

Edda Klipp, Wolfram Liebermeister,
Christoph Wierling and Axel Kowald

Systems Biology

A Textbook

Second Edition



Edda Klipp
Wolfram Liebermeister
Christoph Wierling
Axel Kowald

Systems Biology

*Edda Klipp
Wolfram Liebermeister
Christoph Wierling
Axel Kowald*

Systems Biology

A Textbook

Second, Completely Revised
and Enlarged Edition

WILEY-VCH
Verlag GmbH & Co. KGaA

Authors

Prof. Dr. h.c. Edda Klipp

Theoretical Biophysics
Humboldt-Universität zu Berlin
Invalidenstr. 42
10115 Berlin
Germany

Dr. Wolfram Liebermeister

Institute of Biochemistry
Charité - Universitätsmedizin Berlin
Charitéplatz 1
10117 Berlin
Germany

Dr. Christoph Wierling

Alacris Theranostics GmbH
Fabeckstr. 60-62
14195 Berlin
Germany

and

Max Planck Institute for Molecular Genetics
Ihnestr. 63-73
14195 Berlin
Germany

Dr. Axel Kowald

Theoretical Biophysics
Humboldt University Berlin
Invalidenstr. 42
10115 Berlin
Germany

Cover

Cover design by Wolfram Liebermeister. The cover picture was provided with kind permission by Jörg Bernhardt.

■ 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>>.

© 2016 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: 978-3-527-33636-4

ePDF ISBN: 978-3-527-67566-1

ePub ISBN: 978-3-527-67567-8

Mobi ISBN: 978-3-527-67568-5

Typesetting Thomson Digital, Noida, India

Printed on acid-free paper

Contents

Preface	xi		
Guide to Different Topics of the Book	xiii		
About the Authors	xv		
Part One Introduction to Systems Biology	1		
1 Introduction	3		
1.1 Biology in Time and Space	3		
1.2 Models and Modeling	4		
1.2.1 What Is a Model?	4		
1.2.2 Purpose and Adequateness of Models	5		
1.2.3 Advantages of Computational Modeling	5		
1.3 Basic Notions for Computational Models	6		
1.3.1 Model Scope	6		
1.3.2 Model Statements	6		
1.3.3 System State	6		
1.3.4 Variables, Parameters, and Constants	6		
1.3.5 Model Behavior	7		
1.3.6 Model Classification	7		
1.3.7 Steady States	7		
1.3.8 Model Assignment Is Not Unique	7		
1.4 Networks	8		
1.5 Data Integration	8		
1.6 Standards	9		
1.7 Model Organisms	9		
1.7.1 <i>Escherichia coli</i>	9		
1.7.2 <i>Saccharomyces cerevisiae</i>	11		
1.7.3 <i>Caenorhabditis elegans</i>	11		
1.7.4 <i>Drosophila melanogaster</i>	11		
1.7.5 <i>Mus musculus</i>	12		
References	12		
Further Reading	14		
2 Modeling of Biochemical Systems	15		
2.1 Overview of Common Modeling Approaches for Biochemical Systems	15		
2.2 ODE Systems for Biochemical Networks	17		
2.2.1 Basic Components of ODE Models	18		
2.2.2 Illustrative Examples of ODE Models	18		
References	21		
Further Reading	21		
3 Structural Modeling and Analysis of Biochemical Networks	23		
3.1 Structural Analysis of Biochemical Systems	24		
3.1.1 System Equations	24		
3.1.2 Information Encoded in the Stoichiometric Matrix N	25		
3.1.3 The Flux Cone	27		
3.1.4 Elementary Flux Modes and Extreme Pathways	27		
3.1.5 Conservation Relations – Null Space of N^T	29		
3.2 Constraint-Based Flux Optimization	30		
3.2.1 Flux Balance Analysis	31		
3.2.2 Geometric Interpretation of Flux Balance Analysis	31		
3.2.3 Thermodynamic Constraints	31		
3.2.4 Applications and Tests of the Flux Optimization Paradigm	32		
3.2.5 Extensions of Flux Balance Analysis	33		
Exercises	35		
References	36		
Further Reading	37		
4 Kinetic Models of Biochemical Networks: Introduction	39		
4.1 Reaction Kinetics and Thermodynamics	39		
4.1.1 Kinetic Modeling of Enzymatic Reactions	39		
4.1.2 The Law of Mass Action	40		
4.1.3 Reaction Thermodynamics	40		
4.1.4 Michaelis–Menten Kinetics	42		
4.1.5 Regulation of Enzyme Activity by Effectors	44		
4.1.6 Generalized Mass Action Kinetics	48		
4.1.7 Approximate Kinetic Formats	48		
4.1.8 Convenience Kinetics and Modular Rate Laws	49		
4.2 Metabolic Control Analysis	50		
4.2.1 The Coefficients of Control Analysis	51		

4.2.2	The Theorems of Metabolic Control Theory	53	6.4	Coupled Systems and Emergent Behavior	110
4.2.3	Matrix Expressions for Control Coefficients	55	6.4.1	Modeling of Coupled Systems	111
4.2.4	Upper Glycolysis as Realistic Model Example	58	6.4.2	Combining Rate Laws into Models	113
4.2.5	Time-Dependent Response Coefficients	59	6.4.3	Modular Response Analysis	113
	Exercises	61	6.4.4	Emergent Behavior in Coupled Systems	114
	References	61	6.4.5	Causal Interactions and Global Behavior	115
	Further Reading	62		Exercises	116
5	Data Formats, Simulation Techniques, and Modeling Tools	63		References	117
5.1	Simulation Techniques and Tools	63		Further Reading	119
5.1.1	Differential Equations	63	7	Discrete, Stochastic, and Spatial Models	121
5.1.2	Stochastic Simulations	64	7.1	Discrete Models	122
5.1.3	Simulation Tools	65	7.1.1	Boolean Networks	122
5.2	Standards and Formats for Systems Biology	72	7.1.2	Petri Nets	124
5.2.1	Systems Biology Markup Language	72	7.2	Stochastic Modeling of Biochemical Reactions	127
5.2.2	BioPAX	74	7.2.1	Chance in Biochemical Reaction Systems	127
5.2.3	Systems Biology Graphical Notation	74	7.2.2	The Chemical Master Equation	129
5.3	Data Resources for Modeling of Cellular Reaction Systems	75	7.2.3	Stochastic Simulation	129
5.3.1	General-Purpose Databases	75	7.2.4	Chemical Langevin Equation and Chemical Noise	130
5.3.2	Pathway Databases	76	7.2.5	Dynamic Fluctuations	132
5.3.3	Model Databases	77	7.2.6	From Stochastic to Deterministic Modeling	133
5.4	Sustainable Modeling and Model Semantics	78	7.3	Spatial Models	133
5.4.1	Standards for Systems Biology Models	78	7.3.1	Types of Spatial Models	134
5.4.2	Model Semantics and Model Comparison	78	7.3.2	Compartment Models	135
5.4.3	Model Combination	80	7.3.3	Reaction–Diffusion Systems	136
5.4.4	Model Validity	82	7.3.4	Robust Pattern Formation in Embryonic Development	138
	References	83	7.3.5	Spontaneous Pattern Formation	139
	Further Reading	85	7.3.6	Linear Stability Analysis of the Activator–Inhibitor Model	140
6	Model Fitting, Reduction, and Coupling	87		Exercises	142
6.1	Parameter Estimation	88		References	143
6.1.1	Regression, Estimators, and Maximal Likelihood	88		Further Reading	144
6.1.2	Parameter Identifiability	90	8	Network Structure, Dynamics, and Function	145
6.1.3	Bootstrapping	91	8.1	Structure of Biochemical Networks	146
6.1.4	Bayesian Parameter Estimation	92	8.1.1	Random Graphs	147
6.1.5	Probability Distributions for Rate Constants	94	8.1.2	Scale-Free Networks	148
6.1.6	Optimization Methods	97	8.1.3	Connectivity and Node Distances	149
6.2	Model Selection	99	8.1.4	Network Motifs and Significance Tests	150
6.2.1	What Is a Good Model?	99	8.1.5	Explanations for Network Structures	151
6.2.2	The Problem of Model Selection	100	8.2	Regulation Networks and Network Motifs	152
6.2.3	Likelihood Ratio Test	102	8.2.1	Structure of Transcription Networks	153
6.2.4	Selection Criteria	102	8.2.2	Regulation Edges and Their Steady-State Response	156
6.2.5	Bayesian Model Selection	103	8.2.3	Negative Feedback	156
6.3	Model Reduction	104	8.2.4	Adaptation Motif	157
6.3.1	Model Simplification	104	8.2.5	Feed-Forward Loops	158
6.3.2	Reduction of Fast Processes	105			
6.3.3	Quasi-Equilibrium and Quasi-Steady State	107			
6.3.4	Global Model Reduction	108			

