein structure. This understanding is essential for intelligent use of crystallographic models, whether that use is studying molecular action and interaction, trying to unlock the secrets of protein folding, exploring the possibilities of engineering new protein functions, or interpreting the results of chemical, kinetic, thermodynamic, or spectroscopic experiments on proteins. Indeed, if you use protein models without knowing how they were obtained, you may be treading on hazardous ground. For instance, you may fail to use available information that would give you greater insight into the molecule and its action. Or worse, you may devise and publish a detailed molecular explanation based on a structural feature that is quite uncertain. Fuller understanding of the strengths

#### Preface to the First Edition

and limitations of crystallographic models will enable you to use them wisely and effectively.

If you are part of my intended audience, I do not believe you need to know, or are likely to care about, all the gory details of crystallographic methods and all the esoterica of crystallographic theory. I present just enough about methods to give you a feeling for the experiments that produce crystallographic data. I present somewhat more theory, because it underpins an understanding of the nature of a crystallographic model. I want to help you follow a logical thread that begins with diffraction data and ends with a colorful picture of a protein model on the screen of a graphics computer. The novice crystallographer, or the student pondering a career in crystallography, may find this book a good place to start, a means of seeing if the subject remains interesting under closer scrutiny. But these readers will need to consult more extensive works for fine details of theory and method. I hope that reading this book makes those texts more accessible. I assume that you are familiar with protein structure, at least at the level presented in an introductory biochemistry text.

I wish I could teach you about crystallography without using mathematics, simply because so many readers are apt to throw in the towel upon turning the page and finding themselves confronted with equations. Alas (or hurrah, depending on your mathematical bent), the real beauty of crystallography lies in the mathematical and geometric relationships between diffraction data and molecular images. I attempt to resolve this dilemma by presenting no more math than is essential and taking the time *to explain in words what the equations imply*. Where possible, I emphasize geometric explanations over equations.

If you turn casually to the middle of this book, you will see some forbidding mathematical formulas. Let me assure you that I move to those bushy statements step-by-step from nearby clearings, making minimum assumptions about your facility and experience with math. For example, when I introduce periodic functions, I tell you how the simplest of such functions (sines and cosines) "work," and then I move slowly from that clear trailhead into the thicker forest of complicated wave equations that describe X-rays and the molecules that diffract them. When I first use complex numbers, I define them and illustrate their simplest uses and representations, sort of like breaking out camping gear in the dry safety of a garage. Then I move out into real weather and set up a working camp, showing how the geometry of complex numbers reveals essential information otherwise hidden in the data. My goal is to help you see the relationships implied by the mathematics, not to make you a calculating athlete. My ultimate aim is to prove to you that the structure of molecules really does lie lurking in the crystallographic data-that, in fact, the information in the diffraction pattern implies a unique structure. I hope thereby to remove the mystery about how structures are coaxed from data.

If, in spite of these efforts, you find yourself flagging in the most technical chapters (4 and 7), please do not quit. I believe you can follow the arguments of these chapters, and thus be ready for the take-home lessons of Chapters 8 and 11, even if the equations do not speak clearly to you. Jacob Bronowski once described the verbal argument in mathematical writing as analogous to melody

#### xxiv

#### Preface to the First Edition

in music, and thus a source of satisfaction in itself. He likened the equations to musical accompaniment that becomes more satisfying with repeated listening. If you follow and retain the melody of arguments and illustrations in Chapters 4 through 7, then the last chapters and their take-home lessons should be useful to you.

I aim further to enable you to read primary journal articles that announce and present new protein structures, including the arcane sections on experimental methods. In most scientific papers, experimental sections are directed primarily toward those who might use the same methods. In crystallographic papers, however, methods sections contain information from which the quality of the model can be roughly judged. This judgment should affect your decision about whether to obtain the model and use it, and whether it is good enough to serve as a guide in drawing the kinds of conclusions you hope to draw. In Chapter 8, to review many concepts, as well as to exercise your new skills, I look at and interpret experimental details in literature reports of a recent structure determination.

Finally, I hope you read this book for pleasure—the sheer pleasure of turning the formerly incomprehensible into the familiar. In a sense, I am attempting to share with you my own pleasure of the past ten years, after my mid-career decision to set aside other interests and finally see how crystallographers produce the molecular models that have been the greatest delight of my teaching. Among those I should thank for opening their labs and giving their time to an old dog trying to learn new tricks are Professors Leonard J. Banaszak, Jens Birktoft, Jeffry Bolin, John Johnson, and Michael Rossman.

I would never have completed this book without the patience of my wife, Pam, who allowed me to turn part of our home into a miniature publishing company, nor without the generosity of my faculty colleagues, who allowed me a sabbatical leave during times of great economic stress at the University of Southern Maine. Many thanks to Lorraine Lica, my Acquisitions Editor at Academic Press, who grasped the spirit of this little project from the very beginning and then held me and a full corps of editors, designers, and production workers accountable to that spirit throughout.

Gale Rhodes Portland, Maine August 1992

#### Phase

These still days after frost have let down the maple leaves in a straight compression to the grass, a slight wobble from circular to

the east, as if sometime, probably at night, the wind's moved that way—surely, nothing else could have done it, really eliminating the *as* 

*if*, although the *as if* can nearly stay since the wind may have been a big, slow one, imperceptible, but still angling

off the perpendicular the leaves' fall: anyway, there was the green-ribbed, yellow, flat-open reduction: I just now bagged it up.

A. R. Ammons<sup>1</sup>

<sup>&</sup>lt;sup>1</sup> "Phase," from *The Selected Poems, Expanded Edition* by A. R. Ammons. Copyright @ 1987, 1977, 1975, 1974, 1972, 1971, 1970, 1966, 1965, 1964, 1955 by A. R. Ammons. Reprinted by permission of W. W. Norton & Company, Inc.

### ► Chapter 1

### Model and Molecule

Proteins perform many functions in living organisms. For example, some proteins regulate the expression of genes. One class of gene-regulating proteins contains structures known as *zinc fingers*, which bind directly to DNA. Figure 1.1*a* shows a complex composed of a double-stranded DNA molecule and three zinc fingers from the mouse protein Zif268 (PDB 1zaa).

The protein backbone is shown as a yellow ribbon. The two DNA strands are red and blue. Zinc atoms, which are complexed to side chains in the protein, are green. The green dotted lines near the top center indicate two hydrogen bonds in which nitrogen atoms of arginine-18 (in the protein) share hydrogen atoms with nitrogen and oxygen atoms of guanine-10 (in the DNA), an interaction that holds the sharing atoms about 2.8 Å apart. Studying this complex with modern graphics software, you could zoom in, as in Fig. 1.1b, measure the hydrogen-bond lengths, and find them to be 2.79 and 2.67 Å. From a closer study, you would also learn that all of the protein–DNA interactions are between protein *side chains* and DNA *bases*; the protein backbone does not come in contact with the DNA. You could go on to discover all the specific interactions between side chains of Zif268 and base pairs of DNA. You could enumerate the additional hydrogen bonds and other contacts that stabilize this complex and cause Zif268 to recognize a specific sequence of bases in DNA. You might gain some testable insights into how the protein finds the correct DNA sequence amid the vast amount of DNA in the nucleus of a cell. The structure might also lead you to speculate on how alterations in the sequence of amino acids in the protein might result in affinity for different DNA sequences, and thus start you thinking about how to design other DNA-binding proteins.

Now look again at the preceding paragraph and examine its *language* rather than its content. The language is typical of that in common use to describe molecular structure and interactions as revealed by various experimental methods, including single-crystal X-ray crystallography, the primary subject of this book. In fact, this

Chapter 1 Model and Molecule



**Figure 1.1** (*a*) Divergent stereo image of Zif268/DNA complex (N. P. Pavletich and C. O. Pabo, *Science* **252**, 809, 1991). (*b*) Detail showing hydrogen bonding between arginine-18 of the protein and guanine-10 of the DNA. Atomic coordinates for preparing this display were obtained from the Protein Data Bank (PDB), which is described in Chapter 7. The PDB file code is 1zaa. To allow easy access to all models shown in this book, I provide file codes in this format: PDB 1zaa. Image created by DeepView (formerly called Swiss-PdbViewer), rendered by POV-Ray. To obtain these programs, see the CMCC home page at http://www.usm.maine.edu/~rhodes/CMCC/index.html. For help with viewing stereo images, see Appendix, page 293.

language is shorthand for more precise but cumbersome statements of what we learn from structural studies.

First, Fig. 1.1, of course, shows not molecules, but *models* of molecules, in which structures and interactions are *depicted*, not shown. Second, in this specific case, the models are of molecules not in solution, but in the crystalline state, because the models are derived from analysis of X-ray diffraction by crystals of the Zif268/DNA complex. As such, these models depict the average structure of somewhere between  $10^{13}$  and  $10^{15}$  complexes throughout the crystals that

#### Chapter 1 Model and Molecule

were studied. In addition, the structures are averaged over the time of the X-ray experiment, which may range from minutes to days.

To draw the conclusions found in the first paragraph requires bringing additional knowledge to bear upon the graphics image, including a more precise knowledge of exactly what we learn from X-ray analysis. The same could be said for structural models derived from spectroscopic data or any other method. In short, the graphics image itself is incomplete. It does not reveal things we may know about the complex from other types of experiments, and *it does not even reveal all that we learn from X-ray crystallography*.

For example, how accurately are the relative positions of atoms known? Are the hydrogen bonds precisely 2.79 and 2.67 Å long, or is there some tolerance in those figures? Is the tolerance large enough to jeopardize the conclusion that hydrogen bonds join these atoms? Further, do we know anything about how rigid this complex is? Do parts of these molecules vibrate, or do they move with respect to each other? Still further, in the aqueous medium of the cell, does this complex have the same structure as in the crystal, which is a solid? As we examine this model, are we really gaining insight into cellular processes? Two final questions may surprise you: First, does the model fully account for the chemical composition of the crystal? In other words, are any of the known contents of the crystal missing from the model? Second, does the crystallographic data suggest additional crystal contents that have not been identified, and thus are not shown in the model?

The answers to these questions are not revealed in the graphics image, which is more akin to a cartoon than to a molecule. Actually, the answers vary from one model to the next, and from one region of a model to another region, but they are usually available to the user of crystallographic models. Some of the answers come from X-ray crystallography itself, so the crystallographer does not miss or overlook them. They are simply less accessible to the noncrystallographer than is the graphics image.

Molecular models obtained from crystallography are in wide use as tools for revealing molecular details of life processes. Scientists use models to learn how molecules "work": how enzymes catalyze metabolic reactions, how transport proteins load and unload their molecular cargo, how antibodies bind and destroy foreign substances, and how proteins bind to DNA, perhaps turning genes on and off. It is easy for the user of crystallographic models, being anxious to turn otherwise puzzling information into a mechanism of action, to treat models as everyday objects seen as we see clouds, birds, and trees. But the informed user of models sees more than the graphics image, recognizing it as a static depiction of dynamic objects, as the average of many similar structures, as perhaps lacking parts that are present in the crystal but not revealed by the X-ray analysis, as perhaps failing to show as-yet unidentified crystal contents, and finally, as a fallible interpretation of data. The informed user knows that the crystallographic model is richer than the cartoon.

In the following chapters, I offer you the opportunity to become an informed user of crystallographic models. Knowing the richness and limitations of models requires an understanding of the relationship between data and structure. In Chapter 2, I give an overview of this relationship. In Chapters 3 through 7, the heart of the crystallography in this book, I simply expand Chapter 2 in enough detail to produce an intact chain of logic stretching from diffraction data to final model. Topics come in roughly the same order as the tasks that face a crystallographer pursuing an important structure.

As a practical matter, informed use of a model requires evaluating its quality, which may entail using online *model validation tools* to assess model quality, as well as reading the crystallographic papers and data files that report the new structure, in order to extract from them criteria of model quality. In Chapter 8, I discuss these criteria and provide guided exercises in extracting them from model files themselves and from the literature. Chapter 8 includes an annotated version of a published structure determination and its supporting data, as well as an introduction to online validation tools. Equipped with the background of previous chapters and experienced with the real-world exercises of using validation tools and taking a guided tour through a recent publication, you should be able to read new structure publications in the scientific literature, understand how the structures were obtained, and be aware of just what is known—and what is still unknown—about the molecules under study. Then you should be better equipped to use models wisely.

Chapter 9, "Other Kinds of Macromolecular Methods," builds on your understanding of X-ray crystallography to help you understand other methods in which diffraction provides insights into the structure of large molecules. These methods include fiber diffraction, neutron diffraction, electron diffraction, and various forms of X-ray spectroscopy. These methods often seem very obscure, but their underlying principles are similar to those of X-ray crystallography.

In Chapter 10, "Other Kinds of Macromolecular Models," I discuss alternative methods of structure determination: NMR spectroscopy and various forms of theoretical modeling. Just like crystallographic models, NMR and theoretical models are sometimes more, sometimes less, than meets the eye. A brief description of how these models are obtained, along with some analogies among criteria of quality for various types of models, can help make you a wiser user of all types of models.

For new or would-be users of models, I present in Chapter 11 an introduction to molecular modeling, demonstrating how modern graphics programs allow users to display and manipulate models and to perform powerful structure analysis, as well as model validation, on desktop computers. I also provide information on how to use the World Wide Web to obtain graphics programs and learn how to use them. Finally, I introduce you to the Protein Data Bank (PDB), a World Wide Web resource from which you can obtain most of the available macromolecular models.

There is an additional chapter that does not lie between the covers of this book. It is the Crystallography Made Crystal Clear (CMCC) home page on the World Wide Web at www.usm.maine.edu/~rhodes/CMCC. This web site is devoted to making sure that you can find all the Internet resources mentioned here. Because even major Internet resources and addresses may change (the Protein Data Bank

#### Chapter 1 Model and Molecule

moved while I was writing the second edition of this book), I include only one web address in this book. For all web resources that I describe, I refer you to the CMCC home page. At that web address, I maintain links to all resources mentioned here or, if they disappear or change markedly, to new ones that serve the same or similar functions. For easy reference, the address of the CMCC home page is shown on the cover and title page of this book.

Today's scientific textbooks and journals are filled with stories about the molecular processes of life. The central character in these stories is often a protein or nucleic acid molecule, a thing never seen in action, never perceived directly. We see models of molecules in books and on computer screens, and we tend to treat them as everyday objects accessible to our normal perceptions. In fact, models are hard-won products of technically difficult data collection and powerful but subtle data analysis. And they are richer and more informative than any single image, or even a rotating computer image, can convey. This book is concerned with where our models of structure come from and how to use them wisely. This Page Intentionally Left Blank

### ► Chapter 2

### An Overview of Protein Crystallography

### 2.1 Introduction

The most common experimental means of obtaining a detailed model of a large molecule, allowing the resolution of individual atoms, is to interpret the diffraction of X-rays from many identical molecules in an ordered array like a crystal. This method is called *single-crystal X-ray crystallography*. As of January 2005, the Protein Data Bank (PDB), the world's largest repository of macromolecular models obtained from experimental data (called experimental models), contains roughly 25,000 protein and nucleic-acid models determined by X-ray crystallography. In addition, the PDB holds roughly 4500 models, mostly proteins of fewer than 200 residues, that have been solved by nuclear magnetic resonance (NMR) spectroscopy, which provides a model of the molecule in solution, rather than in the crystalline state. (Because many proteins appear in multiple forms-for example, wild types and mutants, or solo and also as part of protein-ligand or multiprotein complexes-the number of unique proteins represented in the PDB is only a fraction of the almost 30,000 models.) Finally, there are theoretical models, either built by analogy with the structures of known proteins having similar sequence, or based on simulations of protein folding. (Theoretical models are available from databases other than the PDB.) All methods of obtaining models have their strengths and weaknesses, and they coexist happily as complementary methods. One of the goals of this book is to make users of *crystallographic* models aware of the strengths and weaknesses of X-ray crystallography, so that users expectations of the resulting models are in keeping with the limitations of crystallographic methods. Chapter 10 provides, in brief, complementary information about other types of models.

In this chapter, I provide a simplified overview of how researchers use the technique of X-ray crystallography to obtain models of macromolecules. Chapters 3 through 8 are simply expansions of the material in this chapter. In order to keep the language simple, I will speak primarily of proteins, but the concepts I describe apply to all macromolecules and macromolecular assemblies that possess ordered structure, including carbohydrates, nucleic acids, and nucleoprotein complexes like ribosomes and whole viruses.

#### 2.1.1 Obtaining an image of a microscopic object

When we see an object, light rays bounce off (are diffracted by) the object and enter the eye through the lens, which reconstructs an image of the object and focuses it on the retina. In a simple microscope, an illuminated object is placed just beyond one focal point of a lens, which is called the *objective* lens. The lens collects light diffracted from the object and reconstructs an image beyond the focal point on the opposite side of the lens, as shown in Fig. 2.1.

For a simple lens, the relationship of object position to image position in Fig. 2.1 is (OF)(IF') = (FL)(F'L). Because the distances FL and F'L are constants (but not necessarily equal) for a fixed lens, the distance OF is inversely proportional to the distance IF'. Placing the object just beyond the focal point F results in a magnified image produced at a considerable distance from F' on the other side of the the lens, which is convenient for viewing. In a compound microscope,



**Figure 2.1** Action of a simple lens. Rays parallel to the lens axis strike the lens and are refracted into paths passing through a focus (*F* or *F'*). Rays passing through a focus strike the lens and are refracted into paths parallel to the lens axis. As a result, the lens produces an image at *I* of an object at *O* such that (OF)(IF') = (FL)(F'L).

#### Section 2.1 Introduction

the most common type, an additional lens, the *eyepiece*, is added to magnify the image produced by the objective lens.

#### 2.1.2 Obtaining images of molecules

In order for the object to diffract light and thus be visible under magnification, the wavelength ( $\lambda$ ) of the light must be, roughly speaking, no larger than the object. Visible light, which is electromagnetic radiation with wavelengths of 400–700 nm (nm = 10<sup>-9</sup> m), cannot produce an image of individual atoms in protein molecules, in which bonded atoms are only about 0.15 nm or 1.5 angstroms (Å = 10<sup>-10</sup> m) apart. Electromagnetic radiation of this wavelength falls into the X-ray range, so X-rays are diffracted by even the smallest molecules. X-ray analysis of proteins seldom resolves the hydrogen atoms, so the protein models described in this book include elements on only the second and higher rows of the periodic table. The positions of all hydrogen atoms can be deduced on the assumption that bond lengths, bond angles, and conformational angles in proteins are just like those in small organic molecules.

Even though individual atoms diffract X-rays, it is still not possible to produce a focused image of a single molecule, for two reasons. First, X-rays cannot be focused by lenses. Crystallographers sidestep this problem by measuring the directions and strengths (intensities) of the diffracted X-rays and then using a computer to simulate an image-reconstructing lens. In short, the computer acts as the lens, computing the image of the object and then displaying it on a screen (Fig. 2.2).

Second, a single molecule is a very weak scatterer of X-rays. Most of the X-rays will pass through a single molecule without being diffracted, so the diffracted beams are too weak to be detected. Analyzing diffraction from crystals, rather than individual molecules, solves this problem. A crystal of a protein contains many ordered molecules in identical orientations, so each molecule diffracts identically, and the diffracted beams for all molecules augment each other to produce strong, detectable X-ray beams.

#### 2.1.3 A thumbnail sketch of protein crystallography

In brief, determining the structure of a protein by X-ray crystallography entails growing high-quality crystals of the purified protein, measuring the directions and intensities of X-ray beams diffracted from the crystals, and using a computer to simulate the effects of an objective lens and thus produce an image of the crystal's contents, like the small section of a molecular image shown in Fig. 2.3*a*. Finally, the crystallographer must *interpret* that image, which entails displaying it by computer graphics and building a molecular model that is consistent with the image (Fig. 2.3*b*).

The resulting model is often the only product of crystallography that the user sees. It is therefore easy to think of the model as a real entity that has been directly observed. In fact, our "view" of the molecule is quite indirect. Understanding just how the crystallographer obtains models of protein molecules from diffraction measurements is essential to fully understanding how to use models properly.


**Figure 2.2**  $\triangleright$  Crystallographic analogy of lens action. X-rays diffracted from the object are received and measured by a detector. The measurements are fed to a computer, which simulates the action of a lens to produce a graphics image of the object. Compare Fig. 2.2 with Fig. 2.1 and you will see that to magnify molecules, you merely have to replace the light bulb with a synchrotron X-ray source (175 feet in diameter), replace the glass lens with the equivalent of a 5- to 10-megapixel camera, and connect the camera output to a computer running some of the world's most complex and sophisticated software. Oh, yes, and you will need to spend somewhere between a few days and the rest of your life getting your favorite protein to form satisfactory crystals. No, it's not quite as simple as microscopy.

## 2.2 Crystals

#### 2.2.1 The nature of crystals

Under certain circumstances, many molecular substances, including proteins, solidify to form crystals. In entering the crystalline state from solution, individual molecules of the substance adopt one or a few identical orientations. The resulting crystal is an orderly three-dimensional array of molecules, held together by noncovalent interactions. Figure 2.4 depicts such a crystalline array of molecules.

The lines in the figure divide the crystal into identical *unit cells*. The array of points at the corners or vertices of unit cells is called the *lattice*. The unit cell is the smallest and simplest volume element that is completely representative of the whole crystal. If we know the exact contents of the unit cell, we can imagine the whole crystal as an efficiently packed array of many unit cells stacked beside and on top of each other, more or less like identical boxes in a warehouse.



**Figure 2.3** (*a*) Small section of molecular image displayed on a computer. (*b*) Image (*a*) is *interpreted* by building a molecular model to fit within the image. Computer graphics programs allow the crystallographer to add parts to the model and adjust their positions and conformations to fit the image. The protein shown here is adipocyte lipid binding protein (ALBP, PDB 1alb).

From crystallography, we obtain *an image of the electron clouds* that surround the molecules in the average unit cell in the crystal. We hope this image will allow us to locate all atoms in the unit cell. The location of an atom is usually given by a set of three-dimensional Cartesian coordinates, x, y, and z. One of the vertices (a lattice point or any other convenient point) is used as the origin of the unit cell's coordinate system and is assigned the coordinates x = 0, y = 0, and z = 0, usually written (0, 0, 0) (Fig. 2.5).

#### 2.2.2 Growing crystals

Crystallographers grow crystals of proteins by slow, controlled precipitation from aqueous solution under conditions that do not denature the protein. A number of substances cause proteins to precipitate. Ionic compounds (salts) precipitate proteins by a process called "salting out." Organic solvents also cause precipitation, but they often interact with hydrophobic portions of proteins and thereby denature them. The water-soluble polymer polyethylene glycol (PEG) is widely used



**Figure 2.4**  $\triangleright$  Six unit cells in a crystalline lattice. Each unit cell contains two molecules of alanine (hydrogen atoms not shown) in different orientations.



**Figure 2.5**  $\blacktriangleright$  One unit cell from Fig. 2.4. The position of an atom in the unit cell can be specified by a set of spatial coordinates *x*, *y*, *z*.

because it is a powerful precipitant and a weak denaturant. It is available in preparations of different average molecular masses, such as PEG 400, with average molecular mass of 400 daltons.

One simple means of causing slow precipitation is to add denaturant to an aqueous solution of protein until the denaturant concentration is just below that required to precipitate the protein. Then water is allowed to evaporate slowly, which gently raises the concentration of both protein and denaturant until precipitation occurs. Whether the protein forms crystals or instead forms a useless amorphous solid depends on many properties of the solution, including protein concentration, temperature, pH, and ionic strength. Finding the exact conditions to produce good crystals of a specific protein often requires many careful trials, and is perhaps more art than science. I will examine crystallization methods in Chapter 3.

## 2.3 Collecting X-ray data

Figure 2.6 depicts the collection of X-ray diffraction data. A crystal is mounted between an X-ray source and an X-ray detector. The crystal lies in the path of a narrow beam of X-rays coming from the source. The simplest source is an X-ray tube, and the simplest detector is X-ray film, which when developed exhibits dark spots where X-ray beams have impinged. These spots are called *reflections* because they emerge from the crystal as if reflected from planes of atoms.

Figure 2.7 shows the complex diffraction pattern of X-ray reflections produced on a detector by a protein crystal. Notice that the crystal diffracts the source beam into many discrete beams, each of which produces a distinct reflection on the film. The greater the intensity of the X-ray beam that reaches a particular position, the darker the reflection.



**Figure 2.6**  $\triangleright$  Crystallographic data collection. The crystal diffracts the source beam into many discrete beams, each of which produces a distinct spot (reflection) on the film. The positions and intensities of these reflections contain the information needed to determine molecular structures.



**Figure 2.7** Diffraction pattern from a crystal of the MoFe (molybdenum-iron) protein of the enzyme nitrogenase from *Clostridium pasteurianum*, recorded on film. Notice that the reflections lie in a regular pattern, but their intensities (darkness of spots) are highly variable. The hole in the middle of the pattern results from a small metal disk (*beam stop*) used to prevent the direct X-ray beam, most of which passes straight through the crystal, from destroying the center of the film. Photo courtesy of Professor Jeffrey Bolin.

An optical scanner precisely measures the position and the intensity of each reflection and transmits this information in digital form to a computer for analysis. The position of a reflection can be used to obtain the direction in which that particular beam was diffracted by the crystal. The intensity of a reflection is obtained by measuring the optical absorbance of the spot on the film, giving a measure of the strength of the diffracted beam that produced the spot. The computer program that reconstructs an image of the molecules in the unit cell requires these two parameters, the relative intensity and direction, for each diffracted beam that produces a reflection at the detector. The intensity is simply a number that tells how dark the reflection is in comparison to the others. The beam direction, as I will describe shortly, is specified by a set of three-dimensional coordinates h, k, and l for each reflection.

#### Section 2.4 Diffraction

Although film for data collection has almost completely been replaced by devices that feed diffraction data (positions and intensities of each reflection) directly into computers, I will continue, in this overview, to speak of the data as if collected on film because of the simplicity of that format, and because diffraction patterns are usually published in a form identical to their appearance on film. I will discuss modern methods of collecting data in Chapter 4.

## 2.4 Diffraction

#### 2.4.1 Simple objects

You can develop some visual intuition for the information available from X-ray diffraction by examining the diffraction patterns of simple objects like spheres or arrays of spheres (Figs. 2.8–2.11). Figure 2.8 depicts diffraction by a single sphere, shown in cross section on the left. The diffraction pattern, on the right, exhibits high intensity at the center, and smoothly decreasing intensity as the diffraction angle increases.<sup>1</sup>

For now, just accept the observation that diffraction by a sphere produces this pattern, and think of it as the diffraction signature of a sphere. In a sense, you are already equipped to do very simple structure determination; that is, you can now recognize a simple sphere by its diffraction pattern.



**Figure 2.8** ► Sphere (cross-section, on left) and its diffraction pattern (right). Images for Figs. 2.8–2.11 were generously provided by Dr. Kevin Cowtan.

<sup>&</sup>lt;sup>1</sup>The images shown in Figs. 2.8–2.11 are computed, rather than experimental, diffraction patterns. Computation of these patterns involves use of the Fourier transform (Section 2.6.5).

#### 2.4.2 Arrays of simple objects: Real and reciprocal lattices

Figure 2.9 depicts diffraction by spheres in a crystalline array, with a cross section of the crystalline lattice on the left, and its diffraction pattern on the right.

The diffraction pattern, like that produced by crystalline nitrogenase (Fig. 2.7), consists of reflections (spots) in an orderly array on the film. The spacing of the reflections varies with the spacing of the spheres in their array. Specifically, observe that although the lattice spacing of the crystal is smaller vertically, the diffraction spacing is smaller horizontally. In fact, there is a simple inverse relationship between the spacing of unit cells in the crystalline lattice, called the *real lattice*, and the spacing of reflections in the lattice on the film, which, because of its inverse relationship to the real lattice, is called the *reciprocal lattice*.

Because the real lattice spacing is inversely proportional to the spacing of reflections, crystallographers can calculate the dimensions, in angstroms, of the unit cell of the crystalline material from the spacings of the reciprocal lattice on the X-ray film (Chapter 4). The simplicity of this relationship is a dramatic example of how the macroscopic dimensions of the diffraction pattern are connected to the submicroscopic dimensions of the crystal.

#### 2.4.3 Intensities of reflections

Now look carefully at the intensities of the reflections in Fig. 2.9. Some are intense ("bright"), whereas others are weak or perhaps missing from the otherwise evenly spaced pattern. These variations in intensity contain important information. If you blur your eyes slightly while looking at the diffraction pattern, so that you cannot see individual spots, you will see the intensity pattern characteristic of diffraction by a sphere, with lower intensities farther from the center, as in Fig. 2.8. (You just



**Figure 2.9** Lattice of spheres (left) and its diffraction pattern (right). If you look at the pattern and blur your eyes, you will see the diffraction pattern of a sphere. The pattern is that of the average sphere in the real lattice, but it is sampled at the reciprocal lattice points.

#### Section 2.4 Diffraction

determined your first crystallographic structure.) The diffraction pattern of spheres in a lattice is simply the diffraction pattern of the average sphere in the lattice, but this pattern is incomplete. The pattern is *sampled* at points whose spacings vary inversely with real-lattice spacings. The pattern of varied intensities is that of the *average* sphere because all the spheres contribute to the observed pattern. To put it another way, the observed pattern of intensities is actually a superposition of the many identical diffraction patterns of all the spheres.

#### 2.4.4 Arrays of complex objects

This relationship between (1) diffraction by a single object and (2) diffraction by many identical objects in a lattice holds true for complex objects also. Figure 2.10 depicts diffraction by six spheres that form a planar hexagon, like the six carbon atoms in benzene. Notice the starlike six-fold symmetry of the diffraction pattern. Again, just accept this pattern as the diffraction signature of a hexagon of spheres. (Now you know enough to recognize *two* simple objects by their diffraction patterns.) Figure 2.11 depicts diffraction by these hexagonal objects in a lattice of the same dimensions as that in Fig. 2.9.

As before, the spacing of reflections varies reciprocally with lattice spacing, but if you blur your eyes slightly, or compare Figs. 2.10 and 2.11 carefully, you will see that the starlike signature of a single hexagonal cluster is present in Fig. 2.11. From these simple examples, you can see that the reciprocal-lattice spacing (the spacing of reflections in the diffraction pattern) is characteristic of (inversely related to) the spacing of identical objects in the crystal, whereas the reflection intensities are characteristic of the shape of the individual objects. From the reciprocal-lattice spacing in a diffraction pattern, we can compute the dimensions of the unit cell. From the intensities of the reflections, we can learn the shape of the individual molecules that compose the crystal. It is actually advantageous that the object's



#### **Figure 2.10** ► A planar hexagon of spheres (left) and its diffraction pattern (right).

|     |      | 2.0       |           |      |  |
|-----|------|-----------|-----------|------|--|
|     | 32   | 32        |           |      |  |
|     | 2.6  | <b>76</b> | <b>74</b> |      | a liter later and a strategy of the second s |
|     | 32   | 32        | 32        |      |  |
|     |      |           |           |      |  |
|     | 22   | 32        | 22        |      |  |
|     |      |           |           |      |  |
|     |      |           | 2.5       |      |  |
|     |      |           |           |      |  |
|     | 1.5  | 1.1       | 1.5       | 1.1  |  |
|     |      |           |           |      |  |
| 1.0 | - 63 | - 63      | - 63      | 10.3 |  |
|     |      |           |           |      |  |
| 1.2 | 1.2  | 6.3       | 10.0      | 2.3  |  |
|     |      |           |           |      |  |

**Figure 2.11**  $\blacktriangleright$  Lattice of hexagons (left) and its diffraction pattern (right). If you look at the pattern and blur your eyes, you will see the diffraction pattern of a hexagon. The pattern is that of the average hexagon in the real lattice, but it is sampled at the reciprocal lattice points.

diffraction pattern is sampled at reciprocal-lattice positions. This sampling reduces the number of intensity measurements we must take from the film and makes it easier to program a computer to locate and measure the intensities.

#### 2.4.5 Three-dimensional arrays

Unlike the two-dimensional arrays in these examples, a crystal is a threedimensional array of objects. If we rotate the crystal in the X-ray beam, a different cross section of objects will lie perpendicular to the beam, and we will see a different diffraction pattern. In fact, just as the two-dimensional arrays of objects I have discussed are cross sections of objects in the three-dimensional crystal, each two-dimensional array of reflections (each diffraction pattern recorded on film) is a cross section of a three-dimensional lattice of reflections. Figure 2.12 shows a hypothetical three-dimensional diffraction pattern, with the reflections that would be produced by all possible orientations of a crystal in the X-ray beam.

Notice that only one plane of the three-dimensional diffraction pattern is superimposed on the film. With the crystal in the orientation shown, reflections shown in the plane of the film (solid spots) are the only reflections that produce spots on the film. In order to measure the directions and intensities of all additional reflections (shown as hollow spots), the crystallographer must collect diffraction patterns from all unique orientations of the crystal with respect to the X-ray beam. The direct result of crystallographic data collection is a list of intensities for each point in the three-dimensional reciprocal lattice. This set of data is the raw material for determining the structures of molecules in the crystal.

(*Note*: The spatial relationship involving beam, crystal, film, and reflections is more complex than shown here. I will discuss the actual relationship in Chapter 4.)



**Figure 2.12** Crystallographic data collection, showing reflections measured at one particular crystal orientation (solid, on the film) and those that could be measured at other orientations (hollow, within the sphere but not on the film). Each reflection is located by its three-dimensional coordinates h, k, and l. The relationship between measured and unmeasured reflections is more complex than shown here (see Chapter 4).

#### 2.5

### Coordinate systems in crystallography

Each reflection can be assigned three coordinates or *indices* in the imaginary threedimensional space of the diffraction pattern. This space, the strange land where the reflections live, is called *reciprocal space*. Crystallographers usually use h, k, and l to designate the position of an individual reflection in the reciprocal space of the diffraction pattern. The central reflection (the round solid spot at the center of the film in Fig. 2.12) is taken as the origin in reciprocal space and assigned the coordinates (h, k, l) = (0, 0, 0), usually written hkl = 000. (The 000 reflection is not measurable because it is always obscured by X-rays that pass straight through the crystal, and are blocked by the beam stop.) The other reflections are assigned whole-number coordinates counted from this origin, so the indices h, k, and l are *integers*. Thus the parameters we can measure and analyze in the X-ray diffraction pattern are (1) the position hkl and (2) the intensity  $I_{hkl}$  of each reflection. The position of a reflection is related to the angle by which the diffracted beam diverges from the source beam. For a unit cell of known dimensions, the angle of divergence uniquely specifies the indices of a reflection, as I will show in Chapter 4.

Alternatively, actual distances, rather than reflection indices, can be measured in reciprocal space. Because the dimensions of reciprocal space are the inverse of dimensions in the real space of the crystal, distances in reciprocal space are expressed in the units Å<sup>-1</sup> (called *reciprocal angstroms*). Roughly speaking, the inverse of the reciprocal-space distance from the origin out to the most distant measurable reflections gives the potential resolution of the model that we can obtain from the data. So a crystal that gives measurable reflections out to a distance of 1/(3 Å) from the origin is said to yield a model with a resolution of 3 Å.

Crystallographers work back and forth between two different coordinate systems. I will review them briefly. The first system (see Fig. 2.5, p. 12) is the unit cell (real space), where an atom's position is described by its coordinates x, y, z. A vertex of the unit cell, or any other convenient position, is taken as the origin, with coordinates x, y, z = (0, 0, 0). Coordinates in real space designate real spatial positions within the unit cell. Real-space coordinates are usually given in angstroms or nanometers, or in fractions of unit cell dimensions. The second system (see Fig. 2.12, p. 19) is the three-dimensional diffraction pattern (reciprocal space), where a reflection's position is described by its indices *hkl*. The central reflection is taken as the origin with the index *hkl* = 000 (round black dot at center of sphere). The position of a reflection is designated by counting reflections from 000, so the indices h, k, and l are integers. Distances in reciprocal space, expressed in reciprocal angstroms (Å<sup>-1</sup>) or reciprocal nanometers (nm<sup>-1</sup>), are used to judge the potential resolution of the model that the diffraction data can yield.

Like Alice's looking-glass world, reciprocal space may seem strange to you at first (Fig. 2.13). We will see, however, that some aspects of crystallography are actually easier to understand, and some calculations are more convenient, in reciprocal space than in real space (Chapter 4).

## 2.6 The mathematics of crystallography: A brief description

The problem of determining the structure of objects in a crystalline array from their diffraction pattern is, in essence, a matter of converting the experimentally accessible information in the reciprocal space of the diffraction pattern to otherwise inaccessible information about the real space inside the unit cell. Remember that a computer program that makes this conversion is acting as a simulated lens to reconstruct an image from diffracted radiation. Each reflection is produced by a beam of electromagnetic radiation (X-rays), so the computations entail treating the reflections as waves and recombining these waves to produce an image of the molecules in the unit cell.



**Figure 2.13** ► Fun in reciprocal space. © The New Yorker Collection, 1991. John O'Brien, from cartoonbank.com. All rights reserved.

#### 2.6.1 Wave equations: Periodic functions

Each reflection is the result of diffraction from complicated objects, the molecules in the unit cell, so the resulting wave is complicated also. Before considering how the computer represents such an intricate wave, I will consider mathematical descriptions of the simplest waves (Fig. 2.14).

A simple wave, like that of visible light or X-rays, can be described by a periodic function, for instance, an equation of the form

$$f(x) = F \cos 2\pi \left(hx + \alpha\right) \tag{2.1}$$

or

$$f(x) = F \sin 2\pi (hx + \alpha). \tag{2.2}$$

In these functions, f(x) specifies the vertical height of the wave at any horizontal position x along the wave. The variable x and the constant  $\alpha$  are angles expressed in fractions of the wavelength; that is, x = 1 implies a position of one full wavelength ( $2\pi$  radians or  $360^{\circ}$ ) from the origin. The constant F specifies the *amplitude* of the wave (the height of crests from the horizontal wave axis). For example, the crests of the wave  $f(x) = 3 \cos 2\pi x$  are three times as high and the troughs are three times as deep as those of the wave  $f(x) = \cos 2\pi x$  (compare b with a in Fig. 2.14).



**Figure 2.14** Graphs of four simple wave equations  $f(x) = F \cos 2\pi (hx + \alpha)$ . (a)  $F = 1, h = 1, \alpha = 0$ :  $f(x) = \cos 2\pi (x)$ . (b)  $F = 3, h = 1, \alpha = 0$ :  $f(x) = 3 \cos 2\pi (x)$ . Increasing *F* increases the amplitude of the wave. (c)  $F = 1, h = 3, \alpha = 0$ :  $f(x) = \cos 2\pi (3x)$ . Increasing *h* increases the frequency (or decreases the wavelength  $\lambda$ ) of the wave. (d)  $F = 1, h = 1, \alpha = 1/4$ :  $f(x) = \cos 2\pi (x + 1/4)$ . Changing  $\alpha$  changes the phase (position) of the wave.

The constant h in a simple wave equation specifies the frequency or wavelength of the wave. For example, the wave  $f(x) = \cos 2\pi (3x)$  has three times the frequency (or one-third the wavelength) of the wave  $f(x) = \cos 2\pi x$  (compare c with a in Fig. 2.14). (In the wave equations used in this book, h takes on integral values only.)

Finally, the constant  $\alpha$  specifies the phase of the wave, that is, the position of the wave with respect to the origin of the coordinate system on which the wave is



**Figure 2.15** Visualizing a one-dimensional function, the average daily high temperature, requires two dimensions. The height of the curve f(x) (black) represents the average temperature on day x. To illustrate the notion of phase difference, note that the phase  $\alpha$ of this wave is shifted with respect to a plot of day length, with maximum value around June 20 and minimum around December 20 (green curve).

plotted. For example, the position of the wave  $f(x) = \cos 2\pi (x + 1/4)$  is shifted by one-quarter of  $2\pi$  radians (or one-quarter of a wavelength, or 90°) from the position of the wave  $f(x) = \cos 2\pi x$  (compare *d* with *a* in Fig. 2.14). Because the wave is repetitive, with a repeat distance of one wavelength or  $2\pi$  radians, a phase of 1/4 is the same as a phase of  $1\frac{1}{4}$ , or  $2\frac{1}{4}$ , or  $3\frac{1}{4}$ , and so on. In radians, a phase of 0 is the same as a phase of  $2\pi$ , or  $4\pi$ , or  $6\pi$ , and so on. (This use of the term *phase* is different from common parlance, in which, for example, the new moon is called a phase of the lunar month, or autumn is thought of as a phase or time of the year. Mathematically, *phase* gives the *position of the entire wave* with respect to a specified origin, not merely a *location* on that wave; location is given by *x*.)

These equations describe one-dimensional waves, in which a property (in this case, the height of the wave) varies in one direction. Visualizing a *one*-dimensional function f(x) requires a *two*-dimensional graph, with the second dimension used to represent the numerical value of f(x) (Fig. 2.15). For example, if f(x) describes the average daily high temperature over a year's time, the x-axis represents time in days, and the height of the curve f(x) on day x represents the average temperature on that day. The temperature f(x) is in no real sense perpendicular to the time x, but it is convenient to use the perpendicular direction to show the numerical value of the temperature. In general, visualizing a function in n dimensions requires n + 1 dimensions.

#### 2.6.2 Complicated periodic functions: Fourier series and sums

As discussed in Sec. 2.6.1, p. 21, any simple sine or cosine wave can be described by three constants—the amplitude F, the frequency h, and the phase  $\alpha$ . It is less obvious that far more complicated waves can also be described with this same simplicity. The French mathematician Jean Baptiste Joseph Fourier (1768–1830) showed that even the most intricate periodic functions can be described as the sum of simple sine and cosine functions whose wavelengths are integral fractions of the wavelength of the complicated function. Such a sum is called a *Fourier series* and each simple sine or cosine function in the sum is called a *Fourier term*.

Figure 2.16 shows a periodic function, called a step function, and the beginning of a Fourier series that describes it. A method called *Fourier synthesis* is used to compute the sine and cosine terms that describe a complex wave, which I will call the "target" of the synthesis. I will discuss the results of Fourier synthesis, but not the method itself. In the example of Fig. 2.16, the first four terms produced by Fourier synthesis are shown individually ( $f_0$  through  $f_3$ , on the left), and each is added sequentially to the Fourier sum (on the right). Notice that the first term in the series,  $f_0 = 1$ , simply displaces the sums upward so that they have only positive values like the target function. (Sine and cosine functions themselves have both positive and negative values, with average values of zero.) The second term  $f_1 = \cos 2\pi x$ , has the same wavelength as the step function, and wavelengths of subsequent terms are simple fractions of that wavelength. (It is equivalent to say, and it is plain in the equations, that the frequencies h are simple multiples of the frequency of the step function.) Notice that the sum of only the first few Fourier terms merely approximates the target. If additional terms of shorter wavelength are computed and added, the fit of the approximated wave to the target improves, as shown by the sum of the first six terms. Indeed, using the tenets of Fourier theory, it can be proved that such approximations can be made as similar as desired to the target waveform, simply by including enough terms in the series.

Look again at the components of this Fourier series, functions  $f_0$  through  $f_3$ . The low-frequency terms like  $f_1$  approximate the gross features of the target wave. Higher-frequency terms like  $f_3$  improve the approximation by filling in finer details, for example, making the approximation better in the sharp corners of the target function. We would need to extend this series infinitely to reproduce the target perfectly.

#### 2.6.3 Structure factors: Wave descriptions of X-ray reflections

Each diffracted X-ray that arrives at the film to produce a recorded reflection can also be described as the sum of the contributions of all scatterers in the unit cell. The sum that describes a diffracted ray is called a *structure-factor equation*. The computed sum for the reflection *hkl* is called the *structure factor*  $F_{hkl}$ . As I will show in Chapter 4, the structure-factor equation can be written in several different ways. For example, one useful form is a sum in which each term describes diffraction by one atom in the unit cell, and thus the sum contains the same number of terms as the number of atoms.

If diffraction by atom A in Fig. 2.17 is represented by  $f_A$ , then one diffracted ray (producing one reflection) from the unit cell of Fig. 2.17 is described by a structure-factor equation of this form:

$$F_{hkl} = f_A + f_B + \dots + f_{A'} + f_{B'} + \dots + f_{F'}.$$
(2.3)



**Figure 2.16** Beginning of a Fourier series to approximate a target function, in this case, a step function or square wave.  $f_0 = 1$ ;  $f_1 = \cos 2\pi(x)$ ;  $f_2 = (-1/3) \cos 2\pi(3x)$ ;  $f_3 = (1/5) \cos 2\pi(5x)$ . In the left column are the target and terms  $f_1$  through  $f_3$ . In the right column are  $f_0$  and the succeeding sums as each term is added to  $f_0$ . Notice that the approximation improves (that is, each successive sum looks more like the target) as the number of Fourier terms in the sum increases. In the last graph, terms  $f_4$ ,  $f_5$ , and  $f_6$  are added (but not shown separately) to show further improvement in the approximation.



**Figure 2.17**  $\blacktriangleright$  Every atom contributes to every reflection in the diffraction pattern, as described for this unit cell by Eq. 2.3.

The structure-factor equation implies, and correctly so, that each reflection on the film is the result of diffractive contributions from all atoms in the unit cell. That is, every atom in the unit cell contributes to every reflection in the diffraction pattern. The structure factor  $F_{hkl}$  is a wave created by the superposition of many individual waves  $f_j$ , each resulting from diffraction by an individual atom. So the structure factor is the sum of many wave equations, one for diffraction by each atom. In that sense, the structure factor equation is a Fourier *sum* (sometimes called a Fourier *summation*, but I prefer one syllable to three), but not a Fourier *series*. In a Fourier *series*, each succeeding term can be generated from the previous one by some repetitive formula, as in Fig. 2.16.

#### 2.6.4 Electron-density maps

To be more precise about diffraction, when we direct an X-ray beam toward a crystal, the *actual diffractors of the X rays are the clouds of electrons* in the molecules of the crystal. Diffraction should therefore reveal the distribution of electrons, or the *electron density*, of the molecules. Electron density, of course, reflects the molecule's shape; in fact, you can think of the molecule's boundary as a van der Waals surface, the surface of a cloud of electrons that surrounds the molecule. Because, as noted earlier, protein molecules are ordered, and because, in a crystal, the molecules are in an ordered array, the electron density in a crystal can be described mathematically by a periodic function.

If we could walk through the crystal depicted in Fig. 2.4, p. 12, along a linear path parallel to a cell edge, and carry with us a device for measuring electron density, our device would show us that the electron density varies along our path in a complicated periodic manner, rising as we pass through molecules, falling in the space between molecules, and repeating its variation identically as we pass through each unit cell. Because this statement is true for linear paths parallel to all three cell edges, the electron density, which describes the surface features and overall shapes of all molecules in the unit cell, is a three-dimensional periodic



**Figure 2.18** Small volume element *m* within the unit cell, one of many elements formed by subdividing the unit cell with planes parallel to the cell edges. The average electron density within *m* is  $\rho_m(x, y, z)$ . Every volume element contributes to every reflection in the diffraction pattern, as described by Eq. 2.4.

function. I will refer to this function as  $\rho(x, y, z)$ , implying that it specifies a value  $\rho$  for electron density at every position x, y, z in the unit cell. A graph of the function is an image of the electron clouds that surround the molecules in the unit cell. The most readily interpretable graph is a contour map—a drawing of a surface along which there is constant electron density (refer to Fig. 2.3, p. 11). The graph is called an *electron-density map*. The map is, in essence, a fuzzy image of the molecules in the unit cell. The goal of crystallography is to obtain the mathematical function whose graph is the desired electron-density map.

#### 2.6.5 Electron density from structure factors

Because the electron density we seek is a complicated periodic function, it can be described as a Fourier sum. Do the many structure-factor equations, each a sum of wave equations describing one reflection in the diffraction pattern, have any connection with the Fourier function that describes the electron density? As mentioned earlier, each structure-factor equation can be written as a sum in which each term describes diffraction from one atom in the unit cell. But this is only one of many ways to write a structure-factor equation. Another way is to imagine dividing the electron density in the unit cell into many small volume elements by inserting planes parallel to the cell edges (Fig. 2.18).

These volume elements can be as small and numerous as desired. Now because the true diffractors are the clouds of electrons, each structure-factor equation can be written as a Fourier sum in which each term describes diffraction by the electrons in one volume element. In this sum, each term contains the average numerical value of the desired electron density function  $\rho(x, y, z)$  within one volume element. If the cell is divided into *n* elements, and the average electron density in volume element *m* is  $\rho_m$ , then one diffracted ray from the unit cell of Fig. 2.18 is described by a structure-factor equation, another Fourier sum, of this form:

$$F_{hkl} = f(\rho_1) + f(\rho_2) + \dots + f(\rho_m) + \dots + f(\rho_n).$$
(2.4)

Each reflection is described by an equation like this one, giving us a large number of equations describing reflections in terms of the electron density. Is there any way to solve these equations for the function  $\rho(x, y, z)$  in terms of the measured reflections? After all, structure factors like Eq. 2.4 describe the reflections in terms of  $\rho(x, y, z)$ , which is precisely the function the crystallographer is trying to learn. I will show in Chapter 5 that a mathematical operation called the *Fourier transform* solves the structure-factor equations for the desired function  $\rho(x, y, z)$ , just as if they were a set of simultaneous equations describing  $\rho(x, y, z)$  in terms of the amplitudes, frequencies, and phases of the reflections.

The Fourier transform describes precisely the mathematical relationship between an object and its diffraction pattern. In Figs. 2.8–2.11 (pp. 15–18), the diffraction patterns are the Fourier transforms of the corresponding objects or arrays of objects. To put it another way, the Fourier transform is the lenssimulating operation that a computer performs to produce an image of molecules (or more precisely, of electron clouds) in the crystal. This view of  $\rho(x, y, z)$  as the Fourier transform of the structure factors implies that if we can measure three parameters—amplitude, frequency, and phase—of *each* reflection, then we can add them together to obtain the function  $\rho(x, y, z)$ , graph the function, and "see" a fuzzy image of the molecules in the unit cell.

#### 2.6.6 Electron density from measured reflections

Are all three of these parameters accessible in the data that reaches our detectors? I will show in Chapter 5 that the measurable intensity  $I_{hkl}$  of one reflection gives the amplitude of one Fourier term in the series that describes  $\rho(x, y, z)$ , and that the position *hkl* specifies the frequency for that term. But the phase  $\alpha$  of each reflection is not recorded on any kind of detector. In Chapter 6, I will show how to obtain the phase of each reflection, completing the information we need to calculate  $\rho(x, y, z)$ .

A final note: Even though we cannot measure phases by simply collecting diffraction patterns, we can compute them from a known structure, and we can depict them by adding color to images like those of Figs. 2.8–2.11. In his innovative World Wide Web *Book of Fourier*, Kevin Cowtan illustrates phases in diffraction patterns in this clever manner. For example, Fig. 2.19 shows a simple group of atoms, like the carbon atoms in ethylbenzene. Figure 2.19*b* is the computed Fourier transform of (*a*). Image (*c*) depicts a lattice of the objects in (*a*), and (*d*) is the corresponding Fourier transform.

Because patterns (b) and (d) were *computed* from objects of known structure, rather than measured experimentally from real objects, the phases are included in the calculated results, and thus are known. The phase of each reflection is depicted by its color, according to the color wheel (f). The phase can be expressed as an



**Figure 2.19** Simple asymmetric object, alone (*a*) and in a lattice (*c*), and the computed Fourier transforms of each (*b* and *d*). Phases in *b* and *d* are depicted by color. Darkness of color indicates the intensity of a reflection. The phase angle of a region in *b* or a reflection in *d* corresponds to the angle of its color on the color wheel (*f*). Experimental diffraction patterns do not contain phase information, as in (*e*). Images computed and generously provided by Dr. Kevin Cowtan. For additional vivid illustrations of Fourier transforms as they apply to crystallography, direct your web browser to the CMCC home page and select Kevin Cowtan's *Book of Fourier*.

angle between 0° and 360° [this is the angle  $\alpha$  in Eqs. (2.1) or (2.2)]. In Fig. 2.19, the phase angle of each region (in *b*) or reflection (in *d*) is the angle that corresponds to the angle of its color on the color wheel (*f*). For example, red corresponds to a phase angle of 0°, and green to an angle of about 135°. So a dark red reflection has a high intensity (dark color) and a phase angle of 0° (red). A pale green reflection has a low intensity (faint color) and a phase angle of about 135° (green). With the addition of color, these Fourier transforms give a full description of each reflection, including the phase angle that we do not learn from diffraction experiments, which would give us only the intensities, as shown in (*e*). In a sense then, Figs. 2.8–2.11 show *diffraction* patterns, whereas Fig. 2.19*b* and *d* show *structure-factor patterns*, which depict the structure factors fully. Note again that (*d*) is a sampling of (*b*) at points corresponding to the reciprocal lattice of the lattice in (*c*). In other words, the diffraction pattern (*d*) still contains the diffraction signature, including both intensities and phases, of the object in (*a*).

In these terms, I will restate the central problem of crystallography: In order to determine a structure, we need a full-color version of the diffraction pattern that is, a full description of the structure factors, including amplitude, frequency, and phase. But diffraction experiments give us only the black-and-white version (*e* in Fig. 2.19), the positions and intensities of the reflections, but no information about their phases. We must learn the phase angles from further experimentation, as described fully in Chapter 6.

#### 2.6.7 Obtaining a model

Once we obtain  $\rho(x, y, z)$ , we graph the function to produce an electron-density map, an image of the molecules in the unit cell. Finally, we interpret the map by building a model that fits it (Fig. 2.3*b*, p. 11). In interpreting the molecular image and building the model, a crystallographer takes advantage of all current knowledge about the protein under investigation, as well as knowledge about protein structure in general. The most important element of this *prior knowledge* is the sequence of amino acids in the protein. In a few rare instances, the amino-acid sequence has been learned from the crystallographic structure. But in almost all cases, crystallographers know the sequence to start with, from the work of chemists or molecular biologists, and use it to help them interpret the image obtained from crystallography. In effect, the crystallographer starts with knowledge of the chemical structure, but without knowledge of the conformation. Interpreting the image amounts to finding a chemically realistic conformation that fits the image precisely.

A crystallographer interprets a map by displaying it on a graphics computer and building a graphics model within it. The final model must be (1) consistent with the image and (2) chemically realistic; that is, it must possess bond lengths, bond angles, conformational angles, and distances between neighboring groups that are all in keeping with established principles of molecular structure and stereochemistry. With such a model in hand, the crystallographer can begin to explore the model for clues about its function.

In Chapters 3–7, I will take up in more detail the principles introduced in this chapter.

## ► Chapter 3

# **Protein Crystals**

## 3.1 Properties of protein crystals

#### 3.1.1 Introduction

As the term *X-ray crystallography* implies, the sample being examined is in the crystalline state. Crystals of many proteins and other biomolecules have been obtained and analyzed in the X-ray beam. A few macromolecular crystals are shown in Fig. 3.1.

In these photographs, the crystals appear much like inorganic materials such as sodium chloride. But there are several important differences between protein crystals and ionic solids.

#### 3.1.2 Size, structural integrity, and mosaicity

Whereas inorganic crystals can often be grown to dimensions of several centimeters or larger, it is frequently impossible to grow protein crystals as large as 1 mm in their shortest dimension. In addition, larger crystals are often twinned (two or more crystals grown into each other at different orientations) or otherwise imperfect and not usable. Roughly speaking, protein crystallography requires a crystal of at least 0.2 mm in its shortest dimension, although modern methods of data collection can sometimes succeed with smaller crystals, and modern software can sometimes decipher data from twinned crystals.

Inorganic crystals derive their structural integrity from the electrostatic attraction of fully charged ions. On the other hand, protein crystals are held together by weaker forces, primarily hydrogen bonds between hydrated protein surfaces. In other words, proteins in the crystal stick to each other primarily by hydrogen bonds through intervening water molecules. Protein crystals are thus much more fragile than inorganic crystals; gentle pressure with a needle is enough to crush the hardiest protein crystal. Growing, handling, and mounting crystals for analysis



**Figure 3.1** Some protein crystals grown by a variety of techniques and using a number of different precipitating agents. They are (*a*) deer catalase, (*b*) trigonal form of fructose-1,6-diphosphatase from chicken liver, (*c*) cortisol binding protein from guinea pig sera, (*d*) concanavalin B from jack beans, (*e*) beef liver catalase, (*f*) an unknown protein from pineapples, (*g*) orthorhombic form of the elongation factor Tu from *Escherichia coli*, (*h*) hexagonal and cubic crystals of yeast phenylalanine tRNA, (*i*), monoclinic laths of the gene 5 DNA unwinding protein from bacteriophage fd, (*j*) chicken muscle glycerol-3-phosphate dehydrogenase, and (*k*) orthorhombic crystals of canavalin from jack beans. From A. McPherson, in *Methods in Enzymology* **114**, H. W. Wyckoff, C. H. W. Hirs, and S. N. Timasheff, eds., Academic Press, Orlando, Florida, 1985, p. 114. Photo generously provided by the author; photo and caption reprinted with permission.



**Figure 3.2**  $\triangleright$  Crystals are not perfectly ordered. They consist of many small arrays in rough alignment with each other. As a result, reflections are not points, but are spherical or ovoid, and must be measured over a small angular range.

thus require very gentle techniques. If possible, protein crystals are often harvested, examined, and mounted for crystallography within their *mother liquor*, the solution in which they formed.

The textbook image of a crystal is that of a perfect array of unit cells stretching throughout. Real macroscopic crystals are actually mosaics of many submicroscopic arrays in rough alignment with each other, as illustrated in Fig. 3.2. The result of mosaicity is that an X-ray reflection actually emerges from the crystal as a narrow cone rather than a perfectly linear beam. Thus the reflection must be measured over a very small range of angles, rather than at a single, well-defined angle. In protein crystals, composed as they are of relatively flexible molecules held together by weak forces, this mosaicity is more pronounced than in crystals of rigid organic or inorganic molecules, and the reflections from protein crystals therefore suffer greater *mosaic spread* than do those from more ordered crystals.

#### 3.1.3 Multiple crystalline forms

In efforts to obtain crystals, or to find optimal conditions for crystal growth, crystallographers sometimes obtain a protein or other macromolecule in more than one crystalline form. Compare, for instance, Figs. 3.1*a* and *e*, which show crystals of the enzyme catalase from two different species. Although these enzymes are almost identical in molecular structure, they crystallize in different forms. In Fig. 3.1*h*, you can see that highly purified yeast phenylalanyl tRNA (transfer ribonucleic acid) crystallizes in two different forms. Often, the various crystal forms differ in quality of diffraction, in ease and reproducibility of growth, and perhaps in other properties. The crystallographer must ultimately choose the best form with which to work. Quality of diffraction is the most important criterion, because it determines the ultimate quality of the crystallographic model. Among forms that diffract equally well, more symmetrical forms are usually preferred because they require less data collection (see Chapter 4).

#### 3.1.4 Water content

Early protein crystallographers, proceeding by analogy with studies of other crystalline substances, examined dried protein crystals and obtained no diffraction patterns. Thus X-ray diffraction did not appear to be a promising tool for analyzing proteins. In 1934, J. D. Bernal and Dorothy Crowfoot (later Hodgkin) measured diffraction from pepsin crystals still in the mother liquor. Bernal and Crowfoot recorded sharp diffraction patterns, with reflections out to distances in reciprocal space that correspond in real space to the distances between atoms. The announcement of their success was the birth announcement of protein crystallography.

Careful analysis of electron-density maps usually reveals many ordered water molecules on the surface of crystalline proteins (Fig. 3.3). Additional disordered water is presumed to occupy regions of low density between the ordered particles. *Ordered water* refers to water molecules that occupy the same site on every protein molecule in every unit cell (or a high percentage of them) and thus show up clearly in electron-density maps. *Disordered water* refers to bulk water molecules that occupy the spaces between protein molecules, are in different arrangements in each unit cell, and thus show up only as uniform regions of low electron density. The quantity of water varies among proteins and even among different crystal forms of the same protein. The number of detectable ordered water molecules averages about one per amino-acid residue in the protein. Both the ordered and disordered water are essential to crystal integrity, so drying destroys the crystal structure. For this reason, protein crystals are subjected to X-ray analysis in a very humid atmosphere or in a solution that will not dissolve them, such as the mother liquor or a protective harvest buffer.



Figure 3.3 ► Model (stereo) of one molecule of crystalline adipocyte lipid-binding protein (ALBP, PDB 1alb), showing ordered water molecules on the surface and within a molecular cavity where lipids are usually bound. Protein is shown as a ball-and-stick model with carbon dark gray, oxygen red, and nitrogen blue. Ordered water molecules, displayed as space-filling oxygen atoms, are green. Image: DeepView/POV-Ray.

34

NMR analysis of protein structure suggests that the ordered water molecules seen by X-ray diffraction on protein surfaces have very short residence times in solution. Thus most of these molecules may be of little importance to an understanding of protein function. However, ordered water is of great importance to the crystallographer. As the structure determination progresses, ordered water becomes visible in the electron-density map. For example, in Fig. 2.3, p. 11, water molecules are implied by small regions of disconnected density. Positions of these molecules are indicated by red crosses. Assignment of water molecules to these isolated areas of electron density improves the overall accuracy of the model, and for reasons I will discuss in Chapter 7, improvements in accuracy in one area of the model give accompanying improvements in all other regions.

3.2

# Evidence that solution and crystal structures are similar

Knowing that crystallographers study proteins in the crystalline state, you may be wondering if these molecules are altered when they crystallize, and whether the structure revealed by X-rays is pertinent to the molecule's action in solution. Crystallographers worry about this problem also, and with a few proteins, it has been found that crystal structures are in conflict with chemical or spectroscopic evidence about the protein in solution. These cases are rare, however, and the large majority of crystal structures appear to be identical to the solution structure. Because of the slight possibility that crystallization will alter molecular structure, an essential part of any structure determination project is an effort to show that the crystallized protein is not significantly altered.

#### 3.2.1 Proteins retain their function in the crystal

Probably the most convincing evidence that crystalline structures can safely be used to draw conclusions about molecular function is the observation that many macromolecules are still functional in the crystalline state. For example, substrates added to suspensions of crystalline enzymes are converted to product, albeit at reduced rates, suggesting that the enzyme's catalytic and binding sites are intact. The lower rates of catalysis can be accounted for by the reduced accessibility of active sites within the crystal, in comparison to solution.

In a dramatic demonstration of the persistence of protein function in the crystalline state, crystals of deoxyhemoglobin shatter in the presence of oxygen. Hemoglobin molecules are known to undergo a substantial conformational change when they bind oxygen. The conformation of oxyhemoglobin is apparently incompatible with the constraints on deoxyhemoglobin in crystalline form, and so oxygenation disrupts the crystal.

It makes sense, therefore, after obtaining crystals of a protein and before embarking on the strenuous process of obtaining a structure, to determine whether the protein retains its function in the crystalline state. If the crystalline form is functional, the crystallographer can be confident that the model will show the molecule in its functional form.

## 3.2.2 X-ray structures are compatible with other structural evidence

Further evidence for the similarity of solution and crystal structures is the compatibility of crystallographic models with the results of chemical studies on proteins. For instance, two reactive groups in a protein might be linked by a cross-linking reagent, demonstrating their nearness. The groups shown to be near each other by such studies are practically always found near each other in the crystallographic model.

In a growing number of cases, both NMR and X-ray methods have been used to determine the structure of the same molecule. Figure 3.4 shows the alpha-carbon backbones of two models of the protein thioredoxin. The blue model was obtained by X-ray crystallography and the red model by NMR. Clearly the two methods produce similar models. The models are most alike in the pleated-sheet core and the alpha helices. The greatest discrepancies, even though they are not large, lie in the surface loops at the top and bottom of the models. This and other NMR-derived models confirm that protein molecules are very similar in crystals and in solution. In some cases, small differences are seen and can usually be attributed to crystal packing. Often these packing effects are detectable in the crystallographic



**Figure 3.4**  $\blacktriangleright$  Models (stereo) of the protein thioredoxin (human, reduced form) as obtained from X-ray crystallography (blue, PDB 1ert) and NMR (red, PDB 3trx). Only backbone alpha carbons are shown. The models were superimposed by least-squares minimization of the distances between corresponding alpha carbons, using DeepView. Image: DeepView/POV-Ray.

model itself. For instance, in the crystallographic model of cytoplasmic malate dehydrogenase (PDB file 4mdh), whose functional form is a symmetrical dimer, an external loop has different conformations in the two molecules of one dimer. On examination of the dimer in the context of neighboring dimers, it can be seen that one molecule of each pair lies very close to a molecule of a neighboring pair. It was thus inferred that the observed difference between the oligomers in a dimer is due to crystal packing, and further, that the unaffected molecule of each pair is probably more like the enzyme in solution.

#### 3.2.3 Other evidence

In a few cases, the structure of a protein has been obtained from more than one type of crystal. The resulting models were identical, suggesting that the molecular structure was not altered by crystallization.

Recall that stable protein crystals contain a large amount of both ordered and disordered water molecules. As a result, the proteins in the crystal are still in the aqueous state, subject to the same solvent effects that stabilize the structure in solution. Viewed in this light, it is less surprising that proteins retain their solution structure in the crystal.

## 3.3 Growing protein crystals

#### 3.3.1 Introduction

Crystals suffer damage in the X-ray beam, primarily due to free radicals generated by X-rays. For this reason and others discussed later, a full structure determination project usually consumes many crystals. I will now consider the problem of developing a reliable, reproducible source of protein crystals. This entails not only growing good crystals of the pure protein, but also obtaining *derivatives*, or crystals of the protein in complex with various nonprotein components (loosely called *ligands*). For example, in addition to pursuing the structures of proteins themselves, crystallographers also seek structures of proteins in complexes with ligands such as cofactors, substrate analogs, inhibitors, and allosteric effectors. Structure determination then reveals the details of protein-ligand interactions, giving insight into protein function.

Another vital type of ligand is a heavy-metal atom or ion. Crystals of protein/ heavy-metal complexes, often called *heavy-atom derivatives*, are usually needed in order to solve the phase problem mentioned in Sec. 2.6.6, p. 28. I will show in Chapter 6 that, for the purpose of obtaining phases, it is crucial that crystals of heavy-atom derivatives be *isomorphic* with crystals of the pure protein. This means that derivatives must possess the same unit-cell dimensions and symmetry, and the same protein conformation, as the pure protein, which in discussions of derivatives are called *native crystals*. So in most structure projects, the crystallographer must produce both native and derivative crystals under the same or very similar circumstances. Modern methods of obtaining phases can often succeed with proteins in which residues of the amino acid methionine are replaced by selenomethionine, in which selenium replaces the usual sulfur of methionine. This substitution provides selenium as built-in heavy atoms that usually do not alter protein conformation or unit-cell structure. I will discuss the production of crystals of heavy-atom and so-called *selenomet* derivatives after describing general procedures for crystallization.

#### 3.3.2 Growing crystals: Basic procedure

Crystals of an inorganic substance can often be grown by preparing a hot, saturated solution of the substance and then slowly cooling it. Polar organic compounds can sometimes be crystallized by similar procedures or by slow precipitation from aqueous solutions by addition of organic solvents. If you work with proteins, just the mention of these conditions probably makes you cringe. Proteins, of course, are usually denatured by heat or exposure to organic solvents, so techniques used for small molecules are not appropriate. In the most common methods of growing protein crystals, purified protein is dissolved in an aqueous buffer containing a precipitant, such as ammonium sulfate or polyethylene glycol, at a concentration, [precipitant], just below that necessary to precipitate the protein. Then water is removed by controlled evaporation to raise both [protein] and [precipitant], resulting in precipitation. Slow precipitation is more likely to produce larger crystals, whereas rapid precipitation may produce many small crystals, or worse, an amorphous solid.

In theory, precipitation should occur when the combination of [protein] and [precipitant] exceeds threshold values, as shown in the phase diagram of Fig. 3.5*a*. Crystal formation occurs in two stages, *nucleation*, and *growth*. Nucleation, the initial formation of molecular clusters from which crystals grow, requires protein and/or precipitant concentrations higher than those optimal for slow precipitation (Fig 3.5*a*, blue region). In addition, nucleation conditions, if they persist, result in the formation of many nuclei, and as a result, either an amorphous precipitate or many small crystals instead of a few larger ones. An ideal strategy (Fig. 3.5*b*) would be to start with conditions corresponding to the blue region of the phase diagram, and then, when nuclei form, move into the green region, where growth, but not additional nucleation, can occur.

One widely used crystallization technique is *vapor diffusion*, in which the protein/precipitant solution is allowed to equilibrate in a closed container with a larger aqueous reservoir whose precipitant concentration is optimal for producing crystals. One of many examples of this technique is the hanging-drop method (Fig. 3.6).

Less than 25  $\mu$ L of the solution of purified protein is mixed with an equal amount of the reservoir solution, giving precipitant concentration about 50% of that required for protein crystallization (conditions represented by the red circle in Fig. 3.5b). This solution is suspended as a droplet underneath a cover slip, which is sealed onto the top of the reservoir with grease. Because the precipitant



**Figure 3.5** (a) Phase diagram for crystallization mediated by a precipitant. The red region represents concentrations of protein and precipitant at which the solution is not saturated with protein, so neither nucleation nor growth occurs. The green and blue regions represent unstable solutions that are supersaturated with protein. Conditions in the blue region support both nucleation and growth, while conditions in the green support growth only. (b) An ideal strategy for growing large crystals is to allow nucleation to occur under conditions in the blue region, then to move to conditions in the green region until crystal growth ceases.



**Figure 3.6**  $\blacktriangleright$  Growing crystals by the hanging-drop method. The droplet hanging under the cover slip contains buffer, precipitant, protein, and, if all goes well, protein crystals.

is the major solute present, vapor diffusion (evaporation and condensation) in this closed system results in net transfer of water from the protein solution in the drop to the reservoir, until the precipitant concentration is the same in both solutions. Because the reservoir is much larger than the protein solution, the final concentration of the precipitant in the drop is nearly equal to that in the reservoir. When the system comes to equilibrium, net transfer of water ceases, and the protein solution is maintained at constant precipitant concentration. At this point, drop shrinkage has increased both [precipitant] and [protein], moving conditions diagonally into the nucleation region (blue circle in Fig. 3.5*b*). In this way, the precipitant concentration in the protein solution rises to the level required for nucleation and remains there without overshooting because, at equilibrium, the vapor pressure in the closed system equals the inherent vapor pressure of both protein solution and reservoir. As nuclei form, the protein concentration decreases, moving the conditions vertically into the growth region (green circle in Fig. 3.5*b*).

Frequently the crystallographer obtains many small crystals instead of a few that are large enough for diffraction measurements. If many crystals grow at once, the supply of dissolved protein will be depleted before crystals are large enough to be useful. Small crystals of good quality can be used as seeds to grow larger crystals. The experimental setup is the same as before, except that each hanging droplet is seeded with a few small crystals. Seed crystals are sometimes *etched* before use by brief soaking in buffer with precipitant concentration lower than that of the mother liquor. This soak dissolves outer layers of the seed crystal, exposing fresh surface on which crystallization can proceed. Seeds may also be obtained by crushing small crystals or by stroking a crystal with a hair and passing the hair through the crystallization droplet (it is reported that animal whiskers are best—really). Whatever the seeding method, crystals may grow from seeds up to ten times faster than they grow anew, so most of the dissolved protein goes into only a few crystals.

#### 3.3.3 Growing derivative crystals

Crystallographers obtain the derivatives needed for phase determination and for studying protein-ligand interactions by two methods: cocrystallizing protein and ligand, and soaking preformed protein crystals in mother-liquor solutions containing ligand.

It is sometimes possible to obtain crystals of protein-ligand complexes by crystallizing protein and ligand together, a process called *cocrystallization*. For example, a number of NAD<sup>+</sup>-dependent dehydrogenase enzymes readily crystallize as NAD<sup>+</sup> or NADH complexes from solutions containing these cofactors. Cocrystallization is the only method for producing crystals of proteins in complexes with large ligands, such as nucleic acids or other proteins.

A second means of obtaining crystals of protein-ligand complexes is to soak protein crystals in mother liquor that contains ligand. As mentioned earlier, proteins retain their binding and catalytic functions in the crystalline state, and ligands can diffuse to active sites and binding sites through channels of water in the crystal. Soaking is usually preferred over cocrystallization when the crystallographer plans

#### Section 3.3 Growing protein crystals

to compare the structure of a pure protein with that of a protein-ligand complex. Soaking preformed protein crystals with ligands is more likely to produce crystals of the same form and unit-cell dimensions as those of pure protein, so this method is recommended for first attempts to make isomorphic heavy-atom derivatives.

Making selenomet derivatives requires taking advantage of modern methods of molecular biology, in which the gene encoding a desired protein is introduced (for example, on a plasmid) into a specially designed strain of bacterium or other microbe, which is called an *expression vector*. The microbe, in turn, *expresses* the gene, which means that it produces messenger RNA from the gene and synthesizes the desired protein. To produce a selenomet derivative, the gene for the desired protein is expressed in a mutant microbe that cannot make its own methionine, and thus can live only in a growth medium that provides methionine. If the growth medium provides selenomethionine instead of methionine, the microbe usually grows normally, and expression results in incorporation of selenomethionine wherever methionine would normally appear. Purification and crystallization of the selenomet derivative usually follow the same procedures as for the native protein.

Finally, some proteins naturally contain metal ions that can serve the same purpose in phasing as introduced heavy-atom compounds. For example, hemoglobin contains iron (II) ions that can be used to obtain phase information. For such proteins, there is often no need to produce heavy-atom or selenomet derivatives.

#### 3.3.4 Finding optimal conditions for crystal growth

The two most important keys to success of a crystallographic project are purity and quantity of the macromolecule under study. Impure samples will not make suitable crystals, and even for proteins of the highest purity, repeated trials will be necessary before good crystals result.

Many variables influence the formation of macromolecular crystals. These include obvious ones like protein purity, concentrations of protein and precipitant, pH, and temperature, as well as more subtle ones like cleanliness, vibration and sound, convection, source and age of the protein, and the presence of ligands. Clearly, the problem of developing a reliable source of crystals entails controlling and testing a large number of parameters. The difficulty and importance of obtaining good crystals has prompted the invention of crystallization robots that can be programmed to set up many trials under systematically varied conditions.

The complexity of this problem is illustrated in Fig. 3.7, which shows the effects of varying just two parameters, the concentrations of protein (in this case, the enzyme lysozyme) and precipitant (NaCl). Notice the effect of slight changes in concentration of either protein or precipitant on the rate of crystallization, as well as the size and quality of the resulting crystals.

A sample scheme for finding optimum crystallization conditions is to determine the effect of pH on precipitation with a given precipitant, repeat this determination at various temperatures, and then repeat these experiments with different precipitating agents. Notice in Fig. 3.7 that the region of [protein] versus [precipitant] that gives best crystals is in the shape of an arc, like the arc-shaped growth region



Figure 3.7 ► Schematic map of crystallization kinetics as a function of lysozyme and NaCI concentration obtained from a matrix of dishes. Inserts show photographs of dishes obtained one month after preparation of solutions. From G. Feher and X. Kam, in *Methods in Enzymology* 114, H. W. Wyckoff, C. H. W. Hirs, and S. N. Timasheff, eds., Academic Press, Orlando, Florida, 1985, p. 90. Photo and caption reprinted with permission.

of Fig. 3.5*a*. It turns out that if these same data are plotted as [protein] versus ([protein]  $\times$  [precipitant]), this arc-shaped region becomes a rectangle, which makes it easier to survey the region systematically. For such surveys of crystal-lization conditions, multiple batches of crystals can be grown conveniently by the hanging-drop or other methods in crystallization plates of 24, 48, or 96 wells (Fig. 3.8), each with its own cover. This apparatus has the advantage that the growing crystals can be observed through the cover slips with a dissecting microscope. Then, once the ideal conditions are found, many small batches of crystals can be grown at once, and each batch can be harvested without disturbing the others.

Crystallographers have developed sophisticated schemes for finding and optimizing conditions for crystal growth. One approach, called a response-surface procedure, begins with the establishment of a scoring scheme for results, such as giving higher scores for lower ratios of the shortest to the longest crystal dimension. This method gives low scores for needles and higher scores for cubes.



**Figure 3.8**  $\triangleright$  Well-plate, in which 24 sitting-drop crystallization trials can be carried out. Each well contains a pedestal with a concave top, in which the drop sits. Vapor diffusion occurs between drop and reservoir in the bottom of the well.

Then crystallization trials are carried out, varying several parameters, including pH, temperature, and concentrations of protein, precipitant, and other additives. The results are scored, and the relationships between parameters and scores are analyzed. These relationships are fitted to mathematical functions (like polynomials), which describe a complicated multidimensional surface (one dimension for each variable or for certain revealing combinations of variables) over which the score varies. The crystallographer wants to know the location of the "peaks" on this surface, where scores are highest. Such peaks may lie at sets of crystallization conditions that were not tried in the trials and may suggest new and more effective conditions for obtaining crystals. Finding peaks on such surfaces is just like finding the maximum or minimum in any mathematical function. You take the derivative of the function, set it equal to zero, and solve for the values of the parameters. The sets of values obtained correspond to conditions that lie at the top of mountains on the surface of crystal scores.

