UNIVERSITY OF TARTU Faculty of Biology and Geography Institute of Molecular and Cell Biology

Hedi Peterson

The discovery and characterization of regulatory motifs in *Saccharomyces cerevisiae* Master's Thesis

Supervisor: Jaak Vilo, PhD

Tartu 2006

"Computers are useless. They can only give you answers."

Pablo Picasso

Contents

1	Intr	roduction 1		
	1.1	Backg	round	1
	1.2	Object	tive	3
2	Bac	kgrour	ıd	4
	2.1	Biolog	ical background	4
		2.1.1	Gene Ontology	5
		2.1.2	Regulatory motifs	8
		2.1.3	Protein-protein interactions	8
	2.2	Bioinf	ormatics approaches	10
		2.2.1	Pattern Discovery	10
		2.2.2	Phylogenetic footprinting	10
3	Ma	terial a	and methods	12
	3.1	Data		12
		3.1.1	Sequences	13
		3.1.2	Previously known transcription factors and regulatory	
			motifs	14
		3.1.3	GO groups	16
		3.1.4	Protein-protein interaction data	18

	3.2	Pattern discovery	19
	3.3	Statistical evaluation	19
	3.4	Expansion of groups by PPI data	21
4	Res	ults	24
	4.1	Randomization	24
	4.2	Pattern discovery for known transcription factors' target groups	27
	4.3	Pattern discovery in GO dataset	28
	4.4	Expansion of GO groups by PPI	32
	4.5	From protein-protein interactions to regulatory motifs and GO	
		annotations	32
		4.5.1 Pipeline example	33
	4.6	Web-tool Gviz-PPI	36
	4.7	Data update to BiGeR	36
5	Disc	cussion	41
6	Conclusions		
Su	ımma	ary	45
Su	ımma	ary in Estonian	47
A	cknov	vledgements	49
Re	References 5		
A	ppen	dix	56

List of Figures

2.1	Example of a GO molecular function subtree. Here GO:0003887	
	is a leave and $GO:0003674$ is a root node. The parent of	
	GO:0003887 is GO:0016779 and it holds genes from both of	
	its children (GO:0003887 and GO:0003964). In this example	
	only is_a relationships occur	6
3.1	Phylogenetic tree of yeast species derived from 25S rDNA se-	
	quences. The species used in the current thesis are underlined.	
	The picture originates from Dujon $et al 2004$	13
3.2	Histogram of known transcription factors per gene. \ldots .	16
3.3	The distribution of the size of the GO groups. The size of the	
	groups is highly variable	17
4.1	p-values from real and random data with the input size of 160	
	genes	25
4.2	P-values from random data for each GO group size \ldots .	26
4.3	The connected graph of input ORFs. The arrows indicate the	
	known interactions between genes	34
4.4	The GO annotations for the 5 ORFs with GOSt tool (Reimand	
	2006). The bottom line on the figure describes the most rele-	
	vant GO annotation for the input set	35

- 4.5 The clustering based on stress response (Gasch 2000) 38
- 4.6 The expansion of GO:0008652 using Gcn4 motif TGACTC.The rectangles are ORFs belonging to the GO group, the circles are interacting proteins. The red fulfil indicates that the TGACTC pattern is in the upstream region of the ORF . . . 39
- 4.7 BiGeR output of all patterns matching known Gcn4 binding site TGACTC. The top motifs come from various *in vivo* and *in vitro* analyzes. The last two rows describe motifs from ChIP-on-chip and Gene Ontology specific data 40

List of Tables

3.1	Size and number of ORFs of yeasts	14
3.2	The distribution of the number of orthologs for $Saccharomyces$	
	cerevisiae ORFs	15
3.3	Protein-protein interaction datasets overview	18
4.1	Patterns with the smallest p-values from GO groups	30
1	TF target genes and their relative motifs	58

Abbreviations

AD	Activation Domain
bp	base pair
CC	Cellular Component
ChIP-on-chip	Chromatin ImmunoPrecipitation on chip
DBD	DNA Binding Domain
DNA	DeoxyriboNucelic Acid
EP	Expression Profiler
GO	Gene Ontology
MF	Molecular Function
ORF	Open Reading Frame
PPI	Protein-protein interaction
RNA	RiboNucleic Acid
SGD	Saccharomyces Genome Database
TAP	Tandem Affinity Purification
TF	Transcription Factor
TFBS	Transcription Factor Binding Site

Chapter 1

Introduction

1.1 Background

The full genome of *Saccharomyces cerevisiae* was sequenced already ten years ago and the studies of genes and regulatory regions of baker's yeast has been ongoing for even longer (Goffeau *et al.* 1996; Liti & Louis 2005). The current number of *Saccharomyces cerevisiae* open reading frames (ORFs) in Saccharomyces Genome Database (SGD) is 6604, however only approximately 70% of them have gene names and a similar proportion have at least one Gene Ontology (GO) annotation attached to it (Hong *et al.*; Ashburner *et al.* 2000; 2001; Harris *et al.* 2004). It means that we know very little about the remaining 30% of the ORFs. Of course some of the ORFs can be false positives and not actually code proteins, but there is still a lot to study about these nearly 2000 ORFs.

Even if a gene has a GO annotation, it does not mean that the function or biological process is known. Almost every fourth ORF described with some GO has the annotation 'molecular function unknown', and approximately 5% are connected with the GO annotations 'biological process unknown' or 'cellular component unknown'. These facts increase the number of ORFs that should and could be described up to almost 3000.

The most reliable sources of gene function come from *in vivo* and *in vitro* experiments. The bioinformatics approaches of analyzing data from chromatin-immunoprecipitation on chip (ChIP-on-chip) or gene expression, and the predictions of putative transcription factor binding sites (TFBS) have contributed to the overall knowledge as well.

The function predictions often begin with the comparison of similarly regulated genes. The regulatory regions have been studied extensively in *Saccharomyces cerevisiae* (Stormo 2000; Vilo *et al.* 2000; Tompa 2001; van Helden, André, & Collado-Vides 1998). The sequencing of several additional yeasts have made available the comparative genomics approach to study *Saccharomyces cerevisiae* more closer. The phylogenetic footprinting approach helps to find regulatory motifs that are conserved in different species and therefore increase the trustworthiness of the predicted motifs (Duret & Bucher 1997).

Previous studies have shown that using additional yeasts, it is possible to discover stronger regulatory motifs in *Saccharomyces cerevisiae* (Kellis *et al.* 2003; Dujon *et al.* 2004). A previous study also showed that genes interacting with each other, especially if they have co-expression, are often functionally related (Kemmeren *et al.* 2002).

The abundance of various data sources, bioinformatics methods and the growing computational power have made a good basis for starting large-scale analysis. With the ever-growing data relationships it is a huge challenge to be able to tie together data from Gene Ontology, protein-protein interactions and known transcription factor target genes. Connecting this multidimensional dataset will help to create new knowledge.

1.2 Objective

The aim of the thesis was to develop methods for large-scale pattern discovery and characterization of regulatory motifs in *Saccharomyces cerevisiae*. Usually the regulatory signals are studied for one specific dataset, for example similarly expressed genes or ChIP-on-chip data. In this project we aimed to detect and describe putative transcription factor binding sites from different large-scale datasets.

Three big datasets i.e. known targets of transcription factors, Gene Ontology annotations and protein-protein interactions were combined. Different combinations of the data gives several possibilities to predict new knowledge about functional annotations, possible regulatory motifs or target genes.

Chapter 2

Background

2.1 Biological background

The main emphasis of the current molecular biology is to understand the complex mechanisms acting inside and between the cells. The protein functions depend on the structure, the solution they are in, and the interactions they make with other proteins or macromolecules. Studies for understanding the possible function of a protein starts from comparing the sequence similarity to other previously known proteins, comparisons of homologous proteins from other species, pattern discovery and matching in regulatory regions, just to mention few starting points. To describe a protein from the scratch with *in vivo* or *in vitro* methods is time and money consuming. With the help of bioinformatics the possible functions can be predicted and therefore the search range can be reduced.

Lots of work has been done with more primitive organisms like bacteria, lower eukaryotes and plants. Huge amount of knowledge has been collected to large variety of databases describing gene functions, protein structures, regulatory sequences etc. One of the most widely used data source of biological data is Gene Ontology. GO consortium builds and maintains a vocabulary of different kind of descriptors to relate genes with their functions or cellular localizations in regard of the actual species (Ashburner *et al.* 2000; 2001; Harris *et al.* 2004).

2.1.1 Gene Ontology

The Gene Ontology consists of three ontologies, that are described as a tree structured vocabulary (Ashburner *et al.* 2000). The three ontologies describe gene products with their associated biological processes (BP), molecular functions (MF) or cellular components (CC). The ontology trees are described as directed acyclic graphs, where each annotation can have one or more parent annotations as well as several children annotations. There can be more than one path from a leave up to the root, but there can be no path leading from a node to itself. The relationships between graph nodes are '*part_of*', which states that the child is a structural part of its parent, or '*is_a*' stating that a child is an instance of the parent. A small example of a molecular function subtree is shown in Figure 2.1.

The ontologies describe the following:

- Molecular function describes the activity of the gene at the molecular level, e.g. *nucleic acid binding GO:0003676*.
- Biological process is series of events accomplished by one or more ordered assemblies of molecular function, e.g. *biopolymer metabolism* GO:0043283
- Cellular component describes the location of the gene product in the cell, e.g. *membrane-bound organelle GO:0043227*



Figure 2.1: Example of a GO molecular function subtree. Here GO:0003887 is a leave and GO:0003674 is a root node. The parent of GO:0003887 is GO:0016779 and it holds genes from both of its children (GO:0003887 and GO:0003964). In this example only is_a relationships occur.

The knowledge about the function or cellular location of a gene can originate from various sources. For example data about transcription regulation can come from ChIP-on-chip experiments, DNA footprinting methods or from bioinformatics analyses. Each of the approaches have its pros and cons and ranked in GO by evidence codes. The evidence codes are the following, given by suggested reliability hierarchy:

- Inferred by Curator (IC)
- Traceable Author Statement (TAS)
- Inferred from Direct Assay (IDA)
- Inferred from Genetic Interaction (IGI)
- Inferred from Mutant Phenotype (IMP)
- Inferred from Physical Interaction (IPI)
- Inferred from Expression Pattern (IEP)
- Inferred from Sequence or Structural Similarity (ISS)
- Non-traceable Author statement (NAS)
- Inferred from Reviewed Computational Analysis (RCA)
- Inferred from Electronic Annotation (IEA)
- No biological Data available (ND)
- Not Recorded (NR)

2.1.2 Regulatory motifs

The large variety of proteins in an eukaryotic cell are transcribed by RNA polII polymerase. The RNA polII polymerase needs transcription factors to be bound to their specific binding sites in the regulatory regions to recognize the correct starting point for the transcription. Around 2% of the 6604 ORFs in *Saccharomyces cerevisiae* are so far described as transcription factors. For these approximately 150 TFs roughly a thousand binding motifs have been distinguished and described (Wingender *et al.* 2000; Dwight *et al.* 2002; Teixeira *et al.* 2006; Peterson 2004). Transcription factors can have a number of a bit different binding sites and some sites can be bound by various transcription factors. For example transcription factor Gcn4 binds both to TGATTCAT and TGACTA motifs. The later motif is also bound by Yap1. The multiple-valued relationships between the binding sites and TFs makes it complicate to always infer the correct relations between factors and target sites.

