

# INTRODUCTION TO PROTEIN STRUCTURE PREDICTION

**WILEY SERIES ON BIOINFORMATICS:  
COMPUTATIONAL TECHNIQUES AND ENGINEERING**

**Series Editors, Yi Pan & Albert Zomaya**

**Knowledge Discovery in Bioinformatics: Techniques, Methods and Applications / Xiaohua Hu & Yi Pan**

**Grid Computing for Bioinformatics and Computational Biology / Albert Zomaya & El-Ghazali Talbi**

**Analysis of Biological Networks / Björn H. Junker & Falk Schreiber**

**Bioinformatics Algorithms: Techniques and Applications / Ion Mandoiu & Alexander Zelikovsky**

**Machine Learning in Bioinformatics / Yanqing Zhang & Jagath C. Rajapakse**

**Biomolecular Networks / Luonan Chen, Rui-Sheng Wang, & Xiang-Sun Zhang**

**Computational Systems Biology / Huma Lodhi**

**Computational Intelligence and Pattern Analysis in Biology Informatics / Ujjwal Maulik, Sanghamitra, & Jason T. Wang**

**Mathematics of Bioinformatics: Theory, Practice, and Applications / Matthew He**

**Introduction to Protein Structure Prediction: Methods and Algorithms / Huzefa Rangwala & George Karypis**

---

# INTRODUCTION TO PROTEIN STRUCTURE PREDICTION

Methods and Algorithms

---

Edited by

**HUZEFA RANGWALA**  
**GEORGE KARYPIS**

 **WILEY**

A JOHN WILEY & SONS, INC., PUBLICATION

Copyright © 2010 by John Wiley & Sons, Inc. All rights reserved

Published by John Wiley & Sons, Inc., Hoboken, New Jersey  
Published simultaneously in Canada

No part of this publication may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, electronic, mechanical, photocopying, recording, scanning, or otherwise, except as permitted under Section 107 or 108 of the 1976 United States Copyright Act, without either the prior written permission of the Publisher, or authorization through payment of the appropriate per-copy fee to the Copyright Clearance Center, Inc., 222 Rosewood Drive, Danvers, MA 01923, (978) 750-8400, fax (978) 750-4470, or on the web at [www.copyright.com](http://www.copyright.com). Requests to the Publisher for permission should be addressed to the Permissions Department, John Wiley & Sons, Inc., 111 River Street, Hoboken, NJ 07030, (201) 748-6011, fax (201) 748-6008, or online at <http://www.wiley.com/go/permission>.

**Limit of Liability/Disclaimer of Warranty:** While the publisher and author have used their best efforts in preparing this book, they make no representations or warranties with respect to the accuracy or completeness of the contents of this book and specifically disclaim any implied warranties of merchantability or fitness for a particular purpose. No warranty may be created or extended by sales representatives or written sales materials. The advice and strategies contained herein may not be suitable for your situation. You should consult with a professional where appropriate. Neither the publisher nor author shall be liable for any loss of profit or any other commercial damages, including but not limited to special, incidental, consequential, or other damages.

For general information on our other products and services or for technical support, please contact our Customer Care Department within the United States at (800) 762-2974, outside the United States at (317) 572-3993 or fax (317) 572-4002.

Wiley also publishes its books in a variety of electronic formats. Some content that appears in print may not be available in electronic formats. For more information about Wiley products, visit our web site at [www.wiley.com](http://www.wiley.com).

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

Rangwala, Huzefa.

Introduction to protein structure prediction : methods and algorithms / Huzefa Rangwala, George Karypis.

p. cm.—(Wiley series in bioinformatics; 14)

Includes bibliographical references and index.

