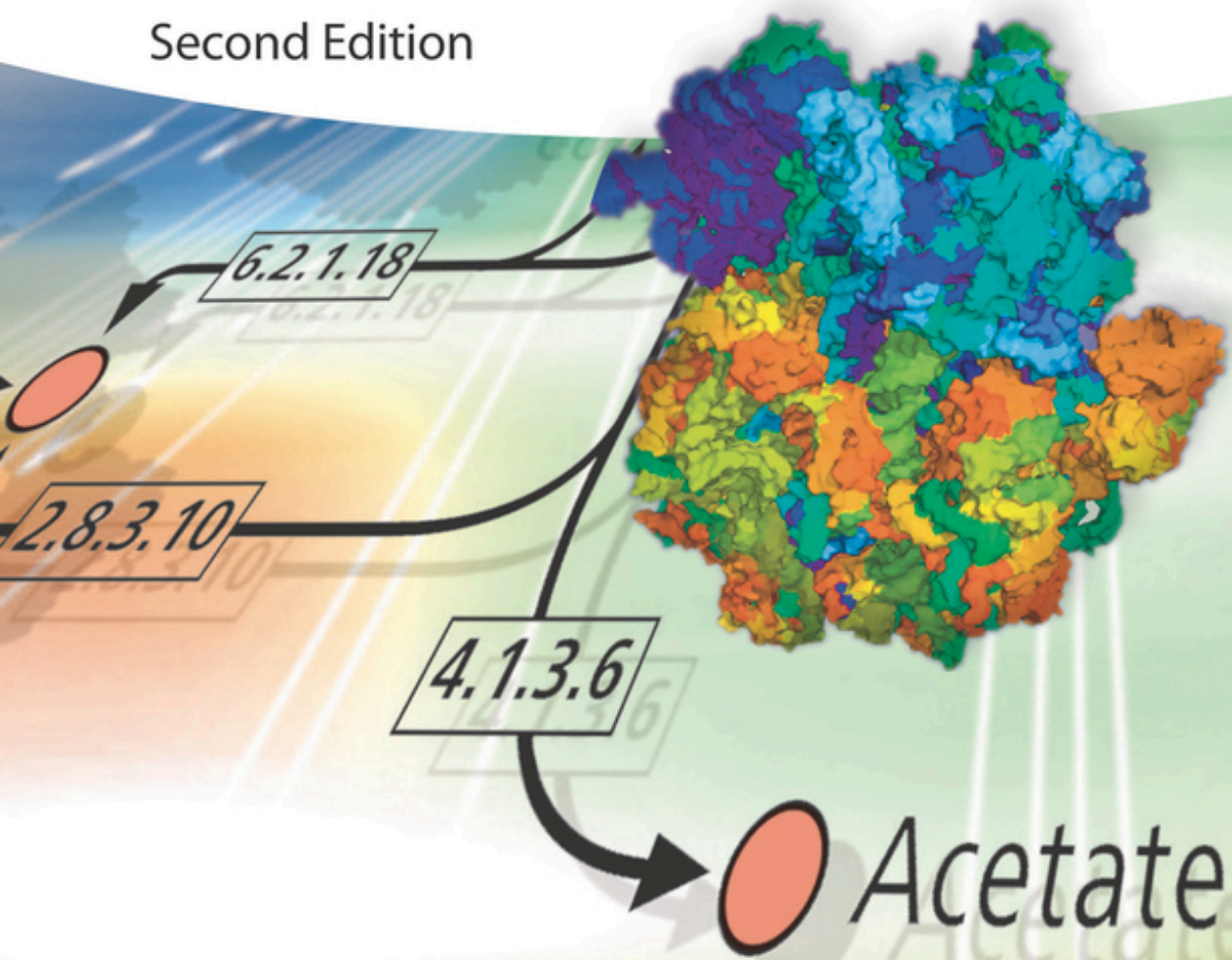


Volkhard Helms

# Principles of Computational Cell Biology

From Protein Complexes to Cellular Networks

Second Edition





## **Principles of Computational Cell Biology**



# **Principles of Computational Cell Biology**

From Protein Complexes to Cellular Networks

*Volkhard Helms*

Second Edition

**WILEY-VCH**

**Author**

*Volkhard Helms*

Universität des Saarlandes  
Zentrum für Bioinformatik  
66041 Saarbrücken  
Germany

■ All books published by **Wiley-VCH** are carefully produced. Nevertheless, authors, editors, and publisher do not warrant the information contained in these books, including this book, to be free of errors. Readers are advised to keep in mind that statements, data, illustrations, procedural details or other items may inadvertently be inaccurate.

**Library of Congress Card No.:**  
applied for

**British Library Cataloguing-in-Publication Data**

A catalogue record for this book is available from the British Library.

**Bibliographic information published by the Deutsche Nationalbibliothek**

The Deutsche Nationalbibliothek lists this publication in the Deutsche Nationalbibliografie; detailed bibliographic data are available on the Internet at <<http://dnb.d-nb.de>>.

© 2019 Wiley-VCH Verlag GmbH & Co. KGaA, Boschstr. 12, 69469 Weinheim, Germany

All rights reserved (including those of translation into other languages). No part of this book may be reproduced in any form – by photoprinting, microfilm, or any other means – nor transmitted or translated into a machine language without written permission from the publishers. Registered names, trademarks, etc. used in this book, even when not specifically marked as such, are not to be considered unprotected by law.

**Print ISBN:** 978-3-527-33358-5

**ePDF ISBN:** 978-3-527-81033-8

**ePub ISBN:** 978-3-527-81032-1

**Typesetting** SPi Global, Chennai, India  
**Printing and Binding**

Printed on acid-free paper

10 9 8 7 6 5 4 3 2 1

## Contents

**Preface of the First Edition** *xv*  
**Preface of the Second Edition** *xvii*

<b>1</b>	<b>Networks in Biological Cells</b>	<b>1</b>
1.1	Some Basics About Networks	1
1.1.1	Random Networks	2
1.1.2	Small-World Phenomenon	2
1.1.3	Scale-Free Networks	3
1.2	Biological Background	4
1.2.1	Transcriptional Regulation	5
1.2.2	Cellular Components	5
1.2.3	Spatial Organization of Eukaryotic Cells into Compartments	7
1.2.4	Considered Organisms	8
1.3	Cellular Pathways	8
1.3.1	Biochemical Pathways	8
1.3.2	Enzymatic Reactions	11
1.3.3	Signal Transduction	11
1.3.4	Cell Cycle	12
1.4	Ontologies and Databases	12
1.4.1	Ontologies	12
1.4.2	Gene Ontology	13
1.4.3	Kyoto Encyclopedia of Genes and Genomes	13
1.4.4	Reactome	13
1.4.5	Brenda	14
1.4.6	DAVID	14
1.4.7	Protein Data Bank	15
1.4.8	Systems Biology Markup Language	15
1.5	Methods for Cellular Modeling	17
1.6	Summary	17
1.7	Problems	17
	Bibliography	18