An example of this approach is illustrated in Fig. 3.9. The graph in the center is a two-dimensional slice of a four-dimensional surface over which [protein], ([protein] × [precipitant]), pH, and temperature were varied, in attempts to find optimal crystallization conditions for the enzyme tryptophanyl-tRNA synthetase. Note that this surface samples the rectangular region [protein] versus ([protein] × [precipitant]), mentioned earlier. The height of the surface is the score for the crystallization. Surrounding the graph are photos of typical crystals obtained in multiple trials of each set of conditions. None of the trial conditions were near the peak of the surface. The photos labeled Opt1 and Opt2 are of crystals obtained from conditions defined by the surface peak. In this instance, the response-surface approach predicted conditions that produced better crystals than any from the trials that pointed to these conditions.



**Figure 3.9**  $\triangleright$  Optimization of conditions for crystallization of tryptophanyl-tRNA synthetase. Photo insets show crystals obtained from various conditions represented by points on the surface. Coordinates of the surface are protein concentration (PROTEIN), product of protein concentration and precipitant concentration (PRO\_PPNT), and the shape of the crystal as reflected by the ratio of its two smallest dimensions, width and length (WL\_RATIO). From C. W. Carter, in *Methods in Enzymology* **276**, C. W. Carter and R. M. Sweet, eds., Academic Press, New York, 1997, p. 75. Reprinted with permission.

So if you decide to try to grow some of your own crystals, how should you proceed? Theoretical studies like those described above, as well as the recorded experience of myriad crystallization successes and failures, have led to development of commercial screening kits that can often streamline the pursuit of crystals. Typical kits are sets of 24, 48, or 96 solutions containing various buffers, salts, and precipitants, representing a wide variety of potential crystallization conditions. After establishing appropriate protein concentration for screening (there is a kit for that, too), you would set up one trial with each of the screen solutions in cells of crystallization plates like the one shown in Fig. 3.8 (\$ome kit\$ even come with prefilled well plate\$). If a particular screen solution produces promising crystals, you can then try to optimize the conditions by varying pH, [salt] or [precipitant] around the values of the screen solution.

Another way to tap accumulated wisdom about crystallization is through online databases. For example, at the combined Biological Macromolecule Crystallization Database and NASA Archive for Protein Crystal Growth Data (see CMCC home page), you can search for successful crystallization conditions for thousands

#### Section 3.3 Growing protein crystals

of macromolecules. You can search by many criteria, including molecule name, source species, prosthetic groups, molecular weight, space groups, as well as specific precipitants, methods, or conditions. Conditions that have succeeded with proteins similar to your target may be good starting points.

When varying the more conventional parameters fails to produce good crystals, the crystallographer may take more drastic measures. Sometimes limited digestion of the protein by a proteolytic enzyme removes a disordered surface loop, resulting in a more rigid, hydrophilic, or compact molecule that forms better crystals. A related measure is adding a ligand, such as a cofactor, that is known to bind tightly to the protein. The protein-ligand complex may be more likely to crystallize than the free protein, either because the complex is more rigid than the free protein or because the cofactor induces a conformational change that makes the protein more amenable to crystallizing. Desperation has even prompted addition of coffee (usually readily at hand in research labs) to precipitant mixtures, but I am aware of no successes from this measure.

Many membrane-associated proteins will not dissolve in aqueous buffers and tend to form amorphous precipitates instead of crystals. The intractability of such proteins often results from hydrophobic domains or surface regions that are normally associated with the interior of membranes. Such proteins have sometimes been crystallized in the presence of detergents, which coat the hydrophobic portion and decorate it with ionic groups, thus rendering it more soluble in water. A small number of proteins have been diffused into crystalline phases of lipid to produce ordered arrays that diffracted well and yielded structures. In some cases, limited proteolysis of membrane-associated proteins has removed exposed hydrophobic portions, leaving crystallizable fragments that are more like a typical water-soluble protein. Membrane proteins are greatly under-represented in the Protein Data Bank, due to their resistance to crystallization. The search for widely applicable conditions for crystallizing membrane proteins is one of crystallography's holy grails. The announcement of a model of a new membrane protein is usually greeted with much attention, and the first question is usually, "How did they crystallize it?"

The effects of modifications of the target protein, as well as the potential crystallizability of a newly purified protein, can be tentatively assessed before crystallization trials begin, through analysis of laser light scattering by solutions of the macromolecule. Simple, rapid light-scattering experiments (see Sec. 9.3, p. 219) can reveal much about the nature of the substance in solutions of varied composition, pH, and temperature, including estimates of average molecular mass of the particles, radius of gyration (dependent on shape of particles), rates of diffusion through the solution, and range and distribution of particle sizes (degree of *polydispersity*). Some of the measured properties correlate well with crystallizability. In particular, *monodisperse* preparations—those containing particles of uniform size—are more promising candidates for crystallization than those in which the protein is polydisperse. In many cases, polydispersity arises from non-specific interactions among the particles, which at higher concentrations is likely to result in random aggregation rather than orderly crystallization.
When drastic measures like proteolysis are required to yield good crystals, the crystallographer is faced with the question of whether the resulting fragment is worthy of the arduous effort to determine its structure. This question is similar to the basic issue of whether a protein has the same structure in crystal and in solution, and the question must be answered in the same way. Specifically, it may be possible to demonstrate that the fragment maintains at least part of the biological function of the intact molecule, and further, that this function is retained after crystallization.

# 3.4 Judging crystal quality

The acid test of a crystal's suitability for structure determination is, of course, its capacity to give sharp diffraction patterns with clear reflections at large angles from the X-ray beam. Using equipment typical of today's crystallography laboratories, researchers can collect preliminary diffraction data quickly and decide whether to obtain a full data set. However, a brief inspection of crystals under a low-power light microscope can also provide some insight into quality and can help the crystallographer pick out the most promising crystals.

Desirable visible characteristics of crystals include optical clarity, smooth faces, and sharp edges. Broken or twinned crystals sometimes exhibit dark cleavage planes within an otherwise clear interior. Depending on the lattice type (Chapter 4) and the direction of viewing relative to unit-cell axes, some crystals strongly rotate plane-polarized light. This property is easily observed by examining the crystal between two polarizers, one fixed and one rotatable, under a microscope. Upon rotation of the movable polarizer, a good-quality crystal will usually brighten and darken sharply.

Once the crystallographer has a reliable source of suitable crystals, data collection can begin.

# 3.5

# Mounting crystals for data collection

The classical method of mounting crystals is to transfer them into a fine glass capillary along with a droplet of the mother liquor. The capillary is then sealed at both ends and mounted onto a goniometer head (see Fig. 4.25, p. 81, and Sec. 4.3.4, p. 80), a device that allows control of the crystal's orientation in the X-ray beam. The droplet of mother liquor keeps the crystal hydrated.

For many years, crystallographers have been aware of the advantages of collecting X-ray data on crystals at very low temperatures, such as that of liquid nitrogen (boiling point -196°C). In theory, lowering the temperature should increase molecular order in the crystal and improve diffraction. In practice, however, early

#### Section 3.5 Mounting crystals for data collection

attempts to freeze crystals resulted in damage due to formation of ice crystals. Then crystallographers developed techniques for flash freezing crystals in the presence of agents like glycerol, which prevent ice from forming. Crystallography at low temperatures is called *cryocrystallography* and the ice-preventing agents are called *cryoprotectants*. Other cryoprotectants include xylitol or sugars such as glucose. Some precipitants, for example, polyethylene glycol, also act as cryoprotectants, and often it is only necessary to increase their concentration in order to achieve protection from ice formation.

If the crystal was not grown in cryoprotectant, preparation for cryocrystallography typically entails placing it in a cryoprotected mother liquor for 5–15 seconds to wash off the old mother liquor (this liquid is sometimes called a *harvest buffer*). If sudden exposure to cryoprotectant damages the crystal, it might be serially transferred through several solutions of gradually increasing cryoprotectant concentration. After transfer into protectant, the crystal is picked up in a small (<1 mm) circular loop of glass wool or synthetic fiber, where it remains suspended in a thin film of solvent, sort of like the soap film in a plastic loop for blowing soap bubbles. The crystal is then flash frozen by dipping the loop into liquid nitrogen. If flash-freezing is successful, the liquid film in the loop freezes into a glass and remains clear (if it is frosty, crystalline water has formed, usually destroying the crystal in the process). For data collection, the loop is mounted onto the goniometer (see Fig. 4.25*b*, p. 81), where it is held in a stream of cold nitrogen gas coming from a reservoir of liquid nitrogen. A temperature of  $-100^{\circ}$ C can be maintained in this manner.

In addition to better diffraction, other benefits of cryocrystallography include reduction of radiation damage to the crystal and hence the possibility of collecting more data—perhaps an entire data set—from a single crystal; reduction of X-ray scattering from water (resulting in cleaner backgrounds in diffraction patterns) because the amount of water surrounding the crystal is far less than that in a droplet of mother liquor in a capillary; and the possibility of safe storage, transport, and reuse of crystals. Crystallographers can take or ship loop-mounted flash-frozen crystals, in liquid-nitrogen-filled insulated containers, to sites of data collection, minimizing handling of crystals at the collection site. With all these benefits, it is not surprising that cryocrystallography is now common practice.

This Page Intentionally Left Blank

# ► Chapter 4

# **Collecting Diffraction Data**

# 4.1 Introduction

In this chapter, I will discuss the geometric principles of diffraction, revealing, in both the real space of the crystal's interior and in reciprocal space, the conditions that produce reflections. I will show how these conditions allow the crystallographer to determine the dimensions of the unit cell and the symmetry of its contents and how these factors determine the strategy of data collection. Finally, I will look at the devices used to produce and detect X-rays and to measure precisely the intensities and positions of reflections.

# 4.2 Geometric principles of diffraction

W. L. Bragg showed that the angles at which diffracted beams emerge from a crystal can be computed by treating diffraction as if it were reflection from sets of equivalent, parallel planes of atoms in a crystal. (This is why each spot in the diffraction pattern is called a *reflection*.) I will first describe how crystallographers denote the planes that contribute to the diffraction pattern.

## 4.2.1 The generalized unit cell

The dimensions of a unit cell are designated by six numbers: the lengths of three unique edges **a**, **b**, and **c**; and three unique angles  $\alpha$ ,  $\beta$ , and  $\gamma$  (Fig. 4.1, p. 50). Notice the use of bold type in naming the unit cell edges or the axes that correspond to them. I will use bold letters (**a**, **b**, **c**) to signify the edges or axes themselves, and letters in italics (*a*, *b*, *c*) to specify their length. Thus *a* is the length of unit cell edge **a**, and so forth.

#### Chapter 4 Collecting Diffraction Data



**Figure 4.1**  $\blacktriangleright$  General (triclinic) unit cell, with edges **a**, **b**, **c** and angles  $\alpha$ ,  $\beta$ ,  $\gamma$ .

A cell in which  $a \neq b \neq c$  and  $\alpha \neq \beta \neq \gamma$ , as in Fig. 4.1, is called *triclinic*, the simplest *crystal system*. If  $a \neq b \neq c$ ,  $\alpha = \gamma = 90^{\circ}$ , and  $\beta > 90^{\circ}$ , the cell is *monoclinic*. If a = b,  $\alpha = \beta = 90^{\circ}$  and  $\gamma = 120^{\circ}$ , the cell is hexagonal. For cells in which all three cell angles are  $90^{\circ}$ , if a = b = c, the cell is *cubic*; if  $a = b \neq c$ , the cell is *tetragonal*; and if  $a \neq b \neq c$ , the cell is *orthorhombic*. The possible crystal systems are shown in Fig. 4.2. The crystal systems form the basis for thirteen unique *lattice types*, which I will describe later in this chapter.

The most convenient coordinate systems for crystallography adopt coordinate axes based on the directions of unit-cell edges. For cells in which at least one cell angle is not 90°, the coordinate axes are not the familiar orthogonal (mutually perpendicular) x, y, and z. In this book, for clarity, I will emphasize unit cells and coordinate systems with orthogonal axes ( $\alpha = \beta = \gamma = 90^\circ$ ), and I will use orthorhombic systems most often, making it possible to distinguish the three cell edges by their lengths. In such systems, the **a** edges of the cell are parallel to the x-axis of an orthogonal coordinate system, edges **b** are parallel to y, and edges **c** are parallel to z. Bear in mind, however, that the principles discussed here can be generalized to all unit cells.

## 4.2.2 Indices of the atomic planes in a crystal

The most readily apparent sets of planes in a crystalline lattice are those determined by the faces of the unit cells. These and all other regularly spaced planes that can be drawn through lattice points can be thought of as sources of diffraction and can be designated by a set of three numbers called *lattice indices* or *Miller indices*. Three indices *hkl* identify a particular set of equivalent, parallel planes. The index *h* gives the number of planes in the set per unit cell in the *x* direction or, equivalently, the number of parts into which the set of planes cut the **a** edge

#### Section 4.2 Geometric principles of diffraction



**Figure 4.2**  $\triangleright$  Crystal systems beginning with the most symmetric (cubic, upper left), and ending with the least symmetric (triclinic, lower right).

of each cell. The indices k and l specify how many such planes exist per unit cell in the y and z directions. An equivalent way to determine the indices of a set of planes is to start at any lattice point and move out into the unit cell away from the plane cutting that lattice point. If the first plane encountered cuts the **a** edge at some fraction  $1/n_a$  of its length, and the same plane cuts the **b** edge at some fraction  $1/n_b$  of its length, then the h index is  $n_a$  and the k index is  $n_b$  (examples are given later). Indices are written in parentheses when referring to the set of planes; hence, the planes having indices hkl are the (hkl) planes.

In Fig. 4.3, each face of an orthorhombic unit cell is labeled with the indices of the set of planes that includes that face. (The crossed arrows lie on the labeled face, and parallel faces have the same indices.)



**Figure 4.3** ► Indices of faces in an orthorhombic unit cell.

The set of planes including and parallel to the **bc** face, and hence normal to the x-axis, is designated (100) because there is one such plane per lattice point in the x direction. In like manner, the planes parallel to and including the **ac** face are called (010) planes (one plane per lattice point along y). Finally, the **ab** faces of the cell determine the (001) planes. (To recognize these planes easily, notice that if you think of the index as (**abc**), the zeros tell you the location of the plane: the zeros in (010) occupy the **a** and **c** positions, so the plane corresponds to the **ac** face.) In the Bragg model of diffraction as reflection from parallel sets of planes, any of these sets of planes can be the source of one diffracted X-ray beam. (Remember that an entire set of parallel planes, not just one plane, acts as a single diffractor, the number of diffracted beams would be small, and the information obtainable from diffraction would be very limited.

In Fig. 4.4, an additional set of planes, and thus an additional source of diffraction, is indicated. The lattice (solid lines) is shown in section parallel to the **ab** faces or the xy plane. The dashed lines represent the intersection of a set of equivalent, parallel planes that are perpendicular to the xy plane of the paper. Note that the planes cut each **a** edge of each unit cell into two parts and each **b** edge into one part, so these planes have indices 210. Because all (210) planes are parallel to the z axis (which is perpendicular to the plane of the paper), the l index is zero. Or equivalently, because the planes are infinite in extent, and are coincident with **c** edges, and thus do not cut edges parallel to the z axis, there are zero (210) planes per unit cell in the z direction. As another example, for any plane in the set shown in Fig. 4.5, the first plane encountered from any lattice point cuts that unit cell at a/2 and b/3, so the indices are 230.

All planes perpendicular to the xy plane have indices hk0. Planes perpendicular to the xz plane have indices h0l, and so forth. Many additional sets of planes are not perpendicular to x, y, or z. For example, the (234) planes cut the unit cell edges **a** into two parts, **b** into three parts, and **c** into four parts (Fig. 4.6).



**Figure 4.4** ► (210) planes in a two-dimensional section of lattice.



**Figure 4.5** ► (230) planes in a two-dimensional section of lattice.



**Figure 4.6**  $\triangleright$  The intersection of three (234) planes with a unit cell. Note that the (234) planes cut the unit-cell edges **a** into two parts, **b** into three parts, and **c** into four parts.

Finally, Miller indices can be negative as well as positive. Sets of planes in which all indices have opposite signs are identical. For example, the (210) planes are the same as (-2 - 1 0), which is commonly written ( $\overline{210}$ ). The ( $\overline{210}$ ) or ( $\overline{210}$ ) planes are identical (all signs opposite), but tilt in the opposite direction from the (210) planes (Fig. 4.7). To determine the sign of indices the *h* and *k* indices for a set of planes that cut the *xy* plane, look at the direction of a line perpendicular to the planes (green arrows in Fig. 4.7), and imagine the line passing through the origin of an *xy* coordinate system. If this perpendicular line points into the (++) and (--) quadrants of the xy plane, then the indices *h* and *k* are both positive or both negative. If the perpendicular points into the (+-) and (-+) quadrants, then *h* and *k* have opposite signs. This mnemonic works in three dimensions as well: a perpendicular to the planes in Fig. 4.6 points into the (+ + +) and (- - -) octants of an *xyz* coordinate system. You will see why this mnemonic works when I show how to construct the reciprocal lattice.

In Bragg's way of looking at diffraction as reflection from sets of planes in the crystal, each set of parallel planes described here (as well as each additional set of planes interleaved between these sets) is treated as an independent diffractor and produces a single reflection. This model is useful for determining the geometry of data collection. Later, when I discuss structure determination, I will consider another model in which each atom or each small volume element of electron density is treated as an independent diffractor, represented by one term in a Fourier sum that describes each reflection. What does the Fourier sum model add to the Bragg model? Bragg's model tells us where to look for the data. The Fourier sum model tells us what the data has to say about the molecular structure, that is, about where the atoms are located in the unit cell.



**Figure 4.7**  $\blacktriangleright$  The (210) and ( $\overline{210}$ ) planes are identical. They tilt in the opposite direction from ( $\overline{210}$ ) and ( $\overline{210}$ ) planes.

# 4.2.3 Conditions that produce diffraction: Bragg's law

Notice that the different sets of equivalent parallel planes in the preceding figures have different interplanar spacing *d*. Among sets of planes (*hkl*), interplanar spacing decreases as any index increases (more planes per unit cell means more closely spaced planes). W. L. Bragg showed that a set of parallel planes with index *hkl* and interplanar spacing  $d_{hkl}$  produces a diffracted beam when X-rays of wavelength  $\lambda$  impinge upon the planes at an angle  $\theta$  and are reflected at the same angle, *only* if  $\theta$  meets the condition

$$2d_{hkl}\sin\theta = n\lambda,\tag{4.1}$$

where *n* is an integer. The geometric construction in Fig. 4.8 demonstrates the conditions necessary for producing a strong diffracted ray. The red dots represent two parallel planes of lattice points with interplanar spacing  $d_{hkl}$ . Two rays  $R_1$  and  $R_2$  are reflected from them at angle  $\theta$ .

Lines *AC* are drawn from the point of reflection *A* of  $R_1$  perpendicular to the ray  $R_2$ . If ray  $R_2$  is reflected at *B*, then the diagram shows that  $R_2$  travels the same distance as  $R_1$  plus an added distance 2*BC*. Because *AB* in the small triangle *ABC* is perpendicular to the atomic plane, and *AC* is perpendicular to the



**Figure 4.8**  $\triangleright$  Conditions that produce strong diffracted rays. If the additional distance traveled by the more deeply penetrating ray  $R_2$  is an integral multiple of  $\lambda$ , then rays  $R_1$  and  $R_2$  interfere constructively.

incident ray, the angle *CAB* equals  $\theta$ , the angle of incidence (two angles are equal if corresponding sides are perpendicular). Because *ABC* is a right triangle, the sine of angle  $\theta$  is *BC*/*AB* or *BC*/*d*<sub>*hkl*</sub>. Thus *BC* equals *d*<sub>*hkl*</sub> sin  $\theta$ , and the additional distance 2*BC* traveled by ray  $R_2$  is 2*d*<sub>*hkl*</sub> sin  $\theta$ .

If this difference in path length for rays reflected from successive planes  $(2d_{hkl}\sin\theta)$  is equal to an integral number of wavelengths  $(n\lambda)$  of the impinging X rays (that is, if  $2d_{hkl} \sin \theta = n\lambda$ ), then the rays reflected from successive planes emerge from the crystal in phase with each other, interfering constructively to produce a strong diffracted beam. For other angles of incidence  $\theta'$  (where  $2d_{hkl} \sin \theta'$  does not equal an integral multiple of  $\lambda$ ), waves emerging from successive planes are out of phase, so they interfere destructively, and no beam emerges at that angle. Think of it this way: If X-rays impinge at an angle  $\theta'$  that does not satisfy the Bragg conditions, then for every reflecting plane p, there will exist, at some depth in the crystal, another parallel plane p' producing a wave precisely  $180^{\circ}$  out of phase with that from p, and thus precisely cancelling the wave from p. So all such waves will be cancelled by destructive interference, and no diffracted ray will emerge at the angle  $\theta'$ . Diffracted rays reflect from (*hkl*) planes of spacing  $d_{hkl}$  only at angles  $\theta$  for which  $2d_{hkl}\sin\theta = n\lambda$ . Notice that what I am calling the diffraction angle  $\theta$  is the angle of incidence *and* the angle of reflection. So the actual angle by which this reflection diverges from the *incident* X-ray beam is  $2\theta$ . I should also add that the intensity of this diffracted ray will depend on how many atoms, or much electron density, lies on this set of planes in the unit cell. If electron density on this set of planes is high, the ray will be strong (high intensity).

If there is little electron density of this set of planes, this ray, although allowed by Bragg's law, will be weak or undetectable.

Notice that the angle of diffraction  $\theta$  is inversely related to the interplanar spacing  $d_{hkl}$  (sin  $\theta$  is proportional to  $1/d_{hkl}$ ). This implies that large unit cells, with large spacings, give small angles of diffraction and hence produce many reflections that fall within a convenient angle from the incident beam. On the other hand, small unit cells give large angles of diffraction, producing fewer measurable reflections. In a sense, the number of measurable reflections depends on how much information is present in the unit cell. Large cells contain many atoms and thus more information, and they produce more information in a diffraction pattern of the same size. Small unit cells contain fewer atoms, and diffraction from them contains less information.

It is not coincidental that I use the variable names h, k, and l for both the indices of planes in the crystal and the indices of reflections in the diffraction pattern (Sec. 2.5, p. 19). I will show later that in fact the electron density on (or parallel to) the set of planes (hkl) produces the reflection hkl of the diffraction pattern. In the terms I used in Chapter 2, each set of parallel planes in the crystal produces one reflection, or one term in the Fourier sum that describes the electron density within the unit cell. The intensity of that reflection depends on the electron distribution and density along the planes that produce the reflection.

#### 4.2.4 The reciprocal lattice

Now let us consider the Bragg conditions from another point of view: in reciprocal space—the space occupied by the reflections. Before looking at diffraction from this vantage point, I will define and tell how to construct a new lattice, the *reciprocal lattice*, in what will at first seem an arbitrary manner. But I will then show that the points in this reciprocal lattice are the locations of all the Bragg reflections, and thus they are guides that tell the crystallographer the angles at which all reflections will occur.

Figure 4.9*a* shows an **ab** section of lattice with an arbitrary lattice point *O* chosen as the origin of the reciprocal lattice I am about to define. This point is thus the origin for both the real and reciprocal lattices. Each red + in the figure is a real lattice point.

Through a neighboring lattice point N, draw one plane from each of the sets (110), (120), (130), and so forth. These planes intersect the **ab** section in lines labeled (110), (120), and (130) in Fig. 4.9*a*. From the origin, draw a line normal to the (110) plane. Make the length of this line  $1/d_{110}$ , the inverse of the interplanar spacing  $d_{110}$ . Define the reciprocal lattice point 110 as the point at the end of this line (green dot). Now repeat the procedure for the (120) plane, drawing a line from O normal to the (120) plane, and of length  $1/d_{120}$ . Because  $d_{120}$  is smaller than  $d_{110}$  (recall that d decreases as indices increase), this second line is longer than the first. The end of this line defines a second reciprocal lattice point, with indices 120 (green dot). Repeat for the planes (130), (140), and so forth. Notice that the reciprocal lattice points lie on a straight line.

Now continue this operation for planes (210), (310), (410), and so on, defining reciprocal lattice points 210, 310, 410, and so on (Fig. 4.9*b*). Note that the points





**Figure 4.9** (*a*) Construction of reciprocal lattice. Real-lattice points are red + signs, and reciprocal lattice points are green dots. Notice the real cell edges **b** and a reciprocal cell edge **b**\*. (*b*) Continuation of (*a*). Notice the real cell edges **a** and a reciprocal cell edge **a**\*.

b



Figure 4.10 ► Reciprocal unit cells of large and small real cells.

defined by continuing these operations form a lattice, with the arbitrarily chosen real lattice point as the origin (indices 000). This new lattice (green dots) is the reciprocal lattice. The planes hk0, h0l, and 0kl correspond, respectively, to the xy, xz, and yz planes. They intersect at the origin and are called the *zero-level planes* in this lattice. Other planes of reciprocal-lattice points parallel to the zero-level planes are called *upper-level planes*.

We can also speak of the reciprocal unit cell in such a lattice (Fig. 4.10). If the angles  $\alpha$ ,  $\beta$ , and  $\lambda$  in the real cell (red) are 90°, the reciprocal unit cell (green) has axes **a**\* lying along (colinear with) real unit cell edge **a**, **b**\* lying along **b**, and **c**\* along **c**. The lengths of edges **a**\*, **b**\*, and **c**\* are reciprocals of the lengths of corresponding real cell edges **a**, **b**, and **c**: **a**\* = 1/a, and so forth, so small real cells have large reciprocal cells and vice versa. If axial lengths are expressed in angstroms, then reciprocal-lattice spacings are in the unit 1/Å or  $\text{Å}^{-1}$  (reciprocal angstroms).

For real unit cells with nonorthogonal axes, the spatial relationships between the real and reciprocal unit-cell edges are more complicated. Examples of monoclinic real and reciprocal unit cells are shown in Fig. 4.11 with a brief explanation. I will make no further use of nonorthogonal unit cells in this book (much to your relief, I expect).

Now envision this lattice of imaginary points surrounding the crystal in space. For a small real unit cell, interplanar spacings  $d_{hkl}$  are small, and hence the lines from the origin to the reciprocal lattice points are long. Therefore, the reciprocal unit cell is large, and lattice points are widely spaced. On the other hand, if the real unit cell is large, the reciprocal unit cell is small, and reciprocal space is densely populated with reciprocal lattice points.





**Figure 4.11** Example of real and reciprocal cells in the monoclinic system (stereo). As in Fig. 4.10, the real monoclinic unit cell is red, and its reciprocal unit cell is green. Even when unit cell angles are not  $90^{\circ}$ , the following relationship always holds: **a**\* is perpendicular to the real-space plane **bc**, **b**\* is perpendicular to **ac**, and **c**\* is perpendicular to **ab**. In this monoclinic cell, **b**\* and **b** are colinear (both perpendicular to **ac**), but **a**\* and **c**\* are not colinear with corresponding real axes **a** and **c**.

The reciprocal lattice is spatially linked to the crystal because of the way the lattice points are defined, so if we rotate the crystal, the reciprocal lattice rotates with it. So now when you think of a crystal, and imagine the many identical unit cells stretching out in all directions (real space), imagine also a lattice of points in reciprocal space, points whose lattice spacing is inversely proportional to the interplanar spacings within the crystal.

You are now in a position to understand the logic behind the mnemonic device for signs of Miller indices (p. 54). Recall that reciprocal lattice points are constructed on normals or perpendiculars to Miller planes. As a result, if the real and reciprocal lattices are superimposed as they were during the construction of Fig. 4.9, p. 58, a normal to a Miller plane points toward the corresponding reciprocal-lattice point. If the normal points toward the (++) quadrant of a two-dimensional coordinate system, then the corresponding reciprocal-lattice point lies in that (++) quadrant of reciprocal space, and thus its corresponding plane is assigned positive signs for both of its indices. The same normal also points toward the (--) quadrant, where the corresponding (--) reciprocal-lattice point lies, so the same set of planes can also be assigned negative signs for both indices.

### 4.2.5 Bragg's law in reciprocal space

Now I will look at diffraction from within reciprocal space. I will show that the reciprocal-lattice points give the crystallographer a convenient way to compute the direction of diffracted beams from all sets of parallel planes in the crystalline lattice (real space). This demonstration entails showing how each reciprocal-lattice point must be arranged with respect to the X-ray beam in order to satisfy Bragg's

law and produce a reflection from the crystal. It will also show how to predict the direction of the diffracted ray.

Figure 4.12*a* shows an  $\mathbf{a}^*\mathbf{b}^*$  plane of reciprocal lattice. Assume that an X-ray beam (arrow *XO*) impinges upon the crystal along this plane. Point *O* is arbitrarily chosen as the origin of the reciprocal lattice. (Remember that *O* is also a *reallattice* point in the crystal.) Imagine the X-ray beam passing through *O* along the line *XO* (arrow). Draw a circle of radius  $1/\lambda$  having its center *C* on *XO* and passing through *O*. This circle represents the wavelength of the X-rays in reciprocal space. (If the wavelength is  $\lambda$  in real space, it is  $1/\lambda$  in reciprocal space.) Rotating the crystal about *O* rotates the reciprocal lattice about *O*, successively bringing reciprocal lattice points like *P* and *P'* into contact with the circle. In Fig. 4.12*a*, *P* (whose indices are *hkl*) is in contact with the circle, and the lines *OP* and *BP* are drawn. The angle *PBO* is  $\theta$ . Because the triangle *PBO* is inscribed in a semicircle, it is a right triangle and

$$\sin \theta = \frac{OP}{OB} = \frac{OP}{2/\lambda}.$$
(4.2)

Rearranging Eq. 4.2 gives

$$2\frac{1}{OP}\sin\theta = \lambda. \tag{4.3}$$

Because *P* is a reciprocal lattice point, the length of the line *OP* is  $1/d_{hkl}$ , where *h*, *k*, and *l* are the indices of the set of planes represented by *P*. (Recall from the construction of the reciprocal lattice that the length of a line from *O* to a reciprocal-lattice point *hkl* is  $1/d_{hkl}$ .) So  $1/OP = d_{hkl}$  and

$$2d_{hkl}\sin\theta = \lambda,\tag{4.4}$$

which is Bragg's law with n = 1. So Bragg's law is satisfied, and reflection occurs, when a reciprocal-lattice point touches this circle.

In Fig. 4.12*b*, the crystal, and hence the reciprocal lattice, has been rotated clockwise about origin *O* until *P'*, with indices h'k'l', touches the circle. The same construction as in Fig. 4.12*a* now shows that

$$2d_{h'k'l'}\sin\theta = \lambda. \tag{4.5}$$

We can conclude that whenever the crystal is rotated about origin O so that a reciprocal-lattice point comes in contact with this circle of radius  $1/\lambda$ , Bragg's law is satisfied and a reflection occurs.

What direction does the diffracted beam take? Recall (from construction of the reciprocal lattice) that the line defining a reciprocal-lattice point is normal to the set of planes having the same indices as the point. So BP, which is perpendicular to OP, is parallel to the planes that are producing reflection P in Fig. 4.12a.



**Figure 4.12**  $\triangleright$  Diffraction in reciprocal space. (a) Ray *R* emerges from the crystal when reciprocal lattice point *P* intersects the circle. (b) As the crystal rotates clockwise around origin *O*, point *P'* intersects the circle, producing ray *R'*.

#### Section 4.2 Geometric principles of diffraction

If we draw a line (red) parallel to *BP* and passing through *C*, the center of the circle, this line (or any other line parallel to it and separated from it by an integral multiple of  $d_{hkl}$ ) represents a plane in the set that reflects the X-ray beam under these conditions. The beam impinges upon this plane at the angle  $\theta$ , is reflected at the same angle, and so diverges from the beam at *C* by the angle  $2\theta$ , which takes it precisely through the point *P*. So *CP* gives the direction of the reflected ray *R* in Fig. 4.12*a*. In Fig. 4.12*b*, the reflected ray *R'* follows a different path, the line *CP'*.

The conclusion that reflection occurs in the direction *CP* when reciprocal latticepoint *P* comes in contact with this circle also holds for all points on all circles produced by rotating the circle of radius  $1/\lambda$  about the X-ray beam. The figure that results, called the *sphere of reflection*, or the *Ewald sphere*, is shown in Fig. 4.13 intersecting the reciprocal-lattice planes *h0l* and *h1l*. In the crystal orientation shown, reciprocal-lattice point 012 is in contact with the sphere, so a diffracted ray *R* is diverging from the source beam in the direction defined by *C* and point 012. This ray would be detected as the 012 reflection.

As the crystal is rotated in the X-ray beam, various reciprocal-lattice points come into contact with this sphere, each producing a beam in the direction of a line from the center of the sphere of reflection through the reciprocal-lattice point that is in contact with the sphere. The reflection produced when reciprocal-lattice point  $P_{hkl}$ contacts the sphere is called the *hkl* reflection and, according to Bragg's model, is caused by reflection from the set of equivalent, parallel, real-space planes (*hkl*).

This model of diffraction implies that the directions of reflection, as well as the number of reflections, depend only upon unit-cell dimensions and not on the contents of the unit cell. As stated earlier, the *intensity* of reflection *hkl* depends upon the amount of electron density, or the average value of  $\rho(x, y, z)$ , on planes (*hkl*).



**Figure 4.13**  $\triangleright$  Sphere of reflection. When reciprocal lattice point 012 intersects the sphere, ray *R* emerges from the crystal as reflection 012.

I will show (Chapter 5) that the intensities of the reflections give us the structural information we seek.

# 4.2.6 Number of measurable reflections

If the sphere of reflection has a radius of  $1/\lambda$ , and the crystal is rotated about origin 000 on the surface of this sphere, then any reciprocal-lattice point within a distance  $2/\lambda$  of the origin can be rotated into contact with the sphere of reflection (Fig. 4.14). This distance defines the *limiting sphere*. The number of reciprocal lattice points within the limiting sphere is equal to the number of reflections that can be produced by rotating the crystal through all possible orientations in the X-ray beam. This demonstrates that the unit-cell dimensions and the wavelength of the X-rays determine the number of measurable reflections. Shorter wavelengths make a larger sphere of reflection, bringing more reflections into the measurable realm. Larger unit cells mean smaller reciprocal unit cells, which populate the limiting sphere more densely, also increasing the number of measurable reflections.

Because there is one lattice point per reciprocal unit cell (one-eighth of each lattice point lies within each of the eight unit-cell vertices), the number of reflections within the limiting sphere is approximately the number of reciprocal unit cells within this sphere. So the number N of possible reflections equals the volume of the limiting sphere divided by the volume  $V_{\text{recip}}$  of one reciprocal cell. The volume of a sphere of radius r is  $(4\pi/3)r^3$ , and r for the limiting sphere is  $2/\lambda$ , so



 $N = \frac{(4\pi/3) \cdot (2/\lambda)^3}{V_{\text{recip}}}.$ (4.6)

**Figure 4.14**  $\blacktriangleright$  Limiting sphere. All reciprocal-lattice points within the limiting sphere of radius  $2/\lambda$  can be rotated through the sphere of reflection.

#### Section 4.2 Geometric principles of diffraction

The volume V of the real unit cell is  $V_{\text{recip}}^{-1}$ , so

$$N = \frac{33.5 \cdot V}{\lambda^3}.\tag{4.7}$$

Equation 4.7 shows that the number of available reflections depends only on V and  $\lambda$ , the unit-cell volume and the wavelength of the X radiation. For a modest-size protein unit cell of dimensions  $40 \times 60 \times 80$  Å, 1.54-Å radiation can produce  $1.76 \times 10^6$  reflections, an overwhelming amount of data. Fortunately, because of cell and reciprocal-lattice symmetry, not all of these reflections are unique (Sec. 4.3.7, p. 88). Still, getting most of the available information from the diffraction experiment with protein crystals usually requires measuring somewhere between one thousand and one million reflections.

It can be further shown that the limit of resolution in an image derived from diffraction information is roughly equal to 0.707 times  $d_{\min}$ , the minimum interplanar spacing that gives a measurable reflection at the wavelength of the X-rays. For instance, with 1.54-Å radiation, the resolution attainable from all the available data is 0.8 Å, which is more than needed to resolve atoms. A resolution of 1.5 Å, which barely resolves adjacent atoms, can be obtained from about half the available data. Interpretable electron-density maps can usually be obtained with data only out to 2.5 or 3 Å. The number of reflections out to 2.5 Å is roughly the volume of a limiting reciprocal sphere of radius 1/(2.5Å) multiplied by the volume of the real unit cell. For the unit cell in the preceding example, this gives about 50,000 reflections. (For a sample calculation, see Chapter 8.)

#### 4.2.7 Unit-cell dimensions

Because reciprocal-lattice spacings determine the angles of reflection, the spacings of reflections on the detector are related to reciprocal-lattice spacings. (The exact relationship depends on the geometry of recording the reflections, as discussed later.) Reciprocal-lattice spacings, in turn, are simply the inverse of real-lattice spacings. So the distances between reflections on the detector and the dimensions of the unit cell are closely connected, making it possible to measure unit-cell dimensions from reflection spacings. I will discuss the exact geometric relationship in Sec. 4.3.6, p. 86, in the context of data-collection devices, whose geometry determines the method of computing unit-cell size.

### 4.2.8 Unit-cell symmetry

If the unit-cell contents are symmetric, then the reciprocal lattice is also symmetric and certain sets of reflections are equivalent. In theory, only one member of each set of equivalent reflections need be measured, so awareness of unit-cell symmetry can greatly reduce the magnitude of data collection. In practice, redundancy of measurements improves accuracy, so when more than one equivalent reflection is observed (measured), or when the same reflection is observed more than once, the average of these multiple observations is considered more accurate than any single observation.

#### Chapter 4 Collecting Diffraction Data

In this section, I will discuss some of the simplest aspects of unit-cell symmetry. Crystallography in practice requires detailed understanding of these matters, but users of crystallographic models need only understand their general importance. As I will show later (Sec. 4.3.7, p. 88, and Chapter 5), the crystallographer can determine the unit-cell symmetry from a limited amount of X-ray data and thus can devise a strategy for data collection that will control the redundancy of observations of equivalent reflections. While today's data-collection software can often make such decisions automatically, sometimes it can do no more than reduce the decision to several alternatives. The crystallographer must be able to check the decisions of the software, and know how to make the correct choice if the software offers alternatives.

When it comes to data collection, the *internal* symmetry of the unit cell is fundamental, and the equalities of Fig. 4.2, p. 51 are actually identities. That is, when Fig. 4.2 says  $\mathbf{a} = \mathbf{b}$ , it means not only that the axis lengths are the same, but that the contents of the unit cell must be identical along those axes. So the equalities must reflect symmetry within the unit cell. For example, a cell in which edges  $\mathbf{a}$ ,  $\mathbf{b}$ , and  $\mathbf{c}$  are equal, and cell angles  $\alpha$ ,  $\beta$ , and  $\lambda$  are 90° to within the tolerance of experimental measurement, would appear to be cubic. But it might actually be triclinic if it possesses no internal symmetry.

The symmetry of a unit cell and its contents is described by its *space group*, which describes the cell's internal symmetry elements. Space group is designated by a cryptic symbol (like  $P2_12_12_1$ ), in which a capital letter indicates the *lattice type* and the other symbols represent *symmetry operations* (defined below) that can be carried out on the unit cell without changing its appearance. Mathematicians in the late 1800s showed that there are exactly 230 possible space groups. The unit cells of a few lattice types in cubic crystal systems are shown in Fig. 4.15. *P* designates a *primitive* lattice, containing one lattice point at each comer or vertex of the cell. Because each lattice point is shared among eight neighboring unit cells, a primitive lattice contains eight times one-eighth or one lattice point per unit cell. Symbol *I* designates a body-centered or *internal* lattice, with an additional lattice point in the center of the cell, and thus two lattice points per unit cell.



**Figure 4.15** Frimitive (*P*), body-centered or internal (*I*), and face-centered (*F*) unit cells.



**Figure 4.16**  $\triangleright$  Of the several possible ways to divide a lattice into unit cells, the preferred choice is the unit cell that is most symmetrical. Centered unit cells are chosen when they are more symmetrical than any of the possible primitive cells.

Symbol *F* designates a *face-centered* lattice, with additional lattice points (beyond the primitive ones) on the centers of some or all faces.

There are usually several ways to choose the unit cell in any lattice (Fig 4.16). By convention, the choice is the unit cell that is most symmetrical. Body-centered and face-centered unit cells are chosen when they have higher symmetry than any of the primitive unit-cell choices. In Fig. 4.16, the first three choices are primitive, but are not as symmetrical as the centered cell on the right. With the addition of body- and face-centered lattices to the primitive lattices of Fig. 4.2, p. 51, there are just 13 allowable lattices, known as the *Bravais lattices*.

After defining the lattice type, the second part of identifying the space group is to describe the internal symmetry of the unit cell, using symmetry operations. To illustrate with a familiar object, one end of rectangular table looks just like a mirror reflection of the other end (Fig. 4.17*a*). We say that the table possesses a *mirror plane* of symmetry, cutting the table perpendicularly across its center (actually, it has two mirror planes; where is the other one?). In addition, if the rectangular table looks just the same as it did before rotation (ignoring imperfections such as coffee stains). We say that the table also possesses a twofold rotation axis because, in rotating the table one full circle about this axis, we find two positions that are equivalent:  $0^{\circ}$  and  $180^{\circ}$ . The mirrors and the twofold axis are examples of a *symmetry elements*.

Protein molecules are inherently asymmetric, being composed of chiral aminoacid residues coiled into larger chiral structures such as right-handed helices or twisted beta sheets. If only one protein molecule occupies a unit cell, then the cell itself is chiral, and there are no symmetry elements, as in Fig. 4.17*b*, when one chiral object is placed anywhere on the table, destroying all symmetry



**Figure 4.17**  $\triangleright$  (*a*) Table with two symmetry elements, mirror plane and twofold rotation axis. (*b*) Placing a chiral object anywhere on the table destroys all symmetry, but (*c*) if two identical chiral objects are properly placed, they restore twofold rotational symmetry. No placement of them will restore mirror symmetry.

in the figure. This situation is rare; in most cases, the unit cell contains several identical molecules or oligomeric complexes in an arrangement that produces symmetry elements. In the unit cell, the largest aggregate of molecules that possesses no symmetry elements, but can be juxtaposed on other identical entities by symmetry operations, is called the *asymmetric unit*. In the simplest case for proteins, the asymmetric unit is a single protein molecule.

The simplest symmetry operations and elements needed to describe unit-cell symmetry are translation, rotation (element: rotation axis), and reflection (element: mirror plane). Combinations of these elements produce more complex symmetry elements, including centers of symmetry, screw axes, and glide planes (discussed later). Because proteins are inherently asymmetric, mirror planes and more complex elements involving them are not found in unit cells of proteins. For example, notice in Fig. 4.17*c* that proper placement of two identical chiral objects restores twofold rotational symmetry, but no placement of the two objects will restore mirror symmetry. Because of protein chirality, symmetry elements in protein crystals include only translations, rotations, and screw axes, which are rotations and translations combined. This limitation on symmetry of unit cells containing asymmetrical objects reduces the number of space groups for chiral molecules from 230 to 65.

Now let us examine some specific symmetry operations. *Translation* simply means movement by a specified distance. For example, by the definition of *unit cell*, movement of its contents along one of the unit-cell axes by a distance equal to the length of that axis superimposes the atoms of the cell on corresponding atoms in the neighboring cell. This translation by one axial length is called a *unit translation*. Unit cells often exhibit symmetry elements that entail translations by a simple fraction of axial length, such as a/4.

In the space-group symbols, rotation axes such as the twofold axis of the table in Fig. 4.17*a* or *c* are represented in general by the symbol *n* and specifically by a number. For example, 4 means a fourfold rotation axis. If the unit cell possesses this symmetry element, then it has the same appearance after each 90° rotation around the axis.

The screw axis results from a combination of rotation and translation. The symbol  $n_m$  represents an *n*-fold screw axis with a translation of m/n of the unit translation. For example, Fig. 4.18 shows models of the amino acid alanine on a 3<sub>1</sub> screw axis in a hypothetical unit cell. On the screw axis, each successive molecule is rotated by 120° (360°/3) with respect to the previous one, and translated one-third of the axial length in the direction of the rotation axis.

Figure 4.19 shows alanine in hypothetical unit cells of two space groups. A triclinic unit cell (Fig. 4.19*a*) is designated *P*1, being a primitive lattice with only a one-fold axis of symmetry (that is, with no symmetry). *P*2<sub>1</sub> (Fig. 4.19*b*) describes a primitive unit cell possessing a twofold screw axis parallel to **c**, which points toward you as you view the figure. Notice that along any 2<sub>1</sub> screw axis, successive alanines are rotated 180° and translated one-half the axis length. A cell in space group *P*2<sub>1</sub>2<sub>1</sub>2<sub>1</sub> possesses three perpendicular twofold screw axes.

#### Chapter 4 Collecting Diffraction Data



**Figure 4.18**  $\blacktriangleright$  Three alanine molecules on a 3<sub>1</sub> screw axis in a hypothetical unit cell (stereo).

Because crystallographers deal with unit cell contents by specifying their x, y, and z coordinates, one of the most useful ways to describe unit-cell symmetry is by *equivalent positions*, locations in the unit cell that are superimposed on each other by the symmetry operations. In a P2<sub>1</sub> cell with the screw axis on or parallel to cell axis **b**, for an atom located at (x, y, z), an identical atom can always be found at (-x, y + 1/2, -z) (more commonly written  $(\overline{x}, y + 1/2, \overline{z})$ , because the operation of a 2<sub>1</sub> screw axis interchanges these positions. So a P2<sub>1</sub> cell has the equivalent positions (x, y, z) and  $(\overline{x}, y+1/2, \overline{z})$ . (The 1/2 means one-half of a unit translation along **b**, or a distance b/2 along the y-axis.) To see an example, look again at Fig. 4.19*b*, and focus on the carboxyl carbon atom of the alanine at bottom front left of the figure. If the **b**-axis is perpendicular to the page, with **a** horizontal and **c** vertical, you can see that moving this atom from position (x, y, z) to position (-x, y + 1/2, -z) will superimpose it on the next alanine carboxy carbon along the **b**-axis.

Lists of equivalent positions for the 230 space groups can be found in *International Tables for X-ray Crystallography*, a reference series that contains an enormous amount of practical information that crystallographers need in their daily work. So the easiest way to see how asymmetric units are arranged in a cell of complex symmetry is to look up the space group in *International Tables*. Each entry contains a list of equivalent positions for that space group, and several types of diagrams of the unit cell. The entry for space group  $P2_1$  is shown in Fig. 4.20.



**Figure 4.19** Alanine molecules in P1(a) and  $P2_1(b)$  unit cells (stereo).

Certain symmetry elements in the unit cell announce themselves in the diffraction pattern by causing specific reflections to be missing (intensity of zero). For example, a twofold screw axis (2<sub>1</sub>) along the **b** edge causes all 0k0 reflections having odd values of k to be missing. Notice in Fig. 4.20 that "Reflection conditions" in a P2<sub>1</sub> cell includes the condition that, along the (0k0) axis, k = 2n, meaning that only the even-numbered reflections are present along the k-axis. So the missing reflections include 010, 030, 050, and so forth. As another example, body-centered (I) lattices show missing reflections for all values of hkl where the sum of h, k, and l is odd. International Tables will list among "Reflections in the (0kl) plane are missing, including reflections 010, 001, 030, 003, and so



п

Figure 4.20 ► Entry for space group P21 in International Tables for Crystallography, Brief Teaching Edition of Volume A, Space-Group Symmetry, Theo Hanh, ed., Kluwer Academic Publishers, Norwell, MA, 5th, revised edition 2002, pp. 91-92. Reprinted with kind permission from Kluwer Academic Publishers and the International Union of Crystallography.

forth, giving the reflections in zero-level plane a diamond-shaped pattern. These patterns of missing reflections are called *systematic absences*, and they allow the crystallographer to determine the space group by looking at a few crucial planes of reflections. I will show later in this chapter how symmetry guides the strategy of data collection. In Chapter 5, I will show why symmetry causes systematic absences.

# 4.3 Collecting X-ray diffraction data

# 4.3.1 Introduction

Simply stated, the goal of data collection is to determine the indices and record the intensities of as many reflections as possible, as rapidly and efficiently as possible. One cause for urgency is that crystals, especially those of macromolecules, deteriorate in the beam because X-rays generate heat and reactive free radicals in the crystal. Thus the crystallographer would like to capture as many reflections as possible during every moment of irradiation. Often the diffracting power of the crystal limits the number of available reflections. Protein crystals that produce measurable reflections from interplanar spacings down to about 3 Å or less are usually suitable for structure determination.

In the following sections, I will discuss briefly a few of the major instruments employed in data collection. These include the X-ray sources, which produce an intense, narrow beam of radiation; detectors, which allow quantitative measurement of reflection intensities; and cameras, which control the orientation of the crystal in the X-ray beam, and thus direct reflections having known indices to detectors.

# 4.3.2 X-ray sources

X-rays are electromagnetic radiation of wavelengths 0.1–100 Å. X-rays in the useful range for crystallography can be produced by bombarding a metal target (most commonly copper, molybdenum, or chromium) with electrons produced by a heated filament and accelerated by an electric field. A high-energy electron collides with and displaces an electron from a low-lying orbital in a target metal atom. Then an electron from a higher orbital drops into the resulting vacancy, emitting its excess energy as an X-ray photon.

The metal in the target exhibits narrow *characteristic lines* (specific wavelengths) of emission resulting from the characteristic energy-level spacing of that element. The wavelengths of emission lines are longer for elements of lower atomic number Z. Electronic shells of atoms are designated, starting from the lowest level, as K, L, M, . . . . Electrons dropping from the L shell of copper (Z = 29) to replace displaced K electrons ( $L \rightarrow K$  or  $K_{\alpha}$  transition) emit X-rays of  $\lambda = 1.54$  Å. The M  $\rightarrow K$  transition produces a nearby emission band ( $K_{\beta}$ ) at 1.39 Å (Fig. 4.21*a*, solid curve). For molybdenum (Z = 42),  $\lambda(K_{\alpha}) = 0.71$  Å, and  $\lambda(K_{\beta}) = 0.63$  Å.



**Figure 4.21** (a) Emission (solid red line) and absorption (dashed green line) spectra of copper. (b) Emission spectrum of copper (solid red) and absorption spectrum of nickel (dashed green). Notice that Ni absorbs copper  $K_{\beta}$  more strongly than  $K_{\alpha}$ .

A monochromatic (single-wavelength) source of X-rays is desirable for crystallography because the radius of the sphere of reflection is  $1/\lambda$ . A source producing two distinct wavelengths of radiation gives two spheres of reflection and two interspersed sets of reflections, making indexing difficult or impossible because of overlapping reflections. Elements like copper and molybdenum make good X-ray sources if the weaker K<sub> $\beta$ </sub> radiation can be removed.

At wavelengths away from the characteristic emission lines, each element absorbs X-rays. The magnitude of absorption increases with increasing X-ray wavelength and then drops sharply just at the wavelength of  $K_{\beta}$ . The green curve in Fig. 4.21*a* shows the absorption spectrum for copper. The wavelength of this *absorption edge*, or sharp drop in absorption, like that of characteristic emission lines, increases as Z decreases such that the absorption edge for element Z - 1 lies slightly above the  $K_{\beta}$  emission line of element Z. This makes element Z - 1 an effective  $K_{\beta}$  filter for element Z, leaving almost pure monochromatic  $K_{\alpha}$  radiation. For example, a nickel filter 0.015 mm in thickness reduces Cu– $K_{\beta}$  radiation to about 0.01 times the intensity of Cu– $K_{\alpha}$ . Figure 4.21*b* shows the copper emission spectrum (red) and the nickel absorption spectrum (green). Notice that Ni absorbs strongly at the wavelength of Cu– $K_{\beta}$  radiation, but transmits Cu– $K_{\alpha}$ .

There are three common X-ray sources, *X-ray tubes* (actually a cathode ray tube sort of like a television tube), *rotating anode tubes*, and *particle storage rings*, which produce synchrotron radiation in the X-ray region. In the X-ray tube, electrons from a hot filament (cathode) are accelerated by electrically charged



**Figure 4.22** ► (*a*) X-ray tube. (*b*) Rotating anode source.

plates and collide with a water-cooled anode made of the target metal (Fig. 4.22*a*). X-rays are produced at low angles from the anode, and emerge from the tube through windows of beryllium.

Output from X-ray tubes is limited by the amount of heat that can be dissipated from the anode by circulating water. Higher X-ray output can be obtained from rotating anode sources, in which the target is a rapidly rotating metal disk (Fig. 4.22*b*). This arrangement improves heat dissipation by spreading the powerful electron bombardment over a much larger piece of metal. Rotating anode sources are more than ten times as powerful as tubes with fixed anodes.

Particle storage rings, which are associated with the particle accelerators used by physicists to study subatomic particles, are the most powerful X-ray sources. In these giant rings, electrons or positrons circulate at velocities near the speed of light, driven by energy from radio-frequency transmitters and maintained in circular motion by powerful magnets. A charged body like an electron emits energy (synchrotron radiation) when forced into curved motion, and in accelerators, the energy is emitted as X-rays. Accessory devices called *wigglers* cause additional bending of the beam, thus increasing the intensity of radiation. Systems of focusing mirrors and monochromators tangential to the storage ring provide powerful monochromatic X-rays at selectable wavelengths.

A particle-storage ring designed expressly for producing X-rays, part of the National Synchrotron Light Source (NSLS) at Brookhaven National Laboratory on Long Island in New York, is shown in Fig. 4.23. In the interior floor plan

# Chapter 4 Collecting Diffraction Data



a



b

**Figure 4.23** ► National Synchrotron Light Source, Brookhaven National Laboratory, Brookhaven NY. (*a*) Aerial view of exterior. (*b*) Interior floor plan.

(Fig. 4.18*b*), paths of particle-storage rings are shown as dark circles. The largest of the three rings in the building is the X-ray ring, which is 54 meters (177 feet) in diameter. The dark lines tangent to rings represent beam lines for work in experimental stations (green blocks). NSLS began operations in 1984. Crystallographers can apply to NSLS and other synchrotron sources for grants of time for

data collection. For a detailed virtual tour of NSLS and other synchrotron sources, see the CMCC home page.

Although synchrotron sources are available only at storage rings and require the crystallographer to collect data away from the usual site of work, there are many advantages that compensate for the inconvenience. X-ray data that requires several hours of exposure to a rotating anode source can often be obtained in seconds or minutes at a synchrotron source like NSLS. In a day or two at a synchrotron source, a crystallographer can collect data that might take weeks to acquire with conventional sources. Another advantage, as I will show in Chapter 6, is that X-rays of selectable wavelength can be helpful in solving the phase problem.

Whatever the source of X-rays, the beam is directed through a *collimator*, a narrow metal tube that selects and reflects the X-rays into parallel paths, producing a narrow beam. After collimation, beam diameter can be further reduced with systems of metal plates called *focusing mirrors*. In the ideal arrangement of source, collimators, and crystal, all points on the crystal can "see" through the collimator and mirrors to all line-of-sight points on the X-ray source.

SAFETY NOTE: X-ray sources pose dangers that the crystallographer must consider in daily work. X-ray tubes require high-voltage power supplies containing large condensers that can produce a dangerous shock even after equipment is shut off. The X-rays themselves are relatively nonpenetrating, but can cause serious damage to surface tissues. Even brief exposure to weak X-rays can damage eyes, so protective goggles are standard attire in the vicinity of X-ray sources. The direct beam is especially powerful, and sources are electronically interlocked so that the beam entrance shutter cannot be opened while the user is working with the equipment. The beam intensity is always reduced to a minimum during alignment of collimating mirrors or cameras. During data collection, the direct beam is blocked just beyond the crystal by a piece of metal called a *beam stop*, which also has the beneficial effect of preventing excessive radiation from reaching the center of the detector, thus obscuring low-angle reflections. In addition, the entire source, camera, and detector are usually surrounded by Plexiglas to block scattered radiation from the beam stop or collimators but to allow observation of the equipment. As a check on the efficacy of measures to prevent X-ray exposure, the prudent crystallographer wears a dosage-measuring ring or badge during all work with X-ray equipment. These devices are periodically sent to radiation-safety labs for measurement of the X-ray dose received by the worker.

## 4.3.3 Detectors

Reflection intensities can be measured by *scintillation counters*, which in essence count X-ray photons and thus give quite accurate intensities over a wide range. Scintillation counters contain a phosphorescent material that produces a flash of light (a scintillation) when it absorbs an X-ray photon. A photocell counts the flashes. With simple scintillation counters, each reflection must be measured separately, an arrangement that was the basis of *diffractometry*, which is now used only for small-molecule crystallography, where the number of reflections per data

set is very small. Macromolecular crystallography requires means to collect many reflections at once. Such devices are called *area detectors*.

The simplest X-ray area detector, and for years the workhorse of detectors, is X-ray-sensitive film (for example, see Fig. 2.7, p. 14), but film has been replaced by image plates and CCD detectors, which I will describe below. Various types of cameras (next section) can direct reflections to area detectors in useful arrangements, allowing precise determination of indices and intensities for thousands of reflections from a single image.

Image plate detectors are somewhat like reusable films that can store diffraction images reversibly and have a very wide dynamic range, or capacity to record reflections of widely varying intensity. Image plates are plastic sheets with a coating of small crystals of a phosphor, such as BaF:Eu<sup>2+</sup>. The crystals can be stimulated by X-rays into a stable excited state in which Eu<sup>2+</sup> loses an electron to the F<sup>-</sup> layer, which contains electron vacancies introduced by the manufacturing process. Further stimulation by visible light causes the electrons to drop back to the Eu layer, producing visible light in proportion to the intensity of the previously absorbed X-rays. After X-ray exposure, data are read from the plate by a scanner in which a fine laser beam induces luminescence from a very small area of the plate, and a photocell records the intensity of emitted light. The intensities are fed to a computer, which can then reconstruct an image of the diffraction pattern. Image plates can be erased by exposure to bright visible light and reused indefinitely.

*Multiwire area detectors* were the first type to combine the accuracy and wide dynamic range of scintillation counting; the simultaneous measurement of many reflections, as with image plates; and the advantage of direct collection of data by computer, without a separate scanning step. In these detectors, some of which are still in use, two oppositely charged, perpendicular sets of parallel wires in an inert gas detect and accurately locate ionization of the gas induced by single X-ray photons, relaying to a computer the positions of X-ray photons almost instantaneously as they strike the detector. The computer records events and their locations to build up an image of the reflections that reach the detector. Multiwire systems are insensitive to X-rays for a short time after recording each event, a period known as *dead time*, which sets a limit on their rate of X-ray photon counting.

The latest designs in area detectors employ *charge-coupled devices* (*CCDs*) as detectors (Fig 4.24). CCDs first found use in astronomy as light collectors of great sensitivity, and have since replaced other kinds of light detectors in many common devices, most notably digital cameras. In effect, CCDs are photon counters, solid-state devices that accumulate charge (electrons) in direct proportion to the amount of light that strikes them. An incident photon raises an electron to a higher energy level, allowing it to move into, and become trapped within, a positively charged region at the center of the pixel (Fig. 4.24*a*). Each pixel in the array can accurately accumulate electrons from several hundred thousand absorption events before reaching its capacity (becoming saturated). Before saturation, its contents must be read out to a computer. At the end of a collection cycle, which produces one frame of data, the charges are read out by a process in which rows of pixel charge



**Figure 4.24** ► (*a*) Schematic diagram of CCD (8 pixels), showing path of data readout. (*b*) Diagram of CCD X-ray detector.

are transferred sequentially, by flipping the voltage applied at the edges, into a serial readout row at one edge of the CCD. The charges in the readout row are transferred serially to an amplifier at the end of the row, and then the next row of pixel charges is transferred into the readout row. Because all data are read out at the end of data collection, a CCD, unlike a multiwire area detector, has no dead time, and thus no practical limit on its rate of photon counting.

CCDs are sensitive to visible light, not X-rays, so phosphors that produce visible light in response to X-ray absorption must be interposed between the diffracted X-rays and the CCD. In addition, the individual pixels of commercially available CCD arrays are smaller and more densely packed than is best for typical X-ray detection geometry, and in addition, some light is lost by reflection from the CCD surface. To overcome all these obstacles, CCD arrays are bonded to a tapered bundle of optical fibers (Fig 4.24*b*). The large end of the optical taper is coated with phosphors that emit visible light in response to X-rays, producing, in effect, a very dense array of scintillation counters. As an example of detector and CCD dimensions, the CCD area detector of the NSLS beamline X-12, which is designed

#### Chapter 4 Collecting Diffraction Data

especially for macromolecular crystallography, consists of four  $1150 \times 1150$ -pixel CCD arrays in a 2 × 2 array, giving a 5.3 megapixel array. The optical taper is about 188 mm (roughly 7.5 inches) on a side at the phosphor-coated face of the detector, and tapers down to about 50 mm (less than 2 inches) at the face of the CCD array. Each pixel can count 450,000 to 500,000 photons between readouts, and readouts take from 1 to 10 seconds, depending on the desired signal-to-noise ratio in the output. These numbers are typical in 2005, but such technology changes rapidly.

### 4.3.4 Cameras

Between the X-ray source and the detector lies a mounted crystal, in the grip of an *X-ray camera*. The term is a misnomer. The word *camera* is from Latin for chamber or vault, referring in common photographic cameras to the darkened chamber inside that prevents stray light from reaching the film (or CCD array). X-ray cameras have no darkened chamber; they are simply mechanical devices, usually a combination of a crystal-holding head on a system of movable circular mounts for orienting and moving crystals with great precision (to within  $10^{-5}$ degrees). The goal in data collection is to use carefully controlled movement of the crystal by the camera to direct all unique reflections to a detector like one of those described in the previous section. In this section, I will describe cameras that can rotate a crystal through a series of known orientations, causing specified reciprocal lattice points to pass through the sphere of reflection and thus produce diffracted X-ray beams.

In all forms of data collection, the crystal is mounted on a *goniometer head*, a device that allows the the crystal orientation to be set and changed precisely. The most complex goniometer heads (Fig. 4.25*a*) consist of a holder for a capillary tube or cryoloop containing the crystal; two arcs (marked by angle scales), which permit rotation of the crystal by up to  $40^{\circ}$  in each of two perpendicular planes; and two dovetailed sledges, which permit small translations of the arcs for centering the crystal on the rotation axis of the head. Simpler heads (Fig. 4.25*b*) contain sledges only. Protein crystals, either sealed in capillary tubes with mother liquor (as in *a*) or flash-frozen in a fiber loop (as in *b*) are mounted on the goniometer head, which is adjusted to center the crystal in the X-ray beam and to allow rotation of the crystal while maintaining centering. Flash-frozen crystals are held in a stream of cold nitrogen gas emerging from a reservoir of liquid nitrogen, and the goniometer head is heated to prevent condensation from forming ice on it.

The goniometer head and crystal are mounted in a system of movable circles called a *goniostat*, which allows automated, highly precise movement of the crystal into almost any orientation with respect to the X-ray beam and the detector. The crystal orientation is specified by a system of angles, whose nomeclature is a vestige of diffractometry, in which individual reflections are directed to a scintillation counter for intensity measurement (diffractometry is still in wide use for crystallography of small molecules). Figure 4.26 shows this system of angles. A complete diffractometer consists of a fixed X-ray source, the goniostat,



**Figure 4.25** (*a*) Full goniometer head, with capillary tube holder at top. The tool (right) is an Allen wrench for adjusting arcs and sledges. Photo courtesy of Charles Supper Company. (*b*) Simple goniometer head, having only sledges. The holder on top is a magnetic disk that accepts a cryoloop holder. The wiring is for heating the head, to prevent ice formation from the nitrogen stream that keeps the crystal at low temperatures. Images courtesy of Hampton Research.

and a movable scintillation-counter detector. The system of circles (Fig. 4.26) allows rotation of the goniometer head (angle  $\phi$ ), movement of the head around a circle centered on the X-ray beam (angle  $\chi$ ), and rotation of the  $\chi$  circle around an axis perpendicular to the beam (angle  $\omega$ ). Furthermore, the detector moves on a circle coplanar with the beam. The axis of this circle coincides with the  $\omega$ -axis. The position of the detector with respect to the beam is denoted by the angle  $2\theta$ . (Why? Think about it, then look at Fig. 4.12, p. 62.) With this arrangement, the crystal can be moved to bring any reciprocal lattice point that lies within the limiting sphere into the plane of the detector and into contact with the sphere of reflection, producing diffracted rays in the detector plane. The detector can be moved into proper position to receive, and measure the intensity of, the resulting diffracted beam.


**Figure 4.26**  $\blacktriangleright$  System of angles in diffractometry. The crystal in the center is mounted on a goniometer head.

Diffractometry gives highly accurate intensity measurement but is slow in comparison with methods that record many reflections at once. In addition, the total irradiation time is long, so crystals may deteriorate and have to be replaced. While one reflection is being recorded, there are usually other unmeasured reflections present, so a considerable amount of diffracted radiation is wasted. Diffractometers teamed up with area detectors give substantial increases in the efficiency of data collection, but are no match for more modern experimental arrangement, which can direct hundreds of reflections to area detectors simultaneously.

Modern, high-speed data collection relies on the *rotation/oscillation* method (Fig. 4.27). An oscillation camera is far simpler than a diffractometer, providing, at the minimum, means to rotate the crystal about an axis perpendicular to the beam ( $\phi$ -axis), as well as to oscillate it back and forth by a few degrees about the same axis. Figure 4.27 shows a diffractometer, X-ray source, and area detector set up for rotation or oscillation photography. The camera provides means for movement through all the diffractometer angles, as shown in the figure labels. In the arrangement shown, the detector and X-ray source (a rotating anode, not shown) are colinear, and rotations about  $\phi$  and  $\chi$  are used for data collection. Rotation of the crystal through a small angle casts large numbers of reflections in a complex pattern onto the area detector, as illustrated in Fig. 4.28.

Figure 4.28*a* shows how rotation of a crystal tips many planes of the reciprocal lattice through the sphere of reflection, producing many reflections. The crystal (dark red) lies in the middle of the Ewald sphere (sphere of reflection, light brown surface). The smaller sphere of gray dots represents the reciprocal lattice. Purple lines show reflections being produced at the current orientation of the crystal, from



**Figure 4.27**  $\triangleright$  X-ray instrumentation used in the 2004 course, X-ray Methods in Structural Biology, Cold Spring Harbor Laboratory. The instrument is actually a four-circle diffractometer using only the  $\phi$ -axis for rotation photography. X-rays come from the left from a rotating anode source (not shown), and emerge through a metal collimator to strike the crystal. A CCD detector (right) captures reflections. The crystal is kept frozen by a stream of cold nitrogen gas coming from a liquid nitrogen supply. The slight cloudiness just to the right of the goniometer head is condensation of moisture in the air as it is cooled by the nitrogen stream. Rotation of the crystal about the  $\phi$ - and  $\chi$  axes directs many reflections to the detector. In the arrangement shown, axes  $\phi$ ,  $2\theta$ , and  $\omega$  are coincident.

reciprocal-lattice points in contact with the Ewald sphere. The other reflections on the detector were produced as the crystal was rotated through a small angle to its current position.

Figure 4.28*b* illustrates the actual geometry of an oscillation image (called a *frame* of data). To produce this figure, a computer program calculated the indices of reflections expected when a crystal is oscillated a few degrees about its **c**-axis. At the expected position of each reflection, the program plotted the indices of that reflection. Only the *l* index of each reflection is shown here, revealing that reflections from many levels of reciprocal space are recorded at once. Although frames of oscillation data are very complex, modern software can index them.

As the crystal oscillates about a fixed starting position, many reciprocal-lattice points pass back and forth through the sphere of reflection, and their intensities are recorded. The amount of data from a single oscillation is limited only by overlap of reflections. The strategy is to collect one frame by oscillating the crystal through a small angle about a starting position of rotation, recording all the



Figure 4.28 ► (a) (Stereo) rotation or oscillation of a crystal by a few degrees tips many reciprocal-lattice planes through the sphere of reflection, sending many reflections to the detector. Image created with the free program XRayView, which allows the user to study diffraction geometry interactively. To obtain the program, see the CMCC home page. Image used with permission of Professor George N. Phillips, Jr. (b) Diagram showing expected positions of reflections in a frame of oscillation data. Diagram courtesy of Professor Michael Rossmann.

#### Section 4.3 Collecting X-ray diffraction data

resulting reflections, and then rotating the crystal to a new starting point such that the new oscillating range overlaps the previous one slightly. From this new position, oscillation produces additional reflections in a second frame. This process is continued until all unique reflections have been recorded.

In classical crystallography, the goniometer head was viewed through a microscope to orient the mounted crystal precisely, using crystal faces as guides. This allowed the crystallographer to make photographs of the reciprocal lattice in specific orientations, such as to record the zero-level planes for measuring unit-cell lengths and angles. Well-formed crystals show distinct faces that are parallel to unit-cell edges, and first attempts to obtain a diffraction pattern were made by placing a crystal face perpendicular to the X-ray beam. Preliminary photos of the diffraction pattern allowed the initial setting to be refined. Photographing the very revealing zero-level planes required a complex movement of crystal and detector, called *precession photography*. Images of zero-level planes are still sometimes called precession photographs. (Figure 2.7, p. 14 is a true precession photograph, taken on film around 1985.) Examination of the zero-level planes of diffraction allowed determination of unit-cell dimensions and space group. The crystallographer then devised a data-collection strategy that would record all unique reflections from a minimum number of crystal orientations, as described later.

None of this is necessary today, because data-collection software can quickly determine the crystal orientation from frames taken over a small range of rotation, making initial orientation of the crystal unimportant (which is good, because orientation is hard to control in cryoloops). The usual procedure is to use the goniometer head sledges to position the randomly oriented crystal so that it stays centered in the beam through all rotations, collect reflections over a few degrees of rotation, and then let the software determine the crystal orientation, index the reflections, determine unit-cell parameters and space group, and devise the collection strategy. The procedure is sometimes called the American method—shoot first, ask questions afterward. The software cannot always make the determination unequivocally, and so sometimes offers choices. The crystallographer's knowledge of crystal systems and symmetry guide the final decision. After full data collection, the software provides means to view and analyze specific sections of the data, such as the zero-level planes.

#### 4.3.5 Scaling and postrefinement of intensity data

The goal of data collection is a set of consistently measured, indexed intensities for as many of the reflections as possible. After data collection, the raw intensities must be processed to improve their consistency and to maximize the number of measurements that are sufficiently accurate to be used.

A complete set of measured intensities often includes many frames, as well as distinct blocks of data obtained from several (or many) crystals in different orientations. Because of variability in the diffracting power of crystals, the difference in the length of the X-ray path through the crystal in different orientations, and the intensity of the X-ray beam, the crystallographer cannot assume that the absolute intensities are consistent from one frame or block of data to the next. An obvious

#### Chapter 4 Collecting Diffraction Data

way to obtain this consistency is to compare reflections of the same index that were measured from more than one crystal or in more than one frame and to rescale the intensities of the two blocks of data so that identical reflections are given identical intensities. This process is called *scaling*. Scaling is often preliminary to a more complex process, *postrefinement*, which recovers usable data from reflections that were only partially measured.

Primarily because real crystals are mosaics of submicroscopic crystals (Fig. 3.2, p. 33), a reciprocal-lattice point acts as a small three-dimensional entity (sphere or ovoid) rather than as an infinitesimal point. As a reciprocal-lattice point moves through the sphere, diffraction is weak at first, peaks when the center of the point lies precisely on the sphere, and then weakens again before it is extinguished. Accurate measurement of intensity thus entails recording the X-ray output during the entire passage of the point through the sphere. Any range of oscillation will record some reflections only partially, but these may be recorded fully at another rotation angle, allowing partial reflections to be discarded from the data. The problem of partial reflections is serious for large unit cells, where smaller oscillation angles are employed to minimize overlap of reflections. In such cases, if partial reflections are discarded, then a great deal of data is lost. Data from partial reflections can be interpreted accurately through postrefinement of the intensity data. This process produces an estimate of the partiality of each reflection. Partiality is a fraction  $p \ (0 \ge p \ge 1)$  that can be used as a correction factor to convert the measured intensity of a partial reflection to an estimate of that reflection's full intensity.

Scaling and postrefinement are the final stages in producing a list of internally consistent intensities for as many of the available reflections as possible.

### 4.3.6 Determining unit-cell dimensions

The unit-cell dimensions determine the reciprocal-lattice dimensions, which in turn tell us where we must look for the data. Methods like oscillation photography require that we (or our software) know precisely which reflections will fall completely and partially within a given oscillation angle so that we can collect as many reflections as possible without overlap. So we need the unit-cell dimensions in order to devise a strategy of data collection that will give us as many identifiable (by index), measurable reflections as possible.

Modern software can search the reflections, measure their precise positions, and subsequently compute unit-cell parameters. This search entails complexities we need not encounter here. Instead, I will illustrate the simplest method for determining unit-cell dimensions: measuring reflection spacings from an orthorhombic crystal on an image of a zero-level plane. Because reciprocal-lattice spacings are the inverse of real-lattice spacings, the unit-cell dimensions are inversely proportional to the spacing of reflections on planes in reciprocal space, and determining unit-cell dimensions from reciprocal-lattice spacings is a remarkably simple geometric problem. Figure 4.29 shows the geometric relationship between reflection spacings on the film and actual reciprocal-lattice spacings.



**Figure 4.29** ► Reflection spacings on the film are directly proportional to reciprocallattice spacings, and so they are inversely proportional to unit-cell dimensions.

The crystal at *C* is precessing about its  $c^*$ -axis, and therefore recording *hk*0 reflections on the detector, with the *h*00 axis horizontal and the 0*k*0 axis vertical. Point *P* is the reciprocal-lattice point 100, in contact with the sphere of reflection, and *O* is the origin. Point *F* is the origin on the detector and *R* is the recording of reflection 100 on the detector. The distance *OP* is the reciprocal of the distance  $d_{100}$ , which is the length of unit-cell edge **a**. Because *CRF* and *CPO* are similar triangles (all corresponding angles equal), and because the radius of the sphere of reflection is  $1/\lambda$ ,

$$\frac{RF}{CF} = \frac{PO}{CO} = \frac{PO}{1/\lambda} = PO \cdot \lambda.$$
(4.8)

Therefore,

$$PO = \frac{RF}{CF \cdot \lambda}.$$
(4.9)

Because  $d_{100} = 1/PO$ ,

$$d_{100} = \frac{CF \cdot \lambda}{RF}.$$
(4.10)

In other words, the axial length a (length of unit-cell edge **a**) can be determined by dividing the crystal-to-detector distance (*CF*) by the distance from the detector origin to the 100 reflection (RF) and multiplying the quotient by the wavelength of X-rays used in taking the photograph.

In like manner, the vertical reflection spacing along 0k0 or parallel axes gives  $1/d_{010}$ , and from it, the length of unit-cell axis **b**. Because we are considering an orthorhombic crystal, which has unit-cell angles of 90°, a second zero-level image, taken after rotating this orthorhombic crystal by 90° about its vertical axis, would record the 00l axis horizontally, giving  $1/d_{001}$ , and the length of **c**.

Of course, the distance from the detector origin to the 100 reflection on a zerolevel image is the same as the distance between any two reflections along this or other horizontal lines, so one zero-level image allows many measurements to determine accurately the *average* spacing of reciprocal-lattice points along two different axes. From accurate average values, unit-cell-axis lengths can be determined with sufficient accuracy to guide a data-collection strategy.

Except perhaps for fun or curiosity, no one today works out this little geometry problem to determine unit-cell dimensions, and I was tempted to omit the topic during this revision. But nothing in crystallography shows more dramatically how simple is the relationship between the submicroscopic dimensions of the unit cell and the macroscopic dimensions of spacing between reflections at the detector. I once sat down with an X-ray photo, used a small ruler under a magnifier to measure the distances between the centers of reflections, and used a pocket calculator to compute the unimaginably small dimensions of unit cells in some of the first crystals I had ever grown. To make a measurement at the molecular level with an ordinary ruler was a magical experience.

### 4.3.7 Symmetry and the strategy of collecting data

Strategy of data collection is guided not only by the unit cell's dimensions but also by its internal symmetry. If the cell and its contents are highly symmetric, then certain sets of crystal orientations produce exactly the same reflections, reducing the number of crystal orientations needed in order to obtain all of the distinct or unique reflections.