ISBN 978-0-470-47059-6 (hardback)

1. Proteins—Structure—Mathematical models. 2. Proteins—Structure—Computer simulation. I. Karypis, G. (George) II. Title.

QP551.R225 2010

572'.633—dc22

2010028352

Printed in Singapore

10 9 8 7 6 5 4 3 2 1

# CONTENTS

---

<b>PREFACE</b>	<b>vii</b>
<b>CONTRIBUTORS</b>	<b>xi</b>
<b>1 INTRODUCTION TO PROTEIN STRUCTURE PREDICTION</b>	<b>1</b>
<i>Huzefa Rangwala and George Karypis</i>	
<b>2 CASP: A DRIVING FORCE IN PROTEIN STRUCTURE MODELING</b>	<b>15</b>
<i>Andriy Kryshchak, Krzysztof Fidelis, and John Moult</i>	
<b>3 THE PROTEIN STRUCTURE INITIATIVE</b>	<b>33</b>
<i>Andras Fiser, Adam Godzik, Christine Orengo, and Burkhard Rost</i>	
<b>4 PREDICTION OF ONE-DIMENSIONAL STRUCTURAL PROPERTIES OF PROTEINS BY INTEGRATED NEURAL NETWORKS</b>	<b>45</b>
<i>Yaoqi Zhou and Eshel Faraggi</i>	
<b>5 LOCAL STRUCTURE ALPHABETS</b>	<b>75</b>
<i>Agnel Praveen Joseph, Aurélie Bornot, and Alexandre G. de Brevern</i>	
<b>6 SHEDDING LIGHT ON TRANSMEMBRANE TOPOLOGY</b>	<b>107</b>
<i>Gábor E. Tusnády and István Simon</i>	
<b>7 CONTACT MAP PREDICTION BY MACHINE LEARNING</b>	<b>137</b>
<i>Alberto J.M. Martin, Catherine Mooney, Ian Walsh, and Gianluca Pollastri</i>	
<b>8 A SURVEY OF REMOTE HOMOLOGY DETECTION AND FOLD RECOGNITION METHODS</b>	<b>165</b>
<i>Huzefa Rangwala</i>	
<b>9 INTEGRATIVE PROTEIN FOLD RECOGNITION BY ALIGNMENTS AND MACHINE LEARNING</b>	<b>195</b>
<i>Allison N. Tegge, Zheng Wang, and Jianlin Cheng</i>	

<b>10</b>	<b>TASSER-BASED PROTEIN STRUCTURE PREDICTION</b>	<b>219</b>
	<i>Shashi Bhushan Pandit, Hongyi Zhou, and Jeffrey Skolnick</i>	
<b>11</b>	<b>COMPOSITE APPROACHES TO PROTEIN TERTIARY STRUCTURE PREDICTION: A CASE-STUDY BY I-TASSER</b>	<b>243</b>
	<i>Amrisha Roy, Sitao Wu, and Yang Zhang</i>	
<b>12</b>	<b>HYBRID METHODS FOR PROTEIN STRUCTURE PREDICTION</b>	<b>265</b>
	<i>Dmitri Mourado, Bostjan Kobe, Nicholas E. Dixon, and Thomas Huber</i>	
<b>13</b>	<b>MODELING LOOPS IN PROTEIN STRUCTURES</b>	<b>279</b>
	<i>Narcis Fernandez-Fuentes, Andras Fiser</i>	
<b>14</b>	<b>MODEL QUALITY ASSESSMENT USING A STATISTICAL PROGRAM THAT ADOPTS A SIDE CHAIN ENVIRONMENT VIEWPOINT</b>	<b>299</b>
	<i>Genki Terashi, Mayuko Takeda-Shitaka, Kazuhiko Kanou and Hideaki Umeyama</i>	
<b>15</b>	<b>MODEL QUALITY PREDICTION</b>	<b>323</b>
	<i>Liam J. McGuffin</i>	
<b>16</b>	<b>LIGAND-BINDING RESIDUE PREDICTION</b>	<b>343</b>
	<i>Chris Kauffman and George Karypis</i>	
<b>17</b>	<b>MODELING AND VALIDATION OF TRANSMEMBRANE PROTEIN STRUCTURES</b>	<b>369</b>
	<i>Maya Schushan and Nir Ben-Tal</i>	
<b>18</b>	<b>STRUCTURE-BASED MACHINE LEARNING MODELS FOR COMPUTATIONAL MUTAGENESIS</b>	<b>403</b>
	<i>Majid Masso and Iosif I. Vaisman</i>	
<b>19</b>	<b>CONFORMATIONAL SEARCH FOR THE PROTEIN NATIVE STATE</b>	<b>431</b>
	<i>Amarda Shehu</i>	
<b>20</b>	<b>MODELING MUTATIONS IN PROTEINS USING MEDUSA AND DISCRETE MOLECULE DYNAMICS</b>	<b>453</b>
	<i>Shuangye Yin, Feng Ding, and Nikolay V. Dokholyan</i>	
	<b>INDEX</b>	<b>477</b>

