

# **Transcription Factors Across Life: Control of Gene Expression**

Matt Weirauch

Center for Autoimmune Genomics and Etiology (CAGE)

Biomedical Informatics, Dev Bio

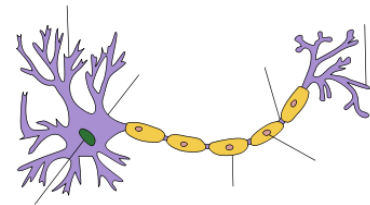
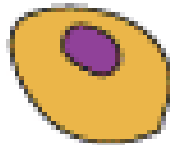
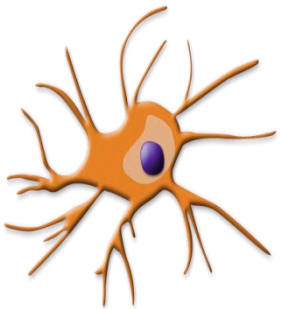
Cincinnati Children's Hospital

# Some perspective...

- You are composed of ~100 trillion cells
- Each cell has an identical copy of your genome
- Within your genome, there are genes
  - Stretches of DNA that encode blueprints for making proteins
  - Humans have ~20,000 of them

# Same genome/different cells

- Every cell has the same genomic DNA
- How can you explain the wide range of cellular form and function?



# Same genes/different forms

- Life is incredibly diverse
- Organisms with different forms and lifestyles use highly similar genes
  - Human/Chimp DNA is >99% identical
  - Human and Mouse have the same basic set of genes



# How is this so?

## It's the utilization of genes

- Genes aren't all used at all times
- Certain genes function only in certain
  - Developmental stages
  - Environments
  - Tissue types
  - Cellular states
- How does the cell “know” whether or not to use a certain gene? – Gene regulation

# Transcription Factors (TFs)

- Proteins that “turn on” and “turn off” the synthesis of genes in response to specific cues

# Transcription Factors (TFs)

- Proteins that “turn on” and “turn off” the synthesis of genes in response to specific cues
  - Humans have ~1700 TFs

# Transcription Factors (TFs)

- Proteins that “turn on” and “turn off” the synthesis of genes in response to specific cues
  - Humans have ~1700 TFs
  - Brain-specific TFs
  - Immune-responsive TFs
  - Developmental stage-specific TFs



# Transcription Factors (TFs)

- Proteins that “turn on” and “turn off” the synthesis of genes in response to specific cues
  - Humans have ~1700 TFs
  - Brain-specific TFs
  - Immune-responsive TFs
  - Developmental stage-specific TFs
- Interact with the genome by binding to short sequences

# Transcription Factors (TFs)

- Proteins that “turn on” and “turn off” the synthesis of genes in response to specific cues
  - Humans have ~1700 TFs
  - Brain-specific TFs
  - Immune-responsive TFs
  - Developmental stage-specific TFs
- Interact with the genome by binding to short sequences
- Each TF recognizes a specific DNA motif in the genome:

AAACCGTTA

TGACCA

GCGGTTA

cCGCAA

# TF regulatory elements

- Usually located near the genes they control
- Humans have ~20,000 genes, all of which are controlled by multiple regulatory elements (Dozens? Hundreds?)
- In order to understand how/when/why genes are turned on and off, we need to know what regulatory elements TFs can recognize (motifs)

# Quick summary

**Gene 1**



**Gene 2**



**Gene 3**



# Quick summary (cont'd)

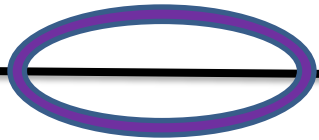
**Gene 1**



**Gene 2**



**Gene 3**



# Quick summary (cont'd)

**..AATGACTCATGATTGACCGTTAACCGGTGACGCATCTC..**

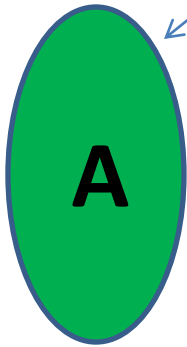
# Quick summary (cont'd)



**..AATGACTCATGATTGACCGTTAACCGGTGACGCATCTC..**

# Quick summary (cont'd)

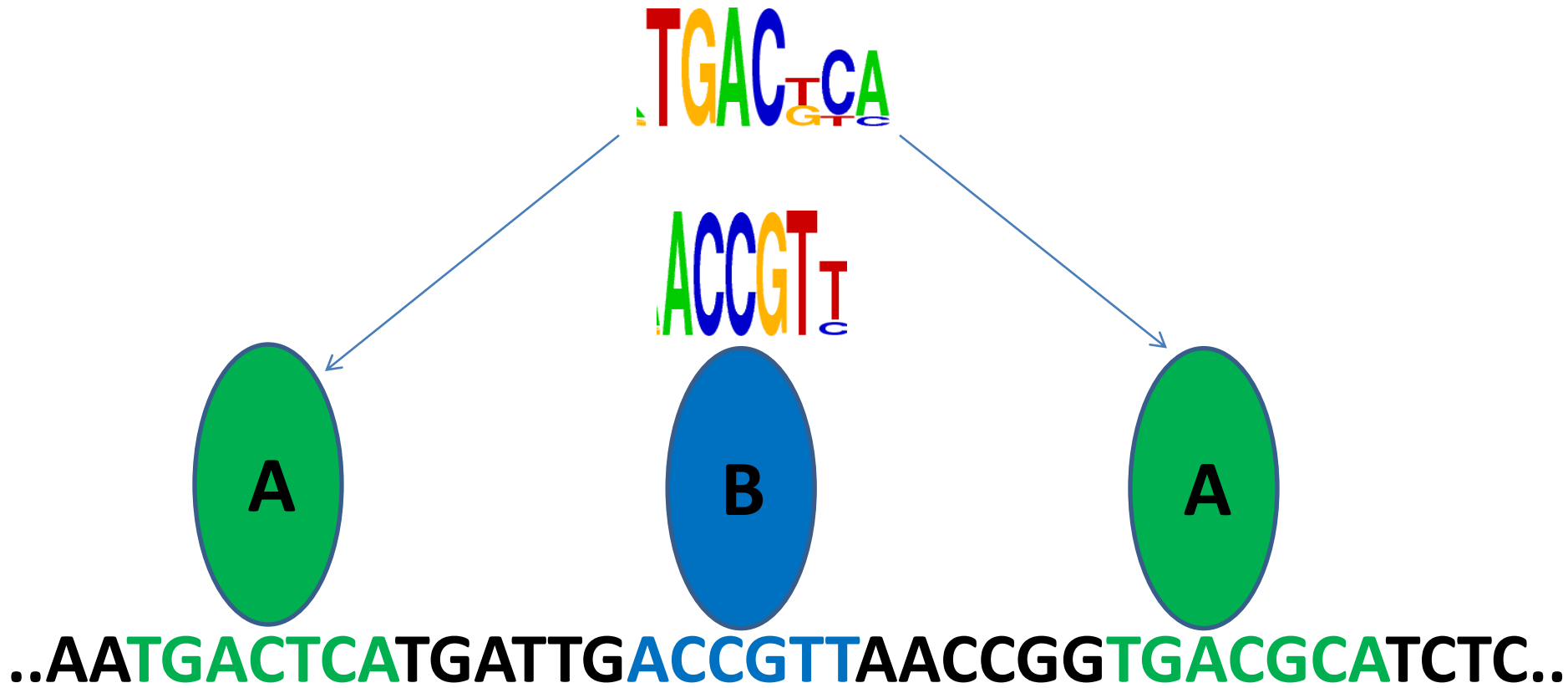
TGAC<sub>T</sub>CA<sub>GTC</sub>



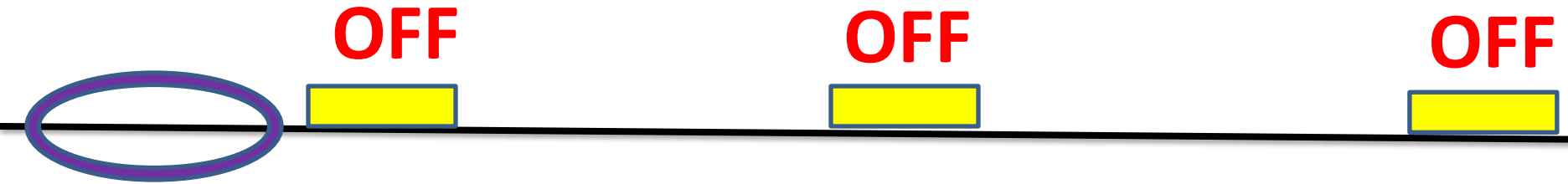
..AATGACTCATGATTGACCGTTAACCGGTGACGCATCTC..



# Quick summary (cont'd)



# Quick summary (cont'd)

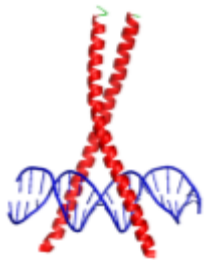


# Quick summary (cont'd)

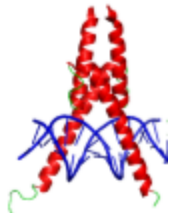


# Quick summary (cont'd)

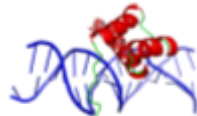




**bZIP (C/EBPalpha)**  
1NWQ [12578822]



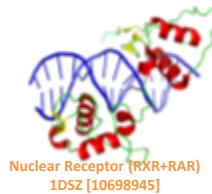
**bHLH (MyoD)**  
1MDY [8181063]



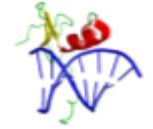
**Homeodomain (VND/NK-2)**  
1NK3 [9154919]



**C<sub>2</sub>H<sub>2</sub> Zinc Finger (Zif268)**  
1A1L [9562555]



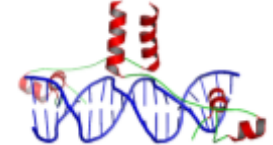
**Nuclear Receptor (RXR+RAR)**  
1DSZ [10698945]



**GATA (GATA-1)**  
1GAT [1567844]



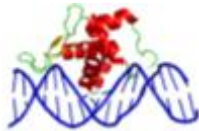
**THAP (THAP)**  
3KDE [20010837]



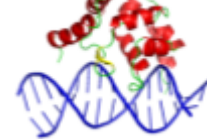
**Zinc cluster (GAL4)**  
1D66 [1557122]



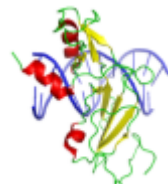
**Myb/SANT (c-Myb)**  
1MSF [7954830]



**Forkhead (Genesis)**  
2HDC [10369754]



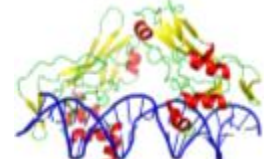
**ARID (Dead ringer)**  
1KQQ [11867548]



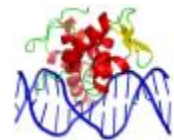
**GCM (GCMa)**  
1ODH [12682016]



**MAD5 box (SRF)**  
1SR5 [7637780]



**T-box (Bracyury/T)**  
1XBR [9349824]



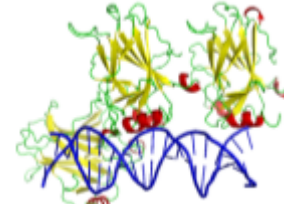
**E2F+DP (E2F4+DP2)**  
1CF7 [100907231]



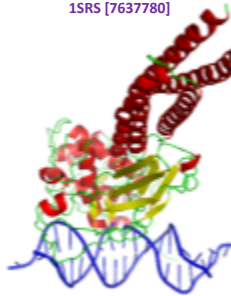
**IRF (IRF-1)**  
1IF1 [9422515]



**Ets (SAP-1)**  
1BC8 [9734357]



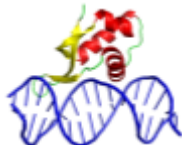
**p53 (p53)**  
1TSR [8023157]



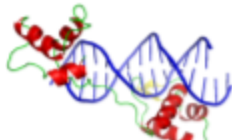
**STAT (STAT1)**  
1BF5 [9630226]



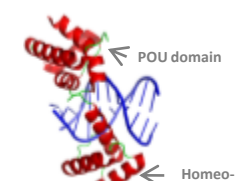
**SMAD (Smad3)**  
1O2J [12686552]



**RFX (RFX1)**  
1DP7 [10706293]



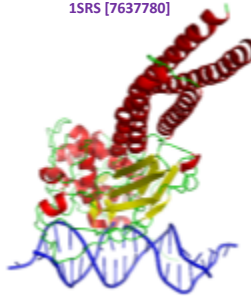
**Paired Box (prd)**  
1PDN [7867071]



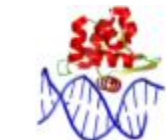
**POU (Oct-1)**  
1E3O [11583619]



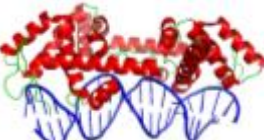
**RHD (NFkB)**  
1VKX [9450761]



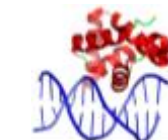
**CSL (C)**  
1TTU [152



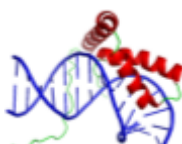
**IBD (IBP39)**  
1PP7 [14622596]



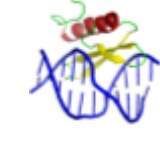
**LFY (LFY)**  
2VY1 [18784751]



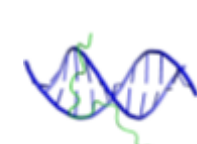
**CUT (SATB1)**  
2O49 [17652321]



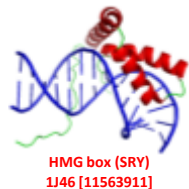
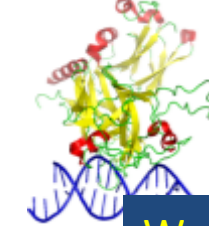
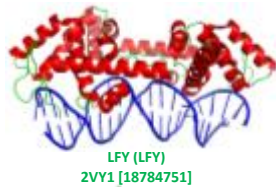
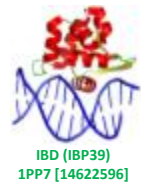
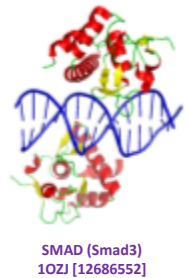
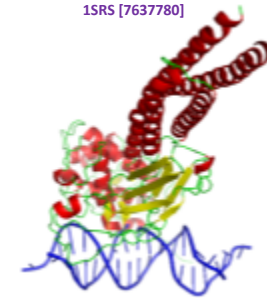
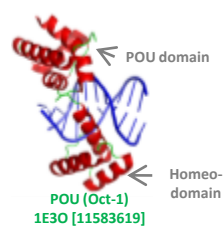
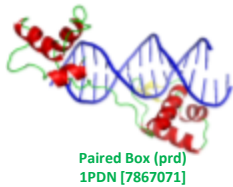
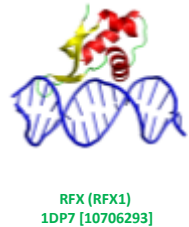
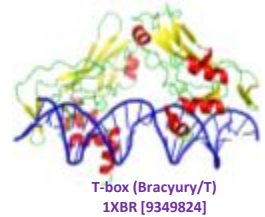
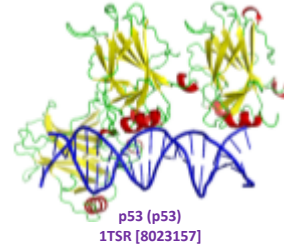
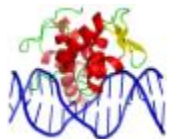
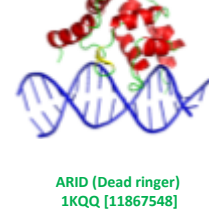
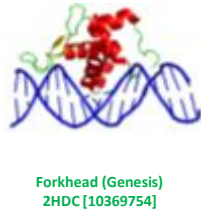
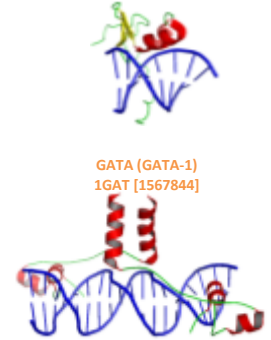
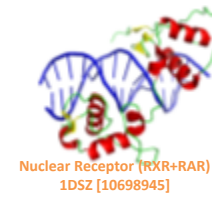
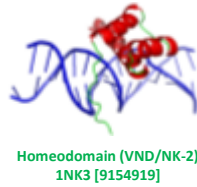
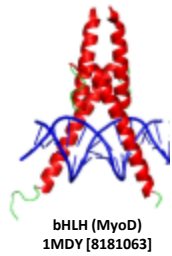
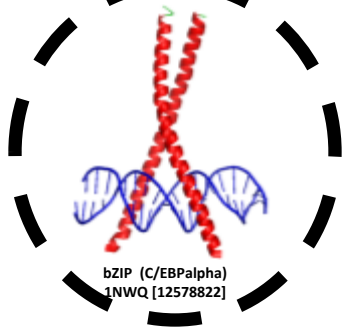
**HMG box (SRY)**  
1J46 [11563911]

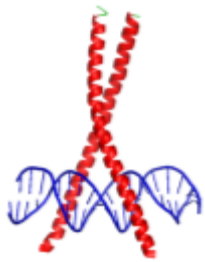


**AP2 (ATERF1)**  
1GCC [9736626]

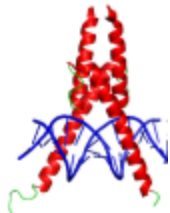


**AT hook (HMG-(Y))**  
2E2D [9253416]

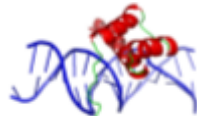




**bZIP (C/EBPalpha)**  
1NWQ [12578822]



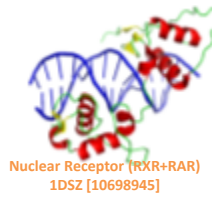
**bHLH (MyoD)**  
1MDY [8181063]



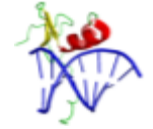
**Homeodomain (VND/NK-2)**  
1NK3 [9154919]



**C<sub>2</sub>H<sub>2</sub> Zinc Finger (Zif268)**  
1A1L [9562555]



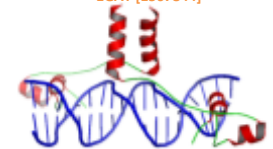
**Nuclear Receptor (RXR+RAR)**  
1DSZ [10698945]



**GATA (GATA-1)**  
1GAT [1567844]



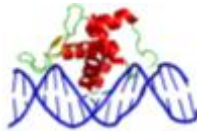
**THAP (THAP)**  
3KDE [20010837]



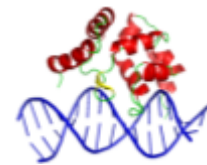
**Zinc cluster (GAL4)**  
1D66 [1557122]



**Myb/SANT (c-Myb)**  
1MSF [7954830]



**Forkhead (Genesis)**  
2HDC [10369754]



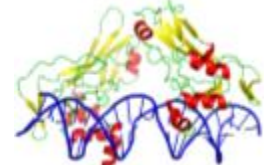
**ARID (Dead ringer)**  
1KQQ [11867548]



**GCM (GCMa)**  
1ODH [12682016]



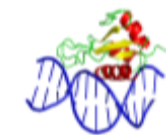
**MAD5 box (SRF)**  
1SR5 [7637780]



**T-box (Bracyury/T)**  
1XBR [9349824]



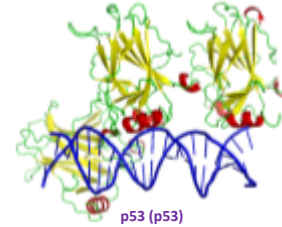
**E2F+DP (E2F4+DP2)**  
1CF7 [100907231]



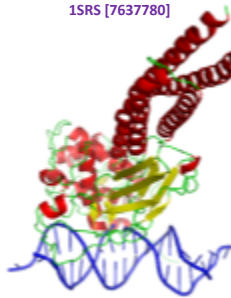
**IRF (IRF-1)**  
1IF1 [9422515]



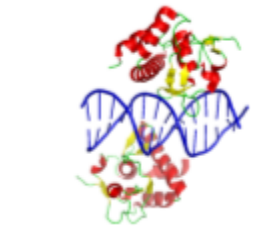
**Ets (SAP-1)**  
1BC8 [9734357]



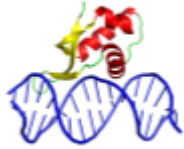
**p53 (p53)**  
1TSR [8023157]



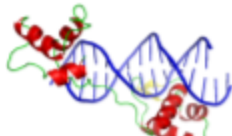
**STAT (STAT1)**  
1BF5 [9630226]



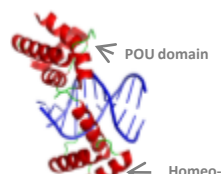
**SMAD (Smad3)**  
1OZJ [12686552]



