

Springer Protocols

Methods in Molecular Biology 541

Computational Systems Biology

Edited by

Jason McDermott
Ram Samudrala
Roger E. Bumgarner
Kristina Montgomery
Reneé Ireton

 Humana Press

METHODS IN MOLECULAR BIOLOGY™

Series Editor
John M. Walker
School of Life Sciences
University of Hertfordshire
Hatfield, Hertfordshire, AL10 9AB, UK

For other titles published in this series, go to
www.springer.com/series/7651

METHODS IN MOLECULAR BIOLOGY™

Computational Systems Biology

Edited by

Jason McDermott

Pacific Northwest National Laboratory, Richland, WA, USA

Ram Samudrala

University of Washington, Seattle, WA, USA

Roger E. Bumgarner

University of Washington, Seattle, WA, USA

Kristina Montgomery

University of Washington, Seattle, WA, USA

René Ireton

Fred Hutchinson Cancer Research Center, Seattle, WA, USA

 **Humana Press**

Editors

Jason McDermott
Pacific Northwest National
Laboratory
Computational Biology
& Bioinformatics Group
Richland, WA, USA
jason.mcdermott@pnl.gov

Ram Samudrala
Department of Microbiology
University of Washington
960 Republican St., WA, USA
ram@compbio.washington.edu

Roger E. Bumgarner
Department of Microbiology
University of Washington
960 Republican St., WA, USA
rogerb@u.washington.edu

Kristina Montgomery
Department of Microbiology
University of Washington
960 Republican St., WA, USA
kmontgom@u.washington.edu

Reneé Ireton
Fred Hutchinson Cancer Research
Center
Public Health Sciences Div.
Molecular Diagnostics Program
WA, USA
rireton@fhcrc.org

ISSN 1064-3745 e-ISSN 1940-6029
ISBN 978-1-58829-905-5 e-ISBN 978-1-59745-243-4
DOI 10.1007/978-1-59745-243-4

Library of Congress Control Number: 2009920380

© Humana Press, a part of Springer Science+Business Media, LLC 2009

All rights reserved. This work may not be translated or copied in whole or in part without the written permission of the publisher (Humana Press, c/o Springer Science+Business Media, LLC, 233 Spring Street, New York, NY 10013, USA), except for brief excerpts in connection with reviews or scholarly analysis. Use in connection with any form of information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed is forbidden.

The use in this publication of trade names, trademarks, service marks, and similar terms, even if they are not identified as such, is not to be taken as an expression of opinion as to whether or not they are subject to proprietary rights.

Cover illustration: Figure 1 in Chapter1.

Printed on acid-free paper

springer.com

Preface

Computational systems biology is the term that we use to describe computational methods to identify, infer, model, and store relationships between the molecules, pathways, and cells (“systems”) involved in a living organism. Based on this definition, the field of computational systems biology has been in existence for some time. However, the recent confluence of high-throughput methodology for biological data gathering, genome-scale sequencing, and computational processing power has driven a reinvention and expansion of this field. The expansions include not only modeling of small metabolic (1–3) and signaling systems (2, 4) but also modeling of the relationships between biological components in very large systems, including whole cells and organisms (5–15). Generally, these models provide a general overview of one or more aspects of these systems and leave the determination of details to experimentalists focused on smaller subsystems. The promise of such approaches is that they will elucidate patterns, relationships, and general features, which are not evident from examining specific components or subsystems. These predictions are either interesting in and of themselves (e.g., the identification of an evolutionary pattern) or interesting and valuable to researchers working on a particular problem (e.g., highlight a previously unknown functional pathway).

Two events have occurred to bring the field of computational systems biology to the forefront. One is the advent of high-throughput methods that have generated large amounts of information about particular systems in the form of genetic studies, gene and protein expression analyses and metabolomics. With such tools, research to consider systems as a whole are being conceived, planned, and implemented experimentally on an ever more frequent and wider scale. The other event is the growth of computational processing power and tools. Methods to analyze large data sets of this kind are often computationally demanding and, as is the case in other areas, the field has benefited from continuing improvements in computational hardware and methods.

The field of computational biology is very much like a telescope with two sequential lenses: one lens represents the biological data and the other represents a computational and/or mathematical model of the data. Both lenses must be properly coordinated to yield an image that reflects biological reality. This means that the design parameters for both lenses must be designed in concert to create a system that yields a model of the organism, which provides both predictive and mechanistic information. The chapters in this book describe the construction of subcomponents of such a system. Computational systems biology is a rapidly evolving field and no single group of investigators has yet developed a complete system that integrates both data generation and data analysis in such a way so as to allow full and accurate modeling of any single biological organism. However, the field is rapidly moving in that direction. The chapters in this book represent a snapshot of the current methods being developed and used in the area of computational systems biology. Each method or database described within represents one or more steps on the path to a complete description of a biological system. How

these tools will evolve and ultimately be integrated is an area of intense research and interest. We hope that readers of this book will be motivated by the chapters within and become involved in this exciting area of research.