8.3	Modularity and Gene Functions	160	10	Variability, Robustness, and Information	209
8.3.1	Cell Functions Are Reflected in Structure, Dynamics, Regulation, and Genetics	160	10.1	Variability and Biochemical Models	210
8.3.2	Metabolic Pathways and Elementary Modes	162	10.1.1	Variability and Uncertainty Analysis	210
8.3.3	Epistasis Can Indicate Functional Modules	163	10.1.2	Flux Sampling	212
8.3.4	Evolution of Function and Modules	163	10.1.3	Elasticity Sampling	213
8.3.5	Independent Systems as a Tacit Model Assumption	165	10.1.4	Propagation of Parameter Variability in Kinetic Models	214
8.3.6	Modularity and Biological Function Are Conceptual Abstractions	165	10.1.5	Models with Parameter Fluctuations	216
	Exercises	166	10.2	Robustness Mechanisms and Scaling Laws	217
	References	167	10.2.1	Robustness in Biochemical Systems	218
	Further Reading	169	10.2.2	Robustness by Backup Elements	219
9	Gene Expression Models	171	10.2.3	Feedback Control	219
9.1	Mechanisms of Gene Expression Regulation	171	10.2.4	Perfect Robustness by Structure	222
9.1.1	Transcription Factor-Initiated Gene Regulation	171	10.2.5	Scaling Laws	224
9.1.2	General Promoter Structure	173	10.2.6	Time Scaling, Summation Laws, and Robustness	227
9.1.3	Prediction and Analysis of Promoter Elements	174	10.2.7	The Role of Robustness in Evolution and Modeling	228
9.1.4	Posttranscriptional Regulation through microRNAs	176	10.3	Adaptation and Exploration Strategies	229
9.2	Dynamic Models of Gene Regulation	180	10.3.1	Information Transmission in Signaling Pathways	230
9.2.1	A Basic Model of Gene Expression and Regulation	180	10.3.2	Adaptation and Fold-Change Detection	230
9.2.2	Natural and Synthetic Gene Regulatory Networks	183	10.3.3	Two Adaptation Mechanisms: Sensing and Random Switching	231
9.2.3	Gene Expression Modeling with Stochastic Equations	186	10.3.4	Shannon Information and the Value of Information	232
9.3	Gene Regulation Functions	187	10.3.5	Metabolic Shifts and Anticipation	233
9.3.1	The Lac Operon in <i>E. coli</i>	187	10.3.6	Exploration Strategies	234
9.3.2	Gene Regulation Functions Derived from Equilibrium Binding	188		Exercises	236
9.3.3	Thermodynamic Models of Promoter Occupancy	189		References	237
9.3.4	Gene Regulation Function of the Lac Promoter	191		Further Reading	239
9.3.5	Inferring Transcription Factor Activities from Transcription Data	192	11	Optimality and Evolution	241
9.3.6	Network Component Analysis	194	11.1	Optimality in Systems Biology Models	243
9.3.7	Correspondences between mRNA and Protein Levels	196	11.1.1	Mathematical Concepts for Optimality and Compromise	245
9.4	Fluctuations in Gene Expression	196	11.1.2	Metabolism Is Shaped by Optimality	248
9.4.1	Stochastic Model of Transcription and Translation	197	11.1.3	Optimality Approaches in Metabolic Modeling	250
9.4.2	Intrinsic and Extrinsic Variability	200	11.1.4	Metabolic Strategies	252
9.4.3	Temporal Fluctuations in Gene Cascades	202	11.1.5	Optimal Metabolic Adaptation	253
	Exercises	203	11.2	Optimal Enzyme Concentrations	255
	References	205	11.2.1	Optimization of Catalytic Properties of Single Enzymes	255
	Further Reading	207	11.2.2	Optimal Distribution of Enzyme Concentrations in a Metabolic Pathway	257
			11.2.3	Temporal Transcription Programs	259
			11.3	Evolution and Self-Organization	261
			11.3.1	Introduction	261
			11.3.2	Selection Equations for Biological Macromolecules	263
			11.3.3	The Quasispecies Model: Self-Replication with Mutations	265

- 11.3.4 The Hypercycle 267
- 11.3.5 Other Mathematical Models of Evolution: Spin Glass Model 269
- 11.3.6 The Neutral Theory of Molecular Evolution 270
- 11.4 Evolutionary Game Theory 271**
 - 11.4.1 Social Interactions 272
 - 11.4.2 Game Theory 273
 - 11.4.3 Evolutionary Game Theory 274
 - 11.4.4 Replicator Equation for Population Dynamics 274
 - 11.4.5 Evolutionarily Stable Strategies 275
 - 11.4.6 Dynamical Behavior in the Rock–Scissors–Paper Game 276
 - 11.4.7 Evolution of Cooperative Behavior 276
 - 11.4.8 Compromises between Metabolic Yield and Efficiency 278
 - Exercises 279
 - References 280
 - Further Reading 283
- 12 Models of Biochemical Systems 285**
 - 12.1 Metabolic Systems 285**
 - 12.1.1 Basic Elements of Metabolic Modeling 286
 - 12.1.2 Toy Model of Upper Glycolysis 286
 - 12.1.3 Threonine Synthesis Pathway Model 289
 - 12.2 Signaling Pathways 291**
 - 12.2.1 Function and Structure of Intra- and Intercellular Communication 292
 - 12.2.2 Receptor–Ligand Interactions 293
 - 12.2.3 Structural Components of Signaling Pathways 295
 - 12.2.4 Analysis of Dynamic and Regulatory Features of Signaling Pathways 304
 - 12.3 The Cell Cycle 307**
 - 12.3.1 Steps in the Cycle 309
 - 12.3.2 Minimal Cascade Model of a Mitotic Oscillator 310
 - 12.3.3 Models of Budding Yeast Cell Cycle 311
 - 12.4 The Aging Process 314**
 - 12.4.1 Evolution of the Aging Process 316
 - 12.4.2 Using Stochastic Simulations to Study Mitochondrial Damage 318
 - 12.4.3 Using Delay Differential Equations to Study Mitochondrial Damage 323
 - Exercises 327
 - References 327
- Part Two Reference Section 331**
- 13 Cell Biology 333**
 - 13.1 The Origin of Life 334**
 - 13.2 Molecular Biology of the Cell 336**
 - 13.2.1 Chemical Bonds and Forces Important in Biological Molecules 336
 - 13.2.2 Functional Groups in Biological Molecules 338
 - 13.2.3 Major Classes of Biological Molecules 338
 - 13.3 Structural Cell Biology 345**
 - 13.3.1 Structure and Function of Biological Membranes 347
 - 13.3.2 Nucleus 349
 - 13.3.3 Cytosol 349
 - 13.3.4 Mitochondria 350
 - 13.3.5 Endoplasmic Reticulum and Golgi Complex 350
 - 13.3.6 Other Organelles 351
 - 13.4 Expression of Genes 351**
 - 13.4.1 Transcription 351
 - 13.4.2 Processing of the mRNA 353
 - 13.4.3 Translation 353
 - 13.4.4 Protein Sorting and Posttranslational Modifications 355
 - 13.4.5 Regulation of Gene Expression 355
 - Exercises 356
 - References 356
 - Further Reading 356
- 14 Experimental Techniques 357**
 - 14.1 Restriction Enzymes and Gel Electrophoresis 358**
 - 14.2 Cloning Vectors and DNA Libraries 359**
 - 14.3 1D and 2D Protein Gels 361**
 - 14.4 Hybridization and Blotting Techniques 362**
 - 14.4.1 Southern Blotting 363
 - 14.4.2 Northern Blotting 363
 - 14.4.3 Western Blotting 363
 - 14.4.4 *In Situ* Hybridization 364
 - 14.5 Further Protein Separation Techniques 364**
 - 14.5.1 Centrifugation 364
 - 14.5.2 Column Chromatography 364
 - 14.6 Polymerase Chain Reaction 365**
 - 14.7 Next-Generation Sequencing 366**
 - 14.8 DNA and Protein Chips 367**
 - 14.8.1 DNA Chips 367
 - 14.8.2 Protein Chips 367
 - 14.9 RNA-Seq 368**
 - 14.10 Yeast Two-Hybrid System 368**
 - 14.11 Mass Spectrometry 369**
 - 14.12 Transgenic Animals 370**
 - 14.12.1 Microinjection and ES Cells 370
 - 14.12.2 Genome Editing Using ZFN, TALENs, and CRISPR 370
 - 14.13 RNA Interference 371**
 - 14.14 ChIP-on-Chip and ChIP-PET 372**
 - 14.15 Green Fluorescent Protein 374**
 - 14.16 Single-Cell Experiments 375**

