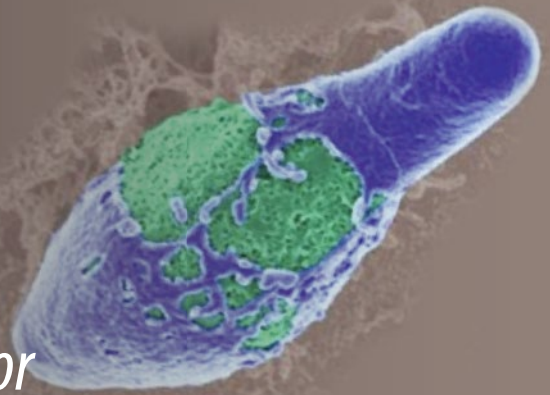


Methods in
Molecular Biology 1140

Springer Protocols



Wayne F. Anderson *Editor*

Structural Genomics and Drug Discovery

Methods and Protocols

 Humana Press

METHODS IN MOLECULAR BIOLOGY

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

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

Structural Genomics and Drug Discovery

Methods and Protocols

Edited by

Wayne F. Anderson

*Center for Structural Genomics of Infectious Diseases,
Midwest Center for Structural Genomics,
Northwestern University Feinberg School of Medicine, Chicago, IL, USA*

 **Humana Press**

Editor

Wayne F. Anderson
Center for Structural Genomics
of Infectious Diseases
Midwest Center for
Structural Genomics
Northwestern University Feinberg
School of Medicine
Chicago, IL, USA

ISSN 1064-3745 ISSN 1940-6029 (electronic)
ISBN 978-1-4939-0353-5 ISBN 978-1-4939-0354-2 (eBook)
DOI 10.1007/978-1-4939-0354-2
Springer New York Heidelberg Dordrecht London

Library of Congress Control Number: 2014931091

© Springer Science+Business Media New York 2014

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed. Exempted from this legal reservation are brief excerpts in connection with reviews or scholarly analysis or material supplied specifically for the purpose of being entered and executed on a computer system, for exclusive use by the purchaser of the work. Duplication of this publication or parts thereof is permitted only under the provisions of the Copyright Law of the Publisher's location, in its current version, and permission for use must always be obtained from Springer. Permissions for use may be obtained through RightsLink at the Copyright Clearance Center. Violations are liable to prosecution under the respective Copyright Law.

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.

While the advice and information in this book are believed to be true and accurate at the date of publication, neither the authors nor the editors nor the publisher can accept any legal responsibility for any errors or omissions that may be made. The publisher makes no warranty, express or implied, with respect to the material contained herein.

Printed on acid-free paper

Humana Press is a brand of Springer
Springer is part of Springer Science+Business Media (www.springer.com)

Dedication

For Caroline, who has been a wonderful partner for nearly 50 years and made this, as well as everything else, possible.

Preface

“Structural Genomics” as an area of investigation arose from the recognition that genome sequence information could be combined with improved methods for macromolecular structure determination to allow high-throughput structure determination. One of the early justifications for developing the field was the potential to make use of the structural information in drug discovery efforts. All three of these areas, genome sequencing, macromolecular structure determination, and structure-aided drug discovery, have seen dramatic improvements in technology and methodology.

This volume focuses on high-throughput structure determination methods and how they can be applied to lay the groundwork for structure-aided drug discovery. The methods and protocols that are described can be applied in any laboratory interested in using detailed structural information to advance the initial stages of drug discovery. Due to the advances in technology and methodology that have occurred during the past 10–15 years, even the nonspecialist can apply structural biology to most biomedical problems. The methods and approaches that distinguish structural genomics from “classical” structural biology have been decreasing as more and more research groups adopt high-throughput methods and apply them to their specific biological research problems.

In some respects, structure-aided drug discovery is very specific to the one particular protein target being studied and the approaches of structural genomics would not seem to be appropriate. However, if one looks at the problem broadly, there often is more than one protein that could be targeted, and when multiple proteins are being investigated, the advantages of carrying out most of the steps in parallel can increase productivity.

The initial chapters deal with bioinformatics and data management because selecting target proteins and planning how the large amount of diverse data will be handled are the first steps. Following these are the chapters on high-throughput methods for cloning, expression and solubility testing, protein production, purification, crystallization screening, and screening for suitability for NMR structure determination. One of the continuing problems faced by structural genomics efforts is the limited success rate that, not surprisingly, accompanies increased throughput and the associated reduction in individual attention to each protein. Although there is no panacea, a number of chapters describe methods that can rescue, or salvage, target proteins that are failing as they proceed through the pipeline. Finally, the concluding chapters describe methods that use the proteins that have been produced in order to identify initial small molecule hits. These hits can then feed into drug discovery efforts. At this point in the process, the number of technically and biologically suitable targets will have been reduced and each protein, together with the hits that have been generated, will require individual attention.

The structural genomics approach provides an efficient initial step toward drug discovery and the methods described will be useful to anyone interested in moving in this direction.

Chicago, IL, USA

Wayne F. Anderson

Contents

<i>Preface</i>	<i>vii</i>
<i>Contributors</i>	<i>xi</i>
1 Data Management in the Modern Structural Biology and Biomedical Research Environment	1
<i>Matthew D. Zimmerman, Marek Grabowski, Marcin J. Domagalski, Elizabeth M. MacLean, Maksymilian Chruszcz, and Wlodek Minor</i>	
2 Structural Genomics of Human Proteins	27
<i>Khan Tanjid Osman and Aled Edwards</i>	
3 Target Selection for Structural Genomics of Infectious Diseases.	35
<i>Corin Yeats, Benoit H. Dessailly, Elizabeth M. Glass, Daved H. Fremont, and Christine A. Orengo</i>	
4 Selecting Targets from Eukaryotic Parasites for Structural Genomics and Drug Discovery	53
<i>Isabelle Q.H. Phan, Robin Stacy, and Peter J. Myler</i>	
5 High-Throughput Cloning for Biophysical Applications	61
<i>Keehwan Kwon and Scott N. Peterson</i>	
6 Expression and Solubility Testing in a High-Throughput Environment	75
<i>Keehwan Kwon and Scott N. Peterson</i>	
7 Protein Production for Structural Genomics Using <i>E. coli</i> Expression	89
<i>Magdalena Makowska-Grzyska, Youngchang Kim, Natalia Maltseva, Hui Li, Min Zhou, Grazyna Joachimiak, Gyorgy Babnigg, and Andrzej Joachimiak</i>	
8 Eukaryotic Expression Systems for Structural Studies	107
<i>Christopher A. Nelson, William H. McCoy, and Daved H. Fremont</i>	
9 Automated Cell-Free Protein Production Methods for Structural Studies	117
<i>Emily T. Beebe, Shin-ichi Makino, John L. Markley, and Brian G. Fox</i>	
10 Parallel Protein Purification	137
<i>Ludmilla Shuvalova</i>	
11 Oxidative Refolding from Inclusion Bodies	145
<i>Christopher A. Nelson, Chung A. Lee, and Daved H. Fremont</i>	
12 High-Throughput Crystallization Screening	159
<i>Tatiana Skarina, Xiaohui Xu, Elena Evdokimova, and Alexei Savchenko</i>	
13 Screening Proteins for NMR Suitability.	169
<i>Adelinda A. Yee, Anthony Semesi, Maite Garcia, and Cheryl H. Arrowsmith</i>	

