

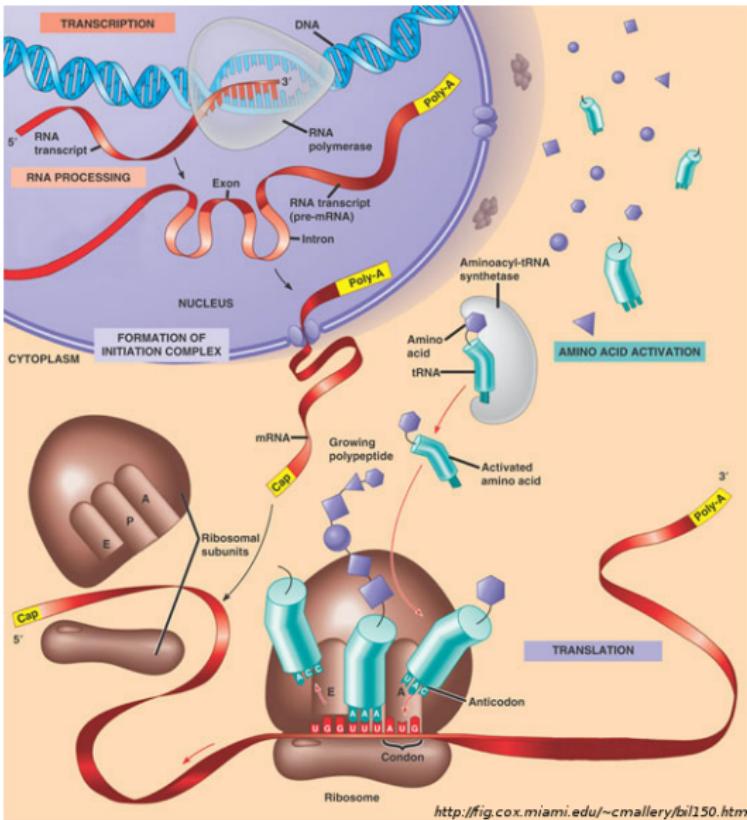
Discovery of regulatory motifs in *Saccharomyces cerevisiae* as a first step of understanding the gene regulation of baker's yeast

Hedi Peterson

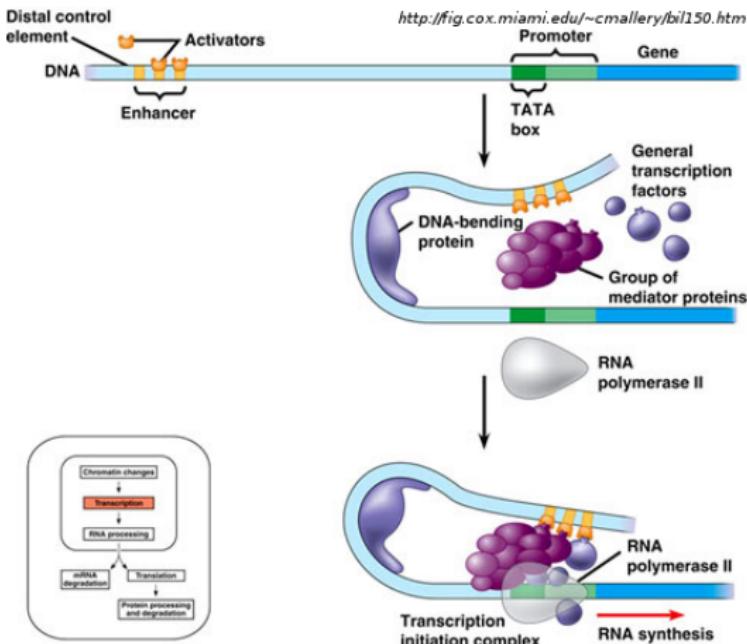
Institute of Molecular and Cell Biology
University of Tartu

29.09.2006
Voore

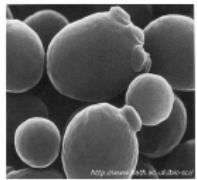
Biology: DNA <-> RNA -> Protein



Transcription (DNA -> RNA)



