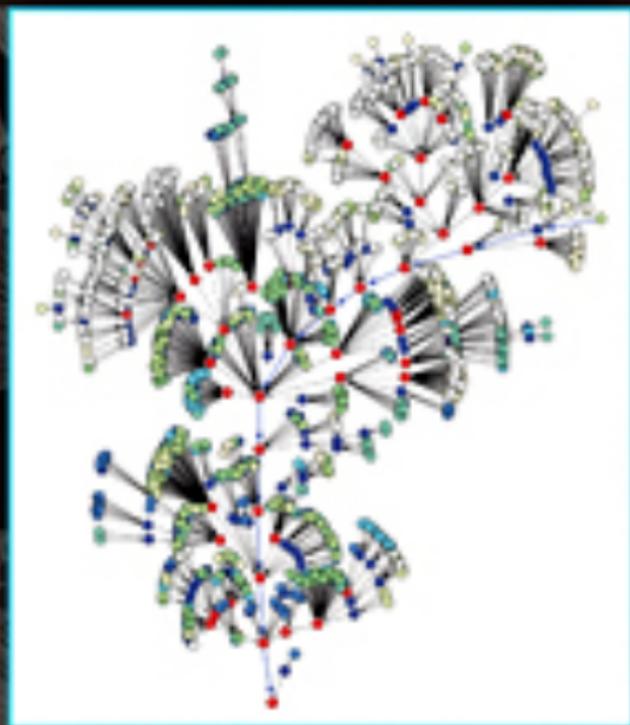


# PROTEIN INTERACTION NETWORKS

Computational Analysis



AIDONG ZHANG

CAMBRIDGE

CAMBRIDGE

[www.cambridge.org/9780521888950](http://www.cambridge.org/9780521888950)

This page intentionally left blank

## Protein Interaction Networks: Computational Analysis

The analysis of protein–protein interactions is fundamental to the understanding of cellular organization, processes, and functions. Proteins seldom act as single isolated species; rather, proteins involved in the same cellular processes often interact with each other. Functions of uncharacterized proteins may be predicted through comparison with the interactions of similar known proteins. Recent large-scale investigations of protein–protein interactions using such techniques as two-hybrid systems, mass spectrometry, and protein microarrays have enriched the available protein interaction data and facilitated the construction of integrated protein–protein interaction networks. The resulting large volume of protein–protein interaction data has posed a challenge to experimental investigation.

This book provides a comprehensive understanding of the computational methods available for the analysis of protein–protein interaction networks. It offers an in-depth survey of a range of approaches, including statistical, topological, data-mining, and ontology-based methods. The author discusses the fundamental principles underlying each of these approaches and their respective benefits and drawbacks, and she offers suggestions for future research.

Aidong Zhang is a professor in the Department of Computer Science and Engineering at the State University of New York at Buffalo and the director of the Buffalo Center for Biomedical Computing (BCBC). She is an author of more than 200 research publications and has served on the editorial boards of the *International Journal of Bioinformatics Research and Applications* (IJBRA), *ACM Multimedia Systems*, the *International Journal of Multimedia Tools and Applications*, the *International Journal of Distributed and Parallel Databases*, and *ACM SIGMOD DiSC* (Digital Symposium Collection). Dr. Zhang is a recipient of the National Science Foundation CAREER Award and SUNY (State University of New York) Chancellor’s Research Recognition Award. Dr. Zhang is an IEEE Fellow.



# PROTEIN INTERACTION NETWORKS

Computational Analysis

**Aidong Zhang**

State University of New York, Buffalo



**CAMBRIDGE**  
UNIVERSITY PRESS

CAMBRIDGE UNIVERSITY PRESS

Cambridge, New York, Melbourne, Madrid, Cape Town, Singapore, São Paulo

Cambridge University Press

The Edinburgh Building, Cambridge CB2 8RU, UK

Published in the United States of America by Cambridge University Press, New York

[www.cambridge.org](http://www.cambridge.org)

Information on this title: [www.cambridge.org/9780521888950](http://www.cambridge.org/9780521888950)

© Aidong Zhang 2009

This publication is in copyright. Subject to statutory exception and to the provision of relevant collective licensing agreements, no reproduction of any part may take place without the written permission of Cambridge University Press.

First published in print format 2007

ISBN-13 978-0-511-53049-4 eBook (Adobe Reader)

ISBN-13 978-0-521-88895-0 hardback

Cambridge University Press has no responsibility for the persistence or accuracy of urls for external or third-party internet websites referred to in this publication, and does not guarantee that any content on such websites is, or will remain, accurate or appropriate.

*To my daughter, Cathy*



# Contents

<i>Preface</i>	<i>page</i> xiii
<b>1 Introduction</b>	1
1.1 Rapid Growth of Protein–Protein Interaction Data	1
1.2 Computational Analysis of PPI Networks	3
1.2.1 Topological Features of PPI Networks	4
1.2.2 Modularity Analysis	5
1.2.3 Prediction of Protein Functions in PPI Networks	6
1.2.4 Integration of Domain Knowledge	7
1.3 Significant Applications	7
1.4 Organization of this Book	9
1.5 Summary	10
<b>2 Experimental Approaches to Generation of PPI Data</b>	11
2.1 Introduction	11
2.2 The Y2H System	11
2.3 Mass Spectrometry (MS) Approaches	13
2.4 Protein Microarrays	15
2.5 Public PPI Data and Their Reliability	15
2.5.1 Experimental PPI Data Sets	15
2.5.2 Public PPI Databases	16
2.5.3 Functional Analysis of PPI Data	17
2.6 Summary	20
<b>3 Computational Methods for the Prediction of PPIs</b>	21
3.1 Introduction	21
3.2 Genome-Scale Approaches	21
3.3 Sequence-Based Approaches	25
3.4 Structure-Based Approaches	26
3.5 Learning-Based Approaches	27
3.6 Network Topology-Based Approaches	29
3.7 Summary	32

