

METHODS IN MOLECULAR BIOLOGY

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

For further volumes:
<http://www.springer.com/series/7651>

The Gene Ontology Handbook

Edited by

Christophe Dessimoz

*Department of Genetics, Evolution & Environment, University College London, London, UK;
Swiss Institute of Bioinformatics, Lausanne, Switzerland; Department of Ecology and Evolution,
University of Lausanne, Lausanne, Switzerland; Center of Integrative Genomics, University of Lausanne,
Lausanne, Switzerland; Department of Computer Science, University College London, London, UK*

Nives Škunca

*Department of Computer Science, ETH Zurich, Zurich, Switzerland; SIB Swiss Institute of Bioinformatics,
Zurich, Switzerland; University College London, London, UK*

Editors

Christophe Dessimoz
Department of Genetics
Evolution and Environment
University College London
London, UK

Swiss Institute of Bioinformatics
Lausanne, Switzerland

Department of Ecology and Evolution
University of Lausanne
Lausanne, Switzerland

Center of Integrative Genomics
University of Lausanne
Lausanne, Switzerland

Department of Computer Science
University College London
London, UK

Nives Škunca
Department of Computer Science
ETH Zurich
Zurich, Switzerland

SIB Swiss Institute of Bioinformatics
Zurich, Switzerland

University College London
London, UK

ISSN 1064-3745

Methods in Molecular Biology

ISBN 978-1-4939-3741-7

DOI 10.1007/978-1-4939-3743-1

ISSN 1940-6029 (electronic)

ISBN 978-1-4939-3743-1 (eBook)

Library of Congress Control Number: 2016943478

© The Editor(s) (if applicable) and The Author(s) 2017. This book is published open access.

Open Access This book is distributed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits use, duplication, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, a link is provided to the Creative Commons license and any changes made are indicated.

The images or other third party material in this book are included in the work's Creative Commons license, unless indicated otherwise in the credit line; if such material is not included in the work's Creative Commons license and the respective action is not permitted by statutory regulation, users will need to obtain permission from the license holder to duplicate, adapt or reproduce the material.

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

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

Printed on acid-free paper

This Humana Press imprint is published by Springer Nature
The registered company is Springer Science+Business Media LLC New York

Preface

The Gene Ontology (GO) is the leading project to organize biological knowledge on genes and their products in a formal and consistent way across genomic resources. This has had a profound impact at several levels. First, such standardization has made possible the integration of multiple resources and sources of knowledge, thereby increasing their discoverability and simplifying their usage. Second, it has greatly facilitated—some might say *exceedingly so*—data mining, aggregate analyses, and other forms of automated knowledge extraction. Third, it has led to an increase in the overall quality of the resources by enforcing minimum requirements across all of them.

Even considering these advantages, the rapid adoption of the GO in the community has been remarkable. In the 15 years since the publication of its introductory article [1], over 100,000 scientific articles containing the keyword “Gene Ontology” have been published and the rate is still increasing (Google Scholar).

However, despite this popularity and widespread use, many aspects of the Gene Ontology remain poorly understood [2], at times even by experts [3]. For instance, unbeknownst to most users, routine procedures such as GO term enrichment analyses remain subject to biases and simplifying assumptions that can lead to spurious conclusions [4].

The objective of this book is to provide a practical, self-contained overview of the GO for biologists and bioinformaticians. After reading this book, we would like the reader to be equipped with the essential knowledge to use the GO and correctly interpret results derived from it. In particular, the book will cover the state of the art of how GO annotations are made, how they are evaluated, and what sort of analyses can and cannot be done with the GO. In the spirit of the *Methods in Molecular Biology* book series in which it appears, there is an emphasis on providing practical guidance and troubleshooting advice.

The book is intended for a wide scientific audience and makes few assumptions about prior knowledge. While the primary target is the nonexpert, we also hope that seasoned GO users and contributors will find it informative and useful. Indeed, we are the first to admit that working with the GO occasionally brings to mind the aphorism “the more we know, the less we understand.”