As mentioned earlier, the unit-cell space group can be determined from systematic absences in the diffraction pattern. With the space group in hand, the crystallographer can determine the space group of the reciprocal lattice, and thus know which orientations of the crystal will give identical data. All reciprocal lattices possess a symmetry element called a *center of symmetry* or *point of inversion* at the origin. That is, the intensity of each reflection *hkl* is identical to the intensity of reflection  $\overline{hkl}$ . To see why, recall from our discussion of lattice indices (Sec. 4.2.2, p. 50) that the the index of the (230) planes can also be expressed as ( $\overline{230}$ ). In fact, the 230 and the  $\overline{230}$  reflections come from opposite sides of the same set of planes, and the reflection intensities are identical. (The equivalence of  $I_{hkl}$  and  $I_{\overline{hkl}}$  is called *Friedel's law*, but there are exceptions, as I will show in Sec. 6.4, p. 128) This means that half of the reflections in the reciprocal lattice axis will capture all unique reflections.

#### Section 4.4 Summary

Additional symmetry elements in the reciprocal lattice allow further reduction in the total angle of data collection. It can be shown that the reciprocal lattice possesses the same symmetry elements as the unit cell, plus the additional point of inversion at the origin. The 230 possible space groups reduce to only 11 different groups, called *Laue groups*, when a center of symmetry is added. For each Laue group, and thus for all reciprocal lattices, it is possible to compute the fraction of reflections that are unique. For monoclinic systems, such as P2, the center of symmetry is the only element added in the reciprocal lattice and the fraction of unique reflections is 1/4. At the other extreme, for the cubic space group P432, which possesses four-, three-, and twofold rotation axes, only 1/48 of the reflections are unique. Determination of the crystal symmetry can greatly reduce the number of reflections that must be measured. It also guides the crystallographer in choosing the best axis about which to rotate the crystal during data collection. In practice, crystallographers collect several times as many reflections as the minimum number of unique reflections. They use the redundancy to improve the signal-to-noise ratio by averaging the multiple determinations of equivalent reflections. They also use redundancy to correct for X-ray absorption, which varies with the length of the X-ray path through the crystal.

### 4.4 Summary

The result of X-ray data collection is a list of intensities, each assigned an index hkl corresponding to its position in the reciprocal lattice. The intensity assigned to reflection hkl is therefore a measure of the relative strength of the reflection from the set of lattice planes having indices hkl. Recall that indices are counted from the origin (indices 000), which lies in the direct path of the X-ray beam. In an undistorted image of the reciprocal lattice, such as an image of a zero-level plane, reflections having low indices lie near the origin, and those with high indices lie farther away. Also recall that as indices increase, there is a corresponding decrease in the spacing  $d_{hkl}$  of the real-space planes represented by the indices. This means that the reflections near the origin come from sets of widely spaced planes, and thus carry information about larger features of the molecules in the unit cell. On the other hand, the reflections far from the origin come from closely spaced lattice planes in the crystal, and thus they carry information about the fine details of structure.

In this chapter, I have shown how the dimensions and symmetry of the unit cell determine the dimensions and symmetry of the diffraction pattern. Next I will show how the molecular *contents* of the unit cell determine the *contents* (that is the reflection intensities) of the diffraction pattern. In the next three chapters, I will examine the relationship between the intensities of the reflections and the molecular structures we seek, and thus show how the crystallographer extracts structural information from the list of intensities.

This Page Intentionally Left Blank

### ► Chapter 5

# From Diffraction Data to Electron Density

### 5.1 Introduction

In producing an image of molecules from crystallographic data, the computer simulates the action of a lens, computing the electron density within the unit cell from the list of indexed intensities obtained by the methods described in Chapter 4. In this chapter, I will discuss the mathematical relationships between the crystallographic data and the electron density.

As stated in Chapter 2, computation of the Fourier transform is the lenssimulating operation that a computer performs to produce an image of molecules in the crystal. The Fourier transform describes precisely the mathematical relationship between an object and its diffraction pattern. The transform allows us to convert a Fourier-sum description of the reflections to a Fourier-sum description of the electron density. A reflection can be described by a structure-factor equation, containing one term for each atom, or for each volume element, in the unit cell. In turn, the electron density is described by a Fourier sum in which each term is a structure factor. The crystallographer uses the Fourier transform to convert the structure factors to  $\rho(x, y, z)$ , the desired electron density equation.

First I will discuss Fourier sums and the Fourier transform in general terms. I will emphasize the form of these equations and the information they contain, in the hope of helping you to interpret the equations—that is, to translate the equations into words and visual images. Then I will present the specific types of Fourier sums that represent structure factors and electron density and show how the Fourier transform interconverts them.

### 5.2 Fourier sums and the Fourier transform

### 5.2.1 One-dimensional waves

j

Recall from Sec. 2.6.1, p. 21, that waves are described by periodic functions, and that simple wave equations can be written in the form

$$f(x) = F \cos 2\pi (hx + \alpha) \tag{5.1}$$

or

$$f(x) = F \sin 2\pi (hx + \alpha), \qquad (5.2)$$

where f(x) specifies the vertical height of the wave at any horizontal position x (measured in wavelengths, where x = 1 implies one full wavelength or one full repeat of the periodic function). In these equations, F specifies the amplitude of the wave (the distance from from horizontal axis to peak or valley), h specifies its frequency (number of wavelengths per radian), and  $\alpha$  specifies its phase (position of the wave, in radians, with respect to the origin). These equations are *one*-dimensional in the sense that they represent a numerical value [f(x), the height of the wave] at all points along *one* axis, in this case, the x-axis. See Fig. 2.14, p. 22 for graphs of such equations.

I also stated in Chapter 2 (see Fig. 2.16, p. 25) that any wave, no matter how complicated, can be described as the sum of simple waves. This sum is called a *Fourier sum* and each simple wave equation in the sum is called a *Fourier term*. Either Eq. (5.1) or (5.2) could be used as a single Fourier term. For example, we can write a Fourier sum of *n* terms using Eq. (5.1) as follows:

$$f(x) = F_0 \cos 2\pi (0x + \alpha_0) + F_1 \cos 2\pi (1x + \alpha_1) + F_2 \cos 2\pi (2x + \alpha_2) + \dots + F_n \cos 2\pi (nx + \alpha_n),$$
(5.3)

or equivalently,

$$f(x) = \sum_{h=0}^{n} F_h \cos 2\pi (hx + \alpha_h).$$
(5.4)

According to Fourier theory, any complicated periodic function can be approximated by such a sum, by putting the proper values of h,  $F_h$ , and  $\alpha_h$  in each term. Think of the cosine terms as basic wave forms that can be used to build any other waveform. Also according to Fourier theory, we can use the sine function or, for that matter, *any* periodic function in the same way as the basic wave for building any other periodic function.

A very useful basic waveform is  $[\cos 2\pi (hx) + i \sin 2\pi (hx)]$ . Here, the waveforms, of cosine and sine are combined to make a *complex number*, whose general

form is a + ib, where *i* is the imaginary number  $(-1)^{1/2}$ . Although the phase  $\alpha$  of this waveform is not shown, it is implicit in the combination of the cosine and sine functions, and it depends only upon the values of *h* and *x*. As I will show in Chapter 6, expressing a Fourier term in this manner gives a clear geometric means of representing the phase  $\alpha$  and allows us to see how phases are computed. For now, just accept this convention as a convenient way to write completely general Fourier terms. In Chapter 6, I will discuss the properties of complex numbers and show how they are used to represent and compute phases.

With the terms written in this fashion, a general Fourier sum looks like this:

$$f(x) = \sum_{h=0}^{n} F_h[\cos 2\pi (hx) + i \sin 2\pi (hx)]$$
(5.5)

In words, this sum consists of *n* simple Fourier terms, one for each integral value of *h* beginning with zero and ending with *n*. Each term is a simple wave with its own amplitude  $F_h$ , its own frequency *h*, and (implicitly) its own phase  $\alpha$ .

Next, we can express the complex number in square brackets as an exponential, using the following equality from complex number theory:

$$\cos\theta + i\sin\theta = e^{i\theta}.\tag{5.6}$$

In our case,  $\theta = 2\pi (hx)$ , so the Fourier sum becomes

$$f(x) = \sum_{h=0}^{n} F_h e^{2\pi i (hx)}$$
(5.7)

or simply

$$f(x) = \sum_{h} F_{h} e^{2\pi i (hx)}, \qquad (5.8)$$

in which the sum is taken over all values of h, and the number of terms is unspecified.

I will write Fourier sums in this form throughout the remainder of the book. This kind of equation is compact and handy, but quite opaque at first encounter. Take the time now to look at this equation carefully and think about what it represents. Whenever you see an equation like this, just remember that it is a Fourier sum, a sum of sine and cosine wave equations, with the full sum representing some complicated wave. The *h*th term in the sum,  $F_h e^{2\pi i (hx)}$ , can be expanded to  $F_h[\cos 2\pi (hx) + i \sin 2\pi (hx)]$ , making plain that the *h*th term is a simple wave of amplitude  $F_h$ , frequency *h*, and implicit phase  $\alpha_h$ .

### 5.2.2 Three-dimensional waves

The Fourier sum that the crystallographer seeks is  $\rho(x, y, z)$ , the three-dimensional electron density of the molecules under study. This function is a wave equation or periodic function because it repeats itself in every unit cell. The waves described in the preceding equations are one-dimensional: they represent a numerical value f(x) that varies in one direction, along the x-axis. How do we write the equations of two-dimensional and three-dimensional waves? First, what do the graphs of such waves look like?

When you graph a function, you must use one more dimension than specified by the function. You use the additional dimension to represent the numerical value of the function. For example, in graphing f(x), you use the y-axis to show the numerical value of f(x). In Fig. 2.14, p. 22, for example, the y-axes are used to represent f(x), the height of each wave at point x. Graphing a *two*-dimensional function f(x, y) requires the *third* dimension to represent the numerical value of the function.

For example, imagine a weather map with mountains whose height at location (x, y) represents the temperature at that location. Such a map graphs a twodimensional function t(x, y), which gives the temperature t at all locations (x, y) on the plane represented by the map. If we must avoid using the third dimension, for instance in order to print a flat map, the best we can do is to draw a contour map on the plane map (Fig. 5.1), with continuous lines (contours, in this case called *isotherms*) representing locations having the same temperature.

Graphing the *three*-dimensional function  $\rho(x, y, z)$  in the same manner would require *four* dimensions, one for each of the spatial dimensions x, y, and z, and a fourth one for representing the value of  $\rho$ . Here a contour map is the only choice. In three dimensions, contours are continuous surfaces (rather than lines) on which the function has a constant numerical value. A contour map of the threedimensional wave  $\rho(x, y, z)$  exhibits surfaces of constant electron density  $\rho$ . You are already familiar with such contour maps. The common drawings of electronic orbitals (such as the 1s orbital of a hydrogen atom, often drawn as a simple sphere) is a contour map of a three-dimensional function. Everywhere on the surface of this sphere, the electron density is the same. Orbital surfaces are often drawn to enclose the region that contains 98% (or some specified value) of the total electron density.

The blue netlike surface in Fig. 2.3, p. 11 is also a contour map of a threedimensional function. It represents a surface on which the electron density  $\rho(x, y, z)$  of adipocyte lipid binding protein (ALBP) is constant. Imagine that the net encloses some specified value, say 98%, of the protein's electron density, and so the net is in essence an image of the protein's surface. The actual value of  $\rho(x, y, z)$  on the plotted surface is specified as the map's *contour level*, usually given in units of  $\sigma$ , the standard deviation of the overall electron density, from the mean electron density. For example, in a map contoured at  $2\sigma$ , the displayed surface is two standard deviations higher than the mean electron density for the whole map.



**Figure 5.1** Seasonable February morning in Maine. Lines of constant temperature (isotherms) allow plotting a two-dimensional function without using the third dimension. This is a contour map of t(x, y), giving the temperature t at all locations (x, y). Along each contour line lie all points having the same temperature. A planar contour map of a function of two variables takes the form of contour lines on the plane. In contrast, a contour map of a function of three variables takes the form of contour surfaces in three dimensions (see Fig. 2.3, p. 11).

I hope the foregoing helps you to imagine three-dimensional waves. What do the equations of such waves look like? A three-dimensional wave has three frequencies, one along each of the x-, y-, and z-axes. So three variables h, k, and l are needed to specify the frequency in each of the three directions. A general Fourier sum for the wave f(x, y, z), written in the compact form of Eq. (5.8) is as follows:

$$f(x, y, z) = \sum_{h} \sum_{k} \sum_{l} F_{hkl} e^{2\pi i (hx + ky + lz)}.$$
 (5.9)

In words, Eq. (5.9) says that the complicated three-dimensional wave f(x, y, z) can be represented by a Fourier sum. Each term in the sum is a simple threedimensional wave whose frequency is h in the x-direction, k in the y-direction, and l in the z-direction. For each possible set of values h, k, and l, the associated wave has amplitude  $F_{hkl}$  and, implicitly, phase  $\alpha_{hkl}$ . The triple sum simply means to add up terms for all possible sets of integers h, k, and l. The range of values for h, k, and l depends on how many terms are required to represent the complicated wave f(x, y, z) to the desired precision.

### 5.2.3 The Fourier transform: General features

Fourier demonstrated that for any function f(x), there exists another function F(h) such that

$$F(h) = \int_{-\infty}^{+\infty} f(x)e^{2\pi i(hx)} dx,$$
 (5.10)

where F(h) is called the *Fourier transform* (FT) of f(x), and the units of the variable *h* are reciprocals of the units of *x*. For example, if *x* is time in seconds (s), then *h* is reciprocal time, or frequency, in reciprocal seconds (s<sup>-1</sup>). So if f(x) is a function of time, F(h) is a function of frequency. Taking the FT of time-dependent functions is a means of decomposing these functions into their component frequencies and is sometimes referred to as *Fourier analysis*. The FT in this form is used in infrared (IR) and nuclear magnetic resonance (NMR) spectroscopy to obtain the frequencies of many spectral lines simultaneously (as I will describe in Chapter 10 on obtaining models from NMR).

On the other hand, if x is a distance or length in Å, h is reciprocal length in Å<sup>-1</sup>. You can thus see that this highly general mathematical form is naturally adapted for relating real and reciprocal space. In fact, as I mentioned earlier, the Fourier transform is a precise mathematical description of diffraction. The diffraction patterns in Figs. 2.8–2.11 (pp. 15–18) are Fourier transforms of the corresponding simple objects and arrays. If these figures give you some intuition about how an object is related to its diffraction pattern, then they provide the same perception about the kinship between an object and its Fourier transform. According to Eq. (5.10), to compute F(h), the Fourier transform of f(x), just multiply the function by  $e^{2\pi i (hx)}$  and integrate (or better, let a computer integrate) the combined functions with respect to x. The result is a new function F(h), which is the FT of f(x). Computer routines for calculating FTs of functions are widely available, and form one of the vital internal organs of crystallographic software.

The Fourier transform operation is reversible. That is, the same mathematical operation that gives F(h) from f(x) can be carried out in the opposite direction, to give f(x) from F(h); specifically,

$$f(x) = \int_{-\infty}^{+\infty} F(h)e^{-2\pi i(hx)}dh$$
 (5.11)

In other words, if F(h) is the transform of f(x), then f(x) is in turn the transform of F(h). In this situation, f(x) is sometimes called the *back-transform* of (h), but this is a loose term that simply refers to the second successive transform that recreates the original function. Notice that the only difference between Eqs. (5.10) and (5.11) is the sign of the exponential term. You can think of this sign change as analogous (very *roughly* analogous!) to the sign change that makes subtraction the reverse of addition. Adding 3 to 5 gives 8: 5 + 3 = 8. To reverse the operation and generate the original 5, you subtract 3 from the previous result: 8 - 3 = 5.

If you think of 8 as a simple transform of 5 made by adding 3, the back-transform of 8 is 5, produced by subtracting 3.

Returning to the visual transforms of Figs. 2.8–2.11 (pp. 15–18), each object (the sphere in Fig. 2.8, for instance) is the Fourier transform (the back-transform, if you wish) of its diffraction pattern. If we build a model that looks like the diffraction pattern on the right, and then obtain its diffraction pattern, we get an image of the object on the left.

There is one added complication. The preceeding functions f(x) and F(h) are one-dimensional. Fortunately, the Fourier transform applies to periodic functions in any number of dimensions. To restate Fourier's conclusion in three dimensions, for any function f(x, y, z) there exists the function F(h, k, l) such that

$$F(h,k,l) = \int_{x} \int_{y} \int_{z} f(x,y,z) e^{2\pi i (hx+ky+lz)} dx dy dz.$$
(5.12)

As before, F(h, k, l) is called the Fourier transform of f(x, y, z), and in turn, f(x, y, z) is the Fourier transform of F(h, k, l) as follows:

$$f(x, y, z) = \int_{h} \int_{k} \int_{l} F(h, k, l) e^{-2\pi i (hx + ky + lz)} dh \, dk \, dl.$$
(5.13)

### 5.2.4 Fourier this and Fourier that: Review

I have used Fourier's name in discussing several types of equations and operations, and I want to be sure that I have not muddled them in your mind. First, a *Fourier* sum is a sum of simple wave equations or periodic functions that describes or approximates a complicated periodic function. Second, constructing a Fourier sum—that is, determining the proper F, h, and  $\alpha$  values to approximate a specific function—is called *Fourier synthesis*. For example, the sum of  $f_0$  through  $f_6$ in Fig. 2.16, p. 25 is at once a Fourier sum and the result of Fourier synthesis. Third, decomposing a complicated function into its components is called *Fourier* analysis. Fourth, the Fourier transform is an operation that transforms a function containing variables of one type (say time) into a function whose variables are reciprocals of the original type [in this case, 1/(time) or frequency]. The function f(x) is related to its Fourier transform F(h) by Eq. (5.10). The term *transform* is commonly used as a noun to refer to the function F(h) and also loosely as a verb to denote the operation of computing a Fourier transform. Finally, a *Fourier series* is an *infinite* sum based on some iterative formula for generating each term (the sum in Fig. 2.16, p. 25 is actually the beginning of an infinite series). So a sum of experimental terms (such as X-ray reflections), or a sum of individual atomic or electron-density contributions to a reflection is a Fourier sum, not a series. (Last, a grammar note: the word *series* is both singular and plural. You must gather from context whether a writer is talking about one series or many series. But from here on, I will be talking about sums, not series.)

### 5.3 Fourier mathematics and diffraction

### 5.3.1 Structure factor as a Fourier sum

I have stated that both structure factors and electron density can be expressed as Fourier sums. A structure factor describes one diffracted X-ray, which produces one reflection received at the detector. A structure factor  $F_{hkl}$  can be written as a Fourier sum in which each term gives the contribution of one atom to the reflection *hkl* [see Fig. 2.17, p. 26 and Eq. (2.3), p. 24]. Here is a single term, called an *atomic structure factor*,  $f_{hkl}$ , in such a series, representing the contribution of the single atom *j* to reflection *hkl*:

$$f_{hkl} = f_j e^{2\pi i (hx_j + ky_j + lz_j)}.$$
(5.14)

The term  $f_j$  is called the *scattering factor* of atom j, and it is a mathematical function (called a  $\delta$  function) that amounts to treating the atom as a simple sphere of electron density. The function is slightly different for each element, because each element has a different number of electrons (a different value of Z) to diffract the X-rays. The exponential term should be familiar to you by now. It represents a simple three-dimensional periodic function having both cosine and sine components. But the terms in parenthesis now possess added physical meaning:  $x_j$ ,  $y_j$ , and  $z_j$  are the coordinates of atom j in the unit cell (real space), expressed as fractions of the unit-cell axis lengths; and h, k, and l, in addition to their role as frequencies of a wave in the three directions x, y, and z, are also the indices of a specific reflection in the reciprocal lattice.

As mentioned earlier, the phase of a diffracted ray is implicit in the exponential formulation of a structure factor and depends only upon the atomic coordinates  $(x_j, y_j, z_j)$  of the atom. In fact, the phase for diffraction by one atom is  $2\pi (hx_j + ky_j + ly_j)$ , the exponent of *e* (ignoring the imaginary *i*) in the structure factor. For its contribution to the 220 reflection, an atom at (0, 1/2, 0) has phase  $2\pi (hx_j + ky_j + lz_j)$  or  $2\pi (2[0] + 2[1/2] + 0[0]) = 2\pi$ , which is the same as a phase of zero. This atom lies on the (220) plane, and all atoms lying on (220) planes contribute to the 220 reflection with phase of zero. [Try the above calculation for another atom at (1/2, 0, 0), which is also on a (220) plane.] This conclusion is in keeping with Bragg's law, which says that all atoms on a set of equivalent, parallel lattice planes diffract in phase with each other.

Each diffracted ray is a complicated wave, the sum of diffractive contributions from *all* atoms in the unit cell. For a unit cell containing *n* atoms, the structure factor  $F_{hkl}$  is the sum of all the atomic  $f_{hkl}$  values for individual atoms. Thus, in parallel with Eq. (2.3), p. 24, we write the structure factor for reflection  $F_{hkl}$  as follows:

$$F_{hkl} = \sum_{j=1}^{n} f_j e^{2\pi i (hx_j + ky_j + lz_j)}.$$
(5.15)

#### Section 5.3 Fourier mathematics and diffraction

In words, the structure factor that describes reflection hkl is a Fourier sum in which each term is the contribution of one atom, treated as a simple sphere of electron density. So the contribution of each atom *j* to  $F_{hkl}$  depends on (1) what element it is, which determines  $f_j$ , the amplitude of the contribution, and (2) its position in the unit cell  $(x_j, y_j, z_j)$ , which establishes the phase of its contribution.

Alternatively,  $F_{hkl}$  can be written as the sum of contributions from each volume element of electron density in the unit cell [see Fig. 2.18, p. 27 and Eq. (2.4), p. 28]. The electron density of a volume element centered at (x, y, z) is, roughly, the average value of  $\rho(x, y, z)$  in that region. The smaller we make our volume elements, the more precisely these averages approach the correct values of  $\rho(x, y, z)$  at all points. We can, in effect, make our volume elements infinitesimally small, and the average values of  $\rho(x, y, z)$  precisely equal to the actual values at every point, by integrating the function  $\rho(x, y, z)$  rather than summing average values. Think of the resulting integral as the sum of the contributions of an infinite number of vanishingly small volume elements. Written this way,

$$F_{hkl} = \int_{x} \int_{y} \int_{z} \rho(x, y, z) e^{2\pi i (hx + ky + lz)} \, dx \, dy \, dz, \tag{5.16}$$

or equivalently,

$$F_{hkl} = \int_{V} \rho(x, y, z) e^{2\pi i (hx + ky + lz)} \, dV, \qquad (5.17)$$

where the integral over V, the unit-cell volume, is just shorthand for the integral over all values of x, y, and z in the unit cell. Each volume element contributes to  $F_{hkl}$  with a phase determined by its coordinates (x, y, z), just as the phase of atomic contributions depend on atomic coordinates.

You can see by comparing Eq. (5.17) with Eq. (5.10) [or Eq. (5.16) with Eq. (5.12)] that  $F_{hkl}$  is the Fourier transform of  $\rho(x, y, z)$ . More precisely,  $F_{hkl}$  is the transform of  $\rho(x, y, z)$  on the set of real-lattice planes (*hkl*). All of the  $F_{hkl}$ s together compose the transform of  $\rho(x, y, z)$  on all sets of equivalent, parallel planes throughout the unit cell.

### 5.3.2 Electron density as a Fourier sum

Because the Fourier transform operation is reversible [Eqs. (5.10) and (5.11)], the electron density is in turn the transform of the structure factors, as follows:

$$\rho(x, y, z) = \frac{1}{V} \sum_{h} \sum_{k} \sum_{l} F_{hkl} e^{-2\pi i (hx + ky + lz)},$$
(5.18)

where V is the volume of the unit cell.

This transform is a triple sum rather than a triple integral because the  $F_{hkl}$ s represent a set of discrete entities: the reflections of the diffraction pattern. The transform

of a discrete function, such as the reciprocal lattice of measured intensities, is a summation of discrete values of the function. The transform of a continuous function, such as  $\rho(x, y, z)$ , is an integral, which you can think of as a sum also, but a sum of an infinite number of infinitesimals. Superficially, except for the sign change (in the exponential term) that accompanies the transform operation, this equation appears identical to Eq. (5.9), a general three-dimensional Fourier sum. But here, each  $F_{hkl}$  is not just one of many simple numerical amplitudes for a standard set of component waves in a Fourier sum. Instead, each  $F_{hkl}$  is a structure factor, itself a Fourier sum, describing a specific reflection in the diffraction pattern.

"Curiouser and curiouser," said Alice.

### 5.3.3 Computing electron density from data

Equation (5.18) tells us, at last, how to obtain  $\rho(x, y, z)$ . We need merely to construct a Fourier sum from the structure factors. By now, you might be wondering about the practical aspects of calculating an electron-density map from this rather abstract equation. Widely available software can turn archived structure factors into electron-density maps for viewing with graphics programs, making possible a first-hand look at the most critical evidence in support of a model. But if you think about the calculation itself, you might wonder about those imaginary terms and what they describe, or what happens to them in the calculations. What physical meaning could be ascribed to imaginary terms in a representation of electron density? I told you earlier that the exponential description  $(e^{i(hx)})$  is computationally more efficient than its trigonometric equivalent  $[\cos(hx) + i\sin(hx)]$ , but when it comes to calculating a concrete image, what happens to those *i* terms?

In short, they go away. The Fourier sum is taken over all values of indices h, k, and l, and for every term containing a positive index h, there is a term containing the same negative index -h. Under these circumstances, sines and cosines behave very differently, because  $\cos(hx) = \cos(-hx)$ , but  $\sin(hx) = -\sin(-hx)$  [try it on your hand calculator, with thirty degrees for (hx)]. So when you use Eq. (5.18) or similar equations to compute anything from structure factors, the cosine terms for positive and negative indices add together, but the sine terms all cancel out precisely. So the actual computation of an electron density map makes use only of the cosine terms, and the imaginaries go away. But that does not mean that the sine terms serve no purpose. As I have said before, in the structure factor equation, the  $[\cos(hx) + i\sin(hx)]$  terms implicitly define the phases of reflections.

A structure factor describes a diffracted ray, and a full description of a diffracted ray, like any description of a wave, must include three parameters: amplitude, frequency, and phase. In discussing data collection, however, I have mentioned only two measurements: the indices of each reflection and its intensity. Looking again at Eq. (5.18), you see that the indices of a reflection play the role of the three frequencies in one Fourier term. The only measurable variable remaining in the equation is  $F_{hkl}$ . Does the measured intensity of a reflection, the only measurement we can make in addition to the indices, completely define  $F_{hkl}$ ? Unfortunately, the answer is "no."

### 5.3.4 The phase problem

Because  $F_{hkl}$  is a periodic function, it possesses amplitude, frequency, and phase. The amplitude of  $F_{hkl}$  is proportional to the square root of the reflection intensity Ihkl, so structure-factor amplitudes are directly obtainable from measured reflection intensities. The three frequencies of this three-dimensional wave function are h, k, and l, the indices of the planes that produce the reflection described by  $F_{hkl}$ . So the frequency of a structure factor is equal to  $1/d_{hkl}$ , making the wavelength the same as the spacing of the planes producing the reflection. But the phase of  $F_{hkl}$  is not directly obtainable from a single measurement of the reflection intensity. In order to compute  $\rho(x, y, z)$  from the structure factors, we must obtain, in addition to the intensity of each reflection, the phase of each diffracted ray. In Chapter 6, I will present an expression for  $\rho(x, y, z)$  as a Fourier series in which the phases are explicit (finally, huh?), and I will discuss means of obtaining phases. This is one of the most difficult problems in crystallography. For now, on the assumption that the phases can be obtained, and thus that complete structure factors are obtainable, I will consider further the implications of Eq. (5.15) (structure factors F expressed in terms of atoms), Eq. (5.16) [structure factors in terms of  $\rho(x, y, z)$ ], and Eq. (5.18) [ $\rho(x, y, z)$  in terms of structure factors].

### 5.4 Meaning of the Fourier equations

### 5.4.1 Reflections as terms in a Fourier sum: Eq. (5.18)

First consider Eq. (5.18) ( $\rho$  in terms of *F*s). Each term in this Fourier-series description of  $\rho(x, y, z)$  is a structure factor representing a single X-ray reflection. The indices *hkl* of the reflection give the three frequencies necessary to describe the Fourier term as a simple wave in three dimensions. Recall from Sec. 2.6.2, p. 23 that any periodic function can be approximated by a Fourier sum, and that the approximation improves as more terms are added to the sum (see Fig. 2.16, p. 25). The low-frequency terms in Eq. (5.18) determine gross features of the periodic function  $\rho(x, y, z)$ , whereas the high-frequency terms improve the approximation by filling in fine details. You can also see in Eq. (5.18) that the low-frequency terms in the Fourier sum that describes our desired function  $\rho(x, y, z)$  are given by reflections with low indices, that is, by reflections near the center of the diffraction pattern (Fig. 5.2). In some crystallographic circles (so to speak), they are called *low-angle* reflections.

The high-frequency terms are given by reflections with high indices, reflections farthest from the center of the pattern (*high-angle reflections*). Thus you can see the importance of how well a crystal diffracts. If a crystal does not produce diffracted rays at large angles from the direct beam (reflections with large indices), the Fourier sum constructed from all the measurable reflections lacks high-frequency terms, and the resulting transform is not highly detailed—the resolution of the resulting



**Figure 5.2** Structure factors of reflections near the center of the diffraction pattern are low-frequency terms in the Fourier sum that approximates  $\rho(x, y, z)$ . Structure factors of reflections near the outer edge of the pattern are high-frequency terms.

image is poor. The Fourier series of Fig. 2.16, p. 25 is *truncated* in just this manner and does not fit the target function in fine details like the sharp corners.

This might sound like a Zen Buddhist question: *What is the meaning of one reflection*? What does the FT of one structure factor tell you? In Fig. 5.3, I have again used a page from Kevin Cowtan's *Book of Fourier* to demonstrate what individual structure factors contribute to the molecular image we are seeking. The first column (*a*) shows the full diffraction pattern (top panel, with phases indicated by color) of a simple molecular model in its unit cell (*a*, bottom panel). The second column (*b*) shows contour plots of the Fourier transforms of individual structure factors, with reflection indices listed to the left of each transform. The FT of the 01 reflection is simply a sinusoidal electron-density function showing electron density on the 01 planes. The density exhibits peaks (red) on the planes and reaches minima (blue) between the planes. The FT of *any* single reflection, like the 10 reflection shown next, will look something like this: merely a repetitive rise and fall of electron density. After all, the only thing you learn from a single reflection is the average electron-density along the set of Miller planes that share their indices with that specific reflection.

But look what happens when you add structure-factor FTs to each other (column c). When the FTs of the 01 and 10 reflections are added together, we see interference between them, in the form of high positive density (intense red, where red crosses red), zero density (white, where red crosses blue), and high negative density (blue crossing blue). The result is a large lobe of positive density that shows us approximately where the molecule lies in this unit cell. In this simple example, only two reflections, when transformed, roughly locate the molecule!



**Figure 5.3**  $\blacktriangleright$  (a) Structure-factor pattern (top) calculated from a simple model (bottom). (b) Fourier transforms of individual reflections from a. Red is positive electron density, blue is negative. (c) Sums of FTs from b. In each square, the FT of one reflection is added to the sum above it.

As we continue adding FTs of structure factors for more reflections, the molecular image becomes more sharply defined, showing what each reflection contributes to the final image. Notice that FTs of reflections with higher indices have higher frequencies (more red-blue repeats per unit distance). These FTs add fine details to the sum. After we have summed only seven reflections, we have located all the atoms approximately. Interference among the full set of FTs gives a sharp image of all atoms (*a*, lower panel), with only very weak negative density, which is due to the finiteness of our data set. Remember that full data sets for biological macromolecules comprise from thousands to millions of structure factors, each telling us the average electron density on a specific set of parallel planes.

There is even more than meets the eye in Fig. 5.3, p. 103. We will visit it again in Chapter 6, when I take up the phase problem. Can't wait? See p. 111.

### 5.4.2 Computing structure factors from a model: Eq. (5.15) and Eq. (5.16)

Equation (5.15) describes one structure factor in terms of diffractive contributions from all *atoms* in the unit cell. Equation (5.16) describes one structure factor in terms of diffractive contributions from all *volume elements* of electron density in the unit cell. Notice in both cases that *all* parts of the structure—every atom or every scrap of electron density—contribute to *every* structure factor. These two equations suggest that we can calculate the full set of structure factors either from an atomic model of the protein or from an electron density function. In short, if we know the structure, either as atoms or as electron density, we can calculate the diffraction pattern, *including the phases of all reflections*. This computation, of course, appears to go in just the opposite direction that the crystallographer desires. It turns out, however, that computing structure factors from a model of the unit cell (back-transforming the model) is an essential part of crystallography, for several reasons.

First, this computation is used in obtaining phases. As I will discuss in Chapter 6, the crystallographer obtains phases by starting from rough estimates of them and then undertaking an iterative process to improve the estimates. This iteration entails a cycle of three steps. In step 1, an estimated  $\rho(x, y, z)$  (that is, a crude model of the structure) is computed using Eq. (5.18) with observed intensities  $(I_{obs})$ and estimated phases ( $\alpha_{calc}$ ). In step 2, the crystallographer attempts to improve the model by viewing the electron-density map [a computer plot of  $\rho(x, y, z)$ ] and identifying molecular features such as the molecule-solvent boundaries or specific groups of atoms (called interpreting the map). Step 3 entails computing new structure factors ( $F_{\text{calc}}$ ), using either Eq. (5.16) with the improved  $\rho(x, y, z)$  model from step 2, or Eq. (5.15) with a partial atomic model of the molecule, containing only those atoms that can be located with some confidence in the electron-density map. Calculation of new  $F_{\text{calc}}$ s in step 3 produces a new (better, we hope) set of estimated phases, and the cycle is repeated: a new  $\rho(x, y, z)$  is computed from the original measured intensities and the newest phases, interpretation produces a more detailed model, and calculation of structure factors from this model produces improved phases. In each cycle, the crystallographer hopes to obtain an improved

#### Section 5.4 Meaning of the Fourier equations

 $\rho(x, y, z)$ , which means a more detailed and interpretable electron-density map, and thus a more complete and accurate model of the desired structure. I will discuss the iterative improvement of phases and electron-density maps in Chapter 7. For now just take note that obtaining the final structure entails both calculating  $\rho(x, y, z)$  from structure factors and calculating structure factors from some preliminary model, either a rough form of  $\rho(x, y, z)$  or a partial atomic model. Note further that when we compute structure factors from a known or assumed model, *the results include the phases for that model*. In other words, the computed results give us the information needed for a "full-color" diffraction pattern, like that shown in Fig. 2.19*d*, p. 29, whereas experimentally obtained diffraction patterns lack the phases and are merely black and white, like Fig. 2.19*e*.

The second use of back-transforms is to assess the progress of structure determination. Equations (5.15) and (5.16) provide means to monitor the iterative process to see whether it is converging toward improved phases and improved  $\rho(x, y, z)$ . The computed structure factors  $F_{calc}$  include both the desired phases  $\alpha_{calc}$  and a new set of intensities. I will refer to these *calculated* intensities as  $I_{calc}$  to distinguish them from the *measured* reflection intensities  $I_{obs}$  taken from the diffraction pattern. As the iteration proceeds, the values of  $I_{calc}$  should approach those of  $I_{obs}$ . So the crystallographer compares the  $I_{calc}$  and  $I_{obs}$  values at each cycle in order to see whether the iteration is converging. When cycles of computation provide no further improvement in correspondence between calculated and measured intensities, then the process is complete, and the model can be improved no further.

### 5.4.3 Systematic absences in the diffraction pattern: Eq. (5.15)

A third application of Eq. (5.15) allows us to understand how systematic absences in the diffraction pattern reveal symmetry elements in the unit cell, thus guiding the crystallographer in assigning the space group of the crystal. Recall from Sec. 4.2.8, p. 65, that if the unit cell possesses symmetry elements, then certain sets of reciprocal-lattice points are equivalent, and so certain reflections in the diffraction pattern are redundant. The crystallographer must determine the unitcell space group (i.e., determine what symmetry elements are present) in order to devise an efficient strategy for measuring as many unique reflections as efficiently as possible. I stated without justification in Chapter 4 that certain symmetry elements announce themselves in the diffraction pattern as *systematic absences*: regular patterns of missing reflections. Now I will use Eq. (5.15) to show how a symmetry element in the unit cell produces systematic absences in the diffraction pattern.

For example, as indicated by the "Reflection conditions" in Fig. 4.20, p. 72, if the **b**-axis of the unit cell is a twofold screw axis, then reflections 010, 030, 050, along with all other 0k0 reflections in which k is an odd number, are missing. We can see why by using the concept of equivalent positions (Sec. 4.2.8, p. 65). For a unit cell with a twofold screw axis along edge **b**, the equivalent positions are (x, y, z) and (-x, y + 1/2, -z). That is, for every atom j with coordinates (x, y, z) in the unit cell, there is an identical atom j' at (-x, y + 1/2, -z). Atoms j and j' are called *symmetry-related* atoms. According to Eq. (5.15), the structure

factor for reflections  $F_{0k0}$  is

$$F_{0k0} = \sum_{j} f_j e^{2\pi i (ky_j)}.$$
(5.19)

The exponential term is greatly simplified in comparison to that in Eq. (5.15) because h = l = 0 for reflections on the 0k0 axis. Now I will separate the contributions of atoms *j* from that of their symmetry-related atoms *j*':

$$F_{0k0} = \sum_{j} f_{j} e^{2\pi i (ky_{j})} + \sum_{j'} f_{j'} e^{2\pi i (ky_{j'})}.$$
(5.20)

Because atoms j and j' are identical, they have the same scattering factor f, and so I can substitute  $f_j$  for  $f_{j'}$  and factor out the f terms:

$$F_{0k0} = \sum_{j} f_j \left( \sum_{j} e^{2\pi i k y_j} + \sum_{j'} e^{2\pi i k y_{j'}} \right).$$
(5.21)

If the *y* coordinate of atom *j* is *y*, then the *y* coordinate of atom j' is y + 1/2. Making these substitutions for  $z_j$  and  $z_{j'}$ ,

$$F_{0k0} = \sum_{j} f_j \left( \sum_{j} [e^{2\pi i k y} + e^{2\pi i k (y+1/2)}] \right).$$
(5.22)

The  $f_j$  terms are nonzero, so  $F_{0k0}$  is zero, and the corresponding 010 reflection is missing, only if all the summed terms in square brackets equal zero. Simplifying one of these terms,

$$e^{2\pi i k y} + e^{2\pi i k (y+1/2)} = e^{2\pi i k y} (1 + e^{\pi i k}).$$
(5.23)

This term is zero, and hence  $F_{0k0}$  is zero, if  $e^{\pi i k}$  equals -1. Converting this exponential to its trigonometric form [see Eq. (5.6)],

$$e^{\pi i k} = \cos(\pi k) + i \sin(\pi k).$$
 (5.24)

The cosine of  $\pi$  radians (180°), or any odd multiple of  $\pi$  radians, is -1. The sine of  $\pi$  radians is 0. Thus  $e^{\pi i k}$  equals -1 for all odd values of k, and  $F_{0k0}$  equals zero if l is odd.

The preceding shows that  $F_{0k0}$  disappears for odd values of k when the **b** edge of a unit cell is a twofold screw axis. But what is going on physically? In short, the diffracted rays from two atoms at (x, y, z) and (-x, y + 1/2, -z) are identical in amplitude  $(f_j = f_{j'})$  but precisely opposite in phase. Recall that the phase of an

106

#### Section 5.5 Summary: From data to density

atom's contribution to  $F_{hkl}$  is  $2\pi (hx + ky + lz)$ . Consider an atom *j* lying at the origin of the unit cell (0, 0, 0), and its symmetry-related atom *j'* at (0, 1/2, 0). The phase of *j*'s contribution to  $F_{010}$  is  $2\pi ([0 \cdot 0] + [1 \cdot 0] + [0 \cdot 0)] = 0$  radians or 0°. The phase for atom *j'* is  $2\pi ([0 \cdot 0] + [1 \cdot (1/2)] + [0 \cdot 0)] = \pi$  radians or 180°, which is precisely 180° out of phase with *j*'s contribution, thus cancelling it to make  $F_{010}$  a missing reflection. So the symmetry-related pair of atoms contributes nothing to  $F_{0k0}$  when *k* is odd, and because every atom has such a symmetry-related partner, this cancellation occurs no matter where the atoms lie. Putting it another way, if the unit cell contains a twofold screw axis along edge **b**, then every atom in the unit cell is paired with a symmetry-related atom that cancels its contributions to all odd-numbered 0*k*0 reflections. (Can you show that two atoms related by a twofold screw axis along **b** diffract *in phase* for the 020 reflection?)

Similar computations have been carried out for all symmetry elements and combinations of elements. Like equivalent positions, systematic absences are tabulated for all space groups in *International Tables*, so the crystallographer can use this reference as an aid to space-group determination. As mentioned above, the *International Tables* entry for space group  $P2_1$  (Fig. 4.20, p. 72), which possess a  $2_1$ axis on edge **b**, indicates that, for reflections 0k0, the "Reflection conditions" are 0k0 : k = 2n. In other words, in this space group, reflections 0k0 are present only if k is even (2 times any integer n), so they are absent if k is odd, as I proved above.

### 5.5 Summary: From data to density

When we describe structure factors and electron density as Fourier sums, we find that they are intimately related. The electron density is the Fourier transform of the structure factors, which means that we can convert the crystallographic data into an image of the unit cell and its contents. One necessary piece of information is, however, missing for each structure factor. We can measure only the intensity  $I_{hkl}$  of each reflection, not the complete structure factor  $F_{hkl}$ . What is the relationship between them? It can be shown that the amplitude of structure factor  $F_{hkl}$  is proportional to  $(I_{hkl})^{1/2}$ , the square root of the measured intensity. So if we know  $I_{hkl}$  from diffraction data, we know the amplitude of  $F_{hkl}$ . Unfortunately, we do not know its phase  $\alpha_{hkl}$ . In focusing light reflected from an object, a lens maintains all phase relationships among the rays, and thus constructs an image accurately. When we record diffraction intensities, we lose the phase information that the computer needs in order to simulate an X-ray-focusing lens. In Chapter 6, I will describe methods for learning the phase of each reflection, and thus obtaining the complete structure factors needed to calculate the electron density.

This Page Intentionally Left Blank

### ► Chapter 6

## **Obtaining Phases**

### 6.1 Introduction

The molecular image that the crystallographer seeks is a contour map of the electron density  $\rho(x, y, z)$  throughout the unit cell. The electron density, like all periodic functions, can be represented by a Fourier sum. The representation that connects  $\rho(x, y, z)$  to the diffraction pattern is

$$\rho(x, y, z) = \frac{1}{V} \sum_{h} \sum_{k} \sum_{l} F_{hkl} e^{-2\pi i (hx + ky + lz)}.$$
(5.18)

Equation (5.18) tells us how to calculate  $\rho(x, y, z)$ : simply construct a Fourier sum using the structure factors  $F_{hkl}$ . For each term in the sum, h, k, and l are the indices of reflection hkl, and  $F_{hkl}$  is the structure factor that describes the reflection. Each structure factor  $F_{hkl}$  is a complete description of a diffracted ray recorded as reflection hkl. Being a wave equation,  $F_{hkl}$  must specify frequency, amplitude, and phase. Its three frequency terms h, k, and l are the indices of the set of parallel planes that produce the reflection. Its amplitude is proportional to  $(I_{hkl})^{1/2}$ , the square root of the measured intensity  $I_{hkl}$  of reflection hkl. Its phase is unknown and is the only additional information the crystallographer needs in order to compute  $\rho(x, y, z)$  and thus obtain an image of the unit cell contents. In this chapter, I will discuss some of the common methods of obtaining phases.

Let me emphasize that *each reflection has its own phase* (see Fig. 2.19*d*, p. 29), so the phase problems must be solved for *every one* of the thousands of reflections used to construct the Fourier sum that approximates  $\rho(x, y, z)$ . Let me also emphasize how crucial this phase information is. In his *Book of Fourier*, Kevin Cowtan illustrates the relative importance of phases and intensities in solving a structure, as shown in Fig. 6.1. Images (*a*) and (*b*) show two simple models, a duck



**Figure 6.1** Relative amounts of information contained in reflection intensities and phases. (*a*) and (*b*) Duck and cat, along with their Fourier transforms. (*c*) Intensity (shading) of the duck transform, combined with the phases (colors) of the cat transform. (*d*) Back-transform of (*c*) produces recognizable image of cat, but not duck. Phases contain more information than intensities. Figure generously provided by Dr. Kevin Cowtan.

#### Section 6.1 Introduction

and a cat, along with their calculated Fourier transforms. As in Fig. 2.19, p. 29, phases are shown as colors, while the intensity of color reflects the amplitude of the structure factor at each location. (Note that these are continuous transforms, because the model is not in a lattice.) Back-transforming each Fourier transform would produce an image of the duck or cat. In (c) the colors from the cat transform are superimposed on the intensities from the duck transform. This gives us a transform in which the intensities come from the duck and the phases come from the cat. In (d) we see the back-transform of (c). The image of the cat is obvious, but you cannot find any sign of the duck. Ironically, the diffraction intensities, which are relatively easy to measure, contain far less information than do the phases, which are much more difficult to obtain.

Before I begin showing you how to obtain phases, I want to give you one additional bit of feeling for their physical meaning, and as well, to confess to a small half-truth I have quietly sustained until now, for simplicity. Take another look at Fig. 5.3, p. 103. If you look very closely at the FT of individual reflections, you will see that the peak of repetitive electron density does not pass precisely through the origin of the unit cell. (The 10 reflection is a good one to check, and you may have to compare its FT with the lower panel of Fig. 5.3a in order to see the unit cell clearly.) This observation implies that the electron density producing, say, the 10 reflection does not peak precisely on the 10 plane. While we commonly say that the 234 reflection comes from the planes whose Miller indices are 234, the truth is that only if its phase is zero does the reflection come precisely from those planes. If its phase is  $\pi$  (180°), then the reflection is coming from halfway between those planes. Physically, the electron density that repeats with the orientation and frequency of the 234 planes may not have its peak on the planes; it is more likely that the peak lies somewhere between them. So to be most accurate, we should say that the 234 reflection comes from repetitive electron density whose orientation and frequency corresponds to that of the 234 planes, but whose peak may lie anywhere between those planes. Its exact position constitutes its phase. Whatever its phase, if that repetitive electron density is high, the 234 reflection will be strong; if low, it will be weak or missing. If you look again at the basic Bragg model (Fig. 4.8, p. 56), you will see that the actual position of the planes is immaterial. Any set of planes with the same orientation and spacing would produce the reflection shown. (It appears that no one ever bothers to add this last little bit of precision to the Bragg model, but there it is. Now I feel better.)

So let us get to phasing. In order to illuminate both the phase problem and its solution, I will represent structure factors as vectors on a two-dimensional plane of complex numbers of the form a + ib, where *i* is the imaginary number  $(-1)^{1/2}$ . This allows me to show geometrically how to compute phases. I will begin by introducing complex numbers and their representation as points having coordinates (a, b) on the complex plane. Then I will show how to represent structure factors as vectors on the same plane. Because we will now start thinking of the structure factor as a vector, I will hereafter write it in boldface ( $\mathbf{F}_{hkl}$ ) instead of the italics used for simple variables and functions. Finally, I will use the vector representation of structure factors to explain a few common methods of obtaining phases.

# 6.2 Two-dimensional representation of structure factors

### 6.2.1 Complex numbers in two dimensions

Structure-factor equations like Eq. (5.15), p. 98, present the structure factor as a sum of terms each containing the exponential element  $e^{2\pi i(hx+ky+lz)}$ . Remember that these exponential elements can also be expressed trigonometrically as  $[\cos 2\pi (hx + ky + lz) + i \sin 2\pi (hx + ky + lz)]$ . In this form, each term in the Fourier sum, and hence the structure factor itself, is a complex number of the form a + ib. Complex numbers can be represented as points in two dimensions (Fig. 6.2). I will use this representation to help you understand the nature of the phase problem and various ways of solving it. The horizontal axis in the figure represents the real-number line. Any real number *a* is a point on this line, which stretches from  $a = -\infty$  to  $a = +\infty$ . The vertical axis is the imaginary-number line, on which lie all imaginary numbers *ib* between  $b = -i\infty$  and  $b = +i\infty$ . A complex number a + ib, which possesses both real (*a*) and imaginary (*ib*) parts, is thus a point at position (*a*, *ib*) on this plane.

### 6.2.2 Structure factors as complex vectors

A representation of structure factors on this plane must include the two properties we need in order to construct  $\rho(x, y, z)$ : amplitude and phase. Crystallographers represent each structure factor as a *complex vector*, that is, a vector (not a point) on the plane of complex numbers. The length of this vector represents the amplitude of the structure factor. Thus the length of the vector representing structure factor  $\mathbf{F}_{hkl}$  is proportional to  $(I_{hkl})^{1/2}$ . The second property, phase, is represented by the



**Figure 6.2**  $\triangleright$  The complex number N = a + ib, represented as a point on the plane of complex numbers.



**Figure 6.3** (a) The structure factor **F**, represented as a vector on the plane of complex numbers. The length of **F** is proportional to  $I^{1/2}$ , the square root of the measured intensity *I*. The angle between **F** and the positive real axis is the phase  $\alpha$ . (b) (Stereo) **F** can be pictured as a complex vector spinning around its line of travel. The projection of the path taken by the head of the vector is the familiar sine wave. The spinning vector reaches the detector pointing in a specific direction that corresponds to its phase.

angle  $\alpha$  that the vector makes with the positive real-number axis when the origin of the vector is placed at the origin of the complex plane, which is the point 0 + i0 (Fig. 6.3*a*).

We can represent a structure factor **F** as a vector  $\mathbf{A} + i\mathbf{B}$  on this plane. The projection of **F** on the real axis is its real part **A**, a vector of length  $|\mathbf{A}|$  (absolute value of **A**) on the real-number line; and the projection of **F** on the imaginary axis is its imaginary part  $i\mathbf{B}$ , a vector of length  $|\mathbf{B}|$  on the imaginary-number line. The length or magnitude (or in wave terminology, the amplitude) of a complex vector is analogous to the absolute value of a real number, so the length of vector  $\mathbf{F}_{hkl}$  is  $|\mathbf{F}_{hkl}|$ ; therefore,  $|\mathbf{F}_{hkl}|$  is proportional to  $(I_{hkl})^{1/2}$ , and if the intensity is known from data collection, we can treat  $|\mathbf{F}_{hkl}|$  as a known quantity. The angle that  $\mathbf{F}_{hkl}$  makes with the real axis is represented in radians as  $\alpha$  ( $0 \le \alpha \le 2\pi$ ), or in cycles as  $\alpha'(0 \le \alpha' \le 1)$ , and is referred to as the *phase angle*.

#### Chapter 6 Obtaining Phases

This representation of a structure factor is equivalent to thinking of a wave as a complex vector spinning around its axis as it travels thorough space (Fig. 6.3b). If its line of travel is perpendicular to the tail of the vector, then a projection of the head of the vector along the line of travel is the familiar sinusoidal wave. When the wave strikes the detector, the vector is pointing in a specific direction that corresponds to its phase. The phase of a structure factor tells us the direction of the vector at some arbitrary origin (in this case, the plane of the detector), and to know the phase of all reflections means to know all their *individual* phase angles with respect to a common origin.

In Sec. 4.3.7, p. 88, I mentioned Friedel's law, that  $\mathbf{I}_{hkl} = \mathbf{I}_{\overline{hkl}}$ . It will be helpful for later discussions to look at the vector representations of pairs of structure factors  $\mathbf{F}_{hkl}$  and  $\mathbf{F}_{\overline{hkl}}$ , which are called *Friedel pairs*. Even though  $\mathbf{I}_{hkl}$  and  $\mathbf{I}_{\overline{hkl}}$  are equal,  $\mathbf{F}_{hkl}$  and  $\mathbf{F}_{\overline{hkl}}$  are not. The structure factors of Friedel pairs have different phases, as shown in Fig. 6.4. Specifically, if the phase of  $\mathbf{F}_{hkl}$  is  $\alpha$ , then the phase of  $\mathbf{F}_{\overline{hkl}}$  is 360° –  $\alpha$ . Another way to put it is that Friedel pairs are reflections of each other in the real axis;  $\mathbf{F}_{\overline{hkl}}$  is the mirror image of  $\mathbf{F}_{hkl}$  with the real axis serving as the mirror.

The representation of structure factors as vectors in the complex plane (that is, complex vectors) is useful in several ways. Because the diffractive contributions of atoms or volume elements to a single reflection are additive, each contribution can be represented as a complex vector, and the resulting structure factor is the vector sum of all contributions. For example, in Fig. 6.5, **F** (green) represents a structure factor of a three-atom structure, in which  $\mathbf{f}_1$ ,  $\mathbf{f}_2$ , and  $\mathbf{f}_3$  (black) are the atomic structure factors. The length of each atomic structure factor  $\mathbf{f}_n$  represents its amplitude, and its angle with the real axis,  $\alpha_n$ , represents its phase. The vector sum  $\mathbf{F} = \mathbf{f}_1 + \mathbf{f}_2 + \mathbf{f}_3$  is obtained by placing the tail of  $\mathbf{f}_1$  at the origin, the tail of



**Figure 6.4** Structure factors of a Friedel pair.  $\mathbf{F}_{hkl}$  is the reflection of  $\mathbf{F}_{hkl}$  in the real axis.



**Figure 6.5**  $\blacktriangleright$  Molecular structure factor **F** (green) is the vector sum of three atomic structure factors (black). Vector addition of  $\mathbf{f}_1$ ,  $\mathbf{f}_2$ , and  $\mathbf{f}_3$  gives the amplitude and phase of **F**.

 $f_2$  on the head of  $f_1$ , and the tail of  $f_3$  on the head of  $f_2$ , all the while maintaining the phase angle of each vector. The structure factor **F** is thus a vector with its tail at the origin and its head on the head of  $f_3$ . This process sums both amplitudes and phases, so the resultant length of **F** represents its amplitude, and the resultant angle  $\alpha$  is its phase angle. (The atomic vectors may be added in any order with the same result.)

In subsequent sections of this chapter, I will use this simple vector arithmetic to show how to compute phases from various kinds of data. In the next section, I will use complex vectors to derive an equation for electron density as a function of reflection intensities and, at last, phases.

### 6.2.3 Electron density as a function of intensities and phases

Figure 6.3 shows how to decompose  $\mathbf{F}_{hkl}$  into its amplitude  $|\mathbf{F}_{hkl}|$ , which is the length of the vector, and its phase  $\alpha_{hkl}$ , which is the angle the vector makes with the real number line. This allows us to express  $\rho(x, y, z)$  as a function of the measurable amplitude of  $\mathbf{F}$  (measurable because it can be computed from the reflection intensity *I*) and the unknown phase  $\alpha$ . For clarity, I will at times drop the subscripts on  $\mathbf{F}$ , *I*, and  $\alpha$ , but remember that these relationships hold for all reflections.

In Fig. 6.3,

$$\cos \alpha = \frac{|\mathbf{A}|}{|\mathbf{F}|}$$
 and  $\sin \alpha = \frac{|\mathbf{B}|}{|\mathbf{F}|}$ , (6.1)

Chapter 6 Obtaining Phases

and therefore

$$|\mathbf{A}| = |\mathbf{F}| \cdot \cos \alpha$$
 and  $|\mathbf{B}| = |\mathbf{F}| \cdot \sin \alpha$ . (6.2)

Expressing **F** as a complex vector  $\mathbf{A} + i\mathbf{B}$ ,

$$\mathbf{F} = |\mathbf{A}| + i |\mathbf{B}| = |\mathbf{F}| \cdot (\cos \alpha + i \sin \alpha).$$
(6.3)

Expressing the complex term in the parentheses as an exponential [Eq. (5.6)],

$$\mathbf{F} = |\mathbf{F}| \cdot e^{i\alpha}.\tag{6.4}$$

Substituting this expression for  $F_{hkl}$  in Eq. (5.18), the electron-density equation (remembering that  $\alpha$  is the phase  $\alpha_{hkl}$  of the specific reflection hkl), gives

$$\rho(x, y, z) = \frac{1}{V} \sum_{h} \sum_{k} \sum_{l} |F_{hkl}| e^{i\alpha_{hkl}} e^{-2\pi i (hx + ky + lz)}.$$
 (6.5)

We can combine the exponential terms more simply by expressing the phase angle as  $\alpha'$ , using  $\alpha = 2\pi \alpha'$ :

$$\rho(x, y, z) = \frac{1}{V} \sum_{h} \sum_{k} \sum_{l} |F_{hkl}| e^{2\pi i \alpha'_{hkl}} e^{-2\pi i (hx + ky + lz)}.$$
 (6.6)

This substitution allows us to combine the exponentials by adding their exponents:

$$\rho(x, y, z) = \frac{1}{V} \sum_{h} \sum_{k} \sum_{l} |F_{hkl}| e^{-2\pi i (hx + ky + lz - \alpha'_{hkl})}.$$
(6.7)

This equation gives the desired electron density as a function of the known amplitudes  $|\mathbf{F}_{hkl}|$  and the unknown phases  $\alpha'_{hkl}$  of each reflection. Recall that this equation represents  $\rho(x, y, z)$  in a now-familiar form, as a Fourier sum, but this time with the phase of each structure factor expressed explicitly. Each term in the series is a three-dimensional wave of amplitude  $|\mathbf{F}_{hkl}|$ , phase  $\alpha'_{hkl}$ , and frequencies *h* along the *x*-axis, *k* along the *y*-axis, and *l* along the *z*-axis.

The most demanding element of macromolecular crystallography (except, perhaps, for dealing with macromolecules that resist crystallization) is the so-called *phase problem*, that of determining the phase angle  $\alpha_{hkl}$  for each reflection. In the remainder of this chapter, I will discuss some of the common methods for overcoming this obstacle. These include the *heavy-atom method* (also called *isomorphous replacement*), *anomalous scattering* (also called *anomalous dispersion*), and *molecular replacement*. Each of these techniques yield only estimates of phases, which must be improved before an interpretable electron-density map

116

can be obtained. In addition, these techniques usually yield estimates for a limited number of the phases, so phase determination must be extended to include as many reflections as possible. In Chapter 7, I will discuss methods of phase improvement and phase extension, which ultimately result in accurate phases and an interpretable electron-density map.

### 6.3 Isomorphous replacement

Each atom in the unit cell contributes to every reflection in the diffraction pattern [Eq. (5.15)]. The contribution of an atom is greatest to the reflections whose indices correspond to lattice planes that intersect that atom, so a specific atom contributes to some reflections strongly, and to some weakly or not at all. If we could add one or a very small number of atoms to identical sites in all unit cells of a crystal, we would expect to see changes in the diffraction pattern, as the result of the additional contributions of the added atom. As I will show later, the slight perturbation in the diffraction pattern caused by an added atom can be used to obtain initial estimates of phases. In order for these perturbations to be large enough to measure, the added atom must be a strong diffractor, which means it must be an element of higher atomic number than most of the other atoms-a so-called heavy atom. In smallmolecule crystallography, a larger atom like sulfur that is already present among a modest number of smaller atoms like carbon and hydrogen can be used for this purpose. This approach is properly called the *heavy-atom method*. For proteins, it is usually necessary to *add* a heavy atom, such as mercury, lead, or gold, in order to see perturbations against the signal of hundreds or thousands of smaller atoms. Addition of one or more heavy atoms to a protein for phasing is properly called *isomorphous replacement*, but is sometimes loosely referred to as the heavy-atom method.

### 6.3.1 Preparing heavy-atom derivatives

After obtaining a complete set of X-ray data and determining that these data are adequate to produce a high-resolution structure, the crystallographer undertakes to prepare one or more *heavy-atom derivatives*. In the most common technique, crystals of the protein are soaked in solutions of heavy ions, for instance ions or ionic complexes of Hg, Pt, or Au. In many cases, such ions bind to one or a few specific sites on the protein without perturbing its conformation or crystal packing. For instance, surface cysteine residues react readily with Hg<sup>2+</sup> ions, and cysteine, histidine, and methionine displace chloride from Pt complexes like PtCl<sub>4</sub><sup>2-</sup> to form stable Pt adducts. The conditions that give such specific binding must be found by simply trying different ionic compounds at various pH values and concentrations. (The same companies that sell crystal screening kits sell heavy-atom screening kits.)
#### Chapter 6 Obtaining Phases

Several diffraction criteria define a promising heavy-atom derivative. First, the derivative crystals must be *isomorphous* with native crystals. At the molecular level, this means that the heavy atom must not disturb crystal packing or the conformation of the protein. Unit-cell dimensions are quite sensitive to such disturbances, so heavy-atom derivatives whose unit-cell dimensions are the same as native crystals are probably isomorphous. The term *isomorphous replacement* comes from this criterion.

The second criterion for useful heavy-atom derivatives is that there must be measurable changes in at least a modest number of reflection intensities. These changes are the handle by which phase estimates are pulled from the data, so they must be clearly detectable, and large enough to measure accurately.

Figure 6.6 shows precession photographs for native and derivative crystals of the MoFe protein of nitrogenase. Underlined in the figure are pairs of reflections whose relative intensities are altered by the heavy atom. If examining heavy-atom data frames by eye, the crystallographer would look for pairs of reflections whose relative intensities are reversed. This distinguishes real heavy-atom perturbations from simple differences in overall intensity of two photos. For example, consider the leftmost underlined pairs in each photograph. In the native photo (a), the reflection on the right is the darker of the pair, whereas in the derivative photo (b), the reflection on the left is darker. Several additional differences suggest that this derivative might produce good phases. In modern practice, software computes intensity differences between native and heavy-atom reflections, and gives a quantitative measure of the potential *phasing power* of a heavy-atom derivative.



**Figure 6.6**  $\triangleright$  Precession photographs of the *hk*0 plane in native (*a*) and heavy-atom (*b*) crystals of the MoFe protein from nitrogenase. Corresponding underlined pairs in the native and heavy-atom patterns show reversed relative intensities. Photos courtesy of Professor Jeffrey Bolin.

#### Section 6.3 Isomorphous replacement

Finally, the derivative crystal must diffract to reasonably high resolution, although the resolution of derivative data need not be as high as that of native data. Methods of phase extension (Chapter 7) can produce phases for higher-angle reflections from good phases of reflections at lower angles.

Having obtained a suitable derivative, the crystallographer faces data collection again. Because derivatives must be isomorphous with native crystals, the strategy is the same as that for collecting native data. You can see that the phase problem effectively multiplies the magnitude of the crystallographic project by the number of derivative data sets needed. As I will show, at least two, and often more, derivatives are required in isomorphous replacement.

#### 6.3.2 Obtaining phases from heavy-atom data

Consider a single reflection of amplitude  $|\mathbf{F}_P|$  (P for protein) in the native data, and the corresponding reflection of amplitude  $|\mathbf{F}_{PH}|$  (PH for protein plus heavy atom) in data from a heavy-atom derivative. Because the diffractive contributions of all atoms to a reflection are additive, the difference in amplitudes  $(|\mathbf{F}_{PH}| - |\mathbf{F}_P|)$  is the amplitude contribution of the heavy atom alone. If we compute a diffraction pattern in which the amplitude of each reflection is  $(|\mathbf{F}_{PH}| - |\mathbf{F}_P|)^2$ , the result is the diffraction pattern of the heavy atom *alone* in the protein's unit cell, as shown in Fig. 6.7. In effect, we have subtracted away all contributions from the protein atoms, leaving only the heavy-atom contributions. Now we see the diffraction pattern of one atom (or only a small number of atoms) rather than the far more complex pattern of the protein.

In comparison to the protein structure, this "structure"—a sphere (or very few spheres) in a lattice—is very simple. It is usually easy to "determine" this structure, that is, to find the location of the heavy atom in the unit cell. Before considering how to locate the heavy atom, I will show how finding it helps us to solve the phase problem.

Suppose we are able to locate a heavy atom in the unit cell of derivative crystals. Recall that Eq. (5.15) gives us the means to calculate the structure factors  $\mathbf{F}_{hkl}$  for a known structure. This calculation gives us not just the amplitudes but the complete structure factors, including each of their phases. So we can compute the amplitudes and phases of our simple structure, the heavy atom in the protein unit cell (Fig. 6.7*e*). Now consider a single reflection *hkl* as it appears in the native and derivative data. Let the structure factor of the native reflection be  $\mathbf{F}_{PL}$ . Let the structure factor of the corresponding derivative reflection be  $\mathbf{F}_{PH}$ . Finally, let  $\mathbf{F}_{H}$  be the structure factor for the heavy atom itself, which we can compute if we can locate the heavy atom.

Figure 6.8 shows the relationship among the vectors  $\mathbf{F}_{P}$ ,  $\mathbf{F}_{PH}$ , and  $\mathbf{F}_{H}$  on the complex plane. (Remember that we are considering this relationship for a specific reflection, but the same relationship holds, and the same kind of equation must be solved, for all reflections.) Because the diffractive contributions of atoms are additive vectors,

$$\mathbf{F}_{\rm PH} = \mathbf{F}_{\rm H} + \mathbf{F}_{\rm P}.\tag{6.8}$$



**Figure 6.7** Heavy-atom method. (*a*) Protein in unit cell, and its diffraction pattern. (*b*) Heavy-atom derivative, and its diffraction pattern. Can you find slight differences in relative intensities of reflections in (*a*) and (*b*)? Taking the difference between diffraction patterns in (*a*) and (*b*) gives (*c*), the diffraction pattern of the heavy atom alone. (*d*) Interpretation of (*c*) by Patterson methods locates the heavy atom (dark gray) in the unit cell of the protein. FT of (*d*) gives (*e*), the structure factors of the heavy atom, *with phases*. These  $\mathbf{F}_{H}$  terms allow solution of Eq. (6.9).

That is, the structure factor for the heavy-atom derivative (red in the figure) is the vector sum of the structure factors for the protein alone (green) and the heavy atom alone (blue). For each reflection, we wish to know  $\mathbf{F}_{\rm P}$ . (We already know that its length is obtainable from the measured reflection intensity  $I_{\rm P}$ , but we want to learn its phase angle.) According to the previous equation,

$$\mathbf{F}_{\mathrm{P}} = \mathbf{F}_{\mathrm{PH}} - \mathbf{F}_{\mathrm{H}}.\tag{6.9}$$

We can solve this vector equation for  $\mathbf{F}_{P}$ , and thus obtain the phase angle of the structure factor, by use of a *Harker diagram* which represents the equation in the complex plane (Fig. 6.9).

We know  $|\mathbf{F}_{PH}|$  and  $|\mathbf{F}_{P}|$  from measuring reflection intensities  $I_{PH}$  and  $I_P$ . So we know the length of the vectors  $\mathbf{F}_{PH}$  and  $\mathbf{F}_P$ , but not their directions or phase angles. We know  $\mathbf{F}_H$ , *including its phase angle*, from locating the heavy atom and calculating all its structure factors (Fig. 6.7). To solve Eq. (6.9) for  $\mathbf{F}_P$  and thus obtain its phase angle, we place the vector  $-\mathbf{F}_H$  at the origin and draw a circle of radius  $|\mathbf{F}_{PH}|$  centered on the head of vector  $-\mathbf{F}_H$  (Fig. 6.9*a*). All points on this



**Figure 6.8** A structure factor  $\mathbf{F}_{PH}$  for the heavy-atom derivative is the sum of contributions from the native structure ( $\mathbf{F}_P$ ) and the heavy atom ( $\mathbf{F}_H$ ). COLOR KEY: In this and all subsequent diagrams of this type (called Harker diagrams), the structure factor whose phase we are seeking (in this case, that of the protein) is green. The structure factor of the heavy atom (or its equivalent in other types of phasing) is blue, and the structure factor of the heavy atom derivative of the protein is red. Circles depicting possible orientations of structure factors carry the same colors.



**Figure 6.9**  $\triangleright$  Vector solution of Eq. (6.9). (a) All points on the red circle equal the vector sum  $|\mathbf{F}_{PH}| - \mathbf{F}_{H}$ . (b) Vectors from the origin to intersections of the two circles are solutions to Eq. (6.9).

circle equal the vector sum  $|\mathbf{F}_{PH}| - \mathbf{F}_{H}$ . In other words, we know that the head of  $\mathbf{F}_{PH}$  lies somewhere on this circle of radius  $|\mathbf{F}_{PH}|$ . Next, we add a circle of radius  $|\mathbf{F}_{PH}|$  centered at the origin (Fig. 6.9*b*). We know that the head of the vector  $\mathbf{F}_{P}$  lies somewhere on this circle, but we do not know where because we do not know its phase angle. Equation (6.9) holds only at points where the two circles intersect. Thus the phase angles of the two vectors  $\mathbf{F}_{P}^{a}$  and  $\mathbf{F}_{P}^{b}$  that terminate at the points of intersection of the circles are the only possible phases for this reflection.

Our heavy-atom derivative allows us to determine, for each reflection hkl, that  $\alpha_{hkl}$  has one of two values. How do we decide which of the two phases is correct? In some cases, if the two intersections lie near each other, the average of the two phase angles will serve as a reasonable estimate. I will show in Chapter 7 that certain phase improvement methods can sometimes succeed with such phases from only one derivative, in which case the structure is said to be solved by the method of *single isomorphous replacement* (SIR). More commonly, however, a second heavy-atom derivative must be found and the vector problem outlined previously must be solved again. Of the two possible phase angles found by using the second derivative, one should agree better with one of the two solutions from the first derivative, as shown in Fig. 6.10.

Figure 6.10*a* shows the phase determination using a second heavy-atom derivative;  $\mathbf{F}'_{\rm H}$  is the structure factor for the second heavy atom. The radius of the red circle is  $|\mathbf{F}'_{\rm PH}|$ , the amplitude of  $\mathbf{F}'_{\rm PH}$  for the second heavy-atom derivative. For this derivative,  $\mathbf{F}_{\rm P} = \mathbf{F}'_{\rm PH} - \mathbf{F}'_{\rm H}$ . Construction as before shows that the phase angles of  $\mathbf{F}_{\rm P}^{\rm c}$  and  $\mathbf{F}_{\rm P}^{\rm d}$  are possible phases for this reflection. In Fig. 6.10*b*, the circles from



**Figure 6.10** (a) A second heavy-atom derivative indicates two possible phases, one of which corresponds to  $\mathbf{F}^{a}$  in Fig. 6.9*b*. (b)  $\mathbf{F}_{P}$ , which points from the origin to the common intersection of the three circles, is the solution to Eq. (6.9) for both heavy-atom derivatives. Thus  $\alpha$  is the correct phase for this reflection.



**Figure 6.11**  $\blacktriangleright$  The MIR solution for this structure factor gives phase of high uncertainty.

Figs. 6.9*b* and 6.10*a* are superimposed, showing that  $\mathbf{F}_{P}^{c}$  is identical to  $\mathbf{F}_{P}^{a}$ . This common solution to the two vector equations is  $\mathbf{F}_{P}$ , the desired structure factor. The phase of this reflection is therefore the angle labeled  $\alpha$  in the figure, the only phase compatible with data from both derivatives.

In order to resolve the phase ambiguity from the first heavy-atom derivative, the second heavy atom must bind at a different site from the first. If two heavy atoms bind at the same site, the phases of  $\mathbf{F}_{\rm H}$  will be the same in both cases, and both phase determinations will provide the same information. This is true because the phase of an atomic structure factor depends only on the location of the atom in the unit cell, and not on its identity (Sec. 5.3.1, p. 98). In practice, it sometimes takes three or more heavy-atom derivatives to produce enough phase estimates to make the needed initial dent in the phase problem. Obtaining phases with two or more derivatives is called the method of multiple isomorphous replacement (MIR). For many years, MIR was one of the most successful methods in macromolecular crystallography.

To compute a high-resolution structure, we must ultimately know the phases  $\alpha_{hkl}$  for all reflections. High-speed computers can solve large numbers of these vector problems rapidly, yielding an estimate of each phase along with a measure of its precision.<sup>1</sup> For many phases, the precision of the first phase estimate is so low that the phase is unusable. For instance, in Fig. 6.11, the circles graze each other rather than intersecting sharply, so there is a large uncertainty in  $\alpha$ . In some cases,

$$\mathbf{F}_{\rm PH}^2 = \mathbf{F}_{\rm P}^2 + \mathbf{F}_{\rm H}^2 + 2\mathbf{F}_{\rm P}\mathbf{F}_{\rm H}\cos(\alpha_{\rm P} - \alpha_{\rm H})$$

The pairs of solutions for two heavy-atom derivatives should have one solution in common.

<sup>&</sup>lt;sup>1</sup>Computer programs calculate phases for each derivative numerically (rather than geometrically) by obtaining two solutions to the equation

because of inevitable experimental errors in measuring intensities, the circles do not intersect at all. This situation is referred to as *lack of closure*, and there are computer algorithms for making phase estimates when it occurs.

Computer programs for calculating phases also compute statistical parameters representing attempts to judge the quality of phases. Some parameters, usually called *phase probabilities*, are measures of the uncertainty of individual phases. Other parameters, including figure of merit, closure errors, phase differences, and various *R*-factors are attempts to assess the quality of groups of phases obtained by averaging results from several heavy-atom derivatives (or results from other phasing methods). In most cases, these parameters are numbers between 0 (poor phases) and 1 (perfect phases). No single one of these statistics is an accurate measure of the goodness of phases. Crystallographers often use two or more of these criteria simultaneously in order to cull out questionable phases. In short, until correct phases are obtained (see Chapter 7), there is no sure way to measure the quality of estimates. The acid test of phases is whether they give an interpretable electron-density map. In Chapter 7, I will say more about the most modern methods of improving phases at all stages in crystallography.

When promising phases are available, the crystallographer carries out Fourier summation [Eq. (6.7)] to calculate  $\rho(x, y, z)$ . Each Fourier term is multiplied by the probability of correctness of the associated phase. This procedure gives greater weight to terms with more reliable phases. Every phase that defies solution or is too uncertain (and for that matter every intensity that is too weak to measure accurately) forces the crystallographer to omit one term from the Fourier series when calculating  $\rho(x, y, z)$ . Each omitted term lowers the accuracy of the approximation to  $\rho(x, y, z)$ , degrading the quality and resolution of the resulting map. In practice, a good pair of heavy-atom derivatives may allow us to estimate only a small percentage of the phases. We can enlarge our list of precise phases by iterative processes mentioned briefly in Sec. 5.4.2, p. 104, which I will describe more fully in Chapter 7. For now, I will complete this discussion of isomorphous replacement by considering how to find heavy atoms, which is necessary for calculating  $\mathbf{F}_{\mathrm{H}}$ .

#### 6.3.3 Locating heavy atoms in the unit cell

Before we can obtain phase estimates by the method described in the previous section, we must locate the heavy atoms in the unit cell of derivative crystals. As I described earlier, this entails extracting the relatively simple diffraction signature of the heavy atom from the far more complicated diffraction pattern of the heavy-atom derivative, and then solving a simpler "structure," that of one heavy atom (or a few) in the unit cell of the protein (Fig. 6.7). The most powerful tool in determining the heavy-atom coordinates is a Fourier sum called the *Patterson function* P(u, v, w), a variation on the Fourier sum I have described for computing  $\rho(x, y, z)$  from structure factors. The coordinates (x, y, z) locate a point in a *Patterson map*, in the same way that coordinates (x, y, z) locate a point in an electron-density map. The Patterson function or Patterson synthesis is a Fourier sum without phases. The amplitude of each term is the square of one structure factor, which is proportional to the measured reflection intensity.

124

#### Section 6.3 Isomorphous replacement

Thus we can construct this series from intensity measurements, even though we have no phase information. Here is the Patterson function in general form:

$$P(u, v, w) = \frac{1}{V} \sum_{h} \sum_{k} \sum_{l} |\mathbf{F}_{hkl}|^2 e^{-2\pi i (hu + kv + lw)}.$$
 (6.10)

To obtain the Patterson function solely for the heavy atoms in derivative crystals, we construct a *difference Patterson function*, in which the amplitudes are  $(\Delta \mathbf{F})^2 = (|\mathbf{F}_{PH}| - |\mathbf{F}_{P}|)^2$ . The difference between the structure-factor amplitudes with and without the heavy atom reflects the contribution of the heavy atom alone. The difference Patterson function is

$$\Delta P(u, v, w) = \frac{1}{V} \sum_{h} \sum_{k} \sum_{l} \Delta \mathbf{F}_{hkl}^2 e^{-2\pi i (hu+kv+lw)}.$$
(6.11)

In words, the difference Patterson function is a Fourier sum of simple sine and cosine terms. (Remember that the exponential term is shorthand for these trigonometric functions.) Each term in the series is derived from one reflection hkl in both the native and derivative data sets, and the amplitude of each term is  $(|\mathbf{F}_{\rm PH}| - |\mathbf{F}_{\rm P}|)^2$ , which is the amplitude contribution of the heavy atom to structure factor  $\mathbf{F}_{\rm PH}$ . Each term has three frequencies h in the *u*-direction, k in the *v*-direction, and l in the *w*-direction. Phases of the structure factors are not included; at this point, they are unknown. (If we knew them, we wouldn't have to do all this.)