Background



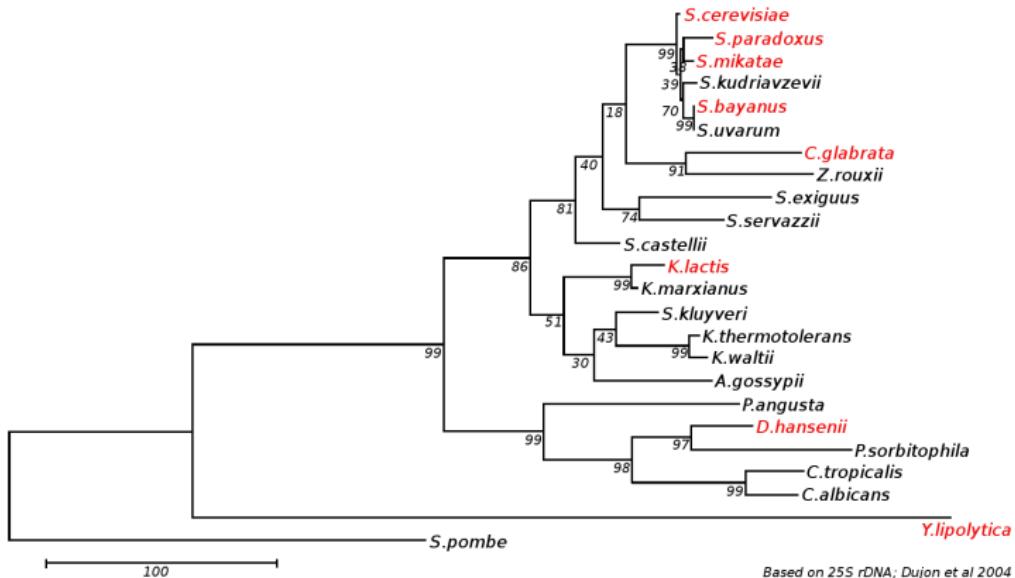
Saccharomyces cerevisiae

- 6604 predicted open reading frames
- 4403 described, 1377 undescribed and 824 dubious (SGD)
- Approximately one third of genes have no functionality or no regulator attached

- Predict regulatory motifs
- Using different gene groups and regulatory motifs, find similar genes

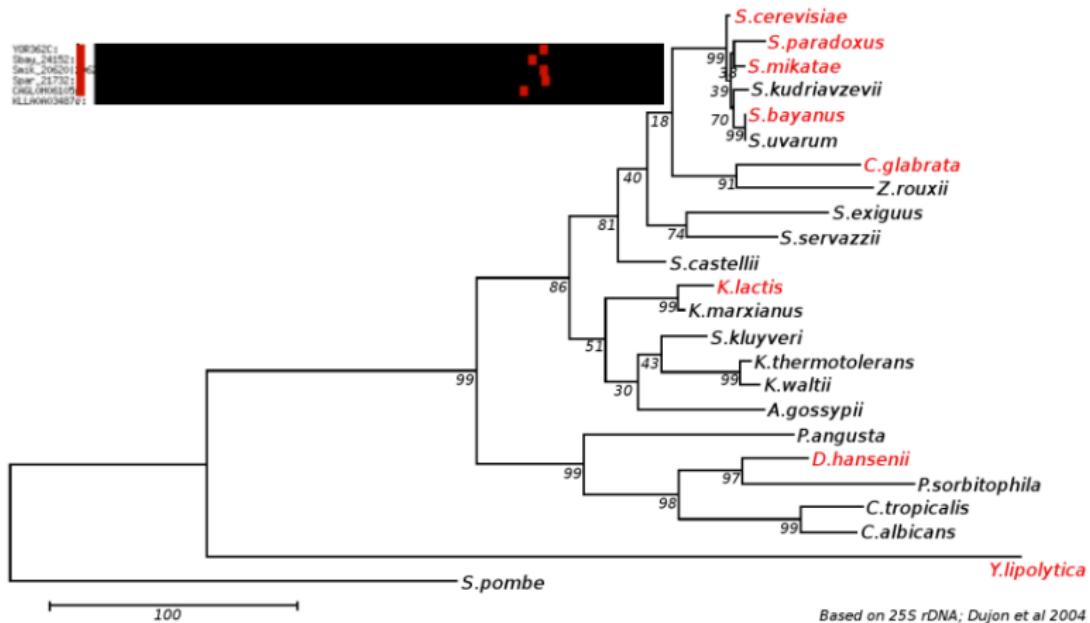
Species used

- *S.cerevisiae*, *S.bayanus*, *S.paradoxus*, *S.mikatae*,
C.glabrata, *K.lactis*, *D.hansenii*, *Y.lipolytica*



Based on 25S rDNA; Dujon et al 2004

Species used



Regulatory motifs

- Transcription factor binding sites
 - discrete motifs
 - alphabet {A,C,G,T}
 - length: 6-20 bp
 - location: 50-500 bp from ATG
- 600bp as promoter length