In Saccharomyces cerevisiae and other lower eukaryotes the distances between ORFs are relatively short and therefore the regulatory regions are short as well. Usually the regulatory region or upstream is considered to be 600 or more rarely 1000 base pair (bp) long. Most of the transcription factor binding sites (TFBSs) are 6-20 bp long and located in the close proximity of ORFs (Zhu & Zhang 1999; Qiu 2003; Vilo 2002). However, some TFBSs can be found near the 3' end of the ORF or even in the coding regions.

2.1.3 Protein-protein interactions

Most proteins need interactions with other proteins to be active in the cell. The structural units like proteasome or ribosome are big PPI complexes. The members of such complexes are quite often regulated by the same transcription factors. For example the proteasome genes are regulated by Rpn4 (Xie & Varshavsky 2001). Some transcription factors, like Gcn4, have DNA binding domain (DBD) and activation domain (AD). The domains are used to combine the interactions between different transcription factors. Gcn4 binds to DNA with its DBD and the transcription is activated when another protein, for example MBF1, binds to GCN4's activation domain.

The protein-protein interactions have been studied in *Saccharomyces cerevisiae* for a long time. The most widely used methods for PPI detection are yeast two-hybrid and tandem affinity purification (TAP) methods (Fields 2005; Puig *et al.* 2001).

In yeast two-hybrid method the functionality of protein is studied using DNA binding and activity domains. A hybrid DBD binding to the DNA leads to gene expression of the reporter gene only when a AD hybrid protein binds to it. The method is mainly used to study DNA binding proteins. The method is reviewed by Fields 2005.

The tandem affinity purification technique uses the fusion of specific tag to the target protein, usually to the 3' end. The construct, after insertion into the host cell, is expressed from the regular promoter. The tagged protein forms a complex with its target proteins. The complexes are concentrated and fractionated on a denaturing gel. The purified complexes are further analyzed with the mass spectrometry to detect the interacting proteins. For a detailed overview of the method look at Puig *et al* 2001.

2.2 Bioinformatics approaches

2.2.1 Pattern Discovery

Pattern discovery is one of the widely used bioinformatics approaches for transcription regulation studies. Several tools have been made to search for the motifs that occur in the input sequences more frequently than in background or random sequences (Vilo 2002; van Helden, André, & Collado-Vides 1998; Brazma *et al.* 1998). The biological motivation behind these extensive searches is the fact that similar genes (either similarly expressed, similarly located in the cell or with similar functions) are often regulated by the same transcription factors.

2.2.2 Phylogenetic footprinting

The term *phylogenetic footprinting* was first introduced in 1988 by Tagle and is defined as a phylogenetic comparison, that reveals evolutionally conserved functional elements from homologous genes (Duret & Bucher 1997). The method is based on different mutation patterns, which can be found in DNA. Genomes change constantly in time, but the results of particular mutation depends on its phenotypic effect. Most of the mutations, whose outcome is negative to the host, will be removed by natural selection. Sequences that are highly conserved in time are probably functional (Duret & Bucher 1997).

There are two main problems for choosing the organisms/sequences to compare. The first one: if one chooses species which are very close in phylogenetic tree then highly conserved elements cannot be differentiated from the overall sequence. The evolutionary time has not been long enough to fix the mutations to DNA by natural selection (Duret & Bucher 1997; Cliften *et al.* 2003; Lenhard *et al.* 2003). The second problem: if one chooses species which are too far away from each other in phylogenetic tree, then the species have been diverged too much to carry any highly conserved regions in their sequences or the species may already have different regulatory processes (Duret & Bucher 1997; Cliften *et al.* 2003). Substitutions in neutral positions in DNA happen with probability 0.5% for million years (Li, Luo, & Wu). From that we can conclude that sequences which diverged 300 million years back should have 30% of similarity if they are not under purifying selection. If there are highly conserved regions after such long time, then it refers to strong natural selection and to important functional element (Duret & Bucher 1997).

The evolutionary distance between widely compared and analyzed *Homo* sapiens and *Mus musculus* is approximately the same distance between *Sac*charomyces cerevisiae and *Saccharomyces bayanus*. The coding regions of the relative pairs are highly conserved and the non-coding regions tend to stay conserved as well. Adding more distant species to comparison we can expect that the non-functional regions are less conserved and functional regions are detectable more easily.

Chapter 3

Material and methods

3.1 Data

The *Saccharomyces cerevisiae* has been a model organism for molecular biology already a long time. The methods for exploring the baker's yeast have evolved simultaneously with the knowledge and growing datasets from wetlabs. The quickly increasing data sources need large-scale analysis pipelines.

The main research object of the project was budding yeast *Saccharomyces* cerevisiae. Saccharomyces cerevisiae was chosen because the main mechanisms and central transcription factors are well described. However there is still a lot ORFs not described and the complexity of gene regulation is not well known.

Additionally seven species were used for comparative pattern analysis. Three out of other seven species are closely related to *Saccharomyces cerevisiae* i.e. *sensu stricto* yeasts: *Saccharomyces bayanus*, *Saccharomyces paradoxus*, *Saccharomyces mikatae* (Kellis *et al.* 2003). In addition, phylogenetically more distant *Candida glabrata*, *Debaryomyces hansenii*, *Kluyveromyces lactis* and *Yarrowia lipolytica* were chosen (Dujon *et al.* 2004). The *sensu* *stricto* group we denote as Kellis dataset and the other four we refer as Pasteur data.

The wide variety of yeasts were used to cover species diverged at various timepoints in evolution and to make use of the idea that functionally active non-coding sequences tend to stay conserved during the evolution (Duret & Bucher 1997; Cliften *et al.* 2003). The evolutionary relationships between the studied yeasts are shown in Figure 3.1.



Figure 3.1: Phylogenetic tree of yeast species derived from 25S rDNA sequences. The species used in the current thesis are underlined. The picture originates from Dujon *et al* 2004.

3.1.1 Sequences

All the experiments were done using 600 base pair (bp) sequences immediately upstream from predicted or verified open reading frames (thereinafter upstream sequences or upstreams). The ORFs and upstream sequences

	Table 3.1. Bine and i		<u>01 </u>	
Species	Genome size (Mbp)	Chromosomes	ORFs*	Sequenced by
S. cerevisiae	12.2	16	6713	(Goffeau et al. 1996)
S.paradoxus	11.8	16	4788	
S.mikatae	12.1	16	4525	(Kellis et al. 2003)
S.bayanus	11.5	16	4492	
C.glabrata	12.3	13	5283	
K.lactis	10.6	6	5329	(Duion at al 2004)
D.hansenii	12.2	7	6906	(Dujon <i>et al.</i> 2004)
Y.lipolytica	20.5	6	6703	

Table 3.1: Size and number of ORFs of yeasts

* The estimated number of protein coding genes for Pasteur sequences.

for *Saccharomyces cerevisiae* were obtained from Saccharomyces Genome Database (SGD). The *sensu stricto* ORFs were predicted by (Kellis *et al.* 2003) and the four distant species were predicted by (Dujon *et al.* 2004).

For most of the *Saccharomyces cerevisiae* ORFs at least one ortholog has been predicted and about one third has 5 orthologs. The orthologs from Pasteur data have quality descriptions '*weakly similar to ORF*' or '*similar* to ORF' or '*highly similar to ORF*'. In current analysis all the predictions were treated as equal. The number of orthologs are given in the Table 3.2.

3.1.2 Previously known transcription factors and regulatory motifs

Information about transcription factors and their binding sites are gathered into many databases (Wingender *et al.* 2000; Dwight *et al.* 2002; Zhu & Zhang 1999; Teixeira *et al.* 2006; Peterson 2004). These databases describe

Number of orthologs	How many ORFs
0	162
1	116
2	168
3	444
4	798
5	1734
6	1378
7	506

 Table 3.2: The distribution of the number of orthologs for Saccharomyces

 cerevisiae ORFs

the relations between motifs and transcription factors, as well as relations between genes and regulators. Most of the databases incorporate only *in vivo* or *in vitro* verified binding sites and do not hold *in silico* predicted motifs. However the BiGeR database contains the *in silico* predicted motifs as well.

The known transcription factor target gene sets were extracted from YEASTRACT database (Teixeira *et al.* 2006). For 143 different transcription factors there was 12328 relations to 4248 different target genes. The average number of target genes per transcription factor is 84.8. Maximum number of target genes were connected to Arr1 (712) and least to Spt23, Rds1, Otu1, Hpc2 (1). In average there is 2.9 relations to TFs per gene, with maximum 20 TFs connected to YGR088W. Figure 3.2 displays a histogram for known regulators per gene.

In the current project we used motifs from YEASTRACT and BiGeR databases to compare the predicted motifs to previously known motifs (Teix-



Figure 3.2: Histogram of known transcription factors per gene.

eira et al. 2006; Peterson 2004).

3.1.3 GO groups

In the current study we used all GO annotations except the very general ones that are either the top nodes of the GO tree or have more than 670 genes associated to it. The limit was set to 670 because larger GO groups are very general, but smaller groups, like 'regulation of biological process GO:0050789', may already share common regulatory motifs. We excluded 42 annotations e.g. biological process GO:0008150, biological process unknown GO:0000004. The gene groups describing the annotations consist of the genes belonging to the annotation itself plus the genes that belong to the annotation's children. By default, not all the genes having specific annotations have the more general annotations given. We used bottom-up annotation addition to get full annotations for each gene and GO node. All annotations irrespective of the evidence codes were used to annotate each GO group with maximum number of genes.

We used 3977 different GO groups with the size ranging from 1 to 658 genes. For the GO group size distribution, see Figure 3.3.



Figure 3.3: The distribution of the size of the GO groups. The size of the groups is highly variable

3.1.4 Protein-protein interaction data

The interactions between proteins form stable complexes, many of these act in transcription control or are parts of a cellular machinery like for example the proteasome. In current thesis PPI data is used to relate proteins with unknown or poorly described functions with more specific annotations. The PPI data comes usually from yeast two-hybrid or tandem affinity purification experiments. We use data from both types of experiments and do not prefer one to another.