The book is structured in six main parts. Part I introduces the reader to the fundamental concepts underlying the Gene Ontology project, with primers on ontologies in general (Chapter 1), on gene function (Chapter 2), and on the Gene Ontology itself (Chapter 3).

To become proficient GO users, we need to know where the GO data comes from. Part II reviews how the GO annotations are made, be it via manual curation of the primary literature (Chapter 4), via computational methods of function inference (Chapter 5), via literature text mining (Chapter 6), or via crowdsourcing and other contributions from the community (Chapter 7).

But can we trust these annotations? In Part III, we consider the problem of evaluating GO annotations. We first provide an overview of the different approaches, the challenges associated with them, but also some successful initiatives (Chapter 8). We then focus on the more specific problem of evaluating enzyme function predictions (Chapter 9). Last, we

reflect on the achievements of the Critical Assessment of protein Function Annotation (CAFA) community experiment (Chapter 10).

Having made and validated GO annotations, we proceed in Part IV to use the GO resource. We consider the various ways of retrieving GO data (Chapter 11), how to quantify the functional similarity of GO terms and genes (Chapter 12), or perform GO enrichment analyses (Chapter 13)—all the while avoiding common biases and pitfalls (Chapter 14). The part ends with a chapter on visualizing GO data (Chapter 15) as well as a tutorial on GO analyses in the programming language Python (Chapter 16).

Part V covers two advanced topics: annotation extensions, which make it possible to express relationships involving multiple terms (Chapter 17), and the evidence code ontology, which provides a more precise and expressive specification of supporting evidence than the traditional GO annotation evidence codes (Chapter 18).

Part VI goes beyond the GO, by considering complementary sources of functional information such as KEGG and Enzyme Commission numbers (Chapter 19), and by considering the potential of integrating GO with controlled clinical nomenclatures (Chapter 20).

The final part concludes the book with a perspective by Suzi Lewis on the past, present, and future of the GO (Chapter 21).

London, UK
Zurich, Switzerland

Christophe Dessimoz
Nives Škunca

References

1. Ashburner M, Ball CA, Blake JA et al (2000) Gene ontology: tool for the unification of biology. The Gene Ontology Consortium. *Nat Genet* 25:25–29
2. Thomas PD, Wood V, Mungall CJ et al (2012) On the use of gene ontology annotations to assess functional similarity among orthologs and paralogs: a short report. *PLoS Comput Biol* 8:e1002386
3. Dessimoz C, Škunca N, Thomas PD (2013) CAFA and the open world of protein function predictions. *Trends Genet* 29:609–610
4. Tipney H, Hunter L (2010) An introduction to effective use of enrichment analysis software. *Hum Genomics* 4:202–206

Acknowledgements

We thank all chapter authors for their contributions to the book. We are particularly indebted to Pascale Gaudet and Ruth Lovering for contributing multiple chapters and demonstrating unabated enthusiasm throughout the process. All chapters were reviewed by at least two independent peers, which represents a considerable effort. Peer reviews were mostly contributed by chapter authors, but also by the following people: Adrian Altenhoff, Natasha Glover, Debra Klopfenstein, Chris Mungall, Prudence Mutowo, Marc Robinson-Rechavi, Kimberly Van Auken, and Haibao Tang. Last but not least, we thank support by John Walker, our series editor, and Patrick Marton from Springer.

Funding for the Open Access charges was generously provided by the University College London Library, the Swiss Institute of Bioinformatics, the Agassiz Foundation, and the Foundation for the University of Lausanne.