**RFX (RFX1)**  
1DP7 [10706293]



**Paired Box (prd)**  
1PDN [7867071]



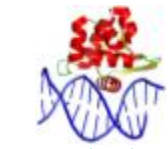
**POU (Oct-1)**  
1E3O [11583619]



**RHD (NFkB)**  
1VKX [9450761]



**CSL (C)**  
1TTU [152



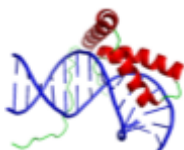
**IBD (IBP39)**  
1PP7 [14622596]



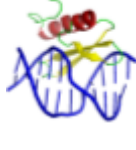
**LFY (LFY)**  
2VY1 [18784751]



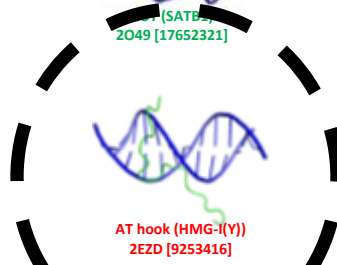
**AT hook (HMG-(Y))**  
2E2D [9253416]



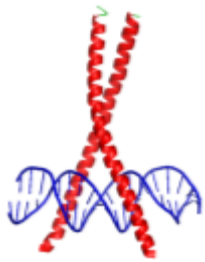
**HMG box (SRY)**  
1J46 [11563911]



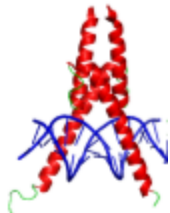
**AP2 (ATERF1)**  
1GCC [9736626]



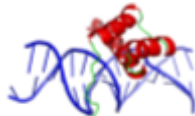
**AT hook (HMG-(Y))**  
2E2D [9253416]



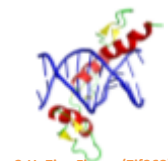
**bZIP (C/EBPalpha)**  
1NWQ [12578822]



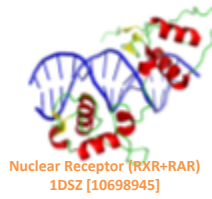
**bHLH (MyoD)**  
1MDY [8181063]



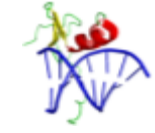
**Homeodomain (VND/NK-2)**  
1NK3 [9154919]



**C<sub>2</sub>H<sub>2</sub> Zinc Finger (Zif268)**  
1A1L [9562555]



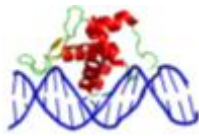
**Nuclear Receptor (RXR+RAR)**  
1DSZ [10698945]



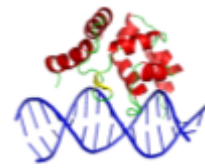
**GATA (GATA-1)**  
1GAT [1567844]



**Myb/SANT (c-Myb)**  
1MSF [7954830]



**Forkhead (Genesis)**  
2HDC [10369754]



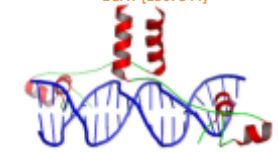
**ARID (Dead ringer)**  
1KQQ [11867548]



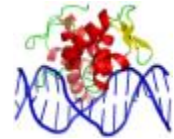
**GCM (GCMa)**  
1ODH [12682016]



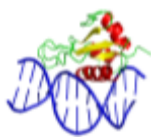
**THAP (THAP)**  
3KDE [20010837]



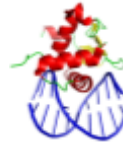
**Zinc cluster (GAL4)**  
1D66 [1557122]



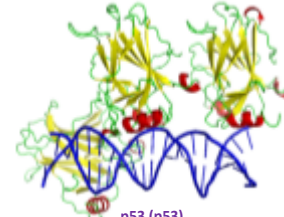
**E2F+DP (E2F4+DP2)**  
1CF7 [100907231]



**IRF (IRF-1)**  
1IF1 [9422515]



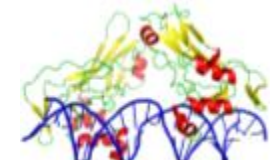
**Ets (SAP-1)**  
1BC8 [9734357]



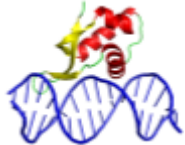
**p53 (p53)**  
1TSR [8023157]



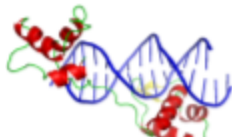
**MADS box (SRF)**  
1SR5 [7637780]



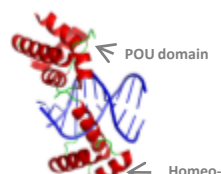
**T-box (Bracyury/T)**  
1XBR [9349824]



**RFX (RFX1)**  
1DP7 [10706293]



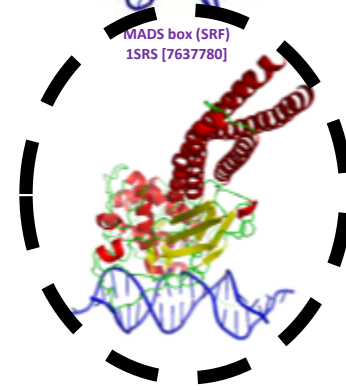
**Paired Box (prd)**  
1PDN [7867071]



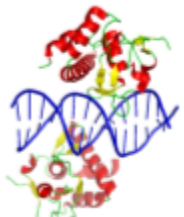
**POU (Oct-1)**  
1E3O [11583619]



**RHD (NFkB)**  
1VKX [9450761]



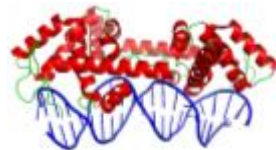
**STAT (STAT1)**  
1BF5 [9630226]



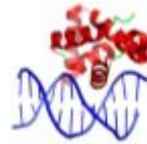
**SMAD (Smad3)**  
1O2J [12686552]



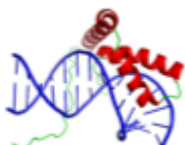
**IBD (IBP39)**  
1PP7 [14622596]



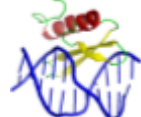
**LFY (LFY)**  
2VY1 [18784751]



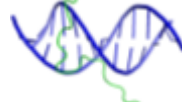
**CUT (SATB1)**  
2O49 [17652321]



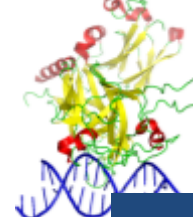
**HMG box (SRY)**  
1J46 [11563911]



**AP2 (ATERF1)**  
1GCC [9736626]



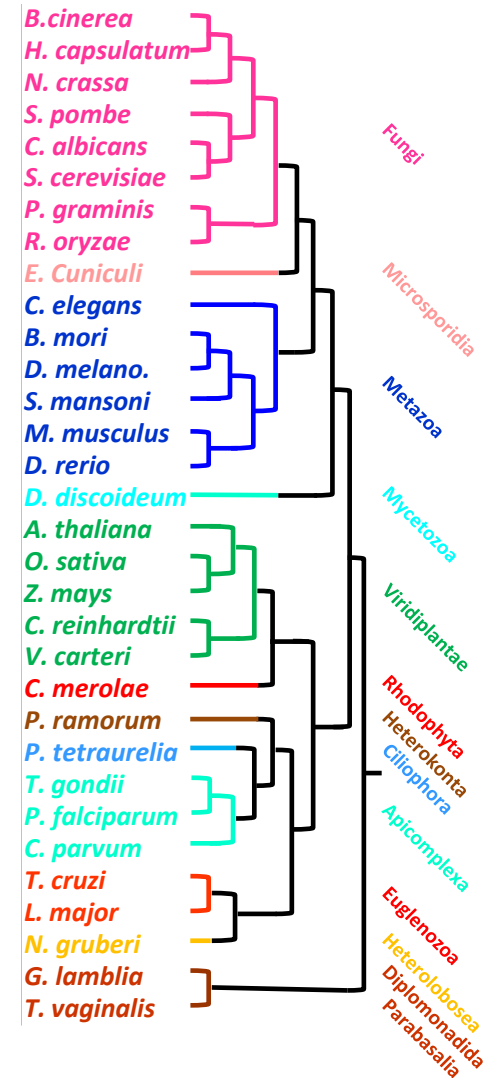
**AT hook (HMG-(Y))**  
2E2D [9253416]