- 14.17 Surface Plasmon Resonance 376**
 - Exercises 377
 - References 377
- 15 Mathematical and Physical Concepts 381**
 - 15.1 Linear Algebra 381**
 - 15.1.1 Linear Equations 381
 - 15.1.2 Matrices 384
 - 15.2 Dynamic Systems 386**
 - 15.2.1 Describing Dynamics with Ordinary Differential Equations 386
 - 15.2.2 Linearization of Autonomous Systems 388
 - 15.2.3 Solution of Linear ODE Systems 388
 - 15.2.4 Stability of Steady States 388
 - 15.2.5 Global Stability of Steady States 390
 - 15.2.6 Limit Cycles 390
 - 15.3 Statistics 391**
 - 15.3.1 Basic Concepts of Probability Theory 391
 - 15.3.2 Descriptive Statistics 396
 - 15.3.3 Testing Statistical Hypotheses 399
 - 15.3.4 Linear Models 401
 - 15.3.5 Principal Component Analysis 404
 - 15.4 Stochastic Processes 405**
 - 15.4.1 Chance in Physical Theories 405
 - 15.4.2 Mathematical Random Processes 406
 - 15.4.3 Brownian Motion as a Random Process 407
 - 15.4.4 Markov Processes 409
 - 15.4.5 Markov Chains 410
 - 15.4.6 Jump Processes in Continuous Time 410
 - 15.4.7 Continuous Random Processes 411
 - 15.4.8 Moment-Generating Functions 412
 - 15.5 Control of Linear Dynamical Systems 412**
 - 15.5.1 Linear Dynamical Systems 413
 - 15.5.2 System Response and Linear Filters 414
 - 15.5.3 Random Fluctuations and Spectral Density 415
 - 15.5.4 The Gramian Matrices 415
 - 15.5.5 Model Reduction 416
 - 15.5.6 Optimal Control 416
 - 15.6 Biological Thermodynamics 417**
 - 15.6.1 Microstate and Statistical Ensemble 417
 - 15.6.2 Boltzmann Distribution and Free Energy 418
 - 15.6.3 Entropy 419
 - 15.6.4 Equilibrium Constant and Energies 421
 - 15.6.5 Chemical Reaction Systems 422
 - 15.6.6 Nonequilibrium Reactions 424
 - 15.6.7 The Role of Thermodynamics in Systems Biology 425
 - 15.7 Multivariate Statistics 426**
 - 15.7.1 Planning and Designing Experiments for Case-Control Studies 426
 - 15.7.2 Tests for Differential Expression 427
 - 15.7.3 Multiple Testing 428
 - 15.7.4 ROC Curve Analysis 429
 - 15.7.5 Clustering Algorithms 430
 - 15.7.6 Cluster Validation 435
 - 15.7.7 Overrepresentation and Enrichment Analyses 436
 - 15.7.8 Classification Methods 438
 - Exercises 441
 - References 443
- 16 Databases 445**
 - 16.1 General-Purpose Data Resources 445**
 - 16.1.1 PathGuide 445
 - 16.1.2 BioNumbers 446
 - 16.2 Nucleotide Sequence Databases 446**
 - 16.2.1 Data Repositories of the National Center for Biotechnology Information 446
 - 16.2.2 GenBank/RefSeq/UniGene 446
 - 16.2.3 Entrez 447
 - 16.2.4 EMBL Nucleotide Sequence Database 447
 - 16.2.5 European Nucleotide Archive 447
 - 16.2.6 Ensembl 447
 - 16.3 Protein Databases 448**
 - 16.3.1 UniProt/Swiss-Prot/TrEMBL 448
 - 16.3.2 Protein Data Bank 448
 - 16.3.3 PANTHER 448
 - 16.3.4 InterPro 448
 - 16.3.5 iHOP 449
 - 16.4 Ontology Databases 449**
 - 16.4.1 Gene Ontology 449
 - 16.5 Pathway Databases 449**
 - 16.5.1 KEGG 450
 - 16.5.2 Reactome 450
 - 16.5.3 ConsensusPathDB 451
 - 16.5.4 WikiPathways 451
 - 16.6 Enzyme Reaction Kinetics Databases 451**
 - 16.6.1 BRENDA 451
 - 16.6.2 SABIO-RK 452
 - 16.7 Model Collections 452**
 - 16.7.1 BioModels 452
 - 16.7.2 JWS Online 452
 - 16.8 Compound and Drug Databases 452**
 - 16.8.1 ChEBI 453
 - 16.8.2 Guide to PHARMACOLOGY 453
 - 16.9 Transcription Factor Databases 453**
 - 16.9.1 JASPAR 453
 - 16.9.2 TRED 453
 - 16.9.3 Transcription Factor Encyclopedia 454
 - 16.10 Microarray and Sequencing Databases 454**
 - 16.10.1 Gene Expression Omnibus 454
 - 16.10.2 ArrayExpress 454
 - References 455

17	Software Tools for Modeling	457
17.1	13C-Flux2	458
17.2	Antimony	458
17.3	Berkeley Madonna	459
17.4	BIOCHAM	459
17.5	BioNetGen	459
17.6	Biopython	459
17.7	BioTapestry	460
17.8	BioUML	460
17.9	CellDesigner	460
17.10	CellNetAnalyzer	460
17.11	Copasi	461
17.12	CPN Tools	461
17.13	Cytoscape	461
17.14	E-Cell	461
17.15	EvA2	461
17.16	FEniCS Project	462
17.17	Genetic Network Analyzer (GNA)	462
17.18	Jarnac	462
17.19	JDesigner	463
17.20	JSim	463
17.21	KNIME	463
17.22	libSBML	464
17.23	MASON	464
17.24	Mathematica	464
17.25	MathSBML	465
17.26	Matlab	465
17.27	MesoRD	465
17.28	Octave	465
17.29	Omix Visualization	466
17.30	OpenCOR	466
17.31	Oscill8	466
17.32	PhysioDesigner	466
17.33	PottersWheel	467
17.34	PyBioS	467
17.35	PySCeS	467
17.36	R	468
17.37	SAAM II	468
17.38	SBMLeditor	468
17.39	SemanticSBML	468
17.40	SBML-PET-MPI	469
17.41	SBMLsimulator	469
17.42	SBMLsqueezer	469
17.43	SBML Toolbox	470
17.44	SBtoolbox2	470
17.45	SBML Validator	470
17.46	SensA	470
17.47	SmartCell	471
17.48	STELLA	471
17.49	STEPS	471
17.50	StochKit2	471
17.51	SystemModeler	472
17.52	Systems Biology Workbench	472
17.53	Taverna	472
17.54	VANTED	473
17.55	Virtual Cell (VCell)	473
17.56	xCellerator	473
17.57	XPPAUT	473
	Exercises	474
	References	474
	Index	475

Preface

Systems biology is the scientific discipline that studies the systemic properties and dynamic interactions in a biological object, be it a cell, an organism, a virus, or an infected host, in a qualitative and quantitative manner and by combining experimental studies with mathematical modeling. Scientists can describe the inner processes of stars a thousand light years away with great accuracy. But how a tiny cell under our microscope grows and divides remains puzzling in many ways. We see kids growing, people aging, plants blooming, and microbes degrading their remains. We use yeast for brewery and bakery, and doctors prescribe drugs to cure diseases. But do we understand how processes of life work?

Starting in the nineteenth century, such processes have no longer been explained by referring to special “life forces,” but by laws of physics and chemistry. By studying the structure and dynamics of living systems in finer and finer details, researchers from different disciplines have revealed how life processes arise from the structure and functional organization of cells, how tens of thousands of biochemical components interact in orchestrated ways, and how these systems are regulated by genetic information and continuously adapted through mutations and selection. With this conceptual shift, new questions became central in biology: How does an organism’s phenotype emerge from the genotype, as encoded in the organism’s DNA, and how is it shaped by environmental factors? Initially, such questions were approached by statistics, for example, by studying what mutations are associated with specific inheritable diseases. But the task, now, is to understand the mechanistic details.

We can easily understand the effects of gene disruptions when gene products have simple, specific functions. However, most gene mutations have only weak or quantitative effects on physiology, and many genetic diseases are multifactorial. Tracing the effects of multiple mutations, of mutations affecting gene regulation, or of

drugs requires a deep, quantitative, and dynamical understanding of cell physiology. In recent years, high-throughput experiments, time series experiments, and imaging techniques with high resolution have provided us with a detailed picture of the cellular machinery. We can observe how physical structures are built, maintained, and reproduced, how the metabolic state is changing, and how signaling and regulation systems allow cells to adapt to their environment. However, to understand how all these systems act together – and how cells can work as complex, robust systems – cataloging and understanding single-cell components is not enough. Instead, we need to capture the global dynamics between these components. This is where mathematical models come into play.

Mathematical modeling has a long, though relatively marginal, tradition in biology, and has influenced the field in many ways. Models can be used to test hypotheses and to yield quantitative predictions or reveal gaps or inconsistencies in previous arguments, thus helping us to improve our understanding of biochemical processes. Inspired by the ideas of cybernetics in the sixties and seventies, dynamical systems theory and control theory have been increasingly applied to biochemical pathways. Thanks to powerful experimental techniques in genomics and proteomics, a wealth of biological data has accumulated and computational models of cells are now within reach. Systems biology, the discipline devoted to developing such models, uses biochemical networks as a main concept. It studies biological systems by investigating the network components and their interactions with the help of experimental high-throughput techniques and dedicated small-scale investigations and by integrating these data into networks and dynamical simulation models.