Contents

<i>Preface</i>	<i>v</i>
<i>Contributors</i>	<i>xi</i>
PART I FUNDAMENTALS	
1 Primer on Ontologies <i>Janna Hastings</i>	3
2 The Gene Ontology and the Meaning of Biological Function <i>Paul D. Thomas</i>	15
3 Primer on the Gene Ontology <i>Pascale Gaudet, Nives Škunca, James C. Hu, and Christophe Dessimoz</i>	25
PART II MAKING GENE ONTOLOGY ANNOTATIONS	
4 Best Practices in Manual Annotation with the Gene Ontology <i>Sylvain Poux and Pascale Gaudet</i>	41
5 Computational Methods for Annotation Transfers from Sequence <i>Domenico Cozzetto and David T. Jones</i>	55
6 Text Mining to Support Gene Ontology Curation and Vice Versa <i>Patrick Ruch</i>	69
7 How Does the Scientific Community Contribute to Gene Ontology? <i>Ruth C. Lovering</i>	85
PART III EVALUATING GENE ONTOLOGY ANNOTATIONS	
8 Evaluating Computational Gene Ontology Annotations <i>Nives Škunca, Richard J. Roberts, and Martin Steffen</i>	97
9 Evaluating Functional Annotations of Enzymes Using the Gene Ontology <i>Gemma L. Holliday, Rebecca Davidson, Eyal Akiva, and Patricia C. Babbitt</i>	111
10 Community-Wide Evaluation of Computational Function Prediction <i>Iddo Friedberg and Predrag Radivojac</i>	133
PART IV USING THE GENE ONTOLOGY	
11 Get GO! Retrieving GO Data Using AmiGO, QuickGO, API, Files, and Tools <i>Monica Munoz-Torres and Seth Carbon</i>	149
12 Semantic Similarity in the Gene Ontology <i>Catia Pesquita</i>	161

13	Gene-Category Analysis	175
	<i>Sebastian Bauer</i>	
14	Gene Ontology: Pitfalls, Biases, and Remedies.	189
	<i>Pascale Gaudet and Christophe Dessimoz</i>	
15	Visualizing GO Annotations	207
	<i>Fran Supek and Nives Škunca</i>	
16	A Gene Ontology Tutorial in Python	221
	<i>Alex Warwick Vesztrocy and Christophe Dessimoz</i>	
PART V ADVANCED GENE ONTOLOGY TOPICS		
17	Annotation Extensions	233
	<i>Rachael P. Huntley and Ruth C. Lovering</i>	
18	The Evidence and Conclusion Ontology (ECO): Supporting GO Annotations	245
	<i>Marcus C. Chibucos, Deborah A. Siegele, James C. Hu, and Michelle Giglio</i>	
PART VI BEYOND THE GENE ONTOLOGY		
19	Complementary Sources of Protein Functional Information: The Far Side of GO	263
	<i>Nicholas Furnham</i>	
20	Integrating Bio-ontologies and Controlled Clinical Terminologies: From Base Pairs to Bedside Phenotypes.	275
	<i>Spiros C. Denaxas</i>	
PART VII CONCLUSION		
21	The Vision and Challenges of the Gene Ontology.	291
	<i>Suzanna E. Lewis</i>	
	<i>Index</i>	303