The three datasets used are Kemmeren (Kemmeren *et al.* 2002), Gavin (Gavin *et al.* 2006) and Krogan (Krogan *et al.* 2006) that describe interactions between 3334 different ORFs. The Gavin and Krogan datasets are from two independent tandem affinity purification experiments (Krogan *et al.* 2006; Gavin *et al.* 2006). The Kemmeren dataset, on the contrary, covers several datasets from yeast two-hybrid and TAP experiments that are verified by expression analysis (Kemmeren *et al.* 2002). For the data source and ORF numbers overview look at Table 3.3

Dataset	Experiment type	Data source	Nr of ORFs
Kemmeren	Two-hybrid; TAP	(Gavin <i>et al.</i> 2002; Uetz <i>et al.</i> 2000; Ito	1309
		et al. 2001; Ho et al. 2002; Hughes et	
		al. 2000; Spellman et al. 1998; Roberts	
		et al. 2000; Chu et al. 1998; Travers et	
		al. 2000; Gasch et al. 2000)	
Gavin	ТАР	(Gavin <i>et al.</i> 2006)	1709
Krogan	ТАР	(Krogan et al. 2006)	2186

Table 3.3: Protein-protein interaction datasets overview

3.2 Pattern discovery

In this work we used SPEXS algorithm from Expression Profiler tool-set to find overrepresented patterns in all of the analyzed datasets (Expression Profiler ; Vilo 2002; Vilo *et al.* 2000). SPEXS is a program for finding common patterns from input sequences using user defined background sequences.

While using the SPEXS we set the following parameters:

- -ms 1, the motif has to occur at least once in our input set
- -genorder 2, the patterns were generated by the most frequent first
- -binomial_prob 1.0e-04, the output pattern has to have a probability less than 1.0e-04

The parameters were set to a low end with the aim of not losing any true results. The primary results were later depleted using different p-value and occurrence filters. The background sequences in the current project were all *Saccharomyces cerevisiae* upstreams or, if we looked for evolutionary conserved patterns, all the upstreams from eight yeasts.

3.3 Statistical evaluation

The statistical evaluation of pattern discovery is highly needed to exclude the false positive motifs from the predictions. We used randomized data to get the p-values of randomly occurring patterns and used these for filtering our original results.

Each of the result sets were filtered with the p-value threshold to minimize the false positive results. The threshold was calculated with expectation of 0.01% of results to be random. For most of the GO groups we believe to have few if any real random motifs.

Algorithm 1 Calculating p-value threshold T_q for input size q

Require: r > 0 {Number of calculation runs}

Require: g > 0 {Number of genes in input set}

Require: $U \{ The set of all upstreams \}$

Require: *p* { *The pre-defined p-value threshold for SPEXS* }

 $P = [] \{ Declare \ empty \ array \ for \ storing \ values \}$

for all $i \in (1, 2, ..., r)$ do

 $U_{g,i} = \operatorname{rand}(g, U) \{ Create \ g \ random \ upstreams \ from \ all \ upstreams \ U \}$

 $P_{g,i} = \text{SPEXS}(U_{g,i}, p = 1.0) \{ \text{Run SPEXS with random input and } p$ -value threshold = 1.0 $\}$

 $P_{g,all} = \min(P_{g,i}) \{ Store \ the \ minimum \ p-value \ from \ SPEXS \ output \}$ end for

 $T_g = \operatorname{avg}(P_{g,all})$ {The average p-value from r runs for the input size g} return T_g {Return p-value threshold T_g for input size g} After having the p-value threshold from randomized data we filter the SPEXS output according to algorithm 2.

Algorithm 2 Filtering the SPEXS output
Require: $g > 0$ { <i>The group size has to be bigger than zero</i> }
Require: $1 \le T_g \le 0$ {Significance threshold for given group size g}
for group with size g do
use T_g threshold
for all patterns in group g do
if pattern $p - value < T_g$ then
keep the pattern
else
remove it
end if
end for
end for
return g {Return all the important patterns}

3.4 Expansion of groups by PPI data

We know that protein complexes, like proteasome, are regulated by the same factor. Knowing the protein complex specific regulatory motif we are able to find other proteins that could be members of these complexes. A contrary hypothesis can be stated: if genes interact and share a common motif then they can be regulated by the same factor. To study this idea we propose an expansion algorithm 3.

This approach can be applied at various datasets e.g. GO groups or known transcription factor target sets. With this method we are able to find

Algorithm 3 Expansion of gene groups by PPI data

Require: $G \{GO \ group\}$ **Require:** I {All PPI pairs} **Require:** $m \{GO \text{ specific regulatory motif}\}$ $G' = [] \{ Declare empty array for expanded GO group \} \}$ C = [] {Declare empty array for genes that interact with genes inside GO and have the same m pattern} for all $g \in G$ do {For all genes in GO group} $g \to G'$ $\{Add \ g \ to \ expanded \ GO \ group\}$ if $(g,i) \in I$ then {If there exists an interaction between g and i} $i \to G \{Add \ i \ to \ expanded \ GO \ group\}$ if $upstream(i) = \sim /m/$ then *{i has GO specific motif in the upstream}* $i \to C \{Add \ i \ to \ possible \ GO \ candidate \ genes\}$ end if end if end for **return** (C, G') {Return expanded GO group G' and group of interacting

ORFs that share input motif C

putative members of input sets based on protein-protein interactions and common regulatory motifs.

Chapter 4

Results

The regulatory complexity of an organism can be studied step by step. One of the first steps is to find relatively simple relations between regulators and their targets. The regulators can be transcription factors that act alone or through protein-protein interactions. The targets can be genes having previously described binding with the factor or genes acting similarly in expression, having similar GO annotations or interacting with each other. In the current thesis we look at the both type of, known and non-verified, targets.

4.1 Randomization

It is easy to make lots of predictions on functions, binding sites, relations between genes etc with *in silico* methods. The important step is evaluation of the predictions to exclude as many false-positive results as possible, keeping still most of the true-positives. The extensively used approach is to estimate the false-positive rate by using random data.

In the current thesis we randomized the upstream regions into groups of various sizes. The group sizes were taken the same as the size of real target sets. Using the pattern discovery approach on the random groups we calculated the thresholds to filter the primary results. The limit between the random and non-random motifs were calculated from the average p-value of the best motifs from ten motif discovery runs. The method for randomization is described in previous section 3.3 with algorithm 1.



Figure 4.1: p-values from real and random data with the input size of 160 genes

The comparison of the p-values for real and random set is given on the Figure 4.1. The given figure illustrates the distribution of p-values for motifs discovered from group size of 160. From the graph we can conclude that the best motifs from the real dataset have approximately 10^{20} smaller p-values than the motifs from random data. We do believe that such a difference between the random and real data supports the trustworthiness of the best

motifs from real datasets.

The random p-values are between 9.1e - 05 and 9.3e - 06. The p-values for different group sizes do not differ much, with an average of 3.9e - 05. The distribution of the p-values for different group sizes is illustrated in Figure 4.2.



Figure 4.2: P-values from random data for each GO group size

The p-values for each group sizes were used in the following steps to filter out the non-random signals. In the next steps of analysis we considered true motifs only the signals that had p-values lower than the random threshold for given group size.

4.2 Pattern discovery for known transcription factors' target groups

Many of the transcription factors have already been identified with their relative binding sites (Dwight *et al.* 2002). The *in vivo* and *in vitro* methods like DNA footprinting or ChIP-on-chip have given more or less specific binding sites for most of the factors. Still there are TFs with a number of known target genes and uncharacterized binding sites. These groups of targets and TFs are very challenging datasets for bioinformatics and especially for pattern discovery methods.

During the project we used target gene sets for 143 transcription factors as described in previous section 3.1.2. Each of the group was analyzed with two datasets. The first set consisting only *Saccharomyces cerevisiae* upstreams and the second with the addition of orthologous upstreams. The motifs were filtered with the random threshold and only putatively non-random signals were analyzed further. For 138 TFs we could find at least one motif with the p-value smaller than the random threshold.

We found 7381 patterns for the 143 TFs, 2020 of them are distinct ¹. For almost half of the TFs we are able to find the known motif as the strongest from our prediction. In some cases we were able to find a motif for a transcription factor that did not had any motif related before.

The best example is Gat3 with very strong motif TACTTCGAAGC in Saccharomyces cerevisiae (p-value 1.6e - 26) that is also conserved in orthologs (p-value 2.1e - 31). The motif is not characterized in the databases before. From the Gene Ontology datasets we find that the motif overlaps

¹Distinct - edit distance is at least one and motif is not a sub-motif for another motif (e.g. CACGTG and CTCGTG are distinct but CACGTG and CACGT are not)

well with telomerase-independent telomere maintenance (BP) and helicase activity (MF) specific motifs. The Gat3 is not well characterized transcription factor. Based on our results we propose that it can be related to the regulation of the genes with helicase activity taking part in telomere maintenance.

We were able to find motifs for other weakly characterized transcription factors as well. The newly found signals quite often overlap with already known motifs belonging to well described TFs. This suggests that these transcription factors can act as protein-protein complexes for example. For more results look the Table 1 in Appendix.

4.3 Pattern discovery in GO dataset

The Gene Ontology is a major knowledge base of biological data. The GO annotations relate genes with similar molecular function, genes participating in the same biological process and genes having the same location in the cell. All the annotations describe a small subset of all the genes. The subsets can often be regulated by the same transcription factors.

The fact that functionally related genes can be regulated by same regulators has been used earlier to find function specific regulatory motifs and to relate these to known transcription factors (Kellis *et al.* 2003; Cliften *et al.* 2003). In the current project we look for overrepresented patterns for each GO group. The groups are analyzed similarly to TF target sets. Both *Saccharomyces cerevisiae* and orthologous sequences are used to detect interesting motifs.

For 2704 GO groups out of 3977, we found at least one pattern below the SPEXS threshold 1.0e - 04 and 2006 groups remain after filtering out the
GO groups with all patterns having better (i.e. smaller) p-value than the random threshold.

The GO groups with patterns having the smallest p-values are often covering a small subtree from the GO tree. This was however expected with our bottom-up annotations, i.e. if the most specific GO group has highly overrepresented patterns, then these patterns occur more frequently than random in the more general parent GO groups as well. Usually in this case the p-value gets worse (i.e. increases) while going up to root in the GO tree, indicating that the motif was child-node specific.

In the set of top 100 GO groups with the top p-values the following GO annotations are represented in the Table 4.1. The GO groups showed in the Table 4.1, ranking between 10 and 100, are the first groups for each TF motif. If there is several groups with the same (sub)motif, then only the first GO group is shown.