Used datasets

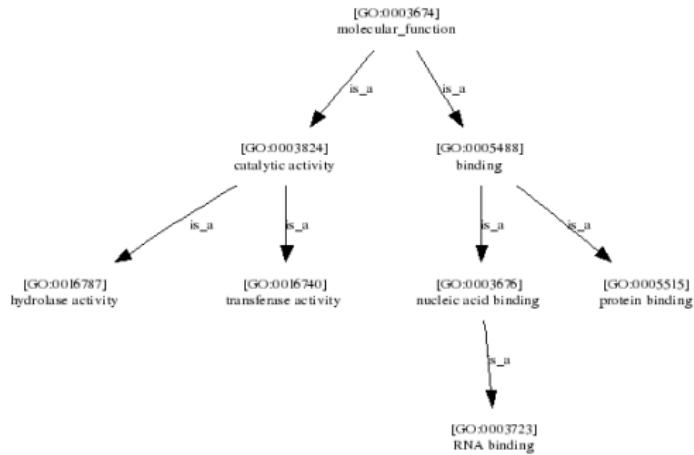
- Experimentally validated transcription factor binding sites
- Gene Ontology (GO)
- Protein-protein interactions (PPI)

TF target genes

- Experimentally validated
- 143 different transcription factors
- 4248 targets
- On average 84.8 targets per TF
- On average 2.9 TF per gene

Gene Ontology

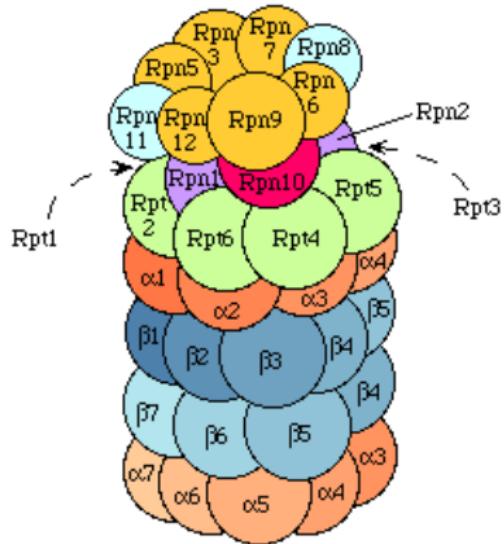
- All annotations are equal
- 3977 different GO groups
- size 1 to 658 genes in a GO group



Protein-protein interactions

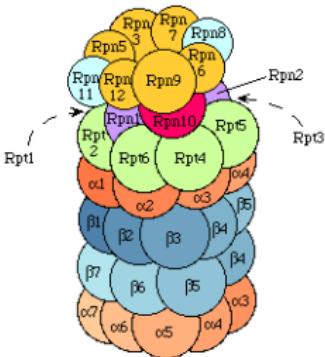
Three experimental datasets

- Kemmeren *et al* 1309
- Krogan *et al* 2186
- Gavin *et al* 1709
- In total 3334 different genes



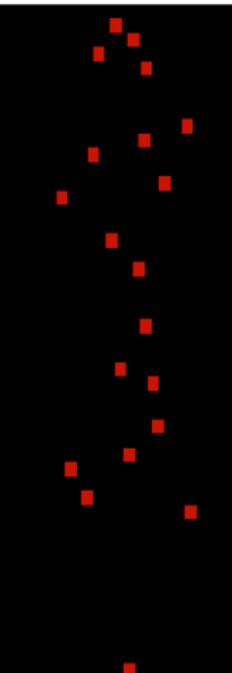
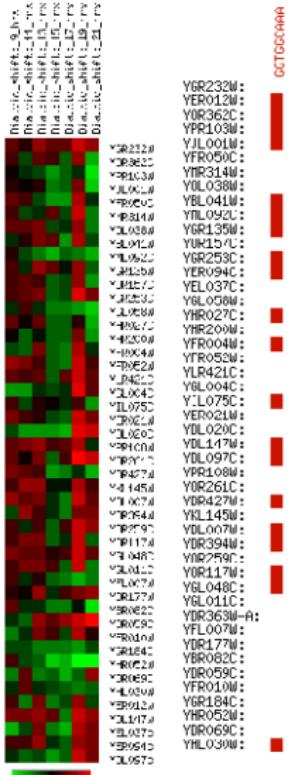
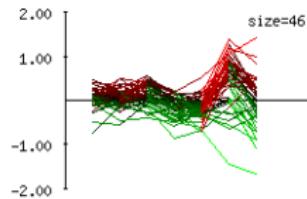
26S Proteasome (*Saccharomyces cerevisiae*)