14	Salvage or Recovery of Failed Targets by In Situ Proteolysis	179
	<i>Yufeng Tong, Aiping Dong, Xiaohui Xu, and Amy Wernimont</i>	
15	Salvage of Failed Protein Targets by Reductive Alkylation	189
	<i>Kemin Tan, Youngchang Kim, Catherine Hatzos-Skintges, Changsoo Chang, Marianne Cuff, Gekleng Chhor, Jerzy Osipiuk, Karolina Michalska, Boguslaw Nocek, Hao An, Gyorgy Babnigg, Lance Bigelow, Grazyna Joachimiak, Hui Li, Jamey Mack, Magdalena Makowska-Grzyska, Natalia Maltseva, Rory Mulligan, Christine Tesar, Min Zhou, and Andrzej Joachimiak</i>	
16	Salvage or Recovery of Failed Targets by Mutagenesis to Reduce Surface Entropy	201
	<i>Lukasz Goldschmidt, David Eisenberg, and Zygmunt S. Derewenda</i>	
17	Data Collection for Crystallographic Structure Determination.	211
	<i>Kanagalaghatta Rajashankar and Zbigniew Dauter</i>	
18	Structure Determination, Refinement, and Validation	239
	<i>George Minasov and Wayne F. Anderson</i>	
19	Virtual High-Throughput Ligand Screening	251
	<i>T. Andrew Binkowski, Wei Jiang, Benoit Roux, Wayne F. Anderson, and Andrzej Joachimiak</i>	
20	Ligand Screening Using Fluorescence Thermal Shift Analysis (FTS)	263
	<i>Chi-Hao Luan, Samuel H. Light, Sara F. Dunne, and Wayne F. Anderson</i>	
21	Ligand Screening Using Enzymatic Assays	291
	<i>Kiira Ratia, Shahila Mehboob, and Hyun Lee</i>	
22	Ligand Screening Using NMR	305
	<i>Benjamin E. Ramirez, Aleksandar Antanasijevic, and Michael Caffrey</i>	
23	Screening Ligands by X-ray Crystallography	315
	<i>Douglas R. Davies</i>	
24	Case Study—Structural Genomics and Human Protein Kinases	325
	<i>Jonathan M. Elkins</i>	
	<i>Index</i>	337

Contributors

- HAO AN • *Biosciences Division, Midwest Center for Structural Genomics, Argonne National Laboratory, Argonne, IL, USA*
- WAYNE F. ANDERSON • *Center for Structural Genomics of Infectious Diseases, Northwestern University Feinberg School of Medicine, Chicago, IL, USA; Midwest Center for Structural Genomics, Northwestern University Feinberg School of Medicine, Chicago, IL, USA*
- ALEKSANDAR ANTANASIJEVIC • *Department of Biochemistry and Molecular Genetics, University of Illinois at Chicago, Chicago, IL, USA*
- CHERYL H. ARROWSMITH • *Division of Cancer Genomics and Proteomics, Northeast Structural Genomics Consortium (NESG), Ontario Cancer Institute, Toronto, ON, Canada*
- GYORGY BABNIGG • *Center for Structural Genomics of Infectious Diseases, Computational Institute, University of Chicago, Chicago, IL, USA; Biosciences Division, Midwest Center for Structural Genomics, Argonne National Laboratory, Argonne, IL, USA*
- EMILY T. BEEBE • *Department of Biochemistry, University of Wisconsin-Madison, Madison, WI, USA*
- LANCE BIGELOW • *Biosciences Division, Midwest Center for Structural Genomics, Argonne National Laboratory, Argonne, IL, USA*
- T. ANDREW BINKOWSKI • *Computation Institute, Center for Structural Genomics of Infectious Diseases, University of Chicago, Chicago, IL, USA*
- MICHAEL CAFFREY • *Department of Biochemistry and Molecular Genetics, University of Illinois at Chicago, Chicago, IL, USA*
- CHANGSOO CHANG • *Biosciences Division, Midwest Center for Structural Genomics, Argonne National Laboratory, Argonne, IL, USA; Biosciences Division, Structural Biology Center, Argonne National Laboratory, Argonne, IL, USA*
- GEKLENG CHHOR • *Biosciences Division, Midwest Center for Structural Genomics, Argonne National Laboratory, Argonne, IL, USA*
- MAKSYMILIAN CHRUSZCZ • *Department of Chemistry and Biochemistry, University of South Carolina, Columbia, SC, USA*
- MARIANNE CUFF • *Biosciences Division, Midwest Center for Structural Genomics, Argonne National Laboratory, Argonne, IL, USA; Biosciences Division, Structural Biology Center, Argonne National Laboratory, Argonne, IL, USA*
- ZBIGNIEW DAUTER • *Synchrotron Radiation Research Section, Macromolecular Crystallography Laboratory, National Cancer Institute, Argonne National Laboratory, Argonne, IL, USA*
- DOUGLAS R. DAVIES • *Emerald Bio, Bainbridge Island, WA, USA*
- ZYGMUNT S. DEREWENDA • *Department of Molecular Physiology and Biological Physics, University of Virginia School of Medicine, Charlottesville, VA, USA*
- BENOIT H. DESSAILLY • *Center for Structural Genomics of Infectious Diseases, Department of Structural and Molecular Biology, University College London, London, UK*