<b>4</b>	<b>Basic Properties and Measurements of Protein Interaction Networks</b>	<b>33</b>
4.1	Introduction	33
4.2	Representation of PPI Networks	33
4.3	Basic Concepts	34
4.4	Basic Centralities	35
4.4.1	Degree Centrality	35
4.4.2	Distance-Based Centralities	35
4.4.3	Current-Flow-Based Centrality	37
4.4.4	Random-Walk-Based Centrality	40
4.4.5	Feedback-Based Centrality	41
4.5	Characteristics of PPI Networks	44
4.6	Summary	49
<b>5</b>	<b>Modularity Analysis of Protein Interaction Networks</b>	<b>50</b>
5.1	Introduction	50
5.2	Useful Metrics for Modular Networks	51
5.2.1	Cliques	51
5.2.2	Cores	51
5.2.3	Degree-Based Index	52
5.2.4	Distance (Shortest Paths)-Based Index	53
5.3	Methods for Clustering Analysis of Protein Interaction Networks	53
5.3.1	Traditional Clustering Methods	54
5.3.2	Nontraditional Clustering Methods	55
5.4	Validation of Modularity	56
5.4.1	Clustering Coefficient	56
5.4.2	Validation Based on Agreement with Annotated Protein Function Databases	57
5.4.3	Validation Based on the Definition of Clustering	59
5.4.4	Topological Validation	60
5.4.5	Supervised Validation	61
5.4.6	Statistical Validation	61
5.4.7	Validation of Protein Function Prediction	62
5.5	Summary	62
<b>6</b>	<b>Topological Analysis of Protein Interaction Networks</b>	<b>63</b>
	<i>With Woo-chang Hwang</i>	
6.1	Introduction	63
6.2	Overview and Analysis of Essential Network Components	64
6.2.1	Error and Attack Tolerance of Complex Networks	64
6.2.2	Role of High-Degree Nodes in Biological Networks	67
6.2.3	Betweenness, Connectivity, and Centrality	69
6.3	Bridging Centrality Measurements	73
6.3.1	Performance of Bridging Centrality with Synthetic and Real-World Networks	75
6.3.2	Assessing Network Disruption, Structural Integrity, and Modularity	77

6.4	Network Modularization Using the Bridge Cut Algorithm	84
6.5	Use of Bridging Nodes in Drug Discovery	87
6.5.1	Biological Correlates of Bridging Centrality	88
6.5.2	Results from Drug Discovery-Relevant Human Networks	92
6.5.3	Comparison to Alternative Approaches: Yeast Cell Cycle State Space Network	94
6.5.4	Potential of Bridging Centrality as a Drug Discovery Tool	95
6.6	PathRatio: A Novel Topological Method for Predicting Protein Functions	97
6.6.1	Weighted PPI Network	97
6.6.2	Protein Connectivity and Interaction Reliability	98
6.6.3	PathStrength and PathRatio Measurements	99
6.6.4	Analysis of the PathRatio Topological Measurement	100
6.6.5	Experimental Results	101
6.7	Summary	108
<b>7</b>	<b>Distance-Based Modularity Analysis</b>	<b>109</b>
7.1	Introduction	109
7.2	Topological Distance Measurement Based on Coefficients	109
7.3	Distance Measurement by Network Distance	112
7.3.1	PathRatio Method	112
7.3.2	Averaging the Distances	113
7.4	Ensemble Method	114
7.4.1	Similarity Metrics	115
7.4.2	Base Algorithms	116
7.4.3	Consensus Methods	116
7.4.4	Results of the Ensemble Methods	118
7.5	UVCLUSTER	118
7.6	Similarity Learning Method	120
7.7	Measurement of Biological Distance	124
7.7.1	Sequence Similarity-Based Measurements	124
7.7.2	Structural Similarity-Based Measurements	125
7.7.3	Gene Expression Similarity-Based Measurements	127
7.8	Summary	128
<b>8</b>	<b>Graph-Theoretic Approaches to Modularity Analysis</b>	<b>130</b>
8.1	Introduction	130
8.2	Finding Dense Subgraphs	130
8.2.1	Enumeration of Complete Subgraphs	130
8.2.2	Monte Carlo Optimization	131
8.2.3	Molecular Complex Detection	132
8.2.4	Clique Percolation	133
8.2.5	Merging by Statistical Significance	134
8.2.6	Super-Paramagnetic Clustering	136
8.3	Finding the Best Partition	137
8.3.1	Recursive Minimum Cut	137
8.3.2	Restricted Neighborhood Search Clustering (RNSC)	138

8.3.3	Betweenness Cut	140
8.3.4	Markov Clustering	140
8.3.5	Line Graph Generation	143
8.4	Graph Reduction-Based Approach	144
8.4.1	Graph Reduction	144
8.4.2	Hierarchical Modularization	146
8.4.3	Time Complexity	147
8.4.4	$k$ Effects on Graph Reduction	147
8.4.5	Hierarchical Structure of Modules	149
8.5	Summary	150
<b>9</b>	<b>Flow-Based Analysis of Protein Interaction Networks</b>	<b>152</b>
9.1	Introduction	152
9.2	Protein Function Prediction Using the FunctionalFlow Algorithm	153
9.3	CASCADE: A Dynamic Flow Simulation for Modularity Analysis	155
9.3.1	Occurrence Probability and Related Models	156
9.3.2	The CASCADE Algorithm	158
9.3.3	Analysis of Prototypical Data	160
9.3.4	Significance of Individual Clusters	162
9.3.5	Analysis of Functional Annotation	164
9.3.6	Comparative Assessment of CASCADE with Other Approaches	169
9.3.7	Analysis of Robustness	175
9.3.8	Analysis of Computational Complexity	175
9.3.9	Advantages of the CASCADE Method	176
9.4	Functional Flow Analysis in Weighted PPI Networks	177
9.4.1	Functional Influence Model	178
9.4.2	Functional Flow Simulation Algorithm	179
9.4.3	Time Complexity of Flow Simulation	180
9.4.4	Detection of Overlapping Modules	181
9.4.5	Detection of Disjoint Modules	189
9.4.6	Functional Flow Pattern Mining	191
9.5	Summary	198
<b>10</b>	<b>Statistics and Machine Learning Based Analysis of Protein Interaction Networks</b>	<b>199</b>
	<i>With Pritam Chanda and Lei Shi</i>	
10.1	Introduction	199
10.2	Applications of Markov Random Field and Belief Propagation for Protein Function Prediction	200
10.3	Protein Function Prediction Using Kernel-Based Statistical Learning Methods	207
10.4	Protein Function Prediction Using Bayesian Networks	211