Because the Patterson function contains no phases, it can be computed from any raw set of crystallographic data, but what does it tell us? A contour map of  $\rho(x, y, z)$ displays areas of high density (peaks) at the locations of atoms. In contrast, it can be proven that a Patterson map, which is a contour map of P(u, v, w), displays peaks at locations corresponding to *vectors between atoms*. (This is a strange idea at first, but the following example will make it clearer.) Of course, there are more vectors between atoms than there are atoms, so a Patterson map is more complicated than an electron-density map. But if the structure is simple, like that of one or a few heavy atoms in the unit cell, the Patterson map may be simple enough to allow us to locate the atom(s). You can see now that the main reason for using the difference Patterson function instead of a simple Patterson using  $\mathbf{F}_{PHS}$ is to eliminate the enormous number of peaks representing vectors between light atoms in the protein.

I will show, in a two-dimensional example, how to construct the Patterson map from a simple crystal structure and then how to use a calculated Patterson map to deduce a structure (Fig. 6.12). The simple molecular structure in Fig. 6.12*a* contains three atoms (red circles) in each unit cell. To construct the Patterson map, first draw all possible vectors between atoms in one unit cell, including vectors between the same pair of atoms but in opposite directions. (For example, treat  $1 \rightarrow 2$  and  $2 \rightarrow 1$  as distinct vectors.) Two of the six vectors  $(1 \rightarrow 3 \text{ and } 3 \rightarrow 2)$ 



**Figure 6.12** Construction and interpretation of a Patterson map. (*a*) Structure of unit cell containing three atoms. Two of the six interatomic vectors are shown. (*b*) Patterson map is constructed by moving the tails of all interatomic vectors to the origin. Patterson "atoms" [peaks (purple) in the contour map] occur at the head of each vector. (*c*) Complete Patterson map, containing all peaks from (*b*) in all unit cells. Peak at origin results from self-vectors. Image of original structure is present (red peaks) amid other peaks. (*d*) Trial solution of map (*c*). If origin and Patterson atoms **a** and **b** were the image of the real unit cell, the interatomic vector  $\mathbf{a} \rightarrow \mathbf{b}$  would produce a peak in the small green box. Absence of the peak disproves this trial solution.

are shown in the figure. Then draw empty unit cells around an origin (Fig. 6.12*b*), and redraw all vectors with their tails at the origin. The head of each vector is the location of a peak in the Patterson map, sometimes called a Patterson "atom" (purple circles in *b*, *c*, and *d*). The coordinates (u, v, w) of a Patterson atom representing a vector between atom 1 at  $(x_1, y_1, z_1)$  and atom 2 at  $(x_2, y_2, z_2)$  are  $(u, v, w) = (x_1 - x_2, y_1 - y_2, z_1 - z_2)$ . The vectors from Fig. 6.12*a* are redrawn in Fig. 6.12*b*, along with all additional Patterson atoms produced by this procedure. Finally, in each of the unit cells, duplicate the Patterson map of the structure in Fig. 6.12*a*. In this case, there are six Patterson atoms in each unit cell. You can easily prove to yourself that a real unit cell containing *n* atoms will give a Patterson unit cell containing n(n - 1) Patterson atoms.

Now let's think about how to go from a computed Patterson map to a structure that is, how to locate real atoms from Patterson atoms. A computed Patterson map exhibits a strong peak at the origin because this is the location of all vectors

#### Section 6.3 Isomorphous replacement

between an atom and itself. Notice in Fig. 6.12*c* that the origin and two of the Patterson atoms (red circles) reconstruct the original arrangement of atoms. Finding six peaks (ignoring the peak at the origin) in each unit cell of the calculated Patterson map, we infer that there are three real atoms per unit cell. [Solve the equation n(n-1) = 6.] We therefore know that the origin and two peaks reconstruct the relationship among the three real atoms, but we do not know which two peaks to choose. To solve the problem, we pick a set of peaks—the origin and two others—as a trial solution, and follow the rules described earlier to generate the expected Patterson map for this arrangement of atoms. If the trial map has the same peaks as the calculated map, then the trial arrangement of atoms is correct. By trial and error, we can determine which pair of Patterson atoms, along with an atom at the origin, would produce the remaining Patterson atoms. Figure 6.12*d* shows an incorrect solution (the origin plus green peaks **a** and **b**). The vector  $\mathbf{a} \rightarrow \mathbf{b}$  is redrawn at the origin to show that the map does not contain the Patterson atom  $\mathbf{a} \rightarrow \mathbf{b}$ , and hence that this solution is incorrect.

You can see that as the number of real atoms increases, the number of Patterson atoms, and with it the difficulty of this problem, increases rapidly. Computer programs can search for solutions to such problems and, upon finding a solution, can refine the atom positions to give the most likely arrangement of heavy atoms.

Unit-cell symmetry can also simplify the search for peaks in a three-dimensional Patterson map. For instance, in a unit cell with a  $2_1$  axis (twofold screw) on edge **b**, recall (equivalent positions, Sec. 4.2.8, p. 65) that each atom at (x, y, z) has an identical counterpart atom at (-x, y + 1/2, -z). The vectors connecting such symmetry-related atoms will all lie at (u, v, w) = (2x, [1/2], 2z) in the Patterson map (just subtract one set of coordinates from the other, and realize that the position v = -1/2 is the same as v = 1/2), which means they all lie in the plane that cuts the Patterson unit cell at w = 1/2—the u[1/2]w plane. Such planes, which contain the Patterson vectors for symmetry-related atoms, are called *Harker sections* or *Harker planes*. If heavy atoms bind to the protein at equivalent positions, heavy-atom peaks in the Patterson vectors that all lie upon a line, called *a Harker line*, rather than on a plane.)

Given the location of a first heavy-atom Patterson peak on a Harker section, what is the location of the heavy atom in the real unit cell? In the  $P2_1$  unit cell, the equality u = 2x means that x = u/2, giving the x-coordinate of the atom. In like manner, z = w/2. The location of the peak on a Harker section at v = 1/2 places no restrictions on the y-coordinate, so it can be given a convenient arbitrary value like y = 0. Thus the heavy-atom coordinates in the real unit cell are (u/2, 0, w/2). Additional heavy-atom sites will have their y-coordinates specified relative to this arbitrary assignment.

There is an added complication (in crystallography, it seems there always is): the arrangement of heavy atoms in a protein unit cell is often enantiomeric. For example, if heavy atoms are found along a threefold screw axis, the screw may be left- or right-handed. The Patterson map does not distinguish between mirrorimage, or more accurately, inverted-image, arrangements of heavy atoms because you cannot tell whether a Patterson vector **ab** is  $\mathbf{a} \rightarrow \mathbf{b}$  or  $\mathbf{b} \rightarrow \mathbf{a}$ . But the phases obtained by calculating structure factors from the inverted model are incorrect and will not lead to an interpretable map. Crystallographers refer to this difficulty as the *hand problem* (although hands are mirror images and the two solutions described here are inversions). If derivative data are available to high resolution, the crystallographer simply calculates two electron-density maps, one with phases from each enantiomer of the heavy-atom structure. With luck, one of these maps will be distinctly clearer than the other. If derivative data are available only at low resolution, this method may not determine the hand with certainty. The problem may require the use of anomalous scattering methods, discussed in Sec. 6.4.6, p. 135.

Having located the heavy atom(s) in the unit cell, the crystallographer can compute the structure factors  $\mathbf{F}_{\rm H}$  for the heavy atoms alone, using Eq. (5.15). This calculation yields both the amplitudes and the phases of structure factors  $\mathbf{F}_{\rm H}$ , giving the vector quantities needed to solve Eq. (6.9) for the phases  $\alpha_{hkl}$  of protein structure factors  $\mathbf{F}_{\rm P}$ . This completes the information needed to compute a first electron-density map, using Eq. (6.7), p. 116. This map requires improvement because these first phase estimates contain substantial errors. I will discuss improvement of phases and maps in Chapter 7.

## 6.4 Anomalous scattering

#### 6.4.1 Introduction

A second means of obtaining phases from heavy-atom derivatives takes advantage of the heavy atom's capacity to absorb X-rays of specified wavelength. As a result of this absorption, Friedel's law (Sec. 4.3.7, p. 88) does not hold, and the reflections hkl and  $\overline{hkl}$  are not equal in intensity. This inequality of symmetry-related reflections is called *anomalous scattering* or *anomalous dispersion*.

Recall from Sec. 4.3.2, p. 73 that elements *absorb* X-rays as well as emit them, and that this absorption drops sharply at wavelengths just below their characteristic emission wavelength  $K_{\beta}$  (Fig. 4.21, p. 74). This sudden change in absorption as a function of  $\lambda$  is called an *absorption edge*. An element exhibits anomalous scattering when the X-ray wavelength is near the element's absorption edge. Absorption edges for the light atoms in the unit cell are not near the wavelength of X-rays used in crystallography, so carbon, nitrogen, and oxygen do not contribute to anomalous scattering. However, absorption edges of heavy atoms are in this range, and if X-rays of varying wavelength are available, as is often the case at synchrotron sources, X-ray data can be collected under conditions that maximize anomalous scattering by the heavy atom.

#### 6.4.2 Measurable effects of anomalous scattering

When the X-ray wavelength is near the heavy-atom absorption edge, a fraction of the radiation is absorbed by the heavy atom and reemitted with altered phase.



**Figure 6.13** Real and imaginary anomalous-scattering contributions alter the magnitude and phase of the structure factor. COLOR KEY: By analogy to colors used in MIR illustrations, the vector whose phase we want to find, in this case,  $F_{PH}^{\lambda 1}$ , is green; the anomalous scattering contributions, whose role is analogous to that of a heavy atom, are shades of blue; and  $F_{PH}^{\lambda 2}$ , the "anomalous-dispersion derivative" that takes the role of the heavy atom derivative in MIR, is red.

The effect of this anomalous scattering on a given structure factor  $\mathbf{F}_{PH}$  in the heavy-atom data is depicted in vector diagrams as consisting of two perpendicular contributions, one real ( $\Delta \mathbf{F}_r$ ) and the other imaginary ( $\Delta \mathbf{F}_i$ ).

In Fig. 6.13,  $\mathbf{F}_{PH}^{\lambda 1}$  (green) represents a structure factor for the heavy-atom derivative measured at wavelength  $\lambda_1$ , where anomalous scattering does not occur.  $\mathbf{F}_{HP}^{\lambda 2}$  (red) is the same structure factor measured at a second X-ray wavelength  $\lambda_2$  near the absorption edge of the heavy atom, so anomalous scattering alters the heavy-atom contribution to this structure factor. The vectors representing anomalous scattering contributions are  $\Delta \mathbf{F}_r$  (real, blue) and  $\Delta \mathbf{F}_i$  (imaginary, cyan). From the diagram, you can see that

$$\mathbf{F}_{\rm PH}^{\lambda 2} = \mathbf{F}_{\rm PH}^{\lambda 1} + \Delta \mathbf{F}_{\rm r} + \Delta \mathbf{F}_{\rm i}. \tag{6.12}$$

Figure 6.14 shows the result of anomalous scattering for a Friedel pair of structure factors, distinguished from each other in the figure by superscripts + and –. Recall that for Friedel pairs in the absence of anomalous scattering,  $|\mathbf{F}_{hkl}| = |\mathbf{F}_{hkl}|$  and  $\alpha_{hkl} = -\alpha_{\overline{hkl}}$ , so  $\mathbf{F}_{PH}^{\lambda 1-}$  is the reflection of  $\mathbf{F}_{PH}^{\lambda 1+}$  in the real axis. The real contributions  $\Delta \mathbf{F}_{r}^{+}$  and  $\Delta \mathbf{F}_{r}^{-}$  to the reflections of a Friedel pair are, like the structure factors themselves, reflections of each other in the real axis. On the other hand, it can be shown (but I will not prove it here) that the imaginary contribution to  $\mathbf{F}_{PH}^{\lambda 1-}$  is the inverted reflection of that for  $\mathbf{F}_{PH}^{\lambda 1+}$ . That is,  $\Delta \mathbf{F}_{i}^{-}$  is obtained by reflecting  $\Delta \mathbf{F}_{i}^{+}$  in the real axis and then reversing its sign or pointing it in the opposite direction. Because of this difference between the imaginary contributions



**Figure 6.14** Under anomalous scattering, at wavelength  $\lambda_2$ ,  $\mathbf{F}_{\overline{hkl}}$  is no longer the mirror image of  $\mathbf{F}_{hkl}$ .

to these reflections, under anomalous scattering, the two structure factors are no longer precisely equal in intensity, nor are they precisely opposite in phase. It is clear from Fig. 6.14 that  $\mathbf{F}_{PH}^{\lambda 2-}$  is not the mirror image of  $\mathbf{F}_{PH}^{\lambda 2+}$ . From this disparity between Friedel pairs, the crystallographer can extract phase information.

#### 6.4.3 Extracting phases from anomalous scattering data

The magnitude of anomalous scattering contributions  $\Delta F_r$  and  $\Delta F_i$  for a given element are constant and roughly independent of reflection angle  $\theta$ , so these quantities can be looked up in tables of crystallographic information. The phases of  $\Delta F_r$  and  $\Delta F_i$  depend only upon the position of the heavy atom in the unit cell, so once the heavy atom is located by Patterson methods, the phases can be computed. The resulting full knowledge of  $\Delta F_r$  and  $\Delta F_i$  allows Eq. (6.12) to be solved for the vector  $\mathbf{F}_{PH}^{\lambda 1}$ , thus establishing its phase. Crystallographers obtain solutions by computer, but I will solve the general equation using Harker diagrams (Fig. 6.15), and thus show that the amount of information is adequate to solve the problem. First consider the structure factor  $\mathbf{F}_{PH}^{\lambda 1+}$  in Fig. 6.14. Applying Eq. (6.12) and solving for  $\mathbf{F}_{PH}^{\lambda 1+}$  gives

$$\mathbf{F}_{\mathrm{PH}}^{\lambda 1+} = \mathbf{F}_{\mathrm{PH}}^{\lambda 2+} - \Delta \mathbf{F}_{\mathrm{r}}^{+} - \Delta \mathbf{F}_{\mathrm{i}}^{+}.$$
(6.13)

To solve this equation (see Fig. 6.15), draw the vector  $-\Delta \mathbf{F}_r^+$  with its tail at the origin, and draw  $-\Delta \mathbf{F}_i^+$  with its tail on the head of  $-\Delta \mathbf{F}_r^+$ . With the head of  $-\Delta \mathbf{F}_i^+$  as center, draw a circle of radius  $|\mathbf{F}_{PH}^{\lambda 2+}|$ , representing the amplitude of this reflection in the anomalous scattering data set. The head of the vector  $\mathbf{F}_{PH}^{\lambda 2+}$  lies somewhere on this circle. We do not know where, because we do not know the phase of the reflection. Now draw a circle of radius  $|\mathbf{F}_{PH}^{\lambda 1+}|$  with its center at the



**Figure 6.15**  $\blacktriangleright$  Vector solution of Eq. (6.13).  $\Delta \mathbf{F}_r$  and  $\Delta \mathbf{F}_i$  play the same role as  $\mathbf{F}_H$  in Figs. 6.9, p. 121, and 6.10, p. 122.

origin, representing the structure-factor amplitude of this same reflection in the nonanomalous scattering data set. The two points of intersection of these circles satisfy Eq. (6.13), establishing the phase of this reflection as either that of  $\mathbf{F}_a$  or  $\mathbf{F}_b$ . As with the SIR method, we cannot tell which of the two phases is correct.

The Friedel partner of this reflection comes to the rescue. We can obtain a second vector equation involving  $\mathbf{F}_{PH}^{\lambda 1+}$  by reflecting  $\mathbf{F}_{PH}^{\lambda 2-}$  and all its vector components across the real axis (Fig. 6.16*a*).

After reflection,  $\mathbf{F}_{PH}^{\lambda 1-}$  equals  $\mathbf{F}_{PH}^{\lambda 1+}$ ,  $\Delta \mathbf{F}_{r}^{-}$  equals  $\Delta \mathbf{F}_{r}^{+}$ , and  $\Delta \mathbf{F}_{i}^{-}$  equals  $-\Delta \mathbf{F}_{i}^{+}$ . The magnitude of  $\mathbf{F}_{PH}^{\lambda 2-}$  is unaltered by reflection across the real axis. If we make these substitutions in Eq. (6.13), we obtain

$$\mathbf{F}_{\mathrm{PH}}^{\lambda 1+} = \left| \mathbf{F}_{\mathrm{PH}}^{\lambda 2-} \right| - \Delta \mathbf{F}_{\mathrm{r}}^{+} - (-\Delta \mathbf{F}_{\mathrm{i}}^{+}). \tag{6.14}$$

We can solve this equation in the same manner that we solved Eq. (6.13), by placing the vectors  $-\Delta \mathbf{F}_{r}^{+}$  and  $\Delta \mathbf{F}_{i}^{+}$  head-to-tail at the origin, and drawing a circle of radius  $|\mathbf{F}_{PH}^{\lambda 2}|$  centered on the head of  $\Delta \mathbf{F}_{i}^{+}$  (Fig. 6.16*b*). The circles intersect at the two solutions to Eq. (6.14). Although the circles graze each other and give two phases with considerable uncertainty, one of the possible solutions corresponds to  $\mathbf{F}_{a}$  in Fig. 6.15, and neither of them is close to the phase of  $\mathbf{F}_{b}$ .

So the disparity between intensities of Friedel pairs in the anomalous scattering data set establishes their phases in the nonanomalous scattering data set. The reflection whose phase has been established here corresponds to the vector  $\mathbf{F}_{PH}$  in Eq. (6.9). Thus the amplitudes and phases of two of the three vectors in the Eq. (6.9) are known: (1)  $\mathbf{F}_{PH}$  is known from the anomalous scattering computation just



**Figure 6.16** Reflection of  $\mathbf{F}^-$  components across the real axis gives a second vector equation involving the desired structure factor. (*a*) All reflected components are labeled with their equivalent contributions from  $\mathbf{F}^+$ . (*b*) Vector solution of Eq. (6.13). These solutions are compatible only with  $\mathbf{F}_a$  in Fig. 6.15.

shown, and (2)  $\mathbf{F}_{H}$  is known from calculating the heavy-atom structure factors after locating the heavy atom by Patterson methods. The vector  $\mathbf{F}_{P}$ , then, is simply the vector difference  $\mathbf{F}_{PH} - \mathbf{F}_{H}$ , establishing the phase of this reflection in the native data.

#### 6.4.4 Summary

Under anomalous scattering, the members of a Friedel pair can be used to establish the phase of a reflection in the heavy atom derivative data, thus establishing the phase of the corresponding reflection in the native data. Let me briefly review the entire project of obtaining the initial structure factors by SIR with anomalous scattering (called SIRAS). First, we collect a complete data set with native crystals, giving us the amplitudes  $|\mathbf{F}_{\rm P}|$  for each of the native reflections. Then we find a heavy-atom derivative and collect a second data set at the same wavelength, giving amplitudes  $|\mathbf{F}_{\rm PH}|$  for each of the reflections in the heavy-atom data. Next we collect a third data set at a different X-ray wavelength, chosen to maximize anomalous scattering by the heavy atom. We use the nonequivalence of Friedel pairs in the anomalous scattering data to establish phases of reflections in the heavy-atom

#### Section 6.4 Anomalous scattering

data, and we use the phased heavy-atom derivative structure factors to establish the native phases. (Puff-puff!)

In practice, several of the most commonly used heavy atoms (including uranium, mercury, and platinum) give strong anomalous scattering with Cu–K<sub> $\alpha$ </sub> radiation. In such cases, crystallographers can measure intensities of Friedel pairs in the heavy-atom data set. In phase determination (refer to Figs. 6.14–6.16), the average of  $|\mathbf{F}_{hkl}|$  and  $|\mathbf{F}_{\overline{hkl}}|$  serves as both  $|\mathbf{F}_{PH}^{\lambda 1+}|$  and  $|\mathbf{F}_{PH}^{\lambda 1-}|$ , while  $|\mathbf{F}_{hkl}|$  and  $|\mathbf{F}_{\overline{hkl}}|$  separately serve as  $|\mathbf{F}_{PH}^{\lambda 2+}|$  and  $|\mathbf{F}_{PH}^{\lambda 2-}|$ , so only one heavy atom data set is required. This method is called single isomorphous replacement with anomalous scattering, or SIRAS.

Like phases from the MIR method, each anomalous scattering phase can only serve as an initial estimate and must be weighted with some measure of phase probability. The intensity differences between Friedel pairs are very small, so measured intensities must be very accurate if any usable phase information is to be derived. To improve accuracy, crystallographers collect intensities of Friedel partners from the same crystal, and under very similar conditions. In rotation/oscillation photography, crystallographers can alter the sequence of frame collection, so that frames of Friedel pairs are measured in succession, thus minimizing artifactual differences between Friedel pairs due to crystal deterioration or changes in the X-ray beam. For example, if we call a frame of data  $F_n$  and the frame of matching Friedel pairs  $F'_n$ , then minimizing artifactual differences between Friedel pairs entails collecting  $F_1$ , rotating the crystal 180° and collecting frame  $F'_1$  of matching Friedel pairs, then collecting  $F'_2$  starting at the end of the rotation for  $F'_1$ , and finally, rotating back 180° to collect  $F_2$ .

#### 6.4.5 Multiwavelength anomalous diffraction phasing

Three developments—variable-wavelength synchrotron X rays, cryocrystallography, and the production of proteins containing selenomethionine instead of the normal sulfur-containing methionine—have recently allowed rapid progress in maximizing the information obtainable from anomalous dispersion. For proteins that naturally contain a heavy atom, such as the iron in a globin or cytochrome, the native heavy atom provides the source of anomalous dispersion. Proteins lacking functional heavy atoms can be expressed in *Escherichia coli* containing exclusively selenomethionine. The selenium atoms then serve as heavy atoms in a protein that is essentially identical to the "native" form. Isomorphism is, of course, not a problem with these proteins, because the same protein serves as both the native and derivative forms.

The power of multiwavelength radiation is that data sets from a heavy-atom derivative at different wavelengths are in many respects like those from distinct heavy-atom derivatives. Especially in the neighborhood of the absorption maximum of the heavy atom [see, for example, the absorption spectra of copper and nickel (dotted lines) in Fig. 4.21, p. 74], the real and imaginary anomalous scattering factors  $\Delta \mathbf{F}_r$  and  $\Delta \mathbf{F}_i$  vary greatly with X-ray wavelength. At the absorption maximum,  $\Delta \mathbf{F}_i$  reaches its maximum value, whereas at the ascending

#### Chapter 6 Obtaining Phases

inflection point or edge,  $\Delta \mathbf{F}_r$  reaches a minimum and then increases farther from the absorption peak. So data sets taken at the heavy-atom absorption maximum, the edge, and at wavelengths distance from the maximum all have distinct values for the real and imaginary contributions of anomalous dispersion. Thus each measurement of a Freidel pair at a specific wavelength provides the components of distinct sets of phasing equations like those solved in Figs. 6.15 and 6.16. In addition to wavelength-dependent differences between Friedel pairs, individual reflection intensities vary slightly with wavelength (called *dispersive differences*), and these differences also contain phase information, which can be extracted by solving equations much like those for isomorphous replacement. All told, data sets at different wavelengths from a single crystal can contain sufficient phasing information to solve a structure if the molecule under study contains one or more atoms that give anomalous dispersion. This method is called *multiwavelength anomalous dispersion*, or MAD, phasing, and it has become one of the most widely used phasing methods.

The principles upon which MAD phasing are based have been known for years. But the method had to await the availability of variable-wavelength synchrotron X-ray sources. In addition, the intensity differences the crystallographer must measure are small, and until recently, these small signals were difficult to measure with sufficient precision. To maximize accuracy in the measured differences between Friedel pairs, corresponding pairs must be measured on the same crystal and at nearly the same time so that crystal condition and instrument parameters do not change between measurements. And even more demanding, crystal condition and instrument parameters should be sufficiently constant to allow complete data sets to be taken at several different wavelengths. These technical demands are met by cryocrystallography and synchrotron sources. Flash-freezing preserves the crystal in essentially unchanged condition through extensive data collection. Synchrotron sources provide X-rays of precisely controllable wavelength and also of high intensity, which shortens collection time.

The first successes of MAD phasing were small proteins that contained functional heavy atoms. Production of selenomethionine proteins opened the door to MAD phasing for nonmetalloproteins. But larger proteins may contain many methionines. Can there be too many heavy atoms? Recall that to solve SIR and anomalous-dispersion phase equations, we must know the position(s) of the heavy atom(s) in the unit cell. For MAD phasing, heavy atoms can be located by Patterson methods (Sec. 6.3.3, p. 124), which entails trial-and-error comparisons of the Patterson map with calculated Pattersons for various proposed models of heavy-atom locations. But Patterson maps of proteins containing many heavy atoms may require so many trials that they resist solution. On the other hand, locating a relatively large number (say, tens) of heavy atoms is similar in complexity to determine the structure of a "small" molecule. And indeed, as I will describe in Sec. 6.4.7, p. 135, the direct-phasing methods used in small-molecule crystallography have come to the rescue in locating numerous Se atoms in selenomethione proteins, allowing application of MAD phasing to larger and larger proteins, with no end in sight.

#### Section 6.4 Anomalous scattering

#### 6.4.6 Anomalous scattering and the hand problem

As I discussed in Sec. 6.3.3, p. 124, Patterson methods do not allow us to distinguish between enantiomeric arrangements of heavy atoms, and phases derived from heavy-atom positions of the wrong hand are incorrect. When high-resolution data are available for the heavy-atom derivative, phases and electron-density maps can be calculated for both enantiomeric possibilities. The map calculated with phases from the correct enantiomer will sometimes be demonstrably sharper and more interpretable. If not, and if anomalous scattering data are available, SIR *and* anomalous scattering phases can be computed for both hands, and maps can be prepared from the two sets of phases. The added phase information from anomalous scattering sometimes makes hand selection possible when SIR phases alone do not.

The availability of two heavy-atom derivatives, one with anomalous scattering, allows a powerful technique for establishing the hand, even at quite low resolution. We locate heavy atoms in the first derivative by Patterson methods, and choose one of the possible hands to use in computing SIR phases. Then, using the same hand assumption, we compute anomalous-scattering phases. For the second heavy-atom derivative, instead of using Patterson methods, we compute a difference Fourier between the native data and the second derivative data, using the SIR phases from the first derivative. Then we compute a second difference Fourier, adding the phases from anomalous scattering. Finally, we compute a third difference Fourier, just like the second except that the signs of all anomalous-scattering contributions are reversed, which is like assuming the opposite hand. The first Fourier should exhibit electron-density peaks at the positions of the second heavy atom. If the initial hand assumption was correct, heavy-atom peaks should be stronger in the second Fourier. If it was incorrect, heavy-atom peaks should be stronger in the third Fourier.

#### 6.4.7 Direct phasing: Application of methods from small-molecule crystallography

Methods involving heavy atoms apply almost exclusively to large molecules (500 or more atoms, not counting hydrogens). For small molecules (up to 200 atoms), phases can be determined by what are commonly called *direct methods*. One form of direct phasing relies on the existence of mathematical relationships among certain combinations of phases. From these relationships, a sufficient number of initial phase estimates can be obtained to begin converging toward a complete set of phases. One such relationship, called a triplet relationship, relates the phases and indices of three reflections as follows:

$$\alpha_{hkl} + \alpha_{h'k'l'} \simeq \alpha_{(h+h')(k+k')(l+l')} \tag{6.15}$$

In words, if the *indices* of the reflections on the left sum to the indices of the reflection on the right, then their *phases* sum approximately. As a specific example, the phase of reflection  $\overline{120}$  is always approximately the sum of the

phases of reflections 110 and  $\overline{2}10$ . If the number of reflections is not too large, so that a large percentage of these relationships can be examined simultaneously, they might put enough constraints on phases to produce some initial estimates. For macromolecules, the number of reflections is far too large to make this method useful.

Another form of direct phasing, executed by a program called Shake-and-Bake, in essence tries out random arrangements of atoms, simulates the diffraction patterns they would produce, and compares the simulated patterns with those obtained from the crystals. Even though the trial arrangements are limited to those that are physically possible (for example, having no two atoms closer than bonding or van der Waals forces allow), the number of trial arrangements can be too large for computation if the number of atoms is large. But this "try-everything" method is enormously powerful for any number of atoms that computation can handle, and of course, this number grows with the rapidly growing capacity of computers.

Direct methods work if the molecules, and thus the unit cells and numbers of reflections, are relatively small. Isomorphous replacement works if the molecules are large enough that a heavy atom does not disturb their structures significantly. The most difficult structures for crystallographers are those that are too large for direct methods and too small to remain isomorphous despite the intrusion of a heavy atom. If a medium-size protein naturally contains a heavy atom, like iron or zinc, or if a selenomethionine derivative can be produced, the structure can often be solved by MAD phasing (Sec. 6.4.5, p. 133). [NMR methods (see Chapter 10) are also of great power for small and medium-size molecules.]

The Shake-and-Bake style of direct phasing apparently has the potential to solve the structures of proteins of over 100 residues if they diffract exceptionally well (to around 1.0 Å). Fewer than 10% of large molecules diffract well enough to qualify. But in a combined process that shows great promise, direct phasing has been combined with MAD phasing to solve large structures. Recall that larger proteins may contain too many methionines to allow Patterson or least-squares location of all seleniums in the selenomethionine derivative. Solving the anomalous-diffraction phase equations requires knowing the locations of all the heavy atoms. Shake-and-Bake can solve this problem, even if there are 50 or more seleniums in the protein. One early success of this combined method was a protein of molecular mass over 250,000 containing 65 selenomethionines.

Our last phasing method applies to all molecules, regardless of size, but it requires knowledge that the desired structure is similar to a known structure.

# 6.5 Molecular replacement: Related proteins as phasing models

#### 6.5.1 Introduction

The crystallographer can sometimes use the phases from structure factors of a known protein as initial estimates of phases for a new protein. If this method is

feasible, then the crystallographer may be able to determine the structure of the new protein from a single native data set. The known protein in this case is referred to as a *phasing model*, and the method, which entails calculating initial phases by placing a model of the known protein in the unit cell of the new protein, is called *molecular replacement*.

For instance, the mammalian serine proteases—trypsin, chymotrypsin, and elastase—are very similar in structure and conformation. If a new mammalian serine protease is discovered, and sequence homology with known proteases suggests that this new protease is similar in structure to known ones, then one of the known proteases might be used as a phasing model for determining the structure of the new protein.

Similarly, having learned the crystallographic structure of a protein, we may want to study the conformational changes that occur when the protein binds to a small ligand and to learn the molecular details of protein-ligand binding. We might be able to crystallize the protein and ligand together or introduce the ligand into protein crystals by soaking. We expect that the protein-ligand complex is similar in structure to the free protein. If this expectation is realized, we do not have to work completely from scratch to determine the structure of the complex. We can use the ligand-free protein as a phasing model for the protein-ligand complex.

In Fig. 6.1, p. 110, I showed that phases contain more information than intensities. How, then, can the phases from a different protein help us find an unknown structure? In his Book of Fourier, Kevin Cowtan uses computed transforms to illustrate this concept, as shown in Fig. 6.17. First we see, posing this time as an unknown structure, the cat (a), with its Fourier transform shown in black and white. The colorless transform is analogous to an experimental diffraction pattern because we do not observe phases in experimental data. Next we see a Manx (tailless) cat (b) (along with its transform) posing as a solved structure that we also know (because of, say, sequence homology) to be similar in structure to the cat (a). If we know that the unknown structure, the cat, is similar to a known structure, the Manx cat, are the intensities of the cat powerful enough to reveal the differences between the unknown structure and the phasing model—in this case, the tail? In (c) the phases (colors) of the Manx cat transform are superimposed on the intensities of the unknown cat transform. In (d) we see the back-transform of (c), and although the image is weak, the cat's tail is apparent. The intensities of cat diffraction do indeed provide enough information to show how the cat differs from the Manx cat. In like manner, measured intensities from a protein of unknown structure do indeed have the power to show us how it differs from a similar, known structure used as a phasing model.

#### 6.5.2 Isomorphous phasing models

If the phasing model and the new protein (the target) are isomorphous, as may be the case when a small ligand is soaked into protein crystals, then the phases from the free protein can be used directly to compute  $\rho(x, y, z)$  from native intensities



**Figure 6.17** Structure determination by molecular replacement. (*a*) Unknown structure, cat, and its diffraction pattern (not colored, because phases are unknown). (*b*) Known structure and phasing model, Manx cat, and transform computed from the model (colored, because calculation of transform from a model tells us phases). (*c*) Manx-cat phases combined with unknown-cat intensities. (*d*) Back-transform of (*c*). Intensities contain enough information to reveal differences (the tail) between phasing model and unknown structure.

of the new protein [Eq. (6.16)]:

$$\rho(x, y, z) = \frac{1}{V} \sum_{h} \sum_{k} \sum_{l} \left| \mathbf{F}_{hkl}^{\text{target}} \right| e^{-2\pi i \left( hx + ky + lz - \alpha_{hkl}^{\text{model}} \right)}$$
(6.16)

In this Fourier synthesis, the amplitudes  $|\mathbf{F}_{hkl}^{\text{target}}|$  are obtained from the native intensities of the new protein, and the phases  $\alpha'^{\text{model}}$  are those of the phasing model. During the iterative process of phase improvement (Chapter 7), the phases should change from those of the model to those of the new protein or complex, revealing the desired structure. In Fig. 6.17, we not only knew that our phasing model (the Manx cat) was similar to the unknown (cat with tail), but we had the added advantage of knowing that its orientation was the same. Otherwise, its phases would not have revealed the unknown structure.

#### 6.5.3 Nonisomorphous phasing models

If the phasing model is not isomorphous with the target structure, the problem is more difficult. The phases of atomic structure factors, and hence of molecular structure factors, depend upon the location of atoms in the unit cell. In order to use a known protein as a phasing model, we must superimpose the structure of the model on the structure of the target protein in its unit cell and then calculate phases for the properly oriented model. In other words, we must find, in the new unit cell, the position and orientation of the phasing model that superimposes it on the target protein, and hence, that would give phases most like those of the target. Then we can calculate the structure factors of the properly positioned model and use the phases of these computed structure factors as initial estimates of the desired phases.

Without knowing the structure of the target protein, how can we copy the model into the unit cell with the proper orientation and position? From native data on the new protein, we can determine its unit-cell dimensions and symmetry. Clearly the phasing model must be placed in the unit cell with the same symmetry as the new protein. This places some constraints upon where to locate the model, but not enough to give useful estimates of phases. In theory, it should be possible to conduct a computer search of all orientations and positions of the model in the new unit cell. For each trial position and orientation, we would calculate the structure factors (called  $\mathbf{F}_{calc}$ ) of the model [Eq. (5.15)], and compare their amplitudes  $|\mathbf{F}_{calc}|$  with the measured amplitudes  $|\mathbf{F}_{obs}|$  obtained from diffraction intensities of the new protein. Finding the position and orientation that gives the best match, we would take the computed phases ( $\alpha_{calc}$ ) as the starting phases for structure determination of the new protein.

#### 6.5.4 Separate searches for orientation and location

In practice, the number of trial orientations and positions for the phasing model is enormous, so a brute-force search is impractical, even on the fastest computers (as of this writing, of course). The procedure is greatly simplified by separating

139

the search for the best orientation from the search for the best position. Further, it is possible to search for the best orientation independently of location by using the Patterson function.

If you consider the procedure for drawing a Patterson map from a known structure (Sec. 6.3.3, p. 124), you will see that the final map is independent of the position of the structure in the unit cell. No matter where you draw the "molecule," as long as you do not change its orientation (that is, as long as you do not rotate it within the unit cell), the Patterson map looks the same. On the other hand, if you rotate the structure in the unit cell, the Patterson map rotates around the origin, altering the arrangement of Patterson atoms in a single Patterson unit cell. This suggests that the Patterson map might provide a means of determining the best orientation of the model in the unit cell of the new protein.

If the model and the target protein are indeed similar, and if they are oriented in the same way in unit cells of the same dimensions and symmetry, they should give very similar Patterson maps. We might imagine a trial-and-error method in which we compute Patterson maps for various model orientations and compare them with the Patterson map of the desired protein. In this manner, we could find the best orientation of the model, and then use that single orientation in our search for the best *position* of the model, using the structure-factor approach outlined earlier. A systematic search for the best orientation entails superimposing the origins of the two three-dimensional Patterson maps, and then rotating the map of the phasing model through three angles, as shown in Fig. 6.18. First, let us imagine an orthogonal system of coordinates (x, y, z) established with a fixed relationship to the Patterson unit cell of the desired protein (Fig. 6.18a). Then imagine a system of spherical polar angles  $\phi$ ,  $\varphi$ , and  $\chi$  defined with respect to the orthogonal system such that  $\phi$  and  $\varphi$  give angles of rotation of a directing axis (blue), and  $\chi$  gives the angle of rotation around that axis. The Patterson cell of the phasing model (b, transparent cell containing colored Patterson peaks) is then rotated with respect to the target Patterson cell (black framework with black Patterson peaks) through appropriately small intervals of the angles  $\phi$ ,  $\varphi$ , and  $\chi$ , (b, c, d) while the correlation between Patterson maps for the two models is monitored (next section). We are searching for values of  $\phi$ ,  $\varphi$ , and  $\chi$  that will superimpose the Patterson peaks in the two models (e).

How much computing do we actually save by searching for orientation and location separately? The orientation of the model can be specified by three angles of rotation about orthogonal axes x, y, and z with their origins at the center of the model. Specifying location also requires three numbers, the x, y, and z coordinates of the molecular center with respect to the origin of the unit cell. For the sake of argument, let us say that we must try 100 different values for each of the six parameters. (In real situations, the number of trial values is much larger.) The number of combinations of six parameters, each with 100 possible values is  $100^6$ , or  $10^{12}$ . Finding the orientation as a separate search requires first trying 100 different values for each of three angles, which is  $100^3$  or  $10^6$  combinations. After finding the orientation, finding the location requires trying 100 different values of each of three coordinates, again  $100^3$  or  $10^6$  combinations.

141



**Figure 6.18**  $\triangleright$  Rotation search (b-e) finds the values of angles  $\phi$ ,  $\varphi$ , and  $\chi$  (*a*) that superimpose the Patterson map of the phasing model (transparent cell with colored Patterson peaks) on that of the target (black framework with black Patterson peaks). The orientation (*e*) that gives the highest correlation between the two Patterson maps gives the best orientation of the phasing model in the unit cell of the target.

The total number of trials for separate orientation and location searches is  $10^6 + 10^6$  or  $2 \times 10^6$ . The magnitude of the saving is  $10^{12}/2 \times 10^6$  or 500,000. In this case, the problem of finding the orientation and location separately is smaller by half a million times than the problem of searching for orientation and location simultaneously.

#### 6.5.5 Monitoring the search

Finally, what mathematical criteria are used in these searches? In other words, as the computer goes through sets of trial values (angles or coordinates) for

the model, how does it compare results and determine optimum values of the parameters?

Comparison of orientations is usually done by computing a *rotation function*, which evaluates the correlation between Patterson maps for the target protein and for the phasing model in various orientations. For this orientation search (often called a *rotation search*), the computer is looking for large values of the model Patterson function  $P^{\text{model}}(u, v, w)$  at locations corresponding to peaks in the Patterson maps over the Patterson unit cell for each orientation of the model Patterson with respect to the target Patterson. Where either Patterson has a peak and the other does not, the product is zero. Where the two Pattersons have coincident peaks, the product is large. So the integral of the product will be very large if there are many coincident peaks in the two maps, and a maximum value at the relative orientation of maximum overlap. For this type of search, the rotation function can be expressed as

$$R(\phi,\varphi,\chi) = \int_{u,v,w} P^{\text{target}}(u,v,w) P^{\text{model}}\{(u,v,w) \times [\phi,\varphi,\chi]\} du \, dv \, dw$$
(6.17)

In words, at each set of rotation angles  $\phi$ ,  $\varphi$ , and  $\chi$ , the value of the rotation function *R* is the integral of the product of two Patterson functions: (1) that of the target molecule [ $P^{\text{target}}(u, v, w$ )], and (2) that of the model [ $P^{\text{model}}(u, v, w$ )] with its coordinates (u, v, w) operated on by rotation matrix<sup>2</sup> [ $\phi$ ,  $\varphi$ ,  $\chi$ ] to produce a specific orientation relative to the target Patterson. This function will exhibit maxima where the two Pattersons have many coincident peaks. This maximum should tell us the best orientation for placing the phasing model in the unit cell of the desired protein (Fig. 6.18*e*). Near the maximum, the rotation search can be repeated at smaller angular intervals to refine the orientation.

For the location or translation search, the criterion is the correspondence between the expected structure-factor amplitudes from the model in a given trial location and the actual amplitudes derived from the native data on the desired protein. This criterion can be expressed as the *R*-factor, a parameter we will encounter later as a criterion of improvement of phases in final structure determination. The *R*factor compares overall agreement between the amplitudes of two sets of structure factors, as follows:

$$R = \frac{\sum ||\mathbf{F}_{obs}| - |\mathbf{F}_{calc}||}{\sum |\mathbf{F}_{obs}|}.$$
(6.18)

142

 $<sup>^{2}</sup>$ The operation of a rotation matrix on a set of coordinates produces a set of simultaneous equations whose solution is the new set of coordinates. For an example of such a set of equations, see Eq. (11.1), p. 271.

In words, for each reflection, we compute the difference between the observed structure-factor amplitude from the native data set  $|\mathbf{F}_{obs}|$  and the calculated amplitude from the model in its current trial location  $|\mathbf{F}_{calc}|$  and take the absolute value, giving the magnitude of the difference. We add these magnitudes for all reflections. Then we divide by the sum of the observed structure-factor amplitudes (the reflection intensities).

If, on the whole, the observed and calculated intensities agree with each other, the differences in the numerator are small, and the sum of the differences is small compared to the sum of the intensities themselves, so R is small. For perfect agreement, all the differences equal zero, and R equals zero. No single difference is likely to be larger than the corresponding  $|\mathbf{F}_{obs}|$ , so the maximum value of R is one. For proteins, R-values of 0.3 to 0.4 for the best placement of a phasing model have often provided adequate initial estimates of phases.

#### 6.5.6 Summary of molecular replacement

If we know that the structure of a new protein is similar to that of a known protein, we can use the known protein as a phasing model, and thus solve the phase problem without heavy atom derivatives. If the new crystals and those of the model are isomorphous, the model phases can be used directly as estimates of the desired phases. If not, we must somehow superimpose the known protein upon the new protein to create the best phasing model. We can do this without knowledge of the structure of the new protein by using Patterson-map comparisons to find the best orientation of the model protein and then using structure-factor comparisons to find the best location of the model protein.

6.6

# Iterative improvement of phases (preview of Chapter 7)

The phase problem greatly increases the effort required to obtain an interpretable electron-density map. In this chapter, I have discussed several methods of obtaining phases. In all cases, the phases obtained are estimates, and often the set of estimates is incomplete. Electron-density maps calculated from Eq. (6.7), p. 116, using measured amplitudes and first phase estimates, are often difficult or impossible to interpret. In Chapter 7, I will discuss improvement of phase estimates and extension of phase assignments to as many reflections as possible. As phase improvement and extension proceed, electron-density maps become clearer and easier to interpret as an image of a molecular model. The iterative process of *structure refinement* eventually leads to a structure that is in good agreement with the original data.

This Page Intentionally Left Blank

# ► Chapter 7

# Obtaining and Judging the Molecular Model

## 7.1 Introduction

In this chapter, I will discuss the final stages of structure determination: obtaining and improving the electron-density map that is based on the first phase estimates, interpreting the map to produce an atomic model of the unit-cell contents, and refining the model to optimize its agreement with the original native reflection intensities. The criteria by which the crystallographer judges the progress of the work overlap with criteria for assessing the quality of the final model. These subjects form the bridge from Chapter 7 to Chapter 8, where I will review many of the concepts of this book by guiding you through the experimental descriptions from the description of a structure determination as published in a scientific journal.

With each passing year, crystallography becomes more highly automated, and the methods of this chapter have probably been the most affected by automation. In routine structure determination, many of the judgments and choices that I will describe are now made by software, running sophisticated algorithms that reflect the collective experience and cleverness of the talented men and women who propel this challenging field. In many cases, it is possible to go from a first map based on MIR phases to a model that is more than 90% complete in just a few minutes of computing that alternates between fitting a model to the map and then computationally improving the map. But just as it helps to know a little about the innards of that car that refuses to start, it helps software users to know what is going on under the hoods of those quiet, sleek, crystallography programs. I will speak of the decisions inherent in structure refinement as if they were made by a prudent and discerning human being. Many of the routine decisions can indeed be handed over to software, and lots of drudgery thus evaded. But there are difficult cases and

crucial mop-up work that software cannot handle, requiring the crystallographer to know what the program is trying to do, how it does it, and what the options are when progress stalls.

# 7.2 Iterative improvement of maps and models—overview

In brief, obtaining a detailed molecular model of the unit-cell contents entails calculating  $\rho(x, y, z)$  from Eq. (6.7),

$$\rho(x, y, z) = \frac{1}{V} \sum_{h} \sum_{k} \sum_{l} |F_{hkl}| e^{-2\pi i (hx + ky + lz - \alpha'_{hkl})}.$$
(6.7)

using amplitudes ( $|F_{hkl}|$ ) computed from measured intensities in the native data set and phases ( $\alpha'_{hkl}$ ) computed from heavy-atom data, anomalous scattering, or molecular replacement. Because the phases are rough estimates, the first map may be uninformative and disappointing. Crystallographers improve the map by an iterative process sometimes called *bootstrapping*. The basic principle of this iteration is easy to state but demands care, judgment, and much labor to execute. The principle is the following: any features that can be reliably discerned in, or inferred from, the map become part of a phasing model for subsequent maps. Without the input of new information, the map will not improve.

Whatever crude model of unit cell contents that can be discerned in the map is cast in the form of a simple electron-density function and used to calculate new structure factors by Eq. (5.16):

$$F_{hkl} = \int_{x} \int_{y} \int_{z} \rho(x, y, z) e^{2\pi i (hx + ky + lz)} \, dx \, dy \, dz, \tag{5.16}$$

The phases of these structure factors are used, along with the original native intensities, to add more terms to Eq. (6.7), the Fourier-sum description of  $\rho(x, y, z)$ , in hopes of producing a clearer map. When the map becomes clear enough to allow location of atoms, these are added to the model, and structure factors are computed from this model using Eq. (5.15),

$$F_{hkl} = \sum_{j=1}^{n} f_j e^{2\pi i (hx_j + ky_j + lz_j)}.$$
(5.15)

which contains *atomic* structure factors rather than electron density. As the model becomes more detailed, the phases computed from it improve, and the model, computed from the original native structure-factor amplitudes and the latest phases,

becomes even more detailed. The crystallographer thus tries to bootstrap from the initial rough phase estimates to phases of high accuracy, and from them, a clear, interpretable map and a model that fits the map well.

I should emphasize that the crystallographer cannot get any new phase information without modifying the model in some way. The possible modifications include solvent flattening, noncrystallographic symmetry averaging, or introducing a partial atomic model, all of which are discussed further in this chapter. It is considered best practice to make sure that initial phases are good enough to make the map interpretable. If it is not, then the crystallographer needs to find additional derivatives and collect better data.

The model can be improved in another way: by refinement of the atomic coordinates. This method entails adjusting the atomic coordinates to improve the agreement between amplitudes calculated from the current model and the original measured amplitudes in the native data set. In the latter stages of structure determination, the crystallographer alternates between map interpretation and refinement of the coordinates. At its heart, refinement is an attempt to minimize the differences between (a) measured diffraction intensities and (b) intensities predicted by the current model, which, in intermediate stages, is incomplete and harbors errors that will eventually be removed. Classical refinement algorithms employ the statistical method of least-squares (Sec. 7.5.1, p. 159), but newer methods, including energy refinement (Sec. 7.5.5, p. 163) and methods based on Bayesian statistics (Sec. 7.5.6, p. 164), are allowing greater automation, as well as more effective extraction of structural information in difficult cases.

The block diagram of Fig. 7.1 shows how these various methods ultimately produce a molecular model that agrees with the native data. The vertical dotted line in Fig. 7.1 divides the operations into two categories. To the right of the line are real-space methods, which entail attempts to improve the electron-density map, by adding information to the map or removing noise from it, or to improve the model, using the map as a guide. To the left of the line are reciprocal-space methods, which entail attempts to improve the agreement between reflection intensities computed from the model and the original measured reflection intensities. In real-space methods, the criteria for improvement or removal of errors are found in electron-density maps, in the fit of model to map, or in the adherence of the model to prior knowledge, such as expected bond lengths and angles (all real-space criteria); in reciprocal-space methods, the criteria for improvement or removal of errors involve reliability of phases and agreement of calculated structure factors with measured intensities (all reciprocal-space criteria). The link between real and reciprocal space is, of course, the Fourier transform (FT).

I will return to this diagram near the end of the chapter, particularly to amplify the meaning of *error removal*, which is indicated by dashed horizontal lines in Fig. 7.1. For now, I will illustrate the bootstrapping technique for improving phases, map, and model with an analogy: the method of successive approximations for solving a complicated algebraic equation. Most mathematics education emphasizes equations that can be solved analytically for specific variables. Many realistic problems defy such analytic solutions but are amenable to numerical methods. The method



**Figure 7.1**  $\triangleright$  Block diagram of crystallographic structure determination. Operations on the left are based on reciprocal space criteria (improving phases), while those on the right are based in real space (improving the accuracy of atomic coordinates).

of successive approximations has much in common with the iterative process that extracts a protein model from diffraction data.

Consider the problem of solving the following equation for the variable *y*:

$$\left(1 + \frac{1}{y^2}\right)(y-1) = 1.$$
(7.1)

Attempts to simplify the equation produce a cubic equation in y, giving no straightforward means to a numerical solution. You can, however, easily obtain a numerical solution for y with a hand calculator. Start by solving for y in terms of

 $y^2$  as follows:

$$y = \frac{1}{\left(1 + \frac{1}{y^2}\right)} + 1. \tag{7.2}$$

Then make an arbitrary initial estimate of y, say y = 1. (This is analogous to starting with the MIR phases as initial estimates of the correct phases.) Plug this estimate into the right-hand  $y^2$  term, and calculate y [analogous to computing a crude structure from measured structure-factor amplitudes ( $|\mathbf{F}_{obs}|$ ) and phase estimates]. The result is 1.5. Now take this computed result as the next estimate (analogous to computing new structure factors from the crude structure), plug it into the  $y^2$  term, and compute y again (analogous to computing a new structure from better phase estimates). The result is 1.6923. Repeating this process produces these answers in succession: 1.7412, 1.752, 1.7543, 1.7547, 1.7549, 1.7549, and so on. After a few iterations, the process converges to a solution; that is, the output value of y is the same as the input. This value is a solution to the original equation.

With Eq. (7.2), any first estimate above 1.0 (even one million) produces the result shown. In contrast, for many other equations, the method of successive approximations works only if the initial estimate is close to a correct solution. Otherwise, the successive answers do not converge; instead, they may oscillate among several values (the iteration "hangs up" instead of converging), or they may continually become larger in magnitude (the iteration "blows up"). In order for the far more complex crystallographic iteration to converge to a protein model that is consistent with the diffraction data, initial estimates of many phases must be close to the correct values. Attempts to start from random phases in hopes of convergence to correct ones appear doomed to failure because of the large number of incorrect solutions to which the process can converge.

The following sections describe the crystallographic bootstrapping process in more detail.

### 7.3 First maps

#### 7.3.1 Resources for the first map

Entering the final stages of structure determination, the crystallographer is armed with several sets of data with which to calculate  $\rho(x, y, z)$  as a Fourier sum of structure factors using Eq. (6.7). First is the original native data set, which usually contains the most accurate and complete (highest-resolution) set of measured intensities. These data will support the most critical tests of the final molecular model. Next are data sets from heavy-atom derivatives, which are often limited to lower resolution. Several sets of phases may be available, calculated from heavyatom derivatives and perhaps anomalous dispersion. Because each phase must be calculated from a heavy-atom reflection, phase estimates are not available for

#### Chapter 7 Obtaining and Judging the Molecular Model

native reflections at higher resolution than that of the best heavy-atom derivative. Finally, for each set of phases, there is usually some criterion of precision. These criteria will be used as weighting factors, numbers between 0 and 1, for Fourier terms containing the phases. A Fourier term containing a phase estimate of low reliability (see Fig. 6.11) will be multiplied by a low weighting factor in the Fourier-sum computation of  $\rho(x, y, z)$ . In other words, such a term will be multiplied by a number less than 1.0 to reduce its contribution to the Fourier sum, and thus reduce bias from a reflection whose phase is questionable. Conversely, a term containing a phase of high reliability will be given full weight (weighting factor of 1.0) in the sum.

Here is the Fourier sum that gives the first electron-density map:

$$\rho(x, y, z) = \frac{1}{V} \sum_{h} \sum_{k} \sum_{l} w_{hkl} |\mathbf{F}_{obs}| e^{-2\pi i \left(hx + ky + lz - \alpha'_{calc}\right)}.$$
 (7.3)

In words, the desired electron-density function is a Fourier sum in which term *hkl* has amplitude  $|\mathbf{F}_{obs}|$ , which equals  $(I_{hkl})^{1/2}$ , the square root of the measured intensity  $I_{hkl}$  from the native data set. The phase  $\alpha'_{hkl}$  of the same term is calculated from heavy-atom, anomalous dispersion, or molecular replacement data, as described in Chapter 6. The term is weighted by the factor  $w_{hkl}$ , which will be near 1.0 if  $\alpha'_{hkl}$  is among the most highly reliable phases, or smaller if the phase is questionable. This Fourier sum is called an  $\mathbf{F}_{obs}$  or  $\mathbf{F}_{o}$  synthesis (and the map an  $\mathbf{F}_{o}$  map) because the amplitude of each term *hkl* is  $|\mathbf{F}_{obs}|$  for reflection *hkl*.

The first term in this Fourier series, the  $\mathbf{F}_{000}$  term, should contain  $(I_{000})^{1/2}$ , where  $I_{000}$  is the intensity of reflection 000, which lies at the origin of the reciprocal lattice. Recall that this reflection is never measured because it is obscured by the direct beam. Examination of Eq. (7.3) reveals that  $\mathbf{F}_{000}$  is a real constant (as opposed to a complex or imaginary number). The phase  $\alpha'_{000}$  of this term is assigned a value of zero, with the result that all other phases will be computed relative to this assignment. Then because h = k = l = 0 for reflection 000, the exponent of *e* is zero and the entire exponential term is 1.0. Thus  $\mathbf{F}_{000}$  is a constant, just like  $\mathbf{f}_0$  in Fig. 2.16, p. 25.

All other terms in the sum are simple trigonometric functions with average values of zero, so it is clear that the value assigned to  $\mathbf{F}_{000}$  will determine the overall amplitude of the electron-density map. (In the same manner, the  $\mathbf{f}_0$  term in Fig. 2.16 displaces all the Fourier sums upward, making the sums positive for all values of *x*, like the target function.) The sensible assignment for  $\mathbf{F}_{000}$  is therefore the total number of electrons in the unit cell, making the sum of  $\rho(x, y, z)$  over the whole unit cell equal to the total electron density. In practice, this term can be omitted from the calculation, and the overall map amplitude can be set by means described in Sec. 7.3.3, p. 151.

#### 7.3.2 Displaying and examining the map

Until the middle to late 1980s, the contour map of the first calculated electron density was displayed by printing sections of the unit cell onto Plexiglas or clear

#### Section 7.3 First maps

plastic sheets and stacking them to produce a three-dimensional model, called a *minimap*. Today's computers can display the equivalent of a minimap and allow much more informative first glimpses of the electron density. With computer displays, the contour level of the map (see p. 94) can be adjusted for maximum detail. These first glimpses of the molecular image are often attended with great excitement and expectation. If the phase estimates are sufficiently good, the map will show some of the gross features of unit cell contents. In the rare best cases, with good phases from molecular replacement, and perhaps with enhancement from noncrystallographic averaging (explained further in Sec. 7.3.3, p. 151), first maps are easily interpretable, clearly showing continuous chains of electron density and features like alpha helices-perhaps even allowing some amino-acid side chains to be identified. At the worst, the first map is singularly uninformative, signaling the need for additional phasing information, perhaps from another heavy-atom derivative. Usually the minimum result that promises a structure from the existing data is that protein be distinguishable from bulk water, and that dense features like alpha helices can be recognized. If the boundary of each molecule, the molecular envelope, shows some evidence of recognizable protein structure, then a full structure is likely to come forth.

I will consider the latter case, in which the first map defines a molecular envelope, with perhaps a little additional detail. If more detail can be discerned, the crystallographer can jump ahead to later stages of the map-improvement process I am about to describe. If the molecular envelope cannot be discerned, then the crystallographer must collect more data.

#### 7.3.3 Improving the map

The crude molecular image seen in the  $\mathbf{F}_0$  map, which is obtained from the original indexed intensity data ( $|\mathbf{F}_{obs}|$ ) and the first phase estimates ( $\alpha'_{calc}$ ), serves now as a model of the desired structure. A crude electron density function is devised to describe the unit-cell contents as well as they can be observed in the first map. Then the function is modified to make it more realistic in the light of known properties of proteins and water in crystals. This process is called, depending on the exact details of procedure, *density modification, solvent leveling, solvent flattening*, or *solvent flipping*.

The electron density function devised by density modification may be no more than a fixed, high value of  $\rho(x, y, z)$  for all regions that appear to be within a protein molecule, and a fixed, low value of  $\rho$  for all surrounding areas of bulk solvent. One automated method first defines the molecular envelope by dividing the unit cell into a grid of regularly spaced points. At each point, the value of  $\rho(x, y, z)$  in the  $\mathbf{F}_0$  map is evaluated. At each grid point, if  $\rho$  is negative, it is reassigned a value of zero; if  $\rho$  is positive, it is assigned a value equal to the average value of  $\rho$  within a defined distance of the gridpoint. This procedure smooths the map (eliminates many small, random fluctuations in density) and essentially divides the map into two types of regions: those of relatively high (protein) and relatively low (solvent) density. Numerous variations of this method are in use. In solvent flipping, the density corresponding to solvent is inverted rather than flattened, to enhance contrast between solvent and protein, and to counteract bias in the model. Some software includes a graphical display that allows the user to draw a mask to define this boundary. Next, the overall amplitude of the map is increased until the ratio of high density to low density agrees with the ratio of protein to solvent in the crystal, usually assuming that the crystal is about half water. This contrived function  $\rho(x, y, z)$  is now used to compute structure factors, using Eq. (5.16). From this computation, we learn what the amplitudes and phases of all reflections would be if this very crude new model were correct. We use the phases from this computation, which constitute a new set of  $\alpha'_{hkl}$ s, along with the  $|\mathbf{F}_{obs}|$ s derived from the original measured intensities, to calculate  $\rho(x, y, z)$ again, using Eq. (7.3).

We do not throw out old phases immediately, but continue to weight each Fourier term with some measure of phase quality (sometimes called a *figure of merit*). In this manner, we continue to let the data speak for itself as much as possible, rather than allowing the current model to bias the results. If the new phase estimates are better, then the new  $\rho(x, y, z)$  will be improved, and the electron-density map will be more detailed. The new map serves to define the molecular boundary more precisely, and the cycle is repeated. (Refer again to the block diagram in Fig. 7.1, p. 148.) If we continue to use good judgment in incorporating new phases and new terms into Eq. (7.3), successive Fourier-series computations of  $\rho(x, y, z)$  include more terms, and successive contour maps become clearer and more interpretable. In other words, the iterative process of incorporating phases from successively better and more complete models converges toward a structure that fits the native data better. The phase estimates "converge" in the sense that the output phases that went into computation of the model [Eq. (7.3)].

As this process continues, and the model becomes more detailed, we begin to get estimates for the phases of structure factors at resolution beyond that of the heavyatom derivatives. In a process called *phase extension*, we gradually increase the number of terms in the Fourier sum of Eq. (7.3), adding terms that contain native intensities (as  $|\mathbf{F}_{obs}|$ ) at slightly higher resolution with phases from the current model. This must be done gradually and judiciously, so as not to let incorrect areas of the current model bias the calculations excessively. If the new phase estimates are good, the resulting map has slightly higher resolution, and structure factors computed from Eq. (5.16) give useful phase estimates at still higher resolution. In this manner, low-resolution phases are improved, and phase assignments are extended to higher resolution.

If phase extension seems like getting something from nothing, realize that by using general knowledge about protein and solvent density, we *impose justifiable restrictions* (sometimes referred to as *prior knowledge*) on the model, giving it realistic properties that are not visible in the map. In effect, we are using known crystal properties to increase the resolution of the model. Thus it is not surprising that the phases calculated from the modified model are good to higher resolution than those calculated from an electron-density function that does little more than describe what can be seen in the map.

#### Section 7.4 The Model becomes molecular

Another means of improving the map at this stage depends upon the presence of noncrystallographic symmetry elements in the unit cell. Recall that the intensity of reflections results from many molecules in identical orientations diffracting identically. In a sense, the diffraction pattern is the sum of diffraction patterns from all individual molecules. This is equivalent to taking a large number of weak, noisy signals (each the diffraction from one molecule) and adding them together to produce a strong signal. The noise in the individual signals, which might include the background intensity of the film or the weak signal of stray X-rays, is random, and when many weak signals are added, this random noise cancels out.

In some cases, the strength of this signal can be increased further by averaging the signals from molecules that are identical but have different orientations in the unit cell, such that no two orientations of the crystal during data collection gives the same orientation of these molecules in the X-ray beam. These molecules may be related by symmetry elements that are not aligned with symmetry elements of the entire unit cell. Thus the diffractive contributions of these identical molecules are never added together. In such cases, the unit cell is said to exhibit noncrystallographic symmetry (NCS). By knowing the arrangement of molecules in the unit cell-that is, by knowing the location and type of noncrystallographic symmetry elements—the crystallographer can use a computer to simulate the movement of these sets of molecules into identical orientations and thus add their signals together. The result is improved signal-to-noise ratio and, in the end, a clearer image of the molecules. This method, called *symmetry averaging*, is spectacularly successful in systems with high symmetry, such as viruses. Many virus coat proteins are icosahedral, possessing two-, three-, and five-fold rotation axes. Often one or more two- and threefold axes are noncrystallographic, and fivefold axes are always noncrystallographic, because no unit cell exhibits fivefold symmetry.

Finding noncrystallographic symmetry elements is another application of rotation searches (Sec. 6.5.4, p. 139, and Fig. 6.18, p. 141) and rotation functions (Sec. 6.5.5, p. 141, and Eq. (6.17), p. 142) using Patterson maps. In this case, the target and the model Pattersons are the same, and the rotation function is called a *self-rotation function*. In carrying out this rotation search, we are asking whether the Patterson map of the unit cell is superimposable on itself in a different orientation. The two or more orientations that superimpose the Patterson on itself reveal the orientations of symmetry axes that relate unit-cell contents, but that do not apply to the unit cell as a whole. These axes then become the basis of symmetry averaging.

## 7.4 The Model becomes molecular

#### 7.4.1 New phases from the molecular model

At some critical point in the iterative improvement of phases, the map becomes clear enough that we can trace the protein chain through it. In the worst
circumstances at this stage, we may only be able to see some continuous tubes of density. Various aids may be used at this point to help the viewer trace the protein chain through the map. One is to *skeletonize* the map, which means to draw line segments along lines of maximum density. These so-called ridge lines show the viewer rough lines along which the molecular chain is likely to lie. They can help to locate both the main chain and the branch points of side chains.

In a clearer map, we may be able to recognize alpha helices, one of the densest features of a protein, or sheets of beta structure. Now we can construct a partial molecular model (as opposed to an electron-density model) of the protein, using computer graphics to build and manipulate a stick model of the known sequence within small sections of the map. This technique is called *map fitting*, and is discussed later. From the resulting model, which may harbor many errors and undefined regions, we again calculate structure factors, this time using Eq. (5.15), which treats each atom in the current model as an independent scatterer. In other words, we calculate new structure factors from our current, usually crude, molecular model rather than from an approximation of  $\rho(x, y, z)$ . Additional iterations may improve the map further, allowing more features to be constructed therein.

Here again, as in density modification, we are using *prior knowledge*—known properties of proteins—to improve the model beyond what we can actually see in the map. Thus we are in effect improving the resolution of the model by making it structurally realistic: giving it local electron densities corresponding to the light atoms we know are present and connecting atoms at bond lengths and angles that we know must be correct. So again, our successive models give us phases for reflections at higher and higher resolution. Electron-density maps computed from these phases and, as always, the original native amplitudes  $|\mathbf{F}_{obs}|$  become more and more detailed. As the map becomes clearer, even more specific prior knowledge of protein structure, can be applied. For example, in pleated sheets, successive carbonyl C=O bonds point in opposite directions. Side chains in pleated sheets also alternate directions, and are roughly perpendicular to the carbonyls. Near the ends of pleated sheet strands, if electron density is weakening, this knowledge of pleated-sheet geometry can often guide the placement and orientation of extra residues that are not as well-defined by the current map.

### 7.4.2 Minimizing bias from the model

Conversion to a molecular model greatly increases the hazard of introducing excessive bias from the model into  $\rho(x, y, z)$ . At this point, bias can be decreased by one of several alternative Fourier computations of the electron-density map. As phases from the model begin to be the most reliable, they begin to dominate the Fourier sum. In the extreme, the series would contain amplitudes purely from the intensity data and phases purely from the model. In order to compensate for the increased influence of model phases, and to continue letting the intensity data influence improvement of the model, the crystallographer calculates electron-density maps using various difference Fourier syntheses, in which the amplitude of each term is of the form ( $|n |\mathbf{F}_{obs}| - |\mathbf{F}_{calc}||$ ), which reduces overall model influence by subtracting the calculated structure-factor amplitudes ( $|\mathbf{F}_{calc}|$ ) from some multiple of

the observed amplitudes ( $|\mathbf{F}_{obs}|$ ) within each Fourier term. For n = 1, the Fourier series is called an  $\mathbf{F}_{o} - \mathbf{F}_{c}$  synthesis:

$$\rho(x, y, z) = \frac{1}{V} \sum_{h} \sum_{k} \sum_{l} (|\mathbf{F}_{o}| - |\mathbf{F}_{c}|) e^{-2\pi i \left(hx + ky + lz - \alpha'_{calc}\right)}.$$
 (7.4)

A contour map of this Fourier series is called an  $\mathbf{F}_o - \mathbf{F}_c$  map. How is this map interpreted? Depending on which of  $\mathbf{F}_o$  or  $\mathbf{F}_c$  is larger, Fourier terms can be either positive or negative. The resulting electron-density map contains both positive and negative "density." Positive density in a region of the map implies that the contribution of the observed intensities ( $\mathbf{F}_o$ s) to  $\rho$  are larger than the contribution of the model ( $\mathbf{F}_c$ s), and thus that the unit cell (represented by  $\mathbf{F}_o$ s) contains more electron density in this region than implied by the model (represented by  $\mathbf{F}_c$ s). In other words, the map is telling us that the model should be adjusted to increase the electron density in this region, by moving atoms toward the region. On the other hand, a region of negative density indicates that the model implies more electron density in the region than the unit cell actually contains. The region of negative density is telling us to move atoms away from this region. As an example, if an amino-acid side chain in the model is in the wrong conformation, the  $\mathbf{F}_o - \mathbf{F}_c$  map may exhibit a negative peak coincident with the erroneous model side chain and a nearby positive peak signifying the correct position.

The  $\mathbf{F}_{o} - \mathbf{F}_{c}$  map emphasizes errors in the current model. In effect, it removes the influence of the current model, so that the original data can tell us where the model is wrong. But the  $\mathbf{F}_{o} - \mathbf{F}_{c}$  map lacks the familiar appearance of the molecular surface found in an  $\mathbf{F}_{o}$  map. In addition, if the model still contains many errors, the  $\mathbf{F}_{o} - \mathbf{F}_{c}$  map is "noisy," full of small positive and negative peaks that are difficult to interpret. The  $\mathbf{F}_{o} - \mathbf{F}_{c}$  map is most useful near the end of the structure determination, when most of the model errors have been eliminated. The  $\mathbf{F}_{o} - \mathbf{F}_{c}$ map is a great aid in detecting subtle errors after most of the serious errors are corrected.

A more easily interpreted and intuitively satisfying difference map, but one that still allows undue influence by the model to be detected, is the  $2\mathbf{F}_o - \mathbf{F}_c$  map, calculated as follows

$$\rho(x, y, z) = \frac{1}{V} \sum_{h} \sum_{k} \sum_{l} (2 |\mathbf{F}_{0}| - |\mathbf{F}_{c}|) e^{-2\pi i (hx + ky + lz - \alpha'_{calc})}.$$
 (7.5)

In this map, the model influence is reduced, but not as severely as with  $\mathbf{F}_{o} - \mathbf{F}_{c}$ . Unless the model contains extremely serious errors, this map is everywhere positive, and contours at carefully chosen electron densities resemble a molecular surface. With experience, the crystallographer can often see the bias of an incorrect area of the model superimposed upon the the signal of the correct structure as implied by the original intensity data. For instance, in a well-refined map (see Sec. 7.5, p. 145), backbone carbonyl oxygens are found under a distinct rounded bulge in the backbone electron density. If a carbonyl oxygen in the model is pointing 180° away from the actual position in the molecule, the bulge in the map may be weaker than usual, or misshapen (sometimes cylindrical) and a weak bulge may be visible on the opposite side of the carbonyl carbon, at the true oxygen position. Correcting the oxygen orientation in the model, and then recalculating structure factors, results in loss of the weak, incorrect bulge in the map and intensification of the bulge in the correct location. This may sound like a serious correction of the model, requiring the movement of many atoms, but the entire peptide bond can be flipped 180° around the backbone axis with only slight changes in the positions of neighboring atoms. Such peptide-flip errors are common in early atomic models.

Various other Fourier syntheses are used during these stages in order to improve the model. Some crystallographers prefer a  $3\mathbf{F}_{o} - 2\mathbf{F}_{c}$  map, a compromise between  $\mathbf{F}_{o} - \mathbf{F}_{c}$  and  $2\mathbf{F}_{o} - \mathbf{F}_{c}$ , for the final interpretation. In areas where the maps continue to be ambiguous, it is often helpful to examine the original MIR or molecular replacement maps for insight into how model building in this area might be started off on a different foot. Another measure is to eliminate the atoms in the questionable region and calculate structure factors from Eq. (5.15), so that the possible errors in the region contribute nothing to the phases, and hence do not bias the resulting map, which is called an *omit map*.

Rigorous analysis of sources and distributions of errors in atomic coordinates, intensities, and phases, and of how these introduce bias into maps, has led to the wide use of *figure-of-merit weighted maps*, often called *sigma-A* weighted maps. The synthesis equivalent to  $\mathbf{F}_{o} - \mathbf{F}_{c}$  in this case is  $m\mathbf{F}_{o} - D\mathbf{F}_{c}$ , where *m* is a figure of merit for each model phase, and *D* is an overall estimate of atomic-coordinate errors in the current model. In some treatments of this synthesis, *D* is replaced by a the term  $\sigma_{A}$  (hence the term sigma-A weighted), which is *D* corrected for completeness of the current model. Comparative tests of  $\sigma_{A}$ -weighted maps versus other types indicate that they are quite powerful at revealing molecular features that are in conflict with the current model, and such maps have come into wide use.

Another important type of difference Fourier synthesis, which is used to compare similar protein structures, is discussed in Sec. 8.3.3, p. 198.

### 7.4.3 Map fitting

Conversion to a molecular model is usually done as soon as the map reveals recognizable structural features. This procedure, called *map fitting* or *model building*, entails interpreting the electron-density map by building a molecular model that fits realistically into the molecular surface implied by the map. Map fitting is done by interactive computer graphics, although automation has eliminated much of the manual labor of this process. A computer program produces a realistic threedimensional display of small sections of one or more electron-density maps, and allows the user to construct and manipulate molecular models to fit the map. As mentioned earlier, such programs can draw ridge lines to help the viewer trace the chain through areas of weak density. Rapid and impressive advances have been

### Section 7.4 The Model becomes molecular

made in automated model building, starting from the first electron density maps that show sufficient detail. As before, I will describe this process as if under manual control in order to give insights into the types of problems the software is solving. Some of these decisions are still necessary in areas of the map where the software is unable to make unequivocal decisions.

As the model is built, the viewer sees the model within the map, as shown in Fig. 2.3*b*, *p*. 11. As the model is constructed or adjusted, the program stores current atom locations in the form of three-dimensional coordinates. The crystallographer, while building a model interactively on the computer screen, is actually building a list of atoms, each with a set of coordinates (x, y, z) to specify its location. Software updates the atomic coordinates whenever the model is adjusted. This list of coordinates is the output file from the map-fitting program and the input file for calculation of new structure factors. When the model is correct and complete, this file becomes the means by which the model is shared with the community of scientists who study proteins (Sec. 7.7, p. 173). This list of coordinates *is* the desired model, in a form that is convenient for further study and analysis.

In addition to routine commands for inserting or changing amino-acid residues, moving atoms and fragments, and changing conformations, map-fitting programs contain many sophisticated tools to aid the model builder. Fragments, treated as rigid assemblies of atoms, can be automatically fitted to the map by the method of least squares (see Sec. 7.5.1, p. 159). After manual adjustments of the model, which may result in unrealistic bond lengths and angles, portions of the model can be *regularized*, which entails automatic correction of bond lengths and angles with minimal movement of atoms. In effect, regularization looks for the most realistic configuration of the known sequence cannot be easily fitted to the map, some map-fitting programs can search fragment databases or the Protein Data Bank (see Sec. 7.7, p. 173) for fragments having the same sequence, and then display these fragments so the user can see whether they fit the map.

Following is a somewhat idealized description of how map fitting may proceed, illustrated with views from a modern map-fitting program. The maps and models are from the structure determination of adipocyte lipid binding protein (ALBP), which I will discuss further in Chapter 8.

When the map has been improved to the point that molecular features are revealed, the crystallographer attempts to trace the protein through as much continuous density as possible. At this point, the quality of the map will vary from place to place, being perhaps quite clear in the molecular interior, which is usually more ordered, and exhibiting broken density in some places, particularly at chain termini and surface loops. Because we know that amino-acid side chains branch regularly off alpha carbons in the main chain, we can estimate the positions of many alpha carbons. These atoms should lie near the center of the main-chain density next to bulges that represent side chains. In proteins, alpha carbons are 3.8–4.2 Å apart. This knowledge allows the crystallographer to construct an alpha-carbon model of the molecule (Fig. 7.2) and to compute structure factors from this model.



**Figure 7.2** Alpha-carbon model (stereo) of ALBP built into electron-density map of Fig. 2.3a, p. 11.

Further improvement of the map with these phases may reveal side chains more clearly. Now the trick is to identify some specific side chains so that the known amino-acid sequence of the protein can be aligned with visible features in the map. As mentioned earlier, chain termini are often ill-defined, so we need a foothold for alignment of sequence with map where the map is sharp. Many times the key is a short stretch of sequence containing several bulky hydrophobic residues, like Trp, Phe, and Tyr. Because they are hydrophobic, they are likely to be in the interior where the map is clearer. Because they are bulky, their side-chain density is more likely to be identifiable. From such a foothold, the detailed model building can begin.

Regions that cannot be aligned with sequence are often built with polyalanine, reflecting our knowledge that all amino acids contain the same backbone atoms, and all but one, glycine, have at least a beta carbon (Fig. 7.3). In this manner,



**Figure 7.3** ► Polyalanine model (stereo) of ALBP built into electron-density map of Fig. 2.3.

#### Section 7.5 Structure refinement

we build as many atoms into the model as possible in the face of our ignorance about how to align the sequence with the map in certain areas.

In pleated sheets, we know that successive carbonyl oxygens point in opposite directions. One or two carbonyls whose orientations are clearly revealed by the map can allow sensible guesses as to the positions of others within the same sheet. As mentioned previously with respect to map fitting, we use prior knowledge of protein structure to infer more than the map shows us. If our inferences are correct, subsequent maps, computed with phases calculated from the model, will show enhanced evidence for the inferred features and will show additional features as well, leading to further improvement of the model. Poor inferences degrade the map, so where electron density conflicts with intuition, we follow the density as closely as possible.

With each successive map, new molecular features are added as soon as they can be discerned, and errors in the model, such as side-chain conformations that no longer fit the electron density, are corrected. As the structure nears completion, the crystallographer may simultaneously use maps based on various Fourier syntheses in order to track down the most subtle disagreements between the model and the data.

### 7.5 Structure refinement

### 7.5.1 Least-squares methods

Cycles of map calculation and model building, which are forms of realspace refinement of the model, are interspersed with computerized attempts to improve the agreement of the model with the original intensity data. (Everything goes back to those original reflection intensities, which give us our  $|\mathbf{F}_{obs}|$  values.) Because these computations entail comparison of computed with observed structure factor amplitudes (reciprocal space), rather than examination of maps and models (real space), these methods are referred to as *reciprocal space refinement*. The earliest successful refinement technique was a massive version of least-squares fitting, the same procedure that freshman chemistry students employ to construct a straight line that fits a scatter graph of data. More recently other versions, energy refinement and methods based on Bayesian statistics, have come to the fore. I will discuss least-squares methods first, and then contrast them with more modern methods.