- MARCIN J. DOMAGALSKI • *Department of Molecular Physiology and Biological Physics, University of Virginia, Charlottesville, VA, USA; Center for Structural Genomics of Infectious Diseases (CSGID), University of Virginia, Charlottesville, VA, USA; Midwest Center for Structural Genomics (MCSG), University of Virginia, Charlottesville, VA, USA; New York Structural Genomics Research Consortium (NYSGRC), University of Virginia, Charlottesville, VA, USA; Enzyme Function Initiative (EFI), University of Virginia, Charlottesville, VA, USA*
- AIPING DONG • *Structural Genomics Consortium, University of Toronto, Toronto, ON, Canada*
- SARA F. DUNNE • *Center for Structural Genomics of Infectious Diseases, High Throughput Analysis Laboratory and Department of Molecular Biosciences, Northwestern University, Evanston, IL, USA*
- ALED EDWARDS • *Structural Genomics Consortium, University of Toronto, Toronto, ON, Canada; Center for Structural Genomics of Infectious Diseases, Toronto, ON, Canada*
- DAVID EISENBERG • *UCLA-DOE Institute for Genomics and Proteomics, Howard Hughes Medical Institute, University of California, Los Angeles, CA, USA*
- JONATHAN M. ELKINS • *Nuffield Department of Clinical Medicine, Structural Genomics Consortium, University of Oxford, Oxford, UK*
- ELENA EVDOKIMOVA • *Department of Chemical Engineering and Applied Chemistry, University of Toronto, Toronto, ON, Canada*
- BRIAN G. FOX • *Department of Biochemistry, University of Wisconsin-Madison, Madison, WI, USA*
- DAVED H. FREMONT • *Center for Structural Genomics of Infectious Diseases, St. Louis, MO, USA; Department of Pathology and Immunology, Washington University, St. Louis, MO, USA; Department of Biochemistry and Molecular Biophysics, Washington University, St. Louis, MO, USA*
- MAITE GARCIA • *Division of Cancer Genomics and Proteomics, Northeast Structural Genomics Consortium (NESG), Ontario Cancer Institute, Toronto, ON, Canada*
- ELIZABETH M. GLASS • *Center for Structural Genomics of Infectious Diseases, Mathematics and Computer Science Division, Argonne National Laboratory, Argonne, IL, USA*
- LUKASZ GOLDSCHMIDT • *UCLA-DOE Institute for Genomics and Proteomics, Howard Hughes Medical Institute, University of California, Los Angeles, CA, USA*
- MAREK GRABOWSKI • *Department of Molecular Physiology and Biological Physics, University of Virginia, Charlottesville, VA, USA; Center for Structural Genomics of Infectious Diseases (CSGID), University of Virginia, Charlottesville, VA, USA; Midwest Center for Structural Genomics (MCSG), University of Virginia, Charlottesville, VA, USA; New York Structural Genomics Research Consortium (NYSGRC), University of Virginia, Charlottesville, VA, USA; Enzyme Function Initiative (EFI), University of Virginia, Charlottesville, VA, USA*
- CATHERINE HATZOS-SKINTGES • *Biosciences Division, Midwest Center for Structural Genomics, Argonne National Laboratory, Argonne, IL, USA*
- WEI JIANG • *Argonne Leadership Computing Facility, Argonne National Laboratory, Argonne, IL, USA*
- GRAZYNA JOACHIMIAK • *Biosciences Division, Midwest Center for Structural Genomics, Argonne National Laboratory, Argonne, IL, USA*

- ANDRZEJ JOACHIMIAK • *Center for Structural Genomics of Infectious Diseases, Computational Institute, University of Chicago, Chicago, IL, USA; Biosciences Division, Midwest Center for Structural Genomics, Argonne National Laboratory, Argonne, IL, USA; Biochemistry and Molecular Biology, University of Chicago, Chicago, IL, USA*
- YOUNGCHANG KIM • *Center for Structural Genomics of Infectious Diseases, University of Chicago, Chicago, IL, USA; Midwest Center for Structural Genomics, Argonne National Laboratory, Argonne, IL, USA; Structural Biology Center, Argonne National Laboratory, Argonne, IL, USA*
- KEEHWAN KWON • *Center for Structural Genomics of Infectious Diseases, J. Craig Venter Institute, Rockville, MD, USA*
- HYUN LEE • *Center for Pharmaceutical Biotechnology, University of Illinois at Chicago, Chicago, IL, USA*
- CHUNG A. LEE • *Department of Pathology and Immunology, Washington University, St. Louis, MO, USA*
- HUI LI • *Midwest Center for Structural Genomics, Argonne National Laboratory, Argonne, IL, USA*
- SAMUEL H. LIGHT • *Center for Structural Genomics of Infectious Diseases, Northwestern University Feinberg School of Medicine, Chicago, IL, USA*
- CHI-HAO LUAN • *Center for Structural Genomics of Infectious Diseases, High Throughput Analysis Laboratory, Department of Molecular Biosciences, Northwestern University, Evanston, IL, USA*
- JAMEY MACK • *Biosciences Division, Midwest Center for Structural Genomics, Argonne National Laboratory, Argonne, IL, USA*
- ELIZABETH M. MACLEAN • *Department of Molecular Physiology and Biological Physics, University of Virginia, Charlottesville, VA, USA; Center for Structural Genomics of Infectious Diseases (CSGID), University of Virginia, Charlottesville, VA, USA; Midwest Center for Structural Genomics (MCSG), University of Virginia, Charlottesville, VA, USA; New York Structural Genomics Research Consortium (NYSGRC), University of Virginia, Charlottesville, VA, USA; Enzyme Function Initiative (EFI), University of Virginia, Charlottesville, VA, USA*
- SHIN-ICHI MAKINO • *Department of Biochemistry, University of Wisconsin-Madison, Madison, WI, USA*
- MAGDALENA MAKOWSKA-GRZYSKA • *Center for Structural Genomics of Infectious Diseases, University of Chicago, Chicago, IL, USA*
- NATALIA MALTSEVA • *Center for Structural Genomics of Infectious Diseases, University of Chicago, Chicago, IL, USA*
- JOHN L. MARKLEY • *Department of Biochemistry, University of Wisconsin-Madison, Madison, WI, USA*
- WILLIAM H. MCCOY • *Department of Pathology and Immunology, Washington University, St. Louis, MO, USA*
- SHAHILA MEHBOOB • *Center for Pharmaceutical Biotechnology, University of Illinois at Chicago, Chicago, IL, USA*
- KAROLINA MICHALSKA • *Biosciences Division, Midwest Center for Structural Genomics, Argonne National Laboratory, Argonne, IL, USA; Biosciences Division, Structural Biology Center, Argonne National Laboratory, Argonne, IL, USA*