# PREFACE

---

## PROTEIN STRUCTURE PREDICTION

Proteins play a crucial role in governing several life processes. Stunningly complex networks of proteins perform innumerable functions in every living cell. Knowing the function and structure of proteins is crucial for the development of better drugs, higher yield crops, and even synthetic biofuels. As such, knowledge of protein structure and function leads to crucial advances in life sciences and biology. The motivation behind the structural determination of proteins is based on the belief that structural information provides insights as to their function, which will ultimately result in a better understanding of intricate biological processes.

Breakthroughs in large-scale sequencing have led to a surge in the available protein sequence information that has far outstripped our ability to characterize the structural and functional characteristic of these proteins. Several research groups have been working on determining the three-dimensional structure of the protein using a wide variety of computational methods. The problem of unraveling the relationship between the amino acid sequence of a protein and its three-dimensional structure has been one of the grand challenges in molecular biology. The importance and the far reaching implications of being able to predict the structure of a protein from its amino acid sequence is manifested by the ongoing biennial competition on “Critical Assessment of Protein Structure Prediction” (CASP) that started more than 16 years ago. CASP is designed to assess the performance of current structure prediction methods and over the years the number of groups that have been participating in it continues to increase.

This book presents a series of chapters by authors who are involved in the task of structure determination and using modeled structures for applications involving drug discovery and protein design. The book is divided into the following themes.

## **BACKGROUND ON STRUCTURE PREDICTION**

Chapter 1 provides an introduction to the protein structure prediction problem along with information about databases and resources that are widely used. Chapters 2 and 3 provide information regarding two very important initiatives in the field: (i) the structure prediction flagship competition (CASP), and (ii) the protein structure initiative (PSI), respectively. Since many of the approaches developed have been tested in the CASP competition, Chapter 2 lays the foundation for the need for such an evaluation, the problem definitions, significant innovations, competition format, as well as future outlook. Chapter 3 describes the protein structure initiative, which is designed to determine representative three-dimensional structures within the human genome.

## **PREDICTION OF STRUCTURAL ELEMENTS**

Within each structural entity called a protein there lies a set of recurring substructures, and within these substructures are smaller substructures. Beyond the goal of predicting the three-dimensional structure of a protein from sequence several other problems have been defined and methods have been developed for solving the same. Chapters 4–6 provide the definitions of these recurring substructures called local alphabets or secondary structures and the computational approaches used for solving these problems. Chapter 6 specifically focuses on a class of transmembrane proteins known to be harder to crystallize. Knowing the pairs of residues within a protein that are within contact or at a closer distance provides useful distance constraints that can be used while modeling the three-dimensional structure of the protein. Chapter 7 focuses on the problem of contact map prediction and also shows the use of sophisticated machine learning methods to solve the problem. A successful solution for each of these subproblems assists in solving the overarching protein structure prediction problem.

## **TERTIARY STRUCTURE PREDICTION**

Chapters 8–11 discuss the widely used structure prediction methods that rely on homology modeling, threading, and fragment assembly. Chapters 8–9 discuss the problems of fold recognition and remote homology detection that attempt to model the three-dimensional structure of a protein using known structures. Chapters 10 and 11 discuss a combination of threading-based approaches along with modeling the protein in parts or fragments and usually helps in modeling the structure of proteins known not to have a close homolog within the structure databases. Chapter 12 is a survey of the hybrid methods that use a combination of the computational and experimental methods to achieve high-resolution protein structures in a high-throughput manner.