Like many new fields of research, systems biology started out with great expectations. High-throughput data and computational models were hoped to provide

answers to basic yet difficult biological questions: Why do we age? What processes control cell proliferation, and how? How do neurodegenerative disorders or diseases such as cancer develop? How can we engineer microbes more efficiently to produce valuable chemicals, fuels, or specific drugs? Only few of these goals have been achieved until now, and most of these questions remain on our agenda. Nevertheless, systems biology has greatly contributed to our understanding of cells and is increasingly becoming a standard part of biological research. It has fostered the formulation of new concepts and methods, such as statistical network analysis, the analysis of the robustness and fragility of dynamical systems, and the analysis of molecular noise. Even more importantly, it has enabled experimental biologists to realize that some scientific ideas cannot be easily expressed by words only. Inspired by electrical engineering, biologists now communicate the structure of biochemical systems by network graphics, which can then be translated into dynamical models.

This book gives an overview of systems biology as a rapidly developing field and provides readers with established and emerging tools and methods. You will learn how to formulate mathematical models of biological processes, how to analyze them, how to use experimental data and other types of knowledge to make models more precise, and how to interpret their simulation results. Based on our own experiences in teaching undergraduate and graduate students, the book is designed as an introductory course for students of biology, biophysics, and bioinformatics. It is as well useful for senior scholars who approach systems biology for the first time or seek more information about specific concepts and techniques. In the first chapters, we introduce stoichiometric and kinetic models, the main theoretical frameworks for metabolism and signaling pathways. We continue with methods for model construction (including model fitting, data handling, and model reduction) and related formalisms (spatial, discrete, and stochastic models). Then, we move on to experimental high-throughput techniques and to cellular networks. The analysis of regulation networks leads us to more general perspectives on cell physiology, including modularity, robustness, and optimality. The main part of the book

ends with a chapter on case studies. Addressing readers with different scientific backgrounds, we have added a reference section summarizing some basic knowledge of cell biology and mathematics, followed by a survey of popular biological databases and software tools. Further material is available on an accompanying Web site (<http://www.wiley-vch.de/home/systemsbiology>), which also contains solutions to the exercises presented in the book.

For the second edition of this book, we have updated and expanded the text to reflect advances in the field, and have reorganized the chapters to improve readability. Many of the changes reflect current developments in systems biology. On the one hand, the development of software tools is a very active area, where many new tools are developed, while others drop into oblivion. In the meantime, SBML has become an established exchange format for computational models in systems biology. We also notice that systems biology as a whole has become a mainstream discipline: High-throughput measurements have become an integral part of cell biology, computational models are used in research and teaching, and collaborations between experimentalists and theoreticians are increasingly common. Today, systems biology is perceived as what it is: the endeavor to understand complex processes in living organisms. Not more, but also not less!

We thank our friends and colleagues who helped us write this book. We are especially grateful to Mariapaola Gritti, Bernd Binder, Andreas Hoppe, Dagmar Waltemath, Elad Noor, Avi Flamholz, Terence Hwa, Ron Milo, Jonathan Karr, Ulrich Liebermeister, David Jesinghaus, Martina Fröhlich, and Severin Ehret for reading and commenting on the text. We thank the Max Planck Society for support and encouragement. We are grateful to the European Commission for funding via different European projects (UniCellSys, SystemeTb and FinSysB to EK, HeCaToS 602156), the German Ministry for Education and Research, BMBF (ViroSign, OncoPath, SysToxChip to EK), and the German Research Foundation (GRK 1772 to EK, LI 1676/2-1 to WL).

The book is dedicated to our teacher Prof. Dr. Reinhard Heinrich (1946–2006) whose work on metabolic control theory in the 1970s paved the way to systems biology and who greatly inspired our minds.

Guide to Different Topics of the Book**Biological systems and processes**

Metabolism (2, 3, 4, 11.1, 11.4, 12.1)
 Gene regulatory network (7, 8.2, 9)
 Gene expression regulation (2, 9)
 Signaling systems (8.2, 12.2)
 Cell cycle (12.3)
 Development (7.3)
 Aging (12.4)

Concepts for biological function

Qualitative behavior (3, 7.1)
 Parameter sensitivity and robustness (4.2, 10.2)
 Modularity and functional subsystems (6.4, 8.3)
 Robustness against failure (10.2)
 Information (10.3, 15.6)
 Population heterogeneity (10.3)
 Optimality (3.2, 11.1, 11.2)
 Evolution (11.3)
 Population dynamics and game theory (11.4)

Model types with different levels of abstraction

Statistical particle models (15.6)
 Stochastic biochemical models (5.1.2, 7.2, 9.4, 15.4)
 Kinetic models (4, 5.1.1, 11, 12, 16.7)
 Constraint-based models (3.2)
 Discrete models (7.1)
 Spatial models (7.3)

Mathematical frameworks to describe cell states

Topological (network structures) (3.1, 8)
 Structural stoichiometric models (3)
 Dynamical systems (4, 12, 15.2)
 Deterministic linear models (6.3)
 Deterministic kinetic models (4, 9.1, 12)
 Uncertain parameters (10.1)
 Optimization and control theory (4.2, 11.1, 11.2, 15.5)

Experimental Techniques

Experimental techniques (14)

Modeling skills

Model building (2, 3, 4, 5.1, 7)
 Model annotation (5.4)
 Parameter estimation (6.1)
 Model testing and selection (6.2)
 Local sensitivity analysis/control theory (4.2, 10.1, 10.2)
 Global sensitivity/uncertainty analysis (10.1)
 Model reduction (6.3)
 Model combination (5.4, 6.4)
 Network theory (8)
 Statistics (15.3, 15.7)
 Optimization of model outputs and structure (11.1)
 Optimal temporal control (11.2, 15.5)

Practical issues in modeling

Use of databases (5.3, 16)
 Data formats (5.2, 5.4)
 Data sources (5.3, 16)
 Modeling software (5.1, 17)
 Simulation techniques and tools (5.1)
 Model visualization (5.1)
 Data visualization (8.3)

About the Authors

Edda Klipp (born 1965) studied biophysics at the Humboldt-Universität zu Berlin and obtained a PhD of theoretical biophysics at the HU Berlin and an honorary doctor at Gothenburg University. She is professor for theoretical biophysics at HU Berlin. A founding member of the International Society for Systems Biology, her research interests include mathematical modeling of cellular systems, application of physical principles in biology, systems biology, and development of modeling tools.

Wolfram Liebermeister (born 1972) studied physics in Tübingen and Hamburg and obtained a PhD of theoretical biophysics at the Humboldt University of Berlin. In his work on complex biological systems, he highlights functional aspects such as variability, information, and optimality.

Christoph Wierling (born 1973) studied biology in Münster and holds a PhD in biochemistry obtained from the Free University Berlin. He was leading a research group for systems biology at the Max Planck Institute for Molecular Genetics, Berlin. Currently he is heading the bioinformatics and modeling unit at Alacris Theranostics, a Berlin-based company applying NGS and systems biology approaches for translational research and personalized medicine. His research interests focus on modeling and simulation of biological systems and the development of systems biology software.

Axel Kowald (born 1963) holds a PhD in mathematical biology from the National Institute for Medical Research, London. His current research interests focus on the mathematical modeling of processes involved in the biology of aging and systems biology.

Part One
Introduction to Systems Biology

Introduction

1

1.1 Biology in Time and Space

Biological systems such as organisms, cells, or biomolecules are highly organized in their structure and function. They have developed during evolution and can only be fully understood in this context. To study them and to apply mathematical, computational, or theoretical concepts, we have to be aware of the following circumstances.

The continuous reproduction of cell compounds necessary for living and the respective flow of information is captured by the central dogma of molecular biology, which can be summarized as follows: genes code for mRNA, mRNA serves as template for proteins, and proteins perform cellular work. Although information is stored in the genes encoded by the DNA sequence, it is made available only through the cellular machinery that can decode this sequence and can translate it into structure and function. In this book, we will explain that from various perspectives.

A description of biological entities and their properties encompasses different levels of organization and different time scales. We can study biological phenomena at the level of populations, individuals, tissues, organs, cells, and compartments down to molecules and atoms. Length scales range from the order of meter (e.g., the size of whale or human) to micrometer for many cell types, down to picometer for atom sizes. Time scales include millions of years for evolutionary processes, annual and daily cycles, seconds for many biochemical reactions, and femtoseconds for molecular vibrations. Figure 1.1 gives an overview about scales.

In a unified view of cellular networks, each action of a cell involves different levels of cellular organization, including genes, proteins, metabolism, or signaling pathways. Therefore, the current description of the individual networks must be integrated into a larger framework.

1.1 Biology in Time and Space

1.2 Models and Modeling

- What Is a Model?
- Purpose and Adequateness of Models
- Advantages of Computational Modeling

1.3 Basic Notions for Computational Models

- Model Scope
- Model Statements
- System State
- Variables, Parameters, and Constants
- Model Behavior
- Model Classification
- Steady States
- Model Assignment Is Not Unique

1.4 Networks

1.5 Data Integration

1.6 Standards

1.7 Model Organisms

- *Escherichia coli*
- *Saccharomyces cerevisiae*
- *Caenorhabditis elegans*
- *Drosophila melanogaster*
- *Mus musculus*

References

Further Reading

Many current approaches pay tribute to the fact that biological items are subject to evolution. The structure and organization of organisms and their cellular machinery has developed during evolution to fulfill major functions such as growth, proliferation, and survival under changing conditions. If parts of the organism or of the cell fail to perform their function, the individual might become unable to survive or replicate.

