

# Protein Structure Prediction

# METHODS IN MOLECULAR BIOLOGY™

*John M. Walker, SERIES EDITOR*

419. **Post-Transcriptional Gene Regulation**, edited by Jeffrey Wilusz, 2008
418. **Avidin–Biotin Interactions: Methods and Applications**, edited by Robert J. McMahon, 2008
417. **Tissue Engineering, Second Edition**, edited by Hamsjörg Hauser and Martin Fussenegger, 2007
416. **Gene Essentiality: Protocols and Bioinformatics**, edited by Andrei L. Osterman, 2008
415. **Innate Immunity**, edited by Jonathan Ewbank and Eric Vivier, 2007
414. **Apoptosis in Cancer: Methods and Protocols**, edited by Gil Mor and Ayesha Alvero, 2008
413. **Protein Structure Prediction, Second Edition**, edited by Mohammed Zaki and Chris Bystroff, 2008
412. **Neutrophil Methods and Protocols**, edited by Mark T. Quinn, Frank R. DeLeo, and Gary M. Bokoch, 2007
411. **Reporter Genes for Mammalian Systems**, edited by Don Anson, 2007
410. **Environmental Genomics**, edited by Cristofre C. Martin, 2007
409. **Immunoinformatics: Predicting Immunogenicity In Silico**, edited by Darren R. Flower, 2007
408. **Gene Function Analysis**, edited by Michael Ochs, 2007
407. **Stem Cell Assays**, edited by Mohan C. Vemuri, 2007
406. **Plant Bioinformatics: Methods and Protocols**, edited by David Edwards, 2007
405. **Telomerase Inhibition: Strategies and Protocols**, edited by Lucy Andrews and Trygve O. Tollefsbol, 2007
404. **Topics in Biostatistics**, edited by Walter T. Ambrosius, 2007
403. **Patch-Clamp Methods and Protocols**, edited by Peter Molnar and James J. Hickman, 2007
402. **PCR Primer Design**, edited by Anton Yuryev, 2007
401. **Neuroinformatics**, edited by Chiquito J. Crasto, 2007
400. **Methods in Lipid Membranes**, edited by Alex Dopico, 2007
399. **Neuroprotection Methods and Protocols**, edited by Tiziana Borsello, 2007
398. **Lipid Rafts**, edited by Thomas J. McIntosh, 2007
397. **Hedgehog Signaling Protocols**, edited by Jamila I. Horabin, 2007
396. **Comparative Genomics, Volume 2**, edited by Nicholas H. Bergman, 2007
395. **Comparative Genomics, Volume 1**, edited by Nicholas H. Bergman, 2007
394. **Salmonella: Methods and Protocols**, edited by Heide Schatten and Abe Eisenstark, 2007
393. **Plant Secondary Metabolites**, edited by Harinder P. S. Makkar, P. Siddhuraju, and Klaus Becker, 2007
392. **Molecular Motors: Methods and Protocols**, edited by Ann O. Sperry, 2007
391. **MRSA Protocols**, edited by Yinduo Ji, 2007
390. **Protein Targeting Protocols, Second Edition**, edited by Mark van der Giezen, 2007
389. **Pichia Protocols, Second Edition**, edited by James M. Cregg, 2007
388. **Baculovirus and Insect Cell Expression Protocols, Second Edition**, edited by David W. Murhammer, 2007
387. **Serial Analysis of Gene Expression (SAGE): Digital Gene Expression Profiling**, edited by Kare Lehmann Nielsen, 2007
386. **Peptide Characterization and Application Protocols**, edited by Gregg B. Fields, 2007
385. **Microchip-Based Assay Systems: Methods and Applications**, edited by Pierre N. Floriano, 2007
384. **Capillary Electrophoresis: Methods and Protocols**, edited by Philippe Schmitt-Kopplin, 2007
383. **Cancer Genomics and Proteomics: Methods and Protocols**, edited by Paul B. Fisher, 2007
382. **Microarrays, Second Edition: Volume 2, Applications and Data Analysis**, edited by Jang B. Rampil, 2007
381. **Microarrays, Second Edition: Volume 1, Synthesis Methods**, edited by Jang B. Rampil, 2007
380. **Immunological Tolerance: Methods and Protocols**, edited by Paul J. Fairchild, 2007
379. **Glycoviroplogy Protocols**, edited by Richard J. Sugrue, 2007
378. **Monoclonal Antibodies: Methods and Protocols**, edited by Maher Albitar, 2007
377. **Microarray Data Analysis: Methods and Applications**, edited by Michael J. Korenberg, 2007
376. **Linkage Disequilibrium and Association Mapping: Analysis and Application**, edited by Andrew R. Collins, 2007
375. **In Vitro Transcription and Translation Protocols: Second Edition**, edited by Guido Grandi, 2007
374. **Quantum Dots**: edited by Marcel Bruchez and Charles Z. Hotz, 2007
373. **Pyrosequencing® Protocols**, edited by Sharon Marsh, 2007
372. **Mitochondria: Practical Protocols**, edited by Dario Leister and Johannes Herrmann, 2007
371. **Biological Aging: Methods and Protocols**, edited by Trygve O. Tollefsbol, 2007
370. **Adhesion Protein Protocols, Second Edition**, edited by Amanda S. Coutts, 2007
369. **Electron Microscopy: Methods and Protocols, Second Edition**, edited by John Kuo, 2007
368. **Cryopreservation and Freeze-Drying Protocols, Second Edition**, edited by John G. Day and Glyn Stacey, 2007
367. **Mass Spectrometry Data Analysis in Proteomics**, edited by Rune Matthiesen, 2007
366. **Cardiac Gene Expression: Methods and Protocols**, edited by Jun Zhang and Gregg Rokosh, 2007
365. **Protein Phosphatase Protocols**: edited by Greg Moorhead, 2007