Chapter 17 provides information about the challenges in modeling transmembrane proteins along with a discussion of some of the widely used methods for these sets of proteins.

Chapter 13 describes the loop prediction problem and how the technique can be used for refinement of the modeled structures. Chapters 14 and 15 assess the modeled structures and provide a notion of the quality of structures. This is extremely important from a biologist's perspective who would like to have a metric that describes the goodness of the structure before use. Chapter 19 provides insights into the different conformations that a protein may take and the approaches used to sample the different conformations.

## **FUNCTIONAL INSIGHTS**

Certain parts of the protein structure may be conserved and interact with other biomolecules (e.g., proteins, DNA, RNA, and small molecules) and perform a particular function due to such interactions. Chapter 16 discusses the problem of ligand-binding site prediction and its role in determining the function of the proteins. The approach uses some of the homology modeling principles used for modeling the entire structure. Chapter 18 introduces a computational model that detects the differences between protein structure (modeled or experimentally-determined) and its modeled mutant. Chapter 20 describes the use of molecular dynamic-based approaches for modeling mutants.

## **ACKNOWLEDGEMENTS**

We wish to acknowledge the many people who have helped us with this project. We firstly thank all the coauthors who spent time and energy to edit their chapters and also served as reviewers by providing critical feedback for improving other chapters. Kevin Deronne, Christopher Kauffman, and Rezwan Ahmed also assisted in reviewing several of the chapters and helped the book take a form that is complete on the topic of protein structure prediction and exciting to read. Finally, we wish to thank our families and friends.

We hope that you as a reader benefit from this book and feel as excited about this field as we are.

HUZEFA RANGWALA

GEORGE KARYPIS

# CONTRIBUTORS

---

NIR BEN-TAL, Department of Biochemistry and Molecular Biology, Tel Aviv University, Tel Aviv, Israel

AURÉLIE BORNOT, Institut National de la Santé et de la Recherche Médicale, UMR-S 665, Dynamique des Structures et Interactions des Macromolécules Biologiques (DSIMB), Université Paris Diderot, Paris, France

ALEXANDRE G. DE BREVERN, Institut National de la Santé et de la Recherche Médicale, Université Paris Diderot, Institut National de la Transfusion Sanguine, 75015, Paris, France

JIANLIN CHENG, Computer Science Department and Informatics Institute University of Missouri, Columbia, MO 65211

FENG DING, Department of Biochemistry and Biophysics University of North Carolina—Chapel Hill, NC 27599

NICHOLAS E. DIXON, School of Chemistry, University of Wollongong, NSW 2522, Australia

NIKOLAY V. DOKHOLYAN, Department of Biochemistry and Biophysics, University of North Carolina, Chapel Hill, NC 27599

ESHEL FARAGGI, Indiana University School of Informatics, Indiana University-Purdue University Indianapolis, and Center for Computational Biology and Bioinformatics, Indiana University School of Medicine, Indianapolis, IN 46202

KRZYSZTOF FIDELIS, Protein Structure Prediction Center, Genome Center, University of California, Davis, Davis, CA

ANDRAS FISER, Department of Systems and Computational Biology and Department of Biochemistry, Albert Einstein College of Medicine, Bronx, NY 10461

NARCIS FERNANDEZ-FUENTES, Leeds Institute of Molecular Medicine, University of Leeds, Leeds, UK

ADAM GODZIK, Program in Bioinformatics and Systems Biology, Sanford-Burnham Medical Research Institute, La Jolla, CA 92037

THOMAS HUBER, The University of Queensland, School of Chemistry and Molecular Biosciences, QLD, Australia

AGNEL PRAVEEN JOSEPH, Institut National de la Santé et de la Recherche Médicale, UMR-S 665, Dynamique des Structures et Interactions des Macromolécules Biologiques (DSIMB), Université Paris Diderot, Paris, France

KAZUHIKO KANOU, School of Pharmacy, Kitasato University, Tokyo 108-8641, Japan

GEORGE KARYPIS, Department of Computer Science, University of Minnesota Minneapolis, MN 55455