One consequence of evolution is the similarity of biological organisms of different species. This similarity

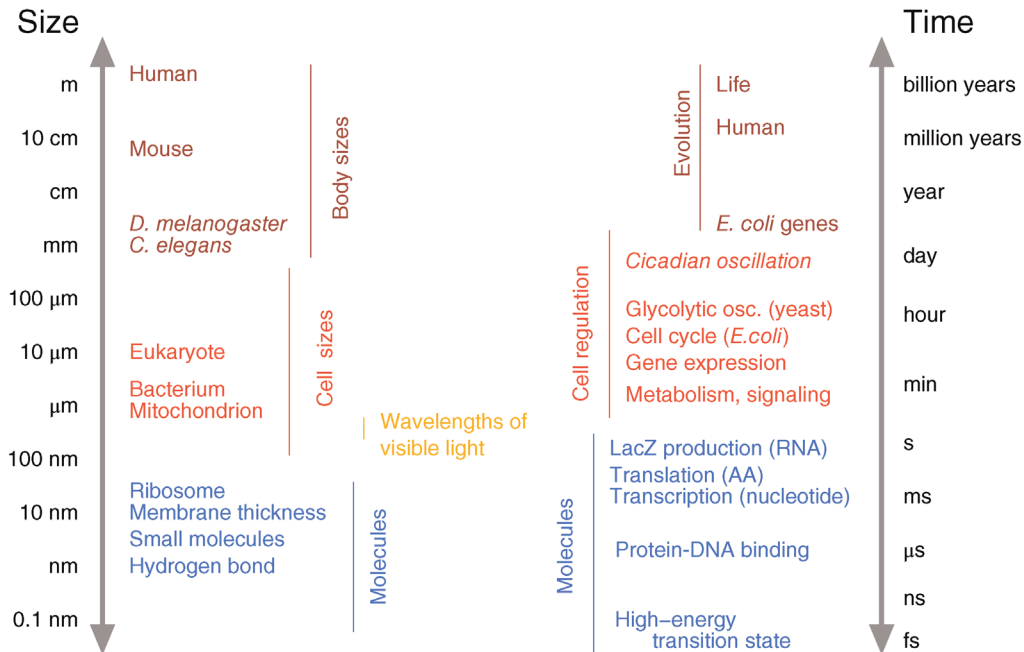


Figure 1.1 Length and time scales in biology. (Data from the BioNumbers database at bionumbers.hms.harvard.edu.)

allows for the use of model organisms and for the critical transfer of insights gained from one cell type to other cell types. Applications include, for example, prediction of protein function from similarity, prediction of network properties from optimality principles, reconstruction of phylogenetic trees, or the identification of regulatory DNA sequences through cross-species comparisons. However, the evolutionary process also leads to genetic variations within species. Therefore, personalized medicine and research is an important new challenge for biomedical research.

1.2 Models and Modeling

If we observe biological phenomena, we are confronted with various complex processes that often cannot be explained from first principles and the outcome of which cannot reliably be foreseen from intuition. Even if general biochemical principles are well established (e.g., the central dogma of transcription and translation or the biochemistry of enzyme-catalyzed reactions), the biochemistry of individual molecules and systems is often unknown and can vary considerably between species. Experiments lead to biological hypotheses about individual processes, but it often remains unclear whether these hypotheses can be combined into a larger coherent picture because it is often difficult to foresee the global

behavior of a complex system from knowledge of its parts. Mathematical modeling and computer simulations can help us to understand the internal nature and dynamics of these processes and to arrive at predictions about their future development and the effect of interactions with the environment.

1.2.1 What Is a Model?

The answer to this question will differ among communities of researchers. In a broad sense, a model is an abstract representation of objects or processes that explains features of these objects or processes (Figure 1.2). A biochemical reaction network can be represented by a graphical sketch showing dots for metabolites and arrows for reactions; the same network could also be described by a system of differential equations, which allows simulating and predicting the dynamic behavior of that network. If a model is used for simulations, it needs to be ensured that it faithfully predicts the system's behavior – at least those aspects that are supposed to be covered by the model. Systems biology models are often based on well-established physical laws that justify their general form, for instance, the thermodynamics of chemical reactions. Besides this, a computational model needs to make specific statements about a system of interest – which are partially justified by experiments and biochemical knowledge, and partially by mere extrapolation from other systems. Such a model can

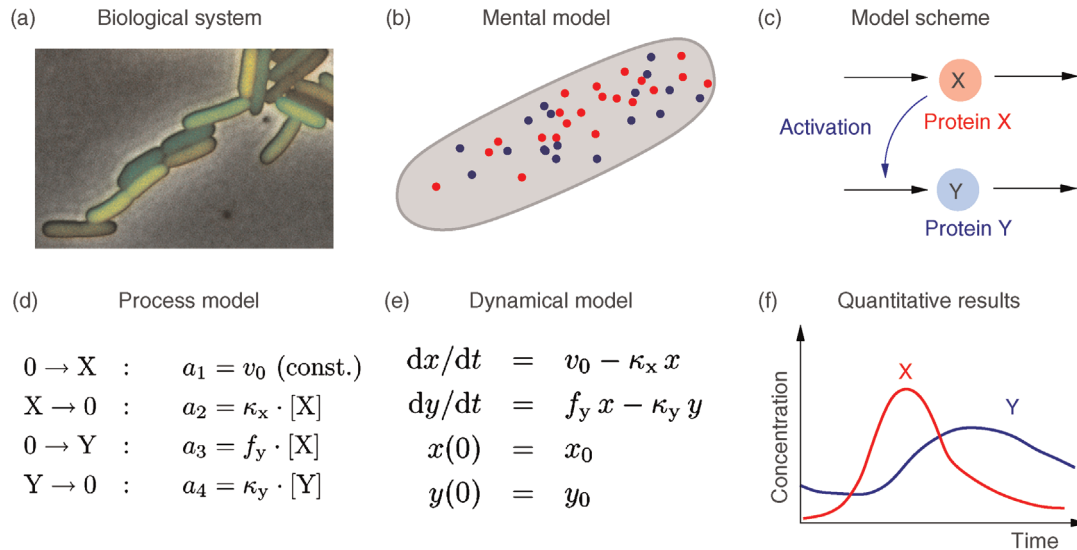


Figure 1.2 Typical abstraction steps in mathematical modeling. (a) *E. coli* bacteria produce thousands of different proteins. If a specific protein type is labeled with a fluorescent marker, cells glow under the microscope according to the concentration of this marker. (Courtesy of M. Elowitz.) (b) In a simplified mental model, we assume that cells contain two enzymes of interest, X (red) and Y (blue), and that the molecules (dots) can freely diffuse within the cell. All other substances are disregarded for the sake of simplicity. (c) The interactions between the two protein types can be drawn in a wiring scheme: each protein can be produced or degraded (black arrows). In addition, we assume that proteins of type X can increase the production of protein Y. (d) All individual processes to be considered are listed together with their rates a (occurrence per time). The mathematical expressions for the rates are based on a simplified picture of the actual chemical processes. (e) The list of processes can be translated into different sorts of dynamic models, in this case, deterministic rate equations for the protein concentrations x and y . (f) By solving the model equations, predictions for the time-dependent concentrations can be obtained. If the predictions do not agree with experimental data, this indicates that the model is wrong or too much simplified. In both cases, the model has to be refined.

summarize established knowledge about a system in a coherent mathematical formulation. In experimental biology, the term “model” is also used to denote a species that is especially suitable for experiments; for example, a genetically modified mouse may serve as a model for human genetic disorders.

1.2.2

Purpose and Adequateness of Models

Modeling is a subjective and selective procedure. A model represents only specific aspects of reality but, if done properly, this is sufficient since the intention of modeling is to answer particular questions. If the only aim is to predict system outputs from given input signals, a model should display the correct input–output relation, while its interior can be regarded as a black box. However, if instead a detailed biological mechanism has to be elucidated, then the system’s structure and the relations between its parts must be described realistically. Some models are meant to be generally applicable to many similar objects (e.g., Michaelis–Menten kinetics holds for many enzymes, the promoter–operator concept is applicable to many genes, and gene regulatory motifs are common), while others are specifically tailored to one

particular object (e.g., the 3D structure of a protein, the sequence of a gene, or a model of deteriorating mitochondria during aging). The mathematical part can be kept as simple as possible to allow for easy implementation and comprehensible results. Or it can be modeled very realistically and be much more complicated. None of the characteristics mentioned above makes a model wrong or right, but they determine whether a model is appropriate to the problem to be solved. The phrase “essentially, all models are wrong, but some are useful” coined by the statistician George Box is indeed an appropriate guideline for model building.

