

BestMasters

Thorsten Will

Predicting Transcription Factor Complexes

A Novel Approach to
Data Integration in Systems Biology



Springer Spektrum

BestMasters

Springer awards „BestMasters“ to the best master’s theses which have been completed at renowned universities in Germany, Austria, and Switzerland.

The studies received highest marks and were recommended for publication by supervisors. They address current issues from various fields of research in natural sciences, psychology, technology, and economics.

The series addresses practitioners as well as scientists and, in particular, offers guidance for early stage researchers.

Thorsten Will

Predicting Transcription Factor Complexes

A Novel Approach to Data
Integration in Systems Biology



Springer Spektrum

Thorsten Will
Saarbrücken, Germany

BestMasters

ISBN 978-3-658-08268-0

ISBN 978-3-658-08269-7 (eBook)

DOI 10.1007/978-3-658-08269-7

Library of Congress Control Number: 2014956553

Springer Spektrum

© Springer Fachmedien Wiesbaden 2015

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.

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

Springer Spektrum is a brand of Springer Fachmedien Wiesbaden
Springer Fachmedien Wiesbaden is part of Springer Science+Business Media
(www.springer.com)

Geleitwort des Betreuers

Die genetische Information jeder biologischen Zelle ist bekanntlich in deren Erbsubstanz kodiert. Insofern gibt es beispielsweise keinerlei Unterschiede zwischen den etwa 200 verschiedenen Zelltypen eines Menschen (Nervenzellen, Muskelzellen, Knochen, Haut etc.). Obwohl diese Zellen natürlich große morphologische Unterschiede aufweisen, enthalten alle Zellen im Prinzip dieselbe Information. Entscheidend für den jeweiligen Zustand einer menschlichen Zelle - d.h. deren Differenzierung in die einzelnen Gewebetypen - ist deshalb nicht, welche Informationen die Erbsubstanz prinzipiell enthält, sondern welche der etwa 22.000 Gene in ihr tatsächlich „abgelesen“ werden und welche nicht. Man kann dies vereinfacht mit dem Lesen eines Buches vergleichen, bei dem auch jeweils eine Seite aufgeschlagen wird. Wer daher in der Zelle diesen Lese-Prozess reguliert, bestimmt quasi das Schicksal der Zelle. Die Entschlüsselung dieser Prinzipien zur Regulierung der Gentranskription ist somit ein sehr wichtiger Schlüssel zum Verständnis von Zellen.

Die wichtigste Rolle bei der Entscheidung, welche Gene abgelesen werden sollen, übernehmen Eiweißmoleküle, sogenannte Transkriptionsfaktoren, die gezielt an Bindungsstellen an der Erbsubstanz binden können. Wenn diese Transkriptionsfaktoren direkt „vor“ dem Beginn einer Gensequenz an die Erbsubstanz binden, kann der Kopierapparat der Zelle, die RNA-Polymerase, gezielt an diese Region rekrutiert werden und beginnt, die sich anschließende Region der Erbsubstanz zu kopieren, d.h. abzulesen.

Man könnte sich nun vorstellen, dass zu jedem einzelnen Gen ein bestimmter Transkriptionsfaktor gehört, der den Ableseprozess dieses Gens reguliert. Jedoch müsste ja auch dieser Transkriptionsfaktor in der Zelle durch Ablesen eines anderen Gens produziert werden, so dass man dann zweimal so viele Gene bräuchte. Und so weiter. Man sieht leicht, dass solch eine Variante nicht funktionieren kann. Eine andere Möglichkeit wäre, dass ein einzelner Transkriptionsfaktor jeweils das Ablesen von vielen anderen Genen reguliert. Dann käme man mit einer wesentlich geringeren Anzahl an Transkriptionsfaktoren aus, könnte aber die Regulation der einzelnen Gene nicht mehr so feinkörnig steuern. Eine dritte Möglichkeit, die nun

tatsächlich in Zellen realisiert wird, ist dass jeweils mehrere Transkriptionsfaktoren gemeinsam den Ableseprozess von einzelnen Genen kontrollieren. Damit reichen wenige hunderte an Transkriptionsfaktoren aus um eine enorme Anzahl an kombinatorischen Varianten zu erzeugen.