Protein-protein interactions



26 S Proteasome (*Saccharomyces cerevisiae*)

KEGG



Steps to do

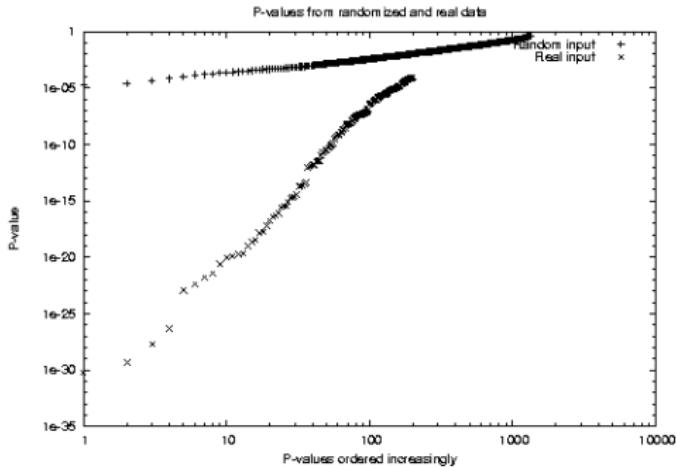
- Pattern discovery using different gene groups
- Expansion of gene groups based on PPI and motifs

Pattern discovery

- Using Sequence Pattern Exhaustive Search (SPEXS) algorithm (*Vilo*)
- Input: promoter sequences of somehow related genes (functionally or regulatory related)
- Background: all promoter sequences (*S.cerevisiae* or from all 8 yeasts)
- p-value threshold: 1.0e-04

Updating confidence threshold

- Input: randomly chosen promoter sequences
- Confidence threshold updated so, that expectance of false-positive motifs is smaller than 0.01%.
- New p-value threshold is between $9.1\text{e-}05$ and $9.3\text{e-}06$



Results - TF

- Known target gene sets
- For 138 out of 143 TF we found at least one motif
- 7381 distinct motifs

TF	Nr	Known TFBS	S.cer mot.	p-value	Ort. mot.	p-value
ACE2	83	TGCTGGT	TGCTGGC	3.95396e-07	AACCAGCA	3.02497e-13
ADR1	43	GGRGK	TTGGGGTA	1.22318e-07	TTGGGGTA	3.88413e-15
ASH1	5	YTGAT	NA	NA	GTCTCCCACATCACCA	2.66303e-17
CBF1	56	RTCACRTG	CACGTGA	4.49396e-39	CACGTG	4.44005e-83
CIN5	133	TTACRTAA	TTACATAA	4.20119e-12	TTACATAAT	3.55377e-14
CUP2	2	TCTTTTGCTG	TTCTTTGCT	2.61691e-10	AATTAGTAAGC	6.91919e-11
GAT3	59	NA	TACTTCGAAGC	1.64527e-26	TACTTCGAAGC	2.10251e-31
IFH1	194	NA	TCCGTACA	2.41217e-27	TACTAAC	2.35358e-50

Results - GO

- For 2006 GO groups from 3977 we found at least one motif

GO description	p-value	pattern	TF
retrotransposon nucleocapsid	2.57319e-136	TGTTGGAATA	Mot3
nucleolus	3.73244e-53	AAAATTTT	
peptidase activity	1.09457e-34	GCAAGGGATTGATAAT	
cytosolic ribosome (sensu Eukaryota)	1.89356e-31	CCGTACA	Rap1
amino acid biosynthesis	5.62027e-31	TGACTCA	Gcn4
proteasome complex (sensu Eukaryota)	5.39292e-28	GGTGGCAA	Rpn4
hydrolase activity, acting on ester bonds	1.75131e-25	ATAATGTAATA	Hcm1?
transferase activity, transferring phosphorus-containing groups	1.10842e-23	GATTGATAATG	
membrane-enclosed lumen	1.74176e-23	GCGATGAG	Esr1/Mec1?
DNA replication	4.16142e-21	ACGCGT	Mbp1
helicase activity	1.53663e-14	CCTCGACTAA	Xbp1
aspartate family amino acid catabolism	1.91911e-11	AGCACGTGAC	Pho4