In the simple least-squares method for functions of one variable, the aim is to find a function y = f(x) that fits a series of observations  $(x_1, y_1), (x_2, y_2), \dots, (x_i, y_i)$ , where each observation is a data point—a measured value of the independent variable x at some selected value y. (For example, y might be the temperature of a gas, and x might be its measured pressure.) The solution to the problem is a function f(x) for which the sum of the squares of distances between the data points and the function itself is as small as possible. In other words, f(x) is the function that minimizes *D*, the sum of the squared differences between observed  $(y_i)$  and calculated  $f(x_i)$  values, as follows:

$$D = \sum_{i} w_i \left[ y_i - f(x_i) \right]^2.$$
(7.6)

The differences are squared to make them all positive; otherwise, for a large number of random differences, D would simply equal zero. The term  $w_i$  is an optional weighting factor that reflects the reliability of observation *i*, thus giving greater influence to the most reliable data. According to principles of statistics,  $w_i$  should be  $1/(\sigma_i)^2$ , where  $\sigma_i$  is the standard deviation computed from multiple measurements of the same data point  $(x_i, y_i)$ . In the simplest case, f(x) is a straight line, for which the general equation is f(x) = mx + b, where m is the slope of the line and b is the intercept of the line on the f(x)-axis. Solving this problem entails finding the proper values of the parameters m and b. If we substitute  $(m_i + b)$ for each  $f(x_i)$  in Eq. (7.6), take the partial derivative of the right-hand side with respect to m and set it equal to zero, and then take the partial derivative with respect to b and set it equal to zero, the result is a set of simultaneous equations in *m* and *b*. Because all the squared differences are to be minimized simultaneously, the number of equations equals the number of observations, and there must be at least two observations to fix values for the two parameters m and b. With just two observations  $(x_1, y_1)$  and  $(x_2, y_2)$ , m and b are determined precisely, and f(x)is the equation of the straight line between  $(x_1, y_1)$  and  $(x_2, y_2)$ . If there are more than two observations, the problem is overdetermined and the values of m and b describe the straight line of best fit to all the observations. So the solution to this simple least-squares problem is a pair of parameters m and b for which the function f(x) = mx + b minimizes D.

### 7.5.2 Crystallographic refinement by least squares

In the crystallographic case, the parameters we seek (analogous to *m* and *b*) are, for all atoms *j*, the positions  $(x_j, y_j, z_j)$  that best fit the observed structure-factor amplitudes. Because the positions of atoms in the current model can be used to calculate structure factors, and hence to compute the *expected* structure-factor amplitudes ( $|\mathbf{F}_{calc}|$ ) for the current model, we want to find a set of atom positions that give  $|\mathbf{F}_{calc}|_s$ , analogous to calculated values  $f(x_i)$ , that are as close as possible to the  $|\mathbf{F}_{obs}|_s$  (analogous to observed values  $y_i$ ). In least-squares terminology, we want to select atom positions that minimize the squares of differences between corresponding  $|\mathbf{F}_{calc}|_s$  and  $|\mathbf{F}_{obs}|_s$ . We define the difference between the observed amplitude  $|\mathbf{F}_{obs}|$  and the measured amplitude  $|\mathbf{F}_{calc}|$  for reflection *hkl* as  $(|\mathbf{F}_0| - |\mathbf{F}_c|)_{hkl}$ , and we seek to minimize the function  $\Phi$ , where

$$\Phi = \sum_{hkl} w_{hkl} \left( |\mathbf{F}_{o}| - |\mathbf{F}_{c}| \right)_{hkl}^{2}.$$
(7.7)

#### Section 7.5 Structure refinement

In words, the function  $\Phi$  is the sum of the squares of differences between observed and calculated amplitudes. The sum is taken over all reflections *hkl* currently in use. Each difference is weighted by the term  $w_{hkl}$ , a number that depends on the reliability of the corresponding measured intensity. As in the simple example, according to principles of statistics, the weight should be  $1/(\sigma_{hkl})^2$ , where  $\sigma$  is the standard deviation from multiple measurements of  $|\mathbf{F}_{obs}|$ . Because the data do not usually contain enough measurements of each reflection to determine its standard deviation, other weighting schemes have been devised. Starting from a reasonable model, the least-squares refinement method succeeds about equally well with a variety of weighting systems, so I will not discuss them further.

### 7.5.3 Additional refinement parameters

We seek a set of parameters that minimize the function  $\Phi$ . These parameters include the atom positions, of course, because the atom positions in the model determine each  $\mathbf{F}_{calc}$ . But other parameters are included as well. One is the *temperature factor*  $B_j$ , or B-factor, of each atom j, a measure of how much the atom oscillates around the position specified in the model. Atoms at side-chain termini are expected to exhibit more freedom of movement than main-chain atoms, and this movement amounts to spreading each atom over a small region of space. Diffraction is affected by this variation in atomic position, so it is realistic to assign a temperature factor to each atom and include the factor among parameters to vary in minimizing  $\Phi$ . From the temperature factors computed during refinement, we learn which atoms in the molecule have the most freedom of movement, and we gain some insight into the dynamics of our largely static model. In addition, adding the effects of motion to our model makes it more realistic and hence more likely to fit the data precisely.

Another parameter included in refinement is the *occupancy*  $n_j$  of each atom j, a measure of the fraction of molecules in which atom j actually occupies the position specified in the model. If all molecules in the crystal are precisely identical, then occupancies for all atoms are 1.00. Occupancy is included among refinement parameters because occasionally two or more distinct conformations are observed for a small region like a surface side chain. The model might refine better if atoms in this region are assigned occupancies equal to the fraction of side chains in each conformation. For example, if the two conformations occur with equal frequency, then atoms involved receive occupancies of 0.5 in each of their two possible positions. By including occupancies among the refinement parameters, we obtain estimates of the frequency of alternative conformations, giving some additional information about the dynamics of the protein molecule. The factor  $|\mathbf{F}_c|$  in Eq. (7.7) can be expanded to show all the parameters included in refinement, as follows:

$$\mathbf{F}_{c} = G \cdot \sum_{j} n_{j} f_{j} e^{2\pi i (hx_{j} + ky_{j} + lz_{j})} \cdot e^{-B_{j} [(\sin \theta)/\lambda]^{2}}.$$
(7.8)

Although this equation is rather forbidding, it is actually the familiar Eq. (5.15) with the new parameters included. Equation (7.8) says that structure factor  $\mathbf{F}_{hkl}$  can be calculated ( $\mathbf{F}_c$ ) as a Fourier sum containing one term for each atom *j* in the current model. *G* is an overall scale factor to put all  $\mathbf{F}_c$ s on a convenient numerical scale. In the *j*th term, which describes the diffractive contribution of atom *j* to this particular structure factor,  $n_j$  is the occupancy of atom *j*;  $f_j$  is its scattering factor, just as in Eq. (5.15);  $x_j$ ,  $y_j$ , and  $z_j$  are its coordinates; and  $B_j$  is its temperature factor. The first exponential term is the familiar Fourier description of a simple three-dimensional wave with frequencies h, k, and l in the directions x, y, and z. The second exponential shows that the effect of  $B_j$  on the structure factor depends on the angle of the reflection  $[(\sin \theta)/\lambda]$ .

### 7.5.4 Local minima and radius of convergence

As you can imagine, finding parameters (atomic coordinates, occupancies, and temperature factors for all atoms in the model) to minimize the differences between all the observed and calculated structure factors is a massive computing task. As in the simple example, one way to solve this problem is to differentiate  $\Phi$  with respect to all the parameters, which gives simultaneous equations with the parameters as unknowns. The number of equations equals the number of observations, in this case the number of measured reflection intensities in the native data set. The parameters are overdetermined only if the number of measured reflections is greater than the number of parameters to be obtained. The complexity of the equations rules out analytical solutions and requires iterative (successive-approximations) methods that we hope will converge from the starting parameters of our current model to a set of new parameters corresponding to a minimum in  $\Phi$ . It has been proved that the atom positions that minimize  $\Phi$  are the same as those found from Eq. (7.3), the Fourier-series description of electron density. So real-space and reciprocal-space methods converge to the same solution.

The complicated function  $\Phi$  undoubtedly exhibits many *local minima*, corresponding to variations in model conformation that minimize  $\Phi$  with respect to other quite similar ("neighboring") conformations. A least-squares procedure will find the minimum that is nearest the starting point, so it is important that the starting model parameters be near the global minimum, the one conformation that gives best agreement with the native structure factors. Otherwise the refinement will converge into an incorrect local minimum from which it cannot extract itself. The greatest distance from the global minimum from which refinement will converge properly is called the *radius of convergence*. The theoretically derived radius is  $d_{\min}/4$ , where  $d_{\min}$  is the lattice-plane spacing corresponding to the reflection of highest resolution used in the refinement. Inclusion of data from higher resolution, while potentially giving more information, decreases the radius of convergence, so the model must be ever closer to its global minimum as more data are included in refinement. There are a number of approaches to increasing the radius of convergence, and thus increasing the probability of finding the global minimum.

#### Section 7.5 Structure refinement

These approaches take the form of additional constraints and restraints on the model during refinement computations. A *constraint* is a fixed value for a certain parameter. For example, in early stages of refinement, we might constrain all occupancies to a value of 1.0. A *restraint* is a subsidiary condition imposed upon the parameters, such as the condition that all bond lengths and bond angles be within a specified range of values. The function  $\Phi$ , with additional restraints on bond lengths and angles, is as follows:

$$\Phi = \sum_{hkl} w_{hkl} (|\mathbf{F}_{o}| - |\mathbf{F}_{c}|)_{hkl}^{2}$$

$$+ \sum_{i}^{\text{bonds}} w_{i} (d_{i}^{\text{ideal}} - d_{i}^{\text{model}})^{2}$$

$$+ \sum_{i} w_{j} (\phi_{j}^{\text{ideal}} - \phi_{j}^{\text{model}})^{2},$$
(7.9)

where  $d_i$  is the length of bond *i*, and  $\phi_j$  is the bond angle at location *j*. Ideal values are average values for bond lengths and angles in small organic molecules, and model values are taken from the current model. In minimizing this more complicated  $\Phi$ , we are seeking atom positions, temperature factors, and occupancies that simultaneously minimize differences between (1) observed and calculated reflection amplitudes, (2) model bond lengths and ideal bond lengths, and (3) model bond angles and ideal bond angles. In effect, the restraints penalize adjustments to parameters if the adjustments make the model less realistic.

### 7.5.5 Molecular energy and motion in refinement

Crystallographers can take advantage of the prodigious power of today's computers to include knowledge of molecular energy and molecular motion in the refinement. In *energy refinement*, least-squares restraints are placed upon the overall energy of the model, including bond, angle, and conformational energies and the energies of noncovalent interactions such as hydrogen bonds. Adding these restraints is an attempt to find the structure of lowest energy in the neighborhood of the current model. In effect, these restraints penalize adjustments to parameters if the adjustments increase the calculated energy of the model.

Another form of refinement employs *molecular dynamics*, which is an attempt to simulate the movement of molecules by solving Newton's laws of motion for atoms moving within force fields that represent the effects of covalent and noncovalent bonding. Molecular dynamics can be turned into a tool for crystallographic refinement by including an energy term that is related to the difference between the measured reflection intensities and the intensities calculated from the model. In effect, this approach treats the model as if its energy decreases as its fit to the native crystallographic data improves. In refinement by *simulated annealing*, the model is allowed to move as if at high temperature, in hopes of lifting it out of local energy minima. Then the model is cooled slowly to find its preferred conformation at the temperature of diffraction data collection. All the while, the computer is searching for the conformation of lowest energy, with the assigned energy partially dependent upon agreement with diffraction data. In some cases, the radius of convergence is greatly increased by this process, a form of molecular dynamics refinement.

### 7.5.6 Bayesian methods: Ensembles of models

Lurking in all the preceding discussions of phasing and refinement has been a tacit but questionable assumption. It is the belief that, at every stage in phasing, it is possible to choose the one best model at the moment, and that it will improve by a linear path through other best-models-at-the-moment, and eventually lead to the best possible model. We make this assumption when we choose a working model and its accompanying set of phases, using procedures like isomorphous replacement, anomalous dispersion, solvent flattening, molecular replacement, and noncrystallographic symmetry. In fact, at each of these stages, there is uncertainty in the choice of models-in the exact values of phases from anomalous dispersion, the exact position of solvent masks, or the exact rotation and translation for a phasing model in molecular replacement. But typically, the crystallographer or the software makes a single choice and proceeds. As you know, the model can bias the subsequent refinement process, at the worst sending it into local minima from which it cannot extricate itself. The gradual release of constraints during least-squares refinement is one means of evading local minima. Weighting structure factors according to the uncertainty in their intensities and associated phases is a means of weakening the model's power to bias the refinement. Finally, simulated annealing is a means of repeatedly getting out of such minima.

A very powerful way to avoid bias is to avoid the best-model-at-the-moment assumption entirely, and instead of refining the model, refine instead a representative sample of all models that are compatible with the current level of ambiguity in the data. The most ambitious movement for permeating crystallography with this kind of thinking is called the Bayesian program, because it is based on Bayes's theorem, which is the foundation of a general and well-established statistical method for decision making in the face of incomplete information. In essence, Bayesian thinking expresses current knowledge about a range of possible hypotheses in the form of prior probabilities. Observation or data collection leads to a new state of knowledge described by *posterior probabilities*. In an iterative process, models with higher prior probability are permuted and confronted with the data to raise prior probabilities further. In its broadest form, the Bayesian program in crystallography is an attempt to combine the brawn of direct phasing methods in small-molecule crystallography with the brains of prior knowledge in macromolecular crystallography. The idea that unifies these two realms is a common means of assessing the correctness of their models. Bringing these two forms of inference together promises to extend to macromolecular

### Section 7.5 Structure refinement

crystallography the kind of automation that has routinized much of small-molecule crystallography.

First, here is a simple example of how Bayesian inference can alter our assessment of possible outcomes. Suppose Dorothy and Max are running for office. Knowing no reason that voters might strongly prefer either candidate (prior knowledge), we assign a prior probability of 0.5 to the election of either. Then we conduct a survey, asking ten randomly chosen voters how they will vote. We find that seven plan to vote for Dorothy. Bayesian inference allows us to compute the probability of Dorothy's success in the light of this survey—the posterior probability of Dorothy's election. The calculated prior probability is almost 0.9 (for calculation, see "Bayesian Inference" at the CMCC home page). In this case, a small amount of data greatly alters our assessment of Dorothy's chances.

Facing the ambiguities typical of many stages of crystallography, Bayesian thinking about improving the current model entails iteration of the following steps:

- 1. Generate an ensemble of models that form a representative sample of all models that are compatible with the current level of ambiguity about the structure factors. Each model, referred to as an hypothesis  $(H_J)$ , is assigned a *prior probability*  $[P_{prior}(H_J)]$  based on knowledge from outside the diffraction measurements. For example, at the stage of solvent flattening, such knowledge might include the percentages of protein and water in the crystal, and the known average densities of protein and bulk solvent. We were assigning a prior probability when we assumed that Dorothy and Max had equal chances of winning the election.
- 2. Compute a *probability distribution* of observed structure factors expected from each model  $[P(\mathbf{F} \mid H_J)]$ , and remove the phases to obtain a probability distribution of diffraction amplitudes expected from each model  $[P(|\mathbf{F}| \mid H_J)]$ . (Read the single vertical line | as "given that" or just "given," so that  $[P(|\mathbf{F}| \mid H_J)]$  says "probability distribution P of a set of structure factor amplitudes  $|\mathbf{F}|$  given hypothesis  $H_J$ .") This entails recognizing that structure factors are not independent of each other, but exhibit characteristic *distributions* in their intensities. (A familiar example of a probability distribution is the "bell curve" or *normal distribution* of a large set of exam scores.) The shapes of distributions of structure factor amplitudes constitute a signal that can be used in assessing the probability of a model's correctness.
- 3. Compute, by comparing the probability distribution of amplitudes given each model with the probability distribution of amplitudes in the actual data [that is, comparing  $P(|\mathbf{F}| | H_J)$  with  $P(|\mathbf{F}_{obs}|)$ ], the *likelihood*  $\Lambda(H_J)$  of each model given the observed data. Likelihoods are related to probabilities as follows:

$$\Lambda(H_J \mid |\mathbf{F}_{obs}|) = P(|\mathbf{F}_{obs}| \mid H_J).$$
(7.10)

4. Use Bayes's theorem to compute a *posterior probability* for each model in the light of the data. Bayes's theorem formulated for this situation is

$$P_{\text{post}}(H_J \mid |\mathbf{F}_{\text{obs}}|) = \frac{P_{\text{prior}}(H_J) \cdot \Lambda(H_J \mid |\mathbf{F}_{\text{obs}}|)}{\sum_K P_{\text{prior}}(H_K) \cdot \Lambda(H_K \mid |\mathbf{F}_{\text{obs}}|)}.$$
(7.11)

In words, the posterior probability of hypothesis  $H_J$  given the data (observed structure factors) equals its prior probability times its likelihood given the data, divided by the sum of all such  $P \cdot \Lambda$  terms for all hypotheses. The denominator assures that the probabilities for all hypotheses add up to 1.00. The posterior probability of a model tells how it stacks up against the other models in their ability to fit the observed data.

5. Use the models of highest posterior probability as the basis for generating a new ensemble of models of higher likelihood, and use these models in step one to begin the next interation. Interating this process can retrieve, in the calculation of structure factors from each hypothesis, missing information, such as uncertain phases or even intensities that are missing or poorly measured.

The founder of the Bayesian program, Gerard Bricogne, has likened this iterative process to the molecular biologist's tool of phage display, a method of "evolving" proteins with specific ligand-binding affinities (Fig. 7.4). The method entails introducing genes into bacteriophages in such a way that expression of the genes results in display of the gene product on the phage surface. This makes it possible to use ligand-affinity chromatography to separate phages that express (and contain) the gene from those that do not. To evolve proteins of higher binding affinity, the molecular biologist generates diversity by producing large numbers of mutated forms of the desired gene, expresses them in phages to produce populations displaying many mutated forms of the protein, and selects the best binders by ligand-affinity chromatography. This process can be interated by generating mutational diversity within populations of the best binders from the previous generation. By analogy, Bricogne refers to the Bayesian method as "phase display," which generates a diverse population of hypothetical models that are compatible with the current state of ambiguity, typically by permuting phases (analogous to mutations); "expresses" these hypotheses by converting them into probability distributions of observable diffraction intensities, and finally, selects by their ability to "bind to the data"-that is, to agree with the experimental data accurately.

How do Bayesian crystallographers tell how they are doing? Recall that, in least-squares refinement (p. 159), the *refinement target* is minimization of differences between (a) measured intensities and (b) intensities calculated from *the* current model. In the searches typical of molecular replacement (rotation and translation searches) and noncrystallographic symmetry averaging (self-rotation), the refinement target is maximization of the correlation between Patterson maps.

166



**Figure 7.4**  $\triangleright$  Phage display analogy of Bayesian approach to crystallography. Both methods entail generation of diversity, followed by expression of diversity in a form that is subject to selection of best results.

In applications of Bayesian methods, by contrast, the refinement target is maximization of either likelihood (when data are plentiful) or posterior probability (a more complex task, required when data are more scarce). An example of a likelihood function is *log-likelihood gain*, expressed as

$$LLG(H_J) = \log \frac{\Lambda(H_J)}{\Lambda(H_0)},$$
(7.12)

which is a measure of the likelihood of a model  $\Lambda(H_J)$  in comparison to the likelihood of some null hypothesis  $\Lambda(H_0)$ . A typical null hypothesis might be a model in which the atoms are distributed in a physically realistic but random manner within the unit cell. As knowledge of phases improves, the surviving models are of higher and higher likelihood.

In essence, Bayesian methods provide rigorous organization, bookkeeping, and quality control for all the models that remain compatible with the current state of knowledge as that knowledge grows. In this respect, the methods formalize the many possibilities that prior knowledge keeps alive in the crystallographer's mind during structure determination. At the same time, these methods "try everything" that remains compatible with the current state of knowledge, reflecting common ground with direct phasing methods.

If least-squares methods are computationally intensive, then certainly refinement methods involving large numbers of hypothetical data sets and coordinate sets are much more so. But what is computationally intensive today typically becomes routine in just a few years—computers get faster and programmers get cleverer. Each year, more of the demands of the Bayesian program can be accomodated by readily available computing. These methods are sure to be lurking in the black boxes of software that ultimately automate more and more of crystallographic structure determination.

### 7.6 Convergence to a final model

### 7.6.1 Producing the final map and model

In the last stages of structure determination, the crystallographer alternates computed, reciprocal-space refinement with map fitting, or real-space refinement. The most powerful software packages can start with the first map that shows any sort of structural detail, and cycle between automated map fitting and refinement cycles until most of the model is complete. In an approach that borrows from the Shake and Bake method, one type of program starts from atoms randomly placed in density, minimizing initial biases in the atomic model. At the end of refinement, such programs point the user to residues that are still uncertain, allowing manual intervention to fix the last details. So in the best cases, manual model building is only needed in a few problem areas. In general, whether the model building is

#### Section 7.6 Convergence to a final model

manual or automated, constraints and restraints are lifted as refinement proceeds, so that agreement with the original reflection intensities is gradually given highest priority. When ordered water becomes discernible in the map, water molecules are added to the model, and occupancies are no longer constrained, to reflect the possibility that a particular water site may be occupied in only a fraction of unit cells. Early in refinement, all temperature factors are assigned a starting value. Later, the value is held the same for all atoms or for groups of similar atoms (like all backbone atoms as one group and all side-chain atoms as a separate group), but the overall value is not constrained. Finally, individual atomic temperature factors are allowed to refine independently. Early in refinement, the whole model is held rigid to refine its position in the unit cell. Then blocks of the model are held rigid while their positions refine with respect to each other. In the end, individual atoms are freed to refine with only stereochemical restraints. This gradual release of the model to refine against the original data is an attempt to prevent it from getting stuck in local minima. Choosing when to relax specific constraints and restraints was once considered to be more art than science, but many of today's crystallographic software packages can now negotiate this terrain and greatly reduce the manual labor of refinement.

Near the end of refinement, the  $\mathbf{F}_{o} - \mathbf{F}_{c}$  map becomes rather empty except in problem areas, and sigma-A weighted maps show good agreement with the model and no empty density. Map fitting becomes a matter of searching for and correcting errors in the model, which amounts to extricating the model from local minima in the reciprocal-space refinement. Wherever model atoms lie outside  $2\mathbf{F}_{o} - \mathbf{F}_{c}$ contours, the  $\mathbf{F}_{o} - \mathbf{F}_{c}$  map will often show the atoms within negative contours, with nearby positive contours pointing to correct locations for these atoms. Many crystalline proteins possess disordered regions, where the maps do not clear up and become unambiguously interpretable. Such regions of structural uncertainty are simply omitted from the model, and this omission is mentioned in published papers on the structure and in the header information of Protein Data Bank files (see Sec. 7.7, p. 173).

At the end of successful refinement, the  $2\mathbf{F}_{o} - \mathbf{F}_{c}$  map almost looks like a space-filling model of the protein. (Refer to Fig. 2.3*b*, p. 11, which is the final model built into the same region shown in Fig. 7.2, p. 158 and Fig. 7.3, p. 158.) The backbone electron density is continuous, and peptide carbonyl oxygens are clearly marked by bulges in the backbone density. Side-chain density, especially in the interior, is sharp and fits the model snugly. Branched side chains, like those of valine, exhibit distinct lobes of density representing the two branches. Rings of histidine, phenylalanine, tyrosine, and tryptophan are flat, and in models of the highest resolution, aromatic rings show a clear depression or hole in the density at their centers. Looking at the final model in the final map, you can easily underestimate the difficulty of interpreting the early maps, in which backbone density is frequently weak and broken, and side chains are missing or shapeless.

You can get a *rough* idea of how refinement gradually reveals features of the molecule by comparing electron-density maps computed at low, medium, and high resolution, as in Fig. 7.5. Each photo in this set shows a section of the final



**Figure 7.5**  $\blacktriangleright$  Electron-density maps at increasing resolution (stereo). Maps were calculated using final phases, and Fourier sums were truncated at the the following resolution limits: (*a*) 6.0 A; (*b*) 4.5 A; (*c*) 3.0 A; (*d*) 1.6 A.

#### Section 7.6 Convergence to a final model

ALBP model in a map calculated with the final phases, but with  $|\mathbf{F}_{obs}|$ s limited to specified resolution. In (*a*), only  $|\mathbf{F}_{obs}|$ s of reflections at resolution 6 Å or greater are used. With this limit on the data (which amounts to including in the  $2\mathbf{F}_{o} - \mathbf{F}_{c}$  Fourier sum only those reflections whose indices *hkl* correspond to sets of planes with spacing *d<sub>hkl</sub>* of 6 Å or greater), the map of this pleated-sheet region of the protein is no more than a featureless sandwich of electron density. As we extend the Fourier sum to include reflections out to 4.5 Å, the map (*b*) shows distinct, but not always continuous, tubes of density for each chain. Extending the resolution to 3.0 Å, we see density that defines the final model reasonably well, including bulges for carbonyl oxygens (red) and for side chains. Finally, at 1.6 Å, the map fits the model like a glove, zigzagging precisely in unison with the backbone of the model and showing well-defined lobes for individual side-chain atoms.

Look again at the block diagram of Fig. 7.1, p. 148, which gives an overview of structure determination. Now I can be more specific about the criteria for error removal or *filtering*, which is shown in the diagram as horizontal dashed lines in real and reciprocal space. Real-space filtering of the map entails removing noise or adding density information, as in solvent flattening or flipping. Reciprocal-space filtering of *phases* entails using only the strongest reflections (for which phases are more accurate) to compute the early maps, and using figures of merit and phase probabilities to select the most reliable phases at each stage. The molecular model can be filtered in either real or reciprocal space. Errors are removed in real space by improving the fit of model to map, and by allowing only realistic bond lengths and angles when adjusting the model (regularization). Here the criteria are structural parameters and congruence to the map (real space). Model errors are removed in reciprocal space (curved arrow in center) by least-squares refinement, which entails adjusting atom positions in order to bring calculated intensities into agreement with measured intensities. Here the criteria are comparative structure-factor amplitudes (reciprocal space). Using the Fourier transform, the crystallographer moves back and forth between real and reciprocal space to nurse the model into congruence with the data.

### 7.6.2 Guides to convergence

Judging convergence and assessing model quality are overlapping tasks. I will discuss criteria of convergence here. In Chapter 8, I will discuss some of the criteria further, particularly as they relate to the quality and usefulness of the final model.

The progress of iterative real- and reciprocal-space refinement is monitored by comparing the measured structure-factor amplitudes  $|\mathbf{F}_{obs}|$  (which are proportional to  $(I_{obs})^{1/2}$ ) with amplitudes  $|\mathbf{F}_{calc}|$  from the current model. In calculating the new phases at each stage, we learn what intensities our current model, if correct, would yield. As we converge to the correct structure, the measured **F**s and the calculated **F**s should also converge. The most widely used measure of convergence is the

residual index, or R-factor (Sec. 6.5.5, p. 141):

$$R = \frac{\sum ||\mathbf{F}_{obs}| - |\mathbf{F}_{calc}||}{\sum |\mathbf{F}_{obs}|}.$$
(7.13)

In this expression, each  $|\mathbf{F}_{obs}|$  is derived from a measured reflection intensity and each  $|\mathbf{F}_{calc}|$  is the amplitude of the corresponding structure factor calculated from the current model. Values of *R* range from zero, for perfect agreement of calculated and observed intensities, to about 0.6, the *R*-factor obtained when a set of measured amplitudes is compared with a set of random amplitudes. An *R*-factor greater than 0.5 implies that agreement between observed and calculated intensities is very poor, and many models with R of 0.5 or greater will not respond to attempts at improvement unless more data are available. An early model with *R* near 0.4 is promising and is likely to improve with the various refinement methods I have presented. A desirable target *R*-factor for a protein model refined with data to 2.5 Å is 0.2. Very rarely, small, well-ordered proteins may refine to R-values as low as 0.1, whereas small organic molecules commonly refine to R values below 0.05.

A more demanding and revealing criterion of model quality and of improvements during refinement is the *free R-factor*,  $R_{\text{free}}$ .  $R_{\text{free}}$  is computed with a small set of randomly chosen intensities, the "test set," which are set aside from the beginning and not used during refinement. They are used only in cross-validation, a quality control process that entails assessing the agreement between calculated (from the model) and observed data. At any stage in refinement,  $R_{\rm free}$  measures how well the current atomic model predicts a subset of the measured intensities that were not included in the refinement, whereas R measures how well the current model predicts the entire data set that produced the model. You can see a sort of circularity in R that is avoided in  $R_{\text{free}}$ . Many crystallographers believe that  $R_{\text{free}}$  gives a better and less-biased measure of overall model. In many test calculations,  $R_{\rm free}$ correlates very well with phase accuracy of the atomic model. In general, during intermediate stages of refinement,  $R_{\rm free}$  values are higher than R, but in the final stages, the two often become more similar. Because incompleteness of data can make structure determination more difficult (and perhaps because the lower values of R are somewhat seductive in stages where encouragement is welcome), some crystallographers at first resisted using  $R_{\text{free}}$ . But most now use both Rs to guide them in refinement, looking for refinement procedures that improve both Rs, and proceeding with great caution when the two criteria appear to be in conflict. In Bayesian methods,  $R_{\rm free}$  is replaced by *free log-likelihood gain*,  $L_{\rm free}$ , calculated over the same test data set as  $R_{\rm free}$ .

In addition to monitoring R- or L- factors as indicators of convergence, the crystallographer monitors various structural parameters that indicate whether the model is chemically, stereochemically, and conformationally reasonable. In a *chemically* reasonable model, the bond lengths and bond angles fall near the expected values for simple organic molecules. The usual criteria applied are the root-mean-square (rms) deviations of all the model's bond lengths and angles from

172

### Section 7.7 Sharing the model

an accepted set of values. A well-refined model exhibits rms deviations of no more than 0.02 Å for bond lengths and 4° for bond angles. Bear in mind, however, that refinement restrains bond lengths and bond angles, making them less informative indicators of model quality than structural parameters that are allowed to refine free of restraints.

A stereochemically reasonable model has no inverted centers of chirality (for instance, no D-amino acids). A conformationally reasonable model meets several criteria. (1) Peptide bonds are nearly planar, and nonproline peptide bonds are *trans*-, except where obvious local conformational constraints produce an occasional *cis*-peptide bond, which is usually followed by proline. (2) The backbone conformational angles  $\Phi$  and  $\Psi$  fall in allowed ranges, as judged from Ramachandran plots of these angles (see Chapter 8). And finally, (3) torsional angles at single bonds in side chains lie within a few degrees of stable, staggered conformations. During the progress of refinement, all of these structural parameters should continually improve. And unlike bond lengths and bond angles, conformational angles are usually not restrained during refinement, so these parameters are better indicators of model quality, as I will discuss in Chapter 8.

In addition to the guides to convergence described in this section, the model validation tools described in Section 8.2.5 can also be used to find potential model errors at the end of refinement. Error correction at this point followed by additional refinement cycles can insure that the published model is as error-free as possible. There is wide variation in the extent to which these additional validation tools are applied.

### 7.7 Sharing the model

If the molecule under study plays an important biological or medical function, an intensely interested audience awaits the crystallographer's final molecular model. The audience includes researchers studying the same molecule by other methods, such as spectroscopy or kinetics, or studying metabolic pathways or diseases in which the molecule is involved. The model may serve as a basis for understanding the properties of the protein and its behavior in biological systems. It may also serve as a guide to the design of inhibitors or to engineering efforts to modify its function by methods of molecular biology. But in the case of models produced by high-throughout crystallography programs, the substance may be of unknown function and not the subject of current study. In this case, the model may simply be deposited, without fanfare or further analysis, into a database. Later, researchers who develop an interest in the substance may find, to their pleasant surprise, that the structure determination has already been done for them. I advise that they worry, however, about whether the final model has received the same scrutiny as one more anxiously awaited.

Most crystallographers would agree that an important part of their work is to make molecular structures available to the larger community of scientists. This belief is reflected in the policies of many journals and funding organizations that require public availability of the structure as a condition of publication or financial support. But in the desire to make sure that they capitalize fully on commercially important leads provided by new models, crystallographers sometimes delay making models available. Research support from industry sometimes carries the stipulation that the full atomic details of new models be withheld long enough to allow researchers to explore and perhaps patent ideas of potential commercial value. Some journals now allow such delays, agreeing to publish announcements and discussions of new models if public access within a reasonable period is assured. This assurance can be enforced by requiring, as a condition of publication, that the researcher submit the model to the Protein Data Bank (discussed later in this section) immediately, but with the stipulation that public access is denied for a time, usually no more than one year. It is becoming more common for publication to require immediate public availability of the model.

Crystallographers share the fruits of their work in the form of lists of atomic coordinates, which can be used to display and study the molecule with molecular graphics programs (Chapter 11). It is becoming more common, with the improvement of resources to use them, for crystallographers to share also the final structure factors, from which electron-density maps can be computed. The audience for structure factors includes other crystallographers developing new techniques of data handling, refinement, or map interpretation. The audience for electron-density maps includes anyone who wants to judge the quality of the final molecular image, and assess the evidence for published conclusions drawn from the model.

Upon request, many authors of published crystallographic structures provide coordinate lists or structure factors by computer mail directly to interested parties. But the great majority of models and structure factors are available through the Protein Data Bank (PDB).<sup>1</sup> Crystallographers can satisfy publication and funding requirements for availability of their structures by depositing coordinates with the PDB. Depositors are required to run a series of checks for errors and inconsistencies in coordinates and format, and a validation check that produces a report on various model properties. These validation checks are not as extensive as can be done with the online validation tools described in Chapter 8 (p. 189). Once the files are processed, the PDB makes them available free over the World Wide Web, in a standard text format. As the size of the PDB has grown, so has the power of its searching tools, and the number of links from individual entries to other databases, including the massive protein and nucleic-acid sequence databases that have grown from worldwide genome projects.

The PDB structure files, which are called *atomic coordinate entries*, can be read within editor or word-processor programs. Almost all molecular graphics programs read PDB files directly or use them to produce their own files in binary

<sup>&</sup>lt;sup>1</sup>The Protein Data Bank is described fully in H. M. Berman, J. Westbrook, Z. Feng, G. Gilliland, T. N. Bhat, H. Weissig, I. N. Shindyalov, P. E. Bourne: The Protein Data Bank. Nucleic Acids Research, 28, pp. 235–242 (2000).

form for rapid access during display. In addition to the coordinate list, a PDB file contains a *header* or opening section with information about published papers on the protein, details of experimental work that produced the structure, and other useful information.

Following is a brief description of PDB file contents, based on requirements for newly deposited files in March of 2005. The line types (also called record types), given in capital letters, are printed at the left of each line in the file. Some of these line types do not apply to all models and may be missing. Typical contents of a coordinate file, in order of appearance, are

- ► HEADER lines, containing the file name ("ID code"), deposition date, and a brief title.
- ► TITLE lines, containing a brief title, usually based on the title of the model publication.
- COMPND lines, containing the name of the protein, including synonyms.
- ► SOURCE lines, giving the organism from which the protein was obtained, and if an engineered protein, the organism in which the protein was expressed.
- ► KEYWDS lines, giving keywords that would guide a search to this file.
- ► EXPDTA lines, giving the experimental method (X-ray diffraction or NMR).
- ► AUTHOR lines, listing the persons who placed this data in the Protein Data Bank.
- ▶ REVDAT lines, listing all revision dates for data on this protein.
- JRNL lines, giving the journal reference to the lead article about this model. When viewed online, these lines are usually linked to the article in PubMed.
- REMARK lines, containing (1) references to journal articles about the structure of this protein and (2) general information about the contents of this file, including many specifics about resolution, refinement methods, and final criteria of model quality (see Chapter 8). The REMARK lines can provide extensive information about the details of structure determination, including lists of missing residues. Ranges of REMARK line numbers are reserved for specific types of information, as detailed in PDB documentation.
- DBREF lines, which cross-reference PDB entries to other databases, particularly sequence databases.
- SEQADV lines, which identify conflicts between the PDB sequence and those of other databases. Such conflicts might arise because of sequencing revisions made during the structure determination, but more commonly, the conflicts reflect engineered sequence changes made to study protein function, or to facilitate structure determination, such as attachment of histidines for purification, or replacement of methionine with selenomethionine.
- ► SEQRES lines, giving the amino-acid sequence of the protein, with amino acids specified by three-letter abbreviation.

### Chapter 7 Obtaining and Judging the Molecular Model

- HET, HETNAM, and FORMUL lines, listing the names, synonyms, and formulas of cofactors, prosthetic groups, or other nonprotein substances present in the structure. Online versions of PDB files often contain links to more information about HET groups, including links to graphics displays of their structures (see Chapter 11).
- ► HELIX, SHEET, TURN, CISPEP, LINK, and SITE lines, listing the elements of secondary structure in the protein, residues involved in *cis*-peptide bonds (almost always involving proline as the second residue), cross-links such as disulfide bonds, and residues in the active site of the protein.
- ► CRYST lines, giving the unit cell dimensions and space group.
- ORIG and SCALE lines, which give matrix transformations relating the orthogonal angstrom coordinates in the entry to the submitted coordinates, which might not have been orthogonal, or might have been in fractions of unit-cell lengths instead of in angstroms.

This concludes what is loosely called the header of the PDB file. The remaining lines provide the atomic coordinates and other information needed to display and analyze the model:

- ATOM lines, containing the atomic coordinates of all protein atoms, plus their structure factors and occupancies. Atoms are listed in the order given in the paragraph following this list.
- ► TER lines, among the ATOM lines, specifying the termini of distinct chains in the model.
- HETATM lines, which contain the same information as ATOM lines for any nonprotein molecules (cofactors, prosthetic groups, and solvent molecules, collectively called *heteromers* or, loosely, *het-groups*) included in the structure and listed in HET, HETNAM, and FORMUL lines above.
- CONECT lines, which list covalent bonds between nonprotein atoms in the file.
- MASTER and END lines, which provide some data for record keeping, and mark the end of the file.

After the header comes a list of model atoms in standard order. Atoms in the PDB file are named and listed according to a standard format in an all-English version of the Greek-letter conventions used by organic chemists. For each amino acid, beginning at the N terminus, the backbone atoms are listed in the order alpha nitrogen N, alpha carbon CA, carbonyl carbon C, and carbonyl oxygen O, followed by the side chain atoms, beta carbon CB, gamma carbon CG, and so forth. In branched side chains (or rings), atoms in the two branches are numbered 1 and 2 after the proper Greek letter. For example, the atoms of aspartic acid, in the order of PDB format, are N, CA, C, O, CB, CG, OE1, and OE2. The terminal atoms of

176

### Section 7.7 Sharing the model

the side chain are followed in the file by atom N of the next residue. There are no markers in the file to tell where one residue begins and another ends; each N marks the beginning of the next residue.

You will find more information on PDB file contents and other available information, as well as tutorials on using the PDB, at their web site (see CMCC home page). The coordinate files are only the first level of information about the model, and related models that you will find at the Protein Data Bank. PDB users may subscribe to a discussion group, providing an email forum for discussion of all aspects of PDB use. Questions to the list often reveal attempts to use the PDB in innovative ways. For example, at the moment, many subscribers are interested in so-called *data mining*—finding ways to extract specific information from many PDB entries automatically.

In this form, as a PDB atomic coordinate entry, a crystallographic structure becomes a matter of public record. The final model of the molecule can then fall before the eyes of anyone equipped with a computer and an appropriate molecular display program. It is natural for the consumer of these files, as well as for anyone who sees published structures in journals or textbooks, to think of the molecule as something someone has seen more or less directly. Having read this far, you know that our crystallographic vision is quite indirect. But you probably still have little intuition about possible limits to the model's usefulness. For instance, just how precise are the relative locations of atoms? How much does molecular motion alter atomic positions? For that matter, how well does the model fit the original diffraction data from which it was extracted? These and other questions are the subject of Chapter 8, in which I will start you off toward becoming a discriminating consumer of the crystallographic product. This entails understanding several criteria of model quality and being able to extract these criteria from published accounts of crystallographic structure determination, or from the model and X-ray data itself, by way of online validation tools. Be aware that the Protein Data Bank does not check or vouch for the quality of models, but a full PDB entry includes much of the information that users need in order to verify that a model is good enough for their purposes. In addition to the coordinate entry, the PDB provides links to many additional tools for the wise user of models.

This Page Intentionally Left Blank

### ► Chapter 8

## A User's Guide to Crystallographic Models

### 8.1 Introduction

If you are not a crystallographer, this chapter may well be the most important one for you. Although structure determination by X-ray crystallography becomes more and more accessible every year, many structural biologists will never determine a protein structure by X-ray crystallography. But many will at some time use a crystallographic model in research or teaching. In research, study of molecular models by computer graphics is an indispensable tool in formulating mechanisms of protein action (for instance, binding or catalysis), searching for modes of interaction between molecules, choosing sites to modify by chemical methods or site-specific mutagenesis, and designing inhibitors of proteins involved in disease. Because protein chemists would like to learn the rules of protein folding, every new model is a potential test for proposed theories of folding, as well as for schemes for predicting conformation from amino-acid sequence. Every new model for which homologous sequences are known is a potential scaffold on which to build homology models (Chapter 11). In education, modern texts in biology and chemistry are effectively and dramatically illustrated with graphics images, sometimes as stereo pairs. Projection monitors allow instructors to show "real-time" graphics displays in the classroom, giving students vivid, animated, three-dimensional views of complex molecules. Any up-to-date structural biology course, such as biochemistry, molecular biology, or bioinformatics, now includes the study of macromolecules using molecular graphics programs on personal computers. Free or inexpensive, yet very powerful, graphics programs for personal computers, combined with easy

online access to the Protein Data Bank, now make it possible for anyone to study any available molecular model.

In all of these applications, there is a tendency to treat the model as a physical entity, as a real object seen or filmed. How much confidence in the crystallographic model is justified? For instance, how precisely does crystallography establish the positions of atoms in the molecule? Are all of the atomic positions equally well established? How does one rule out the possibility that crystallizing the protein alters it in some significant way? The model is a static image of a dynamic molecule, a springy system of atoms that breathes with characteristic vibrations, and tumbles dizzily through solution, as it executes its function. Does crystallography give us any insight into these motions? Are parts of the molecule more flexible than others? Are major movements of structural elements essential to the molecule's action? How does the user decide whether proposed motions of the molecule are reasonable?

More seriously, how does the user confirm that the crystallographic data actually justify published conclusions drawn from the model? As I stated earlier, the Protein Data Bank neither judges nor vouches for the quality of models. The PDB simply requires that authors provide, in a consistent format, the information needed to make these judgments. Is there really any chance that published conclusions are not justified by the data? It is rare, but unfortunately, yes. In any area of science, as in any intellectually or economically competitive endeavor, there are those who gild the lily, take short cuts, or think wishfully. By far most PDB models are deposited with diligent attention to detail and accuracy of representation, but within the anonymity of tens of thousands of entries posted in what is basically an honor system, there are all imaginable shortcomings and all imaginable reasons for them, so you should thoroughly check any model in the PDB before relying on it. If you are just planning to make a picture to illustrate a point in a lecture, then model weaknesses might not affect you. But if you are going to base research directions and conclusions on a published model, or if you are going to use a structure determination as a teaching example (as I will do later in this chapter), you should start by getting to know it and the evidence behind it very well.

In this chapter, I will discuss the *inherent* strengths and limitations of molecular models obtained by X-ray diffraction. My aim is to help you to use crystallographic models wisely and appropriately, and realize just what is known, and what is unknown, about a molecule that has yielded up some of its secrets to crystallographic analysis. Knowing the limitations that are inherent to the crystallographic process will also put you in the position to recognize if a specific model really lives up to its depositor's claims. To help you learn how to make quality assessments for yourself, I will review some online model validation tools that are currently available. I will conclude this chapter by discussing a recent structure publication, as it appeared in a scientific journal. Here my goals are (1) to help you learn to extract criteria of model quality from published structural reports, and (2) to review some basic concepts of protein crystallography from previous chapters (repetition is one of the staples of effective teaching).

# 8.2 Judging the quality and usefulness of the refined model

### 8.2.1 Structural parameters

As discussed in Sec. 7.6.2, p. 171, crystallographers monitor the decrease in some type of *R*-factor (or the increase in a likelihood factor) as an indicator of convergence to a final, refined model, with a general target for *R* of 0.20 or lower for proteins, and adequate additional cycles of refinement to confirm that *R* is not still declining. In addition, various constraints and restraints are relaxed during refinement, and after these restricted values are allowed to refine freely, they should remain in, or converge to, reasonable values. Among these are the root-mean-square (rms) deviations of the model's bond lengths, angles, and conformational angles from an accepted set of values based upon the geometry of small organic molecules. A refined model should exhibit rms deviations of no more than 0.02 Å for bond lengths and 4° for bond angles. These values are routinely calculated during refinement to be sure that all is going well. Because they are restrained in earlier stages, they are not as valuable as quality indicators as parameters that are allowed to vary freely throughout the refinement.

In effect, protein structure determination is a search for the conformation of a molecule whose chemical composition is known. For this reason, conformational angles about single bonds are not constrained during refinement, and they should settle into reasonable values. Spectroscopic evidence abundantly implies that peptide bonds are planar, and some refinements constrain peptide geometry. If unconstrained, peptide bonds should settle down to within a few degrees of planar.

Aside from peptide bonds, the other backbone conformational angles are  $\Phi$ , along the N – C $\alpha$  bond and  $\Psi$ , along the C $\alpha$ –C bond, as shown in Fig. 8.1. In this figure,  $\Phi$  is the torsional angle of the N – C $\alpha$  bond, defined by the atoms C – N – C $\alpha$ –C (C is the carbonyl carbon) and  $\Psi$  is the torsional angle of the C $\alpha$ –C bond, defined by the atoms N – C $\alpha$  – C – N. In the figure,  $\Phi = \Psi = 180^{\circ}$ .

Model studies show that, for each amino acid, the pair of angles  $\Phi$  and  $\Psi$  is greatly restricted by steric repulsion. The allowed pairs of values are depicted on a Ramachandran diagram (Fig. 8.2). A point ( $\Phi$ ,  $\Psi$ ) on the diagram represents the conformational angles  $\Phi$  and  $\Psi$  on either side of the alpha carbon of one residue. Irregular polygons enclose backbone conformational angles that do not give steric repulsion (yellow polygons) or give only modest repulsion (blue polygons). Location of the letters  $\alpha$  and  $\beta$  correspond to conformational angles of residues in right-handed  $\alpha$  helix and in  $\beta$  pleated sheet. Although the differences are relatively small, the shapes of allowed regions are slightly different for each of the 20 common amino acids, and in addition, Ramachandran diagrams from different sources will exhibit slight differences in the shapes of allowed regions.

Every year, someone in my biochemistry class, facing the formidable Ramachandran diagram for the first time, asks me, "Are these diagrams good for anything?" I am happy to say that they are indeed. During the final stages of



**Figure 8.1** ► Backbone conformational angles in proteins (stereo).



**Figure 8.2**  $\triangleright$  Ramachandran diagram for nonglycine amino-acid residues in proteins. Angles  $\Phi$  and  $\Psi$  are as defined in Fig. 8.1. Diagram produced with DeepView.

map fitting and crystallographic refinement, Ramachandran diagrams are a great aid in spotting conformationally unrealistic regions of the model. Now that the majority of residues are built automatically and do not receive individual scrutiny, it is especially important to have easy means of drawing attention to unrealistic features in the model. Crystallographic software packages and map-fitting programs (and some graphics display programs as well) usually contain a routine for computing  $\Phi$  and  $\Psi$  for each residue from the current coordinate list, as well as for generating the Ramachandran diagram and plotting a symbol or residue number at the position ( $\Phi$ ,  $\Psi$ ). Refinement papers often include the diagram, with an explanation of any residues that lie in high-energy ("forbidden") areas. For an example, see Fig. 8.7, p. 209. Glycines, because they lack a side chain, usually account for most of the residues that lie outside allowed regions. If nonglycine residues exhibit forbidden conformational angles, there should be some explanation in terms of structural constraints that overcome the energetic cost of an unusual backbone conformation.

The conformations of amino-acid side chains are unrestrained during refinement. In well-refined models, side-chain single bonds end up in staggered conformations (commonly referred to as *rotamers*). Distributions of side-chain conformations for all amino acids are available among the online structure validation tools that I will discuss later.

### 8.2.2 Resolution and precision of atomic positions

In microscopy, the phrase "resolution of 2 Å," implies that we can resolve objects that are 2 Å apart. If this phrase had the same meaning for a crystallographic model of a protein, in which bond distances average about 1.5 Å, we would be unable to distinguish or resolve adjacent atoms in a 2-Å map. Actually, for a protein refined at 2-Å resolution to an *R*-factor near 0.2, the situation is much better than the resolution statement seems to imply.

In X-ray crystallography, "2-Å model" means that analysis included reflections out to a distance in the reciprocal lattice of 1/(2 Å) from the center of the diffraction pattern. This means that the model takes into account diffraction from sets of equivalent, parallel planes spaced as closely as 2 Å in the unit cell. (Presumably, data farther out than the stated resolution was unobtainable or was too weak to be reliable.) Although the final 2-Å map, viewed as an empty contour surface, may indeed not allow us to discern adjacent atoms, prior knowledge in the form of structural constraints on the model greatly increase the precision of atom positions. The main constraint is that we know we can fit the map with *groups* of atoms amino-acid residues—having known connectivities, bond lengths, bond angles, and stereochemistry.

More than the resolution, we would like to know the precision with which atoms in the model have been located. For years, crystallographers used the Luzzati plot (Fig. 8.3) to estimate the precision of atom locations in a refined crystallographic model. At best, this is an estimate of the upper limit of error in atomic coordinates. The numbers to the right of each smooth curve on the Luzzati plot are theoretical estimates of the average uncertainty in the positions of atoms in the refined model



Figure 8.3 ► Luzzati diagram.

(more precisely, the rms errors in atom positions). The average uncertainty has been shown to depend upon *R*-factors derived from the final model in various resolution ranges. To prepare data for a Luzzati plot, we separate the intensity data into groups of reflections in narrow ranges of 1/d (where d is the spacing of real lattice planes). Then we plot each *R*-factor (vertical axis) versus the midpoint value of 1/d for that group of reflections (horizontal axis). For example, we calculate R using only reflections corresponding to the range 1/d = 0.395 - 0.405 (reflections in the 2.53- to 2.47-Å range) and plot this *R*-factor versus 1/d = 0.400/Å, the midpoint value for this group. We repeat this process for the range 1/d = 0.385 - 0.395, and so forth. As the theoretical curves indicate, the *R*-factor typically increases for lower-resolution data (higher values of 1/d). The resulting curve should roughly fit one of the theoretical curves on the Luzzati plot. From the theoretical curve closest to the experimental R-factor curve, we learn the average uncertainty in the atom positions of the final model. It is now widely accepted that Luzzati plots using  $R_{\rm free}$  (Sec. 7.6.2, p. 171) or cross-validated sigma-A values give better estimates of uncertainty in coordinates. Bayesian methods promise more appropriate models of how phase errors lead to model errors, and hence even more rigorous assessment of uncertainty in atom positions.

Publications of refined structures may include a Luzzati plot or one of its more modern equivalents, allowing the reader to assess very roughly the average uncertainty of atom positions in the model. Alternatively, they may simply report the uncertainty as "determined by the method of Luzzati." For highly refined models, rms errors as low as 0.15 Å are sometimes attained. For example, in Fig. 8.6*a*, p. 207 (to be discussed in Sec. 8.3.3, p. 198), the jagged curve represents the data for the refined model of adipocyte lipid binding protein (ALBP). The position of the curve on the Luzzati plot indicates that rms error for this model is about 0.34 Å, about one-fifth the length of a carbon-carbon bond.

In crystallography, unlike microscopy, the term *resolution* simply refers to the amount of data ultimately phased and used in the structure determination. In contrast, the precision of atom positions depends in part upon the resolution limits of the data, but also depends critically upon the quality of the data, as reflected by such parameters as *R*-factors. Good data can yield atom positions that are precise to within one-fifth to one-tenth of the stated resolution.

### 8.2.3 Vibration and disorder

Notice, however, that the preceding analysis gives only an upper limit and an *average*, or rms value, of position errors, and further, that the errors result from the limits of accuracy in the data. There are also two important physical (as opposed to statistical) reasons for uncertainty in atom positions: thermal motion and disorder. *Thermal motion* refers to vibration of an atom about its rest position. *Disorder* refers to atoms or groups of atoms that do not occupy the same position in every unit cell, in every asymmetric unit, or in every molecule within an asymmetric unit. The temperature factor  $B_j$  obtained during refinement reflects both the thermal motion and the disorder of atom j, making it difficult to sort out these two sources of uncertainty.

Occupancies  $n_j$  for atoms of the protein (but not necessarily its ligands, which may be present at lower occupancies) are usually constrained at 1.0 early in refinement, and in many refinements are never released, so that both thermal motion and disorder show their effects upon the final *B* values. In some cases, after refinement converges, a few *B* values fall far outside the average range for the model. This is sometimes an indication of disorder. Careful examination of sigma-A weighted maps or comparison of  $2\mathbf{F}_0 - \mathbf{F}_c$  and  $\mathbf{F}_0 - \mathbf{F}_c$  maps may give evidence for more than one conformation in such a troublesome region. If so, inclusion of multiple conformations followed by refinement of their occupancies may improve the *R*-factor and the map, revealing the nature of the disorder more clearly.

If  $B_j$  were purely a measure of thermal motion at atom j (and assuming that occupancies are correct), then in the simplest case of purely harmonic thermal motion of equal magnitude in all directions (called *isotropic* vibration),  $B_j$  is related to the magnitude of vibration as follows:

$$B_j = 8\pi^2 \{u_j^2\} = 79\{u_j^2\},\tag{8.1}$$

where  $\{u_j^2\}$  is the mean-square displacement of the atom from its rest position. Thus if the measured  $B_j$  is 79 Å<sup>2</sup>, the total mean-square displacement of atom *j* due to vibration is 1.0 Å<sup>2</sup>, and the rms displacement is the square root of  $\{u_j^2\}$ , or 1.0 Å. The *B* values of 20 and 5 Å<sup>2</sup> correspond to rms displacements of 0.5 and 0.25 Å. But the *B* values obtained for most proteins are too large to be seen as reflecting purely thermal motion and must certainly reflect disorder as well.

With small molecules, it is usually possible to obtain anisotropic temperature factors during refinement, giving a picture of the preferred directions of vibration for each atom. But a description of anisotropic vibration requires six parameters per atom, vastly increasing the computational task. In many cases, the total number of parameters sought, including three atomic coordinates, one occupancy, and six thermal parameters per atom, approaches or exceeds the number of measured reflections. As mentioned earlier, for refinement to succeed, observations (measured reflections and constraints such as bond lengths) must outnumber the desired parameters, so that least-squares solutions are adequately overdetermined. For this reason, anisotropic temperature factors for proteins have not usually been obtained. The increased resolution possible with synchrotron sources and cryocrystallography will make their determination more common. With this development, it may become possible to know coordinate uncertainty for individual atoms, rather than the average uncertainty obtained by methods like that of Luzzati.

Publications of refined structures often include a plot of average isotropic *B* values for side-chain and main-chain atoms of each residue, like that shown in Fig. 8.6*b*, p. 207, for ALBP. Pictures of the model may be color coded by temperature factor, red ("hot") for high values of *B* and blue ("cold") for low values of *B*. Either presentation calls the user's attention to parts of the molecule that are vibrationally active and parts that are particularly rigid. Not surprisingly, side-chain temperature factors are larger and more varied (5–60 Å<sup>2</sup>) than those of main-chain atoms (5–35 Å<sup>2</sup>).

As I mentioned earlier, if occupancies are constrained to values of 1.00, variation in actual occupancies will show up as increased temperature factors. This is true for other types of constraints as well, and so high temperature factors can mask other kinds of errors. For this reason, the temperature factor has been uncharitably referred to as a garbage can for model errors. A less pessimistic assessment is that *B*-values are informative only if other kinds of errors have been carefully removed.

Remember that we see in a crystallographic model an average of all the molecules that diffracted the X-rays. Furthermore, we see a static model representing a stable conformation of a dynamic molecule. It is sobering to realize that the crystallographic model of ALBP exhibits no obvious path for entry and departure of its ligands, which are lipid molecules like oleic acid. Similarly, comparison of the crystallographic models of hemoglobin and deoxyhemoglobin reveals no path for entry of the tiny  $O_2$  molecule. Seemingly simple processes like the binding of small ligands to proteins often involve conformational changes to states not revealed by crystallographic analysis.

Nevertheless, the crystallographic model contributes importantly to solving such problems of molecular dynamics. The refined structure serves as a starting point for simulations of molecular motion. From that starting point, which undoubtedly represents one common conformation of the protein, and from the equations of motion of atoms in the force fields of electrostatic and van der Waals forces, scientists can calculate the normal vibrational motions of the molecules and can simulate random molecular motion, thus gaining insights into how conformational change gives rise to biomolecular function. Even though the crystallographic model is static, it is an essential starting point in revealing the dynamic aspects of structure. As I will show in Chapter 9, time-resolved crystallography offers the potential of giving us highly detailed views of proteins in motion.

### 8.2.4 Other limitations of crystallographic models

The limitations discussed so far apply to all models and suggest questions that the user of crystallographic results should ask routinely. Other limitations are special cases that may or may not apply to a given model. It is important to read the original publications of a structure, as well as its PDB file header, to see whether any of the following limitations apply.

### Low-resolution models

Not all published models are refined to high resolution. For instance, publication of a low-resolution structure may be warranted if it displays an interesting and suggestive arrangement of cofactors or clusters of metal ions, provides possible insights into conformations of a new family or proteins, or displays the application of new imaging methods. In some cases, the published structure is only a crude electron-density model. Or perhaps it contains only the estimated positions of alpha carbons, so-called alpha-carbon or  $C\alpha$  models. Such models may be of limited use for comparison with other proteins, but of course, they cannot support detailed molecular analyses. In alpha-carbon models, there is great deal of uncertainty in the positions, and even in the number, of alpha carbons. Some graphics programs (Chapter 11) will open such files but show no model, giving the impression that the file is empty or damaged. The model appears when the user directs the program to show the alpha-carbon backbone only. It is not unusual for further refinement of these models to reveal errors in the chain tracing. Protein Data Bank headers include important information about model resolution and descriptions of model contents.

### **Disordered regions**

Occasionally, portions of the known sequence of a protein are never found in the electron-density maps, presumably because the region is highly disordered or in motion, and thus invisible on the time scale of crystallography. The usual procedure is simply to omit these residues from the deposited model. It is not at all uncommon for residues at termini, especially the N-terminus, to be missing from a model. Discussions of these structure-specific problems are included in a thorough refinement paper, and lists of missing residues are provided in PDB header.

### **Unexplained density**

Just as the auto mechanic sometimes has parts left over, electron-density maps occasionally show clear, empty density after all known contents of the crystal have been located. Apparent density can appear as an artifact of missing Fourier terms, but this density disappears when a more complete set of data is obtained. Among the possible explanations for density that is not artifactual are ions like phosphate and sulfate from the mother liquor; reagents like mercaptoethanol, dithiothreitol, or detergents used in purification or crystallization; or cofactors, inhibitors, allosteric effectors, or other small molecules that survived the protein purification. Later discovery of previously unknown but important ligands has sometimes resulted in subsequent interpretation of empty density.

### Distortions due to crystal packing

Refinement papers and PDB headers should also mention any evidence that the protein is affected by crystallization. Packing effects may be evident in the model itself. For example, packing may induce slight differences between what are otherwise expected to be identical subunits within an asymmetric unit. Examination of the neighborhood around such differences may reveal that intermolecular contact is a possible cause. In areas where subunits come into direct contact or close contact through intervening water, surface temperature factors are usually lower than at other surface regions.

### Functional unit versus asymmetric unit

The symmetry of functional macromolecular complexes in solution is sometimes important to understanding their functions, as in the binding of regulatory proteins having twofold rotational symmetry to palindromic DNA sequences. As discussed in Sec. 4.2.8, p. 65, in the unit cell of a crystal, the largest aggregate of molecules that possesses no symmetry elements, but can be juxtaposed on other identical entities by symmetry operations, is called the *crystallographic asymmetric unit*. Users of models should be careful to distinguish the asymmetric unit from the *functional unit*, which the Protein Data Bank currently calls the "biological molecule." For example, the functional unit of mammalian hemoglobin is a complex of four subunits, two each of two slightly different polypeptides, called  $\alpha$  and  $\beta$ . We say that hemoglobin functions as an  $\alpha_2\beta_2$  tetramer. In some hemoglobin crystals, the twofold rotational symmetry axis of the tetramer corresponds to a unit-cell symmetry axis, and the asymmetric unit is a single  $\alpha\beta$  dimer. In other cases, the asymmetric unit may contain more than one biological unit.

For technical reasons having to do with data collection strategies, crystal properties, and other processes essential to crystallography itself, the asymmetric unit is often mentioned prominently in papers about new crystallographic models. This discussion is part of a full description of the crystallographic methods for assessment of the work by other crystallographers. It is easy to get the impression that the asymmetric unit is the functional unit, but frequently it is not. Beyond the technical methods sections of a paper, in their interpretations and discussions of the meaning of the model, authors are careful to describe the functional form of the substance under study (if it is known), and this is the form that holds the most interest for users.

It is safe to think of functional-unit symmetry as not necessarily having anything to do with crystallographic symmetry. If the two share some symmetry elements, it is coincidental and may actually be useful to the crystallographer (see, for example, Sec. 7.3.3, p. 151 on noncrystallographic symmetry). But in another crystal form of the same substance, the unit cell and the functional unit may share different symmetry elements or none. For you as a user of crystallographic models, looking

188

at crystal symmetry and packing is primarily of value in making sure that you do not make errors in interpreting the model by not allowing for the relatively rare but possibly disruptive effects of crystallization. Once you are confident that the portions of the molecule that pique your interest are not affected by crystal packing, then you can forget about crystal symmetry and the asymmetric unit, and focus on the functional unit.

As mentioned earlier, the asymmetric unit may be only a part of the functional unit. This sometimes poses a problem for users of crystallographic models because the PDB file for such a crystallographic model may contain only the coordinates of the asymmetric unit. So in the case of hemoglobin, a file may contain only one  $\alpha\beta$  dimer, which is only half of what the user would like to see. The Protein Data Bank routinely takes care of this problem by preparing files containing the coordinates of all atoms in the functional unit (for example, oxy- and deoxyhemoglobin tetramers), and providing online graphics viewers for examining these models and saving the prepared coordinate files. The additional coordinates are computed by applying symmetry operations to the coordinates of the asymmetric unit. In these models, one can study all the important intersubunit interactions of the full tetramer. Another solution to this problem is for users themselves to compute the coordinates of the additional subunits. Many molecular graphic programs provide for such calculations (Sec. 11.3.10, p. 287). Users should also beware that the asymmetric unit, and hence the PDB file, may contain two or several functional units.

### 8.2.5 Online validation tools: Do it yourself!

NOTE: To find all of the validation tools discussed in the following section, see the CMCC home page.

Over the period from 1990 to 2005, the appearance of new crystallographic models turned from a trickle to a flood. As a result, many papers reporting new structures no longer provide detailed information about the crystallographic work, and thus are less helpful in providing measures of model quality. Models produced by high-throughput crystallography may be deposited quietly in the PDB, without detailed structural analysis, and without announcement in journals. A researcher whose interest turns to a previously unstudied gene, and needs to know the structure of its protein product, may find that a model is already present in the PDB, deposited more or less without comment. The researcher may then be the first to explore the new model, and is faced with assessing its quality with perhaps very limited information. In all of these cases, online validation tools can help.

Here is the first step in do-it-yourself model assessment: READ THE HEADER! As mentioned repeatedly in the preceding sections, the inherent limitations of a model should be (usually are) described in the PDB file header. As a subscriber to the PDB Discussion Forum, I know firsthand that many a model user, puzzled by breaks in the graphics display of the model or by chains that begin with residue 13 (both signs of residues not visible in density), or by side chains listed twice (alternative conformations), or by failure of graphics to display a model (C-alpha
model), or by the presence of many superimposed models [not a crystallographic structure at all, but an ensemble of NMR models (Chapter 10)], has sought help from the Forum when they could have resolved their puzzlement by simply reading the header. A conscientiously completed PDB header is a good users' manual for the model. Once you know what the authors are trying to tell you about the model, you are ready to further your assessment on your own. READ THE HEADER!

In my opinion, the current best introduction to concepts and tools for assessing model quality is the tutorial "Model Validation" at Uppsala University, Sweden. Maintained by Research Scientist Gerard Kleywegt, an outspoken proponent of vigiliance in using models, this tutorial introduces all of the classical and modern quality indicators, rates their power in helping you assess models, and guides you in their proper use. The tutorial includes a rogue's gallery of models that fall short in ways both small and large, from bent indole rings in tryptophan and flat beta carbons in valine to whole ligands that fail to show up in the electron-density map. Anyone who graduates from this little validation course will be armed with the know-how and skepticism needed to use models wisely. If this book is your classroom introduction to judging model quality, Kleywegt's tutorial is, at the moment, the best lab practical. For global quality indicators that might help you decide if a model is good enough for comparative modeling, Kleywegt recommends  $R_{\rm free}$ , Ramachandran plots, and packing scores, which tell how "comfortable" each residue is in its protein environment. For assessing local quality, such as whether an active site is good enough for modeling ligand binding, he recommends real-space R- (RSR) factors, which measure how well the electron-density map conforms to a map calculated from the ideal electron density of the model itself, sort of like comparing the map with a very accurate space-filling rendition of the model. He also suggests looking at graphs of main-chain and side-chain torsion-angle com*binations* in comparison to charts of these combinations in high-quality models. He values these measures above widely used ones such as conventional R values, RMS deviation of bond lengths and angles from ideal values, and temperature factors, all of which are more subject to bias from refinement restraints and from the model.

In my opinion, the most powerful way to assess model quality, bar none, is to examine the electron-density map and model. This shows you whether the map supports specific placement of important residues and ligands. When I am curious about a model, my first stop is the Electron Density Server (EDS), also at Uppsala University. If the depositor of the model also deposited structure factors with the PDB, you can use the EDS server to get electron-density maps, which are automatically prepared for all sets of structure factors in the PDB. In addition, the EDS allows you to produce interactive charts such as Ramachandran diagrams and graphs of residue number versus *B* values, packing scores, RSR factors, or other measures. With a click on a residue number in a graph, an online graphics display (a Java applet called Astex Viewer) opens in your browser to allow you to peruse model and map. In short, the EDS site provides almost every imaginable tool for assessing model quality. If you have the proper software (for example, DeepView—see Chapter 11), you can download maps and examine map quality

and map/model conformity in critical areas of the model. If you have the map-fitting program O (that's the name of it—O), you can download a package including model, maps, and most of the graphs mentioned above, all for interactive display within O.

Another useful and eye-opening validation site is the MolProbity Web Service at Duke University. This service uses all-atom contact analysis to help you find and correct bumps (bad contacts) and unrealistic geometry in both side chains and main chains. It provides graphics using KiNG, a modern Java version of the pioneering molecular graphics program Mage, to allow structure analysis within your web browser. Users interested in validating or optimizing a model begin their work by adding hydrogen atoms. Why? Because around 75% of atomic contacts in proteins involve hydrogen. In addition, even though macromolecular crystallography does not usually resolve hydrogens, most hydrogen positions (actually, all but -OH) are well defined by the conformations of the atoms on which hydrogens reside, and thus hydrogens can be added with confidence. With hydrogens present, contact analysis can reveal unrealistic side-chain rotamers, and can also determine the correct rotamers of side chains such as asparagine, glutamine, and histidine, whose rotamer alternatives are often difficult to determine from electron density (do you know why?). Common misplacements of main-chain alpha carbons can also be detected and corrected. These errors can be detected by deviations of beta carbons from their optimum positions, and can be corrected by small rotations about the axis between two alpha carbons, the most useful being correction at residue iby rotation about the axis between  $C\alpha_{i-1}$  and  $C\alpha_{i+1}$ . If structure factors for the model are available, then all of the model corrections at the MolProbity site can be examined within the electron-density map. Of course, this map my be slightly biased by the uncorrected model. A map based on phases from the corrected model would be the ideal reference. The MolProbity site includes a detailed users' manual and an excellent tutorial.

One of the easiest model-validation tools to use is the Biotech Validation Suite at the European Bioinformatics Institute. Users simply submit a coordinate file in PDB format, and receive in return an extensive report on the model, the results of analysis by three programs, PROCHECK, PROVE, and WHAT IF. The report gives results of checks on more than twenty aspects of the model, including verification of bond angles and bond lengths and checks for buried hydrogen-bond donors; bumps (bad contacts); flipped peptides; handedness of chiral atoms; conformational alternatives for HIS, GLN, and ASN; proline puckering; plausibility of water assignments; and atomic occupancy, to name a few. The PROCHECK suite of checks is made automatically for all PDB entries, and is linked to the Structure Explorer page for the entry, under the Analyze menu. At the time of this writing, these validation tools do not provide for online or automated correction of errors, but they probably provide the most extensive and rigorous analysis performed against the most demanding standards. As a result of these high standards, error lists are quite long and may include errors that are of little consequence to most users. Nevertheless, this package of checks leaves no stones unturned, not even the tiny ones.

The prudent crystallographer will use these validation tools before deciding that the model is fully refined, complete, and ready for publication. In other words, it makes sense to use validation tools to find and correct errors during the refinement process, so that subsequent users of the model will find no problems when applying the same validation tools. At the moment, there is no standard way of knowing what validation tools were applied to the model, beyond the standard checks (like PROCHECK) required by the PDB. PDB file headers might provide this information in some cases. Because use of these "accessory" validation tools varies widely among research groups, users of models should assume that validation is necessary.

As mentioned earlier, data mining, the automated extraction of specific information from many PDB entries, is a growing activity. Data miners should also be asking questions about model validation, because in many cases, they collect information from many files on the assumption that all models are error-free. This suggests the need for greater automation in model validation, so that data spanning many PDB files are not skewed by errors in unvalidated models.

And another reminder: For quick access to validation web sites, see the CMCC home page.

### 8.2.6 Summary

Sensible use of a crystallographic model, as with any complex tool, requires an understanding of its limitations. Some limitations, like the precision of atom positions and the static nature of the model, are general constraints on use. Others, like disordered regions, undetected portions of sequence, unexplained density, and packing effects, are model-specific. If you use a protein model from the PDB without reading the header information, or without reading the original publications, you may be missing something vital to the appropriate use of the model. The result may be no more than a crash of your graphics software because of unexpected input like a file containing only alpha carbons. Or more seriously, you may devise and publish a detailed molecular explanation based upon a structural feature that is quite uncertain. In most cases, the PDB model isn't enough. If specific structural details of the model are crucial to a proposed mechanism or explanation, it is prudent to apply the most rigorous validation tools, and to look at the electron-density map in the important region, in order to be sure that the map is well defined there and that the model fits it well.

### 8.3 Reading a crystallography paper

### 8.3.1 Introduction

Original structure publications, especially older ones or ones in the more technical crystallographic journals, often provide enough experimental information to help with assessment of model quality. To help you learn to extract this information

from published papers, as well as to review concepts from the preceding chapters, I will walk you through a detailed structure publication. Following are annotated portions of two papers announcing the structure of adipocyte lipid binding protein (ALBP, PDB 1alb)), a member of a family of hydrophobic-ligand-binding proteins. The first paper<sup>1</sup> appeared in August 1991, announcing the purification and crystallization of the protein, and presenting preliminary results of crystallographic analysis. The second paper,<sup>2</sup> which appeared in April 1992, presented the completed structure with experimental details. In examining these papers and the model that sprang from it, I will focus primarily on the experimental and results sections of the papers and specifically upon (1) methods and concepts treated earlier in this book and (2) criteria of refinement convergence and quality of the model. Although I have reproduced parts of the published experimental procedures here (with the permission of the authors and publisher), you may wish to obtain the full publications and read them before proceeding with this example (see footnotes 1 and 2).

Readers of earlier editions may wonder why I have not replaced these papers with more recent examples. In looking for replacements, I found that most of today's structure publications in mainstream journals provide much less detail about methods, in part because structure determination is far more routine now than in the early 1990s, and in part because many educationally interesting decisions are now made by software. The following papers thus provide insights into crystallographic decision making that are often missing from more recent papers. If you can read papers like the following with insight and understanding, you will be able to extract more than from the much more abbreviated experimental details provided in more recent papers.

In the following material, sections taken from the original papers are presented in indented type. Annotations are in the usual type. For convenience, figures and tables are renumbered in sequence with those of this chapter. For access to references cited in excerpts, see the complete papers. Stereo illustrations of maps and models (not part of the papers) are derived from files kindly provided by Zhaohui Xu. I am indebted to Xu and to Leonard J. Banaszak for allowing me to use their work as an example and for supplying me with an almost complete reconstruction of this structure determination project.

### 8.3.2 Annotated excerpts of the preliminary (8/91) paper

All reprinted parts of this paper (cited in Footnote 1) appear with the permission of Professor Leonard J. Banaszak and the American Society for Biochemistry and Molecular Biology, Inc., publisher of *Journal of Biological Chemistry*.

<sup>&</sup>lt;sup>1</sup>Z. Xu, M. K. Buelt, L. J. Banaszak, and G.A. Bernlohr, Expression, purification, and crystallization of the adipocyte lipid binding protein, *J. Biol. Chem.* **266**, 14367–14370, 1991.

<sup>&</sup>lt;sup>2</sup>Z. Xu, D. A. Bernlohr, and L. J. Banaszak, Crystal structure of recombinant murine adipocyte lipid-binding protein, *Biochemistry* **31**, 3484–3492, 1992.

### Chapter 8 A User's Guide to Crystallographic Models

In the August 5, 1991, issue of Journal of Biological Chemistry, Xu, Buelt, Banaszak, and Berniohr reported the cloning, expression, purification, and crystallization of adipocyte lipid binding protein (ALBP, or rALBP for the recombinant form), along with preliminary results of crystallographic analysis. Even into the mid-1990s, this type of preliminary paper sometimes appeared as soon as a research team had carried a structure project far enough to know that it promised to produce a good model. An important aim of announcing that work was in progress on a molecule was to avoid duplication of effort in other laboratories. Although one might cynically judge that such papers constitute a defense of territory, and a grab for priority in the work at hand, something much more important was at stake. Crystallographic structure determination is a massive and expensive undertanding. The worldwide resources, both equipment and qualified scientists, for structure determination were, and in some respects still are, inadequate for the many molecules we would like to understand. Duplication of effort on the same molecule squanders limited resources in this important field. So generally, as soon as a team had good evidence that they could produce a structure, they alerted the crystallographic community to prevent parallel work from beginning in other labs. Because structure determination has become so much more rapid, it is now common for all elements of a structure project-discovery, expression, purification, crystallization, data collection, structure determination, and structure analysis-to be reported in a single paper.

The following paragraph is an excerpt from the preliminary (8/91) paper, "Experimental Procedures" section:

Crystallization—Small crystals  $(0.05 \times 0.1 \times 0.1 \text{ mm})$  were obtained using the hanging drop/vapor equilibrium method (18). 10-µl drops of 2.5 mg/ml ALBP in 0.05 M Tris, 60% ammonium sulfate, 1 mM EDTA, 1 mM dithiothreitol, 0.05% sodium azide buffer with a pH of 7.0 (crystallization buffer) were suspended over wells containing the same buffer with varying concentrations of ammonium sulfate, from 75 to 85% saturation. Small, well shaped crystals were formed within a month at an 80% saturation and 19°C. These crystals were isolated, washed with mother liquid, and used as seeds by transferring them into a 10-µl drop of 4 mg/ml fresh ALBP in the 80% saturation crystallization buffer over a well containing the same buffer. Large crystals,  $0.3 \times 0.4 \times 0.4$  mm, grew in 2 days at a constant temperature of 19°C.

The precipitant used here is ammonium sulfate, which precipitates proteins by salting out. Notice that Xu and coworkers tried a range of precipitant concentrations, probably after preliminary trials over a wider range. Crystals produced by the hanging drop method (Sec. 3.3.2, p. 38) were too small for X-ray analysis but were judged to be of good quality. The small crystals were used as seeds on which to grow larger crystals under the same conditions that produced the best small crystals. This method, called *repeated seeding*, was also discussed in Chapter 3. The initial unseeded crystallization probably fails to produce large crystals because

many crystals form at about the same rate, and soluble protein is depleted before any crystals become large. The seeded crystallization is probably effective because it decreases the number of sites of crystal growth, causing more protein to go into fewer crystals. Notice also how much faster crystals grow in the seeded drops (2 days) than in the unseeded (1 month). The preformed crystals provide nucleation sites for immediate further growth, whereas the first crystals form by random nucleation events, which are usually rate-limiting in unseeded crystallizations.