METHODS IN MOLECULAR BIOLOGY™

# Protein Structure Prediction

*Second Edition*

Edited by

**Mohammed J. Zaki  
and  
Christopher Bystroff**

*Rensselaer Polytechnic Institute, Troy, New York, USA*

HUMANA PRESS  TOTOWA, NEW JERSEY

©2008 Humana Press Inc.  
999 Riverview Drive, Suite 208  
Totowa, New Jersey 07512

**www.humanapress.com**

All rights reserved. No part of this book may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, electronic, mechanical, photocopying, microfilming, recording, or otherwise without written permission from the Publisher. Methods in Molecular Biology™ is a trademark of The Humana Press Inc.

All papers, comments, opinions, conclusions, or recommendations are those of the author(s), and do not necessarily reflect the views of the publisher.

This publication is printed on acid-free paper.   
ANSI Z39.48-1984 (American Standards Institute) Permanence of Paper for Printed Library Materials

Production Editor: Rhukey Hussain  
Cover design by Karen Schulz

For additional copies, pricing for bulk purchases, and/or information about other Humana titles, contact Humana at the above address or at any of the following numbers: Tel.: 973-256-1699; Fax: 973-256-8341; E-mail: humana@humanapress.com; or visit our Website: www.humanapress.com

**Photocopy Authorization Policy:**

Authorization to photocopy items for internal or personal use, or the internal or personal use of specific clients, is granted by Humana Press Inc., provided that the base fee of US \$30 copy is paid directly to the Copyright Clearance Center at 222 Rosewood Drive, Danvers, MA 01923. For those organizations that have been granted a photocopy license from the CCC, a separate system of payment has been arranged and is acceptable to Humana Press Inc. The fee code for users of the Transactional Reporting Service is: [978-1-58829-752-5/08 \$30].