<b>2</b>	<b>Structures of Protein Complexes and Subcellular Structures</b>	<b>21</b>
2.1	Examples of Protein Complexes	22
2.1.1	Principles of Protein–Protein Interactions	24
2.1.2	Categories of Protein Complexes	27
2.2	Complexome: The Ensemble of Protein Complexes	28
2.2.1	Complexome of <i>Saccharomyces cerevisiae</i>	28
2.2.2	Bacterial Protein Complexomes	30
2.2.3	Complexome of Human	31
2.3	Experimental Determination of Three-Dimensional Structures of Protein Complexes	31
2.3.1	X-ray Crystallography	32
2.3.2	NMR	34
2.3.3	Electron Crystallography/Electron Microscopy	34
2.3.4	Cryo-EM	34
2.3.5	Immunolectron Microscopy	35
2.3.6	Fluorescence Resonance Energy Transfer	35
2.3.7	Mass Spectroscopy	36
2.4	Density Fitting	38
2.4.1	Correlation-Based Density Fitting	38
2.5	Fourier Transformation	40
2.5.1	Fourier Series	40
2.5.2	Continuous Fourier Transform	41
2.5.3	Discrete Fourier Transform	41
2.5.4	Convolution Theorem	41
2.5.5	Fast Fourier Transformation	42
2.6	Advanced Density Fitting	44
2.6.1	Laplacian Filter	45
2.7	FFT Protein–Protein Docking	46
2.8	Protein–Protein Docking Using Geometric Hashing	48
2.9	Prediction of Assemblies from Pairwise Docking	49
2.9.1	CombDock	49
2.9.2	Multi-LZerD	52
2.9.3	3D-MOSAIC	52
2.10	Electron Tomography	53
2.10.1	Reconstruction of Phantom Cell	55
2.10.2	Protein Complexes in <i>Mycoplasma pneumoniae</i>	55
2.11	Summary	56
2.12	Problems	57
2.12.1	Mapping of Crystal Structures into EM Maps	57
	Bibliography	60
<b>3</b>	<b>Analysis of Protein–Protein Binding</b>	<b>63</b>
3.1	Modeling by Homology	63
3.2	Properties of Protein–Protein Interfaces	66
3.2.1	Size and Shape	66
3.2.2	Composition of Binding Interfaces	68

3.2.3	Hot Spots	69
3.2.4	Physicochemical Properties of Protein Interfaces	71
3.2.5	Predicting Binding Affinities of Protein–Protein Complexes	72
3.2.6	Forces Important for Biomolecular Association	73
3.3	Predicting Protein–Protein Interactions	75
3.3.1	Pairing Propensities	75
3.3.2	Statistical Potentials for Amino Acid Pairs	78
3.3.3	Conservation at Protein Interfaces	79
3.3.4	Correlated Mutations at Protein Interfaces	83
3.4	Summary	86
3.5	Problems	86
	Bibliography	86
<b>4</b>	<b>Algorithms on Mathematical Graphs</b>	<b>89</b>
4.1	Primer on Mathematical Graphs	89
4.2	A Few Words About Algorithms and Computer Programs	90
4.2.1	Implementation of Algorithms	91
4.2.2	Classes of Algorithms	92
4.3	Data Structures for Graphs	93
4.4	Dijkstra’s Algorithm	95
4.4.1	Description of the Algorithm	96
4.4.2	Pseudocode	100
4.4.3	Running Time	101
4.5	Minimum Spanning Tree	101
4.5.1	Kruskal’s Algorithm	102
4.6	Graph Drawing	102
4.7	Summary	104
4.8	Problems	105
4.8.1	Force Directed Layout of Graphs	107
	Bibliography	110
<b>5</b>	<b>Protein–Protein Interaction Networks – Pairwise Connectivity</b>	<b>111</b>
5.1	Experimental High-Throughput Methods for Detecting Protein–Protein Interactions	111
5.1.1	Gel Electrophoresis	112
5.1.2	Two-Dimensional Gel Electrophoresis	112
5.1.3	Affinity Chromatography	113
5.1.4	Yeast Two-hybrid Screening	114
5.1.5	Synthetic Lethality	115
5.1.6	Gene Coexpression	116
5.1.7	Databases for Interaction Networks	116
5.1.8	Overlap of Interactions	116
5.1.9	Criteria to Judge the Reliability of Interaction Data	118
5.2	Bioinformatic Prediction of Protein–Protein Interactions	120
5.2.1	Analysis of Gene Order	121
5.2.2	Phylogenetic Profiling/Coevolutionary Profiling	121

5.2.2.1	Coevolution	122
5.3	Bayesian Networks for Judging the Accuracy of Interactions	124
5.3.1	Bayes' Theorem	125
5.3.2	Bayesian Network	125
5.3.3	Application of Bayesian Networks to Protein–Protein Interaction Data	126
5.3.3.1	Measurement of Reliability “Likelihood Ratio”	127
5.3.3.2	Prior and Posterior Odds	127
5.3.3.3	A Worked Example: Parameters of the Naïve Bayesian Network for Essentiality	128
5.3.3.4	Fully Connected Experimental Network	129
5.4	Protein Interaction Networks	131
5.4.1	Protein Interaction Network of <i>Saccharomyces cerevisiae</i>	131
5.4.2	Protein Interaction Network of <i>Escherichia coli</i>	131
5.4.3	Protein Interaction Network of Human	132
5.5	Protein Domain Networks	132
5.6	Summary	135
5.7	Problems	136
5.7.1	Bayesian Analysis of (Fake) Protein Complexes	136
	Bibliography	138
<b>6</b>	<b>Protein–Protein Interaction Networks – Structural Hierarchies</b>	<b>141</b>
6.1	Protein Interaction Graph Networks	141
6.1.1	Degree Distribution	141
6.1.2	Clustering Coefficient	143
6.2	Finding Cliques	145
6.3	Random Graphs	146
6.4	Scale-Free Graphs	147
6.5	Detecting Communities in Networks	149
6.5.1	Divisive Algorithms for Mapping onto Tree	153
6.6	Modular Decomposition	155
6.6.1	Modular Decomposition of Graphs	157
6.7	Identification of Protein Complexes	161
6.7.1	MCODE	161
6.7.2	ClusterONE	162
6.7.3	DACO	163
6.7.4	Analysis of Target Gene Coexpression	164
6.8	Network Growth Mechanisms	165
6.9	Summary	169
6.10	Problems	169
	Bibliography	178
<b>7</b>	<b>Protein–DNA Interactions</b>	<b>181</b>
7.1	Transcription Factors	181
7.2	Transcription Factor-Binding Sites	183