Data Collection and Processing—Crystals were analyzed with the area detector diffractometer from Siemens/Nicolet. A 0.8-mm collimator was used, and the crystal to detector distance was set at 12 cm with the detector midpoint at  $2\theta = 15^{\circ}$ . One  $\phi$  scan totaling 90° and three  $\Omega$  scans of 68° with  $\chi$  at 45° were collected with the Rigaku Ru200 operating at 50 kV and 180 mA. Each frame consisted of a 0.25° rotation taken for 120 s. The diffractometer data were analyzed with the Xengen package of programs (19). Raw data within 50 frames were searched to find about 100 strong reflections which were then indexed, and the cell dimensions were refined by least squares methods. Data from different scans were integrated separately and then merged together.

The angles  $\phi$ ,  $\chi$ ,  $\omega$ , and 26 refer to the diffractometer angles shown in Fig. 4.26, p. 82. The Rigaku Ru200 is the X-ray source, a rotating-anode tube. Each frame of data collection is, in essence, one electronic film on which are recorded all reflections that pass through the sphere of reflection during a 0.25° rotation of the crystal. This rotation size is chosen to collect as many reflections as possible without overlap. As mentioned in Chapter 4, diffractometer measurements are almost fully automated. In this instance, cell dimensions were worked out by a computer program that finds 100 strong reflections and indexes them. Then the program employs a least-squares routine (Sec. 7.6.1, p. 153) to refine the unit-cell dimensions, by finding the cell lengths and angles that minimize the difference between the actual positions of the 100 test reflections and the positions of the same reflections as calculated from the current trial set of cell dimensions. (Least-squares procedures are used in many areas of crystallography in addition to structure refinement.) Using accurate cell dimensions, the program indexed all reflections, and then integrated the X-ray counts received at each location to obtain reflection intensities.

The following excerpt is from the "Results and Discussion" section of the 8/91 paper:

Crystallization experiments using rALBP were immediately successful. With seeding, octahedral crystals of the *apo*-protein grew to a length of 0.4 mm and a height of 0.3 mm. These crystals give diffraction data to 2.4 Å. An entire data set was collected to 2.7-Å resolution using the area detector system. Statistical details of the combined X-ray data set are presented in Table 8.1.

| 0.0426 |
|--------|
| 2.2 Å  |
| 20,478 |
| 5,473  |
| 4.0    |
| 98     |
| 36     |
|        |

TABLE 8.1 
X-Ray Data Collection Statistics for Crystalline ALBP

From Z. Xu et al. (1991) J. Biol. Chem. 266, pp. 14367-14370, with permission.

Xu and colleagues had exceptionally good fortune in obtaining crystals of the recombinant form of ALBP. Efforts to crystallize a desirable protein can give success in a few days, or never, or anything in between. Modern commercial crystallization screens have made quick success more common. The extent of diffraction in preliminary tests (2.4 Å) is a key indicator that the crystals might yield a high-quality structure.

Table 8.1 provides you with a glimpse into the quality of the native data set. The 0.25° frames of data from the area detector are merged into one data set by multiplying all intensities in each frame by a scale factor. A least-squares procedure determines scale factors that minimize the differences between intensities of identical reflections observed on different frames. The merging *R*-factor [see Eq. (7.13)] gives the level of agreement among the different frames of data after scaling. In this type of *R*-factor,  $|\mathbf{F}_{obs}|_s$  are derived from averaged, scaled intensities for all observations of one reflection, and corresponding  $|\mathbf{F}_{calc}|_s$  are derived from scaled intensities for individual observations of the same reflection. The better the agreement between these two quantities throughout the data set, the lower the merging *R*-factor. In this case, individual scaled intensities agree with their scaled averages to within about 4%.

You can see from Table 8.1 that 98% of the reflections available out to 2.7 Å [those lying within a sphere of radius 1/(2.7 Å) centered at the origin of the reciprocal lattice] were measured, and on the average, each reflection was measured four times. Additional reflections were measured out to 2.4 Å. The number of available reflections increases with the third power of the radius of the sampled region in the reciprocal lattice (because the volume of a sphere of radius *r* is proportional to  $r^3$ ), so a seemingly small increase in resolution from 2.7 to 2.4 Å requires 40% more data. [Compare (1/2.4)<sup>3</sup> with (1/2.7)<sup>3</sup>.] For a rough calculation of the number of available reflections at specified resolution, see annotations of the 4/92 paper.

The lattice type was orthorhombic with unit cell dimension of a = 34.4 Å, b = 54.8 Å, c = 76.3 Å. The X-ray diffraction data were examined for systematic absences to determine the space group. Such absences were observed along the **a**<sup>\*</sup>, **b**<sup>\*</sup>, and **c**<sup>\*</sup> axes. Only reflections with *h*, *k*, or l = 2n were observed along the reciprocal

196

axes. This indicated that the space group is  $P_{21}_{21}_{21}$  (25). A unit cell with the dimensions described above has a volume of  $1.44 \times 10^5$  Å<sup>3</sup>. Assuming that half of the crystal volume is water, the volume of protein is approximately  $7.2 \times 10^4$  Å<sup>3</sup>. Considering the space group here, the volume of protein in 1 asymmetric unit would be  $1.8 \times 10^4$  Å<sup>3</sup>. By averaging the specific volume of constituent amino acids, the specific volume of ALBP is 0.715 mL/g. This led to the conclusion that the molecular mass of ALBP is approximately 15 kDa, there is only 1 molecule of ALBP in an asymmetric unit.

Recall from Sec. 5.4.3, p. 105, that for a twofold screw axis along the **b** edge, all odd-numbered 0*k*0 reflections are absent. In the space group  $P2_12_12_1$ , the unit cell possesses twofold screw axes on all three edges, so odd-numbered reflections on all three principal axes of the reciprocal lattice (*h*00, 0*k*0, and 00*l*) are missing. The presence of only even-numbered reflections on the reciprocal-lattice axes announces that the ALBP unit cell has  $P2_12_12_1$  symmetry.

The number of molecules per asymmetric unit can be determined from unit-cell dimensions and a rough estimate of the protein/water ratio. Since this number is an integer, even a rough calculation can give a reliable answer. The assumption that ALBP crystals are 50% water is no more than a guess taken from near the middle of the range for protein crystals (30 to 78%). The unit-cell volume is  $(34.4 \text{ Å}) \times (54.8 \text{ Å}) \times (76.3 \text{ Å}) = 1.44 \times 10^5 \text{ Å}^3$ , and if half that volume is protein, the protein volume is  $7.2 \times 10^4$  Å<sup>3</sup>. In space group  $P2_12_12_1$ , there are four equivalent positions (Sec. 4.2.8, p. 65), so there are four asymmetric units per unit cell. Each one must occupy one-fourth of the protein volume, so the volume of the asymmetric unit is one-fourth of  $7.2 \times 10^4$ , or  $1.8 \times 10^4$  Å<sup>3</sup>. The stated specific volume (volume per gram) of the protein is the weighted average of the specific volumes of the amino-acid residues (which can be looked up), weighted according to the amino-acid composition of ALBP. The molecular mass of one asymmetric unit is obtained by converting the density of ALBP in grams per milliliter (which is roughly the inverse of the specific volume) to daltons per cubic angstrom, and then multiplying by the volume of the asymmetric unit, as follows:

$$\frac{1 \text{ g}}{0.715 \text{ mL}} \cdot \frac{1 \text{ ml}}{\text{cm}^3} \cdot \frac{\text{cm}^3}{(10^8)^3 \text{ Å}^3} \cdot \frac{6.02 \times 10^{23}}{\text{g}} \text{ daltons} \cdot 1.8 \times 10^4 \text{ Å}^3$$
$$= 1.5 \times 10^4 \text{ daltons}.$$

This result is very close to the known molecular mass of ALBP, so there is one ALBP molecule per asymmetric unit. This knowledge is an aid to early map interpretation.

As indicated, ALBP belongs to a family of low molecular weight fatty acid binding proteins. The sequences of the proteins in the family have been shown to be very similar and in particular in the aminoterminal domain where Y19<sup>3</sup> resides. Among them, the structure of myelin P2 and IFABP has been solved. Since the amino acid identity between ALBP and myelin P2 is about 69%, P2 should be a good starting structure to obtain phase information for ALBP using the method of molecular replacement. Preliminary solutions to the rotation and translation functions have been obtained. Seeding techniques will allow us to obtain large crystals for further study of the *holo-* and phosphorylated protein. By comparing the crystal structures of these different forms, it should be possible to structurally determine the effects of protein phosphorylation on ligand binding and ligand binding on phosphorylation.

Because ALBP is related to several proteins of known structure, molecular replacement is an attractive option for phasing. The choice of a phasing model is simple here: just pick the one with the amino-acid sequence most similar to ALBP, which is myelin P2 protein. Solution of rotation and translation functions refers to the search for orientation and position of the phasing model (P2) in the unit cell of ALBP (Section 6.5.4, p. 139). The subsequent paper provides more details.

## 8.3.3 Annotated excerpts from the full structure-determination (4/92) paper

All reprinted parts of this paper (cited in Footnote 2) appear with the permission of Professor Leonard J. Banaszak and the American Chemical Society, publisher of *Biochemistry*.

In April 1992, the structure determination paper appeared in *Biochemistry*. This paper contains a full description of the experimental work, and a complete analysis of the structure. The following is from the 4/92 paper, "Abstract" section:

Adipocyte lipid-binding protein (ALBP) is the adipocyte member of an intracellular hydrophobic ligand-binding protein family. ALBP is phosphorylated by the insulin receptor kinase upon insulin stimulation. The crystal structure of recombinant murine ALBP has been determined and refined to 2.5 Å. The final *R*-factor for the model is 0.18 with good canonical properties.

A 2.5-Å model refined to an *R*-factor of 0.18 should be a detailed model. "Good canonical properties" means good agreement with accepted values of bond lengths, bond angles, and planarity of peptide bonds.

The following excerpts are from the "Materials and Methods" section of the 4/92 paper:

*Crystals and X-ray Data Collection*. Detailed information concerning protein purification, crystallization, and X-ray data collection can be

 $<sup>^{3}</sup>$ Y19 is tyrosine 19, a residue considered important to the function of ALBP.

found in a previous report (Xu *et al.*, 1991) and will be mentioned here in summary form. Recombinant murine *apo*-ALBP crystallizes in the orthorhombic space group  $P_{21}2_{1}2_{1}$  with the following unit cell dimensions a = 34.4 Å, b = 54.8 Å, and c = 76.3 Å. The asymmetric unit contains one molecule with a molecular weight of 14,500. The entire diffraction data set was collected on one crystal. In the resolution range  $\infty$ -2.5 Å, 5115 of the 5227 theoretically possible reflections were measured. Unless otherwise noted the diffraction data with intensities greater than  $2\sigma$  were used for structure determination and refinement. As can be seen in Table 8.2, this included about 96% of the measured data.

This section reviews briefly the results of the preliminary paper. Full data collection from a single crystal was relatively rare at the time of this paper, but with cryocrystallography and more powerful X-ray sources, it is almost the rule today. In the early stages of the work, reflections weaker than two times the standard deviation for all reflections  $(2\sigma)$  were omitted from Fourier syntheses, because of greater uncertainty in the measurements of weak reflections. Table 8.2 is discussed later. The diffractometer software computes the number of reflections available at 2.5-Å resolution by counting the number of reciprocal-lattice points that lie within a sphere of radius [1/(2.5 Å)], centered at the origin of the reciprocal lattice. This number is roughly equal to the number of reciprocal unit cells within the 1/(2.5 Å)sphere, which is, again roughly, the volume of the sphere  $(V_{rs})$  divided by the volume of the reciprocal unit cell  $(V_{\rm rc})$ . The volume of the reciprocal unit cell is the inverse of the real unit-cell volume V. So the number of reflections available at 2.5-Å resolution is approximately  $(V_{rs}) \cdot (V)$ . Because of the symmetry of the reciprocal lattice and of the  $P2_12_12_1$  space group, only one-eighth of the reflections are unique (Sec. 4.3.7, p. 88). So the number of unique reflections is approximately  $(V_{\rm rs}) \cdot (V)/8$ , or

$$\frac{\frac{4}{3}\pi \left(\frac{1}{2.5 \text{ Å}}\right)^3 \left(1.44 \times 10^5 \text{ Å}^3\right)}{8} = 4825 \text{ reflections}$$

The 8% difference between this result and the stated 5227 reflections is due to the approximations made here and to the sensitivity of the calculation to small round-off in unit-cell dimensions.

*Molecular Replacement.* The tertiary structure of crystalline ALBP was solved by using the molecular replacement method incorporated into the XPLOR computer program (Brunger *et al.*, 1987). The refined crystal structure of myelin P2 protein without solvent and fatty acid was used as the probe structure throughout the molecular replacement studies. We are indebted to Dr. A. Jones and his colleagues for permission to use their refined P2 coordinates before publication.

Myelin P2 coordinates were not yet available from the Protein Data Bank and were obtained directly from the laboratory in which the P2 structure was determined. In this project, the search for the best orientation and position of P2 in the ALBP unit cell was divided into three parts: a rotation search to find promising orientations, refinement of the most promising orientations to find the best orientation, and a translation search to find the best position. Here are the details of the search:

(1) *Rotation Search*. The rotation search was carried out using the Patterson search procedures in XPLOR. The probe Patterson maps were computed from structure factors calculated by placing the P2 coordinates into an orthorhombic cell with 100-Å edges. One thousand highest Patterson vectors in the range of 5–15 Å were selected and rotated using the pseudoorthogonal Eulerian angles ( $\theta_+$ ,  $\theta_2$ ,  $\theta_-$ ) as defined by Lattman (1985). The angular search interval for  $\theta_2$  was set to 2.5°; intervals for  $\theta_+$  and  $\theta_-$  are functions of  $\theta_2$ . The rotation search was restricted to the asymmetric unit  $\theta_- = 0 - 180^\circ$ ,  $\theta_2 = 0-90^\circ$ ,  $\theta_+ = 0-720^\circ$  for the  $P2_12_12_1$  space group (Rao, *et al.*, 1980). XPLOR produces a sorted list of the correlation results simplifying final interpretation (Brunger 1990).

XPLOR is a package of refinement programs that includes powerful procedures for energy refinement by simulated annealing, in addition to more traditional tools like least-squares methods and molecular-replacement searches. The package is available for use on many different computer systems.

The P2 phasing model is referred to here as the probe. For the rotation search, the probe was placed in a unit cell of arbitrary size and  $\mathbf{F}_{calc}$ s were obtained from this molecular model, using Eq. (5.15). Then a Patterson map was computed from these  $\mathbf{F}_{calc}$ s using Eq. (6.10), p. 125. Recall that Patterson maps reflect the molecule's orientation but not its position. All peaks in the Patterson map except the strongest 1000 were eliminated. Then the resulting simplified map was compared to a Patterson map calculated from ALBP reflection intensities. The probe Patterson was rotated in a three-dimensional coordinate system to find the orientation that best fit the ALBP Patterson, as illustrated in Fig. 6.18, p. 141 (the system of angles employed here and that shown in the figure are just two of at least nine different ways of specifying angles in a rotation search). The search was monitored by a rotation function like Eq. (6.17), p. 142, producing a plot of the angles versus a criterion of coincidence between peaks in the two Patterson maps (see page 139). Peaks in the rotation function occur at sets of angles where many coincidences occur. The coincidences are not perfect because there is a finite interval between angles tested, and the exact desired orientation is likely to lie between test angles. The interval is made small enough to avoid missing promising orientations altogether.

(2) Patterson Correlation Refinement. To select which of the orientations determined from the rotation search is the correct solution a Patterson correlation refinement of the peak list of the rotation function was performed. This was carried out by minimization against a target

function defined by Brunger (1990) and as implemented in XPLOR. The search model, P2, was optimized for each of the selected peaks of the rotation function.

As discussed later in the "Results" section, the rotation function contains many peaks. The strongest 100 peaks were selected, and each orientation was refined by least squares to produce the best fit to the ALBP Patterson map. For each refined orientation, a correlation coefficient was computed. The orientation giving the highest correlation coefficient was chosen as the best orientation for the phasing model.

(3) *Translation Search*. A translation search was done by using the P2 probe molecule oriented by the rotation function studies and refined by the Patterson correlation method. The translation search employed the standard linear correlation coefficient between the normalized observed structure factors and the normalized calculated structure factors (Funinaga & Read, 1987; Brunger, 1990). X-ray diffraction data from 10–3 Å resolution were used. Search was made in the range x = 0.05, y = 0.05, and z = 0.05, with the sampling interval 0.0125 of the unit cell length.

The last step in molecular replacement was to find the best position for the probe molecule in the ALBP unit cell. The P2 orientation obtained from the rotation search and refinement was tried in all unique locations at intervals of one-eightieth (0.0125) of the unit-cell axis lengths. The symmetry of the  $P_{21}_{21}_{21}$  unit cell allowed this search to be confined to the region bound by one-half of each cell axis. The total number of positions tested was thus (40)(40)(40) or 64,000. For each position,  $\mathbf{F}_{calc}$ s were computed [Eq. (5.15)] from the P2 model and their amplitudes are compared with the  $|\mathbf{F}_{obs}|$ s from the ALBP native data set. An unspecified correlation coefficient, probably similar to an *R*-factor, was computed for each P2 position, and the position that gave P2  $|\mathbf{F}_{calc}|$ s in best agreement with ALBP  $|\mathbf{F}_{obs}|$ s was chosen as the best position for P2 as a phasing model. The starting phase estimates for the refinement were thus the phases of  $\mathbf{F}_{calc}$ s computed [Eq. (5.15)] from P2 in the final orientation and position determined by the three-stage molecular-replacement search.

*Structure Refinement.* The refinement of the structure was based on an energy function approach (Brunger *et al.*, 1987): arbitrary combinations of empirical and effective energy terms describing crystallographic data as implemented in XPLOR. Molecular model building was done on an IRIS Workstation (Silicon Graphics) with the software TOM, a version of FRODO (Jones, 1978).

The initial model of ALBP was built by simply putting the amino acid sequence of ALBP into the molecular structure of myelin P2 protein. After a 20-step rigid-body refinement of the positions and orientations of the molecule, crystallographic refinement with simulated annealing was carried out using a slow-cooling protocol (Brunger *et al.*, 1989,



**Figure 8.4**  $\blacktriangleright$  ALPB electron-density map calculated with molecular-replacement phases before any refinement, shown with the final model. Compare with Fig. 2.3, page 11, which shows the final, much improved, electron-density map in the same region.

1990). Temperature factor refinement of grouped atoms, one for backbone and one for side-chain atoms for each residue, was initiated after the R-factor dropped to 0.249.

The first electron-density map was computed [Eq. (7.03), p. 150] with  $|\mathbf{F}_{obs}|s$  from the ALBP data set and  $\alpha_{calc}s$  from the oriented P2 molecule. Fig. 8.4 shows a small section of this map superimposed upon the *final* model.

An early map like Fig. 8.4, computed from initial phase estimates, harbors many errors, where the map does not agree with the model ultimately derived from refinement. In this section, you can see both false breaks and false connections in the density. For example, there are breaks in density at C- $\beta$  of the phenylalanine residue (side chain ending with six-membered ring) on the right, and along the protein backbone at the upper left. The lobe of density corresponding to the valine side chain (center front) is disconnected and out of place. There is a false connection between density of the carbonyl oxygen (red) at lower left and side chain density above. Subsequent refinement was aimed at improving this map.

Next, the side chains of P2 were replaced with the side chains of ALBP at corresponding positions in the amino-acid sequence to produce the first ALBP model. The position and orientation of this model were refined by least squares, treating the model as a rigid body. Subsequent refinement was by simulated annealing. At first, all temperature factors were constrained at 15.0 Å<sup>2</sup>. After the first round of simulated annealing, temperature factors were allowed to refine for atoms in groups, one value of *B* for all backbone atoms within a residue and another for side-chain atoms in the residue.

The new coordinates were checked and adjusted against a  $(2 |\mathbf{F}_0| - |\mathbf{F}_c|)$  and a  $(|\mathbf{F}_0| - |\mathbf{F}_c|)$  electron density map, where  $|\mathbf{F}_0|$  and  $|\mathbf{F}_c|$  are the observed and calculated structure factor amplitudes. Phases are calculated from the crystal coordinates. The Fourier maps were

calculated on a grid corresponding to one-third of the high-resolution limit of the input diffraction data. All residues were inspected on the graphics system at several stages of refinement. The adjustments were made on the basis of the following criteria: (a) that an atom was located in low electron density in the  $(2 |\mathbf{F}_0| - |\mathbf{F}_c|)$  map or negative electron density in the  $(|\mathbf{F}_0| - |\mathbf{F}_c|)$  map; (b) that the parameters for the  $\Phi$ ,  $\Psi$  angles placed the residue outside the acceptable regions in the Ramachandran diagram. Iterative refinement and model adjustment against a new electron density map was carried out until the *R*-factor appeared unaffected. Isotropic temperature factors for individual atoms were then included in the refinement.

In between rounds of computerized refinement, maps were computed using  $|\mathbf{F}_{obs}|_{s}$  from the ALBP data set and  $\alpha_{calc}s$  from the current model [taken from  $|\mathbf{F}_{calc}|_{s}$  computed by Eq. (5.15)]. The model was corrected where the fit to maps was poor, or where the Ramachandran angles  $\Phi$  and  $\Psi$  were forbidden. Notice that the use of  $2\mathbf{F}_{o} - \mathbf{F}_{c}$  and  $\mathbf{F}_{o} - \mathbf{F}_{c}$  maps [Eq. (7.5) and Eq. (7.4)] is as described in Sec. 7.4.2, p. 154. When alternating rounds of refinement and map fitting produced no further improvement in *R*-factor, temperature factors for each atom were allowed to refine individually, leading to further decrease in *R*.

The next stage of the crystallographic study included the location of solvent molecules. They were identified as well-defined peaks in the electron-density maps within hydrogen-bonding distance of appropriate protein atoms or other solvent atoms. Solvent atoms were assigned as water molecules and refined as oxygen atoms. Those that refined to positions too close to other atoms, ended up located in low electron density, or had associated temperature factors greater than 50 Å<sup>2</sup> were removed from the coordinate list in the subsequent stage. The occupancy for all atoms, including solvent molecules, was kept at 1.0 throughout the refinement. Detailed progress of the crystallographic refinement is given in Table 8.2, p. 204.

Finally, ordered water molecules were added to the model where unexplained electron-density was present in chemically feasible locations for water molecules. Temperature factors for these molecules (treated as oxygen atoms) were allowed to refine individually. If refinement moved these molecules into unrealistic positions or increased their temperature factors excessively, the molecules were deleted from the model. Occupancies were constrained to 1.0 throughout the refinement. This means that *B* values reflect both thermal motion and disorder (Sec. 8.2.3, p. 185). Because all *B* values fall into a reasonable range, the variation in *B* can be attributed to thermal motion. Table 8.2 (p. 204) shows the progress of the refinement. Note that *R* drops precipitously in the first stages of refinement after ALBP side chains replace those of P2. Note also that *R* and the deviations from ideal bond lengths, bond angles, and planarity of peptide bonds decline smoothly throughout the later stages of refinement. The small increase in *R* at the end is due to inclusion of weaker reflections in the final round of simulated annealing.

|                    |                       |          |                   |                  | RN                    | AS deviat              | ions               |
|--------------------|-----------------------|----------|-------------------|------------------|-----------------------|------------------------|--------------------|
| Stage <sup>a</sup> | Number of reflections | R-factor | $B(\text{\AA}^2)$ | Solvent included | Bond<br>length<br>(Å) | Bond<br>angle<br>(deg) | Planarity<br>(deg) |
| 1                  | 2976                  | 0.458    | 15.0              |                  | 0.065                 | 4.12                   | 9.015              |
| 2                  | 2976                  | 0.456    | 15.0              |                  | 0.065                 | 4.12                   | 9.012              |
| 3                  | 4579                  | 0.235    | Group             |                  | 0.019                 | 3.17                   | 1.506              |
| 4                  |                       |          |                   |                  |                       |                        |                    |
| 5                  | 4579                  | 0.220    | Indiv.            |                  | 0.018                 | 3.77                   | 1.408              |
| 6                  |                       |          |                   |                  |                       |                        |                    |
| 7                  | 4579                  | 0.197    | Indiv.            | 31               | 0.018                 | 3.73                   | 1.366              |
| 8                  |                       |          |                   |                  |                       |                        |                    |
| 9                  | 4579                  | 0.172    | Indiv.            | 88               | 0.016                 | 3.47                   | 1.139              |
| 10                 |                       |          |                   |                  |                       |                        |                    |
| 11                 | 4773                  | 0.183    | Indiv.            | 69               | 0.017                 | 3.46                   | 1.070              |

TABLE 8.2 ► Progress of Refinement

Reprinted with permission from Z. Xu et al. (1992). Biochemistry 31, 3484–3492. Copyright 1992 American Chemical Society.

<sup>a</sup>Key to stages of refinement:

Stage Action

1

tage Action Starting model

2 Rigid-body refinement

3 Simulated annealing

4 Model rebuilt using  $(2F_o - F_c)$  and  $(F_o - F_c)$  electron-density maps

5 Simulated annealing

6 Model rebuilt using  $(2\mathbf{F}_0 - \mathbf{F}_c)$  and  $(\mathbf{F}_0 - \mathbf{F}_c)$  electron-density maps, H<sub>2</sub>O included

7 Simulated annealing

8 Model rebuilt using  $(2\mathbf{F}_0 - \mathbf{F}_c)$  and  $(\mathbf{F}_0 - \mathbf{F}_c)$  electron-density maps, H<sub>2</sub>O included

9 Simulated annealing

10 Model rebuilt using  $(2\mathbf{F}_o - \mathbf{F}_c)$  and  $(\mathbf{F}_o - \mathbf{F}_c)$  electron-density maps, H<sub>2</sub>O included

11 Simulated annealing

The following excerpts are from the "Results" section of the 4/92 paper:

*Molecular Replacement.* From the initial rotation search, the 101 highest peaks were chosen for further study. These are shown in Fig. 8.4, p. 206. The highest peak of the rotation function had a value 4.8 times the standard deviation above the mean and 1.8 times the standard deviation above the next highest peak. The orientation was consistently the highest peak for diffraction data within the resolution ranges 10–5, 7–5, and 7–3 Å. Apart from peak number 1, six strong peaks emerged after PC<sup>4</sup> refinement, as can be seen in Fig. 8.4*b*, p. 206. These peaks all corresponded to approximately the same orientation as peak number 1. Three of them were initially away from that orientation and converged to it during the PC refinement.

<sup>&</sup>lt;sup>4</sup>Patterson correlation.

A translation search as implemented in XPLOR was used to find the molecular position of the now oriented P2 probe in the ALBP unit cell. Only a single position emerged at x=0.250, y=0.425, z=0.138 with a correlation coefficient of 0.419. The initial *R*-factor for the P2 coordinates in the determined molecular orientation and position was 0.470 including X-ray data in the resolution range of 10–3 Å. A rigid-body refinement of orientation and position reduced the starting *R*-factor to only 0.456, probably attesting to the efficacy of the Patterson refinement in XPLOR.

In Fig. 8.5*a*, the value of the rotation function, which indicates how well the probe and ALBP Patterson maps agree with each other, is plotted vertically against numbers assigned to the 101 orientations that produced best agreement. Then each of the 101 orientations were individually refined further, by finding the nearby orientation having maximum value of the rotation function. In some cases, different peaks refined to the same final orientation. Each refined orientation of the probe received a correlation coefficient that shows how well it fits the Patterson map of ALBP. The orientation receiving the highest correlation coefficient was taken as the best orientation of the probe, and then used to refine the position of the probe in the ALBP unit cell. The orientation and position of the model obtained from the molecular replacement search was so good that refinement of the model as a rigid body produced only slight improvement in R. The authors attribute this to the effectiveness of the Patterson correlation refinement of model orientation, stage two of the search.

*Refined Structure of apo-ALBP*. The refined ALBP structure has a *R*-factor 0.183 when all observed X-ray data (4773 reflections) between 8.0 and 2.5 Å are included. The rms deviation of bond lengths, bond angles, and planarity from ideality is 0.017 Å, 3.46°, and 1.07°, respectively. An estimate of the upper limit of error in atomic coordinates is obtained by the method of Luzzati (1952). Figure 8.6 summarizes the overall refined model.

The plot presented in Fig. 8.6*a* suggests that the upper limit for the mean error of the refined ALBP coordinates is around 0.34 Å. The mean temperature factors for main-chain and side-chain atoms are plotted in Fig. 8.6*b*.

The final *R*-factor and structural parameters exceed the standards described in Sec. 8.1, p. 179 and attest to the high quality of this model. Atom locations are precise to an average of 0.34 Å, about one-fifth of a carbon-carbon covalent bond length. The plot of temperature factors shows greater variability and range for side-chain atoms, as expected, and shows no outlying values. The model defines the positions of all amino-acid residues in the protein.

Careful examination of  $(2 |\mathbf{F}_0| - |\mathbf{F}_c|)$  and  $(|\mathbf{F}_0| - |\mathbf{F}_c|)$  maps at each refinement step led to the conclusion that no bound ligand was present. There was no continuous positive electron density present near the ligand-binding site as identified in both P2



**Figure 8.5** Rotation function results: P2 into crystalline ALBP. (*a*) Plot of the 101 best solutions to the rotation function, each peak numbered in the horizontal direction (abscissa). The correlation between the Patterson maps of the probe molecule and the measured ALBP X-ray results are shown in the vertical direction (ordinate) and are given in arbitrary units. (*b*) Description of the rotation studies after Patterson correlation refinement. The peak numbers plotted in both panels (*a*) and (*b*) are the same. Reprinted with permission from Z. Xu *et al.* (1992) *Biochemistry* **31**, 3484–3492. Copyright 1992 American Chemical Society.

206



**Figure 8.6** ALBP refinement results. (*a*) Theoretical estimates of the rms positional errors in atomic coordinates according to Luzzati are shown superimposed on the curve for the ALBP diffraction data. The coordinate error estimated from this plot is 0.25 Å with an upper limit of about 0.35 Å.(*b*) Mean values of the main-chain and side-chain temperature factors are plotted versus the residue number. The temperature factors are those obtained from the final refinement cycles. Reprinted with permission from Z. Xu *et al.* (1992) *Biochemistry* **31**, 3484–3492. Copyright 1992 American Chemical Society.

(Jones, *et al.*, 1988) and IFABP (Sacchettini *et al.*, 1989a). The absence of bound fatty acid in crystalline ALBP is consistent with the chemical modification experiment which indicates ALBP purified from *E. coli* is devoid of fatty acid (Xu *et al.*, 1991). The final refined coordinate list includes 1017 protein atoms and 69 water molecules.

The final maps exhibit no unexplained electron density. This is of special concern because ALBP is a ligand-binding protein (its ligand is a fatty acid), and ligands sometimes survive purification and crystallization, and are found in the final electron-density map. It is implied by references to *apo*-protein and *holo*protein that attempts to determine the structure of an ALBP-ligand complex were under way at the time of publication. If it is desired to detect conformational changes upon ligand binding, then it is crucial to know that no ligand is bound to this *apo*-protein, so that conformational differences between *apo*- and *holo*-forms, if found, can reliably be attributed to ligand binding.

To compare *apo*- and *holo*-forms of proteins after both structures have been determined independently, crystallographers often compute difference Fourier syntheses (Sec. 7.4.2, p. 154), in which each Fourier term contains the structure-factor difference  $\mathbf{F}_{holo} - \mathbf{F}_{apo}$ . A contour map of this Fourier series is called a difference map, and it shows only the differences between the *holo*- and *apo*-forms. Like the  $\mathbf{F}_{o} - \mathbf{F}_{c}$  map, the  $\mathbf{F}_{holo} - \mathbf{F}_{apo}$  map contains both positive and negative density. Positive density occurs where the electron density of the *holo*-form is greater than that of the *apo*-form, so the ligand shows up clearly in positive density. In addition, conformational differences between *holo* and *apo*-forms result in positive density where *holo*-protein atoms occupy regions that are unoccupied in the *apo*-form. The standard interpretation of such a map is that negative density locates the same atoms after ligand binding. In regions where the two forms are identical,  $\mathbf{F}_{holo} - \mathbf{F}_{apo} = 0$ , and the map is blank.

Structural Properties of Crystalline ALBP. A Ramachandran plot of the main chain dihedral angle  $\Phi$  and  $\Psi$  is shown in Fig. 8.7. In the refined model, 13 residues have positive  $\Phi$  angles, 9 of which belong to glycine residues. There are 11 glycine residues in ALBP, all associated with good quality electron density.

Most of the residues having forbidden values of  $\Phi$  and  $\Psi$  are glycines, represented by small squares in Fig. 8.7, whereas all other amino acids are represented by + symbols. Succeeding discussion reveals that these unusual conformations are also found in P2 and other members of this protein family, strengthening the argument that these conformations are not errors in the model, and suggesting that they might be important to structure and/or function in this family of proteins.

Figure 2.3, p. 11 shows, at the end of refinement, the same section of map as in Fig. 8.4. By comparing Fig. 2.3 with Fig. 8.4, you can see that the map errors



**Figure 8.7** Ramachandran plot of the crystallographic model of ALBP, generated with DeepView (see Chapter 11). The main-chain torsional angle  $\Phi$  (N-C<sub> $\alpha$ </sub> bond) is plotted versus  $\Psi$  (C-C<sub> $\alpha$ </sub> bond). The following symbols are used: (+) nonglycine residues; ( $\Box$ ) glycine residues. The enclosed areas of the plot show sterically allowed angles for nonglycine residues. The symbols are colored according to their inclusion in secondary structural elements: red, alpha helix; yellow, beta sheet; gray, coil.

described earlier were eliminated, and that the map is a snug fit to a chemically, stereochemically, and conformationally realistic model.

### 8.4 Summary

All crystallographic models are *not* equal. The noncrystallographer can assess model quality by carefully reading original publications of a macromolecular structure and by using the latest online validation tools. The kind of reading and

### Chapter 8 A User's Guide to Crystallographic Models

interpretation implied by my annotations is essential to wise use of models. Don't get me wrong. There is no attempt on the part of crystallographers to hide the limitations of models. On the contrary, refinement papers often represent almost heroic efforts to make plain what the final model says and leaves unsaid. These efforts are in vain if the reader does not understand them, or worse, never reads them. These efforts are often undercut by the simple power of the visual model. The brightly colored stereo views of a protein model, which are in fact more akin to cartoons than to molecules, endow the model with a concreteness that exceeds the intentions of the thoughtful crystallographer. It is impossible for the crystallographer, with vivid recall of the massive labor that produced the model, to forget its shortcomings. It is all too easy for users of the model to be unaware of them. It is also all too easy for the user to be unaware that, through temperature factors, occupancies, undetected parts of the protein, and unexplained density, crystallography reveals more than a single molecular model shows.

Even the highest-quality model does not explain itself. If I showed you a perfect model of a protein of unknown function, it is unlikely that you could tell me what it does, or even pinpoint the chemical groups critical to its action. Using a model to explain the properties and action of a protein means bringing the model to bear upon all the other available evidence. This involves gaining intimate knowledge of the model, a task roughly as complex as learning your way around a small city. In Chapter 11, I will discuss the exploration of macromolecular models by computer graphics. But first, I must carry out two other tasks. In the next chapter, I will build on your understanding of X-ray diffraction to introduce you to other means of structure determination using diffraction, including X-ray diffraction of fibers and powders, and diffraction by neutrons and electrons. Then, in Chapter 10, I will briefly introduce other, noncrystallographic, methods of structure determination, in particular NMR and homology modeling. With each method, I will discuss how to assess the quality of the resulting models, by analogy with the criteria described in this chapter.

### ► Chapter 9

# **Other Diffraction Methods**

### 9.1 Introduction

The same principles that underlie single-crystal X-ray crystallography make other kinds of diffraction experiments understandable. In this chapter, I provide brief, qualitative descriptions of other diffraction methods. First is X-ray diffraction by fibers rather than crystals, which reduces reciprocal space from three dimensions to two. Next we go down to one dimension, with diffraction by amorphous materials like powders and solutions. Then I will look at diffraction using other forms of radiation, specifically neutrons and electrons. I will show how each type of diffraction experiment provides structural information that can complement or supplement information from single-crystal X-ray diffraction. Finally, returning to X-rays and crystals, I will discuss Laue diffraction, which allows the collection of a full diffraction data set from a single brief exposure of a crystal to polychromatic X-radiation. Laue diffraction opens the door to time-resolved crystallography, yielding crystallographic models of intermediate states in chemical reactions. For each of the crystallographic methods, I provide references to research articles that exemplify the method.

### 9.2 Fiber diffraction

Many important biological substances do not form crystals. Among these are most membrane proteins and fibrous materials like collagen, DNA, filamentous viruses, and muscle fibers. Some membrane proteins can be crystallized in matrices of lipid and studied by X-ray diffraction (Sec. 3.3.4, p. 41), or they can be incorporated into lipid films (which are in essence two-dimensional crystals) and studied by

electron diffraction. I will discuss electron diffraction later in this chapter. Here I will examine diffraction by fibers.

Like crystals, fibers are composed of molecules in an ordered form. When irradiated by an X-ray beam perpendicular to the fiber axis, fibers produce distinctive diffraction patterns that reveal their dimensions at the molecular level. Because many fibrous materials are polymeric and of known chemical composition and sequence, their molecular dimensions are sometimes all that is needed to build a feasible model of their structure.

Some materials (for example, certain muscle proteins) form fibers spontaneously or are naturally found in fibrous form. Many other polymeric substances, like DNA, can be induced into fibers by pulling them from an amorphous gel with tweezers or a glass rod. For data collection, the fiber is simply suspended between a well-collimated X-ray source and a detector (Fig. 9.1).

The order in a fiber is one-dimensional (along the fiber) rather than threedimensional, as in a crystal. You can think of molecules in a fiber as being stretched



**Figure 9.1**  $\triangleright$  Fiber diffraction. Molecules in the fiber are oriented parallel to the beam axis but aligned at random. X-ray beams emerging from the fiber strike the detector in layer lines perpendicular to the fiber axis.

212

### Section 9.2 Fiber diffraction

out parallel to the fiber axis but having their termini occurring at random along the fiber, as shown in the expanded detail of the fiber in Fig. 9.1. Because the X-ray beam simultaneously sees all molecules in all possible rotational orientations about the fiber axis, Bragg reflections from a fiber are cylindrically averaged, and irradiation of the fiber by a beam perpendicular to the fiber axis gives a complete, but complex, two-dimensional diffraction pattern from a single orientation of the fiber.

Fibers can be crystalline or noncrystalline. Crystalline fibers are actually composed of long, thin microcrystals oriented with their long axis parallel to the fiber axis. When a crystalline fiber is irradiated with X-rays perpendicular to the fiber axis, the result is the same as if a single crystal were rotated about its axis in the X-ray beam during data collection, sort of like a rotation photograph taken over 360° instead of the usual very small rotation angle (Sec. 4.3.4, p. 80). All Bragg reflections are registered at once, on *layer lines* perpendicular to the fiber axis. All fiber diffraction patterns have two mirror planes, parallel (the meridian) and perpendicular (the equator) to the fiber axis. Many of the reflections overlap, making analysis of the diffraction pattern very difficult.

In noncrystalline fibers, all the *molecules* (as opposed to oriented groups of them in microcrystals) are parallel to the fiber axis but aligned along the axis at random. This arrangement gives a somewhat simpler diffraction pattern, also consisting of layer lines, but with smoothly varying intensity rather than distinct reflections. In the diffraction patterns from both crystalline and noncrystalline materials, spacings of layer lines are related to the periodicity of the individual molecules in the fiber, as I will show.

A simple and frequently occurring structural element in fibrous materials is the helix. I will use the relationship between the dimensions of simple helices and that of their diffraction patterns to illustrate how diffraction can reveal structural information. As a further simplification, I will assume that the helix axis is parallel to the fiber axis. As in all diffraction methods, the diffraction pattern is a Fourier transform of the object in the X-ray beam, averaged over all the orientations present in the sample. In the case of fibers, this means that the transform is averaged cylindrically, around the molecular axis parallel to the fiber axis.

Figure 9.2 shows some simple helices and their transforms. The transform of the helix in Fig. 9.2*a* exhibits an X pattern that is always present in transforms of helices. I will explain the mathematical basis of the X pattern later. Although each layer line looks like a row of reflections, it is actually continuous intensity. This would be apparent if the pattern were plotted at higher overall intensity. The layer lines are numbered with integers from the equator (l = 0). Because of symmetry, the first lines above *and* below the equator are labeled l = 1, and so forth.

Compare helix (*a*) with (*b*), in which the helix has the same radius, but a longer *pitch P* (peak-to-peak distance). Note that the layer lines for (*b*) are more closely spaced. The layer-line spacing is inversely proportional to the helix pitch. The relationship is identical to that in crystals between reciprocal-lattice spacing and unit-cell dimensions [(Eq. (4.10), p. 87]. As a result, precise measurement of layer-line spacing allows determination of helix pitch.



**Figure 9.2**  $\blacktriangleright$  Helices and their Fourier transforms. (*a*) Simple, continuous helix. The first intensity peaks from the centers of each row form a distinctive X pattern. (*b*) Helix with longer pitch than (*a*) gives smaller spacing between layer lines. (*c*) Helix with larger radius than (*a*) gives narrower X pattern. (*d*) Helix of same dimensions as (*a*) but composed of discrete objects gives overlapping X patterns repeated along the meridian.

### Section 9.2 Fiber diffraction

Helix (c) has the same pitch as (a) but a larger radius r. Notice that the X pattern in the transform is narrower than that of helices (a) and (b), which have the same radius. The angle formed by the branches of the X with the meridian, shown as the angle  $\delta$ , is determined by the helix radius. But it appears at first that the relationship between  $\delta$  and helix radius must be more complex, because angles  $\delta$  for helices (a) and (b), which have the same radius, appear to be different. However, we define  $\delta$  as the angle whose tangent is the distance w from the meridian to the center of the first intensity peak divided by the layer-line number l. You can see that the distance w at the tenth layer line is the same in (a) and (b). Defined in this way, and measured at relatively large layer-line numbers (beyond which the tangent of an angle is simply proportional to the length of the side opposite),  $\delta$  is inversely proportional to the helix radius. Because the layer-line spacing is the same in (a) and (c), it is clear that an increase in radius decreases  $\delta$ , and that we can determine the helix radius from  $\delta$ .

Helix (d) has the same pitch and radius as helix (a), but is a helix of discrete objects or "repeats," like a polymeric chain of repeating subunits. The transform appears at first to be far more complex, but it is actually only slightly more so. It is merely a series of overlapping X patterns distributed along the meridian of the transform. To picture how multiple X patterns arise from a helix of discrete objects, imagine that the helix beginning with arbitrarily chosen object number 1 produces the X at the center of the transform. Then imagine that the same helix beginning with object number 2 produces an X of its own. The distance along the meridian between centers of the two Xs is inversely proportional [Eq. (4.10) again] to the distance between successive discrete objects in the helix. So careful measurement of the distance between successive meridional "reflections" (remember that the intensities on a layer line are actually continuous) allows determination of the distance between successive subunits of the helix. In a polymeric helix like that of a protein or nucleic acid, this parameter is called the rise-per-residue or p. Dividing the pitch P by the rise-per-residue p gives the number of residues per helical turn (rise-per-turn divided by rise-per-residue gives residues-per-turn).

In addition, for the simplest type of discrete helix, in which there is an integral number of residues per turn of the helix, this integer P/p is the same as number l of the layer line on which the first meridional intensity peak occurs. Note that helix d contains exactly six residues per turn, and that the first meridional intensity peak above or below the center of the pattern occurs on layer line l = 6. To review, the layer-line spacing Z is proportional to 1/P, and the distance from the origin to the first meridional reflection is proportional to 1/p.

If the number of residues per turn is not integral, then the diffraction pattern is much more complex. For example, a protein alpha helix has 3.6 residues per turn, which means 18 residues in five turns. The diffraction pattern for a discrete helix of simple objects (say, points) with these dimensions has layer lines at all spacings  $Z = (18m + 5\alpha)/5P$ , where *m* and  $\alpha$  are integers, and the diffraction pattern will repeat every 18 layer lines. But of course, protein alpha helix does not contain 3.6 simple points per turn, but instead 3.6 complex groups of atoms per turn. Combined with the rapid drop-off of diffraction intensity at higher diffraction angles, this makes for diffraction patterns that are too complex for detailed analysis.

Now let's look briefly at just enough of the mathematics of fiber diffraction to explain the origin of the X patterns. Whereas each *reflection* in the diffraction pattern of a crystal is described by a Fourier sum of sine and cosine waves, each *layer line* in the diffraction pattern of a noncrystalline fiber is described by one or more *Bessel functions*, graphs that look like sine or cosine waves that damp out as they travel away from the origin (Fig. 9.3). Bessel functions appear when you apply the Fourier transform to helical objects. A Bessel function is of the form



**Figure 9.3** (*a*) Bessel functions J(x) of order  $\alpha = 0$  (red), 1 (green), and 2 (blue), showing positive values of x only. Note that, as the order increases, the distance to the first peak increases. (*b*) Enlargement of a few layer lines from Fig. 9.2*a*, showing the correspondence between diffraction intensities and the squares of Bessel functions having the same order as the layer-line number.

216

### Section 9.2 Fiber diffraction

The variable  $\alpha$  is called the order of the function, and the values of *n* are integers. To plot the Bessel function of order zero, you plug in the values  $\alpha = 0$  and n = 0 and then plot *J* as a function of *x* over some range -x to +x. Next you plug in  $\alpha = 0$  and n = 1, plot again, and add the resulting curve to the one for which n = 0, just as curves were added together to give the Fourier sum in Fig. 2.16, p. 25. Continuing in this way, you find that eventually, for large values of *n*, the new curves are very flat and do not change the previous sum. Bessel functions of orders zero, one, and two, for positive values of *x*, are shown in Fig. 9.3*a*. Notice that as the order increases, the position of the first peak of the function occurs farther from the origin.

Francis Crick showed in his doctoral dissertation that in the transform of a continuous helix, the intensity along a layer line is described by the square of the Bessel function whose order  $\alpha$  equals the number l of the layer line, as shown in Fig. 9.3*b*, which is an enlargement of three layer lines from the diffraction pattern of Fig. 9.2*a*. Thus, the intensity of the central line, layer line zero, varies according to  $[(J_0(x)]^2$ , which is the square of Eq. (9.1) with  $\alpha = 0$  (red). The intensity of the first line above (or below) center varies according to  $[(J_1(x)]^2$  (green), and so forth. This means that, for a helix, the first and largest peak of intensity lies farther out from the meridian on each successive layer line. The first peaks in a series of layer lines thus form the X pattern described earlier. The distance to the first peak in each layer line decreases as the helix radius increases, so thinner helices give wider X patterns.

For helices with a nonintegral number of residues per turn, the intensity functions, like the layer-line spacing, are also more complex, with two or more Bessel functions contributing to the intensities on each layer line. For the  $\alpha$  helix, with 18 residues in five turns, the Bessel functions that contribute to layer line l are those for which  $\alpha$  can be combined with some integral value of m (positive, negative or zero) to make  $l = 18\alpha + 5m$  an integer. For example, for layer line l = 0, one solution to this equation is  $\alpha = m = 0$ , so  $J_0(x)$  contributes to layer line 0. So also does  $J_5(x)$ , because  $\alpha = 5$ , m = -18 also gives l = 0. You can use  $l = 18\alpha + 5m$ also to show that  $J_2(x)$  (m = -7) and  $J_7(x)$  (m = -25) contribute to layer line l = 1, but that  $J_1(x)$  does not.

Probably the most famous fiber diffraction patterns are those of A-DNA and B-DNA obtained by Rosalind Franklin and shown in Fig. 9.4. Franklin's sample of A-DNA was microcrystalline, so its diffraction pattern (*a*) contains discrete reflections, many of them overlapping at the higher diffraction angles. Her B-DNA was noncrystalline, so the intensities in its diffraction pattern (*b*) vary smoothly across each layer line. Considering the B-helix, the narrow spacing between layer lines is inversely proportional to its 34-Å pitch. The distance from the center to the strong meridional intensity near the edge of the pattern is inversely proportional to the 3.4-Å rise per subunit (a nucleotide pair, we now know). Dividing pitch by rise gives 10 subunits per helical turn, as implied by the strong meridional intensity at the tenth layer line. Finally, the angle of the X pattern implies a helix radius of 20 Å.

Francis Crick recognized that Franklin's data implied a helical structure for B-DNA, a wonderful example of Pasteur's dictum: "Chance favors the

### Chapter 9 Other Diffraction Methods



a

b

Figure 9.4 ► Fiber diffraction patterns from A-DNA (left half of figure) and B-DNA (right). A-DNA was microcrystalline and thus gave discrete, but overlapping, Bragg reflections. B-DNA was noncrystalline and thus gave continuous variation in intensity along each layer line. Image kindly provided by Professor Kenneth Holmes.

prepared mind." Using helical parameters deduced from the pattern, knowing the chemical composition of DNA and the structures of DNA's molecular components, and finally, using the tantalizing observation that certain pairs of bases occur in equimolar amounts, James Watson was able to build a feasible model of B-DNA.

Why could Watson and Crick not simply back-transform the diffraction data, get an electron-density map, and fit a model to it? As in any diffraction experiment, we obtain intensities, but not phases. In addition, if we could somehow learn the phases, the back-transform would be the electron density averaged around the molecular axis parallel to the fiber axis, which would show only how electron density varies with distance from the center of the helix. Instead of map interpretation, structure determination by fiber diffraction usually entails inferring the dimensions of chains from the diffraction pattern and then building models using these dimensions plus prior knowledge: the known composition and the constraints on what is stereochemically allowable. As Watson's and Crick's success showed, for a helical substance, being able to deduce from diffraction the pitch, radius, and number of residues per helical turn puts some very strong constraints on a model.

The Fourier transform does, however, provide a powerful means of testing proposed models. With a feasible model in hand, researchers compute its transform and compare it to the pattern obtained by diffraction. If the fit is not perfect, they adjust the model, and again compare its transform with the X-ray pattern. They repeat this process until a model reproduces in detail the experimental diffraction pattern. Thus you can see why a very successful practitioner in this field, who may prefer to remain nameless, said, "Fiber diffraction is not what you'd do if you had a choice." Sometimes it is simply the only way to get structural information from diffraction.

For an example of structure analysis by fiber diffraction, see H. Wang and G. Stubbs, Structure determination of cucumber green mottle mosaic virus by X-ray fiber diffraction. Significance for the evolution of tobamoviruses, *J. Mol. Biol.* **239**, 371–384, 1994.

# 9.3 Diffraction by amorphous materials (scattering)

With fibers, the diffraction pattern, and hence any structural information contained therein, is averaged cylindrically about the molecular axis parallel to the fiber axis. This means that the transform of the diffraction pattern (computation of which would require phases) is an electron-density function showing only how electron density varies with distance from the molecular axis. For a helix of points aligned with the fiber axis, there would be a single peak of density at the radius of the helix. For a polyalanine helix (Fig. 9.5), there would be a density peak for the backbone, because atoms C, O, and CA are all roughly the same distance from the center of the helix (inner circle), and a second peak for the beta carbons, which are farther from the helix axis (outer circle). This averaging greatly reduces the structural information that can be inferred from the diffraction pattern, but it does imply that distance information is present, despite rotational averaging.

Imagine now a sample, such as a powder or solution, in which all the molecules are randomly oriented. Diffraction by such amorphous samples is



**Figure 9.5**  $\triangleright$  Polyalanine  $\alpha$  helix, viewed down the helix axis (stereo). An electrondensity map averaged around this axis would merely show two circular peaks of electron density, one at the distance of the backbone atoms (inner circle) and another at the distance of the  $\beta$  carbons (outer circle).

### Chapter 9 Other Diffraction Methods

usually called *scattering*. The diffraction pattern is averaged in all directions spherically—because the X-ray beam encounters all possible orientations of the molecules in the sample. But the diffraction pattern still contains information about how electron density varies with distance from the center of the molecules that make up the sample. Obviously, for complex molecules, this information would be singularly uninformative. But if the molecule under study contains only a few atoms, or only a few that dominate diffraction (like metal atoms in a protein), then it may be possible to extract useful distance information from the way that scattered X-ray intensity varies with the angle of scattering from the incident beam of radiation. As usual, when we try to extract information from intensity measurements, we work without knowledge of phases.

I have shown that, in simple systems, Patterson functions can give us valuable clues about distances, even when we know nothing about phases (see Sec. 6.3.3, p. 124). Diffraction from the randomly oriented molecules in a solution or powder would give a spherically averaged diffraction pattern, from which we can compute a spherically averaged Patterson map. Is this map interpretable? As in the case of heavy-atom derivatives, we can interpret a Patterson map if there are just a few atoms or a few very strong diffractors. Consider a linear molecule containing only three atoms (Fig. 9.6). We can see what to expect in a Patterson function computed from diffraction data on this molecule by constructing a Patterson function for the known structure. First, construct a Patterson function for the structure shown in Fig. 9.6*a*, using the procedure described in Fig. 6.12, p. 126. The result is (*b*). Spherical averaging of (*b*) gives a set of spherical shells of intensity, with



**Figure 9.6** Radial Patterson function. (*a*) Linear triatomic molecule. (*b*) Patterson function constructed from (*a*). For construction procedure, see Sec. 6.3.3, p. 124 and Fig. 6.12, p. 126. (*c*) Cross section of Patterson function (*b*) averaged by rotation to produce spherical shells. (*d*) Typical presentation of radial Patterson function as calculated from scattering measurements on an amorphous sample.

cross-section shown in (c). This cross section contains information about distances between the atoms in (a). The radius of the first circle is the bond length r, and the radius of the second circle is the length of the molecule (2r). A plot of the magnitude of the Patterson function as a function of distance from the origin (d), called a radial Patterson function, contains peaks that correspond to vectors between atoms. Furthermore, the intensity of each peak will depend on how many vectors of that length are present. In the molecule of Fig. 9.6a, there are four vectors of length rand two vectors of length 2r, so the Patterson peak at r is stronger than the one at 2r. You can see that, for a small molecule, the radial Patterson function computed from scattering intensities may contain enough information to determine distances between atoms. For larger numbers of atoms, the radial Patterson function would contain peaks corresponding to all the interatomic distances.

The radial Patterson function of an amorphous (powder or dilute solution) sample of a protein contains an enormous number of peaks. But imagine a protein containing a cluster of one or two metal ions surrounded by sulfur atoms. These atoms may dominate the powder diffraction data, and the strongest peaks in the radial Patterson function may reveal the distances among the metal ions and sulfur atoms. Remember that we obtain distance information but no geometry because diffraction is spherically averaged and all directional information is lost. Sometimes powder or solution diffraction can be used to extract distance information relating to a cluster of the heavier atoms in a protein. These distances put constraints on models of the cluster. Researchers can compare spherically averaged back-transforms of plausible models with the experimental diffraction data to guide improvements in the model, as in fiber diffraction. In some cases, model building and comparison of model back-transforms with data allows identification of ligand atoms and estimation of bond distances.

At very small diffraction angles, additional information can be obtained about the size and shape of a molecule. A detailed treatment of this method, called *lowangle scattering*, is beyond the scope of this book, but it can be shown that, for very small angles, the variation in scattering intensity is related to the *radius of gyration*,  $R_G$ , of the molecule under study. The radius of gyration is defined as the root-meansquare average of the distance of all scattering elements from the center of mass of the molecule. For two proteins having the same molecular mass, the one with the larger radius of gyration is the more extended or less spherical one. Combined knowledge of the molecular mass and the radius of gyration of a molecule allows an estimate of its shape. The precise relationship between scattering intensity and radius of gyration is

$$\langle I(\theta) \rangle = n_{\rm e}^2 (1 - 16\pi^2 R_{\rm G}^2 \sin^2 \theta / 3\lambda^2).$$
 (9.2)

 $\langle I(\theta) \rangle$  is the intensity of scattering at angle  $\theta$ ,  $n_e$  is the number of electrons in the molecule, and  $\lambda$  is the X-ray wavelength. Equation (9.2) implies that a graph of radiation intensity versus sin<sup>2</sup>  $\theta$  has a slope that is directly proportional to the square of  $R_G$ . Note also that such a graph can be extrapolated to  $\theta = 0$ , where the second term in parentheses disappears, and  $\langle I(\theta) \rangle$  is equal to the square of  $n_e$ , or roughly,

to the square of the molecular mass. So for molecules about which very little is known, measurement of scattering intensity at low angles provides estimates of both molecular mass and shape.

This kind of information about mass, shape, and distance can be obtained on amorphous samples not only from X-ray scattering but also from scattering by other forms of radiation, including light, and as I will discuss in Sec. 9.4, p. 222, neutrons. The choice of radiation depends on the size of the objects under study.

The variable-wavelength X-rays available at synchrotron sources give researchers an additional, powerful way to obtain precise distance information from amorphous samples. At wavelengths near the absorption edge of a metal atom (Sec. 4.3.2, p. 73), there is rapid oscillation of X-ray absorption as a function of wavelength. This oscillation results from interference between diffraction from the absorbing atom and that of its neighbors. The Fourier transform of this oscillation, in comparison with transforms calculated from plausible models, can reveal information on the number, types, and distances of the neighboring atoms. Distance information in favorable cases can be much more precise than atomic distances determined by X-ray crystallography. Thus this information can anticipate or add useful detail to crystallographic models. The measurement of absorbance as a function of wavelength is, of course, a form of spectroscopy. In the instance described here, it is called *X-ray absorption spectroscopy*, or *XAS*. Fourier analysis of near-edge X-ray absorption is called *extended X-ray absorption fine structure*, or *EXAFS*.

For an example of structure analysis by X-ray scattering, see D. I. Svergun, S. Richard, M. H. Koch, Z. Sayers, S. Kuprin, and G. Zaccai, Protein hydration in solution: Experimental observation by X-ray and neutron scattering, *Proc. Natl. Acad. Sci. USA* **95**, 2267–2272, 1998.

### 9.4 Neutron diffraction

The description of diffraction or scattering as a Fourier transform applies to all forms of energy that have wave character, including not only electromagnetic energy like X-rays and light but also subatomic particles, including neutrons and electrons, which have wavelengths as a result of their motion. The *de Broglie equation* (9.3) gives the wavelength  $\lambda$  of a particle of mass *m* moving at velocity *v*:

$$\lambda = h/mv, \tag{9.3}$$

where h is Planck's constant. We can use the de Broglie wavelength to describe diffraction of particles by matter. In this section, I will describe single-crystal neutron diffraction and neutron scattering by macromolecules, emphasizing the type of information obtainable. In the next section, I will apply these ideas to electron crystallography.

### Section 9.4 Neutron diffraction

Recall that X-rays are diffracted by the electrons that surround atoms, and that images obtained from X-ray diffraction show the surface of the electron clouds that surround molecules. Recall also that the X-ray diffracting power of elements in a sample increases with increasing atomic number. *Neutrons are diffracted by nuclei, not by electrons.* Thus a density map computed from neutron diffraction data is not an electron-density map, but instead a map of nuclear mass distribution, a "nucleon-density map" of the molecule (nucleons are the protons and neutrons in atomic nuclei).

The neutron crystallography experiment is much like X-ray crystallography (see Figs. 2.6, p. 13, and 2.12, p. 19). A crystal is held in a collimated beam of neutrons. Diffracted beams of neutrons are detected in a diffraction pattern that is a reciprocal-lattice sampling of intensities, but as usual not phases, of the Fourier transform of the average object in the crystal. Structure determination is, in principle, similar to that in X-ray crystallography, involving estimating the phases, back-transforming a set of structure factors composed of experimental intensities and estimated phases, improving the phases, and refining the structure.

There are two common ways to produce a beam of neutrons. One is steady-state nuclear fission in a reactor, which produces a continuous output of neutrons, some of which sustain fission, while the excess are recovered as a usable neutron beam. The second type is a pulsed source in which a cluster of protons or other charged particles from a linear accelerator are injected into a synchrotron, condensed into a tighter cluster or pulse, and allowed to strike a target of metal, such as tungsten. The high-energy particles drive neutrons from the target nuclei in a process called *spallation*. Neutrons from both fission and spallation carry too much energy for use in diffraction, so they are slowed or cooled ("thermalized") by passage through heavy water (D<sub>2</sub>O) at 300°K, producing neutrons with De Broglie wavelengths ranging from 1 to 2 Å. This wavelength is in the same range as X-rays used in crystallography.

Thermal neutrons then enter a collimator followed by a monochromator, which selects a narrow range of wavelengths to emerge and strike the sample. Monochromators are single crystals of graphite, zinc, or copper. They act like diffraction gratings to direct neutrons of different energies (and hence, wavelengths) in different directions, as a prism does with light. The collimated, monochromatic neutron beam is then delivered to the sample mounted on a goniometer, and diffraction is detected by an area detector (Sec. 4.3.3, p. 77). One common type is an image-plate that employs gadolinium oxide, which absorbs a neutron and emits a gamma ray, which in turn exposes the image plate.

The great advantage of neutron diffraction is that small nuclei like hydrogen are readily observed. By comparison with carbon and larger elements, hydrogen is a very weak X-ray diffractor and is typically not observable in electron-density maps of proteins. But hydrogen and its isotope deuterium (<sup>2</sup>H or D) diffract neutrons very efficiently in comparison with larger elements.

The concept of *scattering length* is used to compare diffracting power of elements. The scattering length b (do not confuse it with temperature factor B or the atomic scattering factor f) is the amplitude of scattering at an angle of zero

degrees to the incident beam and is the absolute measure of scattering power of an atom. Measured diffraction intensities are proportional to  $b^2$ , which is why X-ray diffraction from a single heavy atom can be easily detected above the diffraction of all the small atoms in a massive macromolecule. Table 9.1 gives neutron and X-ray scattering lengths for various elements and allows us to compare the scattering power of elements in neutron and X-ray diffraction.

Note that, even though H and D are weak X-ray diffractors compared to other elements in biomolecules (they have so few electrons around them), they are comparable to other elements when it comes to neutron diffraction. So hydrogen will be a prominent feature in density maps from neutron diffraction. Note also that the sign of b is negative for H. This means that H diffracts with a phase that is opposite to that of other elements. As mentioned before, measured diffraction intensities are related to  $b^2$ , so the negative sign of b has no observable effect on diffraction patterns. But in density maps, H gives negative density, which makes it stand out. In addition, the large magnitude of the difference between scattering lengths for H and D allows some powerful ways to use D as a label in diffraction or scattering experiments.

First, let us consider experiments analogous to single-crystal X-ray crystallography. In most electron-density maps of macromolecules, we cannot observe hydrogens. Thus we cannot distinguish the amide nitrogen and oxygen in glutamine and asparagine side chains [although online validation tools (p. 189) like MolProbity can assign them on the theoretical basis of hydrogen-bonding possibilities]. Nor can we determine the locations of hydrogens on histidine side chains, whose  $pK_a$  values allow both protonated and unprotonated forms at physiological pH. And we cannot detect critical hydrogens involved in possible hydrogen bonding with ligands like cofactors and transition-state analogs. Some hydrogens, including those on hydroxyl OH and amide N, can exchange with hydrogens of the solvent if they are exposed to solvent, and if they are not involved in tight

| Element | X-rays $b \times 10^{13}$ (cm) | Neutrons $b \times 10^{13}$ (cm) |
|---------|--------------------------------|----------------------------------|
| Н       | 2.8*                           | -3.74                            |
| D       | 2.8                            | 6.67                             |
| С       | 16.9                           | 6.65                             |
| Ν       | 19.7                           | 9.40                             |
| 0       | 22.5                           | 5.80                             |
| Р       | 42.3                           | 5.10                             |
| S       | 45                             | 2.85                             |
| Mn      | 70                             | -3.60                            |
| Fe      | 73                             | 9.51                             |
| Pt      | 220                            | 9.5                              |

TABLE 9.1 
X-Ray and Neutron Scattering Lengths of Various Elements

Data from C. R. Cantor and P. R. Schimmel (1980). *Bio-physical Chemistry, Part II: Techniques for the study of Biological Structure and Function.* W. H. Freeman and Company, San Francisco, p. 830. \*Value corrected from original.

### Section 9.4 Neutron diffraction

hydrogen bonds. Distinguishing H from D would mean being able to determine which hydrogens are exchangeable. Neutron diffraction, it would appear, gives us a way to answer these questions.

We can collect diffraction intensities from macromolecular crystals, but can we phase them and thus obtain maps that include clear images of hydrogen atoms? How about heavy atoms for phasing? According to Table 9.1, *there is no such thing as a heavy atom.* In other words, no nucleus diffracts so strongly that we can detect it above all others in a Patterson map. That rules out MIR, SIR with anomalous dispersion, and MAD as possible phasing methods. But if we are trying to find hydrogens in a structure known, or even partially known, from X-ray work, we have a source of starting phases in hand.

A crystal has the same reciprocal lattice for all types of diffraction, because the construction of the reciprocal lattice (Sec. 4.2.4, p. 57) does not depend in any way upon the type of radiation involved. So if we know the positions of all or most of the nonhydrogen atoms from X-ray structure determination, we can compute their contributions to the neutron-diffraction phases. These contributions depend only on atomic positions in the unit cell, and like reciprocal lattice positions, they are independent of the type and wavelength of radiation. We start phasing by assigning the final phases computed from the X-ray model of nonhydrogen atoms to the reflections obtained by neutron diffraction. From a Fourier sum combining neutron-diffraction intensities and X-ray phases, we compute the first map. Then we proceed as usual, alternating cycles of examining the map, building a model (in this case, adding hydrogens whose images we see in the map), back-transforming the model to get better phases, and so forth. In essence, this is isomorphous molecular replacement (Sec. 6.5, p. 136), using a hydrogen-free model that is identical to the fully hydrogenated model we seek.

Neutron diffraction has been used to detect critical hydrogen bonds in proteinligand complexes. For example, the presence of a hydrogen bond between dioxygen and the distal histidine of myoglobin is known because of neutron diffraction studies. In highly refined neutron-diffraction models, hydrogen positions are determined as precisely as positions of other atoms. In the best cases, not only can hydrogens in hydrogen bonds be seen, but it is even possible to tell whether hydrogen-containing groups like methyl or hydroxyl can rotate. If methyl groups are conformationally locked, density maps will show three distinct peaks of density for the three hydrogens of a methyl group, whereas if dynamic rotation is possible, or if alternative static conformations occur on different molecules, then we observe a smooth donut of density around the methyl carbon. For hydroxyls, it is possible to determine the conformation angle H-O-C-H in hydrogen-bonding and nonhydrogen-bonding situations. In the former, neutron diffraction sometimes reveals unexpected eclipsed conformations of the hydroxyl. Finally, the structure of networks of water molecules on the surface of macromolecules have been revealed by neutron diffraction. Being able to image the hydrogens means learning the orientation of each water molecule and the exact pattern of hydrogen bonding. In electron-density maps from X-ray work, we usually learn only the location of the oxygen atom of water molecules.
#### Chapter 9 Other Diffraction Methods

Because of the difference in signs of *b* for H and D, it is easy to distinguish them in density maps. This makes it possible to detect exchangeable hydrogens in proteins by comparing the density maps from crystals in H<sub>2</sub>O and in D<sub>2</sub>O. Amide hydrogens exchange with solvent hydrogens only if the amide group is exposed to solvent and is not involved in tight hydrogen bonding to other atoms. Because X-ray data are taken over a long period of time compared to rapid proton exchange, proton exchange rates can usually only be assigned to three categories: no exchange, slow exchange (10–60% during the time of data collection), and fast exchange (more than 60%).

Turning to scattering by amorphous samples, and to studies at lower resolution than crystallography, the negative scattering length for H makes possible interesting and useful scattering experiments on macromolecular complexes in solution. For example, neutron scattering in solution can be used to measure distances between the various protein components of very large macromolecular complexes like ribosomes and viruses. Understanding these methods requires bringing up a point that I have been able to avoid until now. Specifically, all forms of scattering depend on contrast between the scatterer and its surroundings. Even in singlecrystal X-ray crystallography, we get our first leg up on phases by finding the molecular boundary, in essence distinguishing protein from solvent. This is possible because of the contrast in density between ordered protein and disordered water. If the protein and water had exactly the same scattering power, we could not find this crucial boundary.

The importance of contrast to scattering is analogous to the importance of refractive index to refraction, such as the curving of light beams when they enter water. Light travels in a straight line through a pure liquid, but changes direction abruptly when it crosses a boundary into a new medium having a different refractive index. If two liquids have identical refractive indices, the boundary between them does not bend a light ray. Analogously, if X-rays or neutrons pass from one medium (say, solvent) into another (say, protein), scattering occurs only if the average scattering lengths of the media differ.

Because H and D have scattering lengths of different sign, it is possible to make mixtures of  $H_2O$  and  $D_2O$  with average scattering lengths over a wide range. In addition, by preparing proteins containing varying amounts of D substituted for H (by growing protein-producing cells in H<sub>2</sub>O/D<sub>2</sub>O mixtures), researchers can prepare proteins of variable average scattering length. A protein that scatters identically with the solvent is invisible to scattering. Imagine then reconstituting a bacterial ribosome, which is a complex of 3 RNA molecules and about 50 proteins, from partially D-labeled proteins and RNAs whose scattering power matches the H<sub>2</sub>O/D<sub>2</sub>O mixture in which they are dissolved. This mixture will not scatter neutrons. But if two of the components are unlabeled, only those two will scatter. At low resolution, it is as if the two proteins constitute one molecule made of two large atoms. The variation in intensity of neutron scattering with scattering angle will reveal the distance between the two proteins, just as radial Patterson functions reveal interatomic distances (see Fig. 9.6, p. 220). Repeating the experiment with different unlabeled pairs of proteins gives enough interprotein distances to direct the building of a three-dimensional map of protein locations.

Finally, low-angle neutron scattering can provide information about the shape of these molecules within the ribosome, which may not be the same as their shape when free in solution.

Neutron diffraction experiments are generally more difficult than those involving X-rays. There are fewer neutron sources worldwide. Available beam intensities, along with the low neutron-scattering lengths of all elements, translate into long exposure times. But for most macromolecules, neutron diffraction is the only source of detailed information about hydrogen locations and the exact orientation of hydrogen-carrying atoms.

For an example of neutron diffraction applied to a crystallographic problem, see D. Pignol, J. Hermoso, B. Kerfelec, I. Crenon, C. Chapus, and J. C. Fontecilla-Camps, The lipase/colipase complex is activated by a micelle: Neutron crystallographic evidence, *Chem. Phys. Lipids* **93**, 123–129, 1998.

# 9.5 Electron diffraction and cryo-electron microscopy

Electrons, like X-rays and neutrons, are scattered strongly by matter and thus are potentially useful in structure determination. Electron microscopy (EM) is the most widely known means of using electrons as structural probes. *Scanning* electron microscopes give an image of the sample surface, which is usually coated with a thin layer of metal. Sample preparation techniques for scanning EM are not compatible with obtaining images of molecules at atomic resolution. *Transmission* electron microscopes produce a projection of a very thin sample or section. In the most familiar electron micrographs of cells and organelles, the sample is stained with metals to outline the surfaces of membranes and large multimolecular assemblies like ribosomes. Unfortunately, staining results in distortion of the sample that is unacceptable at high resolution. The most successful methods of applying electron microscopes used to study unstained samples by either *electron diffraction* or *cryo-electron microscopy* (cryo-EM).

The electrons produced by transmission electron microscopes, whose design is analogous to light microscopes, have de Broglie wavelengths of less than 0.1 Å, so they are potentially quite precise probes of molecular structure. Unlike X-rays and neutrons, electrons can be focused (by electric fields rather than glass lenses) to produce an image, although direct images of objects in transmission EM do not approach molecular resolution. However, electron microscopes can be used to collect electron-diffraction patterns from two-dimensional arrays of molecules, such as closely packed arrays of membrane proteins in a lipid layer. Analysis of diffraction by such two-dimensional "crystals" is called *electron crystallography*.

Among the main difficulties with electron crystallography are (1) sample damage from the electron beam (a 0.1-Å wave carries a lot of energy), (2) low contrast between the solvent and the object under study, and (3) weak diffraction from

#### Chapter 9 Other Diffraction Methods

the necessarily very thin arrays that can be studied by this method. Despite these obstacles, cryoscopic methods (Sec. 3.5, p. 46) and image processing techniques have made electron crystallography a powerful probe of macromolecular structure, especially for membrane proteins, many of which resist crystallization.

Transmission electron microscopy is analogous to light microscopy, with visible light replaced by a beam of electrons produced by a heated metal filament, and glass lenses replaced by electromagnetic coils to focus the beam. An image of the sample is projected onto a fluorescent screen or, for a permanent record, onto film or a CCD detector (Sec. 4.3.3, p. 77). Alternatively, an image of the sample's *diffraction pattern* can be projected onto the detectors.

To see how we can observe either an image or a diffraction pattern, look again at Fig. 2.1, p. 8, which illustrates the action of a simple lens, such as the objective (lower) lens of a microscope. Recall that the lens produces an image at I of an object O placed outside the front focal point F of the lens. In a light microscope, an eyepiece (upper) lens is positioned so as to magnify the image at I for viewing. If we move the eyepiece to get an image of what lies at the back focal plane F', we see instead the diffraction pattern of the sample. Analogously, in EM, we can adjust the focusing power of the lenses to project an image of the back focal plane onto the detector, and thus see the sample's diffraction pattern. If the image is nonperiodic (say, a section of a cell) the diffraction pattern is continuous, as in Fig. 6.1, p. 110. If the sample is a periodic two-dimensional array, the diffraction pattern is sampled at reciprocal-lattice points, just as in X-ray diffraction by crystals (Fig. 2.19e, p. 29). The Fourier transform of this pattern is an image of the average object in the periodic array. As with all diffraction, producing such an image requires knowing both diffraction intensities and phases.

The possibility of viewing both the image and the diffraction pattern is unique to electron crystallography and, in favorable cases, can allow phases to be determined directly. For an ordered array (a 2-D crystal) of proteins in a lipid membrane, the direct image is, even at the highest magnification, a featureless gray field. But phase estimates can be obtained from this singularly uninteresting image in the following manner. The image is digitized to pixels at high resolution, producing a two-dimensional table of *image*-intensity values (not diffraction intensities). This table is then Fourier transformed. Computing the FT of a table of values is the same process as used to produce the images in Fig. 2.19, p. 29, Fig. 6.1, p. 110, and Fig. 6.17, p. 138, in which the "samples" are in fact square arrays of pixels with different numerical values. If the gray EM image is actually a periodic array, the result of the FT on the table of pixel values is the diffraction pattern of the average object in the array, sampled at reciprocal-lattice points. But because this transform is computed from "observed" objects, it includes both intensities and phases. The intensities are not as accurate as those that we can measure directly in the diffraction plane of the EM, but the phases are often accurate enough to serve as first estimates in a refinement process.

Single-crystal X-ray crystallography requires measuring diffraction intensities at many closely spaced crystal orientations, so as to measure most of the unique reflections in the reciprocal lattice. The Fourier transform (with correct phases) of

#### Section 9.5 Electron diffraction and cryo-electron microscopy

reflections from a single orientation gives only a projection of the unit cell in a plane perpendicular to the beam. The transform of the full data set gives a threedimensional image. Electron microscopes allow samples to be tilted, which for diffraction is analogous to rotating the crystal. Tilting gives diffraction patterns whose transforms are projections of the sample contents from a different angle. (In the same way, the fall of leaves during a slow, steady wind produces a pattern of fallen leaves that is like a projection of the tree from an angle off the vertical.) With diffraction patterns taken from a wide enough range of angles, you can obtain a three-dimensional image of the sample contents. The EM sample can be tilted about 75° at most, which may or may not allow a sufficiently large portion of the reciprocal lattice to be sampled. It is not unusual for EM data sets to be missing data in parts of reciprocal space that cannot be brought into contact with the sphere of reflection (see Fig. 4.12, p. 62, and Fig. 4.13, p. 63) by tilting. Most common is for a cone of reflections to be missing, with the result that maps computed from these data do not have uniform resolution in all directions.

Armed with sets of data measured at different sample tilt angles, each set including (1) intensities measured at the diffraction plane and (2) phases from the Fourier transform of direct images, the crystallographer can compute a map of the unit cell. Recall that X-rays are scattered by electrons around atoms, producing a map of electron density, and neutrons are scattered by nuclei, producing a map of mass or nucleon density. Electrons, in turn, are scattered by electrostatic interactions, producing maps of electrostatic potential, sometimes called *electron-potential maps*, or just potential maps. Maps from electron diffraction, like all other types, are interpreted by building molecular models to fit them and refined by using partial models as phasing models.