	GO	CC MF BP	GO description	p-value	ORFs	pattern	related TF	
Ĭ	GO:000943	CC	retrotransposon	2.57319e-136	91	TGTTGGAATA	Mot3	
			nucleo capsid					
	GO:0006319	вР	Ty element trans-	2.12078e-135	92	TGTTGGAATA	Mot3	
			position					
	GO:0006313	BP	DNA transposi-	9.80204e-126	105	TGTTGGAATA	Mot3	
			tion					
	GO:0006310	вР	DNA recombina-	1.21499e-95	186	TGTTGGAATA	Mot3	
			tion					
	GO:0003723	MF	RNA binding	4.21101e-75	316	TGTTGGAATA	Mot3	
	GO:0003964	MF	RNA-directed	1.45907e-60	48	TGTTGGAATA	Mot3	
			DNA polymerase					
			activity					
	GO:0005515	MF	protein binding	2.37398e-57	493	TGTTGGAATA	Mot3	
	GO:0006259	вР	DNA metabolism	5.76027e-56	557	GAGGAGAACTTCTA	Mot3*	
	GO:0003676	MF	nucleic acid bind-	1.53731e-53	567	TGTTGGAATA	Mot3	
			ing					
	GO:0005730	CC	nucleolus	3.73244e-53	222	AAATTTT		
	GO:0008233	MF	peptidase activity	1.09457e-34	155	GCAAGGATTGATAAT		
	GO:0005830	MF	cytosolic ribo-	1.89356e-31	160	CCGTACA	Rapl	
			some (sensu					
			Eukaryota)					
	GO:0008652	BP	amino acid	5.62027e-31	102	TGACTCA	Gcn4	
			biosynthesis					
	GO:0000502	CC	proteasome com-	5.39292e-28	407	GGTGGCAAA	Rpn4	
			plex (sensu Eu-					
			karyota)					
	GO:0016788	MF	hydrolase activ-	1.75131e-25	274	ATAATGTAATA	Hcm1?	
			ity, acting on					
			ester bonds					

Table 4.1: Patterns with the smallest p-values from GO groups

Continued on Next Page...

ż	GO	CC MF BP	GO description	p-value	ORFs	pattern	related TF
34	GO:0016772	MF	transferase activ-	1.10842e-23	331	GATTGATAATG	
			ity, transferring				
			phosphorus-				
			containing groups				
36	GO:0031974	сc	membrane-	1.74176e-23	658	GCGATGAG	Esr1/Mec1?
			enclosed lumen				
50	GO:0006260	BP	DNA replication	4.16142e-21	105	ACGCGT	Mbp1
70	GO:0004386	MF	helicase activity	1.53663e-14	06	CCTCGACTAA	Xbp1
98	GO:0009068	вр	aspartate fam-	1.91911e-11	10	AGCACGTGAC	Pho4
			ily amino acid				
			$\operatorname{catabolism}$				

Table 4.1 – Continued

? marks probable regulator, the relation between the motif and TF is not known before

4.4 Expansion of GO groups by PPI

The Gene Ontology groups were expanded based on the protein-protein interaction data from Kemmeren, Gavin and Krogan datasets described in 3.1.4. Expanding the group with interacting proteins gives us hints about the proteins both in and outside of the group. The example of GO:0008652 *amino acid biosynthesis* expansion by PPI and the additional information of Gcn4 binding site TGACTC is illustrated in Figure 4.6

The genes belonging to the GO group are noted as rectangles and proteins interacting with them are shown as circles. If a gene has the TGACTC motif in the upstream, then the figure is colored red. On the picture one can see that there are two proteins interacting with the GO members and share the Gcn4 motif - YLR058C and YDR172W. The first protein is part of a lysine degradation and glycine, serine and threonine metabolism pathways and related to amino acid metabolism. Therefore the protein could be a target of Gcn4 TF and belong to this GO group as well. The second protein has translation termination function and about this protein we have currently no clue if it could be regulated by GCN4 or does it have a function related to amino acid biosynthesis.

4.5 From protein-protein interactions to regulatory motifs and GO annotations

Previously we showed how known GO group can be expanded with PPI and known regulatory motifs. In this step we show how to start analyzing genes from the opposite direction. As we have mentioned earlier, interacting proteins are often regulated by the same transcription factor. We also showed that genes with common regulator share similar motifs in their regulatory regions and often interact with each other. With this knowledge we can show how to examine genes that interact with each other.

Using the Kemmeren protein-protein interaction data we look for groups of genes that interact with each other and form connected graphs (PPI graphs). The ORFs are divided into groups recursively starting from a random ORF.

The groups are then searched for non-random Gene Ontology annotations with GOSt tool set (Reimand 2006). For each of the PPI graphs the GO annotations that have the probability less than random threshold are kept. The sets of ORFs are then taken as input for pattern discovery step with SPEXS to look for potential regulatory motifs. In this step of pattern discovery we look for non-discrete motifs. We allow wildcard² positions in the motifs, at most 2 wildcards per motif. The motifs are sorted by their p-values increasingly and the motif with the smallest p-value at the given input size is taken.

4.5.1 Pipeline example

A set of ORFs without any previous knowledge is analyzed. The input size is 5 ORFs: YER099C, YKL181W, YBL068W, YHL011C, YOL061W. The ORFs are connected with each other in the following way, Figure 4.3. The five ORFs make connected graph because they are pairwise connected.

The input set is analyzed with Gene Ontology tool GOSt (Reimand 2006). The tool is used to find all GO annotations mapping the input ORFs. The annotations with probability value smaller than GOSt analytical threshold is given in the output. From the Figure 4.4 we can see that the best GO

 $^{^{2}}$ wildcard is a special character representing more than one character



Figure 4.3: The connected graph of input ORFs. The arrows indicate the known interactions between genes

annotation is GO:0004749 ribose phosphate diphosokinase activity.

From the pattern discovery we find 3 motifs occurring in all input sequences. The motif with the smallest p-value is taken into further analysis. In this case it is AATG.TTA, where . denotes the wild character i.e. A, C, G or T can be at that position.

The pattern can be matched back to all *Saccharomyces cerevisiae* upstream sequences to check if the pattern is highly specific to the input set or it belongs to some more general GO node.

The expression dataset can be checked to evaluate the potential coexpression of the ORFs belonging to the connected graph set. One still has to keep in mind that the Kemmeren dataset had already verified the data with expression analysis, so a similar expression should be no big surprise, but rather an expected result.

YOLD61W YHL011C YBL068W YKL181W YER039C	P-value	т	Q	Q&T	Q&T/Q	Q&T/T		term ID	te	erm domain and name
A A A A A	4.44e-11	57	5	5	1.000	0.088	ee	GD:0006725	BP	1 aromatic compound metabolism
AAAAA	1.48e-10	72	5	5	1.000	0.069	닫	G0:0046483	BP	heterocycle metabolism
AAAAA	2.67e-15	10	5	5	1.000	0.500	12	G0:0042430	BP	indole and derivative metabolism
	2.6/e-15	10	5	5	1.000	0.500	P P	GU:0042434	BP	Indole derivative metabolism
A A A A A	1.45e-11	46	5	5	1.000	0.109	-	60:0005117	BP	puripe purleotide metabolism
A A A A A	9.01e-12	42	5	5	1.000	0.119	臣	G0:0009259	BP	ribonucleotide metabolism
A A A A A	6.97e-12	40	5	5	1.000	0.125	臣	G0:0009150	BP	purine ribonucleotide metabolism
A A A A A	1.23e-13	19	5	5	1.000	0.263	臣	60:0009123	BP	nucleoside monophosphate metabolism
AAAAA	9.08e-14	18	5	5	1.000	0.278	12	G0:0009126	BP	purine nucleoside monophosphate metabolism
8 8 8 8 8	6.55e-14	17	5	5	1.000	0.294		60:0009161	BP	ribonucieoside monophosphate metabolism
A A A A A	2.12e-14	14	5	5	1.000	0.357		GD:0046040	BP	IMP metabolism
AAAAA	6.46e-08	238	5	5	1.000	0.021	Έ	G0:0006807	BP	nitrogen compound metabolism
AAAAA	2.51e-12	33	5	5	1.000	0.152	臣	60:0009112	BP	nucleobase metabolism
AAAAA	6.55e-14	17	5	5	1.000	0.294	臣	60:0006206	BP	pyrimidine base metabolism
H H H H H	2.14e-07	302	5	5	1.000	0.017		G0:0006082	BP	organic acid metabolism
	2.14e-07	302	Э Б	5	1.000	0.017		60:0019752	BP	carboxyiic acid metabolism
8 8 8 8 8	2.16e-13	21	5	5	1.000	0.238	÷	60:0046112	BP	nucleobase biosynthesis
AAAAA	3.18e-14	15	5	5	1.000	0.333	臣	60:0019856	BP	purimidine base biosynthesis
AAAAA	8.39e-15	12	5	5	1.000	0.417	Έ	60:0006207	BP	'de novo' pyrimidine base biosynthesis
A A A A A	4.85e-11	58	5	5	1.000	0.086	臣	G0:0009165	BP	nucleotide biosynthesis
A A A A A	6.97e-12	40	5	5	1.000	0.125		G0:0009260	BP	ribonucleotide biosynthesis
<u> </u>	6.55e-14	17	5	5	1.000	0.294	1	60:0009124	BP	nucleoside monophosphate biosynthesis
H H H H H	3.18e-14	15	5	5	1.000	0.333		G0:0009156	BP	ribonucleoside monophosphate biosynthesis
8 8 8 8 8	1.02e-11	43	5	5	1.000	0.115	-	GU:0006164	BP	purine nucleotide biosynthesis
A A A A A	5 32e-12	38	5	5	1.000	0.132		60:0009127	BP	purine ribonucleotide biosunthesis
A A A A A	3.18e-14	15	5	5	1.000	0.333	Έ	60:0009168	BP	purine ribonucleoside monophosphate biosynthesis
AAAAA	2.12e-14	14	5	5	1.000	0.357	Έ	60:0006188	BP	IMP biosynthesis
AAAAA	2.12e-14	14	5	5	1.000	0.357	臣	G0:0006189	BP	'de novo' IMP biosynthesis
A A A A A	2.67e-15	10	5	5	1.000	0.500	臣	60:0042435	BP	indole derivative biosynthesis
AAAAA	1.36e-09	111	5	5	1.000	0.045	드	60:0044271	BP	nitrogen compound biosynthesis
H H H H H	6.55e-14	17	5	5	1.000	0.294	1	G0:0019438	BP	aromatic compound biosynthesis
8 8 8 8 8	4.30e-08	111	5	5	1.000	0.023	-	60:0009308	BP	amine metabolism
A A A A A	2.69e-08	200	5	5	1.000	0.025		60:0006519	BP	amine biolognenesis amine acid and derivative metabolism
AAAAA	1.76e-08	184	5	5	1.000	0.027	Ŧ	G0:0006520	BP	amino acid metabolism
A A A A A	2.16e-13	21	5	5	1.000	0.238	臣	G0:0009072	BP	aromatic amino acid family metabolism
A A A A	8.82e-10	102	5	5	1.000	0.049	닫	60:0008652	BP	amino acid biosynthesis
A A A A A	4.63e-14	16	5	5	1.000	0.312	드	60:0009073	BP	aromatic amino acid family biosynthesis
н н н н н о о о о о	2.67e-15	10	5	5	1.000	0.500	T T	G0:0009096	BP	aromatic amino acid family biosynthesis, anthranilate pathway
	2.120-14	14	5	5	1.000	0.337		G0:0009075	DP RD	histidine family amino acid bicsunthesis
A A A A A	2.12e-14	14	5	5	1.000	0.357	臣	60:0006547	BP	histidine metabolism
AAAAA	2.12e-14	14	5	5	1.000	0.357	臣	G0:0000105	BP	histidine biosynthesis
A A A A A	1.04e-12	28	5	5	1.000	0.179	臣	G0:0006575	BP	amino acid derivative metabolism
A A A A A	2.16e-13	21	5	5	1.000	0.238		G0:0042398	BP	amino acid derivative biosynthesis
AAAAA	5.63e-13	25	5	5	1.000	0.200		G0:0006576	BP	biogenic amine metabolism
	1.23e-13	19	5	5	1.000	0.263	E E	GU:0042401	BP	blogenic amine blosynthesis
8 8 8 8 8	2.67e-15	10	5	5	1.000	0.500	臣	60:0006568	BD	indolaikylamine metabolism
AAAAA	2.67e-15	10	5	5	1.000	0.500	臣	60:0046219	BP	indolalkulamine biosunthesis
AAAAA	2.67e-15	10	5	5	1.000	0.500	臣	G0:0000162	BP	tryptophan biosynthesis
PPPPP	2 02e-04	1182	5	5	1 000	0 004	- - -	KECC •00000	Phi	2 KEGG nathways
PPPP	8.55e-13	27	5	5	1.000	0.185	臣	KEGG:00030	PW	Pentose phosphate pathway
PPPP	4.66e-10	90	5	5	1.000	0.056	臣	KEGG:00230	PW	Purine metabolism
	3 18e-14	15	5	5	1.000	0 333	단단	60+0009116	RD	3 nucleoside metabolism
AAAAA	1.33e-15	9	5	5	1.000	0.556	臣	60:0042278	BP	purine nucleoside metabolism
AAAAA	2.67e-15	10	5	5	1.000	0.500	12	G0:0009119	BP	ribonucleoside metabolism
A A A A A	5.93e-16	8	5	5	1.000	0.625	臣	60:0046128	BP	purine ribonucleoside metabolism
A A A A A	6.55e-14	17	5	5	1.000	0.294	닫	60:0043094	BP	metabolic compound salvage
AAAAA	6.36e-17	6	5	5	1.000	0.833		60:0043174	BP	nucleoside salvage
	2.22e-16	7	5	5	1.000	0.714	1	60:0043101	BP	purine salvage
	0.3be-1/	р	9	3	1.000	0.833		PD:0009199	BP	purine riponucieoside salvage
AAAA	1.30e-05	684	5	5	1.000	0.007	면면	60:0016740	MF	4 transferase activity
H H H H A	3.61e-07	335	5	5	1.000	0.015		G0:0016772	MF	transferase activity, transferring phosphorus-containing groups
8 8 8 8 8	2.228-16	109	9 5	5 5	1.000	0.025		GD:0016778	ME	aipnosphotransterase activity kinase activity
A A A A A	4.63e-14	150	5	5	1.000	0.312		60:0019200	MF	carbohudrate kinase activitu
AAAA	1.06e-17	5	5	5	1.000	1.000	臣	G0:0004749	MF	ribose phosphate diphosphokinase activity