7.3	Experimental Detection of TFBS	183
7.3.1	Electrophoretic Mobility Shift Assay	183
7.3.2	DNase Footprinting	184
7.3.3	Protein-Binding Microarrays	185
7.3.4	Chromatin Immunoprecipitation Assays	187
7.4	Position-Specific Scoring Matrices	187
7.5	Binding Free Energy Models	189
7.6	<i>Cis</i> -Regulatory Motifs	191
7.6.1	DACO Algorithm	192
7.7	Relating Gene Expression to Binding of Transcription Factors	192
7.8	Summary	194
7.9	Problems	194
	Bibliography	195
<b>8</b>	<b>Gene Expression and Protein Synthesis</b>	<b>197</b>
8.1	Regulation of Gene Transcription at Promoters	197
8.2	Experimental Analysis of Gene Expression	198
8.2.1	Real-time Polymerase Chain Reaction	199
8.2.2	Microarray Analysis	199
8.2.3	RNA-seq	201
8.3	Statistics Primer	201
8.3.1	<i>t</i> -Test	203
8.3.2	<i>z</i> -Score	203
8.3.3	Fisher's Exact Test	203
8.3.4	Mann–Whitney–Wilcoxon Rank Sum Tests	205
8.3.5	Kolmogorov–Smirnov Test	206
8.3.6	Hypergeometric Test	206
8.3.7	Multiple Testing Correction	207
8.4	Preprocessing of Data	207
8.4.1	Removal of Outlier Genes	207
8.4.2	Quantile Normalization	208
8.4.3	Log Transformation	208
8.5	Differential Expression Analysis	209
8.5.1	Volcano Plot	210
8.5.2	SAM Analysis of Microarray Data	210
8.5.3	Differential Expression Analysis of RNA-seq Data	212
8.5.3.1	Negative Binomial Distribution	213
8.5.3.2	DESeq	213
8.6	Gene Ontology	214
8.6.1	Functional Enrichment	216
8.7	Similarity of GO Terms	217
8.8	Translation of Proteins	217
8.8.1	Transcription and Translation Dynamics	218
8.9	Summary	219
8.10	Problems	220
	Bibliography	224

<b>9</b>	<b>Gene Regulatory Networks</b>	227
9.1	Gene Regulatory Networks (GRNs)	228
9.1.1	Gene Regulatory Network of <i>E. coli</i>	228
9.1.2	Gene Regulatory Network of <i>S. cerevisiae</i>	231
9.2	Graph Theoretical Models	231
9.2.1	Coexpression Networks	232
9.2.2	Bayesian Networks	233
9.3	Dynamic Models	234
9.3.1	Boolean Networks	234
9.3.2	Reverse Engineering Boolean Networks	235
9.3.3	Differential Equations Models	236
9.4	DREAM: Dialogue on Reverse Engineering Assessment and Methods	238
9.4.1	Input Function	239
9.4.2	YAYG Approach in DREAM3 Contest	240
9.5	Regulatory Motifs	244
9.5.1	Feed-forward Loop (FFL)	245
9.5.2	SIM	245
9.5.3	Densely Overlapping Region (DOR)	246
9.6	Algorithms on Gene Regulatory Networks	247
9.6.1	Key-pathway Miner Algorithm	247
9.6.2	Identifying Sets of Dominating Nodes	248
9.6.3	Minimum Dominating Set	249
9.6.4	Minimum Connected Dominating Set	249
9.7	Summary	250
9.8	Problems	251
	Bibliography	254
<b>10</b>	<b>Regulatory Noncoding RNA</b>	257
10.1	Introduction to RNAs	257
10.2	Elements of RNA Interference: siRNAs and miRNAs	259
10.3	miRNA Targets	261
10.4	Predicting miRNA Targets	264
10.5	Role of TFs and miRNAs in Gene-Regulatory Networks	264
10.6	Constructing TF/miRNA Coregulatory Networks	266
10.6.1	TFmiR Web Service	267
10.6.1.1	Construction of Candidate TF–miRNA–Gene FFLs	268
10.6.1.2	Case Study	269
10.7	Summary	270
	Bibliography	270
<b>11</b>	<b>Computational Epigenetics</b>	273
11.1	Epigenetic Modifications	273
11.1.1	DNA Methylation	273
11.1.1.1	CpG Islands	276
11.1.2	Histone Marks	277
11.1.3	Chromatin-Regulating Enzymes	278

11.1.4	Measuring DNA Methylation Levels and Histone Marks Experimentally	279
11.2	Working with Epigenetic Data	281
11.2.1	Processing of DNA Methylation Data	281
11.2.1.1	Imputation of Missing Values	281
11.2.1.2	Smoothing of DNA Methylation Data	281
11.2.2	Differential Methylation Analysis	282
11.2.3	Comethylation Analysis	283
11.2.4	Working with Data on Histone Marks	285
11.3	Chromatin States	286
11.3.1	Measuring Chromatin States	286
11.3.2	Connecting Epigenetic Marks and Gene Expression by Linear Models	287
11.3.3	Markov Models and Hidden Markov Models	288
11.3.4	Architecture of a Hidden Markov Model	290
11.3.5	Elements of an HMM	291
11.4	The Role of Epigenetics in Cellular Differentiation and Reprogramming	292
11.4.1	Short History of Stem Cell Research	293
11.4.2	Developmental Gene Regulatory Networks	293
11.5	The Role of Epigenetics in Cancer and Complex Diseases	295
11.6	Summary	296
11.7	Problems	296
	Bibliography	301
<b>12</b>	<b>Metabolic Networks</b>	<b>303</b>
12.1	Introduction	303
12.2	Resources on Metabolic Network Representations	306
12.3	Stoichiometric Matrix	308
12.4	Linear Algebra Primer	309
12.4.1	Matrices: Definitions and Notations	309
12.4.2	Adding, Subtracting, and Multiplying Matrices	310
12.4.3	Linear Transformations, Ranks, and Transpose	311
12.4.4	Square Matrices and Matrix Inversion	311
12.4.5	Eigenvalues of Matrices	312
12.4.6	Systems of Linear Equations	313
12.5	Flux Balance Analysis	314
12.5.1	Gene Knockouts: MOMA Algorithm	316
12.5.2	OptKnock Algorithm	318
12.6	Double Description Method	319
12.7	Extreme Pathways and Elementary Modes	324
12.7.1	Steps of the Extreme Pathway Algorithm	324
12.7.2	Analysis of Extreme Pathways	328
12.7.3	Elementary Flux Modes	329
12.7.4	Pruning Metabolic Networks: NetworkReducer	331
12.8	Minimal Cut Sets	332
12.8.1	Applications of Minimal Cut Sets	337