One of the special features of electron crystallography is the possibility of detecting charge on specific atoms or functional groups. The probing electrons interact very strongly with negative charge, with the result that negatively charged atoms have negative scattering factors at low resolution. (Recall that the negative scattering length for H in neutron diffraction can make H particularly easy to detect in neutron-density maps.) At high resolution, the negative sign of the scattering factor weakens the signal of negatively charged groups, but not usually enough to be obvious. However, a comparison between maps computed with and without the low-angle data can reveal charged groups. If the two maps are practically identical around a possibly charged group like a carboxyl, then the group is probably neutral. If, however, the map computed without low-angle data shows a stronger density for the functional group than does the map computed from all data, then the functional group probably carries a negative charge. For example, functional groups in proton pumps like rhodopsin are apparently involved in transferring protons across a membrane by way of a channel through the protein. Determining the ionization state of functional groups in the channel is essential to proposing pumping mechanisms.

Cryo-EM combines cryoscopic electron microscopy with another group of powerful structural methods collectively called *image enhancement*. These methods do not involve diffraction directly, but they take advantage of the averaging power

#### Chapter 9 Other Diffraction Methods

of Fourier transforms to produce direct images at low resolution (10–25 Å). Image enhancement can greatly improve the resolution of supramolecular complexes like ribosomes, viruses, and multienzyme complexes that can be seen individually, but at low resolution, as direct EM images. Samples for cryo-EM are unstained. They are flash frozen in small droplets of cryoprotected aqueous buffers, so that the water freezes as a glass, rather than crystalline ice, just as in cryocrystallography.

Imagine a direct EM image of, say, virus particles or ribosomes, all strewn across the viewing field. Particles lay before you in all orientations, and the images show very little detail. In image enhancement, we digitize the individual images and sort them into images that appear to share the same orientation. Then we can compute the Fourier transforms of the images. These transforms should be similar in appearance if indeed the images in a set share the same orientation. Next, we align the transforms, add them together, and back-transform. The result is an averaged image, in which details common to the component images are enhanced, and random differences, a form of "noise," are reduced or eliminated. This process is then repeated for sets of images at different orientations, the transforms are all combined into a three-dimensional set, and the back-transform gives a threedimensional image.

Recall that the Fourier transform of a crystal diffraction pattern gives an image of the average molecule in the crystal. Because any averaging process tends to eliminate random variations and enhance features common to all components, the image is much more highly resolved than if it were derived from a single molecule. In like manner, the final back-transform of Fouriers from many EM images is an image of the average particle in the direct EM field, and the resolution can be greatly enhanced in comparison to any single particle image. In favorable cases, enhanced EM images can be used as phasing models in crystallographic structure determination by molecular replacement, where the particles or even components of the particles can be crystallized.

Another powerful combination of cryo-EM and crystallography is making possible structure determination of supramolecular complexes that are not subject to crystallization, such as bacteriophages and virus-receptor complexes. Components of such complexes are purified and crystallized, and their individual structures are determined at atomic resolution by single-crystal X-ray crystallography. Structures of the intact complexes are determined at resolutions of 10–20 Å by cryo-EM. Then the crystallographic component structures are fitted into the low-resolution complex structure, and their positions refined by real-space fitting of component electron density into the cryo-EM density. If a majority of components can be fitted this way, then subtraction of their density from the cryo-EM density reveals the shape and orientation of remaining components. The resulting low-resolution models of the remaining components might be useful as molecular replacement models in their crystallographic structure determination. Such studies are yielding detailed structures of supramolecular complexes and insights into their action.

The European Bioinformatics Institute (EBI) is host for the Electron Microscopy Database (EMD), a repository of molecular structures determined by electron microscopy. A typical EMD entry contains sample identity, details of the EM experiment and data processing, density maps, and images of the model. To find the EMD, go to the CMCC home page.

For an example of electron diffraction in structure determination of a membrane protein, see Y. Kimura, D. G. Vassylyev, A. Miyazawa, A. Kidera, M. Matsushima, K. Mitsuoka, K. Murata, T. Hirai, and Y. Fujiyoshi, Surface of bacteriorhodopsin revealed by high-resolution electron crystallography, *Nature* **389**, 206–211, 1997.

For an example of crystallography combined with cryo-EM and applied to a supramolecular complex, see M. Rossmann, V. Mesyanzhinov, F. Arisaka, and P. Leiman, The bacteriophage T4 DNA injection machine, *Current Opinion in Structural Biology*, **14**, 171–180, 2004 and references cited therein. To see EMD entries related to this publication, search the EMD for "Rossmann."

# 9.6

# Laue diffraction and time-resolved crystallography

Now I return to X-ray diffraction to describe probably the oldest type of diffraction experiment, but one whose stock has soared with the advent of synchrotron radiation and powerful computer techniques for the analysis of complex diffraction data. The method, Laue diffraction, is already realizing its promise as a means to determine the structures of short-lived reaction intermediates. This method is sometimes called *time-resolved crystallography*, implying an attempt to take snapshots of a chemical reaction or physical change in progress.

Laue crystallography entails irradiating a crystal with a powerful *polychromatic* beam of X rays, whose wavelengths range over a two- to threefold range, for example, 0.5-1.5 Å. The resulting diffraction patterns are more complex than those obtained from monochromatic X-rays, but they sample a much larger portion of the reciprocal lattice from a single crystal orientation. To see why this is so, look again at Fig. 4.12, p. 62, which demonstrates in reciprocal space the conditions that satisfy Bragg's law, and shows the resulting directions of diffracted rays when the incident radiation is monochromatic. Recall that Bragg's law is satisfied when rotation of the crystal (and with it, the reciprocal lattice) brings point *P* (in Fig. 4.12*a*) or *P'* [in (*b*]] onto the surface of the sphere of reflection. The results are diffracted rays *R* and *R'*. Also recall that the sphere of reflection, whose radius is the reciprocal of the X-ray wavelength  $\lambda$ , passes through the origin of the reciprocal lattice, and its diameter is coincident with the X-ray beam.

Figure 9.7 extends the geometric construction of Fig. 4.12, p. 62 to show the result of diffraction when the X-ray beam provides a continuum of wavelengths. Instead of one sphere of reflection, there are an infinite number, covering the gray region of the figure, with radii ranging from  $1/\lambda_{max}$  to  $1/\lambda_{min}$ , where  $\lambda_{max}$  and  $\lambda_{min}$  are the maximum and minimum wavelengths of radiation in the beam. In the figure, spheres of reflection corresponding to  $\lambda_{max}$  and  $\lambda_{min}$  are shown, along with two others that lie in between. The four points lying on the incident beam *X* are the



**Figure 9.7**  $\triangleright$  Geometric construction for Laue diffraction in reciprocal space. Each diffracted ray  $R_1$  through  $R_4$  is shown in the same color as that of the sphere of reflection and the lattice point that produces it.

centers of the four spheres of reflection. Each of the four spheres is shown passing through one reciprocal-lattice point (some pass through others also), producing a diffracted ray ( $R_1$  through  $R_4$ ).

Note, however, that because there are an infinite number of spheres of reflection, *every* reciprocal-lattice point within the gray region lies on the surface of some sphere of reflection, and thus for every such point Bragg's law is satisfied, giving rise to a diffracted ray. Of course, any crystal has its diffraction limit, which depends on its quality. The heavy arc labeled  $d_{max}^*$  represents the resolution limit

#### Section 9.6 Laue diffraction and time-resolved crystallography

of the crystal in this illustration. Diffraction corresponding to reciprocal-lattice points outside this arc (for example,  $R_3$  and  $R_4$ ), even though they satisfy Bragg's law, will not be detected simply because this crystal does not diffract out to those high angles. Because the detector is at the geometric equivalent of an infinite distance from the crystal, we can picture all rays as emerging from a single point, such as the origin. To show the relative directions of all reflections in the Laue pattern from this diagram, we move all the diffracted rays to the origin, as shown for rays  $R_1$  through  $R_4$  in Fig. 9.7b.

Note also that because of the tapering shape of the gray area near the origin, the amount of available data drops off at low angles. One limitation of Laue diffraction is the scarcity of reflections at small angles. An added complication of Laue diffraction is that some rays pass through more than one point. When this occurs, the measured intensity of the ray is the sum of the intensities of both reflections. Finally, at higher resolution, many reflections overlap.

You can see from Fig. 9.8 that a Laue diffraction pattern is much more complex than a diffraction pattern from monochromatic X-rays. But modern software can index Laue patterns and thus allow accurate measurement of many diffraction intensities from a single brief pulse of X-rays through a still crystal. If the crystal



Figure 9.8 ► A typical Laue image. From I. J. Clifton, E. M. H. Duke, S. Wakatsuki, and Z. Ren, *Evaluation of Laue Diffraction Patterns, in Methods in Enzymology 277B*, C. W. Carter and R. M. Sweet, eds., Academic Press, New York, 1997, p. 453. Reprinted with permission.

#### Chapter 9 Other Diffraction Methods

has high symmetry and is oriented properly, a full data set can in theory be collected in a single brief X-ray exposure. In practice, this approach usually does not provide sufficiently accurate intensities because the data lack the redundancy necessary for high accuracy. Multiple exposures at multiple orientations are the rule.

The unique advantage of the Laue method is that data can be collected rapidly enough to give a freeze-frame picture of the crystal's contents. Typical X-ray data are averaged over the time of data collection, which can vary from seconds to days, and over the sometimes large number of crystals required to obtain a complete data set. Laue data has been collected with X-ray pulses shorter than 200 picoseconds. Such short time periods for data collection are comparable to halftimes for chemical reactions, especially those involving macromolecules, such as enzymatic catalysis. This raises the possibility of determining the structures of reaction intermediates.

Recall that a crystallographic structure is the average of the structures of all diffracting molecules, so structures of intermediates can be determined when all molecules in the crystal react in unison or exist in an intermediate state simultaneously. Thus the crystallographer must devise some way to trigger the reaction simultaneously throughout the crystal. Good candidate reactions include those triggered by light. Such processes can be initiated throughout a crystal by a laser pulse. Some reactions of this type are reversible, so the reaction can be run repeatedly with the same crystal.

A strategy for studying enzymatic catalysis entails introducing into enzyme crystals a "caged" substrate—a derivative of the substrate that is prevented from reacting by the presence of a light-labile protective group. If the caged substrate binds at the active site, then the stage is set to trigger the enzymatic reaction by a light pulse that frees the substrate from the protective group, allowing the enzyme to act. Running this type of reaction repeatedly may mean replacing the crystal, or maybe just introducing a fresh supply of caged substrate to diffuse into the crystal for the next run.

Other possibilities are reactions that can be triggered by sudden changes in temperature or pressure. Reactions that are slow in comparison to diffusion rates in the crystal may be triggered by simply adding substrate to the mother liquor surrounding the crystal and allowing the substrate to diffuse in. If there is a long-lived enzyme-intermediate state, then during some interval after substrate introduction, a large fraction of molecules in the crystal will exist in this state, and a pulse of radiation can reveal its structure.

In practice, multiple Laue images are usually necessary to give full coverage of reciprocal space with the required redundancy. Here is a hypothetical strategy for time-resolved crystallography of a reversible, light-triggered process. From knowledge of the kinetics of the reaction under study, determine convenient reaction conditions (for example, length and wavelength of triggering pulse and temperature) and times after initiation of reaction that would reveal intermediate states. These times must be long compared to the shortest X-ray pulses that will do the job of data collection. From knowledge of crystal symmetry, determine the number of orientations and exposures that will produce adequate data

#### Section 9.7 Summary

for structure determination. Then set the crystal to its first orientation, trigger the reaction, wait until the first data-collection time, pulse with X-rays, and collect Laue data. Allow the crystal to equilibrate, which means both letting the reaction come back to equilibrium and allowing the crystal to cool, since the triggering pulse and the X-ray pulse heat the crystal. Then again trigger the reaction, wait until the second data-collection time, pulse with X-rays, and collect data. Repeat until you have data from this crystal orientation at all the desired data-collection times. Move the crystal to its next orientation, and repeat the process. The result is full data sets collected at each of several time intervals during which the reaction was occurring.

For an example of Laue diffraction applied to time-resolved crystallography, see V. Srajer, T. Teng, T. Ursby, C. Pradervand, Z. Ren, S. Adachi, W. Schildkamp, D. Bourgeois, M. Wulff, and K. Moffat, Photolysis of the carbon monoxide complex of myoglobin: Nanosecond time-resolved crystallography, *Science* **274**, 1726–1729, 1996.

# 9.7 Summary

The same geometric and mathematical principles lie at the root of all types of diffraction experiments, whether the samples are powders, solutions, fibers, or crystals, and whether the experiments involve electromagnetic radiation (X-rays, visible light) or subatomic particles (electrons, neutrons). My aim in this chapter was to show the common ground shared by all of these probes of molecular structure. Note in particular how the methods complement each other and can be combined to produce more inclusive models of macromolecules. For example, phases from X-ray work can serve as starting phase estimates for neutron work, and the resulting accurate coordinates of hydrogen positions can then be added to the X-ray model. As another example, direct images obtained from EM work (sometimes after image enhancement) can sometimes be used as molecular-replacement models for X-ray crystallography, or to determine the organization of components in supramolecular complexes.

As a user of macromolecular models, you are faced with judging whether each model really supports the insights it appears to offer. The principles presented in Chapter 7, on how to judge the quality of models, apply to models obtained from all types of diffraction experiments. But today's structural databases also contain a growing number of models obtained by methods other than diffraction. In the next chapter, I will describe the origin of the major types of nondiffraction models and provide some guidance on how to use them wisely.

This Page Intentionally Left Blank

# ► Chapter 10

# Other Kinds of Macromolecular Models

# 10.1 Introduction

When you go looking for models of macromolecules that interest you, crystallographic models are not the only type you will find. As of mid-2005, about 15% of the models in the Protein Data Bank were derived from NMR spectroscopy of macromolecules in solution. Of over 4000 NMR models, fewer than 2% were proteins of more than 200 residues, although the number of larger models is sure to increase. Also proliferating rapidly are homology models, which are built by computer algorithms that work on the assumption that proteins of homologous sequence have similar three-dimensional structures. Massive databases of homology models are constantly growing, with the goal of providing homology models for all known protein sequences that are homologous to structures determined by crystallography or NMR. Although this effort is just beginning, the number of homology models available in one database, the SWISS-MODEL Repository (see the CMCC home page), contains more models than does the Protein Data Bank. Furthermore, it is quite easy for you to determine protein structures by homology modeling on your personal computer, as I will describe in Chapter 11. Finally, there are various means of producing theoretical models, for example, based on attempts to simulate folding of proteins. As with crystallographic models, other types of models vary widely in quality and reliability. The user of these models is faced with deciding whether the quality of the model allows confidence in the apparent implications of the structure.

In this chapter, I will provide brief descriptions of how protein structures are determined by NMR and by homology modeling. In addition, I will provide some guidance on judging the quality of noncrystallographic models, primarily by drawing analogies to criteria of model quality in crystallography. Recall that in all protein-structure determination, the goal is to determine the *conformation* of a molecule whose chemical *composition* (amino-acid content and sequence) is known.

# 10.2 NMR models

# 10.2.1 Introduction

Models of proteins in solution can be derived from NMR spectroscopy. In brief, the process entails collecting highly detailed spectra; assigning spectral peaks (resonances) to all residues by chemical shift and by decoupling experiments; deriving distance restraints from couplings, which (a) reveal local conformations and (b) pinpoint pairs of atoms that are distant from each other in sequence, but near each other because of the way the protein is folded; and computing a chemically, stereochemically, and energetically feasible model that complies with all distance restraints. That's it. There is no image, no "seeing" the model in a map. NMR structure determination amounts to building a model to fit the conformational and distance restraints learned from NMR, and of course, to fit prior knowledge about protein structure. The power of NMR spectroscopy has grown immensely with the development of pulse-Fourier transform instruments, which allow rapid data collection; more powerful magnets, which increase spectral resolution; and sophisticated, computer-controlled pulse sequences, which separate data into subsets (somewhat misleadingly called "dimensions") to reveal through-bond and through-space couplings.

Typical sample volumes for NMR spectroscopy are around 0.5 mL of protein in  $H_2O/D_2O$  (95/5 to 90/10). Typical protein concentrations are 0.5 to 2.0 mM. For a protein of 200 amino-acid residues, this means that each run consumes 5 to 20 mg of protein. The protein must be stable at room temperature for days to weeks, and must not be prone to aggregation at these relatively high concentrations. Depending on the thermal stability of the protein, a sample may be used only once or many times. In addition, proteins must be labeled with NMR-active isotopes such as <sup>13</sup>C and <sup>15</sup>N. Like crystallography, NMR spectroscopy reaps the benefits of modern molecular biology's powerful systems for expression and purification of desired proteins in quantity, labeled with appropriate isotopes from their growth media.

In this section, I provide a simplified physical picture of pulse NMR spectroscopy, including a simple conceptual model to help you understand multidimensional NMR. Then I briefly discuss the problems of assigning resonances and determining distance restraints for molecules as large and complex as proteins, and the methods for deriving a structure from this information. Finally, I discuss the contents of coordinate files from NMR structure determination and provide some hints on judging the quality of models.

## 10.2.2 Principles

I assume that you are conversant with basic principles of <sup>1</sup>H or proton NMR spectroscopy as applied to small molecules. In particular, I assume that you understand the concepts of chemical shift ( $\delta$ ) and spin-spin coupling, classical continuous-wave methods of obtaining NMR spectra, and decoupling experiments to determine pairs of coupled nuclei. If these ideas are unfamiliar to you, you may wish to review NMR spectroscopy in an introductory organic chemistry textbook before reading further.

## Chemical shift and coupling

Many atomic nuclei, notably <sup>1</sup>H, <sup>13</sup>C, <sup>15</sup>N, <sup>19</sup>F, and <sup>31</sup>P, have net nuclear spins as a result of the magnetic moments of their component protons and neutrons. These spins cause the nuclei to behave like tiny magnets, and as a result, to adopt preferred orientations in a magnetic field. A nucleus having a spin quantum number of 1/2 (for example, <sup>1</sup>H) can adopt one of two orientations in a magnetic field, aligned either with or against the field. Nuclei aligned with the field are slightly lower in energy, so at equilibrium, there are slightly more nuclei (about 1 in 10,000) in the lower energy state. The orientations of nuclear spins can be altered by pulses of electromagnetic radiation in the radio-frequency (RF) range.

Nuclei in different chemical environments absorb different frequencies of energy. This allows specific nuclei to be detected by their characteristic absorption energy. This energy can be expressed as an RF frequency (in *hertz*), but because the energy depends on the strength of the magnetic field, it is expressed as a frequency difference between that of the nucleus in question and a standard nucleus (like hydrogen in tetramethylsilane, a common standard for <sup>1</sup>H NMR) divided by the strength of the field. The result, called the *chemical shift*  $\delta$  and expressed in parts per million (ppm), is independent of field strength, but varies informatively with the type of nucleus and its immediate molecular environment. Figure 10.1 shows the <sup>1</sup>H-NMR spectrum of human thioredoxin, <sup>1</sup> a small protein of 105 amino-acid residues. (In humans, thioredoxin plays a role in activating certain transcriptional and translational regulators by a dithiol/disulfide redox mechanism. More familiar to students of biochemistry is the role of thioredoxin in green plants: mediating light activation of enzymes in the Calvin cycle of photosynthesis.) The spectrum is labeled with ranges of  $\delta$  for various types of hydrogen atoms in proteins.

In addition to exhibiting characteristic chemical shifts, nuclear spins interact magnetically, exchanging energy with each other by a process called *spin-spin coupling*. Coupling distributes or splits the absorption signals of nuclei about their characteristic absorption frequency, usually in a distinctive pattern that depends on the number of equivalent nuclei that are coupled, giving the familiar multiplets in NMR spectra of simple molecules. The spacing between signals produced by

<sup>&</sup>lt;sup>1</sup>J. D. Porman-Kay, G. M. Clore, P. T. Wingfield, and A. M. Gronenbom, High-resolution threedimensional structure of reduced recombinant human thioredoxin in solution, *Biochemistry* **30**, 2685, 1991 (PDB files 2trx and 3trx).



**Figure 10.1**  $\triangleright$  <sup>1</sup>H-NMR spectrum of thioredoxin, reduced form. Signal intensity (vertical axis) is plotted against frequency (horizontal axis). Labels show chemical-shift values typical of various hydrogen types in protein chains having random coil conformation. Some signals lie outside these ranges because of specific interactions not present in random coils. Atom labels are as found in PDB coordinate files (p. 176). Spectrum generously provided by Professor John M. Louis.

splitting is called the coupling constant J, which is expressed in hertz because its magnitude does not depend on field strength. Because nuclei must be within a few bonds of each other to couple, this effect can be used to determine which nuclei are neighbors in the molecule under study.

# Absorption and relaxation

You can view the nuclear spins in a sample as precessing about an axis, designated z, aligned with the magnetic field H, as shown in Fig. 10.2, which depicts the spins for a large number of equivalent nuclei in (a), and only the slight equilibrium excess of spins at the lower energy in (b). The slight excess of spins aligned with the field means that there is a net magnetization vector pointing in the positive direction along z at equilibrium. Spins precess at their characteristic RF absorption frequency, but the phases of precession are random. In other words, equivalent, independent spins precess at the same rate, but their positions are randomly distributed about the z-axis.

In Fig. 10.3, a highly simplified sketch of an NMR instrument shows the orientation of the x-, y-, and z-axes with respect to the magnets that apply the field H, the sample, and wire coils that transmit (green) and detect (red) radio-frequency



**Figure 10.2** (a) Nuclear spins precess around the z-axis, which is parallel to the applied field H, indicated by the arrow in the center of the figure. Each spin precesses at a rate that depends on its RF absorption frequency. (b) Excess spins in the lower energy state. At equilibrium, there are slightly more spins aligned with H than against H. (c) Immediately after a 90° pulse, equal numbers of spins are aligned with and against H, and initially all spins lie in the yz-plane, giving a net magnetic vector on the y-axis. This vector precesses in the xy plane, inducing a signal in the receiving coils. As time passes after the pulse, spins precess about z and spread out, due to their different precession frequencies, ultimately after a 180° pulse, or two successive 90° pulses, excess spins shown in (b) are inverted, and lie in the yz-plane as shown. The net magnetization vector in the xy-plane has a magnitude of zero, and thus induces no signal in the receiver coils.

signals. The receiver coils, which encircle the *x*-axis, detect RF radiation in the xy-plane only. A single nuclear spin precessing alone around the *z*-axis has a rotating component in the xy-plane, and transmits RF energy at its characteristic absorption frequency to the receiver coils, in theory revealing its chemical shift. But when many equivalent spins in a real sample are precessing about the *z*-axis with random phase, the xy components of their magnetic vectors cancel each other out, the net magnetic vector in the xy plane has a magnitude of zero, and no RF signals are detected. Detection of RF radiation from the sample requires that a net magnetization vector be moved into the xy-plane. This is accomplished by a wide-band (multifrequency) pulse of RF energy having just the right intensity to



**Figure 10.3**  $\triangleright$  Diagram of an NMR experiment. The sample lies between the poles of a powerful magnet, and is spun rapidly around its long (*x*) axis in order to compensate for any unevenness (inhomogeneity) of the magnetic field. Radio frequency receiver coils (red) form a helix around the *x*-axis, and transmitter coils (green) spiral around the *y*-axis.

tip the net magnetization vector into the *xy*-plane. This pulse is applied along *x*, and it equalizes the number of spins in the higher and lower energy states, while aligning them onto the *yz*-plane (Fig. 10.2*c*). A pulse of this intensity is call a 90° pulse. The result is that a net magnetization vector appears in the *xy*-plane and begins to rotate, generating a detectable RF signal. Because nuclei of different chemical shifts precess at different frequencies, their net magnetization vectors in the *xy*-plane rotate at different frequencies, and the resulting RF signal contains the characteristic absorption frequencies of all nuclei in the sample, from which chemical shifts can be derived.

As shown in Fig. 10.2*c*, the 90° pulse equalizes the number of nuclei aligned with and against the applied field along *z*, giving a net magnetization vector along *z* of zero. This means that the system of spins is no longer at equilibrium. It will return to equilibrium as a result of spins losing energy to their surroundings, a process called *spin-lattice relaxation*, in which the term *lattice* simply refers to the surroundings of the nuclei. This relaxation is a first-order process whose rate constant I will call  $R_L$  (L for lattice). The inverse of  $R_L$  is a time constant I will call  $T_L$ , the spin-lattice relaxation time constant.  $T_L$  is also called the *longitudinal* relaxation time constant, because relaxation occurs along the *z*-axis, parallel to the magnetic field. ( $T_L$  is traditionally called  $T_1$ , but I adopt a more descriptive symbol to reduce confusion with other symbols in the following discussion.)

After a 90° pulse, another relaxation process also occurs. Although all spins have the same phase just after the pulse (aligned along y), the phases of identical nuclei spread out due to exchanges of energy that result from coupling with other spins. The result is that the net magnetization vector in the xy-plane for each distinct set of nuclei diminishes in magnitude, as individual spin magnetizations move into orientations that cancel each other. Furthermore, the RF frequencies

#### Section 10.2 NMR models

of chemically distinct nuclei disappear from the overall signal at different rates, depending on their coupling constants, which reflect the strength of coupling, or in other words, how effectively the nuclei exchange energy. Phases of pairs of nuclei coupled to each other spread out at the same rate because their spin energies are simply being exchanged. The rate constant for this process, which is called *spin-spin relaxation*, is traditionally called  $R_2$ , but I will call it  $R_S$  (S for spin). Its inverse, the time constant  $T_S$ , is the spin-spin relaxation time constant. It is also called the *transverse* relaxation time constant because the relaxation process is perpendicular or transverse to the applied magnetic field *H*. Each set of chemically identical nuclei has a characteristic  $T_S$ . (Often, and confusingly, the rate constants *R* are loosely called rates, and time constants *T* are simply called relaxation times.)

So after a 90° RF pulse applied along x tips the net magnetization vector onto the y-axis, an RF signal appears at the detector because there is net magnetization rotating in the xy-plane. This signal is a composite of the frequencies of all the precessing nuclei, each precessing at its characteristic RF absorption frequency. In addition, this signal is strong at first but decays because of spin-spin relaxation. Signals for different nuclei diminish at different rates that depend on their individual spin-spin relaxation rates. This complex signal is called a *free-induction decay* or *FID*: *free* because it is free of influence from the applied RF field, which is turned off after the pulse; *induction* because the magnetic spins induce the signal in the receiver coil; and *decay* because the signal decays to an equilibrium value of zero. A typical hydrogen NMR FID signal decays in about 300 ms, as shown in Fig. 10.4.

## **One-dimensional NMR**

How do we extract the chemical shifts of all nuclei in the sample from the freeinduction decay signal? The answer is our old friend the Fourier transform. The FID is called a *time-domain* signal because it is a plot of the oscillating and decaying RF intensity *versus* time, as shown in Fig. 10.4 (the time axis is conventionally labeled  $t_2$ , for reasons you will see shortly). Fourier transforming the FID produces a *frequency-domain* spectrum, a plot of RF intensity versus the frequencies present in the FID signal, with the frequency axis labeled  $v_2$  for frequency or  $F_2$  for chemical shift, as shown in Fig. 10.1. So the Fourier transform decomposes the FID into its component frequencies, revealing the chemical shifts of the nuclei in the sample.

NMR spectroscopists collect this type of classical or one-dimensional (1-D) NMR spectra on modern FT-NMR instruments by applying a 90° pulse and collecting the FID signal that is induced in the receiving coil. To make a stronger signal, they collect FIDs repeatedly and add them together. Real signals appear at the same place in all the FIDs and add up to a large sum. On the other hand, random variation or "noise" appears in different places in different FIDs, so their signals cancel each other. The summed FIDs thus have a high signal-to-noise ratio, and FT produces a clean spectrum with well-resolved chemical shifts. From here on, I may not always mention that pulse/FID collection sequences are repeated to



**Figure 10.4** Free induction decay signal, which appears in the receiver coils after a 90° pulse. The FID is a time-domain spectrum, showing RF intensity as a function of time  $t_2$ . It is a composite of the RF absorption frequencies of all nuclei in the sample. The Fourier transform decomposes an FID into its component frequencies, giving a spectrum like that shown in Fig. 10.1. Figure generously provided by Professor John M. Louis.

improve the signal, but you can assume that all sequences I describe are carried out repeatedly, commonly 64 times, for this purpose.

# **Two-dimensional NMR**

Couplings between nuclei influence the rate at which their characteristic frequencies diminish in the FID. If we could measure, in addition to the frequencies themselves, their rates of disappearance, we could determine which pairs of nuclei are coupled, because the signals of coupled pairs fade from the FID at the same rate. This is the basis of *two-dimensional* (2-D) NMR, which detects not only the chemical shifts of nuclei but also their couplings. The 2-D NMR employs computer-controlled and -timed pulses that allow experimenters to monitor the progress of relaxation for different sets of nuclei.

Figure 10.5 illustrates, with a system of four nuclei, the principles of 2-D NMR in its simplest form, when we want to assign pairs of <sup>1</sup>H nuclei that are spin-spin coupled through a small number of bonds, such as hydrogens on adjacent carbon atoms. You may find this figure daunting at first, but careful study as you read the following description will reward you with a clearer picture of just what you are seeing when you look at a 2-D NMR spectrum.

#### Section 10.2 NMR models

The 2-D experiment is much more complex than obtaining a conventional "onedimensional" spectrum. The experimenter programs a sequence of pulses, delays, and data (FID) collections. In each sequence, the program calls for a 90° pulse (the "preparation" pulse), followed by a delay or "evolution" time  $t_1$ , and then a second 90° "mixing" pulse, followed by "detection" or data collection (all repeated 64 times to enhance the signal). In successive preparation/evolution/mixing/detection (PEMD) sequences, the evolution time  $t_1$  is increased in equal steps. In a typical experiment of this type,  $t_1$  might be incremented 512 times in 150- $\mu$ s steps from 0  $\mu$ s up to a maximum delay of about 77 ms, giving a full data set of 512 signalenhanced FIDs, five of which are shown with their FTs in Fig. 10.5*a*. Each FID is a plot of RF signal intensity versus time  $t_2$ , and its FT shows intensity as a function of frequency  $F_2$ . In each of the FTs (except the first one, in which the signals are too weak—see Fig. 10.1*d* to see why), you can see the RF signals of the four nuclei in the sample.

What is the purpose of the second or mixing pulse? Recall that the first pulse tips the precessing magnetic vectors toward the y-axis and aligns their phases, thus putting rotating net magnetic vectors in the xy-plane and an FID signal into the detector. But in 2-D NMR, we do not record this FID. Instead, we wait for a specified evolution time  $t_1$  (between 0 and 77 ms in the example described here), and then pulse again. During this interval-the evolution time-magnetic vectors spread out in the xy-plane, and the RF signal diminishes in intensity. The second pulse tips any magnetization that currently has a component in the yz-plane (as a result of vector spreading) back onto the xy-plane, and aligns the spins on the y-axis. Then relaxation occurs again, and this time we record the FID. In the first sequence, with  $t_1 = 0$ , essentially no relaxation occurs between pulses, no magnetization enters the yz-plane, and after the second pulse, the FID contains weak or no signals. As  $t_1$  is increased, more relaxation occurs between the first and second pulses, magnetization representative of the precession state of each nucleus at time  $t_1$  is present in the yz-plane, and it is detected after the second pulse.

The FIDs are Fourier transformed to produce 512 frequency spectra, each from a different time  $t_1$ , as shown on the right in Fig.10.5*a*. Just as in 1-D NMR, each spectrum gives the intensities of RF signals as a function of frequency  $F_2$  but, in this case, recorded from all nuclei after relaxation for a time  $t_1$ . To assign couplings, we are interested in finding out which pairs of RF frequencies  $F_2$  vary in intensity at the same rate. This can be accomplished by plotting the signal at each frequency  $F_2$  versus time  $t_1$ . The data for such a plot for nucleus 1 is shown within the narrow vertical rectangle (red) in Fig. 10.5*a*, and the plot itself is shown by the FID symbol at the top left in Fig. 10.5*b*. As implied by the symbol for this plot, at a frequency  $F_2$  at which an RF signal occurs, a plot of its intensity versus  $t_1$  is a time-domain spectrum, actually a pattern of interference among all frequencies present. This pattern is mathematically like an FID, and for convenience, it is usually called an FID. Recall that FIDs are composed of one or more component frequencies. If a nucleus is coupled to more than one other nucleus, its FID taken over time



**Figure 10.5** Two-dimensional NMR for a hypothetical system of four nuclei. (*a*) FIDs are collected after each of a series of sequences  $[90^{\circ} \text{ pulse}/t_1 \text{ delay}/90^{\circ} \text{ pulse}]$ , with  $t_1$  varied. Each FID decays over time  $t_2$ . The Fourier transform of an FID gives the RF intensities of each frequency  $F_2$  in the FID at time  $t_1$ . Each signal occurs at the chemical shift  $F_2$  of a nucleus in the sample. Chemical shifts of the nuclei are labeled 1 through 4.

#### Section 10.2 NMR models

**Figure 10.5** (Continued) Intensity at a single frequency  $F_2$  evolves over time  $t_1$ , as shown in the red vertical rectangle for nucleus 1. (b) A plot of this variation is like an FID that decays over time  $t_1$  [top of (b)]. Fourier transforms of  $t_1$ -FIDs reveal the frequencies  $F_1$  present in the signal, such as the red frequencies that make up the signal from nucleus 1. Coupled nuclei have frequencies  $F_1$  in common, as shown for nuclei 1, 2, and 4 in the green horizontal rectangle. They share decay frequencies because they are coupled, which means they are exchanging energy with each other. (c) Contour plot presenting the information of FTs in (b). Peaks on the diagonal correspond to the 1-D NMR spectrum, with  $F_1$  and  $F_2$  both corresponding to chemical shifts. Rows of peaks off the diagonal indicate sets of nuclei that have frequencies  $F_1$  in common due to coupling. The third horizontal row of signals in (c) corresponds to the information in the green rectangle in (b). In this example, the couplings are 1 to 2, 2 to 1 and 4, 3 to 4, and 4 to 2 and 3. From *Principles of Biophysical Chemistry* by Van Holde *et al.*, copyright 1998; adapted by permission of Prentice-Hall, Inc., Upper Saddle River, NJ.

 $t_1$  will contain frequencies (corresponding to the spreading of vectors in the *xy*plane) corresponding to all of its couplings. To determine these frequencies, we once again use the Fourier transform. The FT of this FID will reveal the different frequencies for each coupling. This FT is a frequency-domain spectrum, a plot of RF signal intensity versus phase-spreading frequencies or spin-spin-relaxation frequencies  $F_1$ , shown for the four  $F_2$  signals in the lower part of Fig. 10.5*b*. The frequencies  $F_1$  in each spectrum are relaxation frequencies for a given nucleus. The intensities in these spectra give the strength of coupling, or to put it another way, they give the amount of correlation between nuclei whose  $F_1$  FIDs contain the same frequency.

A 2-D NMR spectrum is a plot of  $F_1$  versus  $F_2$ , that is, it is our set of 512  $F_1$  spectra laid side by side at a spacing corresponding to the frequency  $F_2$  at which each  $F_1$  FID was taken, as in Fig. 10.5*b*. Looking along horizontals across the  $F_1$  spectra in (*b*), we see the frequencies  $F_2$  of sets of nuclei that relaxed at the same rate. For example, the narrow horizontal rectangle (green) in Fig. 10.5*b* shows that spins 1, 2, and 4 share a common relaxation frequency. Thus we conclude that these nuclei are coupled to each other and are on neighboring atoms. The 2-D spectrum is usually presented as shown in Fig. 10.5*c*. The position and intensity of a spot in this spectrum corresponds to the position and intensity of a signal in (*b*).

The off-diagonal or "cross" peaks have finer details, not shown here, that correspond to the simple and familiar, but information-rich, *J*-splitting patterns, such as doublets, triplets, and quartets, seen in conventional <sup>1</sup>H-NMR spectra of small molecules. These patterns tell us how many chemically equivalent nuclei are involved in these couplings. Of course, each signal correlates most strongly with itself over time  $t_1$ , so the strongest signals lie on the diagonal. The spectrum on the diagonal corresponds to the one-dimensional NMR spectrum.

Interpretation of the 2-D spectrum in Fig. 10.5c is simple. The spins of the four nuclei give strong signals on the diagonal (numbered 1 through 4). All off-diagonal signals indicate couplings. Each such signal is aligned horizontally along

 $F_2$  and vertically along  $F_1$  with signals on the diagonal that correspond to the two nuclei, or sets of nuclei, that are coupled. The signals labeled 2–4 denote coupling between nuclei 2 and 4. This type of 2-D NMR spectroscopy, which reveals spin-spin or *J* couplings by exhibiting correlations between spin-spin relaxation times of nuclei, is called *correlation spectroscopy* or *COSY*. In a sense, it spreads the information of the one-dimensional spectrum into two dimensions, keeping coupled signals together by virtue of their correlated or shared relaxation times. By far the most widely useful nucleus for NMR COSY study of proteins is <sup>1</sup>H. Obtaining <sup>1</sup>H spectra (rather than <sup>13</sup>C or <sup>15</sup>N spectra, for example) simply means carrying out the NMR experiments described here in the relatively narrow range of RF frequencies over which all <sup>1</sup>H atoms absorb energy in a magnetic field.

## Nuclear Overhauser effect

The COSY spectra reveal through-bond couplings because these couplings are involved in the relaxation process revealed by this PEMD sequence of operations:  $[90^{\circ} \text{ pulse}/t_1 \text{ delay}/90^{\circ} \text{ pulse}/data \text{ collection}]$ . In a sense, this pulse sequence stamps the  $t_1$  coupling information about relaxation rates onto the  $t_2$  RF absorption signals. The second set of Fourier transforms extracts this information. Other pulse sequences can stamp other kinds of information onto the  $t_2$  signal.

In particular, there is a second form of coupling between nuclear dipoles that occurs *through space*, rather than through bonds. This interatomic interaction is called the nuclear Overhauser effect (NOE), and the interaction can either weaken or strengthen RF absorption signals. Nuclei coupled to each other by this effect will show the effect to the same extent, just as *J*-coupled nuclei share common spin-spin relaxation rates. Appropriate pulse sequences can impress NOE information onto the  $t_2$  absorption signals, giving 2-D NMR spectra in which the off-diagonal peaks represent NOE couplings, rather than *J* couplings, and thus reveal pairs of nuclei that are near each other in space, regardless of whether they are near each other through bonds. This form of 2-D NMR is called *NOESY*. Though the off-diagonal signals represent different kinds of coupling, NOESY spectra look just like COSY spectra.

A partial NOESY spectrum of human thioredoxin is shown in Fig. 10.6. This 2-D spectrum shows the region from  $\delta = 6.5$  to  $\delta = 10$ . First, compare this spectrum with the 1-D spectrum in Fig. 10.1, p. 240. Recall that the signals on the diagonal of a 2-D spectrum are identical to the 1-D spectrum, but they are in the form of a contour map, sort of like looking at the 1-D spectrum from above, or down its intensity axis. For example, notice the two small peaks near  $\delta = 10$  on the 1-D spectrum, and near  $F_1 = F_2 = 10$  on the 2-D spectrum. Notice that there are off-diagonal peaks aligned horizontally and vertically with these peaks. These off-diagonal signals align with other diagonal peaks corresponding to nuclei NOE-coupled to these nuclei. Several coupling assignments that were used to define distance restraints are indicated by pairs of lines—one horizontal, one vertical—on the spectrum. The pairs converge on an off-diagonal peak and diverge to the signals of coupled nuclei on the diagonal.



**Figure 10.6** NOESY spectrum of thioredoxin in the region  $\delta = 6.5$  to  $\delta = 10$ . Pairs of lines—one horizontal, one vertical—converge at off-diagonal peaks indicating NOE couplings and diverge to the signals on the diagonal for the coupled nuclei. Assigning resonances on the diagonal to specific hydrogens in the protein requires greater resolution and simplification of the spectrum than shown here. The off-diagonal peak labeled "F89N, $\delta$ " indicates a specific NOESY interaction that is described and shown in Fig. 10.10, p. 256. Spectrum generously provided by Professor John M. Louis.

# NMR in higher dimensions

For a large protein, even spectra spread into two dimensions by J or NOE coupling are dauntingly complex. Further simplification of spectra can be accomplished by three- and higher-dimensional NMR, in which the information of a complex 2-D spectrum is separated onto sets of spectra that can be pictured as stacked planes of 2-D spectra, as in Fig. 10.7. One technique separates the 2-D spectrum into separate planes, each corresponding to a specific chemical shift of a different NMR-active nucleus, such as <sup>13</sup>C or <sup>15</sup>N, to which the hydrogens are bonded. This method, called <sup>13</sup>C- or <sup>15</sup>N-editing, requires that the protein be heavily labeled with these isotopes, which can be accomplished by obtaining the protein from cells grown with only <sup>13</sup>C- and <sup>15</sup>N-containing nutrients. In <sup>13</sup>C-edited COSY or NOESY spectra, each plane contains crosspeaks only for those hydrogens attached



**Figure 10.7** 3-D and 4-D NMR. (*a*) 2-D NMR of H-C<sub>a</sub>/H-C<sub>b</sub>/H-N<sub>a</sub>/H-N<sub>b</sub>. Chemical shifts of H-C<sub>a</sub> and H-C<sub>b</sub> are identical, as are those of H-N<sub>a</sub> and H-N<sub>b</sub>. so only two peaks appear on the diagonal. One cross-peak could arise from four possible couplings, listed on the figure. (*b*) Spectrum of (*a*) separated onto planes, each at a different <sup>13</sup>C chemical shift  $\delta$ . The  $\delta$ -<sup>13</sup>C axis represents the third NMR dimension. Gray lines represent the position of the diagonal on each plane. The cross-peak occurs only on the plane at  $\delta$ -<sup>13</sup>C<sub>b</sub>, revealing that coupling involves H-C<sub>b</sub> and eliminating two possible couplings. (*c*) Spectrum as in (*b*), separated onto planes, each at a different <sup>15</sup>N chemical shift  $\delta$ . The  $\delta$ -<sup>15</sup>N axis represents the fourth NMR dimension. The cross-peak occurs only on the plane at  $\delta$ -<sup>15</sup>N<sub>a</sub>, revealing that coupling involves H-N<sub>a</sub>. Thus the coupling indicated by the-cross peak in (a) is H-C<sub>b</sub> to H-N<sub>a</sub>.

to <sup>13</sup>C atoms having the chemical shift corresponding to that plane. Because you can picture these planes as stacked along a <sup>13</sup>C chemical-shift axis perpendicular to  $F_1$  and  $F_2$  (as illustrated in Fig. 10.7*b*), this type of spectroscopy is referred to as *three-dimensional NMR*. Further separation of information on these planes onto additional planes, say of <sup>15</sup>N chemical shifts (Fig. 10.7*c*), is the basis of so-called 4-D and higher dimensional NMR. The dimensions referred to are not really spatial dimensions, but are simply subsets of the data observed in 2-D spectra.

Figure 10.7 shows how 3-D and higher dimensional NMR can simplify spectra and reveal the details of couplings. Consider a hypothetical example of four hydrogens—two bonded to nitrogen,  $N_a$ -H and  $N_b$ -H, and two bonded to carbon,  $C_a$ -H and  $C_b$ -H. The two hydrogens of each pair have the same chemical shift, so the 1-D NMR spectrum shows two peaks, shown on the diagonal in Fig. 10.7*a*. In addition, there is one cross-peak, indicating that one H-C is coupled to one H-N. There are four possibilities for this coupling, listed in (*a*). If we separate

these spectra, as shown in (*b*), into planes that show only cross-peaks for protons attached to <sup>13</sup>C atoms having a narrow range of chemical shifts, the cross-peak occurs only on the plane corresponding to  $C_b$ . This means that  $C_b$  carries one of the coupled hydrogens and reduces the possible couplings to those listed as 3 and 4 in (*a*). If we again separate the spectra, as in (*c*), into planes based on <sup>15</sup>N shifts, we find that the cross-peak occurs only on the N<sub>a</sub> plane, eliminating all but possibility 3, coupling of  $C_b$ -H to N<sub>a</sub>-H. Separating the information onto additional planes in this manner does not really produce any new information. Instead it divides the information into subsets that may be easier to interpret.

# 10.2.3 Assigning resonances

The previous section describes methods that can provide an enormous amount of chemical-shift and coupling information about a protein. Recall that our goal is to determine the protein's conformation. We hope that we can use couplings to decide which pairs of hydrogens are neighbors, and that this information will restrict our model's conformations to one or a few similar possibilities. But before we can use the couplings, we must assign all the resonances in the 1-D spectrum to specific protons on specific residues in the sequence. This is usually the most laborious task in NMR structure determination, and I will provide only a brief sketch of it here.

The 1-D NMR spectrum of thioredoxin in Fig. 10.1, p. 240 shows the range of chemical shifts  $\delta$  for various proton types found in a protein. Because there are only about 20 different types of residues present, and because residues have many proton types in common (most have one proton on amide N, one on C<sub> $\alpha$ </sub>, two on C<sub> $\beta$ </sub>, and so forth), the spectrum is composed of several very crowded regions. Obviously, we could never assign specific signals to specific residues without greatly expanding or simplifying the spectrum. A variety of through-bond correlation and NOE experiments allow us to resolve all these signals. Given the necessary resolution, we are faced first with identifying individual residue types, say, distinguishing alanines from valines.

The cross-peaks in through-bond spectra guide us in this task, because each residue type has a characteristic set of splittings that makes it identifiable. This allows us to determine which of the many  $C_{\alpha}$ ,  $C_{\beta}$ , and  $C_{\gamma}$  signals, each in its own characteristic region of the spectrum, belong to the same residue. Thus spin systems that identify specific residue types can be identified and grouped together. This is usually the first stage of analysis: identification of spin systems and thus of groups of signals that together signify individual residues. There is conformational information in these couplings also, because the precise values of *J*-coupling constants depend on the conformational angles of bonds between atoms carrying coupled atoms. So the splitting patterns reveal residue identities, and actual values of splitting constants put some restrictions on local conformations.

Next, of all those spin systems identified as, say, valine, how do we determine which valine is valine-35, and which is valine-128? The next part of the analysis is making this determination, for which we must examine connectivities between residues. You can imagine that we might step from atom to atom through the protein chain by finding successively coupled proton pairs. But because proton spin-spin couplings are usually observed through at most three bonds (for example, H-C-N-H), and because the carbonyl carbon of each residue carries no hydrogens, proton spin-spin couplings do not extend continuously through the chain. In between two main-chain carbonyl carbons, the spin-spin coupling systems are isolated from each other. But there are interresidue spin-spin couplings because the N-H of residue *n* can couple to the  $C_{\alpha}$ -H of residue *n* – 1.

Covalent connectivities between residues can also be determined using proteins labeled fully with <sup>13</sup>C and <sup>15</sup>N, by way of various three-dimensional double and triple resonance experiments, each of which reveals atoms at opposite ends of a series of one-bond couplings. In essence, a sequence of pulses transfers magnetization sequentially from atom to atom. Coupling sequences such as H-N-C<sub> $\alpha$ </sub>, H-N-(CO)-C<sub> $\alpha$ </sub>, H-N-(CO), H-C<sub> $\alpha$ </sub>-(CO), and H-C<sub> $\alpha$ </sub>-(CO)-N can be used to "walk" through the couplings in a protein chain, establishing neighboring residues. Of course, as in crystallography, much of this work can be automated, but the operator must intervene when software cannot make the decisions unequivocally.

Taking NOE connectivities into account brings more potential correlations to our aid. There is a high probability that at least one proton among the N-H,  $C_{\alpha}$ -H, or  $C_{\beta}$ -H of residue *n* will be within NOE distance of the N-H of residue *n* + 1. This means that NOE correlations can also help us determine which residues are adjacent. In the end, this kind of analysis yields a complete set of resonance assignments. In addition, it reveals much about backbone and sidechain conformations because knowing which protons interact by NOE greatly restricts the number of possible conformations. These restrictions are a key to determining the conformation of the protein.

As you might suspect, assigning closely spaced chemical shifts to many similar atoms and assigning myriad correlation peaks to many similar atom pairs is fruitful turf for Bayesian methods (Sec. 7.5.6, p. 164). Indeed, these methods are being applied successfully in sorting out the enormous numbers of chemical-shift and correlation signals that must be assigned in order to proceed to structure determination. Bayesian methods are also useful in assessing the likelihoods of models that fit all or most of these assignments. As you will see in the next section, constructing models that conform to all of the distance restrictions obtained from analysis of signals is the essence of NMR structure determination.

#### 10.2.4 Determining conformation

Recall that our goal (you might well have forgotten by now) is to determine, in detail, the conformation of the protein. From the analysis described in the previous section, we know which signals in the spectra correspond to which residues in the sequence. In addition, the magnitudes of *J* couplings put some restraints on local conformational angles, both main-chain and side-chain, for each residue. What is more, the specific NOE couplings between adjacent residues also restrain local conformations to those that bring the coupled atoms to within NOE distance of each other (the range of distances detected by NOESY can be determined by the details of pulse timing during data collection).

#### Section 10.2 NMR models

Nuclear Overhauser effect correlations must be interpreted with care, because the NOE couples atoms through space, not through bonds. So NOE coupling between protons on two residues may not mean that they are adjacent in sequence. It may instead mean that the protein fold brings the residues near to each other. Distinguishing NOE cross-peaks between sequential neighbors from those between residues distant in the sequence sets the stage for determining the conformation of the protein. Once we have assigned all neighboring-residue NOESY cross-peaks properly, then the remaining cross-peaks tell us which hydrogens of sequentially distant residues are interacting with each other through space. This adds many of the most powerful and informative entries to the list of distance restraints with which our final model must agree.

The end result of analyzing NMR spectra is a list of distance restraints. The list tells us which pairs of hydrogens are within specified distances of each other in space. This is really about all that we learn about the protein from NMR. But as in X-ray crystallography, we already have a lot of *prior knowledge* about the protein. We know its chemical structure in detail because we know the full sequence of amino-acid residues, as well as the full structures of any cofactors present. We know with considerable precision all of the covalent bond lengths and bond angles. We know that most single-bond dihedral angles will lie within a few degrees of the staggered conformations we call *rotamers*. We know that amide single-bond dihedral angles will be close to  $180^{\circ}$ . We know that most main-chain conformational angles  $\Phi$  and  $\Psi$  will lie in the allowed zones of the Ramachandran diagram. We know all these things about any protein having a known sequence. The question is, does all this add up to knowing the three-dimensional conformation of the protein?

Next we would like to build a model of the protein that fits all we know, including the distance restraints we learn from NMR. This is no trivial task. Much research has gone into developing computer algorithms for this kind of model building. One general procedure entails starting from a model of the protein having the known sequence of residues, and having standard bond lengths and angles but random conformational angles. This starting structure will, of course, be inconsistent with most of the distance and conformational restraints derived from NMR. The amount of inconsistency can be expressed as a numerical parameter that should decline in value as the model improves, in somewhat the same fashion as the *R*-factor decreases as a crystallographic model's agreement with diffraction data improves during refinement.

Starting from a random, high-temperature conformation, simulated annealing or some form of molecular dynamics is used to fold our model under the influence of simulated forces that maintain correct bond lengths and angles, provide weak versions of van der Waals repulsions, and draw the model toward allowed conformations, as well as toward satisfying the restraints derived from NMR. Electrostatic interactions and hydrogen-bonding are usually not simulated, in order to give larger weight to restraints based on experimental data; after all, we want to *discover* these interactions in the end, not build them into the model before the data have had their say. The simulated folding process entails satisfying restraints locally at first, and then gradually satisfying them over greater distances. At some points in the procedure, the model is subjected to higher-temperature dynamics, to shake it out of local minima of the consistency parameter that might prevent it from reaching a global best fit to the data. Near the end, the strength of the simulated van der Waals force is raised to a more realistic value, and the model's simulated temperature is slowly brought down to about 25°C.

Next the model is examined for serious van der Waals collisions, and for large deviations from even one distance or conformational restraint. Models that suffer from one or more such problems are judged not to have converged to a satisfactory final conformation. They are discarded. The entire simulated folding process is carried out repeatedly, each time from a different random starting conformation, until a number of models are found that are chemically realistic and consistent with all NMR-derived restraints. If some of these models differ markedly from others, the researchers may try to seek more distance restraints in the NMR spectra that will address specific differences, and repeat the process with additional restraints. It may be possible at various stages in this process to use the current models to resolve previous ambiguities in NOE assignments and to include them in further model building. When the group of models appears to contain the full range of structures that satisfy all restraints, this phase of structure determination is complete.

The result of the NMR structure determination is thus not a single model, but a group or *ensemble* of from half a dozen to over 30 models, all of which agree with the NMR data. The models can be superimposed on each other by least-squares fitting (Sec. 7.5.1, p. 159, and Sec. 11.4.1, p. 288) and displayed as shown in Fig. 10.8 for thioredoxin (PDB 4trx). In this display of 10 of the 33 models



**Figure 10.8** ► Ten of the 33 NMR models of thioredoxin (stereo, PDB 4trx). Image: DeepView/POV-Ray.

#### Section 10.2 NMR models

obtained, it is readily apparent that all models agree with each other strongly in some regions, often in the interior, whereas there is more variation in other regions, usually on the surface.

Researchers interested in understanding this protein will immediately want to use the results of the NMR work as an aid to understanding the function of this protein. Which of these perhaps two or three dozen models should we use? For further studies, as well as for purposes of illustration, it is appealing to desire a single model, but one that will not let us forget that certain regions may be more constrained by data than others, and that certain regions may be different in different models within the ensemble. One way to satisfy the common, but perhaps uninformed (see Sec. 10.2.6, p. 257), desire for a single model is to compute the average position for each atom in the model and to build a model of all atoms in their average position. This model may be unrealistic in many respects. For example, bond lengths and angles involving atoms in their averaged positions may not be the same as standard values. This averaged model is then subjected to restrained energy minimization, which in essence brings bond lengths and angles to standard values, minimizes van der Waals repulsions, and maximizes noncovalent interactions, with minimal movement away from the averaged atomic coordinates. The result is a single model, the restrained minimized average structure. The rms deviation in position from the average of the ensemble is computed for each atom in this model, and the results are recorded in the *B*-factor column of the atomiccoordinate file. Deviations are low in regions where models in the ensemble agree well with each other, and deviations are high where the models diverge. Figure 10.9



**Figure 10.9** Energy-minimized averaged model (stereo, PDB 3trx), colored by rms deviation of atom positions in the ensemble from the average position. For each residue, main-chain colors reflect the average rms deviation for C, O, N, and CA, and all side-chain atoms are colored to show the average rms deviation for atoms in the whole side chain. Image: DeepView/POV-Ray.

shows the final model of thioredoxin (PDB 3trx) derived from the full ensemble partially represented in Fig. 10.8. Colors of the model represent rms deviations of atom positions from the ensemble average, with blue assigned to the smallest deviations, red to the largest, and colors across the visible spectrum for intermediate values.

How much structural information must we obtain from NMR in order to derive models like those shown in Fig. 10.8? As summarized in the PDB 3trx file header for thioredoxin (remember, READ THE HEADER!), the models shown in Fig. 10.8 were determined from 1983 interproton distance restraints derived from NOE couplings. In addition, the researchers used 52 hydrogen-bonding distance restraints defining 26 hydrogen bonds. Detection of these hydrogen bonds was based on H/D exchange experiments, in which spectra in H<sub>2</sub>O and D<sub>2</sub>O were compared, and NOEs involving exchangeable H disappeared or were diminished in D<sub>2</sub>O (deuterium does not give an NMR signal). All hydrogen bonds detected in this manner involved exchangeable amide hydrogens. Once these hydrogen bonds were detected in early model construction, they were included in the restraints used in further calculations. Finally, there were 98  $\Phi$  and 71  $\Psi$  backbone dihedral-angle restraints, and 72  $\chi_1(C_\beta - C_\gamma)$  side-chain dihedral-angle restraints, derived from NOE and J-coupling. Thus the conformation of each of the 33 final thioredoxin models is defined by a total of 2276 restraints. Recall that thioredoxin is a relatively small protein of 105 residues, so these models are based on about 20 restraints per residue. Finally, for an example of the effect of a single distance restraint on the final model, take another look at Fig. 10.6, p. 249, and Fig. 10.10. In Fig. 10.6, p. 249, one of the NOESY off-diagonal signals is labeled "F89N, $\delta$ ." This notation



**Figure 10.10**  $\blacktriangleright$  Detail (stereo) of the averaged model at phenylalanine-89, showing the averaged distance between the two hydrogens inolved in the distrance restraint indicated as "F89N, $\delta$ " in the NOESY spectrum, Fig. 10.6. Image: DeepView/POV-Ray.

means that the correlated nuclei are in phenylalanine-89 (F89), and specifically, that the hydrogen on the amide nitrogen is correlated with the hydrogen on one of the  $C_{\beta}$  carbons, which are on the aromatic ring, *ortho* to the side-chain connection. If you compare this figure with Fig. 10.1, p. 240, you can identify the diagonal signals: the signal at approximately  $\delta = 9.0$  is the amide hydrogen, and the signal at  $\delta = 7.3$  is the aromatic hydrogen. Figure 10.10 shows the distance between these two hydrogens in the final averaged model. You can see how this distance restraint confines the conformation of the phenylalanine side chain in this model.

# 10.2.5 PDB files for NMR models

Like crystallographers, NMR spectroscopists share their models with the scientific community by depositing them in the Protein Data Bank. For NMR models, two PDB files sometimes appear—one containing all models in the ensemble, and one containing the coordinates of the restrained minimized average (although the trend is away from depositing averaged structures). For human thioredoxin, these two files are 4trx (ensemble of 33 models) and 3trx (averaged model). In parallel with deposition of structure factors for crystallographic models, the most responsible research groups also deposit an accessory file listing all distance restraints used in arriving at the final models (PDB id code 3trx.mr for this model).

At first glance, the averaged model would appear to serve most researchers who are looking for a molecular model to help them explain the function of the molecule and rationalize other chemical, spectroscopic, thermodynamic, and kinetic data. On the other hand, you might think that the ensemble and distancerestraint files are of most use to those working to improve structure determination techniques. There are good reasons however, for *all* researchers to look carefully at the ensemble, as discussed in the next section.

As with PDB files for crystallographic models, NMR coordinate files also include headers containing citations to journal articles about the structure determination work, as well as brief descriptions of specific techniques used in producing the model presented in the file, including the numbers and types of restraints. When you view PDB files in web browsers, the literature citations contain convenient live links to PubMed abstracts of the listed articles.

## 10.2.6 Judging model quality

The most obvious criterion of quality of an NMR model is the level of agreement of models in the ensemble. You can make this assessment qualitatively by viewing the superimposed models in a molecular graphics program like the one I will describe the Chapter 11. You should be particularly interested in agreement of the models in regions of functional importance, including catalytically active or ligand-binding sites. The PDB file for the averaged model usually contains information to help you make this assessment. In particular, the data column that contains *B*-factors for crystallographic models contains, for NMR models, the rms deviations from average ensemble coordinate positions. As with *B*-factors, rms deviations are smaller in the main chain than in side chains. The best quality models exhibit main-chain deviations no greater than 0.4 Å, with side-chain values below 1.0 Å.

Many graphics programs will color the averaged model according to these deviations, as illustrated in Fig. 10.9. This coloring is a useful warning to the user, but it is not as informative as the deviations themselves because the graphics program, depending on its settings, may give colors from blue to red relative to the range of values in the PDB file, regardless of how narrow or wide the range. Thus two models colored by rms deviations may appear similar in quality, but one will have all rms deviations smaller than 3 Å, while the other will have much higher values. You can best compare the quality of two different models by coloring rms deviations on an absolute scale, rather than the relative scale of the contents of a single PDB file.

It is also useful to think about why rms deviations might vary from region to region in the averaged model. In crystallographic models, higher *B*-factors in sections of a well-refined model can mean that these sections are dynamically disordered in the crystal, and thus moving about faster than the time scale of the data collection. The averaged image obtained by crystallography, just like a photo of a moving object, is blurred. On the other hand, high *B*-factors may mean static disorder, in which specific side chains or loops take slightly different conformations in different unit cells. Recall that the electron-density map sometimes reveals alternative conformations of surface side chains. Are there analogous reasons for high rms deviations in sections of an NMR model?

The proximate reason for high rms values is the lack of sufficient distance restraints to confine models in the ensemble to a particular conformation. Thus various models converge to different conformations that obey the restraints, but several or many conformations may fill the bill. At the molecular level, there are several reasons why NMR spectroscopy may not provide enough distance restraints in certain regions of the model. The physically trivial reason is spectral resolution. Either resolution may simply not be good enough for all couplings to be resolved and assigned, or some spectral peaks may overlap simply because the structures and environments they represent are practically identical. This problem is worst in the most crowded regions of the spectrum and might persist despite multidimensional simplification. Because a protein NMR spectrum is composed of many similar spin systems, like many valines, some may simply have such similar chemical shifts and couplings that portions of them cannot be distinguished. Alternatively, the width of spectral peaks may be too wide to allow closely spaced peaks to be resolved, and if peak widths exceed coupling constants, correlations cannot be detected. Peak width is inversely related to the rate at which the protein tumbles in solution, and larger proteins tumble more slowly, giving broader peaks. This is unfortunate, of course, because for larger proteins, you need better resolution. But three- and four-dimensional spectra simplify these complex spectra, and the use of one-bond couplings circumvents the correlation problem because one-bond coupling constants are generally much larger than three-bond proton coupling constants. By the late 1990s, top NMR researchers were claiming that it is now potentially feasible to determine protein structures in the 15- to 35-kDa range (135 to 320 residues) at an accuracy comparable to crystallographic models at 2.5-Å resolution. They see the upper limit of applicability of methods described

#### Section 10.3 Homology models

here as probably 60–70 kDa. In early 2005, when I searched the PDB for the NMR models of the largest molecules, I found 10 models of more than 300 residues and 28 of more than 250 residues, out of the approximately 4500 NMR models in the database.

Beyond limitations in resolution are more interesting reasons for high-rms variations within an ensemble of NMR models. Lack of observed distance restraints may in fact mean a lack of structural restraints in the molecule itself. Large variations among the models may be pointing us to the more flexible and mobile parts of the structure in solution. In fact, if we are satisfied that our models are not limited by resolution or peak overlap, we can take seriously the possibility that variations in models indicate real dynamic processes. (Relaxation experiments can be used to determine true flexibility.) Sometimes, mobility may take the form of two or a few different conformations of high-rms regions of the real molecule, each having a long enough lifetime to produce NOE signals. The distance inferred from NOE intensities would be an average figure for the alternative conformations. The ensemble of final models would then reveal conformations that are compatible with averaged distance restraints, and thus point to the longer-lived conformations in solution.

If some or all of the ensemble conformations reveal actual alternative conformations in solution, then these models contain useful information that may be lost in producing the averaged model. If the most important conformations for molecular function are represented in subsets of models within the final ensemble, then an averaged model may mislead us about function. Just like crystallographic models, NMR models do not simply tell us what we would like to know about the inner workings of molecules. Evidence from other areas of research on the molecule are necessary in interpreting what NMR models have to say.

As for other indicators of model quality, we expect that NMR models, just like crystallographic models, should agree with prior knowledge about protein structure. So Ramachandran diagrams and distributions of side-chain conformations should meet the same standards of quality as those for any other type of model.

For an example of a recent NMR structure determination, see NMR structure of Mistic, a membrane-integrating protein for membrane protein expression, T. Roosild, *et al.*, *Science* **307**, 1317 (2005).

# 10.3 Homology models

## 10.3.1 Introduction

I have repeatedly reminded you that protein structure determination is a search for the conformation of a molecule whose chemical composition is known. Much experience supports the conclusion that proteins with similar amino-acid sequences have similar conformations. These observations suggest that we might be able to use proteins of known structure as a basis for building models of proteins for which we know only the amino-acid sequence. This type of structure determination has been called *knowledge-based modeling* but is now commonly known as *comparative protein modeling* or *homology modeling*. We refer to the protein we are modeling as the *target*, and to the proteins used as frameworks as *templates*.

If, in their core regions, two proteins share 50% sequence identity, the alpha carbons of the core regions can almost always be superimposed with an rms deviation of 1.0 Å or less. This means that the core region of a protein of known structure provides an excellent template for building a model of the core region of a target protein having 50% or higher sequence homology. The largest structural differences between homologous proteins, and thus the regions in which homology models are likely to be in error, lie in surface loops. So comparative modeling is easiest and most reliable in the core regions, and it is the most difficult and unreliable in loops. The structures resulting from homology modeling are, in a sense, low-resolution structures, but they can be of great use, for example, in guiding researchers to residues that might be involved in protein function. Hypotheses about the function of these residues can then be tested by looking at the effects of site-directed mutagenesis. Homology models may also be useful in explaining experimental results, such as spectral properties; in predicting the effects of mutations, such as site-directed mutations aimed at altering the properties or function of a protein; and in designing drugs aimed at disrupting protein function.

Because much of the homology modeling process can be automated, it is possible to develop databases of homology models automatically as genome projects produce new protein sequences. Users of such databases should be aware that automatically generated homology models can often be improved by user intervention in the modeling process (see Sec. 10.3.4, p. 263).

## 10.3.2 Principles

Comparative protein modeling entails the following steps: (1) constructing an appropriate template for the core regions of the target, (2) aligning the target sequence with the core template and producing a target core model, (3) building surface loops, (4) adding side chains in mutually compatible conformations, and (5) refining the model. In this section, I will give a brief outline of a typical modeling strategy. Although a variety of such strategies have been devised, my discussion is based on the program ProMod and on procedures employed by SWISS-MODEL, a public service homology modeling tool on the World Wide Web. SWISS-MODEL is one of many tools available on the ExPASy Molecular Biology Server operated by the Swiss Institute of Bioinformatics. You will find a link to ExPASy on the CMCC home page.

# Templates for modeling

The first step in making a comparative protein model is the selection of appropriate templates from among proteins of known structure (that is, experimental models derived from crystallographic or NMR data) that exhibit sufficient homology with the target protein. Even though it might seem that the best template would be the

#### Section 10.3 Homology models

protein having highest sequence homology with the target, usually two or several proteins of high percent homology are chosen. Multiple templates avoid biasing the model toward one protein. In addition, they guide the modeler in deciding where to build loops and which loops to choose. Multiple templates can also aid in the choice of side-chain conformations for the model.

Users find templates by submitting the target sequence for comparison with sequences in databases of known structures. Programs like BLAST or FastA carry out searches for sequences similar to the target. A BLAST score lower than  $10^{-5}$  (meaning a probability lower than one in 100,000 that the sequence similarity is a coincidence), or a FastA score 10 standard deviations above the mean score for random sequences, indicates a potentially suitable template. A safe threshold for automated modeling is 35% homology between target and templates. Below this threshold, alignment of template and target sequences may be unreliable, even though it may turn out that their three-dimensional structures are quite similar. Structural homology can, and often does, lurk at homology percentages too low for homology to be detectable.

Next, programs like SIM align the templates. Because many proteins contain more than one chain, and many chains are composed of more than one domain, modelers use databases, extracted from the Protein Data Bank, of single chains. There are also means to find and extract single domains homologous to the target from within larger protein chains. If only one homolog of known structure can be found, it is used as the template. If several are available, the template will be an averaged structure based on them. The chain most similar to the target is used as the reference, and the others are superimposed on it by means of least-squares fitting, with alignment criteria emphasizing those regions that are most similar (referred to as the "conserved" regions because sequence similarities presumably represent evolutionary conservation of sequence within a protein family). To obtain the best alignment, a program like ProMod starts by aligning alpha carbons from only those regions that share the highest homology with the reference. Alpha carbons from the other (nonreference) templates are added to the aligned, combined template if they lie within a specified distance, say 3.0 Å, of their homologous atoms in the reference. The result is called a structurally corrected multiple-sequence alignment.

## Modeling the target core

To build the core regions of the target protein, its sequence is first aligned with that of the template or, if several templates are used, with the structurally corrected multiple sequence alignment. The procedure aligns the target with all regions, including core and loops, that give high similarity scores, with the result that the core of the target aligns with the core regions of all models, whereas loop regions of the target align only with individual models having very similar loop sequences. This means that subsequent modeling processes will take advantage not only of the general agreement between target and all templates in core regions but also of the specific agreement between target and perhaps a single template that shares a very similar loop sequence.
With this sequence alignment, a backbone model of the target sequence is *threaded*—folded onto the aligned template atoms—producing a model of the target in the core regions and in any loops for which a highly homologous template loop was found. When multiple templates are used, the target atom positions are the best fit to the template atom positions for a target model that keeps correct bond lengths and angles, perhaps weighted more heavily toward the templates of highest local sequence homology with the target. As for residue side chains in the averaged template, in SWISS-MODEL, side-chain atom positions are averaged among different templates and added so that space is occupied more or less as in the templates. After averaging, the program selects, for each side chain, the allowed side-chain conformation or *rotamer* that best matches the averaged atom positions.

#### Modeling loops

At this stage, we have a model of the core backbone in which the atom positions are the average of the atom positions of the templates, and in which some core side chains are included. Perhaps some loops are also modeled, if one or more templates were highly homologous to the target. But most of the surface loops are not yet modeled, and in most cases, the templates give us no evidence about their structures. This is because the most common variations among homologous proteins lie in their surface loops. These variations include differences in both sequence and loop size. The next task is to build reasonable loops containing the residues specified by the target sequence.

One approach is to search among high-resolution ( $\leq 2.5$ -Å) structures in the Protein Data Bank for backbone loops of the appropriate size (number of residues) and end-point geometry. In particular, we are looking for loops whose conformation allows their own neighboring residues to superimpose nicely on the target loop's neighboring residues in the already modeled target core region. To hasten this search, modelers extract and keep loop databases from the PDB, containing the alpha-carbon coordinates of all loops plus those of the four neighboring residues, called "stem" residues, on both ends of each loop.

To build a loop, a modeling program selects loops of desired size and scores them according to how well their stem residues can be superimposed on the stem residues of the target core, aligning each prospective loop as a rigid body. The program might search for loops whose rms stem deviations are less than 0.2 Å and, finding none, might search again with a criterion of 0.4 Å, continuing until a small number, say five, candidates are found. The target loop is then modeled, its coordinates based on the average of five database loops with lowest rms-deviation scores. This process builds only an alpha-carbon model, so the amide groups must be added. This might involve another search through even higher-resolution structures in the PDB, looking for short peptides (not necessarily loops) whose alpha carbons align well with short stretches of the current loop atoms. Again, if several good fits are found, the added atoms are modeled on their average coordinates, thus completing the backbone of our model.

#### Modeling side chains

Our model now consists of a complete backbone, with side-chains only present on those residues where templates and target are identical. In regions of high sequence homology, target side chains nonidentical to templates might be modeled on the template side chains out to the gamma atoms. The remaining atoms and the remaining full side chains called for by the target sequence are then added, using rotamers of the side chains that do not clash with those already modeled, or with each other.

#### Refining the model

Our model now contains all the atoms of the known amino-acid sequence. Because many atom positions are averages of template positions, it is inevitable that the model harbors clashing atoms and less-than-optimal conformations. The end of all homology modeling is some type of structure refinement, including energy minimization. SWISS-MODEL uses the program GROMOS to idealize bond lengths and angles, remove unfavorable atom contacts, and allow the model to settle into lower-energy conformations lying near the final modeled geometry. The number of cycles of energy minimization is limited so that the model does not drift too far away from the modeled geometry. In particular, loops tend to flatten against the molecular surface upon extended energy minimization, probably because the model's energy is being minimized in the absence of simulated interactions with surrounding water.

A modeling process like the one I have described is applied automatically—but with the possibility of some user modifications—to target sequences submitted to the SWISS-MODEL server at ExPASy. This tool also allows for intervention at various stages, during which the user can apply special knowledge about the target protein or can examine and adjust sequence or structural alignments. DeepView (formerly called Swiss-PdbViewer)—an excellent free program for viewing, analyzing, and comparing models—is specifically designed to carry out or facilitate these interventions for SWISS-MODEL users, thus allowing a wider choice of tools for template searching, sequence alignment, loop building, threading, and refinement. You will learn more about DeepView in Chapter 11.

#### 10.3.3 Databases of homology models

As a result of the many genome-sequencing projects now under way, an enormous number of new structural genes (genes that code for proteins) are being discovered. With the sequences of structural genes comes the sequences of their product proteins. Thus new proteins are being discovered far faster than crystallographers and NMR spectroscopists can determine their structures. Although homology models are almost certain to be less accurate than those derived from experimental data, they can be obtained rapidly and automatically. Though they cannot guide detailed analysis of protein function, these models can guide further experimental work on a protein, such as site-directed mutagenesis to pinpoint residues essential to function. The desirability of at least "low-resolution" structures of new proteins has led

to initiatives to automatically model all new proteins as they appear in specified databases. Thus the number of homology models has already exceeded the number of experimentally determined structures in the Protein Data Bank.

One of the first automated modeling efforts was carried out in May of 1998. Called 3D-Crunch, the project entailed submission of all sequences in two major sequence databases, SWISS-PROT and trEMBL, to the SWISSMODEL server. The result was about 64,000 homology models (compared to about 9000 models in the Protein Data Bank at the time), which are now available to the public through the SWISS-MODEL Repository (see the CMCC home page). A versatile searching tool allows users to find and download final models. Alternatively, users can download entire modeling projects, containing the final coordinates of the homology model along with aligned coordinates of the templates. These project files enable the user to carry the modeling project beyond the automated process and to use other tools or special knowledge about the protein to further improve the model. Project files are most conveniently opened in DeepView, which provides both builtin modeling tools and interfaces to additional programs for further improvement of models. Figure 10.11 shows a modeling project returned from SWISS-MODEL. For more discussion of this figure, see the next section and the figure caption.



**Figure 10.11**  $\blacktriangleright$  Homology modeling project returned from SWISS-MODEL (stereo). The target protein is a fragment of FasL, a ligand for the widely expressed mammalian protein Fas. Interaction of Fas with FasL leads to rapid cell death by apoptosis. The template proteins are (1) tumor necrosis factor receptor P55, extracellular domain (PDB 1tnr, black) and (2) tumor necrosis factor-alpha (PDB 2tun, gray). The modeled FasL fragment is shown as ribbon and colored by model *B*-factors. Only the alpha carbons of the templates are shown. Image: DeepView/POV-Ray.

264

#### Section 10.3 Homology models

In addition to general homology-modeling tools like SWISS-MODEL and DeepView, there are specialized modeling tools and servers devoted to specific proteins of wide interest, including antibodies and membrane proteins such as G-protein-coupled receptors. See the CMCC home page for links to such sites.

#### 10.3.4 Judging model quality

Note that the entire comparative protein-modeling process is based on structures of proteins sharing sequence homology with the target, and that no experimental data about the target is included. This means that we have no criteria, such as *R*-factors, that allow us to evaluate how well our model explains experimental observations about the molecule of interest. *The model is based on no experimental observations*. At the worst, in inept hands, a homology modeling program is capable of producing a model of any target from even the most inappropriate template. Even the most respected procedures may produce dubious models. By what criteria do we judge the quality of homology models?

We would like to ask whether the model is *correct*. We could say it is correct if it agrees with the actual molecular structure. But we do not have this kind of assurance about any model, even one derived from experiment. It is more reasonable for us to define *correct* as agreeing to within experimental error with an experimental (X-ray or NMR) structure. Most of the time, we accept a homology model because we do not have an experimental structure, but not always. Researchers working to improve modeling methods often try to model known structures starting from related known structures. Then they compare the model with the known structure to see how the modeling turned out. In this situation, a correct model is one in which the atom positions deviate from those of the experimental model by less than the uncertainty in the experimental coordinates, as assessed by, for example, a Luzzati plot (Sec. 8.2.2, p. 183). Areas in which the model is incorrect are arenas for improving modeling tools.

Typically, however, we want to use a homology model because it is all we have. In some cases, even without an experimental model for comparison, we can recognize incorrectness. A model is incorrect if it is in any sense impossible. What are signs of an incorrect model? One is the presence of hydrophobic side chains on the surface of the model, or buried polar or ionic groups that do not have their hydrogen-bonding or ionic-bonding capabilities satisfied by neighboring groups. Another is poor agreement with expected values of bond lengths and angles. Another is the presence of unfavorable noncovalent contacts or "clashes." Still another is unreasonable conformational angles, as exhibited in a Ramachandran plot (Sec. 8.2.1, p. 181). We know that high-quality models from crystallography and NMR do not harbor these deficiencies, and we should not accept them in a homology model. Many molecular graphics programs can compute deviations from expected bond geometry; highlight clashes with colors, dotted lines, or overlapping spheres; and display Ramachandran diagrams, thus giving us immediate visual evidence of problems with models. We can also say that the model is incorrect if the sequence alignment is incorrect or not optimal. The details of sequence comparison are beyond the scope of this book, but we can test the alignment of target with templates by using different alignment procedures, or by altering the alignment parameters to see if the current alignment is highly sensitive to slight changes in method. If so, it should shake our confidence in the model.