Figure 4.4: The GO annotations for the 5 ORFs with GOSt tool (Reimand 2006). The bottom line on the figure describes the most relevant GO annotation for the input set

From the Figure 4.5 we can see that the expression of the interacting ORFs is very similar in various stress conditions (Gasch *et al.* 2000).

The motif has no relationships to any GO annotations and is not known regulatory motif before. The very similar expression clusters suggest that the proteins might be regulated by a same factor and work in protein complexes.

We checked from the literature that the our analyzed genes are five phosphoribosyl diphosphate synthase-homologous genes. These genes have been previously described as proteins being active only in a complex of at least 3 subunits out of 5 (Hove-Jensen 2004).

4.6 Web-tool Gviz-PPI

The method for expansion a group of genes by addition of interacting proteins is incorporated to a web-tool Gviz-PPI and is publicly accessible at http://bioinf.ebc.ee/u/peterson/gviz/. The tool allows to input a set of *Saccharomyces cerevisiae* ORFs and a motif to find out all the interactions the input gene set has in Kemmeren, Gavin and Krogan datasets. The ORFs that have input motif in their upstream regions are visualized to give a fast and easy overview of the data. It is possible to use regular expressions for the motifs e.g. to use motif like TGA.TC. The tool is implemented in Perl and visualization is enabled by the Simple Web Object Graphics language SWOG (Hansen 2005).

4.7 Data update to BiGeR

The BiGeR database incorporates data about *Saccharomyces cerevisiae* transcription factor binding sites, both predicted by bioinformatics methods and from the wet-lab experiments (Peterson 2004). In the database the TFBS are described as regular DNA motifs, regular expressions, consensus sequences or with position weight matrices. The binding sites are related to known transcription factors and genes they regulate, if the relations are known.

With the previously described analysis a number of binding sites get additional descriptions and relations from GO. The updated motifs contribute to more deeper knowledge in the database and help to distinguish more relevant patterns for each transcription factor.

From GO pattern discovery motifs that had p-value smaller than 1.0e-07 were considered as interesting. The even stricter p-value threshold was chosen because of the large number of motifs with p-values around random threshold. The motifs chosen for BiGeR update should be more reliable. The motifs were then filtered to find distinct motifs. Finally 700 distinct motifs with 5714 connections to 178 different GO annotations were added to the database.

The updated and newly discovered regulatory motifs have been made publicly available in the BiGeR database at http://bioinf.ebc.ee/biger/. An example output of a query describing TGACTC motif is given in Figure 4.7.



Figure 4.5: The clustering based on stress response (Gasch 2000)



Figure 4.6: The expansion of GO:0008652 using Gcn4 motif TGACTC. The rectangles are ORFs belonging to the GO group, the circles are interacting proteins. The red fulfil indicates that the TGACTC pattern is in the upstream region of the ORF

2	signal	start	coord	description	signal type	dbxref	Regulated Gene	Regulating factor
-	TBACTC	-126	108	tor I for writing	digo IR	LF.R00648	1133/YOR202W	GCN4
7 7	ATGACTCAT	102	-92	tase I footprinting direct.get shift used interference	ligo <u>TR</u>	LF.R00649	HIS3/YOR202W	<u>skoi</u>
m	TAATAGTGACTCCGGTAAATT	249	-229	tase I footprinting ged retardation	ligo IR	LF.R00650	HIS4/YCL030C	BAS1
4	TRACTC	350	33	Ize 1 footprinting	lign IK	LF:KUU001	HIS4/VCL//30C	GCN4
н 9	TAATAGTGACTCCGGTAAATTAGTTAATTAA	249	-217	tas el tootprinting gel retarabtion	ligo IR	LF.R00652	HIS4/YCL030C	BAS1
- v	TGACTC	203	-190	tase I to opprinting	ligo IR	LF.R00655	HIS4/YCL030C	<u>GCN4</u>
-	TTGACTCTCtaaaaATGATTCAT	140	-108	tas e la congrateitag gei areurdation.	ligo <u>TR</u>	1F-R02022	RP4/YDR354W	GCN4
	TTGACTCTT	358	-350	las e I foroprinting	ligo IR	LF:R04022	NDE4/YMR300C	GCN4
0	SNSNNNNRTGACTCATNS				onsensus TR	LF:R04800	NDE4/YMR300C	GCN4
10	TAATAGTGACTCCGGTAAATT	249	-229	ence 246.931-935 (1989)	ligo		HIS4/YCL030C	BAS1
1	TCGAACTGACTCTAATAGTGAC	261	-240	(¥3 1992,80 0740-0720	digu		HIS4/YCL0S0C	BAS1
12	TAATAGTGACTCCGGTAAATTAGTTAATTAA	249	-219	ence 246:931-935 (1989)	ligo	H	HIS4/YCL030C	BAS1;PHO2
1 2	TGACTC	199	-194	oc. Natl. Acad. Sci. U SA 83:8516-8520 (1986)	ligo	Pell	HIS4/YCL030C	GCN4
14 14	ACAGTGACTCACGTTT	503	-188	al Cell Biol 1991, 11:364:3651	ligo	PAI	HIS4/VCL030C	GCN4
12	TGACTC	255	-250	oc. Nati. Acad. Sci. USA 83:8516-8520 (1986)	ligo	H	HIS4/YCL030C	GCN4
1	TTGACTCTC	166	-158	LBD 1 0-0411-2625 (1000)	lign		RP4/VDR354W	GCN4
1	TGACTC	184	-179	iol Chem 1997, 272:13343-13354	ligo		ADE5.7/YGL234W	BAS1
18	GCCGACTGACTCGTGTCCTGGT	190	-169	LAS 1992, 896746-6750	ligo		ADE5.7/YGL234W	BAS1
101	TGACTC	217	-212	iol Chem. 1997, 272: 13345-13354	ligo		ADE5.7/YGL234W	BAS1
20	TTCAGTTGACTCGCCCCGGTCGG	333	-202	LAS 1992, 896746-6750	digo		ADE5.7/YGL234W	BAS1
21 1	TTGACTCTT	38	-350	biol. Chem. 266/20453-20466 (1991)	digo	-41	ADE4/YMR300C	GCN4
33	TGACTCA	342	-336	ast 1996,1,11367-1380	ligo		ARG1/YOL058W	GCN4
23 1	TGACTCA	180	-174	ast 1995,11:1367-1380	ligo		ARG8/YOL140W	GCN4
54 74	AGTGATTGACTCTTGCTGACCT	159	-100	LAS 1992,006746-6750	digo		ADE2/YOR 120C	DAS1
22 C	GACAAATGACTCTTGTTGCATG	202	-181	LaS 1992, 8:6746-6750	ligo		ADE2/YOR 128C	BAS1
26 I	TGACTC	122	-117	o.c. Nati. Acad. Sci. USA 83:8516-8520 (1986)	ligo		HIS3/YOR202W	GCN4
27 A	ATGACTCTT	13	-115	al Cell Biol 1995, 15:7059-7066	ligo	Pell	HIS3/YOR202W	GCN4
28 A	ATGACTCWT			rine; MCS 6.1; Best cargoy: ChiP. Gond; CCS: 44; Interpreted:Known: Gond/Bas1	egexp	PHI	HIS3/YOR202W	GCN4
ĺ		ĺ	ĺ		Ì			

Figure 4.7: BiGeR output of all patterns matching known Gcn4 binding site TGACTC. The top motifs come from various *in vivo* and *in vitro* analyzes. The last two rows describe motifs from ChIP-on-chip and Gene Ontology specific data

Chapter 5

Discussion

Combining various data sources like protein-protein interactions with known TF binding sites or with Gene Ontology data we are able to verify the connections between data sources, make new function or annotation predictions and gain additional knowledge about the complex mechanism of a cell. In the current scientific world there is huge data abundance and therefore a lot of effort has been made to combine the data and find the reliable connections between methods and datasets.

In the current thesis we looked at regulatory motifs with the help of protein-protein interactions, Gene Ontology and known transcription factor target sets. Firstly, we were able to find the previously known motifs, which assured that our methods are working well. The large scale analysis on Gene Ontology gave us a great amount of motifs related to one or many ontologies.