12.9	High-Flux Backbone	339
12.10	Summary	341
12.11	Problems	341
12.11.1	Static Network Properties: Pathways	341
	Bibliography	346
<b>13</b>	<b>Kinetic Modeling of Cellular Processes</b>	<b>349</b>
13.1	Biological Oscillators	349
13.2	Circadian Clocks	350
13.2.1	Role of Post-transcriptional Modifications	352
13.3	Ordinary Differential Equation Models	353
13.3.1	Examples for ODEs	354
13.4	Modeling Cellular Feedback Loops by ODEs	356
13.4.1	Protein Synthesis and Degradation: Linear Response	356
13.4.2	Phosphorylation/Dephosphorylation – Hyperbolic Response	357
13.4.3	Phosphorylation/Dephosphorylation – Buzzer	359
13.4.4	Perfect Adaptation – Sniffer	360
13.4.5	Positive Feedback – One-Way Switch	361
13.4.6	Mutual Inhibition – Toggle Switch	362
13.4.7	Negative Feedback – Homeostasis	362
13.4.8	Negative Feedback: Oscillatory Response	364
13.4.9	Cell Cycle Control System	365
13.5	Partial Differential Equations	366
13.5.1	Spatial Gradients of Signaling Activities	368
13.5.2	Reaction–Diffusion Systems	368
13.6	Dynamic Phosphorylation of Proteins	369
13.7	Summary	370
13.8	Problems	372
	Bibliography	373
<b>14</b>	<b>Stochastic Processes in Biological Cells</b>	<b>375</b>
14.1	Stochastic Processes	375
14.1.1	Binomial Distribution	376
14.1.2	Poisson Process	377
14.1.3	Master Equation	377
14.2	Dynamic Monte Carlo (Gillespie Algorithm)	378
14.2.1	Basic Outline of the Gillespie Method	379
14.3	Stochastic Effects in Gene Transcription	380
14.3.1	Expression of a Single Gene	380
14.3.2	Toggle Switch	381
14.4	Stochastic Modeling of a Small Molecular Network	385
14.4.1	Model System: Bacterial Photosynthesis	385
14.4.2	Pools-and-Proteins Model	386
14.4.3	Evaluating the Binding and Unbinding Kinetics	387
14.4.4	Pools of the Chromatophore Vesicle	389
14.4.5	Steady-State Regimes of the Vesicle	389
14.5	Parameter Optimization with Genetic Algorithm	392

14.6	Protein–Protein Association	395
14.7	Brownian Dynamics Simulations	396
14.8	Summary	398
14.9	Problems	400
14.9.1	Dynamic Simulations of Networks	400
	Bibliography	407
<b>15</b>	<b>Integrated Cellular Networks</b>	<b>409</b>
15.1	Response of Gene Regulatory Network to Outside Stimuli	410
15.2	Whole-Cell Model of <i>Mycoplasma genitalium</i>	412
15.3	Architecture of the Nuclear Pore Complex	416
15.4	Integrative Differential Gene Regulatory Network for Breast Cancer Identified Putative Cancer Driver Genes	416
15.5	Particle Simulations	421
15.6	Summary	423
	Bibliography	424
<b>16</b>	<b>Outlook</b>	<b>427</b>
	<b>Index</b>	<b>429</b>



## Preface of the First Edition

This book grew out of a course for graduate students in the first year of the MSc bioinformatics program that the author teaches every year at Saarland University. Also included is some material from a special lecture on cell simulations. The book is designed as a textbook, placing emphasis on transmitting the main ideas of a problem, outlining algorithmic strategies for solving these, and describing possible complications or connections to other parts of the book. The main challenge during the writing of the book was the concentration on conceptual points that may be of general educative value rather than including the latest research results of this fascinating fast-moving field. It is considered more important for a textbook to give a cohesive picture rather than mentioning all possible drawbacks and special cases where particular general guidelines may not apply. We apologize to those whose work could not be mentioned because of space constraints.

The intended audience includes students of bioinformatics and from life science disciplines. Consequently, some basic knowledge in molecular biology is taken for granted. The language used is not very formal. Previous knowledge of computer science is not required, but a certain adeptness in basic mathematics is necessary. The book introduces all of the mathematical concepts needed to understand the material covered. In particular, Chapter 2 introduces mathematical graphs and algorithms on graphs used in classifying protein–protein interaction networks. Chapter 6 introduces linear and convex algebra typically being used in the description of metabolic networks. Chapter 7 discusses ordinary and stochastic differential equations used in the kinetic modeling of signal transduction pathways. Chapter 8 introduces the method of Fourier transformation for protein–protein docking and pattern matching. Also introduced are Bayesian networks in Chapter 4 as a way to judge the reliability of protein–protein interactions and inference techniques to model gene regulatory networks. We note, however, that the emphasis of this book is placed on discrete mathematics rather than on statistical methods. Not included yet are classical network flow algorithms such as Menger’s theorem or the max-flow min-cut theorem as they are currently rarely used in cellular modeling. The book focuses on proteins and the genes coding for them, as well as on metabolites. Less room is given to DNA, RNA, or lipid membranes that would, of course, also deserve a great deal of attention. The main reason for this was to provide a homogenous background for discussing algorithmic concepts.

The author is very grateful to Dr. Tihamér Geyer who coordinated the assignments for the lectures for valuable comments on the manuscript and for many solved examples and problems for this book. The following coworkers from Saarbrücken and elsewhere have provided valuable suggestions on different portions of the text: Kerstin Kunz, Jan Christoph, and Florian Lauck. I thank Dr. Hawoong Jeong, Dr. Julio Collado-Vides, Dr. Agustino Martínez-Antonio, Dr. Ruth Sperling, Dr. James R. Williamson, Dr. Joanna Trylska, Dr. Claude Antony, and Dr. Nicholas Luscombe for sending me high-resolution versions of their graphics. I thank Dr. Andreas Sendtko and the publishing staff at Wiley-VCH for their generous support of this book project, for their seemingly endless patience during the revision stage, and for excellent typesetting.

I also thank the Center of Theoretical Biophysics at the University of California, San Diego, for their hospitality during a sabbatical visit in summer 2007 that finally allowed to complete this work. Finally, this book would not have been possible without the support and patience of my wife Regina and our two daughters.