Beyond criteria for correctness, we can also ask how well the model fits its templates. The rms deviation of model from template should be very small in the core region. If not, we say that the model is *inaccurate*. An inaccurate model implies that the modeling process did not go well. Perhaps the modeling program simply could not come up with a model that aligns well with the coordinates of the template or templates. Perhaps during energy minimization, coordinates of the model drifted away from the template coordinates. Another possibility is poor choice of templates. For instance, occasionally a crystallographic model is distorted by crystal contacts, or an NMR template model is distorted by the binding of a salt ion. If we unwittingly use such models as templates, energy refinement in the absence of the distorting effect would introduce inaccuracy, as defined here, while perhaps actually improving the model. A good rule of thumb is that if the templates share 30-50% homology with the target, rms differences between final positions of alpha carbons in the model and those of corresponding atoms in the templates should be less than 1.5 Å. But it is also essential to look at the template structures and make sure that they are really appropriate. An NMR structure of an enzyme-cofactor complex is likely to be a poor model for a homologous enzyme in the absence of the cofactor.

The rms deviations apply only to corresponding atoms, which means mainly the core regions. Loop regions often cannot be included in such assessment because there is nothing to compare them to. Again, we should demand correctness, that is, the lack of unfavorable contacts or conformations. But beyond this kind of correctness, our criteria are limited. If surface loops contain residues known to be important to function, we must proceed with great caution in using homology models to explain function.

If the model appears to be correct (not harboring impossible regions like clashes) and accurate (fitting its templates well), we can also ask if it is *reasonable*, or in keeping with expectations for similar proteins. Researchers have developed several assessments of reasonableness that can sometimes signal problems with a model or specific regions of a model. One is to sum up the probabilities that each residue should occur in the environment in which it is found in the model. For all Protein Data Bank models, each of the 20 amino acids has a certain probability of belonging to one of the following classes: solvent-accessible surface, buried polar, exposed nonpolar, helix, sheet, or turn. Regions of a model that do not fit expectations based on these probabilities are suspect.

Another criterion of reasonableness is to look at how often pairs of residues interact with each other in the model in comparison to the same pairwise interactions in templates or proteins in general. The sum of pairwise potentials for the model, usually expressed as an "energy" (smaller is better) should be similar to that for the templates. One form of this criterion is called *threading energy*. Threading energy indicates whether the environment of each residue matches what is found for the same residue type in a representative set of protein folds. Such criteria ask, in a sense, whether a particular stretch of residues is "happy" in its three-dimensional setting. If a fragment is "unhappy" by these criteria, then that part of the model may be in error.

To be meaningful, all of these assessments of reasonableness of a model must be compared with the same properties of the templates. After all, the templates themselves, even if they are high-quality experimental structures, may be unusual in comparison to the average protein.

Homology model coordinate files returned from SWISS-MODEL contain, in the B-factor column, a *confidence factor*, which is based on the amount of structural information that supports each part of the model. Actually, it would be better to call this figure an uncertainty factor, or a model B-factor, because a high value implies high uncertainty, or low confidence, about a specific part of the model. (Recall that higher values of the crystallographic *B*-factor imply greater uncertainty in atom positions.) The model B-factor for a residue is higher if fewer template structures were used for that residue. It is also higher for a residue whose alphacarbon position deviates by more than a specified distance from the alpha carbon of the corresponding template residue. This distance is called the *distance trap*. In SWISS-MODEL (or more accurately, in ProMod II, the program that carries out the threading at SWISS-MODEL), the default distance trap is 2.5 Å, but the user can increase or decrease it. However, if the user increases the distance trap, all of the model *B*-factors increase, so they still reflect uncertainty in the model, even if the user is willing to accept greater uncertainty. Finally, all atoms that are built without a template, including loops for which none of the template models had a similar loop size or sequence, are assigned large model *B*-factors, reflecting the lack of template support for those parts of the model. Computer displays of homology models can be colored by these model B-factors to give a direct display of the relative amount of information from X-ray or NMR structures that was used in building the model (Sec. 11.3.5, p. 281). Figure 10.11, p. 264 shows a homology model and its templates. The target model is colored by the model B-factors assigned by SWISS-MODEL. The templates are black and gray. With this color scheme, it is easy to distinguish the parts of the model in which we can have the most confidence. Blue regions were built on more templates and fit the templates better. Red regions were built completely from loop databases, without template contributions. Colors of the visible spectrum between blue and red may align well with none or only a subset of the templates. (For more information about these proteins, see the legend to Fig. 10.11, p. 264.)

### 10.4 Other theoretical models

Structural biologists produce other types of theoretical models as they pursue research in various fields. Work that produces new models includes developing schemes to predict protein conformation from sequence; attempts to simulate folding or other dynamic processes; and attempts to understand ligand binding by building ligands into binding sites and then minimizing the energy of the resulting combined model. The Protein Data Bank once contained some homology and theoretical models, but they were removed in 2002. According to online documentation, "The Protein Data Bank (PDB) is an archive of *experimentally determined* three-dimensional structures . . .." The presence of theoretical models in the Protein Data Bank was only a temporary measure due to the lack of a data bank for homology and other theoretical models. As of this writing, the only database for theoretical models is the SWISS-MODEL Repository mentioned in Sec. 10.3.4, p. 263.

At least two distributed-computing projects promise to produce large numbers of theoretical models. The Folding@Home and Human Proteome Folding projects are both attempts to simulate protein folding and discover folding principles that might eventually make it possible to predict protein conformation from aminoacid sequence. Both projects will produce many theoretical structures, but it is not known at the moment whether results will be available in searchable databases like the PDB or the SWISS-MODEL Repository. Look for links to model databases on the CMCC home page.

My goals in this chapter are to make you aware of the variety of model types available to the structural biologist, to give you a start toward understanding other methods of structure determination, and to guide you in judging the quality of noncrystallographic models, primarily by drawing your attention to analogies between criteria of quality in crystallography.

*Models are not molecules observed*. No matter how they are obtained, before we ask what they tell us, we must ask how well macromolecular models fit with other things we already know. A model is like any scientific theory: it is useful only to the extent that it supports predictions that we can test by experiment. Our initial confidence in it is justified only to the extent that it fits what we already know. Our confidence can grow only if its predictions are verified.

## ► Chapter 11

## Tools for Studying Macromolecules

## 11.1 Introduction

There is an old line about a dog who is finally cured of chasing cars—when it catches one. What now? In this chapter, I discuss what to do when you catch a protein. My main goal is to inform you about some of the tools available for studying protein models and to suggest strategies for learning your way around the unfamiliar terrain of a new protein. I will begin with a very brief glimpse of the computations that underlie molecular graphics displays. Then I will take you on a tour of molecular modeling by detailing the features present on most modeling programs. Finally, I will briefly introduce other computational tools for studying and comparing proteins. My emphasis is on tools that you can use on your personal computer, although today's personal computers do not limit your possibilities very much.

## 11.2 Computer models of molecules

#### 11.2.1 Two-dimensional images from coordinates

Computer programs for molecular modeling provide an interactive, visual environment for displaying and exploring models. The fundamental operation of computer programs for studying molecules is producing vivid and understandable displays convincing images of models. Although the details of programming for graphics displays vary from one program (or programming language, or computer operating system) to another, they all produce an image according to the same geometric principles.

A display program uses a file of atomic coordinates to produce a drawing on the screen. Recall that a PDB coordinate file contains a list of all atoms located by crystallographic, NMR, or theoretical analysis, with coordinates x, y, and zfor each atom. When the model is first displayed, the coordinate system is usually shifted by the modeling program so that the origin is the center of the model. This origin lies at the center of the screen, becoming the origin of a new coordinate system, the *screen coordinates*, which I will designate  $x_s$ ,  $y_s$ , and  $z_s$ . The  $x_s$ -axis is displayed horizontally,  $y_s$  is vertical, and  $z_s$  is perpendicular to the computer screen. (In a right-handed coordinate system, positive  $x_s$  is to the right, positive  $y_s$ is toward the top of the screen, and positive  $z_s$  is toward the viewer.) As the model is moved and rotated, the screen coordinates are continually updated.

The simplest molecular displays are stick models with lines connecting atoms, and atoms simply represented by vertices where lines meet (look ahead to Fig. 11.2, p. 274 for examples of various model displays or *renderings*). It is easy to imagine a program that simply plots a point at each position ( $x_s$ ,  $y_s$ ,  $z_s$ ) and connects the points with lines according to a set of instructions about connectivity of atoms in amino acids.

But the computer screen is two-dimensional. How does the computer plot in three dimensions? It doesn't; it plots a projection of the three-dimensional stick model. Mathematically, projecting the object into two dimensions involves some simple trigonometry, but graphically, projecting is even simpler. The program plots points on the screen at positions  $(x_s, y_s, 0)$ ; in other words, the program does not employ the z<sub>s</sub> coordinate in producing the display. This produces a projection of the molecule on the  $x_s y_s$  plane of the screen coordinate system, which is the computer screen itself (see Fig. 11.1). You can imagine this projection process as analogous to casting a shadow of the molecule on the screen by holding it behind the screen and lighting it from behind. Another analogy is the image of a tree projected onto the ground by its leaves when they fall during a cold windless period. The program may use the z-coordinate to produce shading (by scaling color intensity in proportion to  $z_s$ , thus making foreground objects brighter than background objects) or perspective (by scaling the  $x_s$  and  $y_s$  coordinates in proportion to  $z_s$ , thus making foreground objects larger than background objects), or to allow foreground objects to cover background objects (allowing objects with larger  $z_s$ , to overwrite those with smaller  $z_s$ ), and thus to make the display look three-dimensional.

#### 11.2.2 Into three dimensions: Basic modeling operations

The complexity of a protein model makes it essential to display it as a threedimensional object and move it around (or move our viewpoint around within it). The first step in seeing the model in three dimensions is rotating it, which gives many three-dimensional cues and greatly improves our perception of it. Rotating the model to a new orientation entails calculating new coordinates for all the

270



**Figure 11.1** Geometry of projection. (*a*) Model viewed from off to the side of the screen coordinate system. Each atom is located by screen coordinates  $x_s$ ,  $y_s$ , and  $z_s$ . (*b*) Model projected onto graphics screen. Each atom is displayed at position ( $x_s$ ,  $y_s$ , 0), producing a projection of the model onto the  $x_s y_s$ -plane, which is the plane of the graphics screen.

atoms and redisplaying by plotting on the screen according to the new  $(x_s, y_s)$  coordinates.

The arithmetic of rotation is fairly simple. Consider rotating the model by  $\theta$  degrees around the  $x_s$  (horizontal) axis. It can be shown that this rotation transforms the coordinates of point p,  $[x_s(p), y_s(p), z_s(p)]$ , to new coordinates  $[x'_s(p), y'_s(p), z'_s(p)]$  according to these equations:

$$\begin{aligned} x'_{s}(p) &= x_{s}(p) \\ y'_{s}(p) &= [y_{s}(p) \cdot \cos \theta] - [z_{s}(p) \cdot \sin \theta] \\ z'_{s}(p) &= [y_{s}(p) \cdot \sin \theta] + [z_{s}(p) \cdot \cos \theta]. \end{aligned}$$
(11.1)

Notice that rotating the model about the  $x_s$ -axis does not alter the  $x_s$  coordinates,  $[x'_s(p) = x_s(p)]$ , but does change  $y_s$  and  $z_s$ . A similar set of equations provides for rotation about the  $y_s$ - or  $z_s$ -axis. If we instruct the computer to rotate the model around the  $x_s$ -axis by  $\theta$  degrees, it responds by converting all coordinates  $(x_s, y_s, z_s)$  to new coordinates  $(x'_s, y'_s, z'_s)$  and plotting all points on the screen at the new positions  $(x'_s, y'_s, 0)$ . Graphics programs allow so-called real-time rotation in which the model appears to rotate continuously. This requires fast computer increments  $\theta$  by a small amount, recalculates coordinates of all displayed atoms, using equations such as 11.1, and redraws the screen image. Fast repetition of this process gives the appearance of continuous rotation.

#### 11.2.3 Three-dimensional display and perception

Complicated models become even more comprehensible if seen as threedimensional (3-D) objects even when not in motion. So graphics display programs provide some kind of full-time 3-D display. This entails producing two images like the stereo pairs used in this book, and presenting one image to the left eye and the other to the right eye, which is the function of a stereo viewer. The right-hand view is just like the left-hand view, except that it is rotated about 5° about its  $y_s$ axis (clockwise, as viewed from above). Molecular modeling programs display the two images of a stereo pair side by side on the screen, for viewing in the same manner as printed pairs. (To learn how to see the three-dimensional image using side-by-side pairs, see Appendix, p. 293, or the CMCC home page.)

With proper hardware, modeling programs can produce full-screen 3-D objects. The technique entails flashing full-size left and right views alternately at high speed. The viewer wears special glasses with liquid-crystal lenses that alternate rapidly between opaque and transparent. When the left-eye view is on the screen, the left lens is transparent and the right lens is opaque. When the screen switches to the right-eye view, the right lens becomes transparent and the left becomes opaque. In this manner, the left eye sees only the left view, and the right eye sees only the right view. The alternation is fast, so switching is undetectable, and a full-screen 3-D model appears to hang suspended before the viewer. Several viewers, each wearing the glasses, can see 3-D at the same time, without having to align their eyes with the screen, as is required for stereo pairs.

Both types of stereo presentation mimic the appearance of objects to our two eyes, which produce images on the retina of objects seen from two slightly different viewpoints. The two images are rotated about a vertical axis located at the current focal point of the eyes. From the difference between the two images, called *binocular disparity*, we obtain information about the relative depth of objects in our field of view. For viewers with normal vision, two pictures with a 5° binocular disparity, each presented to the proper eye, gives a convincing three-dimensional image.

Unfortunately, a small but significant percentage of people cannot obtain depth information from binocular disparity alone. In any group of 20 students, there is

a good chance that one or more will not be able to see a three-dimensional image in printed stereo pairs, even with a viewer. In the world of real objects, we decode depth not only by binocular disparity but also by relative motion, size differences produced by perspective, overlap of foreground objects over background objects, effects of lighting, and other means. Molecular modeling programs provide depth information to a wider range of viewers by providing depth cues in the form of shading, perspective, and movement, as described in the next section.

#### 11.2.4 Types of graphical models

Modeling and graphics programs provide many ways to render or represent a molecular model. Figure 11.2 shows several different representations of three strands of beta sheet from human thioredoxin (PDB 1ert). The simplest representation is the wireframe model (a). Although other renderings have greater visual impact, the wireframe model is without question the most useful when you are exploring a model in detail, for several reasons. First, you can see the model, yet see through it at the same time. Foreground and background parts of the model do not block your view of the part of the model on which you center your attention. Second, in wireframe models, you can see atom positions, bond angles, and torsional angles clearly. Third, wireframe models are the simplest and hence the fastest for your computer to draw. This means that, as you rotate or zoom a model, the atom positions are recalculated and the images redrawn faster, so that the model moves more rapidly and smoothly on the screen. Rapid and smooth motion give the model a tangible quality that improves your ability to grasp and understand it. Combined with colorings that represent element types, or other properties like B-factors or solvent accessibility (Sec. 11.3.8, p. 282), wireframe models can be very informative. There is one disadvantage: wireframe models do not contain obvious depth cues; therefore, they are best viewed in stereo.

Other renderings, some of which are also illustrated in Fig. 11.2, have their strengths and weaknesses. All of them make vivid illustrations in textbooks and on web pages, when conveying fine structural details is not necessary. It takes programs longer to draw these more complex models, so they might not move as smoothly on the screen. Ball-and-stick models are similar to wireframe models, but they provide better depth cues, because they are shaded and appear as solid objects that obscure objects behind them. Aside from providing depth cues, the balls primarily produce clutter, but they look nice. Figure 11.2*b* shows a ball-and-stick model of the beta strands shown in (a).

Various cartoon renderings, like the ribbon diagram in Fig. 11.2*c*, show main chains as parallel strands or solid ribbons, perhaps with arrows on strands of pleated sheet to show their direction, and barrel-like alpha helices. These are some of the most vivid renderings and are excellent for giving an overview of protein secondary and tertiary structure. Combined images, for example, with most of the molecule in cartoon form and an important binding site in ball and stick, make nice printed illustrations that guide the reader from the forest of the whole molecule to the trees of functionally important structural details. For example, this particular image also shows more clearly than either (a) or (b) that successive side chains lie



**Figure 11.2**  $\triangleright$  Common types of computer graphics models (stereo), all showing the same three strands of pleated-sheet structure from human thioredoxin (PDB 1ert). (*a*) Wireframe; (*b*) ball and stick; (*c*) ribbon backbone with ball-and-stick side chains; (*d*) space filling. Image: DeepView/POV-Ray.

on opposite sides of a pleated sheet. A disadvantage of more complex images is that rotation on a personal computer can be choppy and slow.

The same atoms are shown again in Fig. 11.2*d* as a space-filling model, which gives a more realistic impression of the density of a model, and of the extent of its surface. The surface of a space-filling model can be either the van der Waals (shown) or the solvent-accessible surface, and can be colored according to various properties, like surface charge. Space-filling models are excellent for examining the details of atomic contacts within and between models. Note, for example, in this figure, how side chains nestle together on the top of this sheet. Space-filling models of an enzyme and its inhibitor appear to fit each other like hand and glove. A disadvantage: you cannot make out exact atom positions and bond angles in space-filling models.

As the renderings of Fig. 11.2 show, there is no single model that shows everything you might like to see. Graphics programs allow you to mix model types and switch quickly among them, to highlight whatever you are trying to see or illustrate. Browsing a model with computer graphics is far superior to viewing even a large number of static images.

# 11.3 Touring a molecular modeling program

As you might infer from a brief description of the computing that underlies molecular graphics, the computer must be fast. But today's personal computers are up to the task of dealing with quite large and complex graphics. The current generation of personal computers can allow you to conduct highly satisfying and informative expeditions into the hearts of the most complex macromolecular models.

The basic operations of projecting and rotating a screen image of the molecular model make the foundation of all molecular graphics programs. Upon these operations are built many tools for manipulating the display. These tools give viewers the feeling of actively exploring a concrete model. Now I will discuss tools commonly found in modeling programs.

Although this is a general discussion of modeling tools, I will use as my example an excellent molecular viewing and modeling program that you can obtain free of charge in Linux, Macintosh, Silicon Graphics, or Windows versions. My example is DeepView, originally and still often called Swiss-PdbViewer, an easy-to-learn yet very powerful molecular viewing and modeling tool. Figure 11.3 provides a picture of a personal computer screen during use of DeepView, and the caption gives a description of its main sets of functions. If you currently own no modeling program and want to learn to use a program that will allow you start easily, yet grow painlessly into very sophisticated modeling and structural analysis, Deep-View is an excellent choice. In fact, it compares favorably with costly commercial programs. As of this writing, the number of free macromolecular graphics programs is growing constantly. But I still know of no single program that does as many things, and as many of them very well, as DeepView. To learn how to download this program from the World Wide Web and to find a complete, self-guiding tutorial for modeling beginners, see the CMCC home page. Whatever the program you decide you use, search its documentation or the World Wide Web for tutorials on how to get started with it. There is no faster way to learn molecular exploration and analysis than by a hands-on tutorial. The CMCC home page also contains links to help you find other molecular graphics programs and tutorials for learning how to use them.

#### 11.3.1 Importing and exporting coordinate files

As indicated earlier, image display and rotation requires rapid computing. Reading coordinates from a text file like a PDB file is slow because every letter or number in the file must be translated into binary or machine language for the computer's internal processing. For this reason, graphics programs work with a machine-language version of the coordinate file, which they can read and recompute faster than a text file. So the first step in exploring a model is usually converting the PDB file to binary form. With almost all modeling programs, this operation is invisible to the user.

Often users produce revealing views of the model and wish to use the coordinates in other programs, such as energy calculations or printing for publication. For this purpose, molecular modeling programs include routines for writing coordinate files in standard formats like PDB, using the current binary model coordinates as input.

Screen shot of DeepView in use on a Macintosh OS X computer. Con-Figure 11.3 ► trols for manipulating the model are at the top of the main graphics window, which can be expanded to fill the screen. The Control Panel lists residues in the model and allows selection of residues for display, coloring, labeling, and surface displays. The Align window shows residue sequence also, including alignment of multiple models if present. The Ramachandran Plot window shows main-chain torsional angles for residues currently selected (red in the Control Panel, purple in the Align window). Dragging dots on the Rama plot changes torsion angles interactively. The Layer Infos window (not shown) allows control of display features for multiple models (layers) in any combination. In the graphics window is a stereo display of the heme region of cytochrome b5 (PDB 1cyo). Selected hydrogen bonds are shown in green, and measured distances in yellow. A TORSION operation is in progress (darkened button, top of graphics window). The side-chain conformation of phenylalanine-58 is being changed. Clashes between the side chain and other atoms are shown in pink. The user can display and read the PDB file of the currently active model by clicking the document icon at the upper right of the graphics windows. Clicking an ATOM line in the PDB file display centers the graphics model on that residue and reduces the display to the residue and its nearest neighbors.



#### 11.3.2 Loading and saving models

The coordinate files produced by graphics programs can be loaded and saved just like any other files. As the model is manipulated or altered, coordinates are updated. At any point, the user can save the current model, or replace it with one saved earlier, just as you can do with a word-processor document. To facilitate studying intermolecular interactions, comparing families of models, or scanning through NMR ensembles, modeling programs can handle several or many models at once, each as separate layers that you can view separately or superimpose. The models currently in memory can be viewed and manipulated individually or together. You can save them with their current relative orientations so that you can resume a complicated project later.

Although almost all macromolecular model coordinates are available in PDB format, modeling programs provide for loading and saving files in a variety of formats. These might include Cambridge Structure Database format, used by the largest database for small molecules, or others from among the dizzying list of small-molecule formats used by organic and inorganic chemists. Other special file formats include those used by energy minimization, dynamics simulation, and refinement programs, many of which are now also available for personal computers. Full-featured modeling programs like DeepView usually provide for exporting coordinates for these powerful programs, shipping them over networks for completion of complex computing tasks, retrieving the resulting model files, and converting them back to PDB format for viewing. As personal computers get faster, these transfers will become unnecessary. Useful adjuncts to graphics programs are online tools for interconverting various types of files. For example, the Dundee PRODRG2 Server allows users to paste or upload coordinates in almost any known format, and then to convert the file to any other format. See the CMCC home page for links to this and other file-conversion web sites.

#### 11.3.3 Viewing models

Standard viewing commands allow users to rotate the model around screen  $x_s$ ,  $y_s$ , and  $z_s$  axes and move (translate) the model for centering on areas of interest. Viewers magnify or zoom the model by moving it along the  $z_s$ -axis toward the viewer, or by using a command that magnifies the image without changing coordinates by simply narrowing the viewing angle. "Clip" or "slab" commands simplify the display by eliminating foreground and background, producing a thin slab of displayed atoms. With most programs, including DeepView, all of these operations are driven by a mouse or other point-and-drag device (trackball, for example), after selecting the desired operation by buttons (top of graphics window in Fig. 11.3) or key commands.

Figure 11.4 shows two views of the small protein cytochrome b5 (PDB 1cyo), an electron-transport protein containing an iron-heme prosthetic group (shown edgeon with gray iron at center). In (a), you are looking into the heme binding pocket, with much of the protein in the background. Even in stereo, the background clutter makes it difficult to see the heme environment clearly. In (b), the background is

278



**Figure 11.4**  $\blacktriangleright$  Heme region of cytochrome *b5* (stereo, PDB 1cyo). (*a*) View without clipping; (*b*) same view after "slab" command to eliminate all except contents of a 12-Å slab in the  $z_s$  direction. Image: DeepView.

removed using a "slab" command to give a clearer view of the heme and protein groups above and below it. The program clips by displaying only those atoms whose current  $z_s$ -coordinates lie within a specified range (12.5 Å, in this case), which is chosen visually by sliding front and rear clipping planes together until unwanted background and foreground atoms disappear.

Centering commands allow the user to select an atom to be made the center of the display. Upon selection by pointing to the atom and clicking a mouse, or by naming the atom, the program moves the model so as to center the desired atom on the screen and within the viewing slab, and also to make it the center of subsequent rotations. For example, DeepView provides a very handy centering feature. Simply pressing the "=" key centers the part of the model currently on display, and resizes it to fit the screen. It is not unusual for novice viewers to accidentally move the model completely out of view and be unable to find it. Nothing is more disconcerting to a beginner than completely losing sight of the model. When the model disappears, it may be off to the side of the display, above or below the display, or still centered, but outside the slab defined by clipping planes. In any event, automatic recentering can often help viewers find the model and regain their bearings. As a last resort, there is usually a reset command, which brings model and clipping planes back to starting positions. The viewer pays a price for resetting, losing the sometimes considerable work of finding a particularly clear orientation for the model, centering on an area of interest, and clipping away obscuring parts of the structure.

Viewing commands usually also include selection of stereo or mono viewing and offer various forms of *depth cueing* to improve depth perception, either by mimicking the effects of perspective (front of model larger than rear), shading (front of model brighter than rear), or rocking the model back and forth by a few degrees of rotation about the  $y_s$ -axis.

#### 11.3.4 Editing and labeling the display

A display of every atom in a protein is often forbidding and incomprehensible. Viewers are interested in some particular aspect of the structure, such as the active site or the path of the backbone chain, and may want to delete irrelevant parts of the model from the display. Display commands allow viewers to turn atoms on and off. Atoms not on display continue to be affected by rotation and translation, so they are in their proper places when redisplayed. Viewers might eliminate specific atoms by pointing to them and clicking a mouse, or they might eliminate whole blocks of sequence by entering residue numbers. They may display only alpha carbons to show the folding of the protein backbone (see Fig. 3.4, p. 36), or only the backbone and certain side chains to pinpoint specific types of interactions.

Editing requires knowledge of how atoms are named in the coordinate file, which is often, but not always, the same as PDB atom labels (Sec. 7.7, p. 173). Thus viewers can produce an alpha-carbon-only model by limiting the display to atoms labeled CA. Each program has its own language for naming atoms, residues (by number or residue name), distinct chains in the model (like the  $\alpha$  and  $\beta$  chains of hemoglobin), and distinct models. Viewers must master this language in order to edit displays efficiently. In DeepView, editing the view is greatly simplified by a Control Panel (Fig. 11.3, right side) that provides a scrollable list of all residues in the PDB file. A menu at the top allows you to switch to other PDB files if you are currently viewing more than one model. You can select, display, label, and color residues, and add surface displays as well, all with mouse operations.

Most programs provide powerful selection tools to allow you to pick specific parts of the model for displaying, labeling, or coloring. For example, DeepView provides, in its Control Panel, the means to select main chain, side chains, or both; individual residues; multiple residues (contiguous or not); individual elements of secondary structure; and individual chains. In addition, DeepView menu commands allow selecting residues by type (for example, select all histidines, or select water, or select heteromeric group), property (select acidic residues), or secondary structure (select helices, beta strands, or coils); surface or buried residues; residues with *cis*-peptide bonds; problem residues in a model, like those with  $\Phi$  and  $\Psi$  values outside the range of allowed values (Sec. 8.2.1, p. 181), or residues making clashes; neighbors of the current selection (select within 4.5 Å of current selection); and residues in contact with another chain, to name just a few.

Even with an edited model, it is still easy for viewers to lose their bearings. Label commands attach labels to specified atoms, signifying element, residue number,

or name. Labels like the one for PHE-58 in Fig. 11.3 float with the atom during subsequent viewing, making it easy to find landmarks in the model.

#### 11.3.5 Coloring

Although you may think that color is merely an attractive luxury, adding color to model displays makes them dramatically more understandable. Most programs allow atoms to be colored manually, by selecting a part of the model and then choosing a color from a color wheel or palette. Additional color commands allow the use of color to identify elements or specific residues, emphasize structural elements, or display properties. For example, DeepView provides, among others, commands for coloring the currently selected residues by CPK color (as in Fig. 11.3), residue type, *B*-factor (sometimes called coloring by temperature), secondary structure (different colors for helix, sheet, and turns), secondary structure in sequence (blue for first helix or beta strand, red for last one, and colors of the visible spectrum for each secondary structural element in between), chain (a different color for each monomer in an oligomeric protein), layer (different color for each model currently on display), solvent accessibility, threading or force-field energies (Sec. 10.3.4, p. 263), and various kinds of model problems.

WARNING: Coloring by *B*-factor uses whatever information is in the *B*-factor column of the PDB file. To interpret the colors that result, it is crucial for you to know what kind of model you are viewing, because this information tells you different things about different kinds of models. Coloring a good crystallographic model by *B*-factor reveals relative uncertainty in atom positions due to static or dynamic disorder (Sec. 8.2.3, p. 185). Coloring an averaged NMR model by "B-factor" actually reveals relative rms deviation of atom positions from the average positions of the corresponding atoms in the ensemble of NMR models (Sec. 10.2.5, p. 257). And coloring a homology model by "B-factor" distinguishes parts of the model according to how much support exists for the model in the form of X-ray or NMR structural information (Sec. 10.3.4, p. 263).

Combined with selecting tools, color commands can be powerful tools simply for finding features of interest. If you are interested in cysteines in a protein, you can select all cysteines, choose a vivid color for them, and immediately be able to find them in the forest of a large protein.

#### 11.3.6 Measuring

Measurements are necessary in identifying interactions within and between molecules. In fact, noncovalent interactions like hydrogen bonds are defined by the presence of certain atoms at specified distances and bond angles from each other. In Fig. 11.3, p. 276, yellow dotted lines connect two carbon atoms, a valine methyl and a heme methyl. The distance between the atoms is displayed and is approximately the distance expected for carbon atoms involved in a hydrophobic contact.

Modeling programs allow display of distances, bond angles, and dihedral angles between bonded and nonbonded atoms. These measurements float on dotted lines within the model (just like labels) and are often active; that is, they are continually updated as the model is changed, as described in the next section.

#### 11.3.7 Exploring structural change

Modeling programs allow the viewer to explore the effects of various changes in the model, including conformational rotation, change in bond length or angle, and movement of fragments or separate chains. Used along with active measurements, these tools allow viewers to see whether side chains can move to new positions without colliding with other atoms, or to examine the range of possible movements of a side chain. In Fig. 11.3, a change of torsion angles is in progress for phenylalanine-58. In its current conformation, the side chain is in collision with the heme, as shown by a pink "clash" line.

Like rotation and translation, changes of model conformation, bond angles, or bond lengths are reflected by changes in the coordinate file. The changes are tentative at first, while users explore various alterations of the model. After making changes, users have a choice of saving the changes, removing the changes, or resetting in order to explore again from the original starting point. In Fig. 11.3, notice that the torsions button is darkened, showing that the operation is in progress, and that the user will have an opportunity to keep or discard the changes being made.

#### 11.3.8 Exploring the molecular surface

Stick models of the type shown in Fig. 11.2*a*, p. 274 are the simplest and fastest type of model to compute and display because they represent the molecule with the smallest possible number of lines drawn on the screen. Stick models are relatively open, so the viewer can see through the outer regions of a complex molecule into the interior or into the interface between models of interacting molecules. But when the viewer wants to explore atomic contacts, a model of the molecular surface is indispensable.

Published structure papers often contain impressive space-filling computer images of molecules, with simulated lighting and realistic shadows and reflections. These images require the computer to draw hundreds of thousands of multicolored lines, and so most computers cannot redraw such images fast enough for continuous movements. Some of these views require from seconds to hours to draw just once. Although such views show contacts between atomic surfaces, they are not practical for exploring the model interactively. They are used primarily as snapshots of particularly revealing views.

How then can you study the surface interactively? The most common compromise is called a dotted surface (Fig. 11.5), in which the program displays dots evenly spaced over the surface of the molecule. This image reveals the surface without obscuring the atoms within and can be redrawn rapidly as the viewer manipulates the model. Several types of surfaces can be computed, each with its own potential uses. One type is the van der Waals surface (Fig. 11.5*a*), in which all dots lie at the van der Waals radius from the nearest atom, the same as the surface



**Figure 11.5**  $\blacktriangleright$  Dotted surface displays of heme in cytochrome *b*5 (stereo, PDB 1cyo). (*a*) Van der Waals surface enclosing the entire heme. (*b*) Solvent-accessible surface of heme, showing only the portion of heme surface that is exposed to surrounding solvent. Most of heme is buried within the protein. Image: DeepView.

of space-filling models. This represents the surface of contact between nonbonded atoms. Any model manipulations in which van der Waals surfaces penetrate each other are sterically forbidden. Van der Waals surfaces show packing of structural elements with each other, but the display is complicated because all internal and external atomic surfaces are shown. Elimination of the internal surfaces produces the *molecular surface*, which is simpler, but still shows some surfaces that are not accessible to even the smallest molecules, like water.

A very useful surface display is the *solvent-accessible surface*, which shows all parts of the molecule that can be reached by solvent (usually water) molecules.

#### Chapter 11 Tools for Studying Macromolecules

This display omits all internal atomic surfaces, including crevices that are open to the outside of the model, but too small for solvent to enter. Some modeling programs, like DeepView, contain routines for calculating this surface, whereas others can take as input the results of surface calculations from widely available programs. Calculating the solvent-accessible surface entails simulating the movement of a sphere, called a *probe*, having the diameter of a solvent molecule over the entire model surface, and computing positions of evenly spaced dots wherever model and solvent come into contact. Figure 11.5*b*, p. 283, shows the solventaccessible surface for the heme group in cytochrome b5. It is clear from this view that most of the heme is buried within the protein and not accessible to the solvent.

Carrying out the same simulation that produces solvent-accessible surface displays, but locating the dots at the center of the probe molecule, produces the *extended surface* of the model. This display is useful for studying intermolecular contacts. If the user brings two models together—one with extended surface displayed, the other as a simple stick model—the points of intermolecular contact are where the extended surface of one model touches the atom centers of the second model.

The default color of the displayed surface is usually the same as the color selected for the underlying atoms. In Fig. 11.5*b*, for example, two oxygens of one of the heme carboxyl groups produces the large red bulge of accessible surface (the other carboxyl is hydrogen bonded to serine-64 and is much less accessible to solvent). Alternatively, color can reflect surface charge (commonly, blue for positive, red for negative, with lighter colors for partial charges) or surface polarity (contrasting colors for hydrophobic and hydrophilic regions). These displays facilitate finding regions of the model to which ligands of specified chemical properties might bind.

The surface of proteins carry many charged and polar functional groups that confer electrostatic potentials between the protein and its surroundings. Graphics programs can display such potentials in two ways, with isopotential surfaces or with potentials mapped onto the molecular surface. Figure 11.6 presents one subunit of the dimeric enzyme acetylcholine esterase (PDB 2ack) with isopotential surfaces (calculated by DeepView) showing its electrostatic properties. The red surface is at a constant (that's what iso means) negative potential and the blue is the same level of positive potential. The height of the isopotential surface above the molecular surface is greatest where the surface potential is highest. As you can see, the surface directly above the substrate binding site (that is, toward the top of the figure), which is marked by a white competitive inhibitor, has a very high negative potential. The substrate, acetylcholine, is positively charged, and thus is strongly attracted by the negatively charged residues that surround the active site region. Attraction of the active site region for the substrate gives this enzyme a very high on-rate for substrate binding. The potential surface gives an almost tactile feeling for the attraction the enzyme would have for a passing substrate molecule; a positively charged substrate would be repelled by the back side of the enzyme, but even a grazing encounter with the face of the enzyme should result in substrate binding.

In Fig. 11.6*b*, the model of (*a*) has been rotated 90° around  $x_s$  toward the viewer, and the electrostatic potential has been mapped onto the molecular surface of the



**Figure 11.6** Electrostatic potentials and surfaces (stereo). (*a*) Electrostatic isopotential surface of one subunit of acetylcholine esterase (PDB 2ack). (*b*) Electrostatic potential mapped onto molecular surface, looking down on the top of the view in (*a*). A competitive inhibitor bound at the active site is visible in both views. Image: DeepView/POV-Ray.

enzyme. You can see the competitive inhibitor (space-filling model) through an opening in the enzyme surface. The shade of the surface reflects the surface potential, with darkest red for highest negative potential, corresponding to locations above which an isopotential surface would be farthest away. The view shows that the substrate is deeply buried inside the enzyme, but the electrostatic potential explains how it quickly finds its way to this secluded site.

Combining models with surfaces can make dramatic illustrations of molecular properties. In Fig. 11.7, the full dimeric structure is shown as an alpha-carbon model, with the isopotential surface for the dimer rendered as a smooth,

Chapter 11 Tools for Studying Macromolecules



**Figure 11.7** ► Acetylcholine esterase dimer (stereo). Alpha-carbon stick model of protein, space filling model of inhibitor, and isopotential surface. Subrates enter at lower left and upper right. Image: OpenGL rendering in DeepView.

transparent surface. This view helps to explain how acetylcholine esterase can bind substrate on almost every encounter—its on-rate for binding is very close to the diffusion rate, the rate at which two species at 1.0 M concentrations are expected to collide with each other in solution.

#### 11.3.9 Exploring intermolecular interactions: Multiple models

Formulating proposed mechanisms of protein action requires investigating how proteins interact with ligands of all kinds, including other proteins. Molecular modeling programs allow the user to display and manipulate several models, either individually or together. With DeepView and many other modeling programs, the number of models is limited only by computer memory and speed. Tools for this purpose usually allow all of the same operations as the viewing tools but permit selection of models affected by the operations. In *docking experiments* (a term taken from satellite docking in the space program), one model can be held still while another is moved into possible positions for intermolecular interaction. Labeling, measurement, and surface tools are used simultaneously during docking to ensure that the proposed interactions are chemically realistic. Some programs include computational docking, in which the computer searches for optimal interaction, usually from a user-specified starting point. Such calculations are quite slow, and usually done by stand-alone programs that produce output coordinate files for

viewing on programs like DeepView. Many graphics programs can prepare files for external docking programs. See the CMCC home page for links to programs and servers for protein-protein and protein-ligand docking.

#### 11.3.10 Displaying crystal packing

Many molecular modeling programs include the capacity to display models of the entire unit cell. All the program needs as input is a set of coordinates for one molecule, the unit-cell dimensions, and a list of equivalent positions for the crystal space group. The user can display one cell or clusters (say,  $2 \times 2 \times 2$ ) of cells. The resulting images, particularly when teamed with surface displays, reveal crystal-packing interactions, allowing the user to see which parts of the crystallographic model might be altered by packing, and might thus be different from the solution structure. For examples of crystal-packing displays, see Fig. 4.18, p. 70, and Fig. 4.19, p. 71. Unit-cell tools usually allow the user to turn equivalent positions on and off individually, making them useful for teaching the topics of equivalent positions and symmetry. For example, DeepView allows the user to create new models by specifying symmetry operations or selecting them from a list or from symmetry lines in PDB files.

#### 11.3.11 Building models from scratch

In addition to taking coordinate files as inputs, modeling programs allow the user to build peptides to specification and to change amino-acid residues within a model. To build new models, users select amino acids from a palate or menu and direct the program to link the residues into chains. Users can specify conformation for the backbone by entering backbone angles  $\Phi$  and  $\Psi$ , by selecting a common secondary structure or by using the tools described earlier for exploring structural change. Model-builder tools are excellent for making illustrations of common structural elements like helices, sheet, and turns. DeepView allows you to start building a model by simply providing a text file of the desired sequence in one-letter abbreviations; it imports the sequence and models it as an alpha helix for compactness. Then you can select parts of the model and either select their secondary structural type or set their Ramachandran angles. Alternatively, you can display a Ramachandran diagram for the model and change torsional angles by dragging residue symbols to new locations on the diagram, watching the model change as you go-a great tool for teaching the meaning of the  $\Phi$  and  $\Psi$  angles. Similar tools are used to replace one or more side chains in a model with side chains of different amino acids, and thus explore the local structural effects of mutation.

#### 11.3.12 Scripts and macros: Automating routine structure analysis

In many areas of structure analysis, you find yourself repeating sequences of operations. For example, in comparing a number of related proteins, you might routinely open a model file, align it with a previously loaded model (your reference model), and color main-chain atoms by rms deviation from the reference. You might even want to capture the results by constructing a table of sequence alignments and rms deviations by residue. Such analysis of dozens of models quickly becomes drudgery. Sounds like a job for a programmer, and some modeling programs make such programming easy by providing for writing and executing command scripts, or for recording and playing back sequences of operations (sometimes called macros). For example, DeepView incorporates a scripting language that is similar to Perl or C++ . You can write scripts with any text editor, and execute them with menu commands within the program. Sample scripts are included with program documentation. Scripting and macros allow you to turn your modeling program loose on large repetitive tasks, and then retrieve the results in easily usable form.

## 11.4 Other tools for studying structure

It is beyond the scope of this little book to cover all the tools available for studying protein structure. I will conclude by listing and briefly describing additional tools, especially ones used in conjunction with modeling on graphics systems.

#### 11.4.1 Tools for structure analysis and validation

In addition to molecular graphics, a complete package of tools for studying protein structure includes many accessory programs for routine structure analysis and judging model quality. The chores executed by such programs include the following:

- Calculating Φ and Ψ angles and using the results to elements of secondary structure as well as to display a Ramachandran diagram, which is useful in finding model errors during structure refinement (crystallography or NMR) or homology modeling. As I mentioned earlier, DeepView has a unique interactive Ramachandran plot window that allows the user to change main-chain conformational angles in the model (Fig. 11.3, p. 276).
- Using distance and angle criteria to search for hydrogen bonds, salt links, and hydrophobic contacts, and producing a list of such interactions. Deep-View calculates and displays hydrogen bonds according to user specifications of distance and angle. A very powerful pair of menu commands in Deep-View allows you to show hydrogen bonds only from selected residues or hetero groups and then to show only residues with visible hydrogen bonds. In these two operations, you can eliminate everything from the view except, for example, a prosthetic group and its hydrogen-bonding neighbors.
- Comparing homologous structures by least-squares superposition of one protein backbone on another. The result is a new coordinate set for one model that best superimposes it on the other model. I used such a tool in DeepView to compare the X-ray and NMR structures of thioredoxin in Fig. 3.4, p. 36 (alignment was instantaneous). DeepView provides superposition tools combined



Figure 11.8 ► X-ray and NMR models of human thioredoxin (stereo, PDB lert and 3trx), aligned by least-squares superposition of corresponding alpha carbons. X-ray model is gray. NMR model is colored by deviation from X-ray model. Greatest deviations are red, and smallest deviations are blue.

with sequence comparison to improve the structural alignment between proteins that are only moderately homologous. In addition, you can color a protein by the rms deviation of its atoms from those of the reference protein on which it was superimposed, giving a vivid picture of areas where the structures are alike and different. For example, in Fig. 11.8, full backbone models of the X-ray and NMR structures of human thioredoxin are superimposed. The X-ray model is gray, and the NMR model is colored by rms deviations of corresponding residues from the X-ray model. Residues for which the two models deviate least are colored blue. Those exhibiting the greatest deviations are red. Residues with intermediate deviations are assigned spectral colors between blue and red. Even if the X-ray model were not shown, it would be easy to see that the most serious disagreement between the two models lies in the surface loop at the bottom, and that the two models agree best in the interior residues.

- ► Building additional subunits using symmetry operations, either to complete a functional unit (Sec. 8.2.4, p. 187) or to examine crystal packing. DeepView allows you to build additional subunits by selecting symmetry operations from a palette, clicking on symmetry operations listed in a PDB file, or typing in the components of a transformation matrix.
- Carrying out homology modeling. As mentioned in Chapter 10, DeepView is a full homology modeling program. Using DeepView, a web browser, and electronic mail, you can obtain template files from SWISS-PROT, align, average, and thread your sequence onto the template, build loops or select them from a database, find and fix clashes, submit modeling projects to SWISS-MODEL, and retrieve them to examine the results or apply other

#### Chapter 11 Tools for Studying Macromolecules



**Figure 11.9** ► Model and portion of electron-density map of bovine Rieske ironsulfur protein (stereo, PDB 1rie). The map is contoured around selected residues only. Image: DeepView/POV-Ray.

modeling tools. GROMOS energy minimization is included in DeepView. Figure 10.11, p. 264, shows a homology model started in DeepView and completed at SWISS-MODEL. The homology model is shown as ribbon colored according to model *B*-factors (see Sec. 10.3.4, p. 265; also see the warning about *B*-factor coloring on p. 281). Two templates used in the modeling are shown as black and gray alpha-carbon displays. (For more information about these proteins, see the legend to Fig. 10.11.)

Displaying an electron-density map and adjusting the models to improve its fit to the map (see Fig. 11.9). DeepView can display maps of several types (CCP4, X-PLOR, DN6). Outside of full crystallographic computation packages, I am aware of no programs currently available for computation of maps from structure factors on personal computers, but I am sure this will change. Structure factors are available for some of the models in the PDB, and can usually be obtained from the depositors of a model. As described in Sec. 8.2.5, p. 189, the Electron Density Server at Uppsala University provides electron-density maps for most models for which structure factors are deposited in the Protein Data Bank.

#### 11.4.2 Tools for modeling protein action

The crystallographic model is used as a starting point for further improvement of the model by energy minimization and for simulations of molecular motion. Additional insight into molecular function can be obtained by calculating charge densities and bond properties by molecular orbital theory. For small molecules, some of these calculations can be done "on the fly" as part of modeling. For the more complex computations, and for larger molecules, such calculations are done outside the graphics program, often as separate tasks on computers whose forte is number crunching rather than graphics. DeepView can write files for export to several widely used programs for energy minimization and molecular dynamics. But as personal computers get faster, these operations will no longer require transfer to specialized machines.

### 11.5 Final note

Making computer images and printed pictures of molecular models endows them with the concreteness of everyday objects. While exploring models, viewers can easily forget the difficult and indirect manner by which they are obtained. I wrote this book in hopes of providing an intellectually satisfying understanding of the origin of molecular models, especially those obtained from single-crystal X-ray crystallography. I also hope to encourage readers to explore the many models now available, but to approach them with full awareness of what is known and what is unknown about the molecules under study. Just as good literature depicts characters and situations in a manner that is "true to life," a sound model depicts a molecule in a manner that is true to the data from which it was derived. But just as real life is more multifarious than the events, settings, and characters of literature, not all aspects of molecular truth (or even of crystallographic, NMR, or modeling truth) are reflected in the colorful model floating before us on the computer screen. Users of models must probe more deeply into the esoterica of structure determination to know just where the graphics depiction is not faithful to the data. The user must probe further still—by using validation tools and by reading wider literature on the molecule-to know whether the model is faithful to other evidence about structure and action. The conversation between structural models and evidence on all sides continually improves models as depictions of molecules.

Perhaps I have stimulated your interest in crystallography itself, and have made you wonder if you might jump in and determine the structure of that interesting protein you are studying. I am happy that I can encourage you by reiterating that crystallography, though still one of structural biology's more challenging callings, is faster and easier than ever before. Screening for crystal growth conditions does not require expensive equipment or chemicals. Most research universities have at least enough crystallographic instrumentation to allow you to assess the diffracting power of your crystals, or perhaps even to collect preliminary data. You can take promising crystals to synchrotron sources for data collection. Finally, you can do the computation on reasonably modern personal computers. Look for a local crystallography research group to help you get started.

Of course, if you pursue crystallography, there are many more details to learn. As your next step toward a truly rigorous understanding of the method, I suggest *Introduction to Macromolecular Crystallography* by Alexander McPherson (Wiley-Liss, 2002); *Crystal Structure Analysis for Chemists and Biologists: Methods in Stereochemical Analysis* by Jenny Glusker, Mitchell Lewis, and Miriam Rossi (John Wiley & Sons, 1994); *X-ray Structure Determination: A Practical Guide*, 2nd edition, by George H. Stout and Lyle H. Jensen (John Wiley and

#### Chapter 11 Tools for Studying Macromolecules

Sons, Inc., 1989); and *Practical Protein Crystallography*, by Duncan E. McRee (Academic Press, Inc., 1993). On the World Wide Web, I recommend *Crystallography 101* by Bernhard Rupp. Finally, if you are really serious about being a crystallographer, and don't mind getting very little sleep for about 16 days, apply for admission to the course *X-Ray Methods in Structural Biology*, taught each fall at Cold Spring Harbor Laboratory. You will find links to this and other crystallography resources at the CMCC home page, www.usm.maine.edu/~rhodes/CMCC.

292

## **Viewing Stereo Images**

To see three-dimensional images using divergent stereo pairs in this book, use a stereo viewer such as item #D8-GEO8570, Carolina Biological Supply Company (1-800-334-5551). Or better, you can view stereo pairs without a viewer by training yourself to look at the left image with your left eye and the right image with your right eye (called divergent viewing). This is a very useful skill for structural biologists, and is neither as difficult nor as strange as it sounds. (According to my ophthamologist, it is not harmful to your eyes, and may in fact be good exercise for eye muscles.)

Try it with this stereo image:



Figure A.1 > Divergent stereo pair of interior, Dom St. Stephan, Passau, Deutschland.

#### Appendix Viewing Stereo Images

Try putting your nose on the page between the two views. With both eyes open, you will see the two images superimposed, but out of focus, because they are too close to your eyes. Slowly move the paper away from your face, trying to keep the images superimposed until you can focus on them. (Keep the line between image centers parallel to the line between your eyes.) When you can focus, you will see three images. The middle one should exhibit convincing depth. Try to ignore the flat images on either side. This process becomes easier and more comfortable with practice. If you have difficulty, try it with a very simple image, such as Fig. 4.18, p. 70.

For additional help with viewing stereo images in books or on computers, click on *Stereo Viewing* at the CMCC home page, http://www.usm.maine.edu/~rhodes/CMCC.

294

## Index

#### **Numerics**

3D-Crunch project, 264

#### A

A-DNA diffraction patterns, 217 absences in diffraction patterns, 73, 105-107 absorption, nuclear spin and, 240-241 absorption edge, 128 absorption of X rays, 74. See also detectors, X-ray anomalous scattering, 128-136 direct phasing methods, 135-136 extracting phase, 130-132 hand problem, 135 measurable effects, 128-130 multiwavelength (MAD), 133-134 amorphous materials, diffraction by, 219-222 amplitude of electronic-density map, 150 amplitude, wave, 21 angles within cells, 49-50 anomalous scattering (anomalous dispersion), 128-136 direct phasing methods, 135-136 extracting phase, 130-132 hand problem, 135 measurable effects, 128-130 multiwavelength (MAD), 133-134

applying prior knowledge to models, 152, 154, 253 area detectors, 78 arrays, diffraction by, 16, 18 assessing model quality. See quality of models assigning NMR model resonances, 251-252. See also NMR models asymmetric units, 68 asymmetric vs. functional units, 188 ATOM lines (PDB file), 176 atom occupancy, 161, 185 atomic coordinates refinement of, 147-148 resolution and precision, 183-185 atomic coordinates entries, 174-176 atomic coordinates, modeling from, 269-270 atomic plane indices, 50-55 atomic structure factors, 98 AUTHOR lines (PDB file), 175 automating routine structure analysis, 287 - 288average spacing of reciprocal-lattice points, 86

#### В

B-DNA diffraction patterns, 217 B-factors (homology models), 267, 281 B-factors (NMR models), 257–258 back-transforms, 96

#### ball-and-stick models, 273 Bayesian refinement models, 164-168 beam stops, 77 Bessel functions, 216 bias, minimizing, 154-156 Bayesian inference, 164 bichromatic X-ray sources, 74 binocular disparity, 272, 293-294 Biotech Validation Suite, 191 BLAST program, 261 body-centered (internal) lattices, 66 bootstrapping density modification, 151-153 error removal, 147, 171 overview of, 146-149 refinement of atomic coordinates, 147-148 Bragg's law, 55-57, 111 polychromatic rays, 231 in reciprocal space, 60-64 Bravais lattices, 67 brightness of reflections, 16-17, 63 electron density as function of, 115-117 scaling and postrefinement, 85-86 broken crystals, recognizing, 45 building models. See map fitting

#### С

<sup>13</sup>C-editing, 250 caged substrates, 234 calculated intensities, 105 cameras, X-ray, 80-85 capillary mounting, 45 cartoon renderings, 273 CCDs (charge-coupled devices), 78-80 cells. See unit cells center of symmetry, 88 characteristic lines, 73 charge-coupled devices, 78-80 charge detection, 229 chemical shift, 239 chemically reasonable models, 172-173 CISPEP lines (PDB file), 176 CMCC web page, 4-5 cocrystallization, 40 cofactors, to grow protein crystals, 45 collimators, 77

coloring models (DeepView program), 281 comparative protein modeling, 260-263 complex numbers, 92-93 in two dimensions, 112 complex objects, diffraction by, 17 complex vectors, structure factors as, 112-115 COMPND lines (PDB file), 175 computer models of molecules, 269-275 conditions for crystal growth, 41-46 CONECT lines (PDB file), 176 confidence factor (homology models), 267 conformation, determining for NMR models, 252-257 conformationally reasonable models, 173, 183 constraints, defined, 163 contour levels, 94 contour maps. See maps of electron density convergence to final model, 168-173 coordinate systems, 19-20, 50 depositing with PDB, 173-177 importing and exporting, 276-277 refinement of atomic coordinates, 147 - 148two-dimensional computer images from, 269-270 correlation spectroscopy (COSY), 248 COSY (correlation spectroscopy), 248 coupling, nuclear, 239-240 Crick, Francis, 217 cross-validation, 172 cryo-electronic microscopy, 227-231 cryocrystallography, 45-46 CRYST lines (PDB file), 176 crystal packing, 188, 287 crystallation techniques. See growing crystals crystalline fibers, 213 crystallographic asymmetric units, 188 crystallographic refinement. See structure refinement Crystallography Made Crystal Clear web page, 4-5 crystallography models. See models crystallography papers, reading, 192-208

#### Index

#### Index

crystals, in general, 10–13. *See also* protein crystals growing, 11–13, 37–46 optimal conditions for, 41–46 mounting for data collection, 46–47 quality of, judging, 46 cubic cells, 50

#### D

data collection, 13-15, 73-89 cameras, 80-85 detectors, 77-80 intensity scaling and postrefinement, 85-86 mounting crystals for, 46-47 symmetry and strategy, 88-89 unit-cell dimension determination, 86-88 X-ray sources, 73-77 data mining, 177 database of proteins. See PDB databases of homology models, 263-265 DBREF lines (PDB file), 175 de Broglie equation, 222 dead time, 78 DeepView program, 275-288 density modification, 151-153 density, unexplained, 187 depositing coordinates with PDB, 174-177 depth cuing, 280 derivative crystals, 37 growing, 40-41 selenomet derivatives, 41 detectors, X-ray, 77-80 difference Patterson functions, 125 diffracted X-ray reflections, 13, 49 as Fourier terms, 101-104 intensities of, 16-17 measurable, number of, 64-65 measured, electronic density from, 28-30 phases of. See phase sphere of reflection, 63-64 symmetry elements and, 71-73 wave descriptions of, 24-26 diffraction, 15-19, 211-235 by amorphous materials (scattering), 219-222

calculating electronic density, 91-107 data collection, 13-15, 73-89 cameras, 80-85 detectors, 77-80 intensity scaling and postrefinement, 85-86 mounting crystals for, 46-47 symmetry and strategy, 88-89 unit-cell dimension determination, 86-88 X-ray sources, 73-77 electron diffraction, 227-231 fiber diffraction, 211-219 geometric principles of, 49-73 judging crystal quality, 46 Laue diffraction, 231-235 neutron diffraction, 222-227 systematic absences in patterns, 73, 105-107 diffraction patterns, 30 diffractometry, 82 digestion of proteins, 45 dimensions of unit cells, 65, 86-88 direct phasing methods, 135-136 disorder, atomic, 185-186 regions of disorder, 187 disordered water, 34-35 dispersive differences, 134 displaying electron-density maps, 150-151 distance trap, 267 distortion from crystal packing, 188 divergent viewing, 272, 293-294 DNA diffraction patterns, 217 dynamic range, 78

#### Ε

edges, cell, 49–50. *See also* unit cells editing display (DeepView program), 280 EDS (Electron Density Server), 190 electron cloud imaging, 11, 26 electron crystallography, 227 electron density calculating from diffraction data, 91–107 determining from measured reflections, 28–30

#### 297
electron density (continued) determining from structure factors, 27 - 28as Fourier sum, 99-100 as function of intensities and phases, 115-117 maps of, 26-27 amplitude of, 150 developing model from, 153-159 final models, 168-173 first maps, 149-153 improving. See bootstrapping structure refinement, 159-168 unexplained, 187 Electron Density Server (EDS), 190 electron diffraction, 227-231 electron-potential maps, 229 EM (electron microscopy), 227-231 EMD (Electron Microscope Database), 230-231 enantiomeric arrangements, 127-128 END lines (PDB file), 176 energy-minimized averaged models, 255 energy refinement, 163 equivalent positions (symmetry), 69 error removal (filtering), 147, 171 etched crystals, 40 evaluating model quality. See quality of models Ewald sphere, 63-64 examining electron-density maps, 150-151 ExPASy tool, 263 EXPDTA lines (PDB file), 175 experimental models, 7 exporting coordinate files, 276-277 expression vectors, 41 extended surfaces, 284

### F

face-centered lattices, 67 FastA program, 261 fiber diffraction, 211–219 FID (free-induction decay), 243. *See also* NMR models figure of merit, 156 figure of merit (phase), 152 film, X-ray-sensitive, 78 filtering, 147, 171 final model, converging to, 168-173 first maps, 149-153 Fo - Fc maps, 155-156, 169 disorder, indications of, 185 focusing mirrors, 77 FORMUL lines (PDB file), 176 Fourier analysis, defined, 96, 97 Fourier series, 23-24, 97 Fourier sums (summations), 26, 92-97 electron density as, 99-100 reflections as terms in, 101-104 structure factors as, 98-99 Fourier synthesis, 24, 97 Fourier terms, defined, 92 Fourier transforms, 28, 96-97 of FID (free-induction decay), 243.247 of helices, 214 testing models with, 217-218 frames of data, 83 free-induction decay (FID), 243. See also NMR models free log-likelihood gain, 172 free R-factor, 172 freezing crystals, 45-46 frequency, wave, 22 Friedel pairs, 114 anomalous scattering, 131-133 Friedel's law, 88, 114 FTs. See Fourier transforms functional vs. asymmetric units, 188 functions of proteins, structure and, 35

### G

generalized unit cells, 49–50 geometric principles of diffraction, 49–73 Bragg's law, 55–57, 111 polychromatic rays, 231 in reciprocal space, 60–64 global minimum (structure refinement), 162 goniometer heads, 80 goniostats, 80 graphical models, types of, 273–275 graphing three-dimensional functions, 94 GROMOS program, 263, 290

#### Index

growing crystals, 11–13, 37–46 optimal conditions for, 41–46 gyration, radius of, 221

### Η

hand problem anomalous scattering, 135 isomorphous replacement, 128 handing-drop method, 38-40 Harker sections (Harker planes), 127 harvest buffers, 46 header, PDB files, 175-176, 189 heavy-atom derivatives, 37. See also isomorphous replacement hand problem, solving, 135 preparing, 117-119 heavy-atom method. See isomorphous replacement heavy atoms, neutron scattering with, 225 helices, in fibrous materials, 213-214 HELIX lines (PDB file), 176 HET and HETNAM lines (PDB file), 176 high-angle reflections, 101 higher-dimensional NMR spectra, 249-251 homology models, 237, 259-267 basic principles, 260-263 databases, 263-265 judging quality of, 265-267 hydrogen bonds, 31

# I

image enhancement with cryo-EM, 229–230 image plate detectors, 78 imaginary numbers. *See entries at* complex imaging microscopic objects, 8 importing coordinate files, 276–277 improving maps and models. *See* bootstrapping impurity of samples, 41 indices (coordinates), 19, 50–55 atomic planes, 52–55 lattice indices, 50–52 integrity, structural, 31–33 intensity of reflections, 16-17, 63 electron density as function of, 115-117 scaling and postrefinement, 85-86 intermolecular interactions, exploring (DeepView program), 286 internal lattices, 66 internal symmetry of unit cells, 65-73 data collection strategies, 88-89 functional-unit, 188 noncrystallographic (NCS), 153 Patterson map searches, 127 International Tables for X-Ray Crystallography, 69 interplanar spacing. See Bragg's law inverse Fourier transforms, 96 isomorphism, 37, 118 isomorphous phasing models, 137-139 isomorphous replacement, 117-128 isotropic vibration, 185 iterative improvement of maps and models. See bootstrapping

# J

JRNL lines (PDB file), 175 judging crystal quality, 46 judging model quality. *See* quality of models justifiable restrictions on model, 152, 154, 253

# Κ

KEYWS lines (PDB file), 175 knowledge-based modeling, 260 knowledge (prior), applying, 152, 154, 253 known properties, improving model with, 152, 154, 253

### L

labeling display (DeepView program), 280 lack of closure (isomorphous replacement), 124 lattice indices, 50–52 lattices, 10 real, 16 Bragg's law, 55–57

lattices (continued) reciprocal, 16 average point spacing, 86 how to construct, 57-60 limiting sphere, 64 types of, 50, 66-67 Laue diffraction, 231-235 Laue groups, 89 layer lines, 213-216 least-squares refinement methods, 159 - 161lenses, 8 ligands, 37 cocrystallization, 40 to grow protein crystals, 45 limited digestion of proteins, 45 limiting sphere, 64 LINK lines (PDB file), 176 liquid nitrogen. See cryocrystallography loading models (DeepView program), 278 local minima (structure refinement), 162-163 location of phasing model, 139-141 log-likelihood gain, 168, 172 longitudinal relaxation, 242 loops, homology models, 262 low-angle reflections, 101 low-angle scattering, 221, 227 low-resolution models, 187 Luzzati plots, 183-184

### Μ

macros with DeepView program, 287-288 MAD (multiwavelength anomalous dispersion) phasing, 133-134 manual model building, 168-169 map fitting, 154, 156-159 convergence to final model, 168-173 with DeepView program, 287 manual, 168-169 maps of electron density, 26-27 amplitude of, 150 developing model from, 153-159 final models, 168-173 first maps, 149-153 improving. See bootstrapping structure refinement, 159-168 unexplained, 187

maps of electron potential, 229 maps of nucleon density, 223 MASTER lines (PDB file), 176 mathematics of crystallography, 20 - 30maximum-density lines, 154 measurements, taking (DeepView program), 281 membrane-associated proteins, 45, 211 microscopic imaging, 8 Miller indices, 60 minima, local (structure refinement), 162-163 minimaps, 151 minimizing bias from models, 154-156 MIR (multiple isomorphous replacement), 123 mirror plane of symmetry, 67 model validation tools, 189-192, 288-290 Model Validation tutorial, 190 modeling after known proteins. See molecular replacement modeling tools, 269-292 computer modeling, 269-275 DeepView program (example), 275 - 288model validation tools, 189-192, 288-290 models computer models, 269-275 improving electronic-density maps density modification, 151-153 error removal, 147, 171 overview of, 146-149 refinement of atomic coordinates, 147-148 judging quality of. See quality of models modifying to obtain phase information, 147 non-crystallographic, 237-268 homology models, 237, 259-267 NMR models, 238-259 other theoretical models, 267-268 obtaining and improving, 30, 145, 153-159 convergence to final model, 168-173

judging convergence, 171-173 minimizing bias, 154-156 model building. See map fitting regularization, 157, 171 structure refinement, 159-168 structure refinement parameters, 161-162, 189-192 phasing models. See molecular replacement reading crystallography papers, 192-208 sharing, 173-177 types of, 7 understanding, in general, 1-3 molecular dynamics, 163, 186 molecular envelopes, 151 molecular models. See models molecular models, types of, 7 molecular replacement, 136-143 imposing prior knowledge, 152, 154, 253 isomorphous phasing models, 137-139 location and orientation searches, 139 - 143nonisomorphous phasing models, 139 molecular surface exploration (DeepView program), 282-286 molecules, obtaining images of, 9 MolProbity Web Service, 191 monochromatic X-ray sources, 74 monoclinic cells, 50, 59-60 monodispersity, 45 mosaicity (mosaic spread), 33 mother liquor, 33 cryoprotected, 46 with ligand, 40-41 motion, molecular, 163 mounting crystals for data collection, 46-47 multidimensional NMR models, 249-251 multiple forms for crystals, 33 multiple isomorphous replacement (MIR), 123 multiwavelength anomalous diffraction (MAD) phasing, 133-134 multiwavelength X-ray beams, 231-235

multiwire area detectors, 78

### Ν

<sup>15</sup>N-editing, 250 native crystals, 37 NCS (noncrystallographic symmetry), 153 negative indices, 54 neutron diffraction, 222-227 NMR models, 238-259 basic principles, 239-251 conformation determination, 252-257 judging quality of, 257-259 PDB files for, 257 resonances, assigning, 251-252 NOE (nuclear Overhauser effect), 248, 252 conformation determination, 252-253 NOESY (NMR model), 248 noncrystalline fibers, 213 noncrystallographic symmetry (NCS), 153 nonisomorphous phasing models, 137-139 NSLS (National Synchrotron Light Source), 75-77 nuclear Overhauser effect, 248 nuclear spin, 239-240 nucleation, 38 nucleon-density maps, 223

# 0

objective lens, 8 obtaining and improving models, 30, 145, 153-159 convergence to final model, 168-173 judging convergence, 171-173 map fitting (model building), 154, 156-159 convergence to final model, 168-173 with DeepView program, 287 manual, 168-169 minimizing bias, 154-156 regularization, 157, 171 structure refinement, 159-168 parameters for, 161-162, 189-192 obtaining phases, 104-105, 109-143 anomalous scattering, 128-136 isomorphous replacement, 117-128

obtaining phases (continued) molecular replacement, 136-143 two-dimensional representation of structure factors, 112-117 occupancy, atom, 161, 185 omit maps, 156 one-dimensional NMR models, 243, 251 one-dimensional waves, 92-93 online databases, 44-45 online model validation tools, 189-192 ordered water, 34-35, 169 orientation of phasing model, 139-141 ORIG lines (PDB file), 176 orthorhombic cells, 50 lattice indices of (example), 51-52 oscillation cameras, 82-84 overall map amplitude, 150 Overhauser effect (nuclear), 248

### Ρ

packing effects, 36-37 packing scores, defined, 190 papers (crystallography), reading, 192-208 parameters for structure refinement, 161-162, 189-192 partial molecular models, building, 154 particle storage rings, 75-76 Patterson functions and maps, 124-128 phasing model location and orientation, 139-143 spherically averaged, 220-221 PDB (Protein Data Bank), 7, 174-177 file headers, 175-176, 189 NMR models, 257 quality of publications in, 180 PEG (polyethylene glycol), 11-12 periodic functions, 21-24 phage display, 166-167 phase, 98, 101 allusions to "Phase" poem, 23, 229, 270 electron density as function of, 115-117 filtering, 171 improving. See bootstrapping mathematics of, 22-23 obtaining, 104-105, 109-143

anomalous scattering, 128-136 isomorphous replacement, 117-128 molecular replacement, 136-143 two-dimensional representation of structure factors, 112-117 quality of (figure of merit), 152 phase angle, 113 phase extension, 152 phase maps, obtaining, 154-156 "Phase" (poem), xxiv allusions to, 23, 229, 270 phase probabilities, 124 phase problem, defined, 116 phasing models. See molecular replacement phasing power, 118 pitch, helix, 213 plane indices, 50-55 plane-polarized light, observing, 45 pleated sheets, 159 point of inversion, 88 polychromatic X-ray beams, 231-235 polydispersity, 45 polyethylene glycol (PEG), 11-12 position of wave, 23. See also phase positive vs. negative indices, 54 posterior possibilities (Bayesian inference), 164, 166 postrefinement of X-ray intensities, 85-86 precession photography, 84 precipitating proteins, 11-12, 38. See also growing crystals precision of atomic positions, judging, 183-185 primitive lattices, 66 prior knowledge restrictions on model, 152, 154 NMR models, 253 prior possibilities (Bayesian inference), 164, 165 probability distributions (Bayesian inference), 165 PROCHECK program, 191 protein action modeling tools, 290 protein conformation, determining (NMR models), 252-257 protein crystallography, sketch of, 9-10

protein crystals, 31-47 growing. See growing crystals mounting crystals for data collection, 46-47 properties of, 31-34 quality of, judging, 46 related, as phasing models. See molecular replacement solution vs. crystal structures, 34-37 Protein Data Bank. See PDB protein function, structure and, 35 protein precipitation, 11-12, 38. See also growing crystals protein structure. See entries at structure publishing models, 173-177 purity of samples, 41

# Q

quality of crystals, judging, 46 quality of models. See also models, obtaining and improving assessing convergence, 171-173 homology models, 265-267 how to judge, 180-192 atomic positions, 183-185 miscellaneous limitations of, 187-189 model validation tools, 189-192 structural parameters, 161-162, 181-183, 189-192 vibration and disorder, 185-187 NMR models, 257-259 reading crystallography papers, 192-208 quantity necessary for growing crystals, 41

### R

R-factor, 142 *R*-factor, 172, 181 Luzzati plots and, 184 radial Patterson functions, 220–221 radius of convergence, 162–163 radius of gyration, 221 Ramachandran diagrams, interpreting, 181–183 reading crystallography papers, 192–208 real lattices, 16 Bragg's law, 55-57 real-space refinement algorithms, 147-148 manual model building, 168-169 phase filtering, 171 reasonableness of models, 172-173, 183 reciprocal angstroms, 20 reciprocal lattices, 16 average point spacing, 86 how to construct, 57-60 limiting sphere, 64 reciprocal space, 19-20, 57-64 Bragg's law, 60-64 reciprocal-space refinement algorithms, 147-148, 159 phase filtering, 171 refinement of atomic coordinates, 147 - 148refinement of map structure, 159-168 parameters for, 161-162, 189-192 refinement target, defined, 166 reflections (X-ray), 13, 49 as Fourier terms, 101-104 intensities of, 16-17 measurable, number of, 64-65 measured, electronic density from, 28 - 30phases of. See phase sphere of reflection, 63-64 symmetry elements and, 71-73 wave descriptions of, 24-26 regularization, 157, 171 related proteins. See molecular replacement relaxation, nuclear, 242-243 REMARK lines (PDB file), 175 removing error from models, 147, 171 rendering computer models. See computer models residual index (R-factor), 172, 181 Luzzati plots and, 184 resolution limit, 65 resolution of atomic positions, judging, 183-185, 187 resonances, assigning in NMR models, 251-252. See also NMR models response-surface procedure, 41-42

restrained minimized average structures, 255 restraints, defined, 163 REVDAT lines (PDB file), 175 reversibility of Fourier transforms, 96 ridge lines, drawing, 154 rise-per-residue (helix), 215 rms deviations (NMR models), 257–258 rotamers, 183 rotating anode tubes, 75 rotation, 68 rotation functions, 142 self-rotation, 153 rotation of computer models, 271–272 rotation/oscillation method, 82–84

# S

safety, X-ray, 77 sampled diffraction patterns, 17 saving models (DeepView program), 278 SCALE lines (PDB file), 176 scaling X-ray data, 85-86 scattering, diffraction by, 219-222 scattering factor, 98 scattering length, 223-224 scintillation counters, 77-78 screen coordinates (computer models), 270 screw axis, 68 scripts with DeepView program, 287-288 seed crystals, 40 selenomet derivatives, 41 self-rotation functions, 153. See also rotation functions SEQADV lines (PDB file), 175 SEQRES lines (PDB file), 175 Shake-and-Bake direct phasing, 136 sharing models, 173-177 SHEET lines (PDB file), 176 side chains, homology models, 263 side chains, identifying, 158 sigma-A weighted maps, 156, 169 disorder, indications of, 185 simple objects, diffraction by, 15-17 simulated annealing, 163 single-crystal C-ray crystallography, 7 single isomorphous replacement (SIR), 118, 131

SIR (single isomorphous replacement), 118 with anomalous scattering (SIRAS), 131 SITE lines (PDB file), 176 size, crystal, 31-33 skeletonizing maps, 154 slow precipitation, 11-12, 38. See also growing crystals solution structure vs. crystalline structure, 35-37 solvent-accessible surfaces, 283-284 solvent leveling/flattening/flipping, 151-153. See also density modification SOURCE lines (PDB file), 175 sources, X-ray, 73-77 space-filling models, 275 space groups, 66-67 spallation, 223 spectroscopy. See NMR models sphere of reflection, 63-64 spheres, diffraction by, 15 spherically averaged diffraction patterns, 220-221 spin-lattice relaxation, 242-243 spin, nuclear, 239-240 spin-spin relaxation, 243 square wave, 24 step function, 24 stereo images, viewing, 272, 293-294 stereochemically reasonable models, 173 strategies for data collection, 88-89 structural change exploration (DeepView program), 282 structural integrity, 31-33 structural parameters, 161-162, 189-192 structure, solution vs. crystalline, 35-37 structure analysis tools, 288-290 structure-factor patterns, 30 structure factors, 24-26 computing from model, 104-105 electron density from, 27-28 as Fourier sums, 98-99, 102-103 Friedel pairs, 88, 114 two-dimensional representation, 112-117

structure files, PDB, 174–176

#### Index

structure refinement, 159-168 parameters for, 161-162, 189-192 summations, Fourier, 26, 92-97 electron density as, 99-100 reflections as terms in, 101-104 structure factors as, 98-99 SWISS-MODEL Repository, 237, 260, 264.268 Swiss-PdbViewer program. See DeepView program symmetry averaging, 153 symmetry elements, 67 symmetry of unit cells, 65-73 data collection strategies, 88-89 functional-unit, 188 noncrystallographic (NCS), 153 Patterson map searches, 127 symmetry-related atoms, 105-106 synchrotrons, 75-77, 222 systematic absences, 73, 105-107

### Т

target core (homology models), 261-262 temperature factors, 161, 169 vibration measurements, 185-186 temperature of mounted crystals, 45-46 templates for homology models, 260-261 TER lines (PDB file), 176 tetragonal cells, 50 theoretical models, 7. See also models thermal motion, 185-186 three-dimensional arrays, diffraction by, 18 three-dimensional computer models, 270-273 three-dimensional NMR models, 250 three-dimensional waves, 94-95 time-domain signals, 243 time-resolved crystallography, 231-235 TITLE lines (PDB file), 175 tools for studying macromolecules, 269-292 computer modeling, 269-275 DeepView program (example), 275-288 model validation tools, 189-192, 288-290 touring molecular modeling programs, 275-288

translation, 68 translation searches, 139-141 transmission electron microscopes, 227 transverse relaxation time, 243 triclinic cells, 50 triplet relationship (direct phasing), 135-136 truncated Fourier series, 102 tubes, X-ray, 75-76 TURN lines (PDB file), 176 twinned crystals, 31, 45 two-dimensional arrays, diffraction by, 18 two-dimensional computer models, 269-270 two-dimensional NMR models, 244-248 twofold rotation axis, defined, 67

# U

uncertainty factor (homology models), 267 unexplained density, 187 unit cells, 10, 49-50 asymmetric vs. functional, 188 dimensions of, 65, 86-88 heavy atoms in, locating, 124-128 reciprocal, 59 symmetry of, 65-73 data collection strategies, 88-89 functional-unit, 188 noncrystallographic (NCS), 153 Patterson map searches, 127 unit translation, 68 unreasonableness of models, 172-173, 183 upper-level planes, 59 user's guide to crystallographic models, 179-210 judging model quality, 180-192 atomic positions, 183-185 miscellaneous limitations of, 187-189 model validation tools, 189-192 structural parameters, 161-162, 181-183, 189-192 vibration and disorder, 185-187 reading crystallography papers, 192-208

V validation tools, 189–192, 288–290 validation tools, online, 189–192 van der Walls forces, 253–254 vapor diffusion, 38 variable-wavelength X rays. *See* synchrotrons vectors, complex, 112–115 vibration, atomic, 185–186 viewing electronic-density maps, 151 viewing models (DeepView program), 278–280

#### W

water content of crystals, 34–35 wave equations, 21–23 one-dimensional waves, 92–93 three-dimensional waves, 94–95 wavelength, 22 wavelengths of X-ray emissions, 73 weighted maps, 156 wigglers, 75 wireframe models, 273

# Х

X-ray absorption, 74. *See also* detectors, X-ray anomalous scattering, 128–136 direct phasing methods, 135–136 extracting phase, 130–132

hand problem, 135 measurable effects, 128-130 multiwavelength (MAD), 133-134 X-ray analysis, 9 X-ray data collection, 13-15, 73-89 cameras, 80-85 detectors, 77-80 intensity scaling and postrefinement, 85-86 mounting crystals for, 46-47 symmetry and strategy, 88-89 unit-cell dimension determination, 86-88 X-ray sources, 73-77 X-ray reflections, 13, 49 as Fourier terms, 101-104 intensities of, 16-17 measurable, number of, 64-65 measured, electronic density from, 28 - 30phases of. See phase sphere of reflection, 63-64 symmetry elements and, 71-73 wave descriptions of, 24-26 X-ray scattering lengths, 224 X-ray tubes, 75-76

# Ζ

zero-level planes, 59 zinc fingers, 1