

Computational Biology

Volume 24

Editors-in-Chief

Andreas Dress

CAS-MPG Partner Institute for Computational Biology, Shanghai, China

Michal Linial

Hebrew University of Jerusalem, Jerusalem, Israel

Olga Troyanskaya

Princeton University, Princeton, NJ, USA

Martin Vingron

Max Planck Institute for Molecular Genetics, Berlin, Germany

Editorial Board

Walter M. Fitch, Irvine, CA, USA

Robert Giegerich, University of Bielefeld, Bielefeld, Germany

Janet Kelso, Max Planck Institute for Evolutionary Anthropology, Leipzig, Germany

Gene Myers, Max Planck Institute of Molecular Cell Biology and Genetics, Dresden, Germany

Pavel Pevzner, University of California, San Diego, CA, USA

Advisory Board

G.M. Crippen, University of Michigan, Ann Arbor, MI, USA

Joseph Felsenstein, University of Washington, Seattle, WA, USA

Dan Gusfield, University of California, Davis, CA, USA

Sorin Istrail, Brown University, Providence, RI, USA

Sam Karlin, Stanford, CA, USA

Thomas Lengauer, Max Planck Institute for Computer Science, Saarbrücken, Germany

Marcella McClure, Montana State University, Bozeman, MO, USA

Martin Nowak, Harvard University, Cambridge, MA, USA

David Sankoff, University of Ottawa, Ottawa, ON, USA

R. Shamir, Tel Aviv University, Tel Aviv, Israel

Mike Steel, University of Canterbury, Christchurch, New Zealand

Gary Stormo, Washington University in St. Louis, St. Louis, MO, USA

Simon Tavaré, University of Cambridge, Cambridge, UK

Tandy Warnow, University of Texas, Austin, TX, USA

The *Computational Biology* series publishes the very latest, high-quality research devoted to specific issues in computer-assisted analysis of biological data. The main emphasis is on current scientific developments and innovative techniques in computational biology (bioinformatics), bringing to light methods from mathematics, statistics and computer science that directly address biological problems currently under investigation.

The series offers publications that present the state-of-the-art regarding the problems in question; show computational biology/bioinformatics methods at work; and finally discuss anticipated demands regarding developments in future methodology. Titles can range from focused monographs, to undergraduate and graduate textbooks, and professional text/reference works.

More information about this series at <http://www.springer.com/series/5769>

Sourav S. Bhowmick · Boon-Siew Seah

Summarizing Biological Networks

 Springer

Sourav S. Bhowmick
School of Computer Science
and Engineering
Nanyang Technological University
Singapore
Singapore

Boon-Siew Seah
School of Computer Science
and Engineering
Nanyang Technological University
Singapore
Singapore

ISSN 1568-2684

Computational Biology

ISBN 978-3-319-54620-9

ISBN 978-3-319-54621-6 (eBook)

DOI 10.1007/978-3-319-54621-6

Library of Congress Control Number: 2017935562

© Springer International Publishing AG 2017

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, express or implied, with respect to the material contained herein or for any errors or omissions that may have been made. The publisher remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Printed on acid-free paper

This Springer imprint is published by Springer Nature

The registered company is Springer International Publishing AG

The registered company address is: Gewerbestrasse 11, 6330 Cham, Switzerland

To our parents and wives.

Preface

Data do not give up their secrets easily. They must be tortured to confess.

Jeff Hopper, Bell Labs

The desire to study biology from a systems perspective has led to an emergence of new science—biological network analysis. Biological network models biological entities (e.g., proteins and genes) and their relationships (e.g., physical and genetic interactions) to characterize their cooperative activity within a system. With the rapid growth of such network data, the information overload problem has become a major stumbling block to analyze these networks, making human interpretation of such data increasingly difficult. Hence, there is a growing need to construct methods for large-scale topological and functional summaries of biological networks to understand the underlying mechanics of biological systems.

This book presents frameworks, as they stand today, that allow biologists to rapidly visualize and comprehend high-level topological and functional summary of the processes that govern biological systems via topological or functional organization *within* a biological network (intra-system processes) and relationships *between* biological networks (inter-system processes). Drawing on well-founded principles in data mining, systems biology, and bioinformatics, we present a multi-resolution and multi-perspective analysis paradigm to address this broad goal. Note that it is reasonable to expect this picture to change with time.