10.5	Improving Protein Function Prediction Using Bayesian Integrative Methods	213
10.6	Summary	214
<b>11</b>	<b>Integration of GO into the Analysis of Protein Interaction Networks</b>	<b>216</b>
	<i>With Young-rae Cho</i>	
11.1	Introduction	216
11.2	GO structure	217
11.2.1	GO Annotations	217
11.3	Semantic Similarity-Based Integration	218
11.3.1	Structure-Based Methods	219
11.3.2	Information Content-Based Methods	220
11.3.3	Combination of Structure and Information Content	221
11.4	Semantic Interactivity-Based Integration	223
11.5	Estimate of Interaction Reliability	223
11.5.1	Functional Co-Occurrence	224
11.5.2	Topological Significance	225
11.5.3	Protein Lethality	226
11.6	Functional Module Detection	227
11.6.1	Statistical Assessment	227
11.6.2	Supervised Validation	229
11.7	Probabilistic Approaches for Function Prediction	231
11.7.1	GO Index-Based Probabilistic Method	231
11.7.2	Semantic Similarity-Based Probabilistic Method	235
11.8	Summary	241
<b>12</b>	<b>Data Fusion in the Analysis of Protein Interaction Networks</b>	<b>243</b>
12.1	Introduction	243
12.2	Integration of Gene Expression with PPI Networks	243
12.3	Integration of Protein Domain Information with PPI Networks	244
12.4	Integration of Protein Localization Information with PPI Networks	245
12.5	Integration of Several Data Sources with PPI Networks	247
12.5.1	Kernel-Based Methods	247
12.5.2	Bayesian Model-Based Method	249
12.6	Summary	249
<b>13</b>	<b>Conclusion</b>	<b>251</b>
	<i>Bibliography</i>	255
	<i>Index</i>	273



# Preface

I am pleased to offer the research community my second book-length contribution to the field of bioinformatics. My first book, *Advanced Analysis of Gene Expression Microarray Data*, was published in 2006 by World Scientific as part of its Science, Engineering, and Biology Informatics (SEBI) series. I first became involved in the study of bioinformatics in 1998 and, over the ensuing decade, have been struck by the enormous quantity of data being generated and the need for effective approaches to its analysis.

The analysis of protein–protein interactions (PPIs) is fundamental to the understanding of cellular organizations, processes, and functions. It has been observed that proteins seldom act as single isolated species in the performance of their functions; rather, proteins involved in the same cellular processes often interact with each other. Therefore, the functions of uncharacterized proteins can be predicted through comparison with the interactions of similar known proteins. A detailed examination of a PPI network can thus yield significant new insights into protein functions. These interactions have traditionally been examined via intensive small-scale investigations of a small set of proteins of interest, each yielding information about a limited number of PPIs. The existing databases of PPIs have been compiled from such small-scale screens, presented in individual research papers. Because these data were subject to stringent controls and evaluation in the peer-review process, they can be considered to be fairly reliable. However, each experiment observes only a few interactions and yields a data set of very limited size. Recent large-scale investigations of PPIs using such techniques as two-hybrid systems, mass spectrometry, and protein microarrays have enriched the available protein interaction data and facilitated the construction of integrated PPI networks. The resulting large volume of PPI data has posed a challenge to experimental investigation. Consequently, computational analysis of the networks has become a necessary tool for the determination of functionally associated proteins.

This book is intended to provide a comprehensive understanding of the computational methods available for the analysis of PPI networks. It offers an in-depth survey of a range of approaches to this analysis, including statistical, topological, data-mining, and ontology-based methods. The fundamental principles underlying each of

these approaches are discussed, along with their respective benefits and drawbacks. Suggestions for future research are also offered. In total, this book is intended to offer bioinformatics researchers a comprehensive and practical guide to the analysis of PPI networks, which will assist and stimulate their further investigation.

Some knowledge on the part of the reader in the fields of molecular biology, data mining, and statistics is assumed. Apart from this, the book is designed to be self-contained, as it includes introductions to the fundamental concepts underlying data generation and analysis. Thus, this book is expected to be of interest to a variety of researchers. It can be used as a textbook for advanced graduate courses in bioinformatics, and most of its content has been tested in the author's graduate-level course in this field. In addition, it can serve as a resource for graduate students seeking topics for investigation. The book will also be useful to researchers involved in computational biology in universities, organizations, and industry. For this audience, it will provide guidance on the techniques available for analysis of PPI networks. Research professionals interested in expanding their knowledge base can draw upon the material presented here to gain an understanding of principles and methods involved in this growing and highly significant field.

## **ACKNOWLEDGMENTS**

I would like to express my deepest thanks to my doctoral students, Pritam Chanda, Young-rae Cho, Woo-chang Hwang, Taehyong Kim, and Lei Shi, for their excellent technical contributions. I am also highly appreciative of the editorial work of Rachel Ramadhyani.

The inspiration for this book was an invitation from Ms. Lauren Cowles, a senior editor from Cambridge University Press. I would like to express my special thanks to her.

Aidong Zhang  
Buffalo, New York

# Introduction

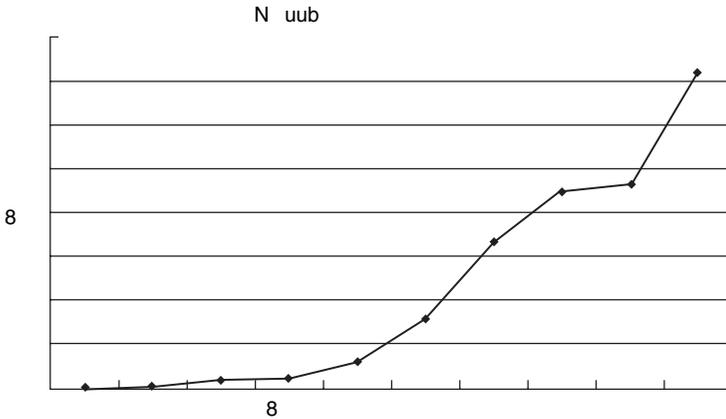
## 1.1 RAPID GROWTH OF PROTEIN–PROTEIN INTERACTION DATA

Since the sequencing of the human genome was brought to fruition [154,310], the field of genetics now stands on the threshold of significant theoretical and practical advances. Crucial to furthering these investigations is a comprehensive understanding of the expression, function, and regulation of the proteins encoded by an organism [345]. This understanding is the subject of the discipline of proteomics. Proteomics encompasses a wide range of approaches and applications intended to explicate how complex biological processes occur at a molecular level, how they differ in various cell types, and how they are altered in disease states.