Der Autor entwickelte in seiner Masterarbeit im Fachgebiet Bioinformatik einen neuen Ansatz um Eiweißkomplexe zu identifizieren, die aus mehreren Transkriptionsfaktoren sowie aus weiteren Proteinen bestehen. Um die Praktikabilität der Methode zu testen, wurde als Modellorganismus die Bäckerhefe (*S. cerevisiae*) ausgewählt, da hierfür besonders gute experimentelle Daten zu paarweisen Proteininteraktionen vorliegen. Der Algorithmus verwendet die in der Informatik oft eingesetzte Baumstruktur zur Aufzählung aller möglichen Komplexe, die an die Erbsubstanz binden und maximal 10 Eiweißmoleküle enthalten. Ein neuartiger Beitrag bestand darin, die grundlegende Datenstruktur für die Interaktionen nicht auf gesamten Proteinen aufzusetzen, sondern auf deren Domänenbausteinen. So konnten zum einen weitere Interaktionsdaten zwischen Proteindomänen eingebunden werden. Zum anderen ergibt sich eine feinere strukturelle Auflösung der in Konkurrenz miteinander stehenden Kontakte. Mit dem neuen Ansatz konnten für Hefe mehr als 10 mal so viele unterschiedliche Proteinkomplexe generiert werden wie mit anderen derzeit verfügbaren Methoden. Der Autor zeigte zudem, dass die Ergebnisse eine bessere Abdeckung der bisher experimentell charakterisierten Komplexe liefern als alle anderen Methoden und dass die vorhergesagten Komplexe eine hohe biologische Plausibilität besitzen.

Diese vielversprechenden Ergebnisse lassen es denkbar erscheinen, ähnliche Methoden auch für komplexere Lebewesen wie die Maus oder sogar den Mensch einzusetzen. Dies wäre ein wichtiger Schritt dabei, die Mechanismen der Genregulation besser zu verstehen, da deren Fehlfunktionen natürlich auch zur Entstehung vieler Krankheiten beitragen. Die Bioinformatik übernimmt bei solchen Projekten meist die wichtige Aufgabe der Datenintegration und ermöglicht es, die Existenz bestimmter Szenarien oder Mechanismen zu postulieren, deren Korrektheit dann im Experiment gezielt getestet werden kann.

Prof. Dr. Volkhard Helms

Institutsprofil

Das Zentrum für Bioinformatik an der Universität des Saarlandes¹ ist eine interfakultäre Einrichtung zwischen der Informatik, Medizin und den Lebenswissenschaften Biologie, Chemie und Pharmazie. Ein wichtiger Schwerpunkt der Saarbrücker Bioinformatik-Forschung ist die medizinische Bioinformatik, d.h. die Anwendung von Bioinformatik-Methoden für die Bearbeitung von biomedizinischen Daten. Man interessiert sich heutzutage zum Beispiel dafür, wie sich die Erbsubstanz in Krebszellen von der Erbsubstanz in daneben liegenden gesunden Zellen unterscheidet. Ein anderes wichtiges Forschungsgebiet ist die Resistenzforschung, wie es Viren und Bakterien durch Veränderung ihres Erbgutes schaffen, die Wirkung von an sich hoch aktiven antiviralen oder antibakteriellen Wirkstoffen ins Leere laufen zu lassen.

Das Zentrum für Bioinformatik organisiert an der Universität des Saarlandes einen grundständigen Bachelorstudiengang Bioinformatik und einen darauf aufbauenden Masterstudiengang. Basierend auf den anerkannten Forschungserfolgen, aber auch befördert durch die Anbindung an die exzellenten Forschungsinstitute der Saarbrücker Informatik und durch die hohe gesellschaftliche Relevanz der bearbeiteten Forschungsthemen, genießt die Saarbrücker Bioinformatik weltweit eine hohe Reputation. Die Absolventen der beiden genannten Studiengänge sind in nationalen und internationalen Unternehmen sowie in Forschungsinstituten sehr stark nachgefragt.