March 2008

*Volkhard Helms*  
Center for Bioinformatics  
Saarland University  
Saarbrücken, Germany

## Preface of the Second Edition<sup>1</sup>

About 10 years after the publication of the first edition, I finally managed to prepare this expanded second edition of this book. Its main spirit remained the same: it is designed as a textbook, placing emphasis on transmitting the main ideas of a problem, outlining algorithmic strategies for solving these, and describing possible complications or connections to other parts of the book. Because of the feedback from colleagues, I have reordered the content, starting now in Chapter 2 with an introduction into the structures of protein–protein complexes before we enter into the world of protein interaction networks. I refrain from listing all the rearrangements here. Usually, I tried to keep subsections intact and simply shifted them around. A few sections were removed from the text because I now felt that they were too specialized. About 50% of new content has been added. In terms of mathematical methods, much more room is now given to statistical methods. In terms of biology, several new chapters now address protein–DNA interactions, epigenetic modifications, and microRNAs. Still not covered are biophysical topics related to intracellular transport, cytoskeletal dynamics, and processes taking place at and across biological membranes. Maybe, there will be a need for a third edition eventually?

In addition to those who contributed to the first edition, the author is very grateful to Thorsten Will and Maryam Nazarieh for solved examples and problems for this book. The following coworkers from Saarbrücken and elsewhere have provided valuable suggestions on different portions of the text: Mohamed Hamed Fahmy, Dania Humaidan, Olga Kalinina, Heiko Rieger, and Thorsten Will. I thank my group members of the past years with whom I had the privilege to work on exciting research projects related to the content of this book and I thank our secretary Kerstin Gronow-Pudelek for technical assistance.

April 2018

*Volkhard Helms*  
Center for Bioinformatics  
Saarland University  
Saarbrücken, Germany

<sup>1</sup> Problems: To really absorb the content of this textbook, it is advisable to also try to solve some of the problems enclosed.



## 1

## Networks in Biological Cells

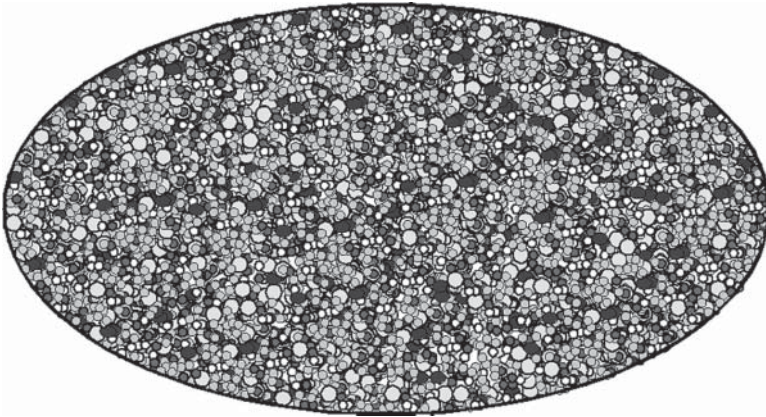
Modern molecular and cell biology has worked out many important cellular processes in more detail, although some other areas are known to a lesser extent. It often remains to understand how the individual parts are connected, and this is exactly the focus of this book. Figure 1.1 displays a cartoon of a cell as a highly viscous soup containing a complicated mixture of many particles. Certainly, several important details are left out here that introduce a partial order, such as the cytoskeleton and organelles of eukaryotic cells. Figure 1.1 reminds us that there is a myriad of biomolecular interactions taking place in biological cells at all times and that it is pretty amazing how a considerable order is achieved in many cellular processes that are all based on pairwise molecular interactions.

The focus of this book is placed on presenting mathematical descriptions developed in recent years to describe various levels of cellular networks. We will learn that many biological processes are tightly interconnected, and this is exactly where many links still need to be discovered in further experimental studies. Many researchers in the field of molecular biology believe that only combined efforts of modern experimental techniques, mathematical modeling, and bioinformatics analysis will be able to arrive at a sufficient understanding of the biological networks of cells and organisms.

In this chapter, we will start with some principles of mathematical networks and their relationship with biological networks. Then, we will briefly look at several biological key players to be used in the rest of this book (cells, compartments, proteins, and pathways). Without going into any further detail, we will directly move into the field of network theory with the amazing “small-world phenomenon.”

### 1.1 Some Basics About Networks

**Network theory** is a branch of applied mathematics and more of physics that uses the concepts of graph theory. Its developments are led by application to real-world examples in the areas of social networks (such as networks of acquaintances or among scientists having joint publications), technological networks (such as the World Wide Web that is a network of web pages and the Internet that is a network of computers and routers or power grids), and biological networks (such as neural networks and metabolic networks).



**Figure 1.1** Is this how we should view a biological cell? The point of this schematic picture is that about 30% of the volume of a biological cell is taken up by millions of individual proteins. Therefore, biological cells are really “full.” However, of course, such pictures do not tell us much about the organization of biological processes. As we will see later in this book, there are many different hierarchies of order in such a cell.

### 1.1.1 Random Networks

In a random network, every possible link between two “vertices” (or nodes) A and B is established according to a given probability distribution irrespective of the nature and connectivity of the two vertices A and B. This is what is “random” about these networks. If the network contains  $n$  vertices in total, the maximal number of undirected edges (links) between them is  $n \times (n - 1)/2$ . This is because we can pick each of the  $n$  vertices as the first vertex of an edge, and there are  $(n - 1)$  other vertices that this vertex can be connected to. In this way, we will actually consider each edge twice, using each end point as the first vertex. Therefore, we need to divide the number of edges by 2.

If every edge is established with a probability  $p \in [0, 1]$ , the total number of edges in an undirected graph is  $p \times n \times (n - 1)/2$ . The mathematics of random graphs was developed and elucidated by two Hungarian mathematicians Erdős and Renyi. However, the analysis of real networks showed that such networks often differ significantly from the characteristics of random graphs. We will turn back to random graphs in Section 6.3.

### 1.1.2 Small-World Phenomenon

The term **small-world phenomenon** was coined to describe the observation that everyone in the world is linked to some other person through a short chain of social acquaintances. In a **small-world experiment**, the psychologist Stanley Milgram found in 1967 that, on average, any two US citizens randomly picked were connected to each other by only six acquaintances. Vertices in a network have short average distances. Usually, the distance between the nodes scales logarithmically with the total number,  $n$ , of the vertices.

In a paper published in the journal *Nature* in 1998, the two mathematicians Duncan J. Watts and Steven H. Strogatz (Watts and Strogatz, 1998) reported

that small-world networks are common in many different areas ranging from neuronal connections of the worm *Caenorhabditis elegans* to power grids.

### 1.1.3 Scale-Free Networks

Only one year after the discovery of Watts and Strogatz, Albert-László Barabási from the Physics Department at the University of Notre Dame introduced an even simpler model for the emergence of the small-world phenomenon (Barabási and Albert 1999). Although Watts and Strogatz's model was able to explain the short average path length and the dense clustering coefficient of a *small world* (all these terms will be introduced in Chapter 6), it did not manage to explain another property that is typical for real-world networks such as the Internet: these networks are **scale-free**. In simple terms, this means that although the vast majority of vertices are weakly connected, there also exist some highly interconnected super-vertices or **hubs**. The term scale-free expresses that the ratio of highly to weakly connected vertices remains the same irrespective of the total number of links in the network. We will see in Section 6.4 that the connectivity of scale-free networks follows a power law. If a network is scale-free, it is also a small world.