**CSL (C)**  
1TTU [152]



# TF families across Eukaryota



# TF families across Eukaryota

#TFs

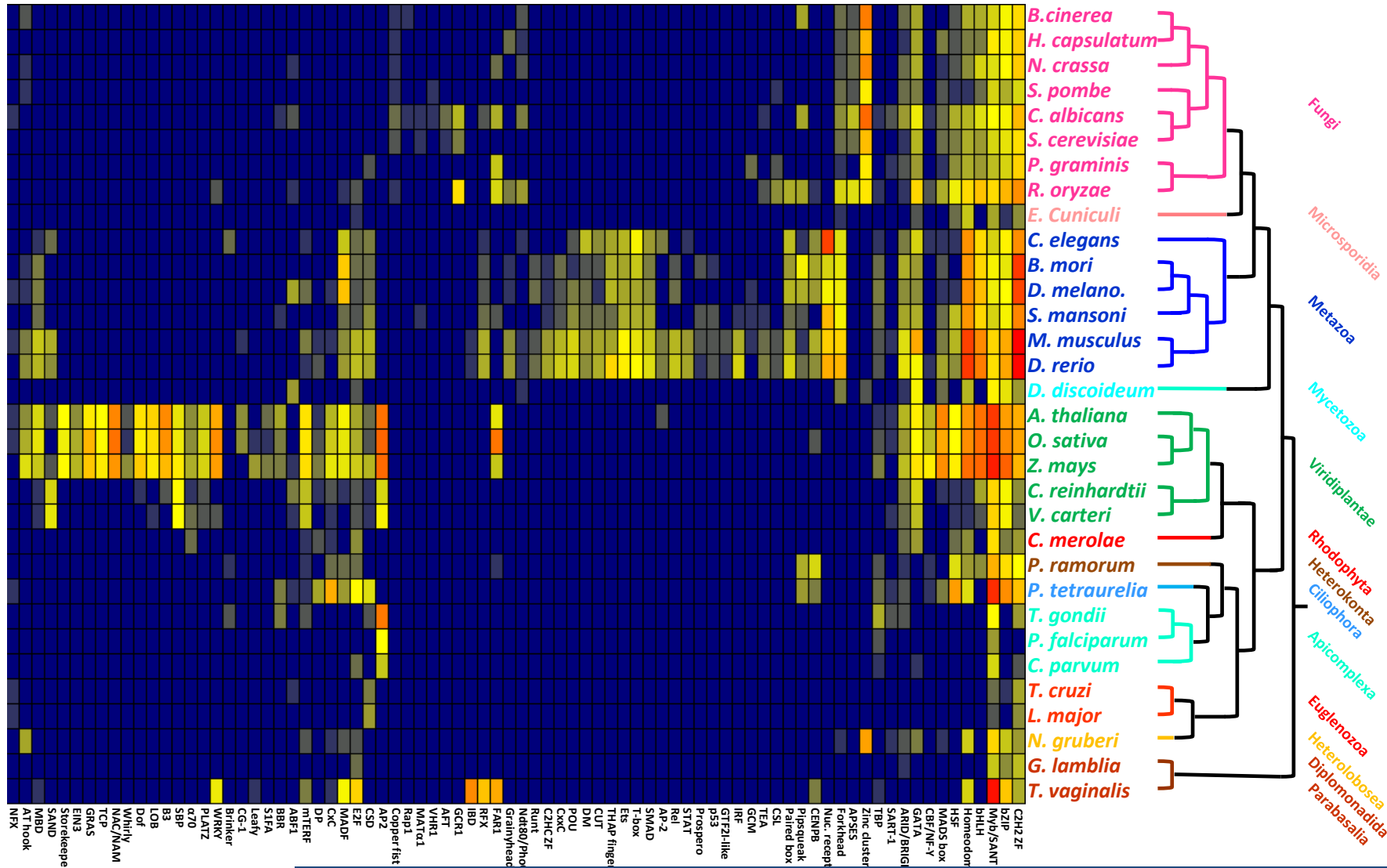
600

100

50

10

0



# TF families across Eukaryota

#TFs

600

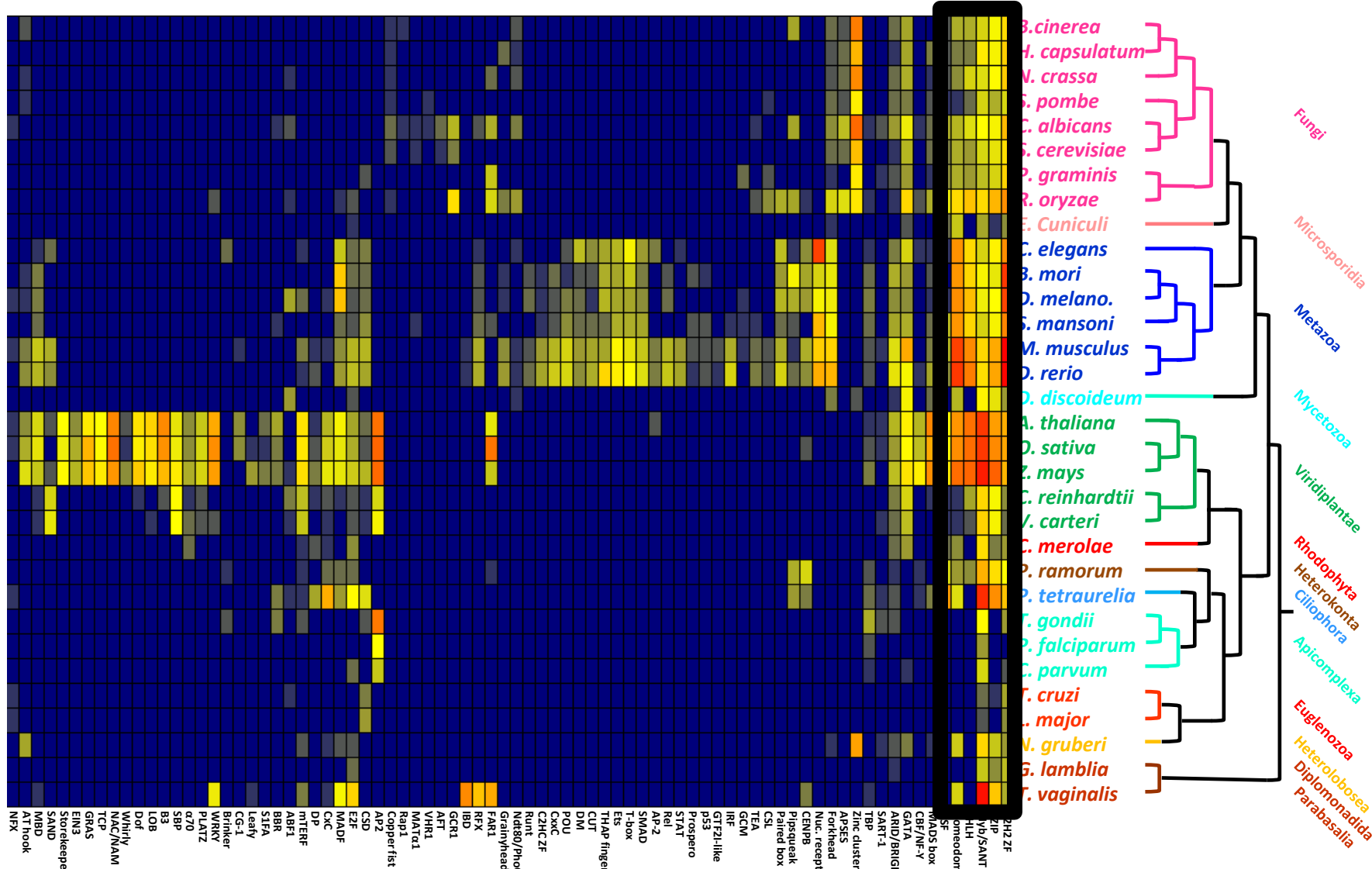
100

50

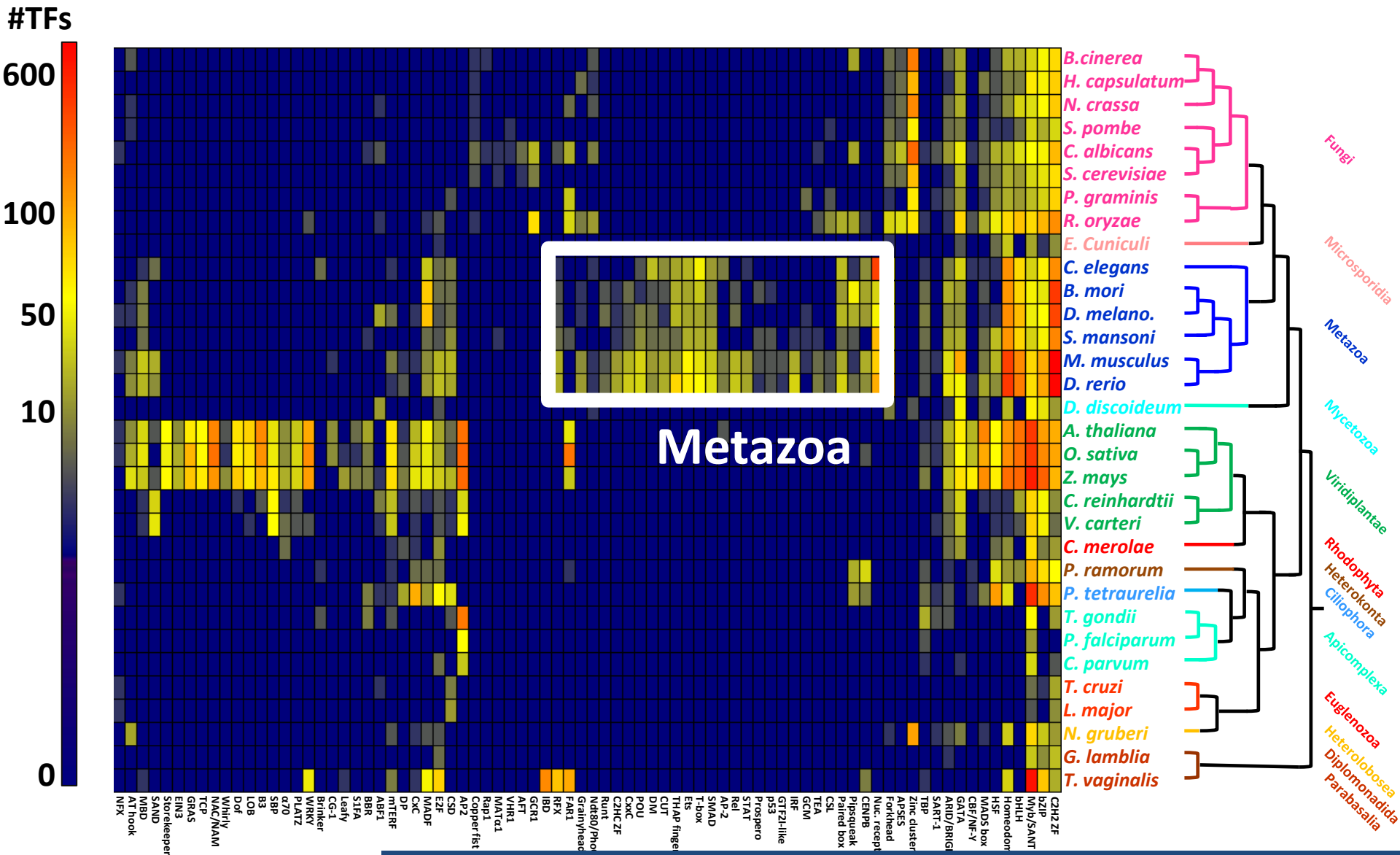
10

0

Eukaryotes



# TF families across Eukaryota



# TF families across Eukaryota

#TFs

600

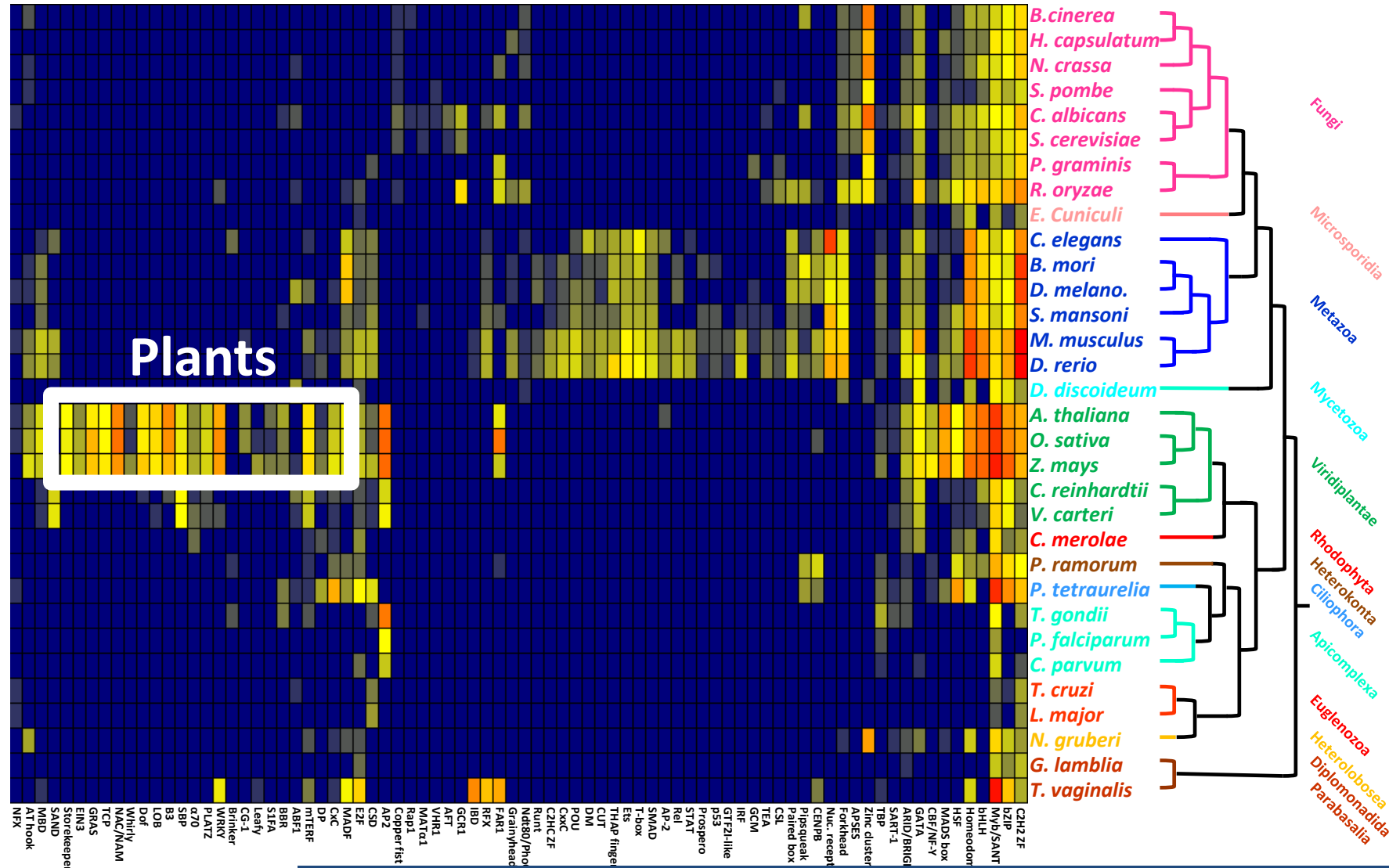
100

50

10

0

Plants



# TF families across Eukaryota

#TFs

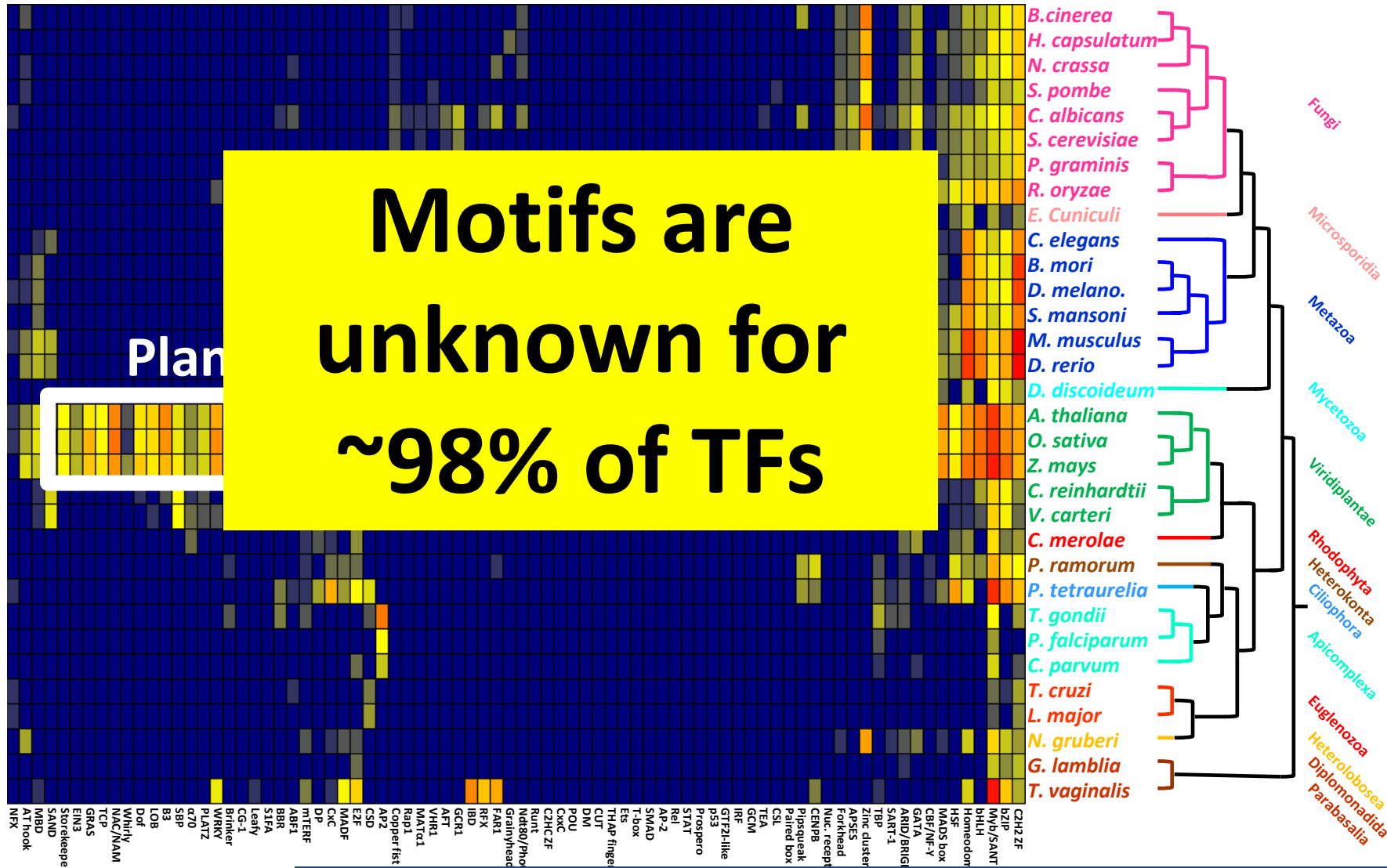
600

100

50

10

0



# Protein-binding microarrays (PBMs)

Make plasmid  
(DNA-binding  
domain)



Bind to DNA  
microarray



Determine  
binding  
preferences

~40,000 probes (35 unique bases each)

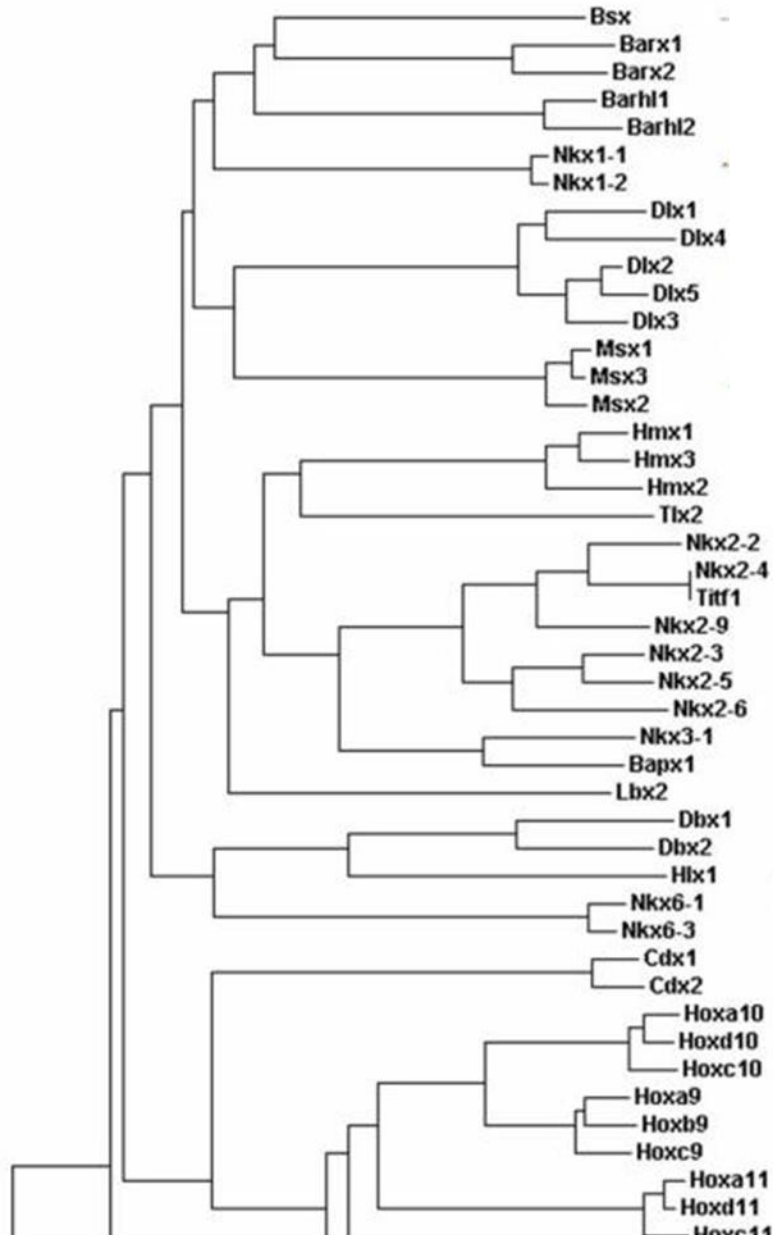
Probes contain all possible 10mers, 32 copies  
of all possible 8mers

Reproducible (8mers  $r \sim 0.70$  across arrays)

Scores track well with binding strength (kD)

*Berger et al. Nature Biotech 2006*

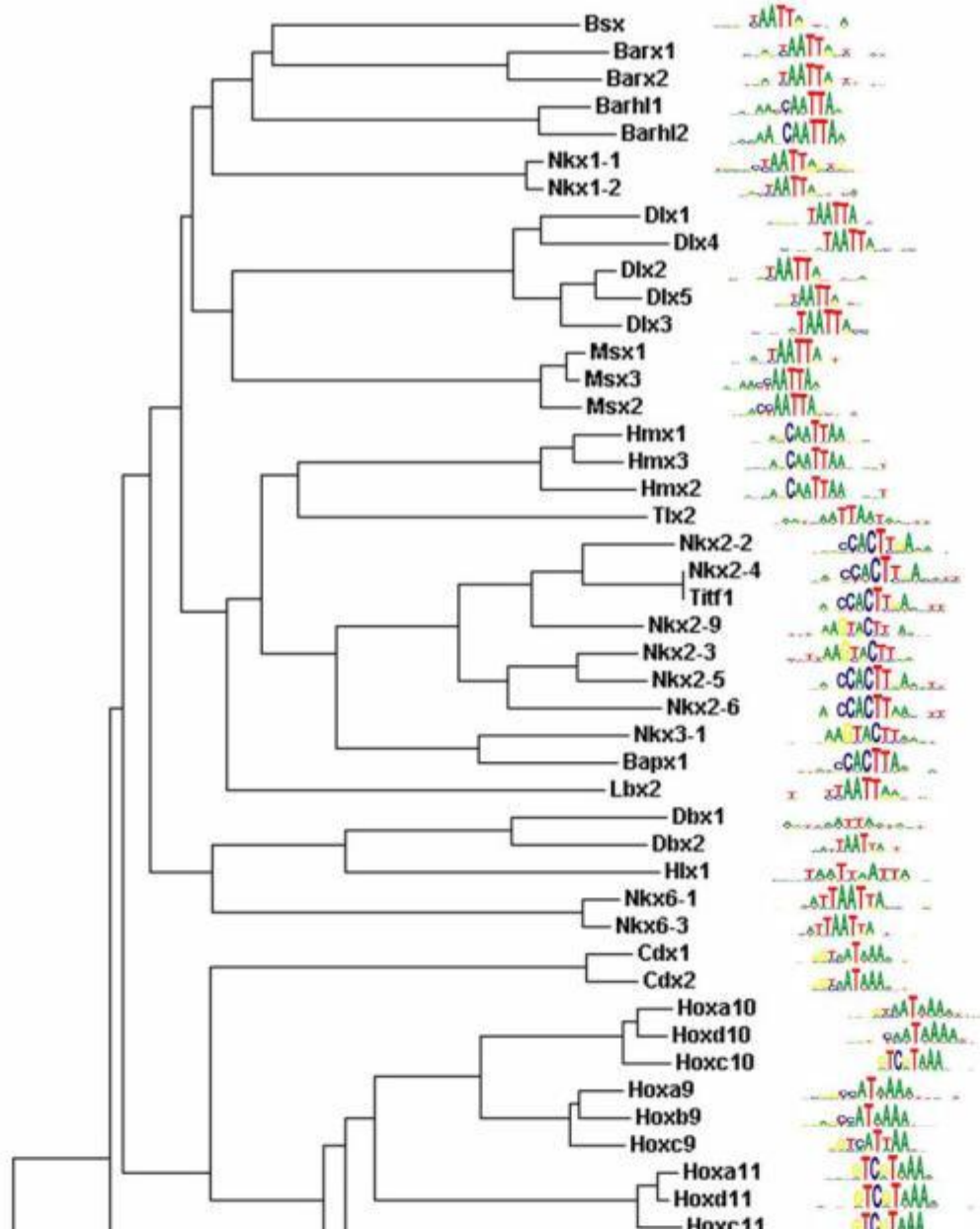
# Similar DBDs recognize similar motifs



**160 Mouse  
Homeodomains  
(Berger *et al.*  
2008, Cell)**



# Similar DBDs recognize similar motifs



# The “65% rule”

- 65% or greater AA DBD sequence identity -> nearly identical motif
  - Independent of evolutionary distance
  - Can use a mouse TF to predict yeast homolog’s binding preference
- Simple rule performs comparably to advanced methods
- Works for RNA binding proteins too!
  - Ray, Kazan, Cook, Weirauch, *et al.* , Nature 2013

# TF “picking” strategy

Choose ~2000 TFs for characterization, in order to:

# TF “picking” strategy

Choose ~2000 TFs for characterization, in order to:

1. Determine motif inference thresholds for the >80 TF families

# TF “picking” strategy

Choose ~2000 TFs for characterization, in order to:

1. Determine motif inference thresholds for the >80 TF families
2. Maximize the total number of motif inferences across all organisms

# TF “picking” strategy

Choose ~2000 TFs for characterization, in order to:

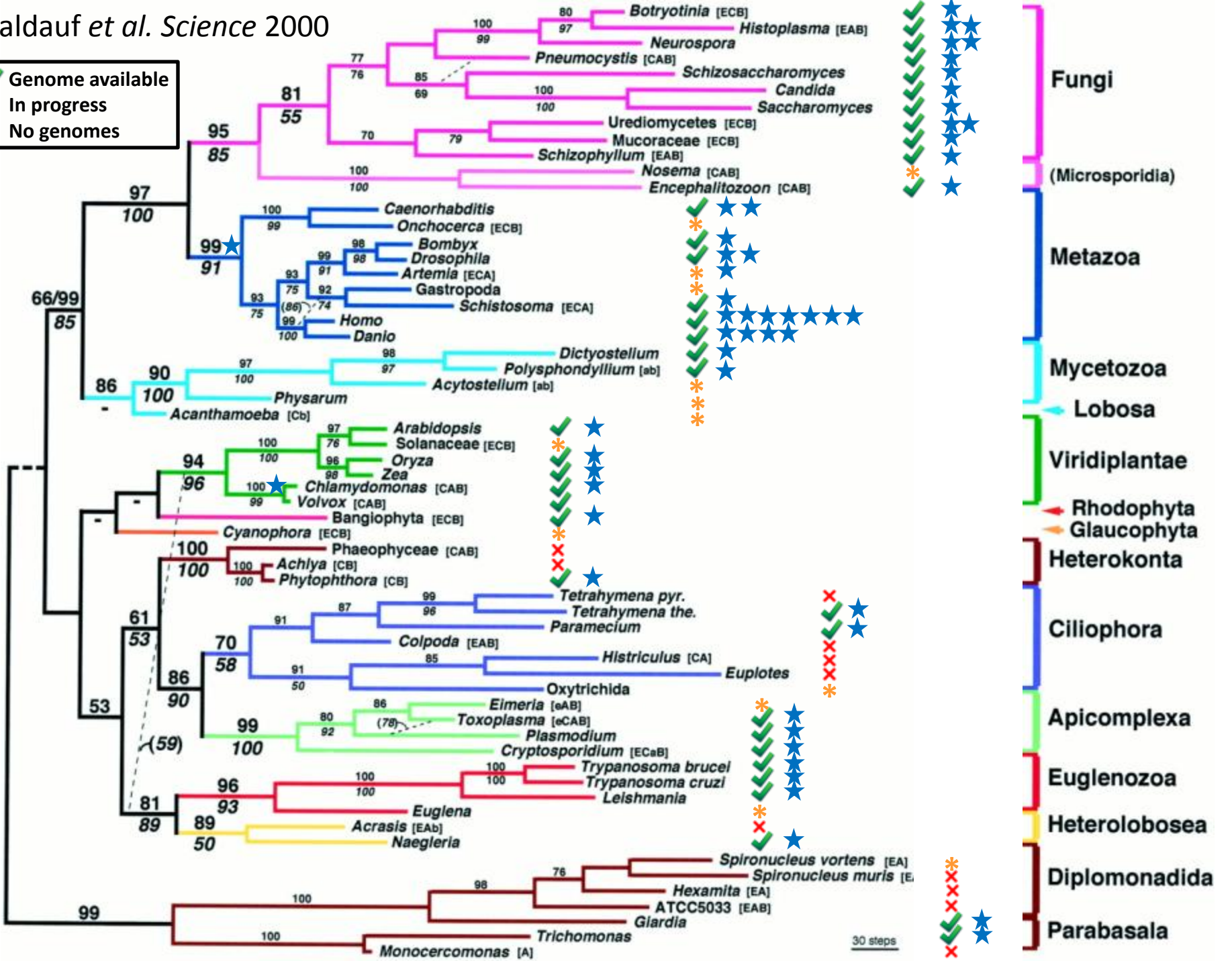
1. Determine motif inference thresholds for the >80 TF families
2. Maximize the total number of motif inferences across all organisms
3. Ensure all major clades are sampled

# TF “picking” strategy

Choose ~2000 TFs for characterization, in order to:

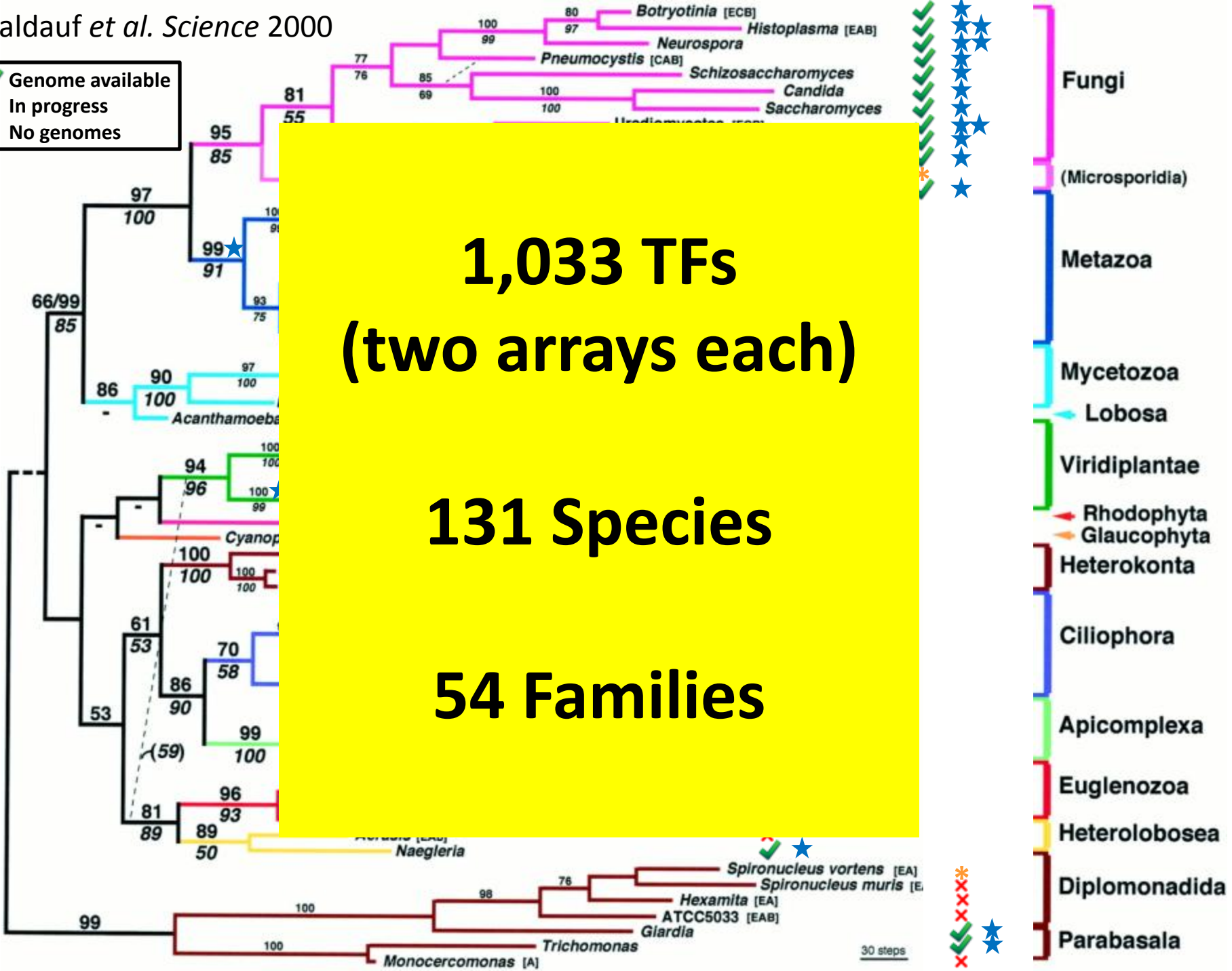
1. Determine motif inference thresholds for the >80 TF families
2. Maximize the total number of motif inferences across all organisms
3. Ensure all major clades are sampled
4. Characterize motifs for select, diverse model organisms
  - *Arabidopsis thaliana* (plant)
  - *Dictyostelium discoideum* (slime mold)
  - *Ostreococcus tauri* (algae)
  - *Neurospora crassa* (fungus)
  - *Mus musculus* (animal)