- GEORGE MINASOV • *Center for Structural Genomics of Infectious Diseases, Northwestern University Feinberg School of Medicine, Chicago, IL, USA; Midwest Center for Structural Genomics, Northwestern University Feinberg School of Medicine, Chicago, IL, USA*
- WLADEK MINOR • *Department of Molecular Physiology and Biological Physics, University of Virginia, Charlottesville, VA, USA; Center for Structural Genomics of Infectious Diseases (CSGID), University of Virginia, Charlottesville, VA, USA; Midwest Center for Structural Genomics (MCSG), University of Virginia, Charlottesville, VA, USA; New York Structural Genomics Research Consortium (NYSGR), University of Virginia, Charlottesville, VA, USA; Enzyme Function Initiative (EFI), University of Virginia, Charlottesville, VA, USA*
- RORY MULLIGAN • *Center for Structural Genomics of Infectious Diseases, Computational Institute, University of Chicago, Chicago, IL, USA*
- PETER J. MYLER • *Seattle Structural Genomics Center for Infectious Disease, Seattle, WA, USA; Seattle Biomedical Research Institute, Seattle, WA, USA; Department of Global Health, University of Washington, Seattle, WA, USA; Department of Medical Education and Biomedical Informatics, University of Washington, Seattle, WA, USA*
- CHRISTOPHER A. NELSON • *Center for Structural Genomics of Infectious Diseases, Washington University, St. Louis, MO, USA; Department of Pathology and Immunology, Washington University, St. Louis, MO, USA*
- BOGUSLAW NOCEK • *Biosciences Division, Midwest Center for Structural Genomics, Argonne National Laboratory, Argonne, IL, USA; Biosciences Division, Structural Biology Center, Argonne National Laboratory, Argonne, IL, USA; Center for Structural Genomics of Infectious Diseases, Computational Institute, University of Chicago, Chicago, IL, USA*
- CHRISTINE A. ORENGO • *Center for Structural Genomics of Infectious Diseases, Department of Structural and Molecular Biology, University College London, London, UK*
- JERZY OSIPIUK • *Biosciences Division, Midwest Center for Structural Genomics, Argonne National Laboratory, Argonne, IL, USA; Biosciences Division, Structural Biology Center, Argonne National Laboratory, Argonne, IL, USA; Center for Structural Genomics of Infectious Diseases, Computational Institute, University of Chicago, Chicago, IL, USA*
- KHAN TANJID OSMAN • *Structural Genomics Consortium, University of Toronto, Toronto, ON, Canada*
- SCOTT N. PETERSON • *Inflammatory and Infectious Disease Center, Sanford-Burnham Medical Research Institute, La Jolla, CA, USA*
- ISABELLE Q.H. PHAN • *Seattle Structural Genomics Center for Infectious Disease, Seattle, WA, USA; Seattle Biomedical Research Institute, Seattle, WA, USA*
- KANAGALAGHATTA RAJASHANKAR • *NE-CAT and Department of Chemistry and Chemical Biology, Cornell University, Argonne National Laboratory, Argonne, IL, USA*
- BENJAMIN E. RAMIREZ • *Center for Structural Biology, University of Illinois at Chicago, Chicago, IL, USA*
- KIIRA RATIA • *Department of Medicinal Chemistry and Pharmacognosy, University of Illinois at Chicago, Chicago, IL, USA; Research Resources Center, University of Illinois at Chicago, Chicago, IL, USA*
- BENOIT ROUX • *Department of Biochemistry and Molecular Biology, University of Chicago, Chicago, IL, USA*
- ALEXEI SAVCHENKO • *Center for Structural Genomics of Infectious Diseases, Department of Chemical Engineering and Applied Chemistry, University of Toronto, Toronto, ON, Canada*

- ANTHONY SEMESI • *Division of Cancer Genomics and Proteomics, Northeast Structural Genomics Consortium (NESG), Ontario Cancer Institute, Toronto, ON, Canada*
- LUDMILLA SHUVALOVA • *Center for Structural Genomics of Infectious Diseases, Northwestern University Feinberg School of Medicine, Chicago, IL, USA; Midwest Center for Structural Genomics, Northwestern University Feinberg School of Medicine, Chicago, IL, USA*
- TATIANA SKARINA • *Department of Chemical Engineering and Applied Chemistry, University of Toronto, Toronto, ON, Canada; CSGID, Toronto, ON, Canada*
- ROBIN STACY • *Seattle Structural Genomics Center for Infectious Disease, Seattle, WA, USA; Seattle Biomedical Research Institute, Seattle, WA, USA*
- KEMIN TAN • *Biosciences Division, Midwest Center for Structural Genomics, Argonne National Laboratory, Argonne, IL, USA; Biosciences Division, Structural Biology Center, Argonne National Laboratory, Argonne, IL, USA; Center for Structural Genomics of Infectious Diseases, Computational Institute, University of Chicago, Chicago, IL, USA*
- CHRISTINE TESAR • *Biosciences Division, Midwest Center for Structural Genomics, Argonne National Laboratory, Argonne, IL, USA*
- YUFENG TONG • *Structural Genomics Consortium, University of Toronto, Toronto, ON, Canada*
- AMY WERNIMONT • *Structural Genomics Consortium, University of Toronto, Toronto, ON, Canada*
- XIAOHUI XU • *Department of Chemical Engineering and Applied Chemistry, University of Toronto, Toronto, ON, Canada; Midwest Center for Structural Genomics, University of Toronto, Toronto, ON, Canada*
- CORIN YEATS • *Center for Structural Genomics of Infectious Diseases, Department of Structural and Molecular Biology, University College London, London, UK*
- ADELINDA A. YEE • *Division of Cancer Genomics and Proteomics, Northeast Structural Genomics Consortium (NESG), Ontario Cancer Institute, Toronto, ON, Canada*
- MIN ZHOU • *Center for Structural Genomics of Infectious Diseases, University of Chicago, Chicago, IL, USA; Midwest Center for Structural Genomics, Argonne National Laboratory, Argonne, IL, USA*
- MATTHEW D. ZIMMERMAN • *Department of Molecular Physiology and Biological Physics, University of Virginia, Charlottesville, VA, USA; Center for Structural Genomics of Infectious Diseases (CSGID), University of Virginia, Charlottesville, VA, USA; Midwest Center for Structural Genomics (MCSG), University of Virginia, Charlottesville, VA, USA; New York Structural Genomics Research Consortium (NYSGRC), University of Virginia, Charlottesville, VA, USA; Enzyme Function Initiative (EFI), University of Virginia, Charlottesville, VA, USA*