Printed in the United States of America. 10 9 8 7 6 5 4 3 2 1  
eISBN 978-1-59745-574-9

Library of Congress Control Number: 2007933144

---

# Preface

For 40 years we have known the essential ingredients for protein folding – an amino acid sequence and water. But the problem of predicting the three-dimensional structure from its sequence has eluded computational biologists even in the age of supercomputers and high throughput structural genomics. Will we ever solve the “protein folding problem”, or will we simply settle for a solution to the “protein prediction problem”? This book covers elements of both the data-driven comparative modeling approach to structure prediction and also recent attempts to simulate folding using explicit or simplified models. Despite the unsolved mystery of how a protein folds, advances are being made in predicting the interactions of proteins with other molecules, such as small ligands, nucleic acids, or other proteins. Also, rapidly advancing are the methods for solving the inverse folding problem, the problem of finding a sequence to fit a structure. This book focuses on the various computational methods for prediction, their successes, and their limitations, from the perspective of their most well-known practitioners. An overview of the chapters in this volume is given below.

## Overview of Protein Structure Prediction

In the first chapter, entitled “A Historical Perspective of Template-Based Protein Structure Prediction,” Jun-tao Guo, Kyle Ellrott, and Ying Xu give a comprehensive, as well as historical, account of protein structure prediction. They touch upon methods spanning threading, fold recognition, homology modeling, ab initio methods, and their hybrids. They also discuss recent progress in the worldwide blind structure prediction evaluation experiments like CASP and its cousin for automated servers, CAFASP.

In the second chapter “The Assessment of Methods for Protein Structure Prediction,” Anna Tramontano, Domenico Cozzetto, Alejandro Giorgetti, and Domenico Raimondo take a critical look at extant methods for protein structure prediction and assess how well they perform. They focus on automatic assessment methods as well as the CASP challenges and discuss their limitations and trade-offs.

## Template-Based Methods

In the third chapter “Aligning Sequences to Structures,” Liam J. McGuffin discusses the current approaches to template-based fold prediction. The goal here is to align new protein sequences to library of known/template folds. Liam also shows a step-by-step guide to template alignment.

In the fourth chapter “Protein Structure Prediction Using Threading,” Jinbo Xu, Feng Jiao, and Libo Yu discuss approaches for protein threading. After setting up the general requirements for protein structure prediction by threading, they specifically focus on their successful new method called RAPTOR, which combines linear programming with machine learning approaches.

## Structure Alignment and Indexing

In the fifth chapter “Algorithms for Multiple Protein Structure Alignment and Structure-Derived Multiple Sequence Alignment,” Maxim Shatsky, Ruth Nussinov, and Haim J. Wolfson, present methods to recognize the structural core common to a set of proteins through multiple structure alignment. They also discuss how to align multiple sequences with knowledge derived from structural alignment.

In the sixth chapter “Indexing Protein Structures using Suffix Trees,” Feng Gao and Mohammed J. Zaki describe a new approach to 3D database searching for protein sub-structures. Given a large set of proteins, they extract local structural features, which are converted into a set of symbols, which can be indexed using a traditional suffix tree. They show how one can rapidly retrieve approximately similar protein substructure matching a query protein.

## Protein Features Prediction

In the seventh chapter “Hidden Markov Models for Prediction of Protein Features,” Christopher Bystroff and Anders Krogh present a comprehensive overview of Hidden Markov Models (HMMs), which are used extensively in protein structure/sequence algorithms. They specifically focus on the applications of HMMs to predict signal peptides, secondary and local structure, and transmembrane helices.

In the eighth chapter “The Pros and Cons of Predicting Protein Contact Maps,” Lisa Bartoli, Emidio Capriotti, Piero Fariselli, Pier Luigi Martelli, and Rita Casadio discuss methods to predict protein contact maps. Contact maps are “simplified” 2D representations of the 3D proteins structure yet, they retain most of the important features for protein folding. They discuss the strengths and weaknesses of the contact map representation and highlight ways to improve contact map predictions.