Contributors

- EYAL AKIVA • *Department of Bioengineering and Therapeutic Sciences, University of California San Francisco, San Francisco, CA, USA*
- PATRICIA C. BABBITT • *Department of Bioengineering and Therapeutic Sciences, University of California San Francisco, San Francisco, CA, USA*
- SEBASTIAN BAUER • *PRIVATE, Berlin, Germany*
- SETH CARBON • *Berkeley Bioinformatics Open-Source Projects, Environmental Genomics and Systems Biology Division, Lawrence Berkeley National Laboratory, Berkeley, CA, USA*
- MARCUS C. CHIBUCOS • *Department of Microbiology and Immunology, Institute for Genome Sciences, University of Maryland School of Medicine, Baltimore, MD, USA*
- DOMENICO COZZETTO • *Bioinformatics Group, Department of Computer Science, University College London, London, UK*
- REBECCA DAVIDSON • *Department of Bioengineering and Therapeutic Sciences, University of California San Francisco, San Francisco, CA, USA*
- SPIROS C. DENAXAS • *Farr Institute of Health Informatics Research, University College London, London, UK; Institute of Health Informatics, University College London, London, UK*
- CHRISTOPHE DESSIMOZ • *Department of Genetics, Evolution & Environment, University College London, London, UK; Swiss Institute of Bioinformatics, Lausanne, Switzerland; Department of Ecology and Evolution, University of Lausanne, Lausanne, Switzerland; Center of Integrative Genomics, University of Lausanne, Lausanne, Switzerland; Department of Computer Science, University College London, London, UK*
- IDD0 FRIEDBERG • *Department of Veterinary Microbiology and Preventive Medicine, Iowa State University, Ames, IA, USA*
- NICHOLAS FURNHAM • *Department of Pathogen Molecular Biology, London School of Hygiene and Tropical Medicine, London, UK*
- PASCALE GAUDET • *CALIPHO Group, Swiss Institute of Bioinformatics, Geneva, Switzerland; Department of Human Protein Sciences, Faculty of Medicine, University of Geneva, Geneva, Switzerland*
- MICHELLE GIGLIO • *Department of Medicine, Institute for Genome Sciences, University of Maryland School of Medicine, Baltimore, MD, USA*
- JANNA HASTINGS • *Cheminformatics and Metabolism, European Molecular Biology Laboratory, European Bioinformatics Institute (EMBL-EBI), Cambridgeshire, UK*
- GEMMA L. HOLLIDAY • *Department of Bioengineering and Therapeutic Sciences, University of California San Francisco, San Francisco, CA, USA*
- JAMES C. HU • *Department of Biochemistry and Biophysics, Texas A&M University and Texas AgriLife Research, College Station, TX, USA*
- RACHAEL P. HUNTLEY • *Functional Gene Annotation Initiative, Centre for Cardiovascular Genetics, Institute of Cardiovascular Science, University College London, London, UK*
- DAVID T. JONES • *Bioinformatics Group, Department of Computer Science, University College London, London, UK*

- SUZANNA E. LEWIS • *Lawrence Berkeley National Laboratory, Berkeley, CA, USA*
- RUTH C. LOVERING • *Functional Gene Annotation Initiative, Centre for Cardiovascular Genetics, Institute of Cardiovascular Science, University College London, London, UK*
- MONICA MUNOZ-TORRES • *Berkeley Bioinformatics Open-Source Projects, Environmental Genomics and Systems Biology Division, Lawrence Berkeley National Laboratory, Berkeley, CA, USA*
- CATIA PESQUITA • *LaSIGE, Faculdade de Ciências, Universidade de Lisboa, Lisbon, Portugal*
- SYLVAIN POUX • *Swiss-Prot Group, SIB Swiss Institute of Bioinformatics, Centre Medical Universitaire, Geneva, Switzerland*
- PREDRAG RADIVOJAC • *Department of Computer Science and Informatics, Indiana University, Bloomington, IN, USA*
- RICHARD J. ROBERTS • *New England Biolabs, Ipswich, MA, USA*
- PATRICK RUCH • *SIB Text Mining, Swiss Institute of Bioinformatics, Geneva, Switzerland; BiTeM Group, HES-SO\HEG Genève, Carouge, Switzerland*
- DEBORAH A. SIEGELE • *Department of Biology, Texas A&M University, College Station, TX, USA*
- NIVES ŠKUNCA • *Department of Computer Science, ETH Zurich, Zurich, Switzerland; SIB Swiss Institute of Bioinformatics, Zurich, Switzerland; University College London, London, UK*
- MARTIN STEFFEN • *Department of Biomedical Engineering, Boston University, Boston, MA, USA; Department of Pathology and Laboratory Medicine, Boston University School of Medicine, Boston, MA, USA*
- FRAN SUPEK • *Division of Electronics, Ruder Boskovic Institute, Zagreb, Croatia; EMBL/CRG Systems Biology Research Unit, Centre for Genomic Regulation (CRG), The Barcelona Institute of Science and Technology, Barcelona, Spain; Universitat Pompeu Fabra (UPF), Barcelona, Spain*
- PAUL D. THOMAS • *Division of Bioinformatics, Department of Preventive Medicine, University of Southern California, Los Angeles, CA, USA*
- ALEX WARWICK VESZTROCY • *Department of Genetics, Evolution and Environment, University College London, London, UK; Swiss Institute of Bioinformatics, Lausanne, Switzerland*