✓ Genome available  
✱ In progress  
✗ No genomes





- ✓ Genome available
- \* In progress
- ✗ No genomes

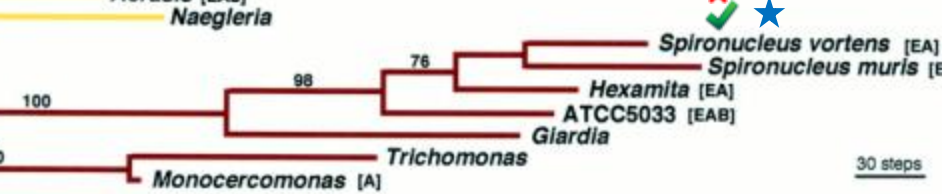
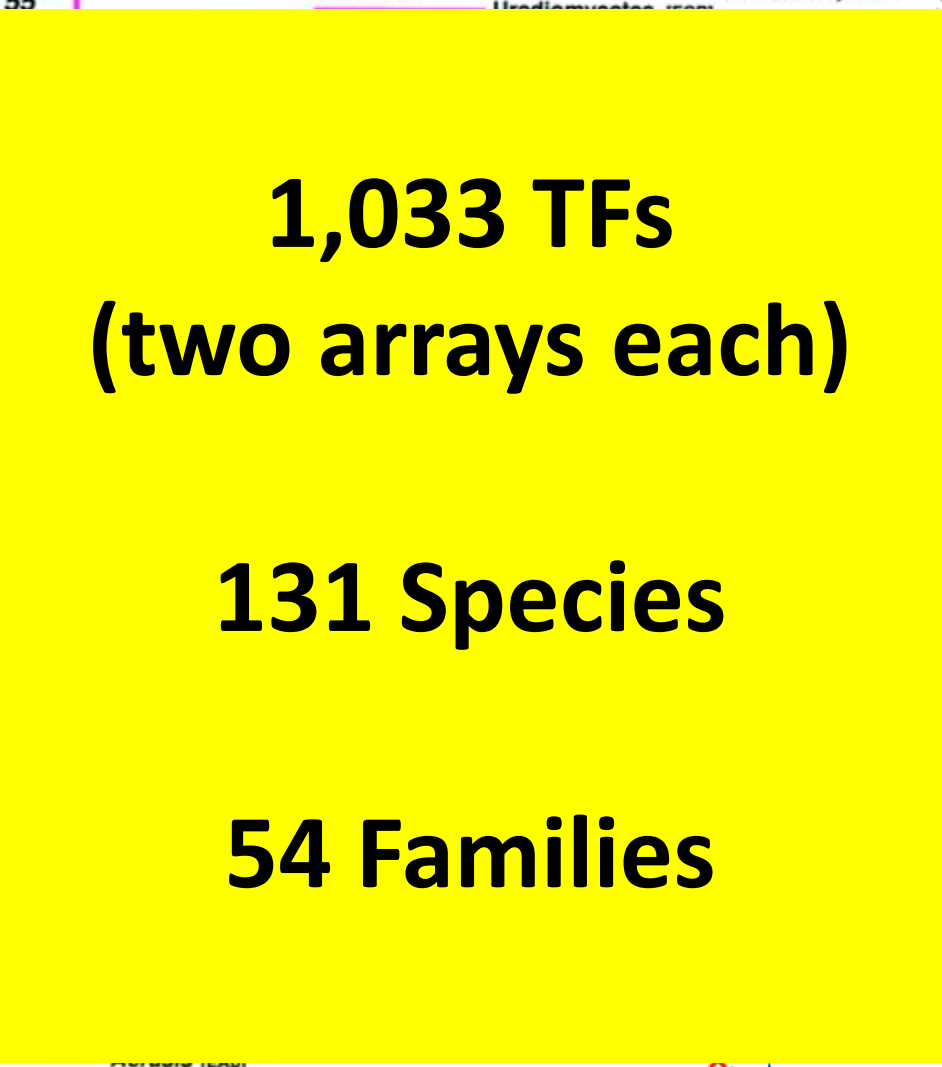
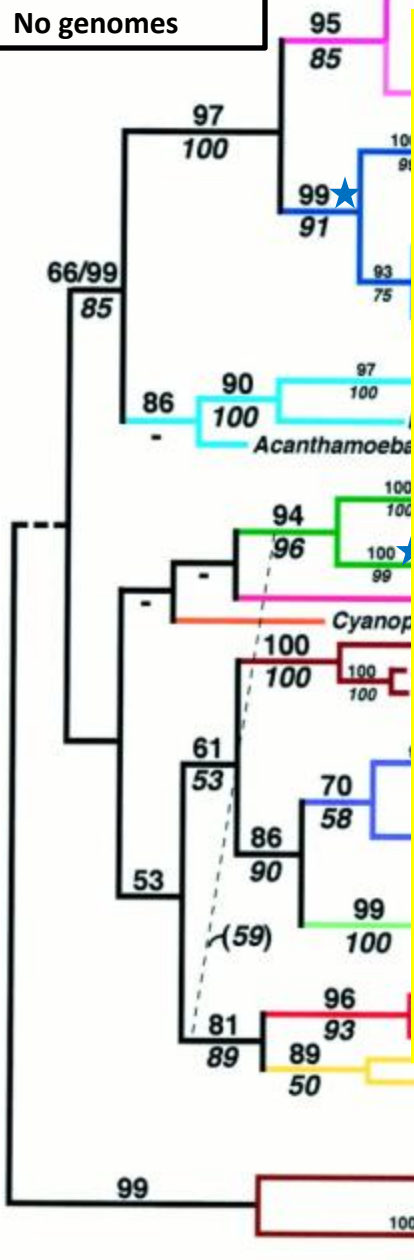


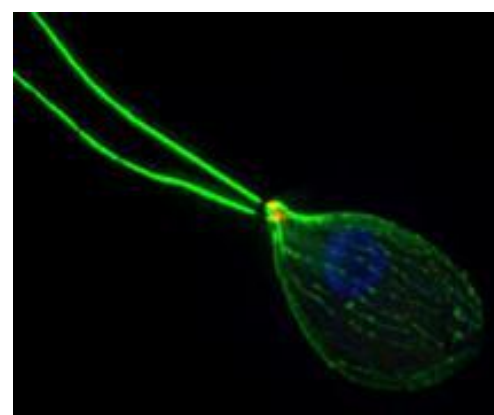
**1,033 TFs**  
**(two arrays each)**

**131 Species**

**54 Families**

- Fungi
- (Microsporidia)
- Metazoa
- Mycetozoa
- Lobosa
- Viridiplantae
- Rhodophyta
- Glaucophyta
- Heterokonta
- Ciliophora
- Apicomplexa
- Euglenozoa
- Heterolobosea
- Diplomonadida
- Parabasala