In the ninth chapter “Road Map Methods for Protein Folding,” Mark Moll, David Schwarz, and Lydia E. Kavraci give a comprehensive survey of “roadmap” approaches to protein folding. Roadmap methods, inspired by motion planning techniques in robotics research, provide a model for understanding and predicting the folding mechanism or pathway.

### **Methods for De Novo Structure Prediction**

In the tenth Chapter “Scoring Functions for De Novo Protein Structure Prediction Revisited,” Shing-Chung Ngan, Ling-Hong Hung, Tianyun Liu, and Ram Samudrala, provide a thorough review of both physics-based and knowledge-based scoring functions for conformational samples in de novo protein structure prediction.

In the eleventh chapter “Protein–Protein Docking: Overview and Performance Analysis,” Kevin Wiehe, Matthew W. Peterson, Brian Pierce, Julian Mintseris, and Zhiping Weng focus on Fast Fourier Transform-based methods for protein docking. They specifically focus on the ZDOCK algorithm and study its performance on benchmark datasets and study its strengths and weaknesses through regression analysis.

In the final chapter “Molecular Dynamics Simulations of Protein Folding,” Angel E. Garcia describes the Replica Exchange Molecular Dynamics (REMD) method for molecular dynamics simulation. He illustrates the effectiveness of the REMD method on the folding of a small protein.

**Mohammed J. Zaki**  
**Chris Bystroff**

---

# Contents

Preface .....	v
Contributors .....	xi

## **PART I: OVERVIEW OF PROTEIN STRUCTURE PREDICTION**

1 A Historical Perspective of Template-Based Protein Structure Prediction <i>Jun-tao Guo, Kyle Ellrott, and Ying Xu</i> .....	3
2 The Assessment of Methods for Protein Structure Prediction <i>Anna Tramontano, Domenico Cozzetto, Alejandro Giorgetti, and Domenico Raimondo</i> .....	43

## **PART II: TEMPLATE-BASED METHODS**

3 Aligning Sequences to Structures <i>Liam James McGuffin</i> .....	61
4 Protein Structure Prediction Using Threading <i>Jinbo Xu, Feng Jiao, and Libo Yu</i> .....	91

## **PART III: STRUCTURE ALIGNMENT AND INDEXING**

5 Algorithms for Multiple Protein Structure Alignment and Structure-Derived Multiple Sequence Alignment <i>Maxim Shatsky, Ruth Nussinov, and Haim J. Wolfson</i> .....	125
6 Indexing Protein Structures Using Suffix Trees <i>Feng Gao and Mohammed J. Zaki</i> .....	147

## **PART IV: PROTEIN FEATURES PREDICTION**

7 Hidden Markov Models for Prediction of Protein Features <i>Christopher Bystroff and Anders Krogh</i> .....	173
8 The Pros and Cons of Predicting Protein Contact Maps <i>Lisa Bartoli, Emidio Capriotti, Piero Fariselli, Pier Luigi Martelli, and Rita Casadio</i> .....	199
9 Roadmap Methods for Protein Folding <i>Mark Moll, David Schwarz, and Lydia E. Kavraki</i> .....	219



**PART V: METHODS FOR DE NOVO STRUCTURE PREDICTION**

- 10 Scoring Functions for De Novo Protein Structure Prediction  
Revisited  
*Shing-Chung Ngan, Ling-Hong Hung, Tianyun Liu, and Ram  
Samudrala* ..... 243
- 11 Protein–Protein Docking: Overview and Performance Analysis  
*Kevin Wiehe, Matthew W. Peterson, Brian Pierce, Julian Mintseris,  
and Zhiping Weng* ..... 283
- 12 Molecular Dynamics Simulations of Protein Folding  
*Angel E. Garcia* ..... 315
- Index* ..... 331