With the help of known facts like Gcn4 is the main regulator of aminoacid biosynthesis and knowing the binding sites for Gcn4, we could verify the known ontology specific motifs and start to predict the unknown.

The new patterns were predicted with SPEXS algorithm and related to GO nodes. The probable motifs were filtered with random p-value distributions. With a very small false-negative approximation we were able to find GO specific motifs for almost half of the input annotations. Few of the groups were analyzed more deeply and verified with the data from previously published studies.

We found several clusters of motifs that differ by one or two bases. The changes of the binding sites may affect the binding affinity of TF. For example changes in proteasome specific motif GTGGCAAA are related to changes in gene expression similarity (data not shown).

The protein-protein interaction datasets gave us a possibility to use the previously known fact that interacting proteins tend to be regulated by the same factors and have alike functions with the pattern discovery and GO annotations to expand the GO groups.

The GO group expansions were made to find genes that do interact with the members of the GO group, but do not have the relationship with the annotation. We used putative regulatory motifs to describe the interacting proteins and add them as possible members of GO groups.

In many cases we saw that same motif belongs to different GO groups. If the groups belong to the same GO subtree, then this is observation is expected. We think it would be interesting to look for motifs that belong to not connected GO annotations.

In the future work we propose to use less discrete patterns to describe the GO groups, because of the variability of binding sites for the TFs. The usage of position weight matrices or regular expressions will make it easier to describe similar motifs in a compact way. The datasets could be filtered according to some quality thresholds, e.g. GO evidence codes or the similarity of orthologous genes. The datasets would be smaller but the predictions could be more reliable. The protein-protein interactions could be filtered according to the expression dataset to have more trustworthy data i.e. to use the same approach as Kemmeren et al (Kemmeren *et al.* 2002).

Chapter 6

Conclusions

The combination of datasets, methods, quality assessments is a real challenge in current bioinformatics. A big number of datasets can give lots of information about the topic, however it can be hard to find the most important data and relations from it. We combined data with evolutionary background, pattern discovery, the functional annotations from GO and protein-protein interactions to develop methods and pipeline for large scale analysis of transcriptional regulation in *Saccharomyces cerevisiae*. The pattern discovery step was quite successful and the usage of random data and phylogenetic data helped to deplete the probable false-positive results even more. The PPI data gave us hints about ontologies that could be expanded and how to predict functions to unknown proteins.

The outcome of the thesis is knowledge about how to combine relevant data sources to understand better the transcription factor binding sites and the complexity of the regulatory mechanism. We propose a few approaches towards the function prediction and describe a large variety of putative regulatory regions.

Summary

We know the whole genome of *Saccharomyces cerevisiae* for 10 years already. The baker's yeast has been studied extensively and thoroughly, but still not all the genes and molecular mechanisms of this fairly simple eukaryote has been described. The challenge for understanding this one cell organism is more and more drifted to bioinformatics. The large datasets covering different aspects of regulatory mechanisms and regions have been published and the main challenge nowadays is to to put all the information together, to connect the small pieces of this huge puzzle.

In this thesis we gave an overview of the possibilities to join the widely used and up-to-date source of Gene Ontology, the sets of protein-protein interactions and pattern discovery methods. The approaches were used to study the connections one can find between the sets, to make predictions of regulatory motifs and widen the Gene Ontology groups.

We were able to show how to find GO specific putative regulatory motifs using PPI data or how to broaden the known GO annotation based on PPI and known regulatory motifs. These experiments helped us to understand the ways how one can start annotating functions or regulatory regions from a set of genes with a little or no previous knowledge.

The thesis gave a glimpse of the complexity of connecting large datasets, analyzing the results and predicting new knowledge. The data sources grow rapidly, so there will be more and more challenges to solve with a similar approach.

Summary in Estonian

Saccharomyces cerevisiae regulatoorsete motiivide ennustamine ja kirjeldus

Saccharomyces cerevisiae genoom sekveneeriti juba kümme aasta tagasi. Hoolimata pagaripärmi põhjalikust teaduslikust uurmisest viimase kahe kümnendi jooksul on siiski paljude geenide ja molekulaarsete mehhanismide funktsioonid teadmata. Eksperimenditulemuste üha suurenev kasv loob võimalusi bioinformaatika laialdaseks kasutamiseks funktsioonide ennustamisel. Tänapäeva bioinformaatika suurimaid väljakutseid on erinevate andmehulkade sidumine nii, et kõik kokku moodustaks tervikliku pildi rakus toimuvast.

Käesolevas töös anti ülevaade kuidas siduda laiadlaselt kasutusel olevat geeniontoloogia andmestikku, valk-valk interaktsioone ning mustrite otsimismeetodeid. Töös loodi seosed andmestike vahel ning kasutati neid regulatoorsete motiivide ennustamiseks ning geeniontoloogiate laiendamiseks.

Töös on toodud näited, kuidas kasutades valk-valk interkatsioone on võimalik leida regulatoorseid motiive. Samuti, kuidas laiendada geeniontoloogiaid spetsiifilisi regulatoorseid motiive ning valk-valk interaktsioone kasutades. Teostatud eksperimendid aitasid leida viise annoteerimaks funktsioone või regulatoorseid regioone kasutades vähest varasemat teadmist sisendgeenide hulga kohta.

Antud töö andis lühiülevaate suurte bioloogiliste andmehulkade sidumise

keerukusest, analüüsivõimalustest ning uue teadmise ennustamisest. Andmehulkade jätkuva kasvu taustal on kindlasti kirjeldatud teadmistest ja meetoditest edaspidises uurimistöös kasu.

Acknowledgements

I would like to express my sincere gratitude to my supervisor dr. Jaak Vilo for introducing the interesting world of bioinformatics to me, for all the support and guidance throughout the years and for believing in me.

I am really grateful to my friend and fellow BIITer Jüri with whom I have shared those long working nights during the last two years. Thanks for all the tips and ideas. It has and will continuously be fun to work with You.

Many thanks go to my friends in BIIT, especially to Asko, Jelena and Lemps.

I wish to thank also Jaanus, Kostja, Meelis, Pavlos and Raivo, who have supported me with tools, helped out with statistics or just given good advice.

Financial support from Estonian Science Foundation grant no 5724, EU FP6 STRE ATD and Kristjan Jaak scholarship foundation is acknowledged.

Last, but not least, I would like to thank my friends from outside of the science world and my family for their continuous support during the studies.

References

Ashburner, M.; Ball, C. A.; Blake, J. A.; Botstein, D.; Butler, H.; Cherry,
J. M.; Davis, A. P.; Dolinski, K.; Dwight, S. S.; Eppig, J. T.; Harris, M. A.;
Hill, D. P.; Issel-Tarver, L.; Kasarskis, A.; Suzanna Lewis, John C. Matese,
J. E. R. M. R.; Rubin, G. M.; and Sherlock, G. 2000. Gene Ontology: tool
for the unification of biology. *Nature genetics* 25(1):25–29.

Ashburner, M.; Ball, C. A.; Blake, J. A.; Botstein, D.; Butler, H.; Cherry,
M. J.; Davis, A. P.; Dolinski, K.; Dwight, S. S.; Eppig, J. T.; Harris, M. A.;
Hill, D. P.; Issel-Tarver, L.; Kasarskis, A.; Lewis, S.; Matese, J. C.; Richardson, J. E.; Ringwald, M.; Rubin, G. M.; and Sherlock, G. 2001. Creating the
Gene Ontology Resource: Design and implementation. *Genome Research* 11:1425–1433.

Brazma, A.; Jonassen, I.; Vilo, J.; and Ukkonen, E. 1998. Predicting gene regulatory elements in silico on a genomic scale. *Genome Research* 8:1202–1215.

Chu, S.; DeRisi, J.; Eisen, M.; Mulholland, J.; Botstein, J.; Brown, O.; and Herskowitz, I. 1998. The transcriptional program of sporulation in buddying yeast. *Science* 282(5389):699-705.

Cliften, P.; Sudarsanam, P.; Desikan, A.; Fulton, L.; Fulton, B.; Majors, J.; Waterston, R.; Cohen, B. A.; and Johnston, M. 2003. Finding functional features in *Saccharomyces* geomes by phylogenetic footprinting. *Science* 301(5629):71-76.

Dujon, B.; Sherman, D.; Fischer, G.; et al. 2004. Genome evolution in yeasts. *Nature* 430(6995):35-44.

Duret, L., and Bucher, P. 1997. Searching for regulatory elements in human noncoding sequences. *Current Opinion in Structural Biology* 7(3):399–406.

Dwight, S.; Harris, M. A.; Dolinski, K.; Ball, C. A.; Binkley, G.; Christie, K. R.; Fisk, D. G.; Issel-Tarver, L.; Schroeder, M.; Sherlock, G.; Sethuraman, A.; Weng, S.; Botstein, D.; and Cherry, J. M. 2002. Saccharomyces Genome Database (SGD) provides secondary gene annotation using the Gene ontology (GO). *Nucleic Acids Research* 30(1):69–72.

Expression Profiler. Webpage http://www.ep.ebi.ac.uk/EP.

Fields, S. 2005. High-throughput two-hybrid analysis. The promise and the peril. *The FEBS Journal* 272(21):5391–5399.

Gasch, A.; Spellman, P.; Kao, C.; Carmel-Harel, O.; Eisen, M.; G.Storz; Botstein, D.; and Brown, P. 2000. Genomic expression programs in the response of yeast cells to environmental changes. *Molecular Biology of Cell* 11(12):4241–4257.

Gavin, A.; Bosche, M.; Krause, R.; Grandi, P.; Marzioch, M.; Bauer, A.; Scultz, J.; Rick, J.; Michon, A.; Cruciat, C.; et al. 2002. Functional organization of the yeast proteome by systematic analysis of protein complexes. *Nature* 415(6868):141–147.

Gavin, A.-C.; Aloy, P.; Grandi, P.; Krause, R.; Boesche1, M.; Marzioch, M.; Rau, C.; Jensen, L. J.; Bastuck, S.; Dümpelfeld, B.; Edelmann, A.; Heurtier, M.-A.; Hoffman, V.; Hoefert, C.; Klein, K.; Hudak, M.; Michon,

A.-M.; Schelder, M.; Schirle, M.; Remor, M.; Rudi, T.; Hooper, S.; Bauer,
A.; Bouwmeester, T.; Casari, G.; Drewes, G.; Neubauer, G.; Rick, J. M.;
Kuster, B.; Bork, P.; Russell, R. B.; and Superti-Furga, G. 2006. Proteome
survey reveals modularity of the yeast cell machiner. *Nature* 440(7084):631–636.

Goffeau, A.; Barrell, B.; Bussey, H.; Davis, R.; Dujon, B.; Feldmann, H.; Galibert, F.; Hoheisel, J.; Jacq, C.; Johnston, M.; Louis, E.; Mewes, H.; Murakami, Y.; Philippsen, P.; Tettelin, H.; and Oliver, S. 1996. Life with 6000 genes. *Science* 274(5287):563-567.