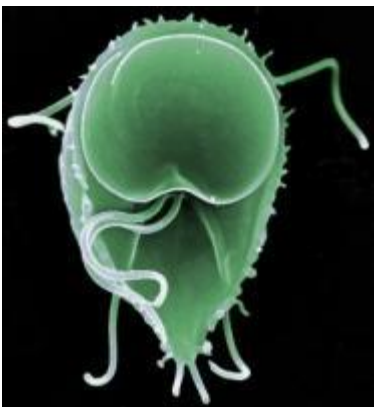
*T. nigroviridis*

*T. pseudonana*

*N. crassa*

*S. commune*

*N. gruberi*



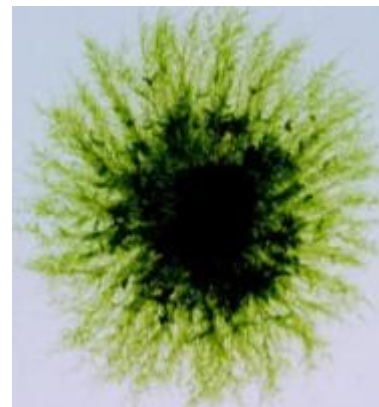
*G. lamblia*

*B. mori*

*D. discoideum*

*M. gallopavo*

*P. tetraurelia*



*O. sativa*

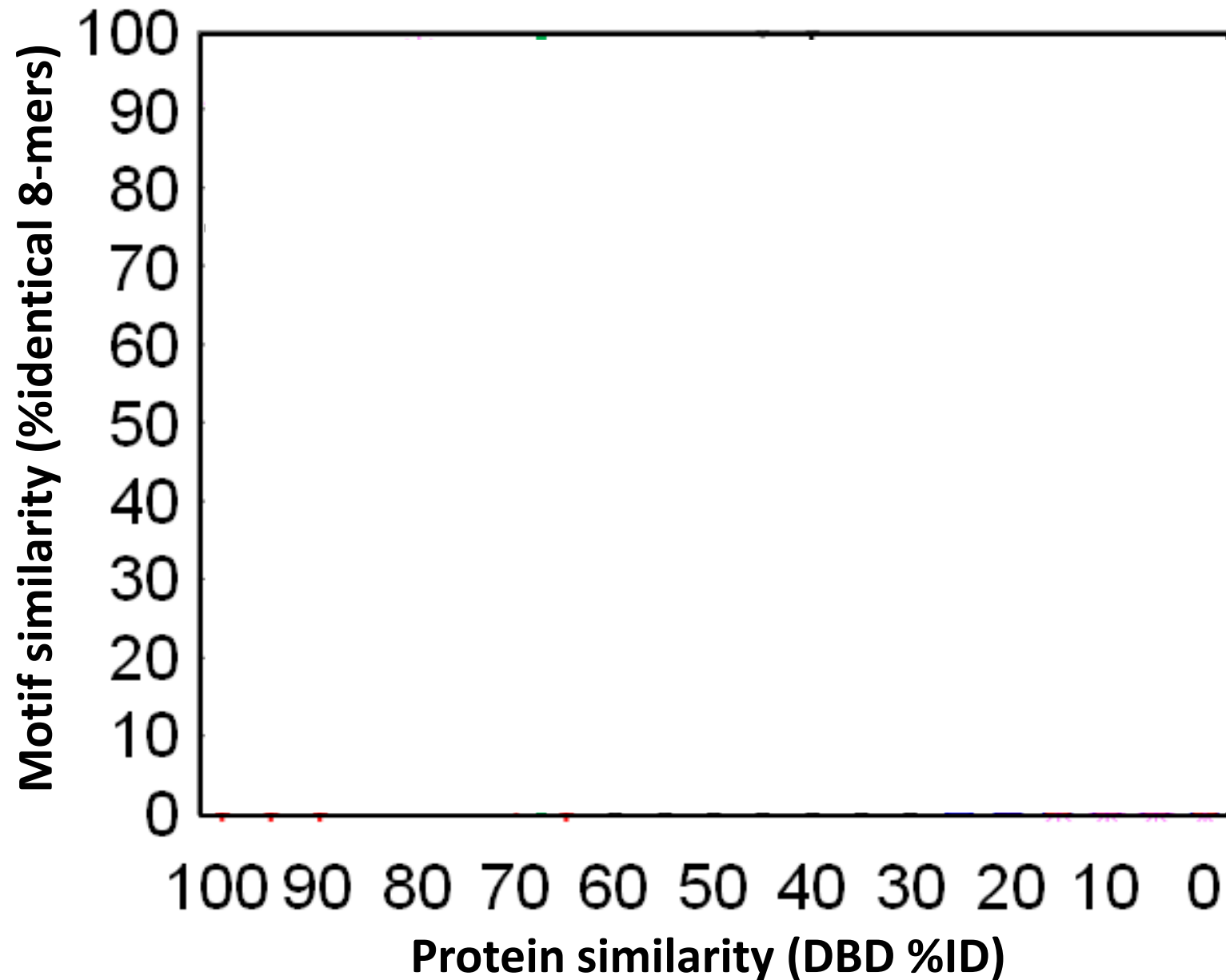
*N. vectensis*

*C. intestinalis*

*P. patens*

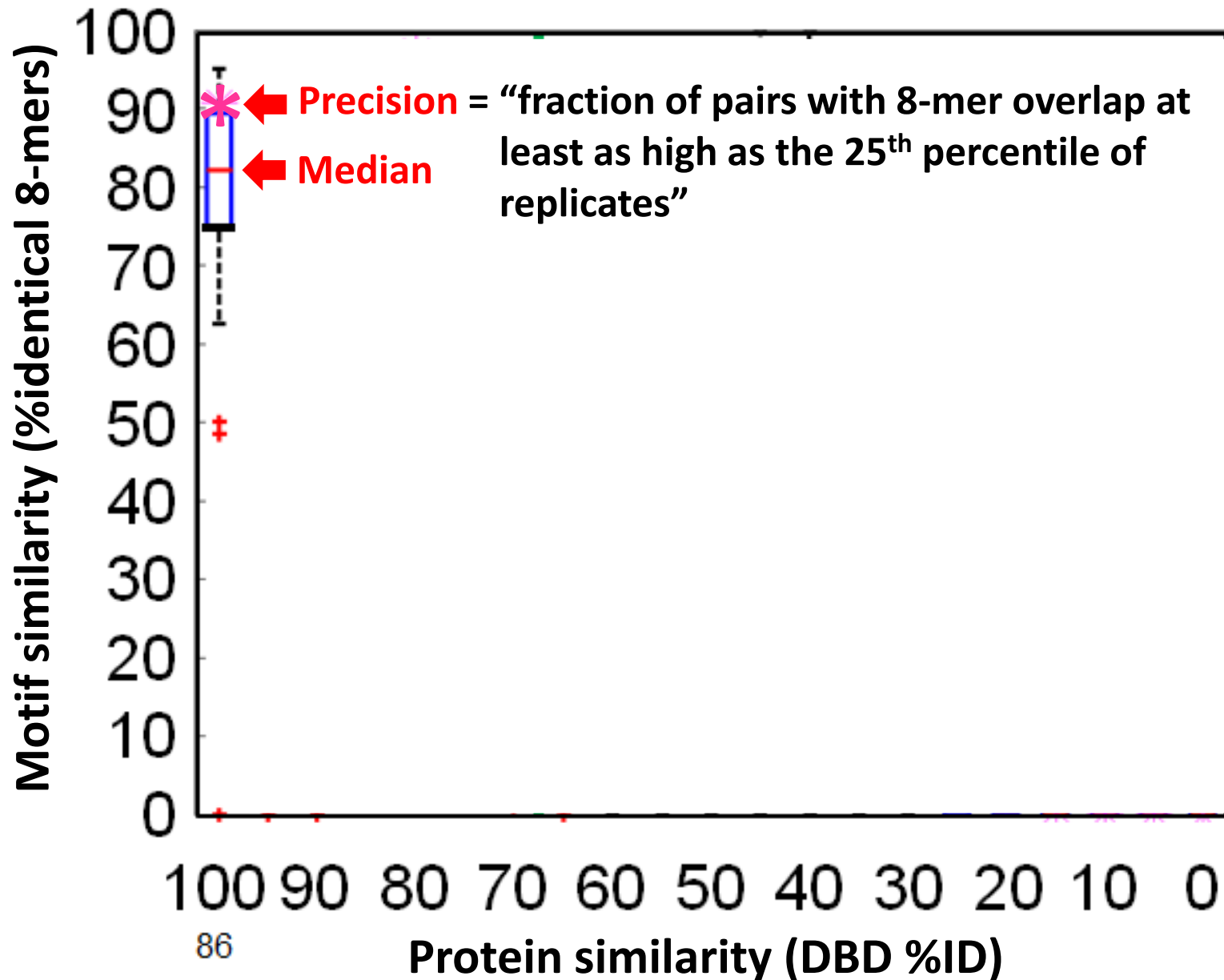
*A. carolinensis*

# 65% rule, revisited



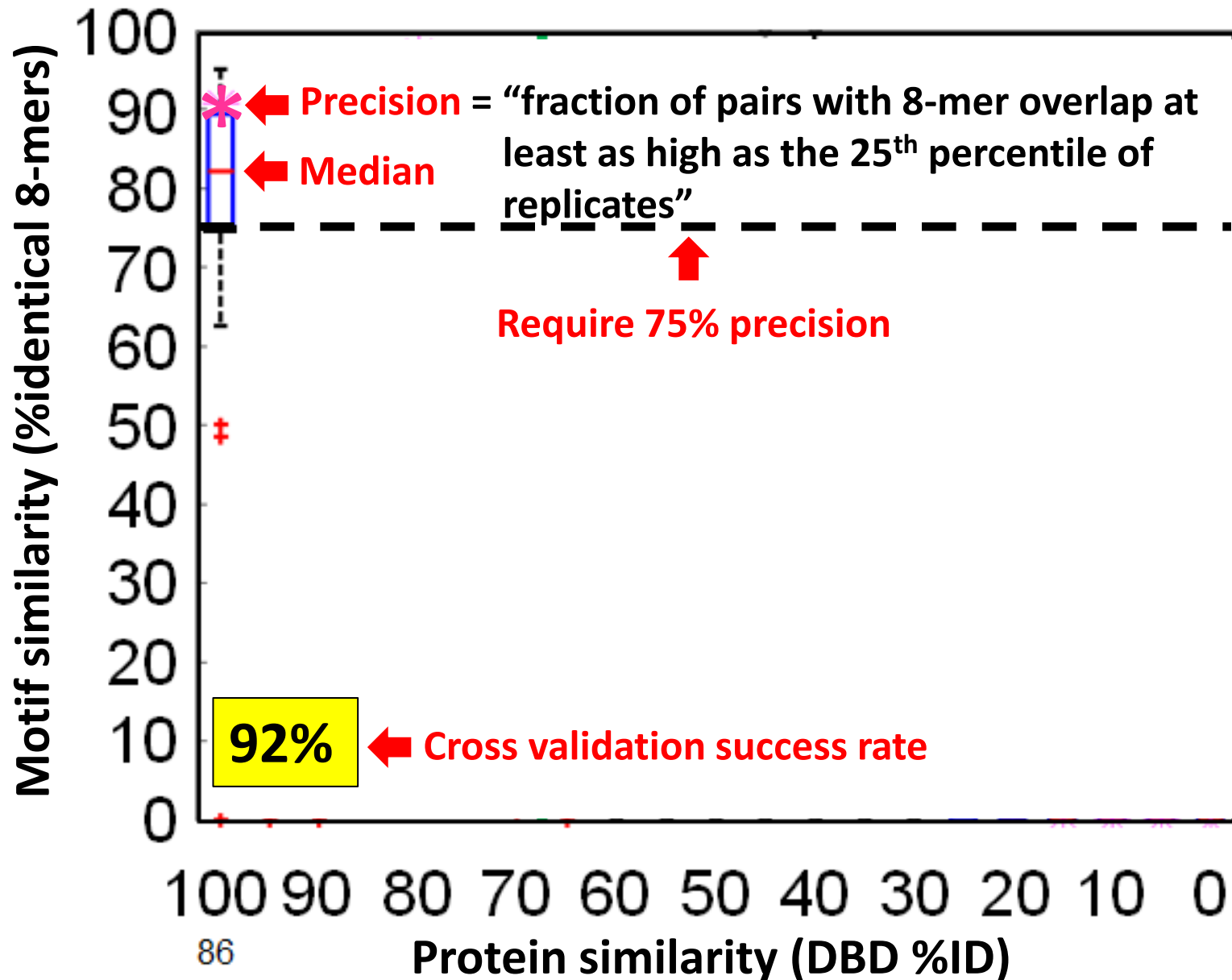


# 65% rule, revisited

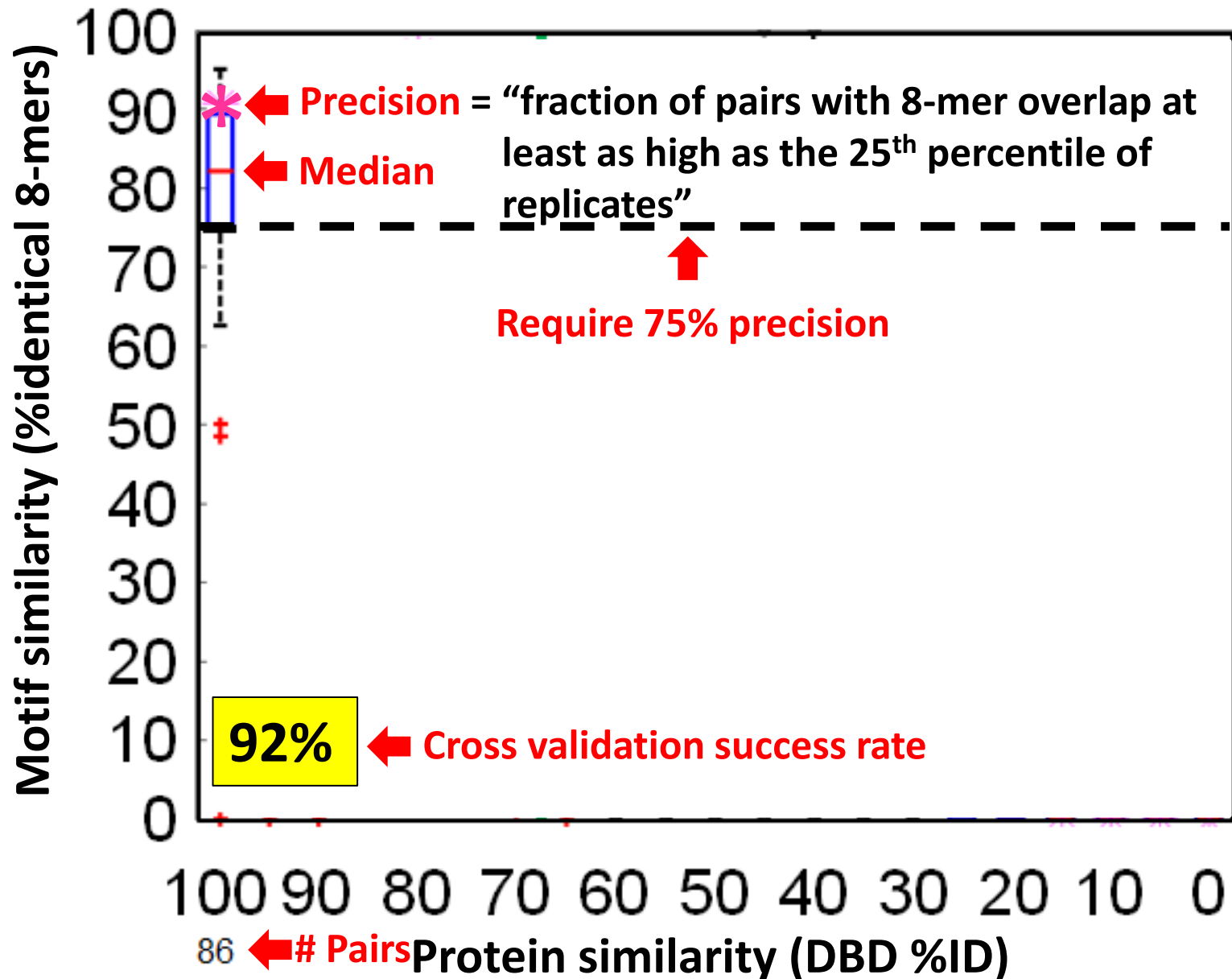




# 65% rule, revisited

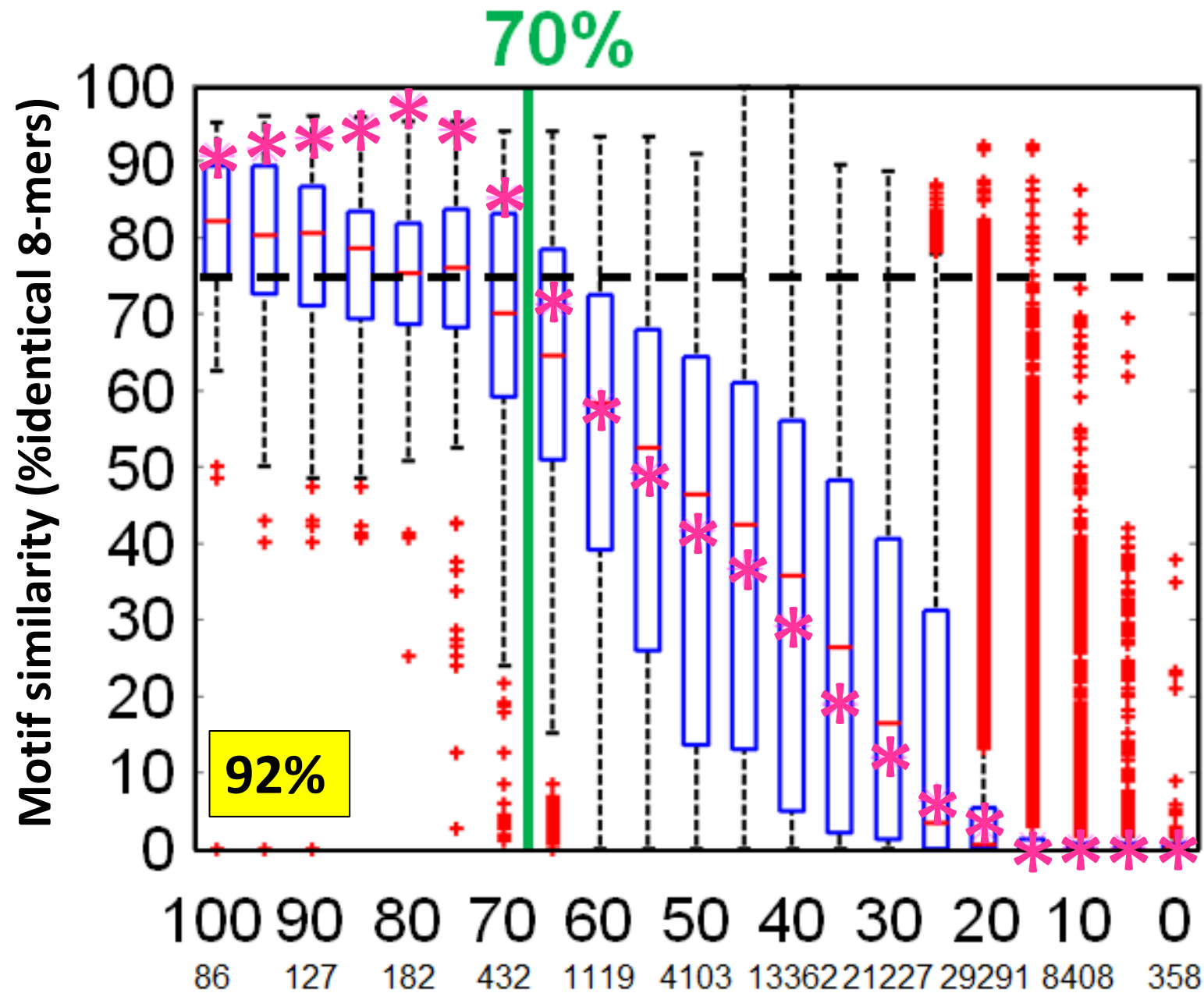


# 65% rule, revisited

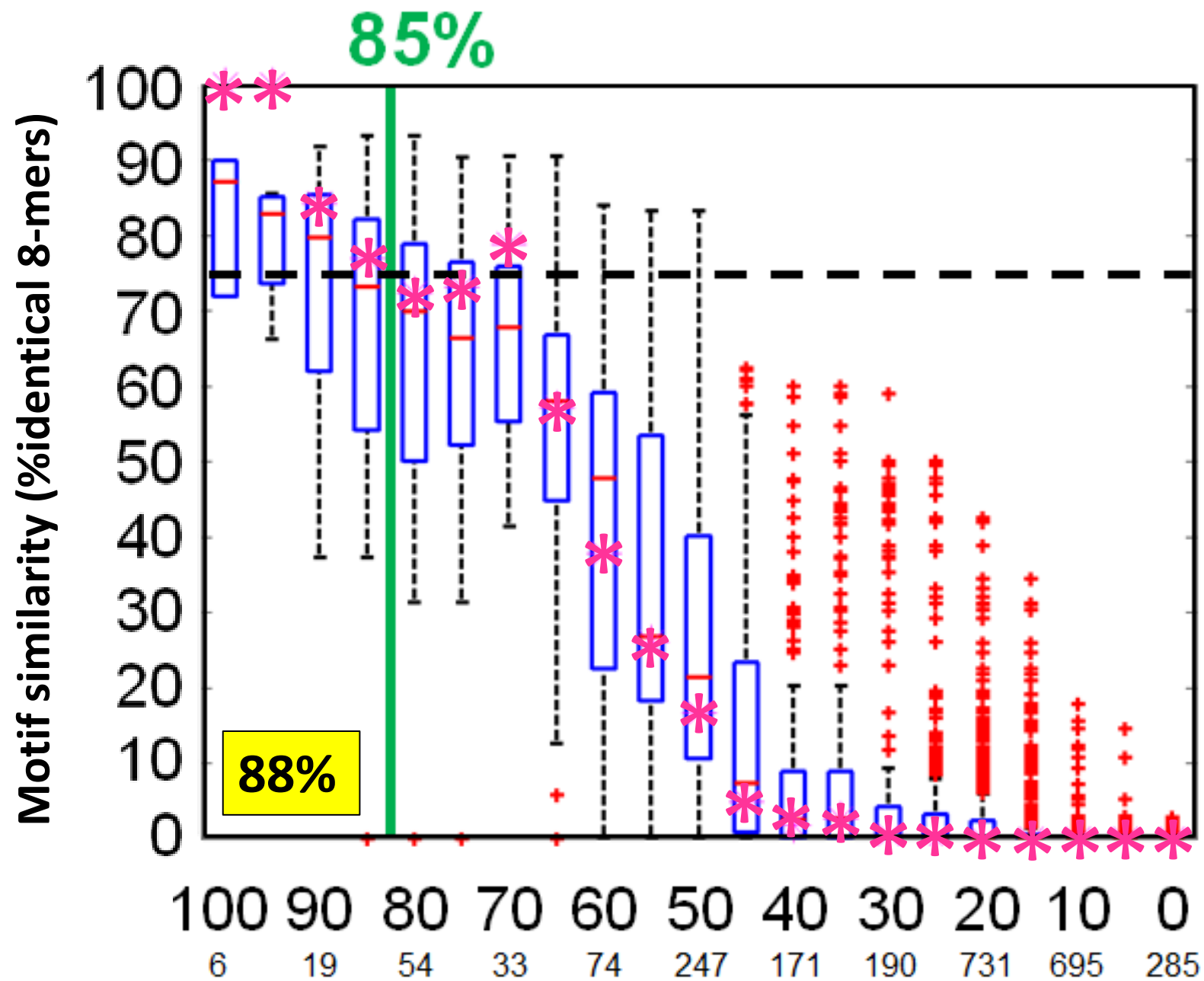




# Homeodomains

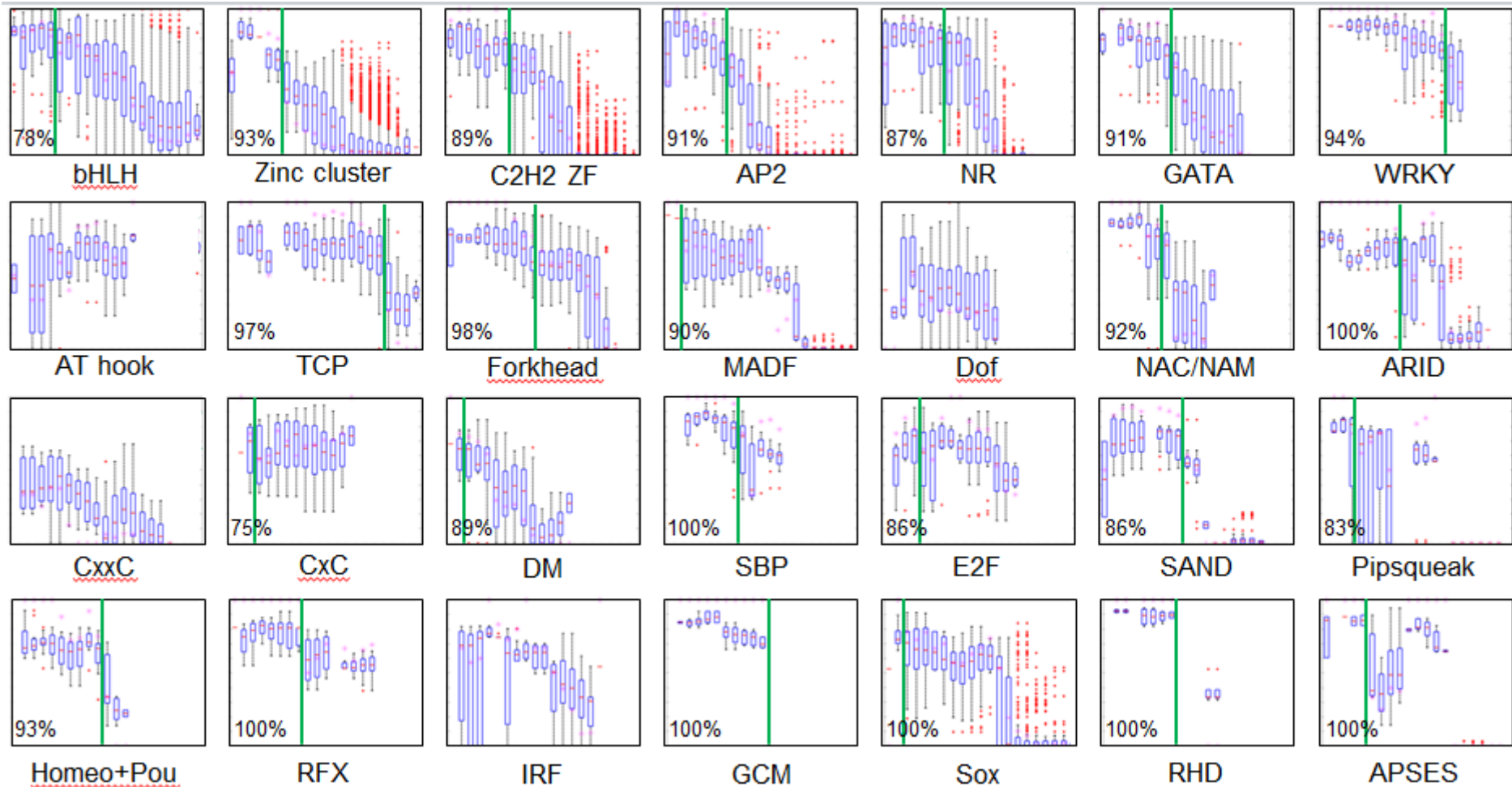


# Myb/SANT



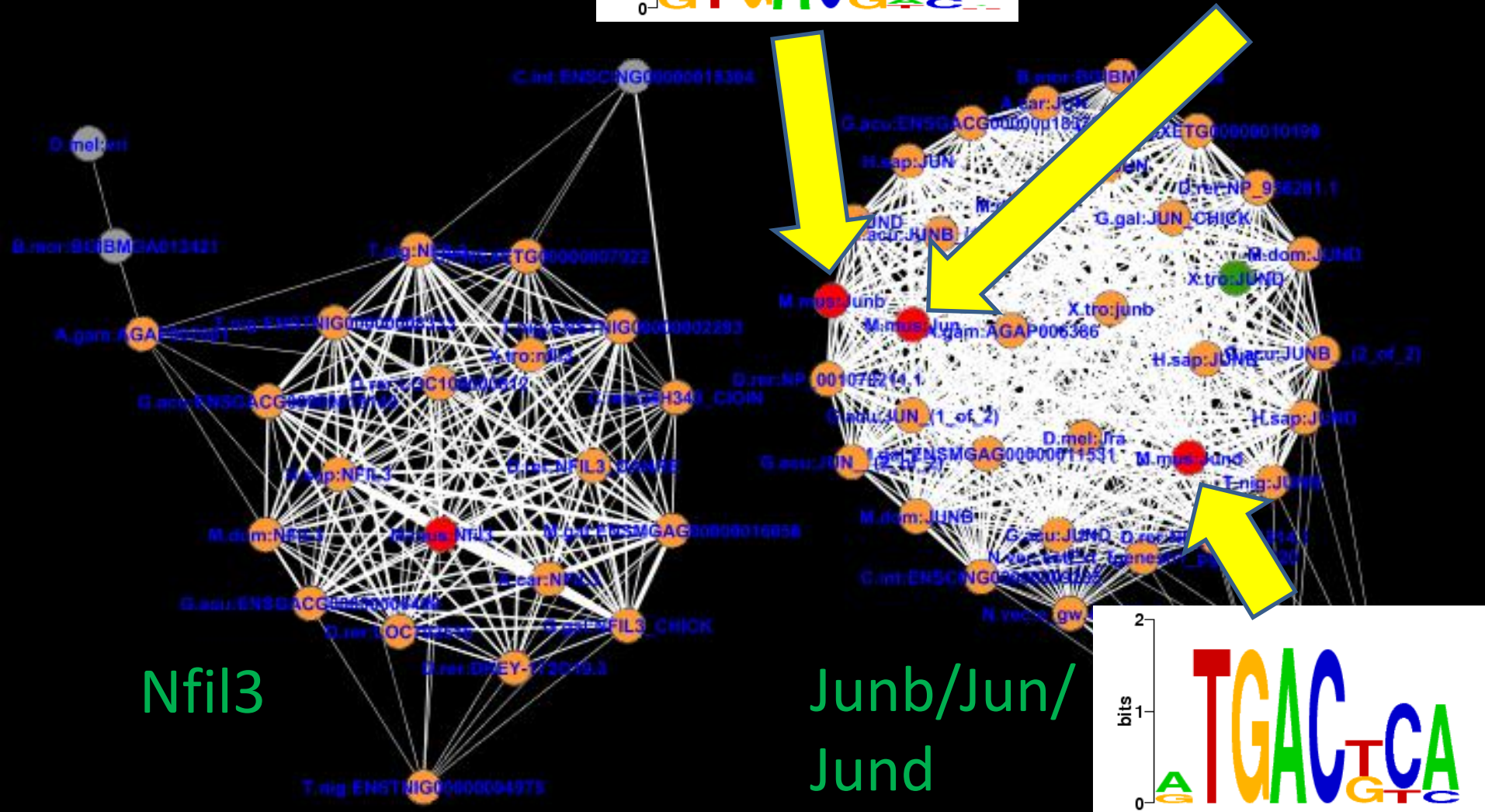


# Thresholds can be drawn for most TF families



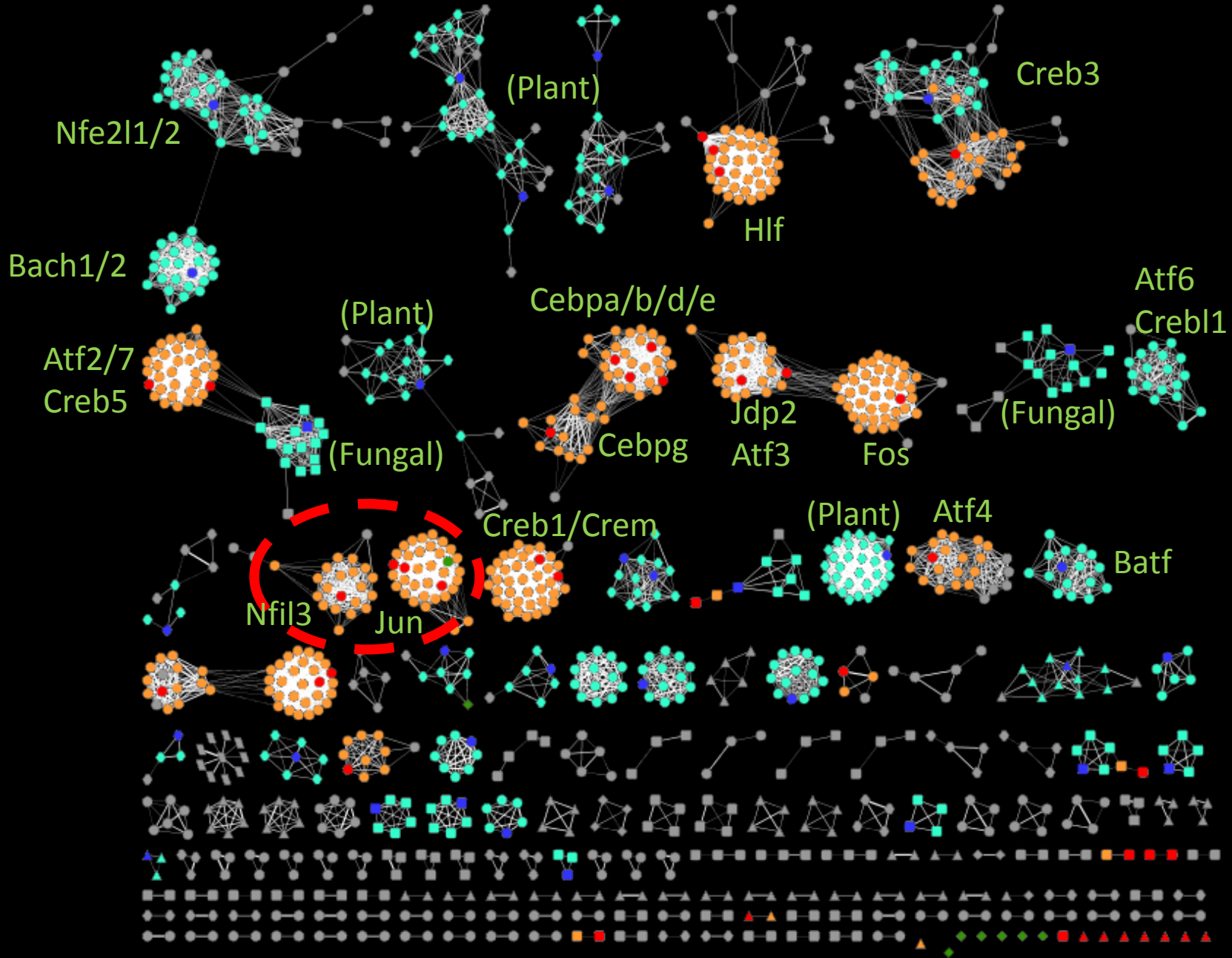
**Overall Cross Validation success rate: 89%**