Organization of the Book

This volume is organized into five major parts: Network Components, Network Inference, Network Dynamics, Function and Evolutionary Systems Biology, and Computational Infrastructure for Systems Biology. Each section is described briefly below.

Part I – Network Components

This section focuses on methods to identify subcomponents of the complete networks. Ultimately, such subcomponents will need to be integrated with each other or used to inform other methods to arrive at a complete description of a biological system. This section begins with two methods for the prediction of transcription factor binding sites. In the first, Chapter 1, Mariño-Ramirez et al. describe a method for the prediction of transcription factor binding sites using a Gibbs sampling approach. In Chapter 2, Liu and Bader show how DNA-binding sites and specificity can be predicted using sophisticated structural analysis. Chapters 3–5 discuss methods to predict protein–protein interaction (PPI) networks, and Chapter 6 builds on predicted PPIs to identify potential regulatory interactions. Finally, Chapter 7 discusses the inherent modularity that is observed in biological networks with a focus on networks of PPIs.

Part II – Network Inference

This section focuses on methodologies to infer transcriptional networks on a genome-wide scale. In general, the methods described within focus on using either mRNA expression data or mRNA expression data coupled with expression quantitative trait locus (eQTL) data. To a large extent, method development in this area is driven primarily by the ubiquitous mRNA expression data that are available in the public domain or that are relatively easily generated within a single laboratory. These methods have been tremendously enabled by the development of array technology and hence predominately model mRNA levels (as that is the most ubiquitous data type). Chapters 8 and 9 present two methods for identifying and modeling transcriptional regulatory networks, while Chapter 10 focuses on inferring mRNA expression networks from eQTL data. Chapter 11 is a review of different methods for inferring and modeling large scale networks from expression and eQTL data.

Part III – Network Dynamics

Systems are not static entities. They change over time and in response to a variety of perturbations. Ultimately, computational systems biology will have to develop methods and corresponding data sets that allow one to infer and model the kinetics and

dynamics of reactions between all the chemical moieties in a cell. The chapters in this section focus on such methods. Chapter 12 discusses methods to infer both static co-expression networks and a finite-state Markov chain model for mimicking the dynamic behavior of a transcriptional network. Chapter 13 focuses on quantitative models of system behavior based on differential equations using biochemical control theory, whereas Chapter 14 focuses on the use of stochastic kinetic simulations. Both approaches have applications where one is superior to the other. At this point in time, it is not clear which methods will turn out to be most useful in dynamically modeling the largest number of biological systems. In general, this is likely the case for most of the technologies described in this book, so it is useful for readers to familiarize themselves with several concepts. Specifically, both Chapters 13 and 14 provide an excellent discussion of a variety of historical approaches to the dynamical modeling of biological systems and the relative merits and downsides to each. Chapter 15 provides an excellent introduction to considerations for the interplay between experimental design and dynamic modeling using lambda phage as an example system. The methods and considerations described within are generally applicable to other biological systems and highlight the importance of integrating the direction of wet bench work and computational modeling to more rapidly refine the models.

Part IV – Function and Evolutionary Systems Biology

The ultimate representation of the function of a given biological moiety is a complete description of all the reactions in which it participates and the relative rates of said reactions. At present, we are quite distant from this goal for most biological molecules or systems. However, we are able to use computational methods to predict the most likely functions of a given protein and even predict which portions and specific sequences of the protein contribute most to that function. This section is focused on methods used to infer protein function and on the relationships between function and evolution.

Ultimately, the reason to study and research “systems” biology is to understand biological function at a given hierarchical level (be it a single catalytic site or entire pathways). The interplay between the detailed atomic study of function and the large-scale study of systems will enable us to achieve this goal. This section contains chapters that address the interdependence of these two aspects: individual algorithms or techniques to understand the functional role of atoms or residues in single molecules (e.g., proteins), which in turn are extrapolated to understand their greater role in terms of biological or organismal function. Conversely and complementarily, the role of larger systems and their influence on single molecules is also explored. Together, all these chapters illustrate the strong dependence between single molecules and entire pathways or systems.