Chapter 1

Data Management in the Modern Structural Biology and Biomedical Research Environment

Matthew D. Zimmerman, Marek Grabowski, Marcin J. Domagalski,
Elizabeth M. MacLean, Maksymilian Chruszcz, and Wlodek Minor

Abstract

Modern high-throughput structural biology laboratories produce vast amounts of raw experimental data. The traditional method of data reduction is very simple—results are summarized in peer-reviewed publications, which are hopefully published in high-impact journals. By their nature, publications include only the most important results derived from experiments that may have been performed over the course of many years. The main content of the published paper is a concise compilation of these data, an interpretation of the experimental results, and a comparison of these results with those obtained by other scientists.

Due to an avalanche of structural biology manuscripts submitted to scientific journals, in many recent cases descriptions of experimental methodology (and sometimes even experimental results) are pushed to supplementary materials that are only published online and sometimes may not be reviewed as thoroughly as the main body of a manuscript. Trouble may arise when experimental results are contradicting the results obtained by other scientists, which requires (in the best case) the reexamination of the original raw data or independent repetition of the experiment according to the published description of the experiment. There are reports that a significant fraction of experiments obtained in academic laboratories cannot be repeated in an industrial environment (Begley CG & Ellis LM, *Nature* 483(7391):531–3, 2012). This is not an indication of scientific fraud but rather reflects the inadequate description of experiments performed on different equipment and on biological samples that were produced with disparate methods. For that reason the goal of a modern data management system is not only the simple replacement of the laboratory notebook by an electronic one but also the creation of a sophisticated, internally consistent, scalable data management system that will combine data obtained by a variety of experiments performed by various individuals on diverse equipment. All data should be stored in a core database that can be used by custom applications to prepare internal reports, statistics, and perform other functions that are specific to the research that is pursued in a particular laboratory.

This chapter presents a general overview of the methods of data management and analysis used by structural genomics (SG) programs. In addition to a review of the existing literature on the subject, also presented is experience in the development of two SG data management systems, UniTrack and LabDB. The description is targeted to a general audience, as some technical details have been (or will be) published elsewhere. The focus is on “data management,” meaning the process of gathering, organizing, and storing data, but also briefly discussed is “data mining,” the process of analysis ideally leading to an understanding of the data. In other words, data mining is the conversion of data into information. Clearly, effective

Matthew D. Zimmerman and Marek Grabowski have contributed equally to this work.

data management is a precondition for any useful data mining. If done properly, gathering details on millions of experiments on thousands of proteins and making them publicly available for analysis—even after the projects themselves have ended—may turn out to be one of the most important benefits of SG programs.

Key words Databases, Data management, Structural biology, LIMS, PSI, CSGID

1 Introduction

1.1 *Data in Structural Biology*

Both structural genomics consortia and individual structural biology laboratories produce tremendous amounts of data, and having accurate, complete and consistent data is critical for reproducibility of biomedical research [1]. A single trip to a synchrotron for data collection by a productive crystallographic lab can generate hundreds of datasets totaling around 2 TB of raw data [2]. Modern data processing software can reduce, on the fly, a raw set of diffraction images into a single file that contains a description of every diffraction peak: Miller indices, intensity, and experimental uncertainty (sigma). These data are further reduced into one relatively small file that contains scaled and merged diffraction intensities. However, each file has to be associated with a particular sample (protein crystal) and the description of the experiment, which is usually written in the header of the diffraction image. These data are further used for structure determination and/or for function–structure relation studies.

To perform these studies the experimenter needs information about the protein (at a minimum, the protein sequence), crystallization conditions, and, for functional studies, protein production details. If this information is available, the process described above is simple to implement. Data harvesting from structure determination is relatively straightforward. The whole process following the placement of a crystal in the X-ray beam can be entirely controlled and captured by computer.

However, while this is very simple in theory, this simplicity has not yet been translated into practice. Analysis of the Protein Data Bank (PDB) [3, 4] shows that the number of data collection parameters marked as “NULL” in the header information (i.e., the detailed description of the experiment) is still significant [5, 6]. Moreover, data in the header are sometimes self-contradictory, contradictory to the experimental description in the paper citing the structure, or both [7, 8]. In that case, contacting the authors of the deposit and paper may be the only way to resolve the arising problems. Taking into account that only a small fraction, about 13 % [9], of structures determined by high-throughput consortia are converted (reduced) to peer-reviewed papers, the correctness of data uploaded to various databases like TargetTrack [10], TargetDB [11], and data banks like PDB is absolutely critical (see below).

1.2 Large-Scale Initiatives Create New Databases: TargetDB/ PepcDB/TargetTrack

Since their inception, many structural genomics efforts have adopted policies that experimental data produced by member consortia should be made available to the community from the moment of target selection. This has been particularly true for the two large initiatives from the National Institutes of Health (NIH): the Protein Structure Initiative (PSI) established in 2000 by the National Institute of General Medical Sciences (NIGMS) and the SG centers focusing on infectious diseases established in 2007 by the National Institute of Allergy and Infectious Diseases (NIAID). Even some partially privately funded SG efforts like the Structural Genomics Consortium (SGC) have established policies to release some experimental data to the general public [12] (typically only after the structure is determined and deposited). In the specific case of the centers funded by NIGMS and NIAID, the NIH established the target registration database, TargetDB [11], and required that all member consortia deposit data on the progress of their targets. Subsequently many other SG centers worldwide have deposited some of their experimental data as well.

Initially, the main purpose of TargetDB was the prevention of duplication of effort between different SG centers and maximization of the structural coverage of the protein fold space. The scope of the data was very modest. It included protein identification information (sequence, organism) and the timeline of changes in experimental status for each target. Status events included target selection, cloning, expression, purification, as well as crystallization, diffraction, determination of crystal structure, and PDB deposition (for targets studied by X-ray crystallography) or obtaining the HSQC spectra, determination of NMR structure, and BMRB/PDB deposition (for targets studied by NMR).