---

# Contributors

- LISA BARTOLI • *CIRB/Department of Biology, University of Bologna, Bologna, Italy*
- CHRISTOPHER BYSTROFF • *Department of Biology, Rensselaer Polytechnic Institute, Troy NY, USA*
- EMIDIO CAPIROTTI • *CIRB/Department of Biology, University of Bologna, Bologna, Italy*
- RITA CASADIO • *CIRB/Department of Biology, University of Bologna, Bologna, Italy*
- DOMENICO COZZETTO • *Department of Biochemical Sciences, University of Rome “La Sapienza”, Rome, Italy*
- KYLE ELLROTT • *Department of Biochemistry and Molecular Biology, University of Georgia, Athens, GA, USA*
- PIERO FARISELLI • *CIRB/Department of Biology, University of Bologna, Bologna, Italy*
- FENG GAO • *Department of Computer Science, Rensselaer Polytechnic Institute, Troy, NY, USA*
- ANGEL E. GARCIA • *Department of Physics, Applied Physics and Astronomy, Rensselaer Polytechnic Institute, Troy, NY 12180*
- ALEJANDRO GIORGETTI • *Department of Biochemical Sciences, University of Rome “La Sapienza”, Rome, Italy*
- JUN-TAO GUO • *Department of Biochemistry and Molecular Biology, University of Georgia, Athens, GA, USA*
- LING-HONG HUNG • *Department of Microbiology, University of Washington School of Medicine, Seattle, WA, USA*
- FENG JIAO • *School of Computer Science, University of Waterloo, Waterloo, Canada*
- LYDIA E. KAVRAKI • *Computer Science Department, Rice University, Houston, TX, USA*
- ANDERS KROGH • *The Bioinformatics Centre, Inst. Mol. Biol. and Physiology, University of Copenhagen, Denmark*
- TIANYUN LIU • *Department of Microbiology, University of Washington School of Medicine, Seattle, WA, USA*

- PIER LUIGI MARTELLI • *CIRB/Department of Biology, University of Bologna, Bologna, Italy*
- LIAM JONES MCGUFFIN • *The University of Reading, Reading, UK*
- JULIAN MINTSERIS • *Bioinformatics Program, Boston University, Boston, MA, USA*
- MARK MOLL • *Information Sciences Institute, University of Southern California, Marina del Rey, CA, USA*
- SHING-CHUNG NGAN, • *Department of Microbiology, University of Washington School of Medicine, Seattle, WA, USA*
- RUTH NUSSINOV • *Sackler Inst. of Molecular Medicine, Tel Aviv University, Tel Aviv, Israel and Basic Research Program, SAIC-Frederick, Inc., Frederick, MD, USA*
- MATTHEW W. PETERSON • *Department of Biomedical Engineering, Boston University, Boston, MA, USA*
- BRIAN PIERCE • *Bioinformatics Program, Boston University, Boston, MA, USA*
- DOMENICO RAIMONDO • *Department of Biochemical Sciences, University of Rome “La Sapienza”, Rome, Italy*
- RAM SAMUDRALA • *Department of Microbiology, University of Washington School of Medicine, Seattle, WA, USA*
- DAVID SCHWARZ • *Computer Science Department, Rice University, Houston, TX, USA*
- MAXIM SHATSKY • *School of Computer Science, Tel Aviv University, Tel Aviv, Israel*
- ANNA TRAMONTANO • *Department of Biochemical Sciences, University of Rome “La Sapienza”, Rome, Italy*
- ZHIPING WENG • *Department of Biomedical Engineering, Boston University, Boston, MA, USA*
- KEVIN WIEHE • *Bioinformatics Program, Boston University, Boston, MA, USA*
- HAIM J. WOLFSON • *School of Computer Science, Tel Aviv University, Tel Aviv, Israel*
- JINBO XU • *Toyota Technological Institute at Chicago, Chicago, IL, USA*
- YING XU • *Department of Biochemistry and Molecular Biology, University of Georgia, Athens, GA, USA*
- LIBO YU • *Bioinformatics Solutions Inc., Waterloo, Canada*
- MOHAMMED J. ZAKI • *Department of Computer Science, Rensselaer Polytechnic Institute, Troy, NY, USA*