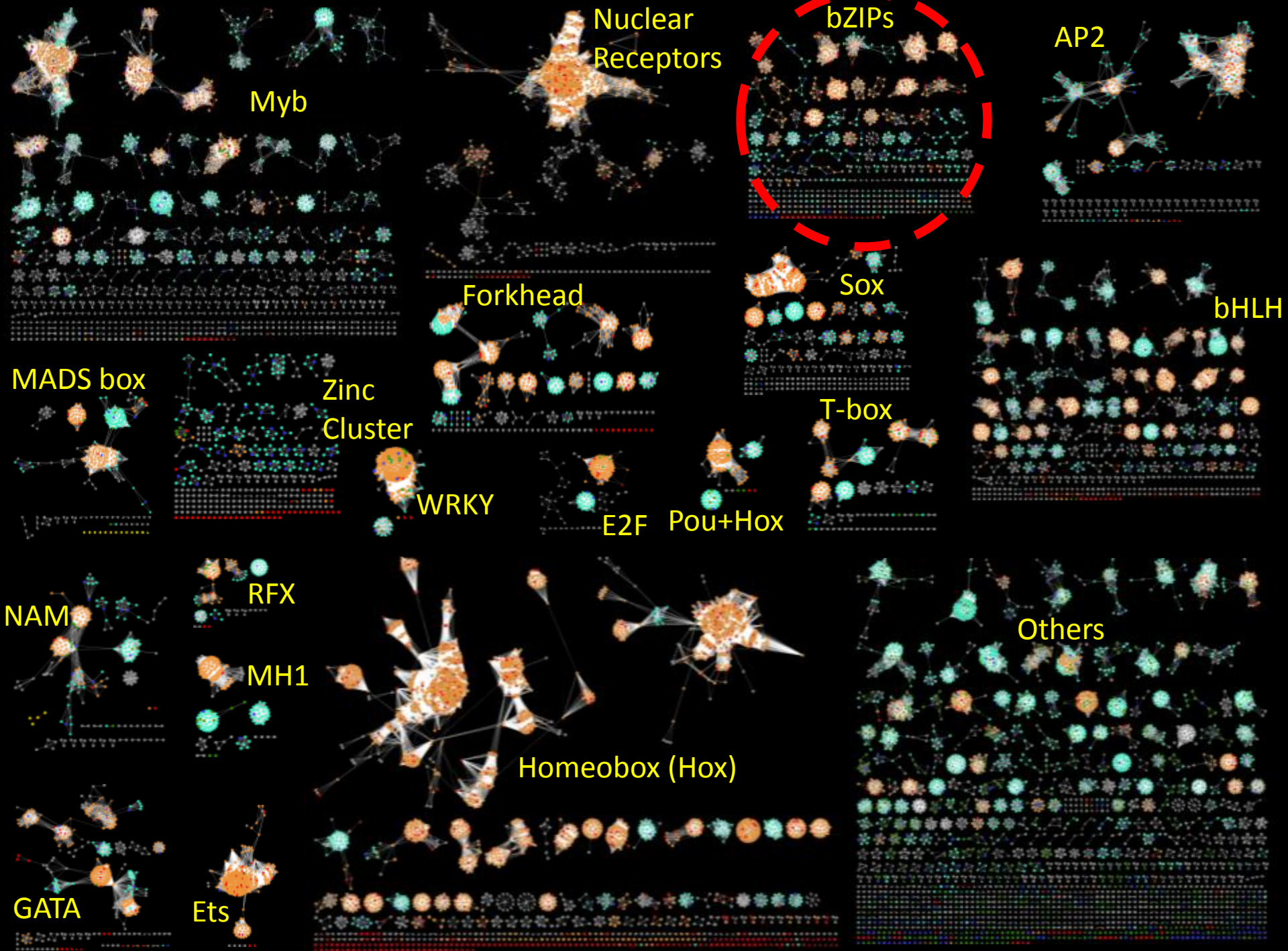




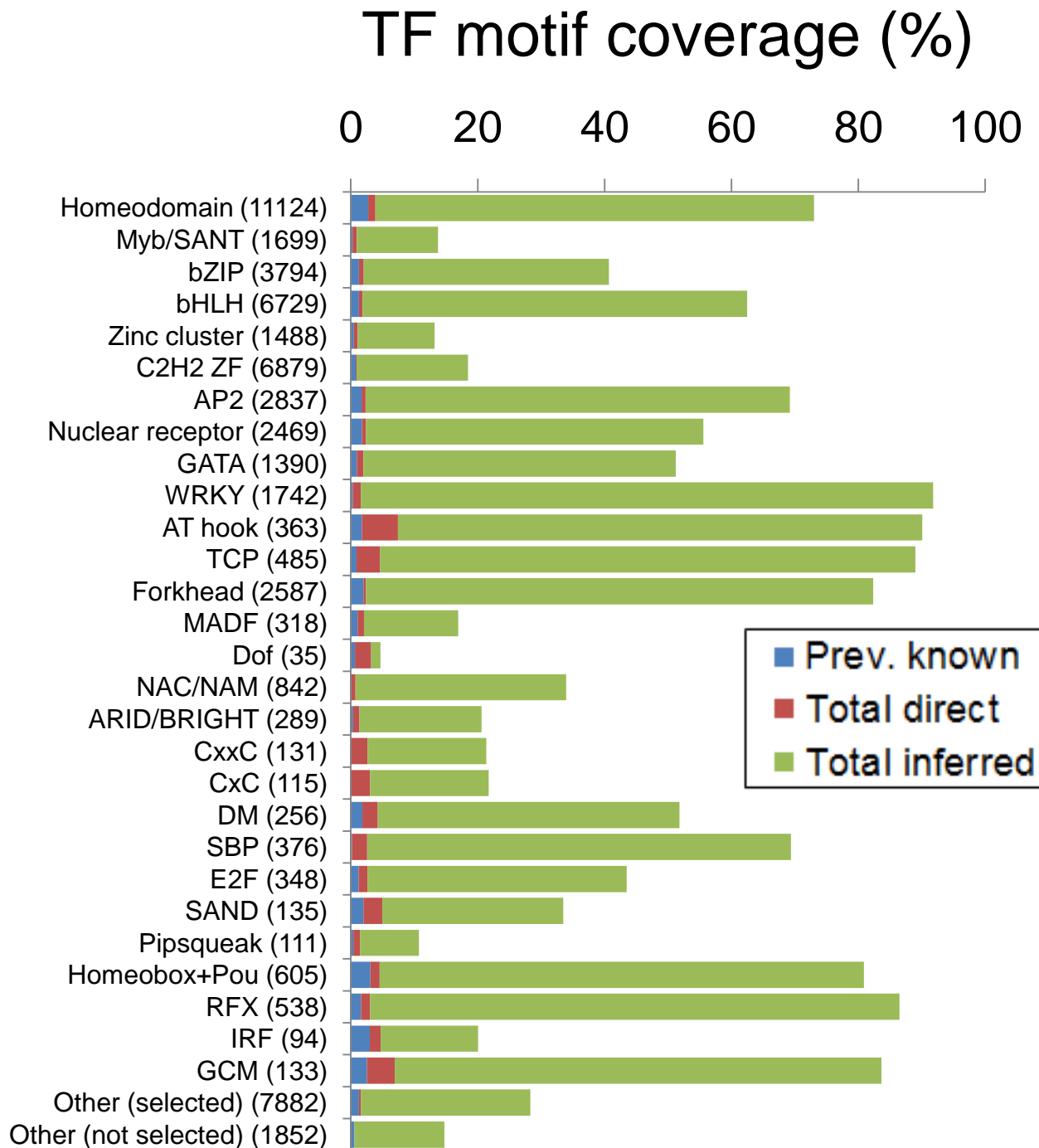
# Basic Leucine Zippers (bZIPs)



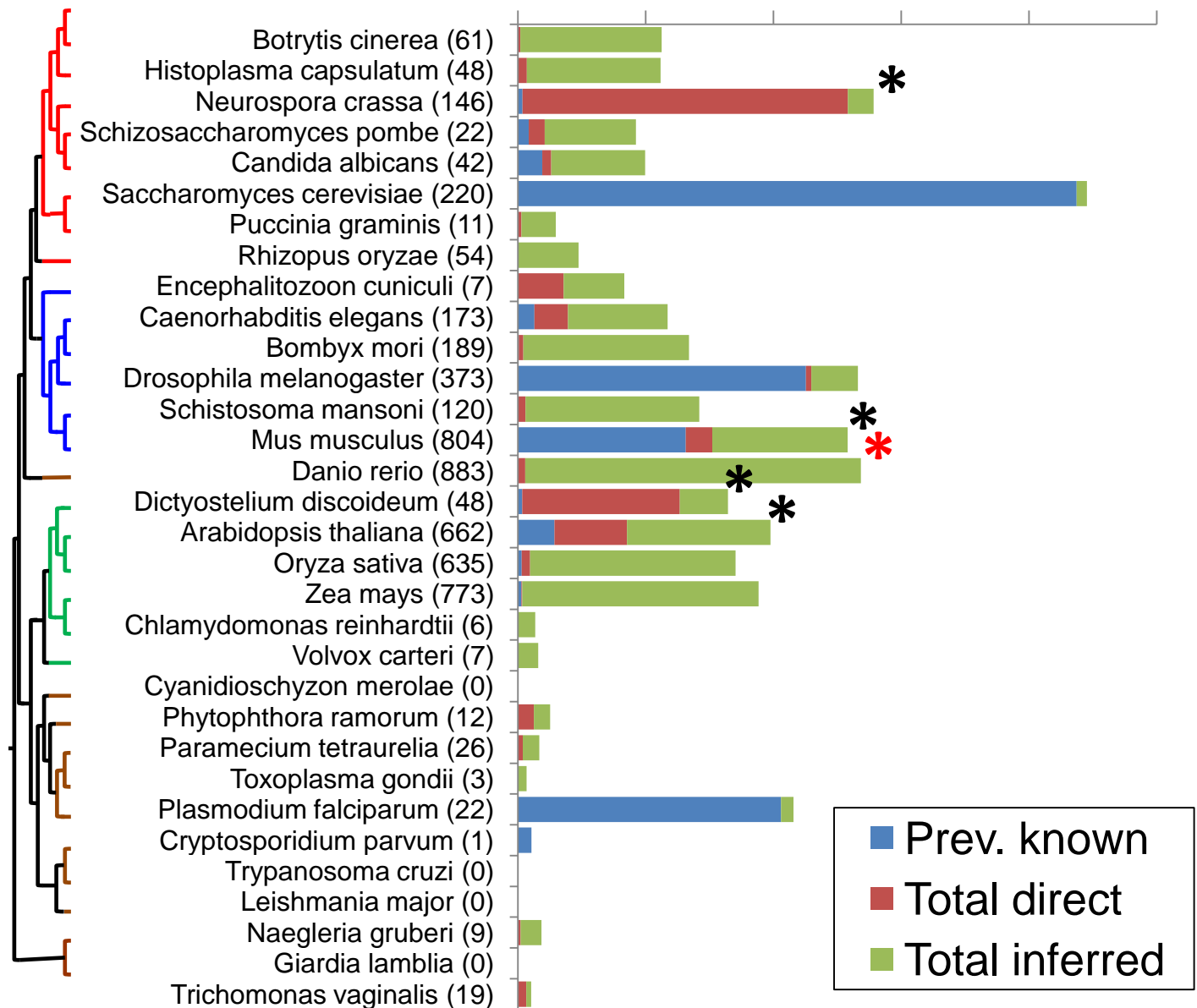




# Coverage by Family



# Coverage by Species



# Coverage summary

|                 |   |
|-----------------|---|
| <b>~162,000</b> | <b>Total estimated number of sequenced eukaryotic TFs</b> |
|-----------------|---|

# Coverage summary

|                                  |   |
|----------------------------------|---|
| <b>~162,000</b>                  | <b>Total estimated number of sequenced eukaryotic TFs</b> |
| <b>~2,000</b><br><b>(&lt;2%)</b> | <b>Number of TF motifs from other studies</b>             |

# Coverage summary

|                                  |  |
|----------------------------------|--|
| <b>~162,000</b>                  | <b>Total estimated number of sequenced eukaryotic TFs</b>            |
| <b>~2,000</b><br><b>(&lt;2%)</b> | <b>Number of TF motifs from other studies</b>                        |
| <b>~63,000</b><br><b>(39%)</b>   | <b>Number of TF motifs we currently have determined or can infer</b> |

# Coverage summary

|                                  |  |
|----------------------------------|--|
| <b>~162,000</b>                  | <b>Total estimated number of sequenced eukaryotic TFs</b>            |
| <b>~2,000</b><br><b>(&lt;2%)</b> | <b>Number of TF motifs from other studies</b>                        |
| <b>~63,000</b><br><b>(39%)</b>   | <b>Number of TF motifs we currently have determined or can infer</b> |
| <b>~80,000</b><br><b>(49%)</b>   | <b>Number of TF motifs (lenient, 50% precision threshold)</b>        |

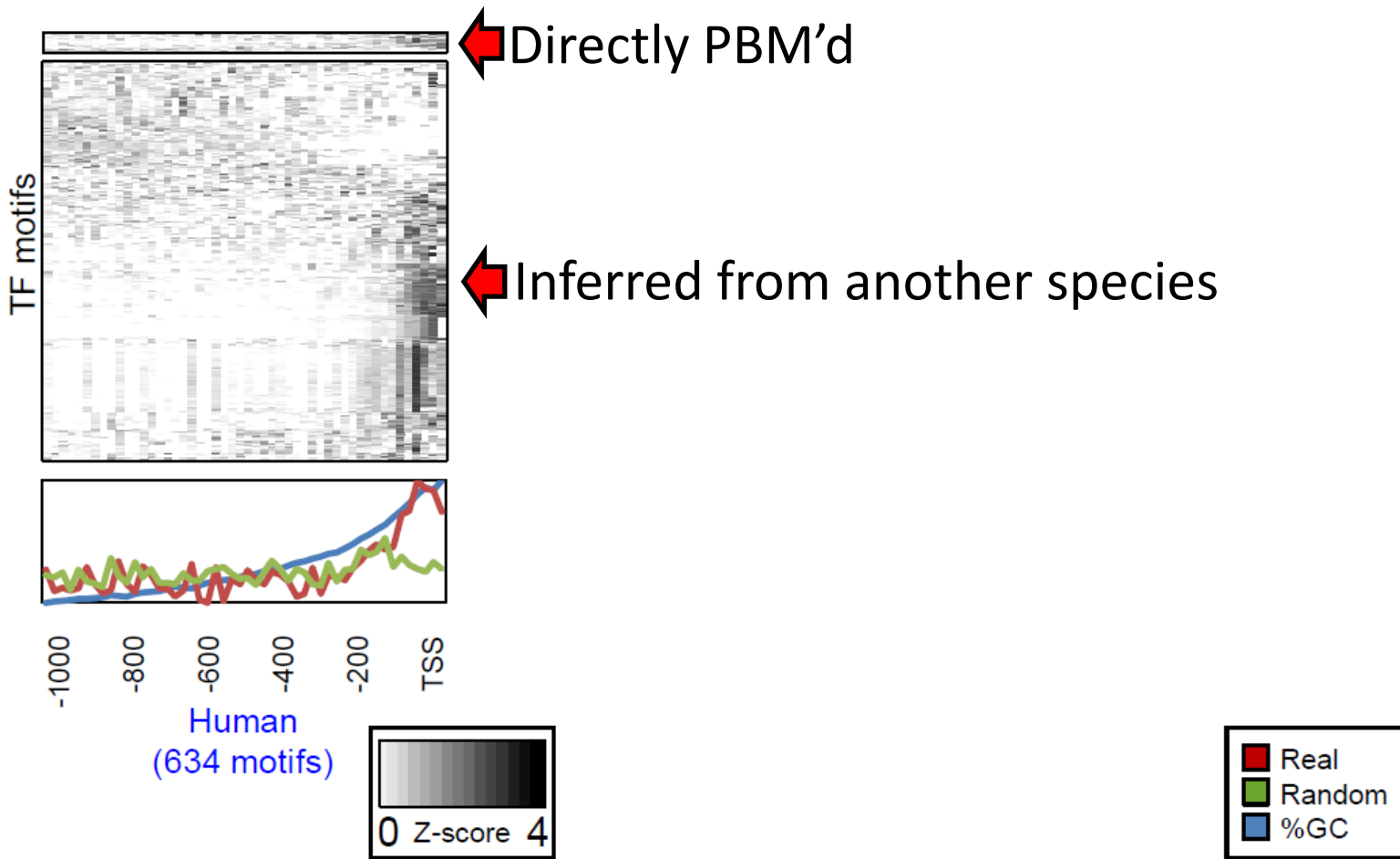
# Coverage summary

|                                  |  |
|----------------------------------|--|
| <b>~162,000</b>                  | <b>Total estimated number of sequenced eukaryotic TFs</b>            |
| <b>~2,000</b><br><b>(&lt;2%)</b> | <b>Number of TF motifs from other studies</b>                        |
| <b>~63,000</b><br><b>(39%)</b>   | <b>Number of TF motifs we currently have determined or can infer</b> |
| <b>~80,000</b><br><b>(49%)</b>   | <b>Number of TF motifs (lenient, 50% precision threshold)</b>        |