¹www.zbi.uni-saarland.de

Preface

Gene regulatory networks are fixed determinants of cellular control and the abundance of differentially expressed regulatory proteins, called transcription factors, their driving signal. In concert with specific epigenetic marks, transcription factors define the active subset of the network to govern distinct cellular states in time and space.

Eukaryotic gene expression is generally controlled through molecular logic circuits combining regulatory signals of several transcription factors. Recently, it has been shown that complexation of regulatory proteins is a prevailing and highly conserved mechanism of signal integration within critical regulatory pathways, like body part formation or differentiation. A knowledge of potential assembly candidates could provide the basic information that is needed to infer possible target genes as well as the exerted mechanism of influence. There already exists a plethora of approaches to predict protein complexes from protein-protein interaction data. However, those are generally designed to detect large self-contained functional complexes and lack the ability to reveal dynamic and highly modular combinatorial complex assemblies, a property of crucial importance for the signal integration exerted by transcription factor complexes.

The method proposed in this thesis combines protein-protein interaction networks and domain-domain interaction networks with the well-known cluster-quality metric cohesiveness. A novel growth algorithm is described that locally maximizes the metric on the holistic level of protein interactions while sophisticated connectivity constraints are preserved. Assuming that each domain can only support one interaction, the domain topology can be utilized to account for the exclusive and thus combinatorial nature of physical interactions between proteins. During the growth process, the complex candidate is thought to be backed by a spanning tree of simultaneously possible domain interactions which restrict further expansion possibilities. Consequently, every addition of a protein requires the choice of an applicable domain interaction which again influences later steps. Often many options have to be taken into account by branching of the algorithm, which naturally allows for the justified prediction of a manifold of transcription factor complexes from a common start.

The proposed approach outperformed popular complex prediction methods by far for the prediction of transcription factor complexes in yeast. The evaluation was based on established benchmarks assessing accordance with several reference complex datasets as well as measures of biological relevance. Additionally, many of the predictions of the proposed method could be associated with target genes and a potential regulatory effect. Furthermore, predicted candidates could be mapped to distinct functions during a defined cellular state and condition by analyzing the expression coherence among their regulated genes for cell cycle expression data. Many findings were backed up by literature evidence.

The results encourage an application to higher eukaryotes where the combinatorial interplay between transcription factors is more pronounced. The knowledge of putative transcription factor complexes - DNA-binding members and recruited potentially regulatory active proteins - offers novel capabilities in the automatized modeling of gene regulatory networks which may assist to surpass nowadays models.

A condensed summary of the novel concept, the main method and the results for yeast was previously published in [1] prior to the production of this book.

Contents

1	Introduction	1
1.1	Motivation	1
1.2	Eukaryotic transcriptional regulation	4
1.2.1	Transcription factors	7
1.2.2	Cis-regulatory modules	13
1.2.3	Complexation is important in development	21
1.3	Outline and Goal	23
2	Related work	25
2.1	Protein complex prediction from networks	25
2.1.1	ClusterONE and cohesiveness	27
2.1.2	Quality measures for complex prediction	30
2.2	Protein complexes beyond plain networks	34
2.2.1	Domain-domain interaction model	37
2.3	Prediction of cooperative TFs	40
2.3.1	Expression coherence scoring	40
3	Materials and Methods	45
3.1	Domain-aware cohesiveness optimization	45
3.1.1	Why seeding from pairs is beneficial	51
3.2	Data sources, their retrieval and preprocessing	53
3.2.1	General yeast protein data	53
3.2.2	Weighted protein-protein interaction data	55
3.2.3	Domain-association and domain-interactions	56
3.2.4	Transcription factors and their binding sites	60
3.2.5	Expression data	64
3.2.6	Reference data	65
3.3	Workflow and implementation	66
3.3.1	Building the domain-domain interaction network . .	68
3.3.2	Domain-aware cohesiveness optimization	69