Defined succinctly, proteomics is the systematic study of the many and diverse properties of proteins with the aim of providing detailed descriptions of the structure, function, and control of biological systems in health and disease [241]. The field has burst onto the scientific scene with stunning rapidity over the past several years. Figure 1–1 shows the trend of the number of occurrences of the term “proteome” found in PubMed bioinformatics citations over the past decade. This figure strikingly illustrates the rapidly increasing role played by proteomics in bioinformatics research in recent years.

A particular focus of the field of proteomics is the nature and role of interactions between proteins. Protein–protein interactions (PPIs) regulate a wide array of biological processes, including transcriptional activation/repression; immune, endocrine, and pharmacological signaling; cell-to-cell interactions; and metabolic and developmental control [9,139,167,184]. PPIs play diverse roles in biology and differ based on the composition, affinity, and lifetime of the association. Noncovalent contacts between residue side chains are the basis for protein folding, protein assembly, and PPI [232]. These contacts facilitate a variety of interactions and associations within and between proteins. Based on their diverse structural and functional characteristics, PPIs can be categorized in several ways [230]. On the basis of their interaction surface, they may be homo- or hetero-oligomeric; as judged by their stability, they may be obligate or nonobligate; and as measured by their persistence, they may be transient or permanent. A given PPI can fall into any combination of these three categorical pairs. An interaction may also require reclassification under certain

## 2 Introduction



**Figure 1–1** Number of results found in PubMed for the term “proteome.” (Reprinted from [200] with permission of John Wiley & Sons, Inc.)

conditions; for example, it may be mainly transient *in vivo* but become permanent under certain cellular conditions.

It has been observed that proteins seldom act as single isolated species while performing their functions *in vivo* [330]. The analysis of annotated proteins reveals that proteins involved in the same cellular processes often interact with each other [312]. The function of unknown proteins may be postulated on the basis of their interaction with a known protein target of known function. Mapping PPIs has not only provided insight into protein function but also facilitated the modeling of functional pathways to elucidate the molecular mechanisms of cellular processes. The study of PPIs is fundamental to understanding how proteins function within the cell. Characterizing the interactions of proteins in a given cellular proteome will be the next milestone along the road to understand the biochemistry of the cell.

The result of two or more proteins interacting with a specific functional objective can be demonstrated in several different ways. The measurable effects of PPIs have been outlined by Phizicky and Fields [254]. PPIs can:

- alter the kinetic properties of enzymes; this may be the result of subtle changes at the level of substrate binding or at the level of an allosteric effect;
- act as a common mechanism to allow for substrate channeling;
- create a new binding site, typically for small effector molecules;
- inactivate or destroy a protein; or
- change the specificity of a protein for its substrate through interaction with different binding partners; for example, demonstrate a new function that neither protein can exhibit alone.

PPIs are much more widespread than once suspected, and the degree of regulation that they confer is large. To properly understand their significance in the cell, one needs to identify the different interactions, understand the extent to which they take place in the cell, and determine the consequences of the interactions.

In recent years, PPI data have been enriched by high-throughput experimental methods, such as two-hybrid systems [155,307], mass spectrometry [113,144], and

protein chip technology [114,205,346]. Integrated PPI networks have been built from these heterogeneous data sources. However, the large volume of PPI data currently available has posed a challenge to experimental investigation. Computational analysis of PPI networks has become a necessary supplemental tool for understanding the functions of uncharacterized proteins.

## 1.2 COMPUTATIONAL ANALYSIS OF PPI NETWORKS

A PPI network can be described as a complex system of proteins linked by interactions. The computational analysis of PPI networks begins with the representation of the PPI network structure. The simplest representation takes the form of a mathematical graph consisting of nodes and edges [314]. Proteins are represented as nodes in such a graph; two proteins that interact physically are represented as adjacent nodes connected by an edge. Based on this graphic representation, various computational approaches, such as data mining, machine learning, and statistical approaches, can be designed to reveal the organization of PPI networks at different levels. An examination of the graphic form of the network can yield a variety of insights. For example, neighboring proteins in the graph are generally considered to share functions (“guilt by association”). Thus, the functions of a protein may be predicted by looking at the proteins with which it interacts and the protein complexes to which it belongs. In addition, densely connected subgraphs in the network are likely to form protein complexes that function as a unit in a certain biological process. An investigation of the topological features of the network (e.g., whether it is scale-free, a small network, or governed by the power law) can also enhance our understanding of the biological system [5].

In general, the computational analysis of PPI networks is challenging, with these major difficulties being commonly encountered:

- *The protein interactions are not reliable.* Large-scale experiments have yielded numerous false positives. For example, as reported in [288], high-throughput yeast two-hybrid (Y2H) assays are ~50% reliable. It is also likely that there are many false negatives in the PPI networks currently under study.
- *A protein can have several different functions.* A protein may be included in one or more functional groups. Therefore, overlapping clusters should be identified in the PPI networks. Since conventional clustering methods generally produce pairwise disjoint clusters, they may not be effective when applied to PPI networks.
- *Two proteins with different functions frequently interact with each other.* Such frequent, random connections between the proteins in different functional groups expand the topological complexity of the PPI networks, posing difficulties to the detection of unambiguous partitions.

Recent studies of complex systems [5,227] have attempted to understand and characterize the structural behaviors of such systems from a topological perspective. Such features as small-world properties [319], scale-free degree distributions [28,29], and hierarchical modularity [261] have been observed in complex systems, elements that are also characteristic of PPI networks. Therefore, topological methods can be

used to address the challenges mentioned earlier and to facilitate the efficient and accurate analysis of PPI networks.

### 1.2.1 Topological Features of PPI Networks

Barabasi and Oltvai [29] introduced the concept of degree distribution,  $P(k)$ , to quantify the probability that a selected node in a network will have exactly  $k$  links. Networks of different types can be distinguished by their degree distributions. For example, a random network follows a Poisson distribution. In contrast, a scale-free network has a power-law degree distribution,  $P(k) \sim k^{-\gamma}$ , indicating that a few hubs bind numerous small nodes. When  $2 \leq \gamma \leq 3$ , the hubs play a significant role in the network [29]. Recent publications have indicated that PPI networks have the features of a scale-free network [121,161,198,313]; therefore, their degree distribution approximates a power law,  $P(k) \sim k^{-\gamma}$ . In scale-free networks, most proteins participate in only a few interactions, while a small set of hubs participate in dozens of interactions.