In this paper, Barabási and Albert presented a strikingly simple and intuitive algorithm that generates networks with a scale-free topology. It has two essential elements:

- *Growth*. The network is started from a small number of (at least two) connected vertices. At every iteration step, a new vertex is added that forms links to  $m$  of the existing vertices.
- *Preferential attachment*. One assumes that the probability of a link between a newly added vertex and an existing vertex  $i$  depends on the degree of  $i$  (the number of existing links between vertex  $i$  and other vertices). The more connections  $i$  has already, the more likely the new vertices will link to  $i$ . This behavior is described by the saying “the rich become richer.” Let us motivate this on the fictitious example of the early days of air traffic. Initially, one needs to build two airports so that a first regular flight connection can be established between them. Eventually, a third airport is established. Most likely, initially, only one new flight will go to either one of the existing airports. Now, the situation is unbalanced. Now, there exists one airport that is connected to two other cities, and the airports of those cities are only connected to one city. There is a certain chance that, after some time, the “missing” connection between the new airport and the other airport would be introduced, which would lead to a balanced situation again. Alternatively, a fourth airport could emerge that would also start by establishing only one flight to one of the existing airports. Now, the airport that already has two connections would have an obvious practical advantage because passengers taking this route simply have more options to carry on. Therefore, the chance that this flight is established is higher than for the other connections. Exactly, this idea is captured by the concept of preferential attachment.

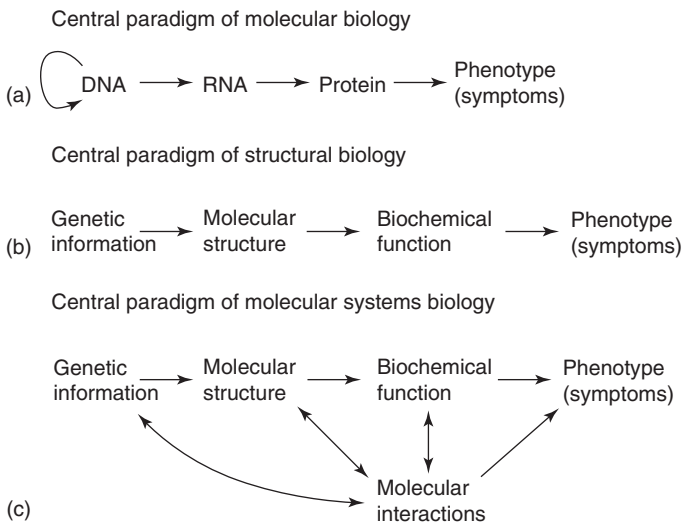
The same growth mechanism applies, for example, to the World Wide Web. Obviously, this network grows constantly over time, and many new pages are

added to it every moment. We know from our own experience that once a new web page is created, its owner will most likely include links to other popular pages (hubs) on the new page so that the second “rule” is also fulfilled.

In the early exciting days of network theory when the study of large-scale networks took off like a storm, it was even suggested that the scale-free network model may be something like a law of nature that controls how natural small-world networks are formed. However, subsequent work on integrated biological networks showed that the concept of scale-free networks may rather be of theoretical value and that it may not be directly applicable to certain biological networks. For the moment, we will consider the idea of network topology (scale-free networks and small-world phenomenon) as a powerful concept that is useful for understanding the mechanism of network growth and vulnerability.

## 1.2 Biological Background

Until recently, the paradigm of molecular biology was that genetic information is read from the genomic DNA by the RNA polymerase complex and is **transcribed** into the corresponding RNA. Ribosomes then bind to messenger RNA (mRNA) snippets and produce amino acid strands. This process is called **translation**. Importantly, the paradigm involved the notion that this entire process is unidirectional, see Figure 1.2.



**Figure 1.2** (a) Since the 1950s, a paradigm was established, whereby the information flows from DNA over RNA to protein synthesis, which then gives rise to particular phenotypes. (b) The emergence of structural biology – the first crystal structure of the protein myoglobin was determined in 1960 – emphasized the importance of the three-dimensional structures of proteins determining their function. (c) Today, we have realized the central role played by molecular interactions that influence all other elements.

### 1.2.1 Transcriptional Regulation

It is now well established that many feedback loops are provided in this system too, e.g. by the proteins known as transcription factors that bind to sequence motifs on the genomic DNA and mediate (activate or repress) transcription of certain genomic segments. Important discoveries of the past 20 years showed that cellular mRNA concentrations are also largely affected by small RNA snippets termed microRNAs and that the chromatin structure is shaped by epigenetic modifications of the DNA and histone proteins that control the accessibility of genomic regions. The cellular network therefore certainly appears much more complicated today than it did 60 years ago.

This brings us to the world of **gene regulatory networks**. Collecting the required information on the regulation of individual genes is a subject of intense active research. For example, the ENCODE project for human cells and the modENCODE project for the model organisms *C. elegans* and *Drosophila melanogaster* mapped the binding sites of hundreds of transcription factors throughout the genomes. Also, the FANTOM initiative started in Japan is a worldwide collaborative project aiming at identifying all the functional elements in mammalian genomes. However, occupancy maps of transcription factors alone are not being considered as compelling evidence of biologically functional regulation. To really prove or disprove which gene is activated or repressed by a particular transcription factor (or microRNA), one could create a knockout organism lacking the gene coding for this transcription factor and see which genes are no longer expressed or are now expressed in excess. Such genome-wide deletion libraries have actually been produced for the model organism *Saccharomyces cerevisiae*. However, in this way, we can only discover those combinations that are not lethal for the organism. Also, pairs or larger assemblies of transcription factors often need to bind simultaneously. It simply appears impossible to discover the full connectivity of this regulatory network by a traditional one-by-one approach. Fortunately, modern microarray and RNAseq experiments probe the expression levels of many genes simultaneously. Ongoing challenges are the noisy nature of the large-scale data and the fact that genes actually do not interact directly with each other. Analysis of gene expression data will be discussed in Chapter 8.

In this book, we will be mostly concerned with the following four types of biological cellular networks: protein–protein interaction networks, gene regulatory networks, signal transduction networks, and metabolic networks. We will discuss them at different hierarchical levels as shown in Figure 1.3 using the example of regulatory networks.