CHRIS KAUFFMAN, Department of Computer Science, University of Minnesota, Minneapolis, MN 55455

BOSTJAN KOBE, The University of Queensland, School of Chemistry and Molecular Biosciences, Brisbane, Australia

ANDRIY KRYSHTAFOVYCH, Protein Structure Prediction Center, Genome Center, University of California, Davis, Davis, CA

ALBERTO J.M. MARTIN, Complex and Adaptive Systems Lab, School of Computer Science and Informatics, UCD Dublin, Ireland

MAJID MASSA, Department of Bioinformatics and Computational Biology, George Mason University, Manassas, VA 20110

LIAM J. MCGUFFIN, School of Biological Sciences, The University of Reading, Reading, UK

CATHERINE MOONEY, Shields Lab, School of Medicine and Medical Science, University College Dublin, Ireland

JOHN MOULT, Institute for Bioscience and Biotechnology Research, University of Maryland, Rockville, MD 20850

DMITRI MOURADOV, The University of Queensland, School of Chemistry and Molecular Biosciences, QLD, Australia

CHRISTINE ORENGO, Department of Structural and Molecular Biology, University College London, London UK

SHASHI BHUSHAN PANDIT, Center for the Study of Systems Biology, School of Biology, Georgia Institute of Technology, Atlanta, GA 30318

GIANLUCA POLLASTRI, Complex and Adaptive Systems Lab, School of Computer Science and Informatics, UCD Dublin, Ireland

HUZEFA RANGWALA, Department of Computer Science, George Mason University, Fairfax, VA 22030

- BURKHARD ROST, Department of Biochemistry and Molecular Biophysics, Columbia University, New York, NY 10032
- AMBRISH ROY, Center for Computational Medicine and Bioinformatics, University of Michigan, Ann Arbor, MI 48109
- MAYA SCHUSHAN, Department of Biochemistry and Molecular Biology, Tel Aviv University, Tel Aviv, Israel
- AMARDA SHEHU, Department of Computer Science, George Mason University, Fairfax, VA 22030
- MAYUKO TAKEDA-SHITAKA, School of Pharmacy, Kitasato University, Tokyo 108-8641, Japan
- ISTVÁN SIMON, Intsitute of Enzymology, BRC, Hungarian Academy of Sciences, Budapest, Hungary
- JEFFREY SKOLNICK, Center for the Study of Systems Biology, School of Biology, Georgia Institute of Technology Atlanta, GA 30318
- ALLISON N. TEGGE, Computer Science Department and Informatics Institute, University of Missouri, Columbia, MO 65211
- GENKI TERASHI, School of Pharmacy, Kitasato University, Tokyo 108-8641, Japan
- GÁBOR E. TUSNADY, Intsitute of Enzymology, BRC, Hungarian Academy of Sciences, Budapest, Hungary
- HIDEAKI UMEYAMA, School of Pharmacy, Kitasato University, Tokyo 108-8641, Japan
- IOSIF I. VAISMAN, Department of Bioinformatics and Computational Biology, George Mason University, Manassas, VA 20110
- IAN WALSH, Complex and Adaptive Systems Lab, School of Computer Science and Informatics, UCD Dublin, Ireland
- ZHENG WANG, Computer Science Department, University of Missouri, Columbia, MO 65211
- SITAO WU, Center for Computational Medicine and Bioinformatics, University of Michigan, Ann Arbor, MI 48109
- SHUANGYE YIN, Department of Biochemistry and Biophysics, University of North Carolina, Chapel Hill, NC 27599
- YANG ZHANG, Center for Computational Medicine and Bioinformatics, University of Michigan, Ann Arbor, MI 48109
- HONGYI ZHOU, Center for the Study of Systems Biology, School of Biology Georgia Institute of Technology, Atlanta, GA 30318

YAOQI ZHOU, Indiana University School of Informatics, Indiana University-Purdue University Indianapolis, and Center for Computational Biology and Bioinformatics, Indiana University School of Medicine, Indianapolis, IN 46202