1.2.3

Advantages of Computational Modeling

Models gain their reference to reality from comparison with experiments, and their benefits therefore depend on the quality of the experiments used. Nevertheless, modeling combined with experimentation has a lot of advantages compared with purely experimental studies:

- Modeling drives conceptual clarification. It requires verbal hypotheses to be made specific and conceptually rigorous.

- Modeling highlights gaps in knowledge or understanding. During the process of model formulation, unspecified components or interactions have to be determined.
- Modeling provides independence of the modeled object.
- Time and space may be stretched or compressed *ad libitum*.
- Solution algorithms and computer programs can be used independently of the concrete system.
- Modeling is cheap compared with experiments.
- Models exert by themselves no harm on animals or plants and help to reduce ethical problems in experiments. They do not pollute the environment.
- Modeling can assist experimentation. With an adequate model, one may test different scenarios that are not accessible by experiment. One may follow time courses of compounds that cannot be measured in an experiment. One may impose perturbations that are not feasible in the real system. One may cause precise perturbations without directly changing other system components, which is usually impossible in real systems. Model simulations can be repeated often and for many different conditions.
- Model results can often be presented in precise mathematical terms that allow for generalization. Graphical representation and visualization make it easier to understand the system.
- Finally, modeling allows for making well-founded and testable predictions.

The attempt to formulate current knowledge and open problems in mathematical terms often uncovers a lack of knowledge and requirements for clarification. Furthermore, computational models can be used to test whether proposed explanations of biological phenomena are feasible. Computational models serve as repositories of current knowledge, both established and hypothetical, about how systems might operate. At the same time, they provide researchers with quantitative descriptions of this knowledge and allow them to simulate the biological process, which serves as a rigorous consistency test.

1.3 Basic Notions for Computational Models

1.3.1 Model Scope

Systems biology models consist of mathematical elements that describe properties of a biological system, for instance, mathematical variables describing the concentrations of

metabolites. As a model can only describe certain aspects of the system, all other properties of the system (e.g., concentrations of other substances or the environment of a cell) are neglected or simplified. It is important – and, to some extent, an art – to construct models in such ways that the disregarded properties do not compromise the basic results of the model.

1.3.2 Model Statements

Alongside the model elements, a model can contain various kinds of statements and equations describing facts about the model elements, most notably, their temporal behavior. In kinetic models, the basic modeling paradigm considered in this book, the dynamics is determined by a set of ordinary differential equations describing the substance balances. Statements in other model types may have the form of equality or inequality constraints (e.g., in flux balance analysis), maximality postulates, stochastic processes, or probabilistic statements about quantities that vary in time or between cells.

1.3.3 System State

In dynamical systems theory, a system is characterized by its *state*, a snapshot of the system at a given time. The state of the system is described by the set of variables that must be kept track of in a model: in deterministic models, it needs to contain enough information to predict the behavior of the system for all future times. Each modeling framework defines what is meant by the state of the system. In kinetic rate equation models, for example, the state is a list of substance concentrations. In the corresponding stochastic model, it is a probability distribution or a list of the current number of molecules of a species. In a Boolean model of gene regulation, the state is a string of bits indicating for each gene whether it is expressed (“1”) or not expressed (“0”). Also, the temporal behavior can be described in fundamentally different ways. In a *dynamical system*, the future states are determined by the current state, while in a *stochastic process*, the future states are not precisely predetermined. Instead, each possible future history has a certain probability to occur.

1.3.4 Variables, Parameters, and Constants

The quantities in a model can be classified as variables, parameters, and constants. A *constant* is a quantity with a fixed value, such as the natural number e or Avogadro’s number (number of molecules per mole). *Parameters* are

quantities that have a given value, such as the K_m value of an enzyme in a reaction. This value depends on the method used and on the experimental conditions and may change. *Variables* are quantities with a changeable value for which the model establishes relations. A subset of variables, the *state variables*, describes the system behavior completely. They can assume independent values and each of them is necessary to define the system state. Their number is equivalent to the dimension of the system. For example, the diameter d and volume V of a sphere obey the relation $V = \pi d^3/6$, where π and 6 are constants, V and d are variables, but only one of them is a state variable since the relation between them uniquely determines the other one.

Whether a quantity is a variable or a parameter depends on the model. In reaction kinetics, the enzyme concentration appears as a parameter. However, the enzyme concentration itself may change due to gene expression or protein degradation, and in an extended model, it may be described by a variable.

1.3.5 Model Behavior

Two fundamental factors that determine the behavior of a system are (i) influences from the environment (input) and (ii) processes within the system. The system structure, that is, the relation among variables, parameters, and constants, determines how endogenous and exogenous forces are processed. However, different system structures may still produce similar system behavior (output); therefore, measurements of the system output often do not suffice to choose between alternative models and to determine the system's internal organization.

1.3.6 Model Classification

For modeling, processes are classified with respect to a set of criteria.

- A structural or *qualitative* model (e.g., a network graph) specifies the interactions among model elements. A *quantitative* model assigns values to the elements and to their interactions, which may or may not change.
- In a *deterministic* model, the system evolution through all following states can be predicted from the knowledge of the current state. *Stochastic* descriptions give instead a probability distribution for the successive states.
- The nature of values that time, state, or space may assume distinguishes a *discrete* model (where values are taken from a discrete set) from a *continuous* model (where values belong to a continuum).

- *Reversible* processes can proceed in a forward and backward direction. Irreversibility means that only one direction is possible.
- *Periodicity* indicates that the system assumes a series of states in the time interval $\{t, t + \Delta t\}$ and again in the time interval $\{t + i\Delta t, t + (i + 1)\Delta t\}$ for $i = 1, 2, \dots$.

1.3.7 Steady States

The concept of stationary states is important for the modeling of dynamical systems. *Stationary states* (other terms are *steady states* or *fixed points*) are determined by the fact that the values of all state variables remain constant in time. The asymptotic behavior of dynamic systems, that is, the behavior after a sufficiently long time, is often stationary. Other types of asymptotic behavior are oscillatory or chaotic regimes.

The consideration of steady states is actually an abstraction that is based on a separation of time scales. In nature, everything flows. Fast and slow processes – ranging from formation and breakage of chemical bonds within nanoseconds to growth of individuals within years – are coupled in the biological world. While fast processes often reach a quasi-steady state after a short transition period, the change of the value of slow variables is often negligible in the time window of consideration. Thus, each steady state can be regarded as a quasi-steady state of a system that is embedded in a larger nonstationary environment. Despite this idealization, the concept of stationary states is important in kinetic modeling because it points to typical behavioral modes of the system under study and it often simplifies the mathematical problems.

Other theoretical concepts in systems biology are only rough representations of their biological counterparts. For example, the representation of gene regulatory networks by Boolean networks, the description of complex enzyme kinetics by simple mass action laws, or the representation of multifarious reaction schemes by black boxes proved to be helpful simplifications. Although being a simplification, these models elucidate possible network properties and help to check the reliability of basic assumptions and to discover possible design principles in nature. Simplified models can be used to test mathematically formulated hypotheses about system dynamics, and such models are easier to understand and to apply to different questions.

1.3.8 Model Assignment Is Not Unique

Biological phenomena can be described in mathematical terms. Models developed during the last few decades range from the description of glycolytic oscillations with

ordinary differential equations to population dynamics models with difference equations, stochastic equations for signaling pathways, and Boolean networks for gene expression. However, it is important to realize that a certain process can be described in more than one way: a biological object can be investigated with different experimental methods and each biological process can be described with different (mathematical) models. Sometimes, a modeling framework represents a simplified limiting case (e.g., kinetic models as limiting case of stochastic models). On the other hand, the same mathematical formalism may be applied to various biological instances: statistical network analysis, for example, can be applied to cellular transcription networks, the circuitry of nerve cells, or food webs.

The choice of a mathematical model or an algorithm to describe a biological object depends on the problem, the purpose, and the intention of the investigator. Modeling has to reflect essential properties of the system and different models may highlight different aspects of the same system. This ambiguity has the advantage that different ways of studying a problem also provide different insights into the system. However, the diversity of modeling approaches makes it also very difficult to merge established models (e.g., for individual metabolic pathways) into larger supermodels (e.g., models of complete cell metabolism).

1.4 Networks

The network is a crucial concept in systems biology. We study protein–protein interaction networks, protein–RNA interaction networks, metabolic networks (see Chapters 3 and 4 and Section 12.1), signaling networks (Section 12.2), guilt-by-association networks, and networks connecting gene defects with diseases or diseases with other diseases via common gene defects [1]. Throughout this book, you will find more examples.