However, even the modest amount of data available in TargetDB permitted interesting analyses of the overall SG structure determination pipeline [13, 14]. In particular, the overall efficiency of the pipeline—the ratio of solved structures to clones—was found to be below 10 % even in the most productive centers. The two steps that contributed most to the failure of a target in the pipeline were production of soluble protein and diffraction-quality crystals. Not surprisingly, the success ratio depended very strongly on the type of protein as well as the methodology used by particular centers. There was not a single overall bottleneck factor. In 2004, TargetDB was extended to the Protein Expression, Purification, and Crystallization Database (PepcDB) [15] which in addition to simple status history included multiple trials, tracking of failed as well as successful experiments, and more detailed descriptions of protocols.

In 2010, PepcDB and TargetDB were merged into a single new database, TargetTrack, part of the new PSI-Structural Biology Knowledge Base (PSI-SBKB) [10, 16]. The new repository

extended the definition of a target to include protein–protein complexes and incorporated tracking of biological assays needed in the PSI:Biography phase. As of January 2013, TargetTrack contained data on over 300,000 targets and over 1,000 protocols.

1.3 Diverse Approaches to Data Management in SG Centers

Development of effective data management systems was a necessity for the large-scale SG centers, not only in order to provide the data to the scientific community but also particularly to effectively handle the huge amounts of experimental data, plan experiments, adjust experimental approaches (e.g., choice of cloning vectors, sequence truncation, crystallization conditions, structure determination procedures), and prioritize targets. These needs required gathering far more data than what was being required by TargetTrack.

In general, two levels of data management are needed in high-throughput, high-output structural biology programs: the *target tracking* level and the *experiment tracking* level. The target tracking level comprises target selection, overall experimental status of each target, center-wide efficiency statistics, and generation of reports to the public and to other databases such as TargetTrack. Almost all SG centers have a separate target-tracking database, though some functionality (e.g., target selection) can be “offloaded” to other specialized databases. The primary audience for the target-tracking level is everyone interested in a “high-level” view of the data produced by the center: the center’s scientists and administrators as well as members of the scientific community with interest in the targeted proteins. This level is typically not designed for uploading new data or providing all details of individual experiments; these tasks are better handled at the experimental tracking level.

The experimental tracking level comprises the tools used to collect the results of experiments performed in the laboratory. This type of tool is generally known as a “laboratory information management system” or LIMS. LIMSs are typically used day to day by the researchers conducting the experimental work of a laboratory and may be highly customized to the protocols and work flow of a particular laboratory. LIMSs may also provide tools to help design experiments, operate laboratory equipment, semiautomatically harvest data, track the use of resources, etc. As a result, the primary audience for the LIMS is composed of those interested in a “low-level” view of the data, the center researchers themselves. As compared to the target-tracking level, it is not uncommon to use more than one LIMS in a single SG center, as different systems may be used in different laboratories.

It should be noted that splitting the data management system of a typical SG center into two distinct levels, “high-level” target tracking and “low-level” experiment tracking, is somewhat arbitrary. Some data are natural candidates to be kept at the LIMS

level only, for example, the location in the freezer where a particular clone is stored or the particular lot of a reagent or a crystallization buffer. Conversely, some data may only apply at the target-tracking level, for example, the number of publications referencing a given protein. In principle, it is possible for a single database and/or data management system to fully implement both levels. However, in practice, it seems that solutions where the two levels are implemented as separate systems/databases appear to be more common, especially for the larger scale projects.

There have been several “top-down” attempts to design a general framework for SG data management systems in the form of data dictionaries [17] or a protein production UML data model [18]. The latter has been implemented by several systems, such as HalX [19] or the Protein Information Management System (PiMS) [20] used by a number of European SG labs. However, most of the SG centers set up data management systems in a more ad hoc, “bottom-up” manner. Initially, some centers attempted to use commercial LIMS, but often these solutions were not flexible enough or even robust enough, and most SG centers developed their own solutions “in-house.” There are exceptions to this rule. For example, the Structural Genomics Consortium uses two commercially available software systems: the Beehive LIMS (Molsoft LLC; <http://www.molsoft.com/beehive.html>) and Electronic Laboratory Notebook (now iLabber; Contur Software; <http://www.contur.com/home/>). It should be noted however that unlike many SG consortia, SGC does not deposit the results of its experiments to PepcDB or TargetTrack. Several of the SG-developed data management systems have been described in the literature [21–23], but to our knowledge, none of these systems have been fully commercialized.

One comprehensive data SG management system that has gained wider use is Sesame, developed by Zsolt Zolnai at Center for Eukaryotic Structural Genomics (CESG) [22]. It has been adopted by a number of labs and specialized centers.

The data management system for the Joint Center for Structural Genomics (JCSG) was developed by the center’s programming team in parallel with the construction of the physical pipeline. The LIMS part of the system functions as a hub of information, recording all pipeline steps from target selection to deposition. The tracking database uses Oracle as its engine and tracks 424 experimental parameters, organized into 130 tables [24]. The tools and interfaces to the database contain approximately 360,000 lines of code, which illustrates the level of complexity of this and similar systems.

The Northeast Structural Genomics (NESG) consortium’s data management system is organized as a “federated database framework,” comprising a set of distributed, interconnecting databases [21]. The main target-tracking database, SPINE, serves

as an analysis system, utilizing data mining and machine learning tools. In particular, decision trees are used for predicting chances for protein solubility, successful purification, and crystallization. These predictions are used in directing targets to X-ray crystallography or NMR studies [14].

The other two large-scale PSI:BiologY centers—the Midwest Center for Structural Genomics (MCSG) and the New York Structural Genomics Research Consortium (NYSGRC)—use the data management system developed in the Minor Lab at the University of Virginia. In both cases, the system is based on a collection of customized LIMS in each site laboratory and a central database (UniTrack, described below) that curates and unifies data obtained by various laboratories. In the case of MCSG, several different LIMSs are used in different laboratories, including LabDB, Mnemosyne, and ANL-DB. In NYSGRC, two different instances of LabDB are used. Similar systems are also deployed in the Center for Structural Genomics of Infectious Diseases (CSGID) and the Enzyme Function Initiative (EFI).