PPI networks also have a characteristic property known as the “small-world effect,” which states that any two nodes can be connected via a short path of a few links. The small-world phenomenon was first investigated as a concept in sociology [217] and is a feature of a range of networks arising in both nature and technology, including the Internet [5], scientific collaboration networks [224], the English lexicon [280], metabolic networks [106], and PPI networks [284,313]. Although the small-world effect is a property of random networks, the path length in scale-free networks is much shorter than that predicted by the small-world effect [74,75]. Therefore, scale-free networks are “ultra-small.” This short path length indicates that local perturbations in metabolite concentrations could permeate an entire network very quickly. In PPI networks, highly connected nodes (hubs) seldom directly link to each other [211]. This differs from the assortative nature of social networks, in which well-connected individuals tend to have direct connections to each other. In contrast, biological networks have the property of disassortativity, in which highly connected nodes are only infrequently linked.

A number of recent publications have proposed the use of centrality indices, including node degree, pagerank, clustering coefficient, betweenness centrality, and bridging centrality metrics, as measurements of the importance of components in a network [47,53,103,110,226,268,319]. For instance, betweenness centrality [225] was proposed to detect the optimal location for partitioning a network [122,145]. The modified betweenness cut approach has been suggested for use with weighted PPI networks that integrate gene expression [61]. Jeong’s group has espoused the degree of a node as a key basis for the identification of essential network components [161]. In this model, power-law networks are very robust to random attacks but highly vulnerable to targeted attacks [7]. Hahn’s group identified differences in degree, betweenness, and closeness centrality between essential and nonessential genes in three eukaryotic PPI networks (yeast, worm, and fly) [131]. Estrada’s group introduced a new subgraph centrality measure to characterize the participation of each node in all subgraphs in a network [103,102]. Palumbo’s group sought to identify lethal nodes by arc deletion, thus facilitating the isolation of network subcomponents [239]. Guimera’s group devised a clustering method to identify functional modules

in metabolic pathways and categorized the role of each component in the pathway according to its topological location relative to detected functional modules [129].

As we will subsequently discuss in greater detail, the unique topological features found to be characteristic of PPI networks will play significant roles in the computational analysis of these networks.

### 1.2.2 Modularity Analysis

The idea of functional modules, introduced in [139], offers a major conceptual tool for the systematic analysis of a biological system. A functional module in a PPI network represents a maximal set of functionally associated proteins. In other words, it is composed of those proteins that are mutually involved in a given biological process or function. A wide range of graph-theoretic approaches have been employed to identify functional modules in PPI networks. However, these approaches have tended to be limited in accuracy due to the presence of unreliable interactions and the complex connectivity of the networks [288]. In particular, the topological complexity of PPI networks, arising from the overlapping patterns of modules and cross talks between modules, poses challenges to the identification of functional modules. Because a protein generally performs different biological processes or functions in different environments, real functional modules are overlapping. Moreover, the frequent, dynamic cross connections between different functions are biologically meaningful and must be taken into account [274].

In an attempt to parse this complexity, the hierarchical organization of modules in biological networks has been recently proposed [261]. The architecture of this model is based on a scale-free topology with embedded modularity. In this model, the significance of a few hub nodes is emphasized, and these nodes are viewed as the determinants of survival during network perturbations and as the essential backbone of the hierarchical structure. This hierarchical network model can plausibly be applied to PPI networks because cellular functionality is typically hierarchical in nature, and PPI networks include a few hub nodes that are biologically lethal.

The identification of functional modules in PPI networks or modularity analysis can be successfully accomplished through the use of cluster analysis. Cluster analysis is invaluable in elucidating network topological structure and the relationships among network components. Typically, clustering approaches focus on detecting densely connected subgraphs within the graphic representation of a PPI network. For example, the maximum clique algorithm [286] is used to detect fully connected, complete subgraphs. To compensate for the high-density threshold imposed by this algorithm, relatively dense subgraphs can be identified in lieu of complete subgraphs, either by using a density threshold or by optimizing an objective density function [56,286]. A number of density-based clustering algorithms using alternative density functions have been presented [12,24,247].

As noted, hierarchical clustering approaches can plausibly be applied to biological networks because of the hierarchical nature of functional modules [261,297]. These approaches iteratively merge nodes or recursively divide a graph into two or more subgraphs. To merge nodes iteratively, the similarity or distance between two nodes or two groups of nodes is measured and a pair is selected for merger in each iteration [17,263]. Recursive division of a graph involves the selection of nodes

or edges to be cut. Partition-based approaches have also been applied to biological networks. One partition-based clustering approach, the Restricted Neighborhood Search Clustering (RNSC) algorithm [180], determines the best partition using a cost function. In addition, other approaches have been applied to biological networks. For example, the Markov Clustering Algorithm (MCL) finds clusters using iterative rounds of expansion and inflation that, respectively, prefer the strongly connected regions and weaken the sparsely connected regions [308]. The line graph generation method [250] transforms a network of proteins connected by interactions into a network of connected interactions and then uses the MCL algorithm to cluster the PPI network. Samantha and Liang [272] applied a statistical approach to the clustering of proteins based on the premise that a pair of proteins sharing a significantly greater number of common neighbors will have a high functional similarity. The recently introduced STM algorithm [148] votes a representative of a cluster for each node.

Topological metrics can be incorporated into the modularity analysis of PPI networks. From our studies, we have observed that the bridging nodes identified in PPI networks serve as the connecting nodes between protein modules; therefore, removing the bridging nodes preserves the structural integrity of the network. Such findings can play an important role in the modularity analysis of PPI networks. Removal of the bridging nodes yields a set of components disconnected from the network. Thus, using bridging centrality to remove the bridging nodes can be an excellent preprocessing procedure to estimate the number and location of modules in the PPI network. Results of this research [151,152] have shown that such approaches can generate larger modules that discard fewer proteins, permitting more accurate functional detection than other current methods.

### **1.2.3 Prediction of Protein Functions in PPI Networks**

Predicting protein function can be, in itself, the ultimate objective of the analysis of a PPI network. Despite the many extensive studies of yeast that have been undertaken, there are still a number of functionally uncharacterized proteins in the yeast database. The functional annotation of human proteins can provide a strong foundation for the complete understanding of cell mechanisms, information that is invaluable for drug discovery and development. The increased interest in and availability of PPI networks have catalyzed the development of computational methods to elucidate protein functions.