As a representative example of biological networks, we utilize protein–protein interaction (PPI) networks in majority part of this book. Our discussion is divided into five parts. First, we have attempted to review, as accurately as possible, a wide spectrum of approaches proposed by the bioinformatics community to cluster PPI networks and highlight their strengths and limitations. The results of such clustering can be considered as a summary of topological or functional modules in the underlying PPI network. In particular, a pervasive desire of this review is to

emphasize the uniqueness of the network clustering problem in the context of PPI networks and highlight why a panoply of generic network clustering algorithms proposed by the data mining community cannot be leveraged to address this problem effectively.

Second, we review a closely related problem to PPI network clustering, functional summarization, which can enable us to make sense out of the information contained in large PPI networks by generating multi-level functional summaries. We discuss a data-driven and generic PPI network summarization framework that constructs higher level functional summary to summarize the underlying PPI network to obtain a concise, interpretable representation of the network. It generates the “best” summary from both interaction and annotation data by maximizing information gain for a specific resolution. We evaluate the performance of this framework on several real-world PPI networks, its superiority over network clustering, and showcase its applicability in comprehending Alzheimer’s disease network.

Third, we discuss a technique that summarizes a PPI network in a multi-perspective manner. This is based on the fact that a biological system can be seen from different functional perspectives (e.g., components in a PPI network can be organized by localization, process, disease, etc.). Each discovered perspective represents a distinct interpretation of how the network can be functionally summarized. The performance of this framework is extensively discussed with several real-world PPI networks highlighting the limitations of network clustering paradigm to generate such multi-perspective summary. We also performed a case study using human autophagy system to illustrate the utility of this framework.

Fourth, we discuss a data-driven effort to construct summaries of differential functional responses of gene interaction networks that undergo “rewiring” after environmental change. Experimental evaluation with real-world dataset demonstrates the superiority of this technique to address the differential network summarization problem.

The last topic consists of several open problems of this young field. The list presented should by no means be considered exhaustive and is centered around challenges and issues currently in vogue. Nevertheless, readers can benefit by exploring the research directions given in this part.

The book is suitable for use in advanced undergraduate- and graduate-level courses on biological networks. It has sufficient material that can be covered as part of a semester-long course, thereby leaving plenty of room for an instructor to choose topics. An undergraduate course in algorithms, graph theory, and basic cell biology should suffice as a prerequisite for most of the chapters. A good knowledge of C++/Java programming language is sufficient to code the algorithms described herein. For completeness, we have provided background information on several topics in Chap. 2: the central dogma of biology, protein–protein interactions, high-throughput experimental techniques to analyze protein–protein interactions,

and annotations of these interactions with Gene Ontology. The knowledgeable reader may omit this chapter and perhaps refer back to comparisons while reading later chapters of this book.

We hope that this book will serve as a catalyst in helping this burgeoning area of biological network summarization grow and have practical impact.

Singapore
December 2016

Sourav S. Bhowmick
Boon-Siew Seah

Acknowledgements

It is a great pleasure for us to acknowledge the assistance and contributions of a large number of individuals to this effort. First, we would like to thank our publisher Springer-Verlag for their support. In particular, we would like to acknowledge the efforts, help, and patience of Melissa Fearon and Jennifer Malat, our primary contacts for this edition.

The work reported in this book grew out of the PANORAMA project at the Nanyang Technological University (NTU), Singapore, and Massachusetts Institute of Technology (MIT), USA, under the auspices of Singapore-MIT Alliance graduate program. In this project, we explored issues on building frameworks that allow biologists to rapidly visualize the processes that govern biological systems. Specifically, the chapters in this book are part of Boon Siew's thesis work under the guidance of Sourav. Some of these chapters are published in reputable journals and conferences in the area of bioinformatics and systems biology.

Dr. C. Forbes Dewey, Jr. of MIT, who was a key collaborator for this project, deserves the first thank you. Not only did he introduce us to interesting and exciting field of network biology and network medicine, but he was also always willing to discuss ideas with us, no matter how strange they were.

In addition, we would also like to express our gratitude to all the group members and collaborators, past and present, in our **Computational Systems Biology** research group (COSBY). In particular, Dr. Huey Eng Chua (NTU), Dr. Jie Zheng (NTU), Dr. Lisa Tucker-Kellogg (Duke-NUS Medical School), Dr. Hanry Yu (NUS), and Mengxuan Chen made substantial contributions to the broader aspect of our research in network biology.

Quite a few people have helped us with the initial vetting of the text for this book. It is our pleasure to acknowledge them all here. We would like to thank Scientific Publishing Services (SPS) for carefully proofreading the complete book in a short span of time and suggesting the changes which have been incorporated.