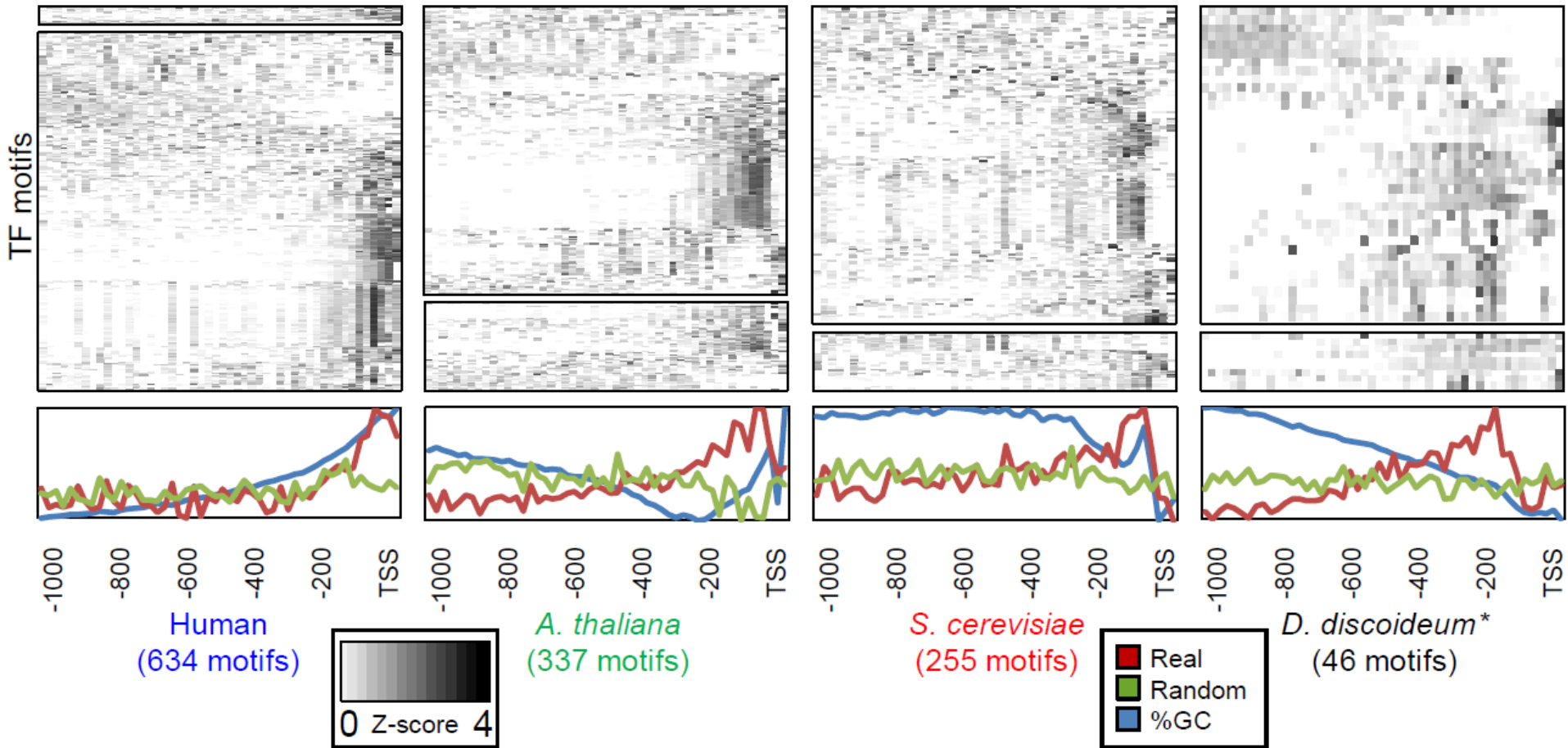
**We recently ordered primers  
for ~3000 more proteins**



# Positional bias of motifs in eukaryotic promoters

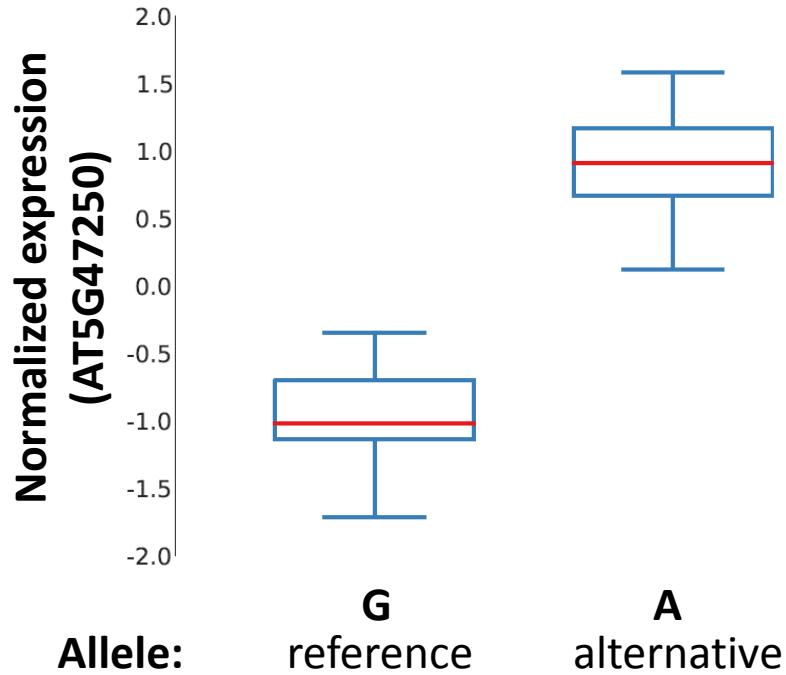


# Positional bias of motifs in eukaryotic promoters

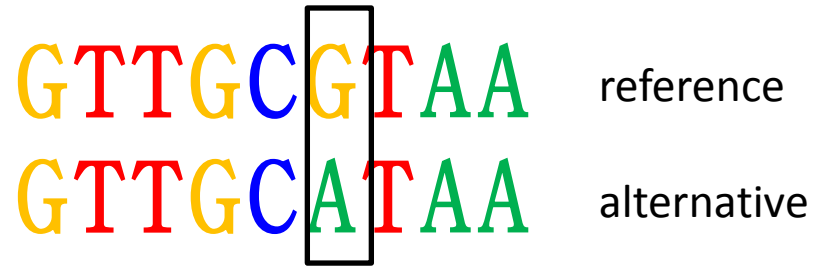
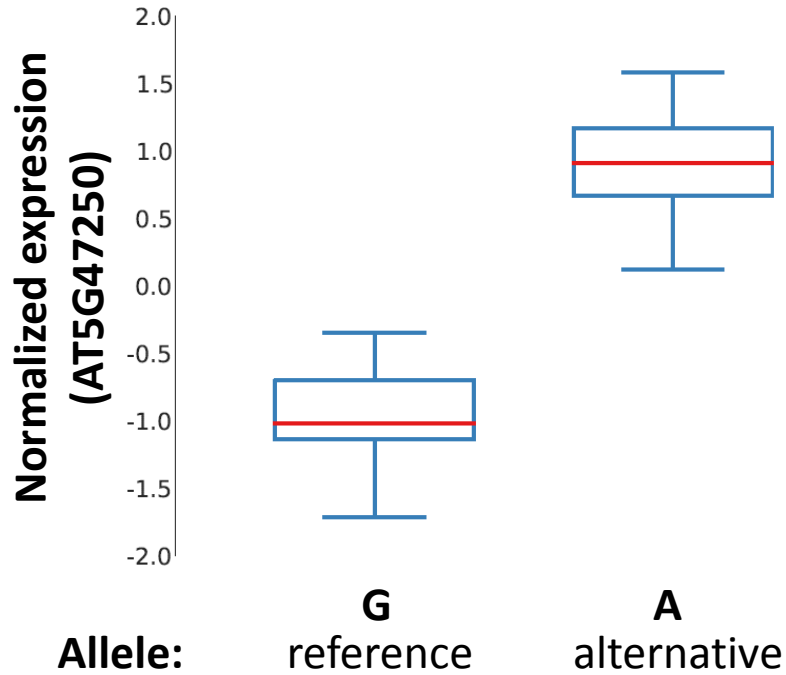


Both direct + inferred motifs “peak” near TSS

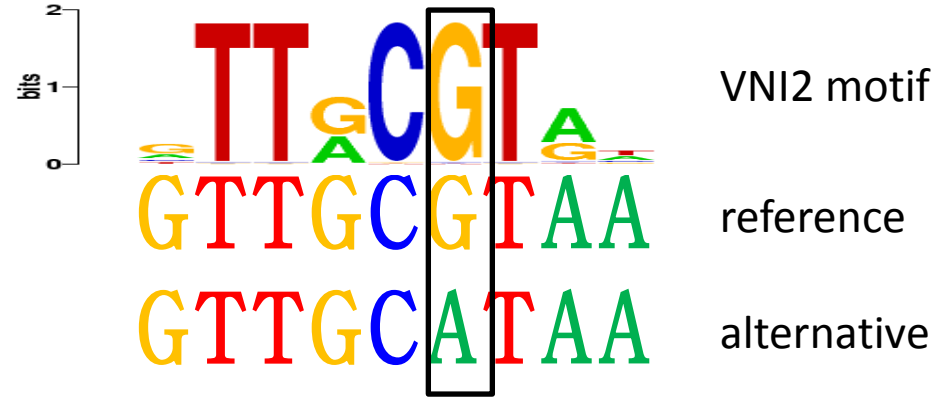
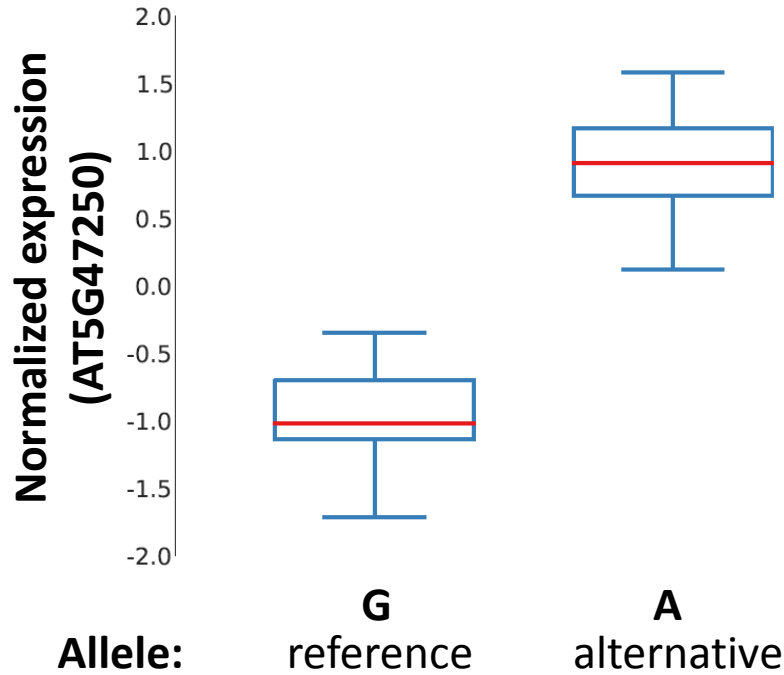
# Motifs tie genotype to expression



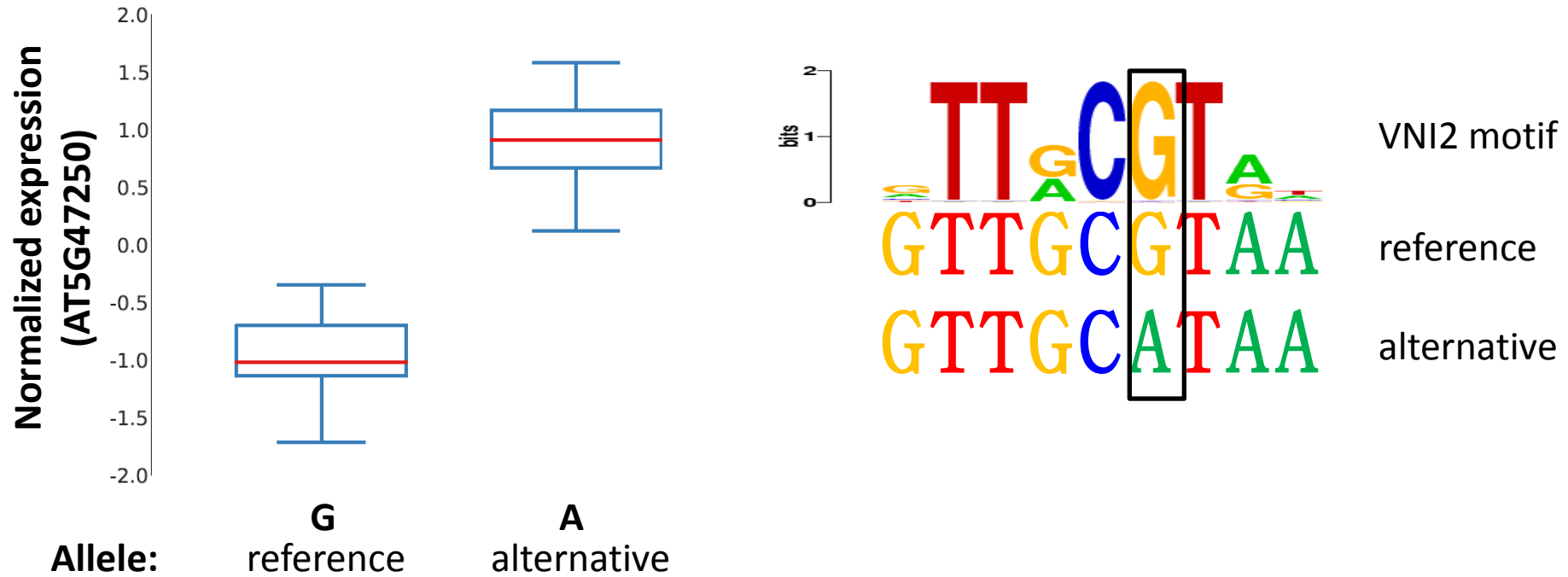
# Motifs tie genotype to expression



# Motifs tie genotype to expression



# Motifs tie genotype to expression



- VNI2 is a repressor (Yamaguchi *et al.* Plant Cell, 2010).
  - The alternative (A) allele “breaks” a binding site for this repressor, causing the de-repression of AT5G47250
- Both genes involved in defense response against the same pathogen

# Identifying TFs affected by human disease SNPs

- **Procedure:**
  - Use all PBM data to “score” risk and non-risk alleles of known disease-associated SNPs
  - Rank all TFs based on the likelihood that their binding would be disrupted by the SNP

# Identifying TFs affected by human disease SNPs

- **Procedure:**
  - Use all PBM data to “score” risk and non-risk alleles of known disease-associated SNPs
  - Rank all TFs based on the likelihood that their binding would be disrupted by the SNP
- **Results:**
  - Applied to set of 16 experimentally confirmed TF/disease SNP pairs

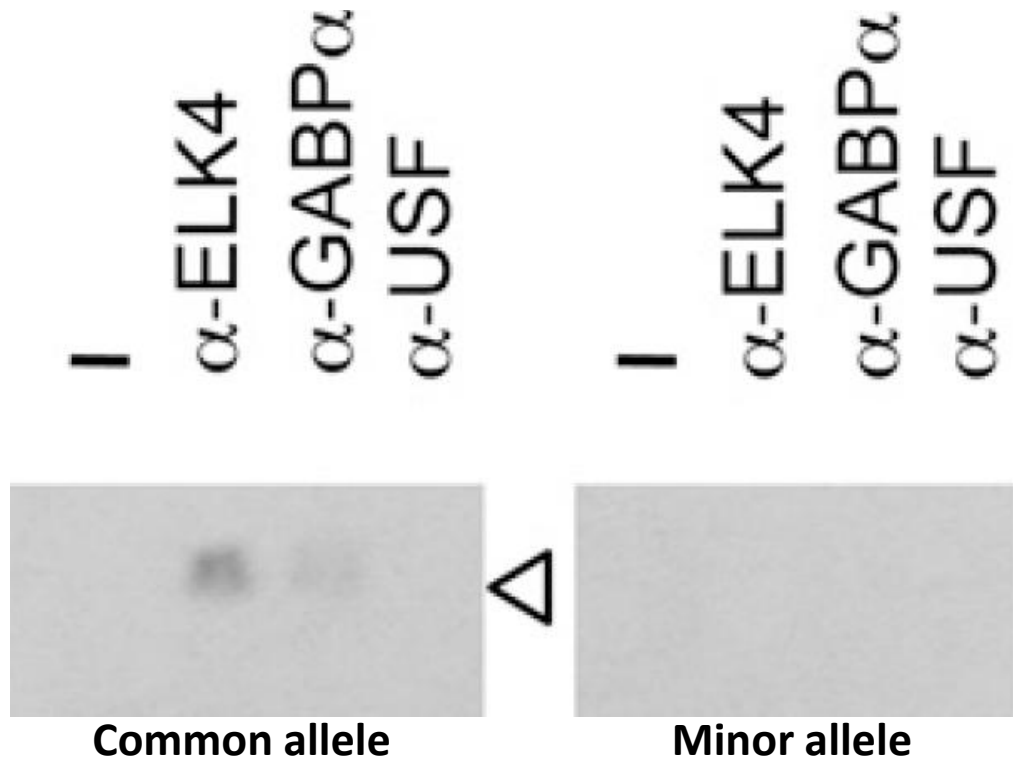


# Identifying TFs affected by human disease SNPs

- **Procedure:**
  - Use all PBM data to “score” risk and non-risk alleles of known disease-associated SNPs
  - Rank all TFs based on the likelihood that their binding would be disrupted by the SNP
- **Results:**
  - Applied to set of 16 experimentally confirmed TF/disease SNP pairs
  - Could rank correct TF (or highly related) in the top five for 10 out of 16

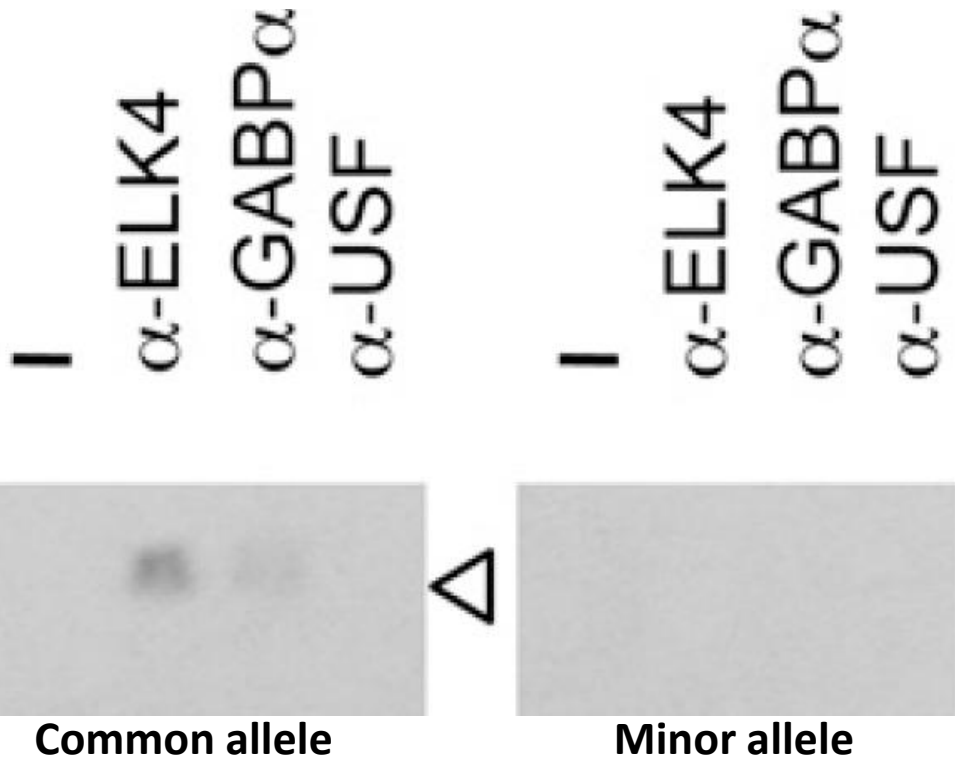
# Example: differential TF binding to rs554219 in breast cancer tumors

## Supershift EMSA



# Example: differential TF binding to rs554219 in breast cancer tumors

## Supershift EMSA



## Computational

| TF               | Common | Minor |
|------------------|--------|-------|
| 1. ELK4          | 0.498  | 0.131 |
| 2. GABP $\alpha$ | 0.496  | 0.329 |

(ranked out of >800 possible TFs)

# Application to lupus: IRF5 and rs4728142

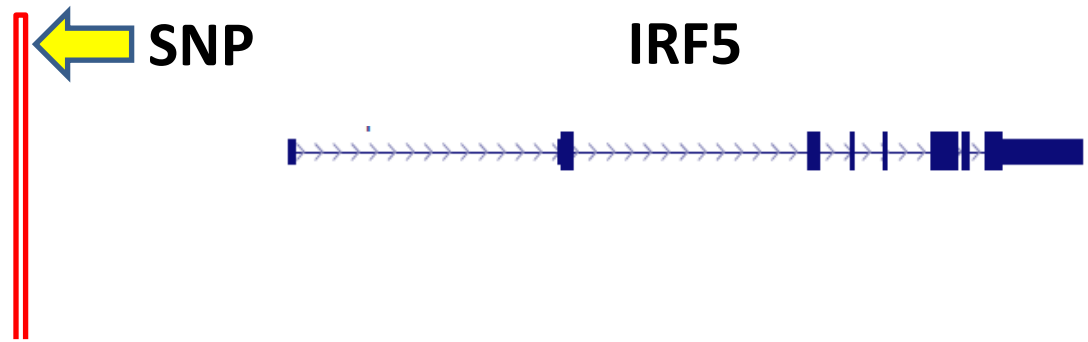
- **IRF5 (Interferon Regulatory Factor 5)**
  - Controls important inflammatory pathways involved in lupus, autoimmunity, and infectious disease
  - Strongly associated with lupus
  - Over-expressed in lupus patients
- **rs4728142**
  - Most likely of six possible causal variants for IRF5's association with lupus (Harley/Kottyan labs, targeted deep sequencing of 1,000s of cases and controls)
  - Risk allele results in increased binding of some protein in blood cell nuclear extracts (Kristjansdottir *et al.*, 2008), and in B cell lines (Harley lab, unpublished data)