4 Results and Discussion	89
4.1 Impact of algorithm engineering on runtime	90
4.2 Common protein complex prediction benchmarks	91
4.2.1 Comparison to reference complexes	94
4.2.2 Assessment of biological relevance	97
4.2.3 Evaluation of postprocessing and thresholds	98
4.3 Analysis in the transcription factor context	100
4.3.1 Estimation of target genes	101
4.3.2 Estimating the modes of action	103
4.3.3 Significance in yeast cell cycle	106
5 Conclusion and Outlook	111
5.1 Conclusion	111
5.2 Outlook	111
A Additional tables	113
A.1 Supplement to complex prediction benchmarks	113
A.2 Supplement to target analysis	119
Bibliography	121

List of Tables

3.1	Databases integrated in InterPro	58
3.2	Mapping motifs to proteins	62
4.1	Predicted TF complexes from various approaches	93
4.2	Results to reference complex benchmarks	95
4.3	Average sizes of predictions	96
4.4	Results of biological relevance checks	98
4.5	Results of added merging procedures	99
4.6	Descriptive yeast GO terms	104
4.7	Significant increase in ECS	110
A.1	Parameter optimization: MCODE	114
A.2	Parameter optimization: MCL	115
A.3	Parameter optimization: ClusterONE	116
A.4	Parameter optimization: ClusterONE with pairs	117
A.5	Size dependency: ClusterONE	118
A.6	Target analysis predictions	119
A.7	Target analysis systematic/random	119
A.8	Target analysis ClusterONE	120

List of Figures

1.1	Combinatorial control in development	3
1.2	Eukaryotic gene architecture	5
1.3	Basal transcription	7
1.4	Proximal promotor activation	10
1.5	Distal promotor activation	10
1.6	TFs and chromatin state	12
1.7	Models of cis-regulatory regulation	14
1.8	Graded signal integration	16
1.9	Single TFs still matter in a cooperative world	19
1.10	Complex assembly is crucial	20
1.11	Collaborative competition	21
2.1	Cohesiveness definitions	28
2.2	Protein interactions are ambiguous	35
2.3	Domain decompositions	38
2.4	Expression coherence score definition	41
2.5	Expression coherence score example	42
3.1	Domain-aware reachability	47
3.2	Detailed algorithm example	49
3.3	Pairwise seeding example	52
3.4	Pairwise seeding scheme	54
3.5	Binding site merge example	64
3.6	Workflow	67
3.7	Domain multiplicity example	72
3.8	State descriptions	76
3.9	ECS Venn-diagram	87
4.1	Optimization and runtimes	90
4.2	Weighted and unweighted variant and depth-thresholds . .	101
4.3	Distribution of targets to TFs	102
4.4	MET4/MET32 refinement	107

List of Algorithms

2.1	Iterative cohesiveness optimization	29
3.1	Domain-aware cohesiveness optimization	48
3.2	BS-sweep	63
3.3	Domain-aware cohesiveness optimization IIa	71
3.4	Domain-aware cohesiveness optimization IIb	72
3.5	Growth manager	74
3.6	Domain-aware cohesiveness optimization III	78
3.7	Iterative merging procedure	79
3.8	Pairwise merge function	80
3.9	Pairwise binding site compatibility	81
3.10	Extend valid binding sites by an additional TF	83
3.11	Distance constraint adjacent binding sites for a set of TFs .	84
3.12	Compute significance of dECS with binding site constraints	86

List of Abbreviations

bp base pair(s)

CRM cis-regulatory module

DDI(N) domain-domain interaction (network)

EC(S) expression coherence (score)

GO Gene Ontology (annotation)

GTF general transcription factor

HMM hidden markov model

MMR maximum matching ratio

ORF open reading frame

PDB Protein Data Bank

Pol II RNA-Polymerase II

PPI(N) protein-protein interaction (network)

SPIN simultaneous protein(-protein) interaction network

TBP TATA-binding protein

TF transcription factor

TSS transcriptional start site

YPA Yeast Promotor Atlas