Networks are best represented by graphs that consist of nodes and edges, which connect the nodes, as illustrated in Figure 1.3. In protein–protein interaction networks, for example, nodes are proteins and edges are their interactions as can for instance be determined by yeast two-hybrid experiments (see Chapter 14). If appropriate, one can introduce different types of nodes for different types of components. For example, the metabolites and converting enzymes in metabolic networks can be represented with bipartite networks, which possess two types of nodes – one for metabolites and the other for enzymes – that are never directly connected by an edge, but only via the other type of node. Petri net type of modeling

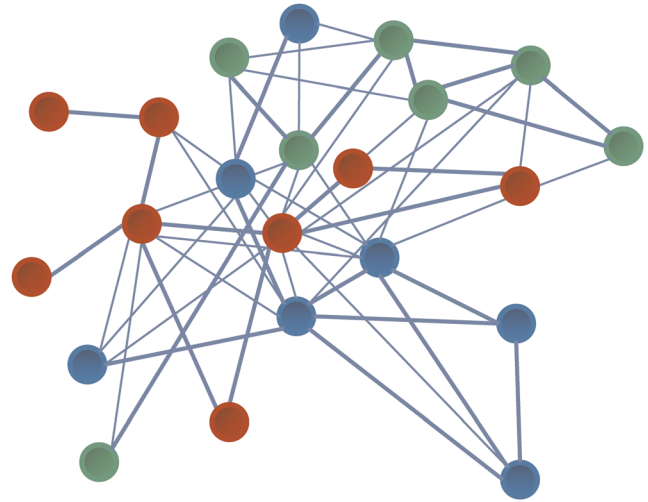


Figure 1.3 Network with nodes (circles) and edges (lines between circles). Different node colors indicate different types of connected components (e.g., proteins, mRNAs, and metabolites).

takes that property into account representing metabolites as places and enzyme-catalyzed reactions as transitions (see Section 7.1). By contrast, classical metabolic modeling considers only one type of node, but different types in different approaches. Systems of ordinary differential equations describing metabolite dynamics take metabolites as nodes and enzymatic reactions as edges (Chapter 4), while flux balance analysis restricts itself to steady states and now focusses on the fluxes through the reactions (now as nodes) that are linked by the stationary metabolites as edges.

1.5 Data Integration

Systems biology has evolved rapidly in the last few years, driven by the new high-throughput technologies. The most important impulse was given by large sequencing projects such as the Human Genome Project, which resulted in the full sequence of the human and other genomes [2,3]. Proteomic technologies have been used to identify the translation status of complete cells (2D gels, mass spectrometry) and to elucidate protein–protein interaction networks involving thousands of components [4]. However, to validate such diverse high-throughput data, one needs to correlate and integrate such information. Thus, an important part of systems biology is data integration.

On the lowest level of complexity, data integration implies common schemes for data storage, data representation, and data transfer. For particular experimental

techniques, this has already been established, for example, in the field of transcriptomics with Minimum Information About a Microarray Experiment [5], Minimum Information for Reporting Next Generation Sequence Genotyping [6], in proteomics with proteomics experiment data repositories [7], and the Human Proteome Organization consortium [8]. On a more complex level, schemes have been defined for biological models and pathways such as Systems Biology Markup Language (SBML) [9], CellML [10], or Systems Biology Graphical Notation (SBGN) [11], which all use an XML-like language style.

Data integration on the next level of complexity consists of data correlation. This is a growing research field as researchers combine information from multiple diverse data sets to learn about and explain natural processes [12,13]. For example, methods have been developed to integrate the results of transcriptome or proteome experiments with genome sequence annotations. In the case of complex disease conditions, it is clear that only integrated approaches can link clinical, genetic, behavioral, and environmental data with diverse types of molecular phenotype information and identify correlative associations. Such correlations, if found, are the key to identifying biomarkers and processes that are either causative or indicative of the disease. Importantly, the identification of biomarkers (e.g., proteins and metabolites) associated with the disease will open up the possibility to generate and test hypotheses on the biological processes and genes involved in this condition. The evaluation of disease-relevant data is a multistep procedure involving a complex pipeline of analysis and data handling tools such as data normalization, quality control, multivariate statistics, correlation analysis, visualization techniques, and intelligent database systems [14]. Several pioneering approaches have indicated the power of integrating data sets from different levels, for example, the correlation of gene membership of expression clusters and promoter sequence motifs [15], the combination of transcriptome and quantitative proteomics data in order to construct models of cellular pathways [13], and the identification of novel metabolite–transcript correlations [16]. Finally, data can be used to build and refine dynamical models, which represent an even higher level of data integration.

1.6 Standards

As experimental techniques generate rapidly growing amounts of data and large models need to be developed and exchanged, standards for both experimental procedures and modeling are a central practical issue in systems biology. Information exchange necessitates a

common language about biological aspects. One seminal example is the Gene Ontology that provides a controlled vocabulary that can be applied to all organisms, even as the knowledge about genes and proteins continues to accumulate. SBML [9] has been established as exchange language for mathematical models of biochemical reaction networks. SBGN [11] defines graphical elements to unambiguously represent biochemical reaction sets and large regulatory networks. A series of “minimum-information-about” statements based on community agreement defines standards for certain types of experiments. Minimum Information Requested in the Annotation of Biochemical Models (MIRIAM) [17] describes standards for this specific type of systems biology models. Minimum Information About a Simulation Experiment (MIASE) [18] helps authors to describe all elements of a computational experiment such that readers can repeat the simulations and create figures as shown in the publication.

1.7 Model Organisms

Model organisms are species that have developed over the years to be extremely popular for scientific investigations. The reasons for such popularity can be manifold. Of great importance is, of course, an easy handling of the organism, that is, culture conditions (temperature, pressure, etc.) that can be set up in the laboratory without much effort and that tolerate some degree of variation, so that results are comparable between groups that use slightly different growth conditions. However, other factors are also important, such as costs (for housing, food, etc.), size (the smaller the size, the more individuals can be studied), or lifespan (short-lived species are more popular for aging studies). Often model organisms are also used to represent important taxonomical properties (prokaryotes, eukaryotes, unicellular organisms, multicellular organisms, vertebrates, and invertebrates), but in all cases the hope is that important biochemical findings made in such model organisms are also of relevance for other species of that taxonomical group or even for humans. Figure 1.4 shows a selection of popular model species, which will be discussed in the next sections. They range from prokaryotic organisms to single and multicellular eukaryotic species up to mammals.

1.7.1 *Escherichia coli*

E. coli is probably the oldest and best studied model organism of all (Figure 1.4a). It is a rod-shaped



Figure 1.4 Popular model organisms for studies of problems in biochemistry and molecular biology. (a) *E. coli* is a rod-like bacterium and the best studied prokaryotic model system. (Public domain image from Wikimedia, http://commons.wikimedia.org/wiki/File:EscherichiaColi_NIAID.jpg.) (b) The yeast *S. cerevisiae* is a simple unicellular eukaryote and is of considerable scientific and industrial interest. (Public domain image from Wikimedia, http://commons.wikimedia.org/wiki/File:S_cerevisiae_under_DIC_microscopy.jpg.) (c) The nematode *C. elegans* is approximately 1 mm and is a popular representative for simple and short-lived multicellular organisms. (“Adult *Caenorhabditis elegans*” by Kbradnam (<http://en.wikipedia.org/wiki/User:Kbradnam>) is licensed under CC BY-SA-2.5, <http://creativecommons.org/licenses/by-sa/2.5>.) (d) The fruit fly *D. melanogaster* is like *C. elegans* a model for simple multicellular organisms and has extensively been studied in developmental biology. (“*Drosophila melanogaster*” by A. Karwath (<http://commons.wikimedia.org/wiki/User:Aka>) is licensed under CC BY-SA-2.5, <http://creativecommons.org/licenses/by-sa/2.5>.) (e) Finally, the mouse *M. musculus* is a popular model species for mammals and is thus also of great relevance for humans. (Public domain image from Wikimedia, http://commons.wikimedia.org/wiki/File:House_mouse.jpg.)

bacterium that is found in the intestines of many organisms, including humans. It is a facultative anaerobic organism, which means that it can grow under aerobic as well as anaerobic conditions. *E. coli* is roughly 2 μm long with a diameter of 0.5 μm . Under laboratory conditions, it can easily be cultivated and doubles its number in less than 30 min. It has been studied for more than 50 years and is the most popular prokaryotic model organism. The genome of the *E. coli* strain K-12 has completely been sequenced in 1997 [19] and contains around 4200 genes dispersed along 4.6 million base pairs (Mbp). It is a very streamlined genome containing very few intergenic sequences. The *E. coli* family consists of a large number of strains, and a comparison of the sequence of more than 60 strains has shown that they contain in total more than 15 500 genes, while

only 6% of this pan-genome is present in each strain [20]. *E. coli* was of pivotal importance for developing many of the experimental techniques described in Chapter 14. Today, a large number of scientific resources regarding this model species are available on the Internet. A good starting point is EcoCyc (ecocyc.org), which provides information about the genome and biochemical machinery of the *E. coli* strain K-12 MG1655. Other websites provide information about protein–protein interactions (bacteriome.org/) and systematic single-gene knockout mutants (<http://ecoli.aist-nara.ac.jp/gb6/Resources/deletion/deletion.html>), or a database of available strains (cgsc.biology.yale.edu). For modelers, the CyberCell Database (ccdb.wishartlab.com/CCDB) is also of interest since it aims at providing enzymatic, genetic, and biological information suitable

for developing mathematical models of all parts of a cell of *E. coli* strain K-12.