Sourav and Boon Siew would like to acknowledge their parents and family members who gave them incredible support throughout the years. A special thanks goes to Dr. Paul Matsudaira (MIT, NUS) and Dr. Hew Choy Leong (NUS), who

were a great motivator during our early days when we were grappling with the new field of systems biology. They were the major force behind our continuous strive for breaking out from the comfort zone of computer science to explore problems that are at the intersection of two or more disparate fields. It has been a great learning experience for us.

Finally, we would like to thank the Singapore-MIT Alliance for the generous resources and financial support provided for the PANORAMA project. We would also like to thank the School of Computer Science and Engineering at the Nanyang Technological University for allowing the use of their resources to help complete this book.

Singapore
December 2016

Sourav S. Bhowmick
Boon-Siew Seah

Contents

1 Introduction	1
1.1 Challenges	4
1.2 Overview of This Book	5
References	7
2 Background	9
2.1 Proteins: The Building Block of Life	9
2.2 Protein-Protein Interaction (PPI)	11
2.3 Methods to Analyze Protein-Protein Interactions	12
2.3.1 Yeast Two-Hybrid (Y2H)	12
2.3.2 Tandem Affinity Purification (TAP)	13
2.3.3 Bimolecular Fluorescence Complementation (BIFC)	13
2.3.4 Noise in High-Throughput Screening Methods	14
2.4 Protein-Protein Interaction Databases	15
2.5 Annotating the Roles of Proteins and Their Interactions	16
2.5.1 The Structure of Gene Ontology	17
2.6 Summary	19
References	19
3 Clustering PPI Networks	23
3.1 PPI Network Clustering Problem	24
3.1.1 Problem Definition	25
3.1.2 Challenges	25
3.1.3 Representative Clustering Measures	27
3.1.4 Overview	29
3.2 PPI Network Clustering Techniques	32
3.2.1 Heuristic-Based Algorithms	32
3.2.2 Flow-Based Algorithms	36
3.2.3 Complete Enumeration Algorithms	38
3.2.4 Random Walks and Message Passing Algorithms	40

3.2.5	Graph-Cut and Hierarchical Clustering Algorithms	42
3.2.6	Multiple Clustering-Based Algorithms	46
3.2.7	Genomic Data-Driven Clustering Algorithms	48
3.3	Cluster Validation Measures	49
3.3.1	Functional Homogeneity-Based Validation	50
3.3.2	MIPS-based Validation	51
3.3.3	Other Measures	52
3.4	Comparative Summary	52
3.5	Conclusions	54
	References	55
4	Functional Summarization	59
4.1	Motivation	59
4.2	Limitations of PPI Clustering Techniques	60
4.3	Overview	63
4.4	Related Work	63
4.5	The Functional Summarization Problem	65
4.5.1	Functional Summary of PPI	65
4.5.2	Problem Statement	69
4.6	The Algorithm FUSE	71
4.7	Experimental Results	74
4.7.1	Evaluation Metrics	75
4.7.2	FUSE Versus Graph Clustering Methods	76
4.7.3	Effects of Different Parameters	82
4.7.4	Runtime and Scalability	87
4.8	Case Study on AD Network	88
4.9	Inferring Functional Cluster Hubs	90
4.10	Conclusions	92
	References	92
5	Multi-faceted Functional Decomposition	95
5.1	Motivation	95
5.2	Related Work	97
5.3	Problem Statement	97
5.3.1	Terminology	97
5.3.2	Multi-faceted Functional Decomposition Problem	98
5.3.3	Problem Definition	101
5.4	FACETS Algorithm	101
5.4.1	The Initialization Phase	101
5.4.2	The Iteration Phase	102
5.5	Experimental Study	105
5.5.1	Experiment Settings	105
5.5.2	Results	106
5.6	Case Study: Human Autophagy System	113

5.7 Conclusions 115

References 115

6 Differential Functional Summarization 117

6.1 Background 117

6.2 Motivation and Overview 119

6.3 Functional Subgraphs in a Differential Network 120

6.3.1 Constructing Differential Networks 120

6.3.2 Functional Subgraphs 123

6.4 Differential Summarization Problem 125

6.5 The DiffNet Algorithm 127

6.6 Experimental Study 128

6.6.1 Functional Analysis of MMS-treated/untreated
dE-MAP Network 129

6.6.2 Comparison with Graph Clustering Algorithms 130

6.6.3 Effect of Various Parameters 132

6.6.4 Running Times 133

6.6.5 Effect of Interaction Noise 135

6.6.6 Effect of Annotation Loss 136

6.7 Conclusions 137

References 137

7 The Road Ahead 139

7.1 Summary 139

7.2 Future Research 140

References 142

Index 143