Protein functions may be predicted on the basis of modularization algorithms. If an unknown protein is included in a functional module, it is expected to contribute toward the function that the module represents. The generated functional modules may thus provide a framework within which to predict the functions of unknown proteins. Each generated module may contain a few uncharacterized proteins along with a larger number of known proteins. It can be assumed that the unknown proteins play a positive role in realizing the function of the generated module. However, predictions arrived at through these means may be inaccurate, since the accuracy of the modularization process itself is typically low. For greater reliability, protein functions should be predicted directly from the topology or connectivity of PPI networks.

Several topology-based approaches that predict protein function on the basis of PPI networks have been introduced. At the simplest level, the “neighbor counting

method” predicts the function of an unknown protein by the frequency of known functions of the immediate neighbor proteins [274]. The majority of functions of the immediate neighbors can be statistically assessed [143]. The function of a protein can be assumed to be independent of all other proteins, given the functions of its immediate neighbors. This assumption gives rise to a Markov random field model [85,196]. Recently, the number of common neighbors of the known protein and the unknown protein has been taken as the basis for the prediction of function [201].

Machine learning has been widely applied to the analysis of PPI networks, and, in particular, to the prediction of protein functions. A variety of methods have been developed to predict protein function on the basis of different information sources. Some of the inputs used by these methods include protein structure and sequence, protein domain, PPIs, genetic interactions, and gene expression analysis. The accuracy of prediction can be enhanced by drawing upon multiple sources of information. The Gene Ontology (GO) database [84] is one example of such semantic integration.

#### 1.2.4 Integration of Domain Knowledge

As noted, the accuracy of results obtained from computational approaches can be compromised by the inclusion of false connections and the high complexity of networks. The reliability of this process can be improved by the integration of other functional information. Initially, the identification of similarities in gene sequence can be a primary indicator of a functional association between two genes. Additionally, genome-level methods for functional inference, such as gene fusion events and phylogenetic profiling, can generate useful data pointing to functional linkages. Beyond this, we know that genes with correlated expression profiles determined through microarray experiments are likely to be functionally related. Many studies [65,66,153,304] have investigated the integration of PPI networks with gene expression data to improve the accuracy of the functional modules identified. Finally, as briefly noted earlier, GO [18,301] can be a useful data source to combine with the PPI networks. GO is currently one of the most comprehensive and well-curated ontology databases in the bioinformatics community. It represents a collaborative effort to address the need for consistent descriptions of genes and gene products. The GO database includes GO terms and their relationships. The former are well-defined biological terms organized into three general conceptual categories that are shared across different organisms: biological processes, molecular functions, and cellular components. The GO database also provides annotations to each GO term, and each gene can be annotated on one or more GO terms. The GO database and its annotations can thus be a significant resource for the discovery of functional knowledge. These tools have been employed to facilitate the analysis of gene expression data [89,105,147] and have been integrated with unreliable PPI networks to accurately predict functions of unknown proteins [84] and identify functional modules [68,70].

### 1.3 SIGNIFICANT APPLICATIONS

The systematic analysis of PPIs can enable a better understanding of cellular organization, processes, and functions. Functional modules can be identified from the

PPI networks that have been derived from experimental data sets. There are many significant applications following this analysis. In this book, the following principal applications to which this analysis can be applied will be discussed:

- *Predicting protein function.* As noted earlier, the most basic application of PPI networks is the use of topological analysis to predict protein function. The generated functional modules can serve as a framework within which to predict the functions of unknown proteins. Each generated module may contain a few uncharacterized proteins. By associating unknown proteins with the known proteins, we can suggest that those proteins participate positively in performing the functions assigned to the modules.
- *Lethality analysis.* The topological analysis of PPI networks can be used to systematically assess the biological importance of bridging and other nodes in a PPI network [65,66,70,148]. Lethality, a crucial factor in characterizing the biological indispensability of a protein, is determined by examining whether a module is functionally disrupted when the protein is eliminated. Information regarding lethality is compiled in most PPI databases. For example, the MIPS database [214] indicates the lethality or viability of each included protein. Such sources allow the researcher to compare the lethality of nodes with high bridging-score values to that associated with other competing network parameters in the PPI networks. These comparisons reveal that nodes with the highest bridging scores are less lethal than both randomly selected nodes and nodes with high degree centrality. However, the average lethality of the neighbors of the nodes with the highest bridging scores is greater than that of a randomly selected subset. Our research has indicated that bridging nodes have relatively low lethality; inter-connecting nodes are characterized by higher lethality; and modular nodes and peripheral nodes have, respectively, the highest and lowest proportion of lethal proteins. These results imply that many of the bridging nodes do not perform tasks critical to biological functions [151,152]. As a result, these nodes would serve as good targets for drugs, as discussed later.
- *Assessing the druggability of molecular targets from network topology.* Translating the societal investments in the Human Genome Project and other similar large-scale efforts into therapies for human diseases is an important scientific imperative in the post-human-genome era. The efficacy, specificity/selectivity, and side-effect characteristics of well-designed drugs depend largely on the appropriate choice of pharmacological target. For this reason, the identification of molecular targets is a very early and critical step in the drug discovery and development process. The goal of the target identification process is to arrive at a very limited subset of biological molecules that will become the principal focus for the subsequent discovery research, development, and clinical trials. Pharmacological targets can span the range of biological molecules from DNA and lipids to metabolites. In fact, though, the majority of pharmacological targets are proteins. Effective pharmacological intervention with the target protein should significantly impact the key molecular processes in which the protein participates, and the resultant perturbation should be successful in modulating the pathophysiological process of interest. Another important consideration that is sometimes overlooked during the target identification step is the potential for side effects.

Ideally, an appropriate balance should be found among efficacy, selectivity, and side effects. In practice, however, compromises are often required in the areas of specificity/selectivity and side effects, since pharmacological interventions with proteins that are central to key processes will likely affect many biological pathways. We have observed that the biological correlates of the nodes with the highest bridging scores indicate that these nodes are less lethal than other nodes in PPI networks. Thus, they are promising drug targets from the standpoints of efficacy and side effects.