Hansen, J. 2005. Graphics language SWOG, Bachelor thesis, University of Tartu.

Harris, M.; Clark, J.; Ireland, A.; Lomax, J.; Ashburner, M.; Foulger, R.;
Eilbeck, K.; Lewis, S.; Marshall, B.; Mungall, C.; Richter, J.; Rubin, G.;
Blake, J.; Bult, C.; Dolan, M.; Drabkin, H.; Eppig, J.; Hill, D.; Ni, L.; Ring-wald, M.; Balakrishnan, R.; Cherry, J.; Christie, K.; Costanzo, M.; Dwight,
S.; Engel, S.; Fisk, D.; Hirschman, J.; Hong, E.; Nash, R.; Sethuraman,
A.; Theesfeld, C.; Botstein, D.; Dolinski, K.; Feierbach, B.; Berardini, T.;
Mundodi, S.; Rhee, S.; Apweiler, R.; Barrell, D.; Camon, E.; Dimmer, E.;
Lee, V.; Chisholm, R.; Gaudet, P.; Kibbe, W.; Kishore, R.; Schwarz, E.;
Sternberg, P.; Gwinn, M.; Hannick, L.; Wortman, J.; Berriman, M.; Wood,
V.; de la Cruz, N.; Tonellato, P.; Jaiswal, P.; Seigfried, T.; and White, R.
2004. The Gene Ontology (GO) database and informatics resource. *Nucleic Acids Research* 32(Database issue):D258–261.

Ho, Y.; Gruhler, A.; Heilbut, A.; Bader, G.; Moore, L.; Adams, S.; Millar, A.; Taylor, P.; Bennet, K.; Boutilier, K.; et al. 2002. Systematic identification of protein complexes in *Saccharomyces cerevisiae* by mass spectometry.

Nature 415(6868):180–183.

Hong, E.; Balakrishnan, R.; Christie, K.; Costanzo, M.; Dwight, S.; Engel,
S.; Fisk, D.; Hirschman, J.; Livestone, M.; Nash, R.; Park, J.; Oughtred, R.;
Skrzypek, M.; Starr, B.; Theesfeld, C.; Andrada, R.; Binkley, G.; Dong, Q.;
Lane, C.; Hitz, B.; Miyasato, S.; Schroeder, M.; Sethuraman, A.; Weng, S.;
Dolinski, K.; Botstein, D.; and Cherry, J. Saccharomyces genome database.
http://www.yeastgenome.org/ 26.03.2006, year = 2006,.

Hove-Jensen, B. 2004. Heterooligomeric phosphoribosyl diphosphate synthase of *Saccharomyces cerevisiae*: combinatorial expression of the five prs genes in *Escherichia coli. J. Biol. Chem* 279(39):40345–40350.

Hughes, T.; Marton, M.; Jones, A.; Robets, C.; Stoughton, R.; Armour, C.; Bennett, H.; Coffey, E.; Dai, H.; He, Y.; et al. 2000. Functional discovery via a compendium of expression profiles. *Cell* 102(1):109–126.

Ito, T.; Chiba, T.; Ozawa, R.; Yoshida, M.; Hattori, M.; and Sakaki, Y. 2001. A comprehensive two-hybrid analysis to explore the yeast protein interactome. In *Proc. Natl. Acad. Sci. USA*, volume 98, 4569–4574.

Kellis, M.; Patterson, N.; Endrizzi, M.; Birren, B.; and Lander, E. S. 2003. Sequencing and comparison of yeast species to identify genes and regulatory elements. *Nature* 423(6937):241–254.

Kemmeren, P.; van Berkum, N.; Vilo, J.; Bijma, T.; Donders, R.; Brazma, A.; and Holstege, F. 2002. Protein interaction verification and functional annotation by integrated analysis of genome-scale data. *Molecular Cell* 9(5):1133–1143.

Krogan, N.; Cagney, G.; Yu, H.; Zhong, G.; Guo, X.; Ignatchenko, A.;Li, J.; Pu, S.; Datta, N.; Tikuisis, A.; Punna, T.; Peregrin-Alvarez, J.;Shales, M.; Zhang, X.; Davey, M.; Robinson, M.; Paccanaro, A.; Bray,

J.; Sheung, A.; Beattie, B.; Richards, D.; Canadien, V.; Lalev, A.; Mena,
F.; Wong, P.; Starostine, A.; Canete, M.; Vlasblom, J.; Wu, S.; Orsi, C.;
Collins, S.; Chandran, S.; Haw, R.; Rilstone, J.; Gandi, K.; Thompson, N.;
Musso, G.; Onge, P. S.; Ghanny, S.; Lam, M.; Butland, G.; Altaf-Ul, A.;
Kanaya, S.; Shilatifard, A.; O'Shea, E.; Weissman, J.; Ingles, C.; Hughes,
T.; Parkinson, J.; Gerstein, M.; Wodak, S.; Emili, A.; and Greenblatt, J.
2006. Global landscape of protein complexes in the yeast *Saccharomyces cerevisiae*. *Nature* 440(7084):637–643.

Lenhard, B.; Sandelin, A.; Mendoza, L.; Engström, P.; Jareborg, N.; and Wasserman, W. W. 2003. Identification of conserved regulatory elements by comparative genome analysis. *Journal of Biology* 2(13).

Li, W.; Luo, C.; and Wu, C. Evolution of DNA sequences. In *Molecular Evolutionary Genetics*.

Liti, G., and Louis, E. J. 2005. Yeast evolution and comparative genomics. Annual Review of Microbiology 59:135–153.

Peterson, H. 2004. Gene regulation database BiGeR, Bachelor thesis, University of Tartu.

Puig, O.; Caspary, F.; Rigaut, G.; Rutz, B.; Bouveret, E.; Bragado-Nilsson,
E.; Wilm, M.; and Séraphin, B. 2001. The tandem affinity purification (TAP) method: A general procedure of protein complex purification. *Methods* 24(3):218-229.

Qiu, P. 2003. Computational approaches for deciphering the transcriptional regulatory network by promoter analysis. *Biosilico* 1(4):125–133.

Reimand, J. 2006. Gene ontology mining tool GOSt. Master's thesis, University of Tartu. Roberts, C.; Nelson, B.; Marton, M.; Stoughton, R.; Meyer, M.; Benet, H.; He, Y.; Dai, H.; Walker, W.; Hughes, T.; et al. 2000. Signaling and circuitry of multiple MAPK pathways revealed by a matrix of global gene expression profiles. *Science* 287(5454):873–880.

Spellman, P.; Sherlock, G.; Zhang, M.; Iyer, V.; Anders, K.; Eisen, M.; Brown, P.; Botstein, D.; and Futcher, B. 1998. Comprehensive identification of cell cycle-regulated genes of the yeast *Saccharomyces cerevisiae* by microarray hybridization. *Molecular Biology of Cell* 9(12):3273–3297.

Stormo, G. D. 2000. DNA binding sites: representation and discovery. Bioinformatics 16(1):16-23.

Teixeira, M.; Monteiro, P.; Jain, P.; Tenreiro, S.; Fernandes, A.; Mira, N.; Alenquer, M.; Freitas, A.; A.L Oliveira, A.; and Sá-Correia, I. 2006. The YEASTRACT database: a tool for the analysis of transcription regulatory associations in *Saccharomyces cerevisiae*. *Nucleic Acids Research* 34(Database issue):D446-451.

Tompa, M. 2001. Identifying functional elements by comparative DNA sequence analysis. *Genome Research* 11(7):1143–1144.

Travers, K.; Patil, C.; Wodicka, L.; Lockhart, D.; Weissman, J.; and Walter,
P. 2000. Functional and genomic analyses reveal an essential coordination
between the unfolded protein response and ER-associated degradation. *Cell* 101(3):249-258.

Uetz, P.; giot, L.; Cagney, G.; Mansfield, T.; Judson, R.; Knight, J.; Lockshon, D.; Narayan, V.; Srinivasan, M.; and Pochart, P. 2000. A comprehensive analysis of protein-protein interactions in *Saccharomyces cerevisiae*. *Nature* 403(6770):623–627.

van Helden, J.; André, B.; and Collado-Vides, J. 1998. Extracting regula-

tory sites from the upstream region of yeast genes by computational analysis of oligonucleotide frequencies. *Journal of Molecular Biology* 281:827–842.

Vilo, J.; Brazma, A.; Jonassen, I.; Robinson, A.; and Ukkonen, E. 2000. Mining for putative regulatory elements in the yeast genome using gene expression data. In *ISMB-2000*, 384–394. AAAI press.

Vilo, J. 2002. Pattern Discovery from Biosequences. Ph.D. Dissertation, University of Helsinki.

Wingender, E.; Chen, X.; Hehl, R.; Karas, H.; Liebich, I.; Matys, V.; Meinhardt, T.; Prüß, M.; Reuter, I.; and Schacherer, F. 2000. TRANSFAC: an integrated system for gene expression regulation. *Nucleic Acids Research* 28(1):316–319.

Xie, Y., and Varshavsky, A. 2001. RPN4 is a ligand, substrate, and transcriptional regulator of the 26s proteasome: a negative feedback circuit. *Proc Natl Acad Sci USA* 98(6):3056–6.

Zhu, J., and Zhang, M. Q. 1999. SCPD: a promoter database of the yeast Saccharomyces cerevisiae. Bioinformatics 15(7/8):607-611.

Appendix A

Tables

ΤF	Nr	Known BS	S.cerevisiae	Prob	Orthologs	p-value	Cons.	Comments
ABF1	102	TCCCCATTAACG	CGCTGAAT	1.27101e-05	VN	VN	ON	weak motif, no
ACA1	2	NA	GCACCTTA	3.73022e-05	NA	NA	ON	weak motif, no
ACE2	s S	төстөөт	TGCTGGC	3.95396e-07	AACCAGCA	3.02497e-13	YES	medium strong motif, good overlap; strong (CA) repeats in S.cerevisiae upstreams
ADR1	43	GGRGK	TTGGGGTA	1.22318e-07	TTGGGGTA	3.88413e-15	YES	medium strong motif, improve- ment of a motif
AFT2	46	YRCACOCR	GGGTGCA	3.07926e-11	GGGTGCA	5.94077e-20	YES	strong motif, improvment of a motif
ARG80	19	CCTCTAAAGG	TTTCACTTA	1.80249e-07	TTTTGCCACCC	3.0292e-12	ON	medium strong motif, no overlap
ARG81	21	AAGTACAGTTAATAACGA	TTTCACTTA	3.11774e-07	TTTTGCCACCC	2.09836e-12	ON	medium strong motif, no overlap
ARO80	30	NA	NA	NA	NA	NA	NA	no motif found
ARR1	712	TTAATAA	ATCCGTACA	6.71412e-06	VN	NA	ON	weak motif, no overlap
ASH1	ы С	YTGAT	NA	NA	GTCTCCCACATCACCA	2.66303e-17	ON	no motif from S.cerevisiae, strong motif from orthologs with overlap

Table 1: TF target genes and their relative motifs

Continued on Next Page...