2 A Centralized Target Management System: UniTrack

The central, public system comprising the target-tracking level of the SG management system developed by the Minor Lab at the University of Virginia is named UniTrack. As mentioned above, the MCSG, NYSGRC, CSGID, and EFI consortia are all driven by variants of the UniTrack system. The system comprises a core abstraction based on 10 years of experience in SG data management, with a common database architecture and set of tools for managing target and experimental data. Each site is based on the UniTrack core but is then highly customized for the needs of the particular center or consortium of research laboratories. In each case, the UniTrack-derived system comprises the central tracking database and a set of auxiliary databases and applications, which collect and integrate experimental data and are provided by distributed LIMSs deployed in participating laboratories (Fig. 1). Experimental data from different LIMSs are combined and incorporated into UniTrack via a standard protocol. In the most basic case, each LIMS generates XML files in a predefined format, which are parsed by UniTrack tools. An alternative (and more efficient) method, where a LIMS directly communicates with the tracking database, has also been developed. The LIMSs can be very diverse; however, they all must be able to provide the minimum set of required data for cloning, expression, purification, and crystallization experiments.

The experimental pipeline starts with target selection and validation, which is specific for a particular center. The validation process is performed automatically and typically involves checking

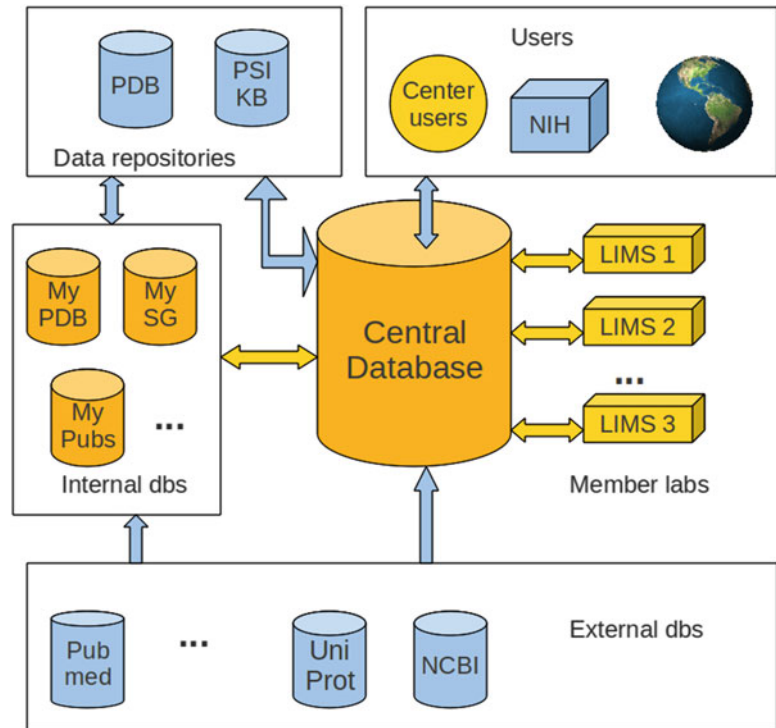


Fig. 1 The architecture of the UniTrack data management system. The central database interacts with LIMSs distributed in member labs. A number of auxiliary databases are used to store data from the PDB, data from other SG centers, and SG publications. The central database is responsible for producing reports for external data repositories such as PSI-SBKB. UniTrack databases are synchronized with external data sources such as NCBI GenBank, UniProt, and PubMed via custom scripts. Users interact with the system via a web interface

the accuracy of the amino acid and the nucleotide sequences as well as checking if the selected protein is homologous to proteins with structures in the PDB or to targets selected by other SG centers. Validated targets are inserted into the tracking database. Protein annotations and related data are automatically imported from external databases such as NCBI GenBank [25], Uniprot [26], PDB, and the PSI-SBKB. Depending on the needs of a particular center, between 30 and 80 attributes of any given protein target are stored in UniTrack.

UniTrack keeps a history and the results of the experiments for each target (Fig. 2). About 400 distinct data attributes are used to describe an experimental trial, from the cloning of a target through the determination of its structure. Almost all protein production and crystallization data can be automatically imported from the local LIMS or equipment database. However, smaller labs that do not have a LIMS deployed can still contribute data to UniTrack by entering it manually using the customized interface. Diffraction