## 1.4 ORGANIZATION OF THIS BOOK

This book is intended to provide an in-depth examination of computational analysis as applied to PPI networks, offering perspectives from data mining, machine learning, graph theory, and statistics. The remainder of this book is organized as follows:

- Chapter 2 introduces the three principal experimental approaches that are currently used for generating PPI data: the Y2H system [121,156,307], mass spectrometry (MS) [113,120,144,187,210,303], and protein microarray methods [114,346].
- Chapter 3 discusses various computational approaches to the prediction of protein interactions, including genomic-scale, sequence-based, structure-based, learning-sequence-based, and network topology-based techniques.
- Chapter 4 introduces the basic properties of and metrics applied to PPI networks. Basic concepts in graphic representation employed to characterize various properties of PPI networks are defined for use throughout the balance of the book.
- Chapter 5 discusses the modularity analysis of PPI networks. Various modularity analysis algorithms used to identify modules in PPI networks are discussed, and an overview of the validation methods for modularity analysis is presented.
- Chapter 6 explores the topological analysis of PPI networks. Various metrics used for assessing specific topological features of PPI networks are presented and discussed.
- Chapter 7 focuses on greater detail on one type of modularity algorithm, specifically, the distance-based modularity analysis of PPI networks.
- Chapter 8 focuses on greater detail on graph-theoretic approaches for modularity analysis of PPI networks.
- Chapter 9 discusses the flow-based analysis of PPI networks.
- Chapter 10 examines statistical- and machine learning-based analysis of PPI networks.
- Chapter 11 discusses the integration of domain knowledge into the analysis of PPI networks.
- Chapter 12 presents some of the more recent approaches that have been developed for incorporating diverse biological information into the explorative analysis of PPI networks.
- Chapter 13 offers a synthesis of the methods and concepts discussed throughout the book and reflections on potential directions for future research and applications.

## **1.5 SUMMARY**

The analysis of PPI networks poses many challenges, given the inherent complexity of these networks, the high noise level characteristic of the data, and the presence of unusual topological phenomena. As discussed in this chapter, effective approaches are required to analyze PPI data and the resulting PPI networks. Recently, a variety of data-mining and statistical techniques have been applied to this end, with varying degrees of success. This book is intended to provide researchers with a working knowledge of many of the advanced approaches currently available for this purpose. (Some of the material in this chapter is reprinted from [200] with permission of John Wiley & Sons, Inc.)

# Experimental Approaches to Generation of Protein–Protein Interaction Data

## 2.1 INTRODUCTION

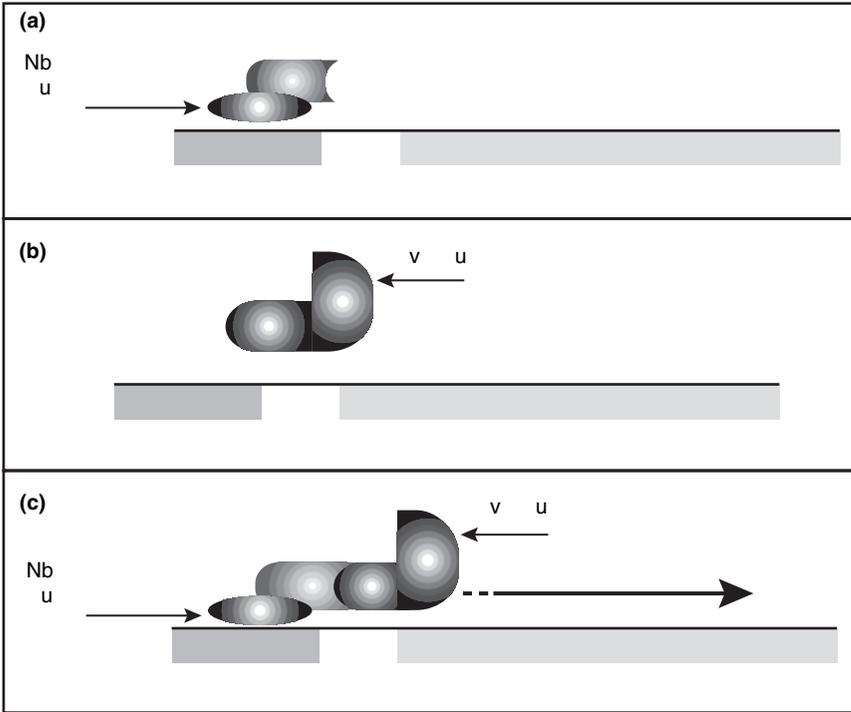
Proteins and their interactions lie at the heart of most fundamental biological processes. Typically, proteins seldom act in isolation but rather execute their functions through interaction with other biomolecular units. Consequently, an examination of these protein–protein interactions (PPIs) is essential to understanding the molecular mechanisms of underlying biological processes [79]. This chapter is intended to provide an overview of the more common experimental methods currently used to generate PPI data.

In the past, PPIs were typically examined via intensive small-scale investigations of restricted sets of proteins of interest, each yielding information regarding a limited number of PPIs. The existing databases of PPIs have been compiled from the results of such small-scale screens presented in individual research papers. Since these data are subject to stringent controls and evaluation in the peer-review process, they can be considered to be fairly reliable. However, each experiment observes only a few interactions and provides a data set of limited size.

Recent high-throughput approaches involve genome-wide detection of protein interactions. Studies using the yeast two-hybrid (Y2H) system [121,156,307], mass spectrometry (MS) [113,120,144,187,210,303], and protein microarrays [114,346] have generated large amounts of interaction data. The Y2H system takes a bottom-up genomic approach to detecting possible binary interactions between any two proteins encoded in the genome of interest. In contrast, mass spectrometric analysis adopts a top-down proteomic approach by analyzing the composition of protein complexes. The protein microarray technology simultaneously captures the expression of thousands of proteins.

## 2.2 THE Y2H SYSTEM

One of the most common approaches to the detection of pairs of interacting proteins *in vivo* is the Y2H system [21,155]. The Y2H system, first introduced in 1989 [107], is a molecular–genetic tool that facilitates the study of PPI. The interaction of two proteins transcriptionally activates a reporter gene, and a color reaction is seen



**Figure 2-1** Y2H system applied to the detection of binary protein interactions. (Reprinted by permission from Macmillan Publishers Ltd: Nature [233], copyright 2000.)

on specific media. This indication can track the interaction between two proteins, revealing “prey” proteins that interact with a known “bait” protein.