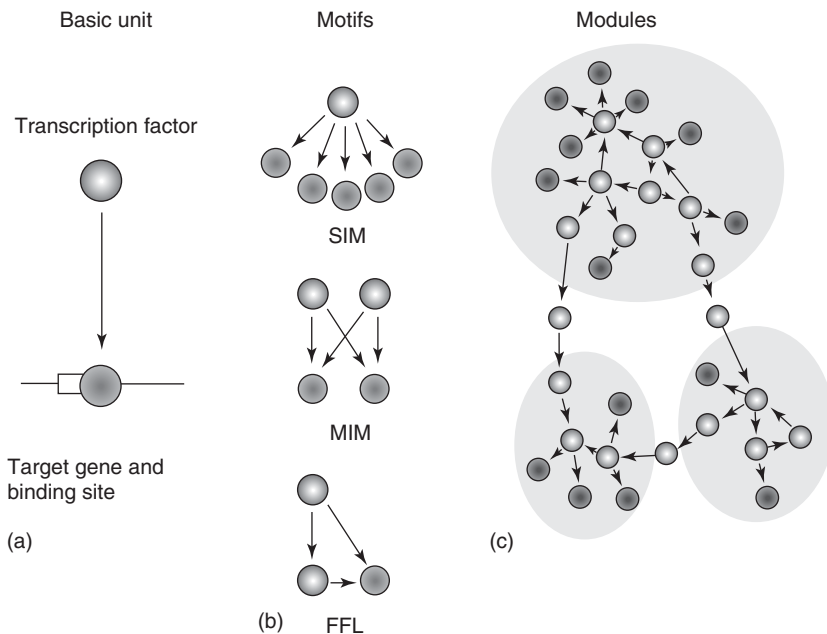
### 1.2.2 Cellular Components

Cells can be described at various levels in detail. We will mostly use three different levels of description:

- (a) *Inventory lists and lists of processes.*
  - Proteins in particular compartments
  - Proteins forming macromolecular complexes

- Biomolecular interactions
  - Regulatory interactions
  - Metabolic reactions
- (b) *Structural descriptions.*
- Structures of single proteins
  - Topologies of protein complexes
  - Subcellular compartments
- (c) *Dynamic descriptions.*
- Cellular processes ranging from nanosecond dynamics for the association of two biomolecules up to processes occurring in seconds and minutes such as the cell division of yeast cells.

We will assume that the reader has a basic knowledge about the organic molecules commonly found within living cells and refer those who do not to basic books on biochemistry or molecular biology. Depending on their role in metabolism, the biomolecules in a cell can be grouped into several classes.



**Figure 1.3** Structural organization of transcriptional regulatory networks. (a) The “basic unit” comprises the transcription factor, its target gene with a DNA recognition site, and the regulatory interaction between them. (b) Units are often organized into network “motifs” that comprise specific patterns of inter-regulation that are overrepresented in networks. Examples of motifs include single-input/multiple output (SIM), multiple input/multiple output (MIM), and feed-forward loop (FFL) motifs. (c) Network motifs can be interconnected to form semi-independent “modules,” many of which have been identified by integrating regulatory interaction data with gene expression data and imposing evolutionary conservation. The next level consists of the entire network (not shown). Source: Babu et al. (2004). Drawn with permission of Elsevier.

1. **Macromolecules** including nucleic acids, proteins, polysaccharides, and certain lipids.
2. The **building blocks** of macromolecules include sugars as the precursors of polysaccharides, amino acids as the building blocks of proteins, nucleotides as the precursors of nucleic acids (and therefore of DNA and RNA), and fatty acids that are incorporated into lipids. Interestingly, in biological cells, only a small number of theoretically synthesizable macromolecules exist at a given time point. At any moment during a normal cell cycle, many new macromolecules need to be synthesized from their building blocks, and this is meticulously controlled by the complex gene expression machinery. Even during a steady state of the cell, there exists a constant turnover of macromolecules.
3. *Metabolic intermediates (metabolites)*. Many molecules in a biological cell have complex chemical structures and must be synthesized in several reactions from specific starting materials that may be taken up as the energy source. In the cell, connected chemical reactions are often grouped into metabolic pathways (Section 1.3).
4. Molecules of **miscellaneous function** including vitamins, steroid hormones, molecules that can store energy storage such as ATP, regulatory molecules, and metabolic waste products.

Almost all biological materials that are needed to construct a biological cell are either synthesized by the RNA polymerase and ribosome machinery of the cell or are taken up from the outside via the cell membrane. Therefore, as a minimum inventory, every cell needs to contain the construction plan (DNA), a processing unit to transcribe this information into mRNA (polymerase), a processing unit to translate these mRNA pieces into protein (ribosome), and transporter proteins inside the cell membrane that transport material through the cell membrane.

### 1.2.3 Spatial Organization of Eukaryotic Cells into Compartments

Organization into various compartments greatly simplifies the temporal and spatial process flow in eukaryotic cells. As mentioned above, at each time point during a cell cycle, only a small subfraction of all potential proteins is being synthesized (and not yet degraded). Also, many proteins are only available in very small concentrations, possibly with only a few copies per cell. However, localizing these proteins to particular spots in the cell, e.g. by attaching them to the cytoskeleton or by partitioning them into lipid rafts, their local concentrations may be much higher. We assume that the reader is vaguely familiar with the compartmentalization of eukaryotic cells involving the lysosome, plasma membrane, cell membrane, Golgi complex, nucleus, smooth endoplasmic reticulum, mitochondrion, nucleolus, rough endoplasmic reticulum, and cytoskeleton.

An important element of cellular organization is the active transport of macromolecules along the microtubules of the cytoskeleton that is carried out by molecular motor proteins such as kinesin and dynein. Here, we will not address the activities of molecular motors because this is rather a research topic in biophysics.

**Table 1.1** Data on the genome length and on the number of protein-coding and RNA genes are taken from the Kyoto Encyclopedia of Genes and Genomes database (April 2018); data on the number of putative transporter proteins are taken from [www.membranetransport.org](http://www.membranetransport.org).

Organism	Length of genome (Mb)	Number of protein-coding genes	Number of RNA genes	Number of transporter proteins
<b>Prokaryotes</b>				
<i>Mycoplasma genitalium</i> G37	0.6	476	43	53
<i>Bacillus subtilis</i> BSN5	4.2	4 145	113	552
<i>Escherichia coli</i> APEC01	4.6	4 890	93	665
<b>Eukaryotes</b>				
<i>Saccharomyces cerevisiae</i> S288C	1.3	6 002	425	341
<i>Drosophila melanogaster</i>	12	13 929	3 209	662
<i>Caenorhabditis elegans</i>	100.2	20 093	24 969	669
<i>Homo sapiens</i>	3 150	20 338	19 201	1 467

## 1.2.4 Considered Organisms

Table 1.1 presents some statistics of the organisms considered in this book.