Part V – Computational Infrastructure for Systems Biology

To represent and organize the large amounts of experimental data and software tools, database frameworks must be created and made available to the larger biological

community. This chapter focuses on computational methods and databases as well as data representations necessary to both integrate and export systems biology information to an end user. The user may be the biologist searching for their gene of interest or they may be the bioinformatician looking for trends in protein function among higher eukaryotes. Several groups are working on this extremely difficult task of providing semantic meaning to the large amounts of underlying biological data collected from single and high-throughput experiments, as well as computational predictions. (As a parenthetical comment, this is a significantly much harder problem than one faced by Internet search engines such as a Google, which at this point do not provide any semantic meaning to a query.) We present only a few such examples in this section (and in this book). One primary focus is on the Bioverse framework, database, and web application, which was developed by the editors of this book. However, we also describe the Biozon as well as the SEBIN and CABIN frameworks. The abstract representations required to model biological systems are still in fruition, and a complement of many tools, technologies, databases, and algorithms will have to be integrated in the future as our knowledge expands.

Acknowledgments

It goes without saying that this volume would not have been possible without the efforts of countless biologists who provided the raw data that enable modeling of biological systems. We first and foremost thank all these researchers who provide the raw data for all the computational modeling that enables the field of “computational” systems biology. We also thank all the researchers who investigate these problems using sophisticated modeling techniques, particularly those described in this volume. We thank two particular editors of this volume, René Iretton and Kristina Montgomery, who have dealt with the thankless job of undertaking the proofreading the chapters in this book and associated bureaucratic requirements. We finally thank the McDermott, Bumgarner, and Samudrala groups for their critical comments as this volume was being prepared, as well as our respective families for their patience. The administrative aspects of this work were in part supported by the NSF CAREER award to Dr. Ram Samudrala.

References

1. Ishii N, Robert, M, Nakayama, Y, Kanai, A and Tomita, M. Toward large-scale modeling of the microbial cell for computer simulation. *J Biotechnol* 2004;113(1-3):281–94.
2. Ekins S, Nikolsky, Y, Bugrim, A, Kirillov, E and Nikolskaya, T. Pathway mapping tools for analysis of high content data. *Methods Mol Biol* 2006;356:319–50.
3. Lafaye A, Junot, C, Pereira, Y, et al. Combined proteome and metabolite-profiling analyses reveal surprising insights into yeast sulfur metabolism. *J Biol Chem* 2005;280(26):24723–30.
4. Stevenson-Paulik J, Chiou, ST, Frederick, JP, et al. Inositol phosphate metabolomics: merging genetic perturbation with modernized radiolabeling methods. *Methods* 2006;39(2):112–21.
5. Ideker T, Thorsson, V, Ranish, JA, et al. Integrated genomic and proteomic analyses of a systematically perturbed metabolic network. *Science* 2001;292(5518):929–34.
6. Pe’er D, Regev, A, Elidan, G and Friedman, N. Inferring subnetworks from perturbed expression profiles. *Bioinformatics* 2001;17 Suppl 1:S215–24.

7. Pilpel Y, Sudarsanam, P and Church, GM. Identifying regulatory networks by combinatorial analysis of promoter elements. *Nat Genet* 2001;29(2):153–9.
8. Ideker T, Ozier, O, Schwikowski, B and Siegel, AF. Discovering regulatory and signalling circuits in molecular interaction networks. *Bioinformatics* 2002;18 Suppl 1:S233–40.
9. Kelley BP, Sharan, R, Karp, RM, et al. Conserved pathways within bacteria and yeast as revealed by global protein network alignment. *Proc Natl Acad Sci U S A* 2003;100(20):11394–9.
10. Shannon P, Markiel, A, Ozier, O, et al. Cytoscape: a software environment for integrated models of biomolecular interaction networks. *Genome Res* 2003;13(11):2498–504.
11. Ideker T. A systems approach to discovering signaling and regulatory pathways—or, how to digest large interaction networks into relevant pieces. *Adv Exp Med Biol* 2004;547:21–30.
12. Schadt EE, Monks, SA, Drake, TA, et al. Genetics of gene expression surveyed in maize, mouse and man. *Nature* 2003;422(6929):297–302.
13. Schadt EE and Lum, PY. Reverse engineering gene networks to identify key drivers of complex disease phenotypes. *J Lipid Res* 2006.
14. McDermott J and Samudrala R. Bioverse: Functional, structural and contextual annotation of proteins and proteomes. *Nucleic Acids Res* 2003;31(13):3736–7.
15. McDermott J, Bumgarner, R and Samudrala, R. Functional annotation from predicted protein interaction networks. *Bioinformatics* 2005;21(15):3217–26.