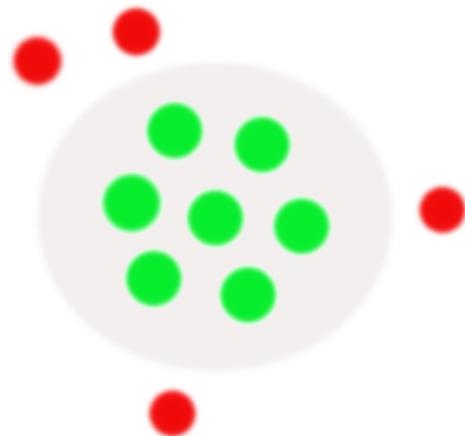
Expansion of gene groups

- For a input geneset (e.g GO:0008652) we found all interacting proteins
- We describe genes with the geneset specific motif (TGACTC)



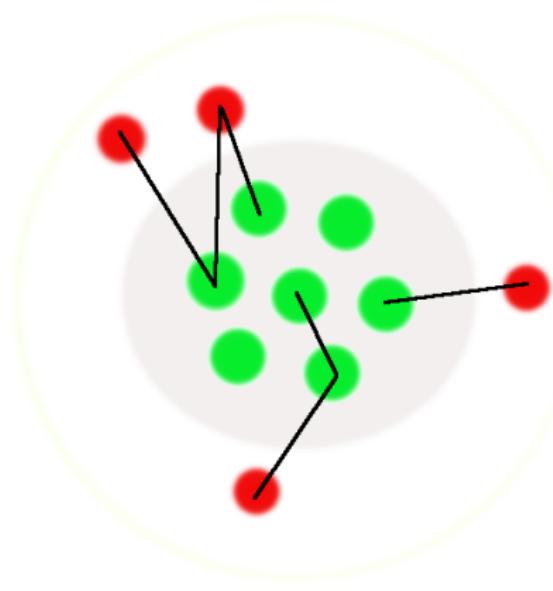
Expansion of gene groups

- For a input geneset (e.g GO:0008652) we found all interacting proteins
- We describe genes with the geneset specific motif (TGACTC)



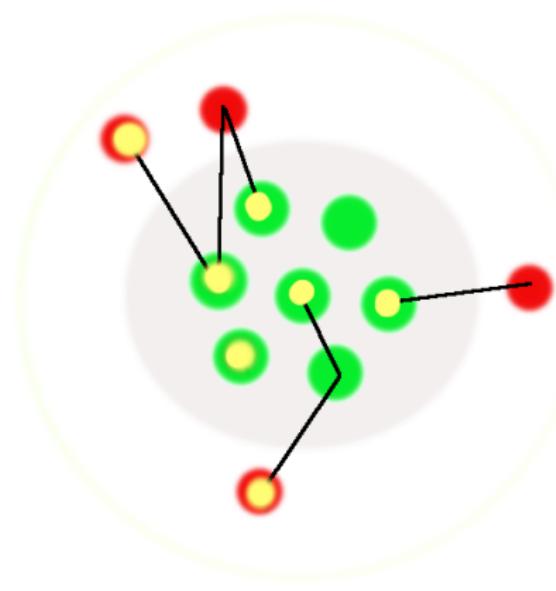
Expansion of gene groups

- For a input geneset (e.g GO:0008652) we found all interacting proteins
- We describe genes with the geneset specific motif (TGACTC)



Expansion of gene groups

- For a input geneset (e.g GO:0008652) we found all interacting proteins
- We describe genes with the geneset specific motif (TGACTC)





Conclusions

- We found regulatory motifs for different input genesets:
 - we improved previously known motifs
 - we found new motifs for transcription factors
 - we found Gene Ontology group specific motifs
- Using protein-protein interactions and regulatory motifs we can expand gene groups with probable candidate genes
Web-tool PPI-Gviz is publicly available
(<http://bioinf.ebc.ee/u/peterson/gviz/>)
- All discovered motifs are in BiGeR database
(<http://bioinf.ebc.ee/biger/>)

- Motif discovery using non-discrete motifs
- Filtering input genesets
 - GO geneset filtering based on evidence codes
 - Using more protein-protein interaction datasets + filter the interactions based on the gene expression similarity.
 - Use only very similar orthologs from more distant yeasts
- Comparison of newly found motifs with previously described motifs in BiGeR database
- Try to discover species specific motifs

Acknowledgements



Jaak Vilo, Jüri Reimand, Priit "Lemps" Adler, Jelena Zaitseva,
Meelis Kull, Hendrik Nigul, Raivo Kolde, Jaanus Hansen ..



ESF grant 5724, EU STRE ATD, EU NoE ENFIN and Kristjan
Jaak scholarship foundation