## 1.3 Cellular Pathways

### 1.3.1 Biochemical Pathways

**Metabolism** denotes the entirety of biochemical reactions that occur within a cell (Figure 1.4). In the past century, many of these reactions have been organized into **metabolic pathways**. Each pathway consists of a sequence of chemical reactions that are catalyzed by specific enzymes, and the outcome of one reaction is the input for the next one. Unraveling the individual enzymatic reactions was one of the big successes of applying biochemical methods to cellular processes. Metabolic pathways can be divided into two broad types. **Catabolic pathways** disintegrate complex molecules into simpler ones, which can be reused for synthesizing other molecules. Also, catabolic pathways provide chemical energy required for many cellular processes. This energy may be stored temporarily as high-energy phosphates (primarily in ATP) or as high-energy electrons (primarily in NADPH). Conversely, **anabolic pathways** synthesize more complex substances from simpler starting reagents by utilizing the chemical energy generated by exergonic catabolic pathways.

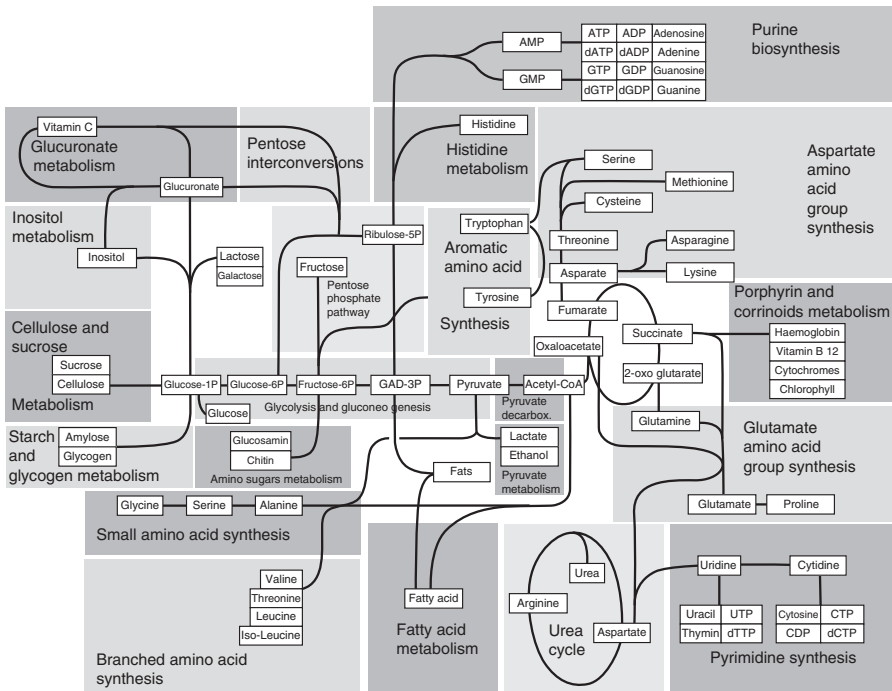
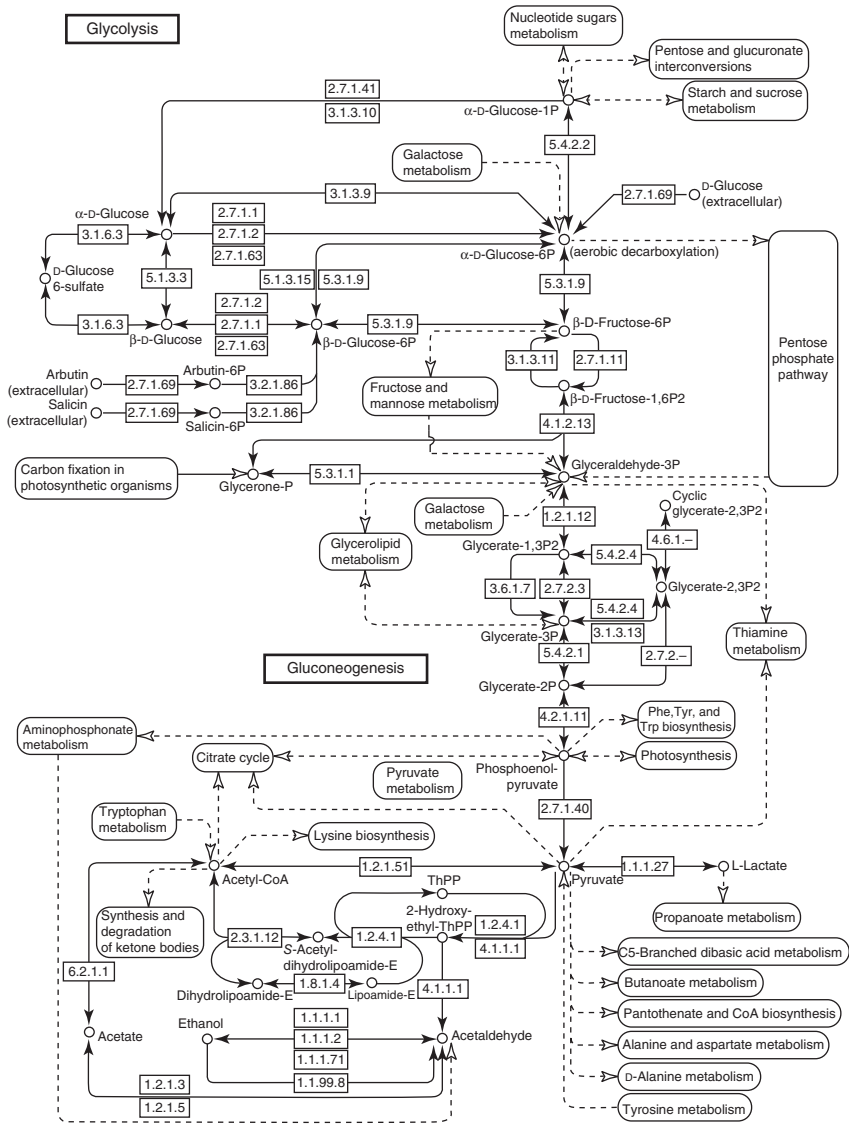


Figure 1.4 Major metabolic pathways.

The traditional biochemical pathways were often derived from studying simple organisms where these pathways constitute a dominating part of the metabolic activity. For example, the **glycolysis** pathway was discovered in yeast (and in muscle) in the 1930s. It describes the disassembly of the nutrient glucose that is taken up by many microorganisms from the outside. Figure 1.5 shows the glycolysis pathway in *Homo sapiens* as represented in the KEGG database (Kanehisa et al. 2016).



00010 8/6/07

**Figure 1.5** The glycolysis pathway as visualized in the KEGG database is connected to many other cellular pathways. Source: From <http://www.genome.ad.jp/kegg>.