1.7.2

Saccharomyces cerevisiae

The yeast *S. cerevisiae* is a unicellular fungus, belonging to the ascomycetes (Figure 1.4b). It is not only a useful organism needed for the production of wine, beer, and bread, but also the best studied eukaryotic model system. The cells are easy to grow and double under optimal conditions every 90–100 min. Like *E. coli*, also *S. cerevisiae* can live under aerobic as well as anaerobic conditions. If oxygen is present, the majority of energy is generated via oxidative phosphorylation at the inner mitochondrial membrane and without oxygen energy is produced via glycolysis and fermentation. The yeast normally propagates as a diploid organism via mitosis. Under stress, however, the diploid cells can undergo sporulation, producing four haploid cells in the process. These haploid cells belong to one of two mating classes (sexes), called “a” and “ α ”. Haploids can either propagate via normal mitosis or mate with other haploids of the different mating class, resulting again in diploid cells. This life cycle makes *S. cerevisiae* interesting for genetic studies; it has also been extensively used by experimental and modeling studies of the cell cycle, glycolysis, osmotic shock, and mating process [21–28]. Cell division occurs in *S. cerevisiae* in an asymmetric fashion called budding and single-cell studies have shown that yeast cells exhibit replicative senescence with a maximum of 30–40 divisions [29]. Since this process is very reminiscent of the replicative senescence known from human fibroblasts [30], *S. cerevisiae* is also employed as a model system for investigations of the aging process. Furthermore, *S. cerevisiae* was also the first eukaryotic organism to be sequenced and its genome consists of about 12 Mbp containing roughly 6000 genes distributed over 16 chromosomes [31]. Homologous recombination (the exchange of sequences between similar strands of DNA) is very efficient in *S. cerevisiae*, which makes the organism also a convenient model for studies of synthetic biology. Using this mechanism, it was possible to replace the complete chromosome 16 with a new, synthetic one through 11 successive rounds of transformation (see Chapter 14) [32]. The synthetic chromosome was streamlined by removing all introns and superfluous tRNA genes and using only two of the three possible stop codons. This opens the possibility to extend the genetic code by a further amino acid once all chromosomes are modified in this way. A good online resource for further information about this model organism is the *Saccharomyces* Genome Database (www.yeastgenome.org).

1.7.3

Caenorhabditis elegans

Of course, model systems for multicellular organisms are also needed and the nematode *C. elegans* (Figure 1.4c) has become such a model since Sidney Brenner introduced it to the research community [33]. Like the other model organisms, it is easy to cultivate (feeding on bacteria or synthetic medium) and thousands of the about 1 mm long animals can live on a large Petri dish. Wild populations of *C. elegans* consist mainly of hermaphrodites together with a few males. Hermaphrodites not only are capable of self-fertilization (leading to natural inbred lines), but can also mate with males. The hermaphrodite then lays eggs that develop into larvae after hatching and after a total of four larval stages (L1–L4) the adult animal emerges. The complete life cycle from egg to egg takes between 2.5 and 5.5 days, depending on the temperature. The total lifespan of *C. elegans* is rather short with 2–3 weeks. This made *C. elegans* another popular model system for the investigation of the aging process [34]. However, the nematode is also an important model for other fields of research such as molecular biology or neurology. RNA interference (RNAi), for instance, is an important experimental technique (Chapter 14) that was developed based on experiments in *C. elegans* [35]. Furthermore, adult nematodes have a fixed number of somatic cells that is identical for all individuals (1031 in the male and 959 in the hermaphrodite), which makes it possible to generate very detailed anatomical models of the worm. The “slidable worm” (www.wormatlas.org/slidableworm.htm), which is a resource available on the webpage of the WormAtlas database, presents the results of such anatomical studies using an easy-to-use interface. *C. elegans* is also the only animal for which the complete wiring diagram (connectome) of the nervous system has been determined (using electron microscopy serial sections) [36,37]. Finally, *C. elegans* has also been the first multicellular organism for which the complete genome sequence has been determined [38,39]. The 97 Mbp contain approximately 19000 genes dispersed over six chromosomes. Good online starting points for more information are WormBase (www.wormbase.org), WormBook (www.wormbook.org/), or WormAtlas (www.wormatlas.org/).

1.7.4

Drosophila melanogaster

The fruit fly *D. melanogaster* (Figure 1.4d) is another, immensely popular, model organism that shares many of the properties of *C. elegans*. The animals are easy to breed in captivity and because of their small size (around 1 mm) it is possible to perform studies involving thousands of

individuals (e.g., for selection or population studies). The generation time (about 7 days at 29°C) and lifespan (about 30 days at 29°C) are very short and depend strongly on the ambient temperature. This facilitates, for example, artificial selection studies, which take several generations [40]. *D. melanogaster* has four chromosomes ($2n = 8$), which can even be studied under the light microscope because of a phenomenon called polyteny. As in many insect larvae, the cells of the salivary glands of *D. melanogaster* undergo multiple rounds of replication without cell division, leading to hundreds of sister chromatids aligned to each other. Polytene chromosomes are found in cells that need to express a large amount of a specific gene product and transcriptionally active areas appear under the microscope as swollen regions, so-called puffs. Although this technique is now outdated regarding the analysis of transcriptional activity, polytene chromosomes are still valuable for taxonomic problems. After staining, the puffs form a specific banding pattern that can be used to identify chromosomal deletions and duplications. This can be used in taxonomy to differentiate and classify closely related subspecies. *D. melanogaster* was arguably the most important model species for investigating developmental processes in multicellular organisms [41], which has led to the discovery of Hox genes [42]. These genes code for a set of transcription factors that contain a common 180 bp motif (the homeo-domain) and control the development of the anterior–posterior axis of the animal. A unique feature of these genes is that they are arranged on the chromosomes in the same linear order as the body region that they affect (called collinearity). Thus, Hox genes at one end of the cluster control the development of the anterior region (head), while the genes at the other end of the cluster influence the development of the posterior region (tail). Although originally found in *Drosophila*, Hox genes have been found in many metazoans, including vertebrates [43]. The complete genome was sequenced in 2000 [44] and somewhat surprisingly the number of genes is with approximately 14 000 clearly smaller than the number of genes in *C. elegans*. Further information, tools, and resources are available at FlyBase (flybase.org) and Ensembl Genome Browser (www.ensembl.org/Drosophila_melanogaster).

1.7.5

Mus musculus

The last model system that we want to introduce here is the house mouse *M. musculus domesticus* (Figure 1.4e). It is clearly the model organism with the largest similarity to humans and is therefore also of great relevance for

human research. Humans and mice are both mammals and thus share a common ancestor roughly 80 million years ago, a rather short time span compared with the other model organisms. Consequently, the genome structure and organization is also very similar. The mouse genome, sequenced in 2002 [45], contains 2.5 Gbp and is thus somewhat smaller than the human genome with 2.9 Gbp [2,3], although both genomes contain approximately 20 000–25 000 genes. The similarity at the gene level is quite amazing insofar that for more than 99% of mouse genes a homolog can also be found in the human genome [3], and vice versa. The mouse is also a popular model system because it is very amenable to genetic manipulations. The first mice were cloned in 1998 [46] and today it is common routine to create transgenic mice by introducing DNA constructs into fertilized egg cells and to study the function of existing genes by knocking them out or down (see Chapter 14). The Knockout Mouse Project (KOMP), for instance, aims at generating and providing mouse embryonic stem cells (and eventually whole mice) with single-gene knockout for every gene in the mouse genome (www.komp.org). Because mice have been used for such a long time as model species, many different inbred strains have been developed, which differ in various aspects of their phenotype (e.g., size, lifespan, and disease susceptibility). Of special interest are the various strains of nude mice that have a deletion of the FOXP1 gene, which prevents the formation of a functioning thymus. Without a thymus, these mice cannot produce mature T lymphocytes and therefore lack most forms of immune response (the lack of fur is a side effect of this mutation). As a consequence, they are valuable tools to study tumor development and are also used for transplantation studies, since they do not reject allo- or xenografts. Useful starting points for further information are, for instance, the Mouse Genome Informatics (www.informatics.jax.org/), the Mouse Atlas Project (www.emouseatlas.org), or the Ensembl Genome Browser (www.ensembl.org/Mus_musculus).

References

- 1 Goh, K.I., Cusick, M.E., Valle, D., Childs, B., Vidal, M., and Barabási, A.L. (2007) The human disease network. *Proc. Natl. Acad. Sci. USA*, 104 (21), 8685–8690.
- 2 Lander, E.S. *et al.* (2001) Initial sequencing and analysis of the human genome. *Nature*, 409, 860–921.
- 3 Venter, J.C. *et al.* (2001) The sequence of the human genome. *Science*, 291, 1304–1351.
- 4 Mering, C. *et al.* (2002) Comparative assessment of large-scale data sets of protein–protein interactions. *Nature*, 417, 399–403.
- 5 Brazma, A. *et al.* (2001) Minimum Information About a Microarray Experiment (MIAME): toward standards for microarray data. *Nat. Genet.*, 29, 365–371.