Contents

<i>Preface</i>	<i>v</i>
<i>Contributors</i>	<i>xiii</i>
<i>Color Plates</i>	<i>xvii</i>

PART I: NETWORK COMPONENTS

1. Identification of <i>cis</i> -Regulatory Elements in Gene Co-expression Networks Using A-GLAM	3
<i>Leonardo Mariño-Ramírez, Kannan Tharakaraman, Olivier Bodenreider, John Spouge, and David Landsman</i>	
2. Structure-Based <i>Ab Initio</i> Prediction of Transcription Factor–Binding Sites	23
<i>L. Angela Liu and Joel S. Bader</i>	
3. Inferring Protein–Protein Interactions from Multiple Protein Domain Combinations	43
<i>Simon P. Kanaan, Chengbang Huang, Stefan Wuchty, Danny Z. Chen, and Jesús A. Izaguirre</i>	
4. Prediction of Protein–Protein Interactions: A Study of the Co-evolution Model	61
<i>Itai Sharon, Jason V. Davis, and Golan Yona</i>	
5. Computational Reconstruction of Protein–Protein Interaction Networks: Algorithms and Issues	89
<i>Eric Franzosa, Bolan Linghu, and Yu Xia</i>	
6. Prediction and Integration of Regulatory and Protein–Protein Interactions	101
<i>Duangdao Wichadakul, Jason McDermott, and Ram Samudrala</i>	
7. Detecting Hierarchical Modularity in Biological Networks	145
<i>Erzsébet Ravasz</i>	

PART II: NETWORK INFERENCE

8. Methods to Reconstruct and Compare Transcriptional Regulatory Networks	163
<i>M. Madan Babu, Benjamin Lang, and L. Aravind</i>	
9. Learning Global Models of Transcriptional Regulatory Networks from Data	181
<i>Aviv Madar and Richard Bonneau</i>	
10. Inferring Molecular Interactions Pathways from eQTL Data	211
<i>Imran Rashid, Jason McDermott, and Ram Samudrala</i>	
11. Methods for the Inference of Biological Pathways and Networks	225
<i>Roger E. Bumgarner and Ka Yee Yeung</i>	

PART III: NETWORK DYNAMICS

12. Exploring Pathways from Gene Co-expression to Network Dynamics	249
<i>Huai Li, Yu Sun, and Ming Zhan</i>	
13. Network Dynamics	269
<i>Herbert M. Sauro</i>	

14.	Kinetic Modeling of Biological Systems	311
	<i>Haluk Resat, Linda Petzold, and Michel F. Pettigrew</i>	
15.	Guidance for Data Collection and Computational Modelling of Regulatory Networks	337
	<i>Adam Christopher Palmer, and Keith Edward Shearwin</i>	
PART IV: FUNCTION AND EVOLUTIONARY SYSTEMS BIOLOGY		
16.	A Maximum Likelihood Method for Reconstruction of the Evolution of Eukaryotic Gene Structure	357
	<i>Liran Carmel, Igor B. Rogozin, Yuri I. Wolf, and Eugene V. Koonin</i>	
17.	Enzyme Function Prediction with Interpretable Models	373
	<i>Umar Syed and Golan Yona</i>	
18.	Using Evolutionary Information to Find Specificity-Determining and Co-evolving Residues	421
	<i>Grigory Kolesov and Leonid A. Mirny</i>	
19.	Connecting Protein Interaction Data, Mutations, and Disease Using Bioinformatics	449
	<i>Jake Y. Chen, Eunseog Youn, and Sean D. Mooney</i>	
20.	Effects of Functional Bias on Supervised Learning of a Gene Network Model	463
	<i>Insuk Lee and Edward M. Marcotte</i>	
PART V: COMPUTATIONAL INFRASTRUCTURE FOR SYSTEMS BIOLOGY		
21.	Comparing Algorithms for Clustering of Expression Data: How to Assess Gene Clusters	479
	<i>Golan Yona, William Dirks, and Shafquat Rahman</i>	
22.	The Bioverse API and Web Application	511
	<i>Michal Guerquin, Jason McDermott, Zach Frazier, and Ram Samudrala</i>	
23.	Computational Representation of Biological Systems	535
	<i>Zach Frazier, Jason McDermott, Michal Guerquin, and Ram Samudrala</i>	
24.	Biological Network Inference and Analysis Using SEBINI and CABIN	551
	<i>Ronald Taylor and Mudita Singhal</i>	
	<i>Index</i>	577