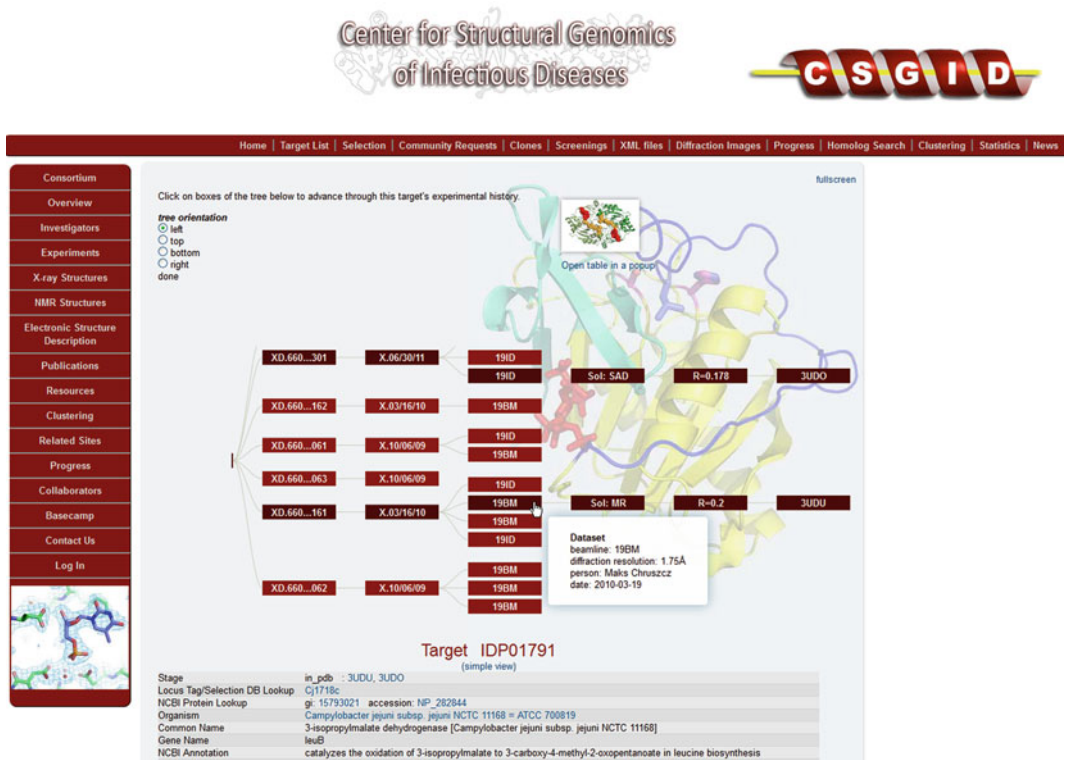


Fig. 2 Fragment of an experiment tree displayed in the UniTrack-based CSGID interface. *Boxes* represent particular experiments: purification (P), crystallization drop (XD), crystal harvest (X), data collection (beamline name), structure solution (Sol), refinement (R factor), and PDB deposit (PDB id). *Paths* in the tree represent trials for a particular sample. The *white box* that appears when the cursor hovers over an item displays additional details about a particular step. In addition, clicking on any of the boxes display all the data known about this step stored in the database

and structure determination data is currently imported automatically only from the LabDB instances that have the *bkldb* module enabled [27]. Researchers in other labs upload scaling logs and refinement files manually via the interface.

The tracking database also generates real-time internal reports and statistics as well as the XML files that are being submitted to the TargetTrack repository. In addition, the periodic reports required by various bodies are generated in real time from the database and accessible to the general public. In some sense, all of the portions of UniTrack that generate publicly accessible web pages serve as reports.

The customized instances of UniTrack for each center drive dynamic parts of the centers' corresponding web portals. The web interfaces are implemented using the Model-View-Controller (MVC) architecture, with separate layers for data retrieval (model), "business logic" (controller), and web page rendering (view).

Even with the use of the CakePHP MVC framework (<http://cakephp.org>) the customized web interfaces for the centers are quite complicated; as an example, the implementation of the CSGID web interface contains over 50,000 lines of source code.

2.1 The LabDB “Super-LIMS”

LabDB is a modular “super-LIMS,” originally developed to track the structure determination pipeline from cloning to structure determination (Fig. 3). The central component of the system is a PostgreSQL database server coupled with a web-based framework, along with two specialized tools: *Xtaldb*, for designing and tracking crystallization experiments, and *hkldb*, a module of the HKL-2000/3000 system [27] for incorporating information from crystallographic data collection and structure determination. *hkldb* and *Xtaldb* can also be used with stand-alone databases.

One of the fundamental design goals of LabDB is to harvest data automatically or semiautomatically from laboratory equipment whenever possible. To that end, the system has modules to import data from a variety of different types of laboratory equipment, including chromatography systems (GE Healthcare AKTA systems), electrophoresis documentation and separation systems

LabDB Macromolecular Crystallography Laboratory Information Management System

[Administrator] Logged in as Matthew Zimmerman (switch user) (logout)

Home File Add Search Batch Stats Admin Options Help

View Project 012178

Edit this project | Delete this project | Add a clone to this project
View thermal shifts

Experimental history	Details
<p>fold unfold</p> <p>012178</p> <ul style="list-style-type: none"> clone: 012178 - pSGC (plate Sino128SGC well B09) <ul style="list-style-type: none"> expression: 0.75 mL PASM-5052 BL21 DE3 RIL (2011-09-08) (pla expression: 0.75 mL PASM-5052 BL21 DE3 RIL (2011-09-20) (pla clone: 12178 - pSGC-His (plate SINO.0149-SGC-PROD well D10) <ul style="list-style-type: none"> expression: 2000 mL PASM-SEMET BL21 (DE3) RIL (2011-10-05) <ul style="list-style-type: none"> purification: Ni SEC (2011-10-11) <ul style="list-style-type: none"> macroprep: 15.71 mg/mL 012178 NpSGC L-1 (2011-10-14) <ul style="list-style-type: none"> plate: 012178 #1 <ul style="list-style-type: none"> crystal: XA003096 plate: 012178 #2 plate: 012178 #3 plate: 012178 #4 clone: C11 - CHS30 (plate SIN0113CHS30 well C11) clone: 012457 - CHS23 (plate Sino Extras CHS23 well C10) <ul style="list-style-type: none"> expression: 2000 mL PSAM-SeMet BL21 DE3 RIL (2011-08-11) (f) <ul style="list-style-type: none"> purification: Ni SEC (2011-08-16) <ul style="list-style-type: none"> macroprep: 18.8 mg/mL 012457 CHS23 L-1 (2011-08-18) <ul style="list-style-type: none"> plate: 012457 #3 plate: 012457 #2 plate: 012457 #1 plate: MSCG 1-3 with NADP #2 plate: MSCG 1-3 with NADP #1 plate: MSCG 1-3 with NADP #3 plate: 012457 #4 <ul style="list-style-type: none"> crystal: XA003134 plate: Opt MCG C5 add #1 expression: 0.75 mL PASM-5052 BL21 DE3 RIL (2011-08-19) (pla 	<p>Description Short-chain dehydrogenase/reductase [Sinorhizobium melloti 1021] (includes experiments from former project 012457)</p> <p>Responsible person Steve Almo</p> <p>Project group PSI/BIO/SINO</p> <p>Organism Sinorhizobium melloti 1021</p> <p>Locus tag SMc04391</p> <p>Uniprot identifier</p> <p>Genbank identifier CAC47875.1</p> <p>Genbank GI 15076322</p> <p>Bridging Project Descriptor</p> <p>Sequence and statistics</p> <p>Sequence <pre> MTKARPVAI VTGPRGIGL GIARALAAG FDIATIGIGD AECVAPVIAE LSQLGARVIF LRKQLADLSS HQATYDAVVA EFGRIIDLIN NAGIASIVRD DFLDLKPENF DTIVGVNLRG TVFFTOAVLK AMLASDARAS RSIINITSVS AVHTSPERLD YCHSKAGLAA P50QLALRLA ETGIAVFEVR PGIIRSDMTA AVSKRYDGLI EESGLVPHRRR GEPEIDIGNIV AGLAGOGGPF ATGVSQIQDG GLSEIGRL </pre> </p> <p>Molecular weight 26418.4 Da</p> <p>Isoelectric point 5.1</p>

Fig. 3 A typical target overview page in the LabDB LIMS