Two-hybrid procedures are typically carried out by screening a protein of interest against a random library of potential protein partners. Figure 2-1 [233] depicts the Y2H process. In Figure 2-1(a), we see that the fusion of the “bait” protein and the DNA-binding domain of the transcriptional activator does not turn on the reporter gene; no color change occurs; and the interaction cannot be tracked. Figure 2-1(b) shows that, similarly, the fusion of the “prey” protein and the activating region of the transcriptional activator is also insufficient to switch on the reporter gene. In Figure 2-1(c), the “bait” and the “prey” associate, bringing the DNA-binding domain and activator region into sufficiently close proximity to switch on the reporter gene. The result is gene transcription and a color change that can be monitored.

The Y2H system enables both highly sensitive detection of PPIs and screening of genome libraries to ascertain the interaction partners of certain proteins. The system can also be used to pinpoint protein regions mediating the interactions [157]. However, the classic Y2H system has several limitations. First, it cannot, by definition, detect interactions involving three or more proteins and those depending on posttranslational modifications (PTMs) except those applied to the budding yeast itself [157]. Second, since some proteins (e.g., membrane proteins) cannot be reconstructed in the nucleus, the Y2H system is not suitable for the detection of interactions involving these proteins. Finally, the method does not guarantee that an interaction

indicated by Y2H actually takes place physiologically. Given these limitations, the Y2H system is most suitable for the detection of binary interactions, particularly those that are transient and unstable.

Despite these drawbacks, the Y2H system has become established as a standard technique in molecular biology and serves as an important method for proteomics analysis [240]. High-throughput Y2H screens have been applied to *Escherichia coli* [31], hepatitis C virus [108], Vaccinia virus [213], *Saccharomyces cerevisiae* [156,307], *Helicobacter pylori* [259], and *Caenorhabditis elegans* [198,315], *Drosophila melanogaster* [121], and *Homo sapiens* [76,266].

Recently, numerous modifications of the Y2H approach have been proposed that characterize PPI networks by screening each protein expressed in a eukaryotic cell [109]. Drees [92] has proposed a variant that includes the genetic information of a third protein. Zhang et al. [342] have suggested the use of RNA for the investigation of RNA–protein interactions. Vidal et al. [311] used the *URA3* gene instead of *GAL4* as the reporter gene; this two-hybrid system can be used to screen for ligand inhibition or to dissociate such complexes. Johnson and Varshavsky [166] have proposed a cytoplasmic two-hybrid system that can be used for screening of membrane protein interactions.

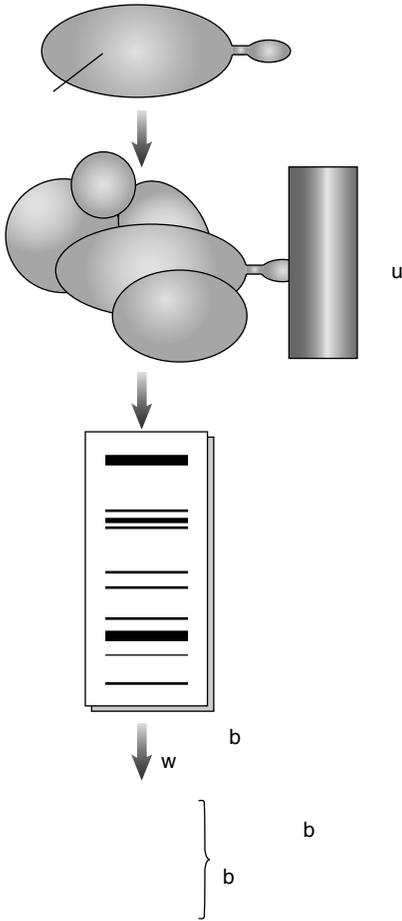
Despite the various limitations of the Y2H system, this approach has revealed a wealth of novel interactions and has helped illuminate the magnitude of the protein interactome. In principle, it could be applied in a more comprehensive fashion to examine all possible binary combinations between the proteins encoded by any single genome.

## 2.3 MASS SPECTROMETRY (MS) APPROACHES

Another traditional approach to PPI detection uses quantitative MS to analyze the composition of a partially purified protein complex together with a control purification in which the complex of interest is not enriched.

Mass spectrometry analysis proceeds in three steps: bait presentation, affinity purification of the complex, and analysis of the bound proteins [2]. Two large-scale studies [113,144], that apply MS analysis to the PPI network in yeast have been published. Each study attempted to identify all the components that were present in “naturally generated” protein complexes, taking as their subject essentially pure preparations of each complex [188]. In both approaches, bait proteins were generated that carried a particular affinity tag. In the case studied by Gavin et al. [113], 1,739 TAP-tagged (Tandem Affinity Purification) genes were introduced into the yeast genome by homologous recombination. Ho et al. [144] expressed 725 proteins modified to carry the FLAG epitope. In both cases, the proteins were expressed in yeast cells, and complexes were purified using a single immunoaffinity purification step. Both groups resolved the components of each purified complex with a one-dimensional denaturing polyacrylamide gel electrophoresis (PAGE) step. From the 1,167 yeast strains generated by Gavin et al. [113], 589 protein complexes were purified, 232 of which were unique. Ho et al. [144] used 725 protein baits and detected 3,617 interactions that involved 1,578 different proteins.

Figure 2–2 illustrates the process of mass spectrometric analysis [188]. In step (1), an “affinity tag” is attached to a target protein (the “bait”). As illustrated in



**Figure 2–2** Mass spectrometric analysis of protein complexes. (Reprinted by permission from Macmillan Publishers Ltd: Nature [188], copyright 2002.)

Figure 2–2(2), bait proteins are systematically precipitated, along with any associated proteins, onto an “affinity column.” In Figure 2–2(3), purified protein complexes are resolved by one-dimensional SDS-PAGE, so that proteins become separated according to mass. Step (4) entails the separating of protein bands by protein size; in step (5), protein bands are digested with trypsin. In steps (6–9), component proteins are detected by MS and bioinformatic analysis.

Mass-spectrometry-based proteomics can be applied not only to identify and quantify individual proteins [77,189,249,318] but also to protein analysis, including protein profiling [192], PTMs [206,207], and, in particular, identification of PPIs.

In general, mass spectrometric analysis is more physiological than the Y2H system. Actual molecular assemblies composed of all combinations of direct and cooperative interactions are analyzed *in vivo*, as opposed to the examination of reconstituted bimolecular interactions *ex vivo* or *in vitro*. MS can detect more complex interactions and is not limited to binary interactions, permitting the isolation of large protein complexes and the detection of networks of interactions. However,