Page.
Next
uo
ontinued

	Comments	very strong over-	allpig motifs from	both sets	weak, not over-	lapping motifs	weak, not over-	lapping motifs	weak	S.cerevisiae	motifs but	medium motif	from orthologs	strong conserved	motif	strong conserved	motif	very weak S .	cerevisiae motif	not very strong	motifs, not over-	lapping	medium-strong	motif with weak	overall with	known motif	medium-strong	overlapping	motifs	very strong con-	served motif	very strong over-	lapping motif
	Cons.	\mathbf{YES}			NO		NO		NO					\mathbf{YES}		\mathbf{YES}		ON		NO			YES				YES			\mathbf{YES}		\mathbf{YES}	
	Prob	1.12426e-33			5.98394e-14		3.7337e-12		1.12228e-13					9.11558e-48		5.25653e-46		1.69306e-21		1.79077e-10			NA				1.16657e-15			2.10251e-31		4.18126e-66	
	Orthologs	CTTATC			AAAAGAAACGG		CCTGCCGTTA		TCGTATAAG					GTAACA		GTAACA		ATCACCACCA		TTATCATCAG			NA				CTTATC			TACTTCGAAGC		TGACTC	
	p-value	7.74645e-14			2.29767e-06		1.35422e-05		3.77335e-06					1.37214e-18		2.15905e-16		2.15306e-05		1.15037e-07			1.51717e-09				1.09203e-07			1.64527e-26		1.06791e-30	
	S.cerevisiae	CTTATC			GGCCCGC		TGACTCGA		CCCACGG					GTAACA		GTAACA		AGGCATCAT		AGTACAAG			GGAAATGTAA				GCTTATC			TACTTCGAAGC		TGACTCA	
ble 1 – Continued	Known BS	GATAA			GAAATTGCGTTT		GATAAG		TCGTATA					GTAAACAA		GTAAACA		TTTGCN97GCAAA		CGTATCGTATAAGGCAACAATAG			CGGN11CCG				AGATAAG			NA		TGASTCA	
Ta	Nr	26			45		15		4					26		136		5		22			62				19			59		291	
	TF	DAL80			DAL81		DAL82		ECM22					FKH1		FKH2		FLO8		FZF1			GAL4				GAT1			GAT3		GCN4	

•
60
ñ,
р.
+-
- S
Z
5
-
ě
2
÷Ξ
E I
8
0

	51							
\mathbf{TF}	Nr	Known BS	S.cerevisiae	p-value	Orthologs	Prob	Cons.	Comments
GCR1	128	GGCTTCCWC	ATATACGGT	1.64991e-18	ATACGGTGTT	6.27441e-19	YES	conserved strong
								motif, not over-
								lapping
GCR2	82	NA	CCGGGCCA	6.18481e-06	ATTGAGCCAATC	1.49211e-13	NO	weak non-
								conserved motifs
GIS1	24	TWAGGGAT	TAAATAGGA	3.00845e-06	TAAGGG	1.63583e-11	NO	weak non-
								conserved motif,
								orthologous motif
								overlapping with
								known motif
GLN3	56	GATAAGA	CTTATC	2.0965e-14	CTTATC	2.10289e-32	\mathbf{YES}	strong conseerved
								overlapping motif
GTS1	14	NA	GTTGATGC	3.13764e-06	NA	NA	NA	weak motif found
GZF3	9	GATAAG	GTAATCAG	4.56157e-08	TCGCTTATC	4.72237e-17	ON	the orthologous
								motif overlaps
								with the known
								motif
HAA1	10	VN	ATAACTAAT	1.65476e-05	TTAATTTTAATT	1.76094e-13	NO	weak non-
								conserved motif
HAP1	56	CGGNNNTANCGG	GCCCC	7.42444e-08	TTGGTTGGTGG	1.54924e-12	NO	non-conserved
								non-overlapping
								motifs
HAP2	127	CCAAT	CCAATCA	8.26735e-12	CCAATCA	1.8314e-34	YES	strong conserved
								overlapping motif
HAP3	127	CCAAT	CCAATCA	8.26735e-12	CCAATCA	1.8314e-34	\mathbf{YES}	strong conserved
								overlapping motif
HAP4	136	GNCCAATCA	CCAATCA	3.74938e-11	CCAATCA	3.3177e-33	YES	strong conserved
								overlapping motif
HAP5	126	CCAAT	CCAATCA	6.92853e-12	CCAATCA	1.22063e-34	YES	strong conserved
								overlapping motif
HCM1	244	WAAYAAACAAW	GTACGCCAAA	5.36295e-06	NA	NA	NA	weak non-
								overlapping
								motif

Table 1 - Continued

Continued on Next Page...

	Comments	weak non-	overlapping	motif	weak non-	overlapping	motif	no motifs found	very strong con-	served and over-	lapping motif	very strong	S.cerevisiae	motif	medium strong,	partly overlap-	ping motif	strong conserved	overlapping motif	strong conserved	overlapping motif	weak non-	conserved motif	weak non-	conserved motif	very strong	conserved non-	overlapping motif	strong, partly	overlapping motif	weak, non-	overlapping motif	no motifs found	no motifs found
	Cons.	ON			NO			NA	YES			NO			ON			YES		YES		ON		NO		YES			YES		NO		NA	NA
	Prob	2.09612e-12			4.07536e-10			NA	1.97021e-64			2.35358e-50			4.7683e-15			3.67712e-32		1.57806e-37		1.27713e-12		7.30219e-11		2.37607e-52			4.07014e-20		1.52828e-11		NA	NA
	Orthologs	GATGGTACAT			AATAATAAGAAG			NA	TTCTAGAA			TACTAAC			TTAGCCGC			TTCACATG		TTCACATG		TTTCTTCAAGA		TTGTTTACCATTT		TGACTCA			AAATTCCG		GGCAGTTTCC		NA	NA
	p-value	5.5389e-05			6.50105e-06			NA	2.85502e-27			2.41217e-27			1.34873e-06			9.42167e-11		1.32319e-13		7.18972e-06		8.87619e-06		8.85301e-18			2.67637e-08		2.85045e-06		NA	NA
	S.cerevisiae	GGGCTC			CTTGATCTG			NA	TTCTAGAA			TCCGTACA			TTTCCCTTTTC			TTCACATG		TTCACATG		AACGATAAG		AGAGCTATT		TGACTCA			AATTCCG		AACAAGACA		NA	NA
ble 1 – Continued	Known BS	NA			NA			NA	cTTCtaGAAgcTTCtaGAAg			NA			TTTTCHHCG			CACATGC		CATGTGAA		NN		NA		CCGGTACCGG			TCCRNYGGA		TTTGCTC		NA	CGGN9CGG
Tai	Nr	10			11			1	114			194			6			76		76		2		2		155			9		39		4	с г
	\mathbf{TF}	HMS1			HOT1			HPC2	HSF1			IFH1			IME1			INO2		INO4		IXR1		KAR4		LEU3			LYS14		MAC1		MAL13	MAL63

Continued on Next Page

	Comments	/ strong con-	red and over-	oing motif	ng conserved	rlapping motif	ng motifs	ng conserved	-overlapping	if	-uou gu	served over-	oing motifs	/ strong mo-		ing conserved	ifs	k, partly over-	oing motifs	k, partly over-	oing motifs	lium-strong,	tly overlap-	g motifs	ng, overlap-	5 S.cerevisiae	if, non-	served	k, non-	served motifs	/ strong	served, over-
	Cons. (YES very	serv	lapi	YES stro	IDAO	NO stro	YES stro	non	mot	NO stro	COII	lapi	NO very	tifs	YES stro	mot	NO wea	lapi	NO wea	lapi	NO med	par	pine	NO stro	ping	mot	con	NO wea	con	YES very	COD
	Prob	1.6667e-67			1.14065e-23		2.75285e-43	4.12471e-27			5.88414e-15			3.36324e-42		1.32824e-23		6.12217e-13		3.30121e-17		3.54364e-17			3.05372e-11				5.06786e-16		1.26991e-34	
	Orthologs	ACGCG			ATTAGGAA		CCACAGTT	CACGTGA			AACTGTGGC			CCACAGTT		TGTTGAGC		AGTTTTCCG		GTTTTCCG		TCCGGGGAAA			AGACTGACA				TCTCCCACATCACCA		CCCCCT	
	p-value	1.42434e-27			5.24175e-09		1.32878e-16	8.70367e-14			5.58093e-08			4.12474e-17		4.38491e-07		1.58549e-06		1.32309e-06		7.982e-07			1.16731e-08				1.20969e-05		1.60505e-19	
	S .cerevisiae	ACGCG			TTTCCTAA		CACGTGA	CACGTGA			CACGTG			CACGTG		TGTTGAGC		CCCAC		GGGGAAGG		CTAACCAG			CGTGCCTT				AACGTAT		9996G	
ble 1 – Continued	Known BS	ACGCGT			CCTAATTAGG		VN	AAACTGTGG			AAACTGTGG			TCACGTG		VN		ATTTTGTGGGG		ATTTTGTGGGG		AATTRTCCGGGG			HAGGYA				N		MAGGGGSGG	
Tab	'n	194			130		34	34			80			36		60		50		22		9			34				14		267	
	TF	MBP1			MCM1		MET28	MET31			MET32			MET4		MGA2		MIG1		MIG2		MIG3			MOT3				MSN1		MSN2	

Continued on Next Page

	5.							
\mathbf{TF}	Nr	Known BS	S.cerevisiae	p-value	Orthologs	Prob	Cons.	Comments
MSN4	264	AAGGGG	AGGGG	4.2888e-17	CCCCT	9.27985e-34	YES	very strong
								conserved, over-
								lapping motifs
MSS11	23	TTTGCN97GCAAA	NA	NA	NA	NA	NA	no motifs found
NDT80	17	GNCRCAAAW	GACACAAA	9.01781e-07	CACAAAA	1.9436e-15	YES	medium-strong
								conserved, over-
								lapping motif
NRG1	222	GGACCCT	CCGCGGA	1.61155e-08	NA	NA	NA	medium-strong
								non-overlapping
								motif
NRG2	144	NA	CGCTCGGA	2.17966e-06	NA	NA	NA	weak motif
OAF1	26	CGGN3TNRN8,12CCG	AACTCCG	1.43239e-05	TAACTCCGA	1.07697e-13	\mathbf{YES}	weak, conserved,
								partly overlap-
								ping motifs
OP11	23	TCGAAYC	CATGTGA	3.29699e-11	TTCACATG	4.89146e-21	YES	conserved non-
								overlapping
								motifs
OTU1	1	NA	NA	NA	NA	NA	NA	no motifs found
PDC2	2	NA	GCTCTTTCA	3.03788e-05	ATTTTGCAT	4.55004e-15	NO	weak non-
								conserved motifs

Table 1 – Continued