# Genomic context of rs4728142

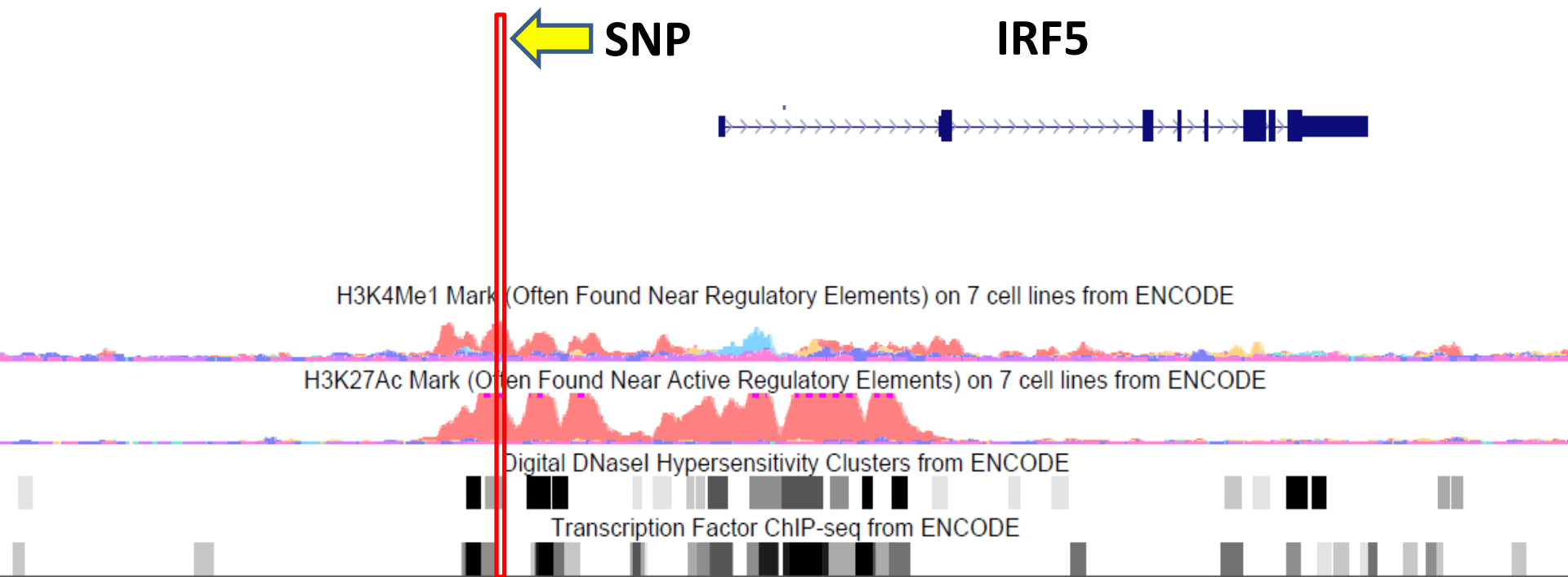
IRF5



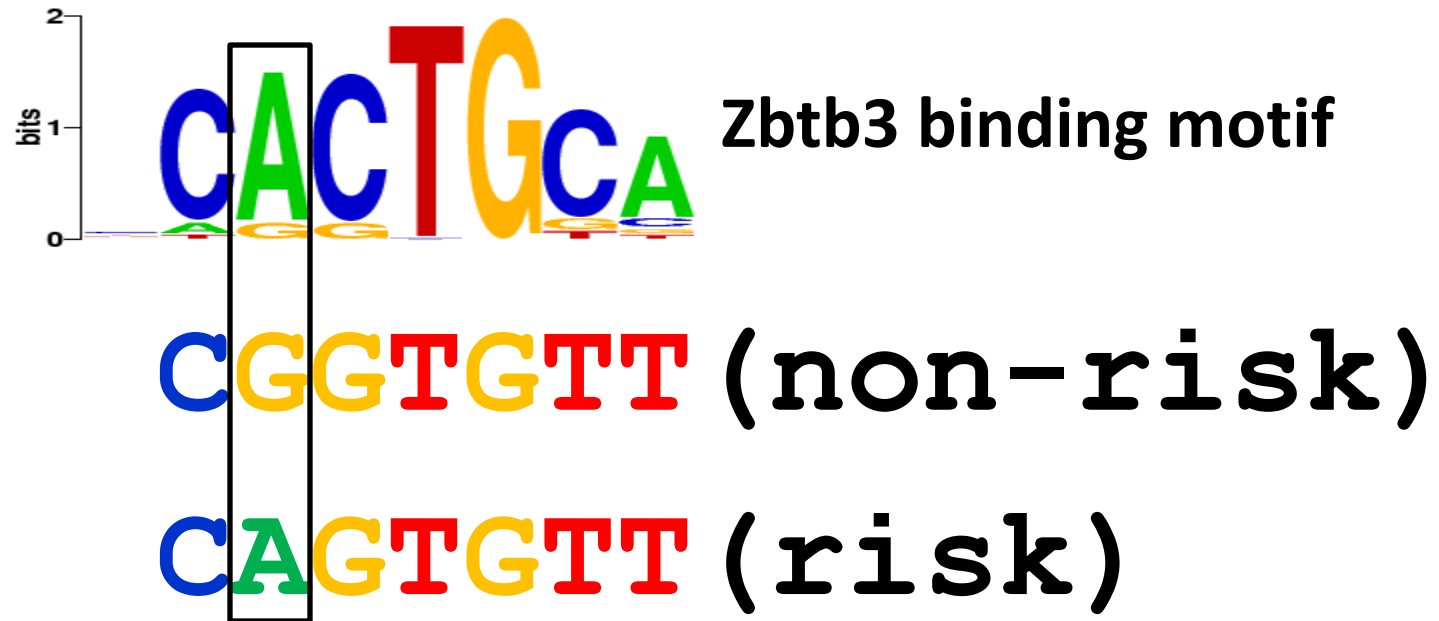
# Genomic context of rs4728142



# Genomic context of rs4728142



# Prediction: stronger binding of Zbtb3



Quantitative score from PBM:

0.03 (won't bind) to 0.45 (strong binding)

Zbtb3 is highly expressed in B cells



**Hypothesis:  
the rs4728142 risk allele  
“creates” a binding site for  
ZBTB3 in B cells, resulting in  
higher expression of the IRF5  
gene in lupus patients**

**Hypothesis:  
the rs4728142 risk allele  
“creates” a binding site for  
ZBTB3 in B cells, resulting in  
higher expression of the IRF5  
gene in lupus patients**

**CONFIRMED VIA SUPERSHIFT EMSA**

# Website frontpage

CIS-BP Database: Catalog of Inferred Sequence Binding Preferences

Home

Tools

View cart

Bulk downloads

Database stats

Contact us

Help

Update Log

FAQ

Links

How to cite



Welcome to CIS-BP, the online library of transcription factors and their DNA binding motifs.

## Search for a TF

By Identifier   
(e.g. Gata\*, YEL009C, ISFTZ\_01)

## Browse TFs / Restrict Search for TFs

By Model Organism

By Any Species

By Domain Type

By Motif Evidence

By Evidence Type

By Study

Database Build




Last updated: Sep 2nd, 2014 Database Build 1.00

Current database contents: 65713 TFs with at least one binding motif (3875 from direct experiments), out of a total of 167077 Eukaryotic TFs from 261 families in 340 species

<http://cisbp.cabr.utoronto.ca/>

# TF pages

 CIS-BP Database: Catalog of Inferred Sequence Binding Preferences



[Home](#)  
[Tools](#)  
[Download cart](#)  
[Bulk downloads](#)  
[Database stats](#)  
[Contact us](#)  
[Help](#)  
[How to cite](#)

## CREB1 (*Homo sapiens*) bZIP







### TF Information

| Pfam ID                          | Interpro ID               | Gene ID                         | CIS-BP ID      | Sequence source                       |
|----------------------------------|---------------------------|---------------------------------|----------------|---------------------------------------|
| <a href="#">PF00170 (bZIP_1)</a> | <a href="#">IPR011616</a> | <a href="#">ENSG00000118260</a> | <b>T086497</b> | <a href="#">Ensembl (2011-Oct-26)</a> |

### Directly determined binding motifs

| Name/Motif ID  | Species             | Forward  | Reverse   | Type/Study/Study ID  | DBD Identity |
|----------------|---------------------|--|---|--|--------------|
| CREB1<br>M1305 | <i>Homo sapiens</i> |  |  | SELEX<br><a href="#">Portales-Casamar et al.(2010)</a><br>MA0018.1 | (Direct)     |

### Motifs from related TFs

| Name/Motif ID                  | Species             | Forward  | Reverse   | Type/Study/Study ID                             | DBD Identity |
|--------------------------------|---------------------|--|---|---|--------------|
| <a href="#">Creb1</a><br>M0352 | <i>Mus musculus</i> |  |  | PBM<br><a href="#">Zoo et al.(0)</a><br>pTH5080 | 1.000        |
| <a href="#">Crem</a><br>M0375  | <i>Mus musculus</i> |  |  | PBM<br><a href="#">Zoo et al.(0)</a><br>pTH5002 | 0.934        |
| <a href="#">Atf1</a>           | <i>Mus musculus</i> |  |  | PBM<br><a href="#">Zoo et al.(0)</a>            | 0.852        |

# Bulk Downloads

CIS-BP Database: Catalog of Inferred Sequence Binding Preferences

| Categories     | Selection            | Logos                               | E-scores                 | Z-scores                 | Probe Intensities        | TF info                             | PWMs                                | Action                    |                                   |
|----------------|----------------------|-------------------------------------|--------------------------|--------------------------|--------------------------|-------------------------------------|-------------------------------------|---------------------------|-----------------------------------|
| By Species     | <input type="text"/> | <input checked="" type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | Download Species Archive! |                                   |
| By Family      | <input type="text"/> | <input checked="" type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | Download Family Archive!  |                                   |
| Entire dataset |                      |                                     |                          |                          |                          |                                     |                                     |                           | Download Entire Datasets Archive! |
| MySQL tables   |                      |                                     |                          |                          |                          |                                     |                                     |                           | Download MySQL tables             |

Home  
Tools  
View cart  
Bulk downloads  
Database stats  
Contact us  
Help  
Update Log  
FAQ  
Links  
How to cite

Last updated: Sep 2nd, 2014 Database Build 1.00

# Tools page

CIS-BP Database: Catalog of Inferred Sequence Binding Preferences

[Home](#)  
[Tools](#)  
[View cart](#)  
[Bulk downloads](#)  
[Database stats](#)  
[Contact us](#)  
[Help](#)  
[Update Log](#)  
[FAQ](#)  
[Links](#)  
[How to cite](#)

### Scan single sequences for TF binding

Scan a DNA sequence for potential binding sites  
*Click for [sample](#)*

Species:

Motif model:  Threshold:

Scan TFs in my cart

### Scan two sequences for differential TF binding

SNP Scan  
*Click for [sample](#)*

Allele 1:

Allele 2:

Species:

Motif model:  Threshold:

8 mers - Escores: 0.45

Scan TFs in my cart


### Protein Scan

Scan an amino acid sequence to predict its DNA binding motif  
*Click for [sample](#)*

### Motif Scan

Compare a given motif to all TFs in the database  
*Click for [PWM Alignment](#), [IUPAC](#)*

Species:



Last updated: Sep 2nd, 2014 Database Build 1.00

Current database contents: 65713 TFs with at least one binding motif (3875 from direct experiments), out of a total of 167077 Eukaryotic TFs from 261 families in 340 species

# There's one for RNA binding proteins too!

CISBP-RNA Database: Catalog of Inferred Sequence Binding Preferences of RNA binding proteins

Home  
Tools  
Download cart  
Bulk downloads  
Database stats  
Contact us  
Help  
How to cite

## CISBP-RNA

Welcome to CIS-BP-RNA, the online library of RNA binding proteins and their motifs.

**Search for an RBP**

By Identifier   
(e.g. Puf\*, YOR359W, RNCMPT00046)

**Browse RBPs / Restrict Search for RBPs**

By Model Organism

By Any Species

By Domain Type

By Motif Evidence

By Evidence Type

By Study

Database Build   
*Latest build: 0.6*

**GO!**

Last updated: -- Database Build 0.6

Current database contents: 8056 RBP binding motifs(250 from direct experiments), out of a total of 62587 Eukaryotic RBPs from 55 families in 289 species

Ray, Weirauch *et al.*, Nature 2013  
<http://cisbp-rna.cabr.utoronto.ca/>

## U. Toronto

- Tim Hughes
- Ally Yang
- Mihai Albu
- Atina Cote
- Hong Zheng
- Hamed Shateri Najafabadi
- Sam Lambert
- Ishminder Mann
- Kate Cook
- Harm van Bakel
- Shaheynoor Talukder
- Andrew Vorobyov
- Anton van der Ven
- Wilfred de Vega
- Nicole Park
- Geanany Rasanathan
- Yogesh Hooda
- Sanie Mnaimneh
- Kenneth Chu
- Oliver Boright
- Jerry Li
- Agnieszka Janska

## Elsewhere

- Joseph Ecker
- Luis Larrondo
- Gunnar Ratsch
- Marian Walhout
- Alejandro Montenegro-Montero
- Philipp Drewe
- Francois-Yves Bouget
- Gadi Shaulsky
- Jean-Claude Lozano
- Mary Galli
- Mathew Lewsey
- Eryong Huang
- Tuhin Mukherjee
- Xiaoting Chen
- John Reece-Hoyes
- Sridhar Govindarajan
- Oliver Stegle

Ref:

Weirauch *et al.* Cell, 2014.

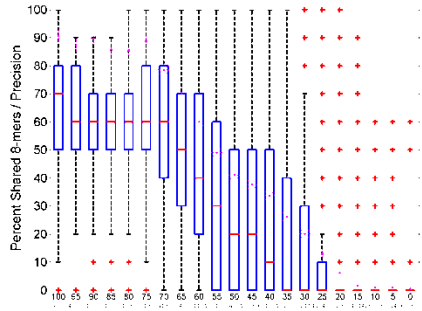
Contact:



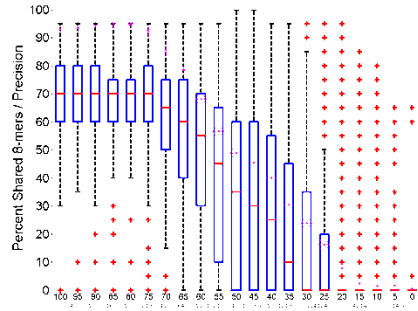


# Homeodomains

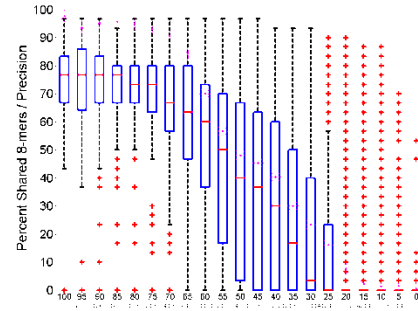
## (8 different measures of 8-mer similarity)



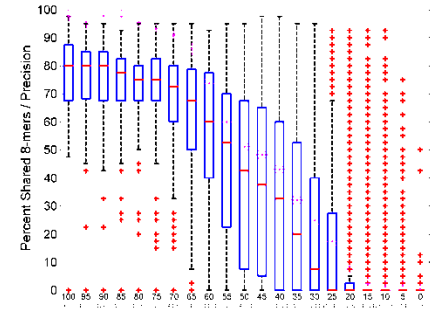
Zscore10



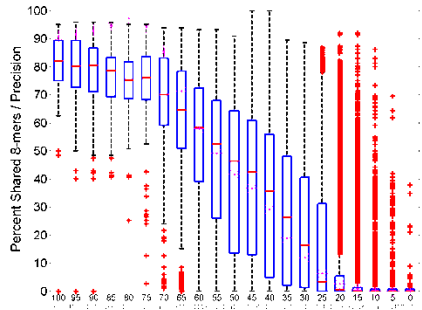
Zscore20



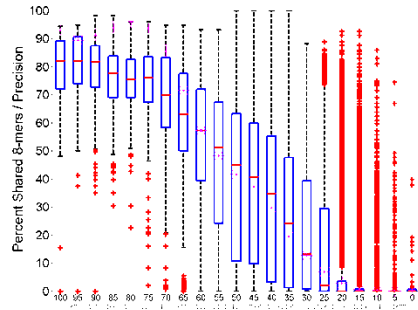
Zscore30



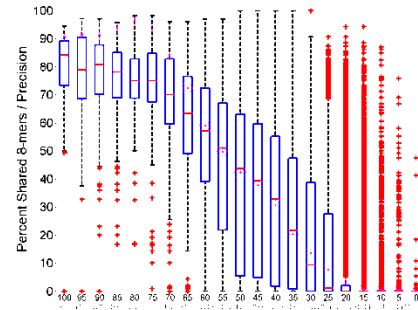
Zscore40



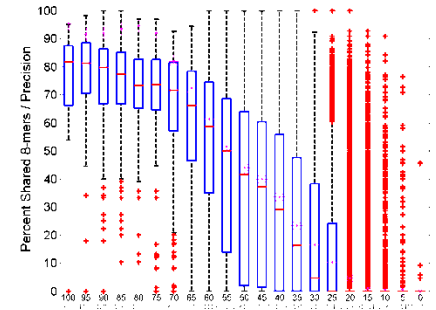
Escore45



Escore46

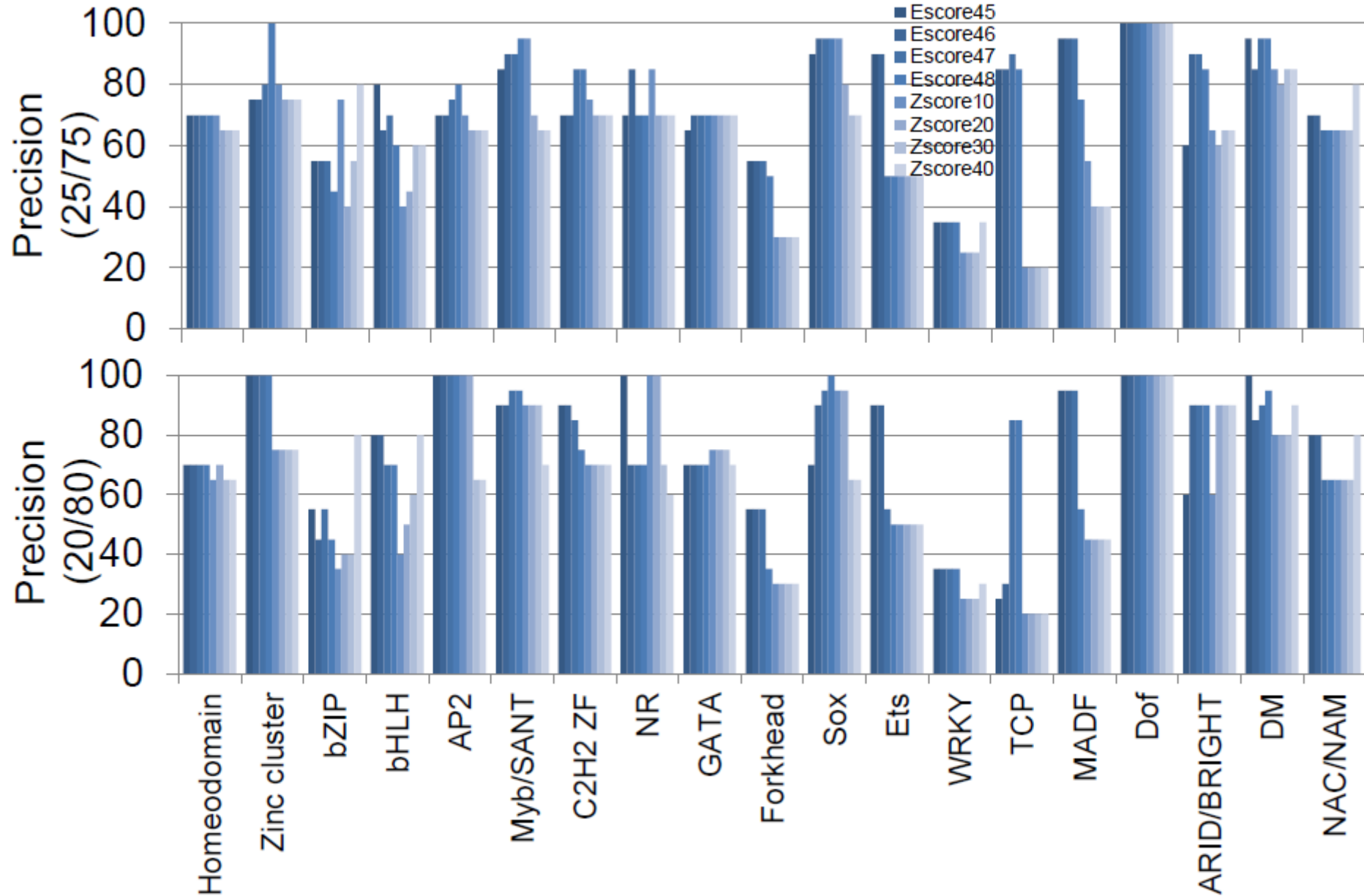


Escore47

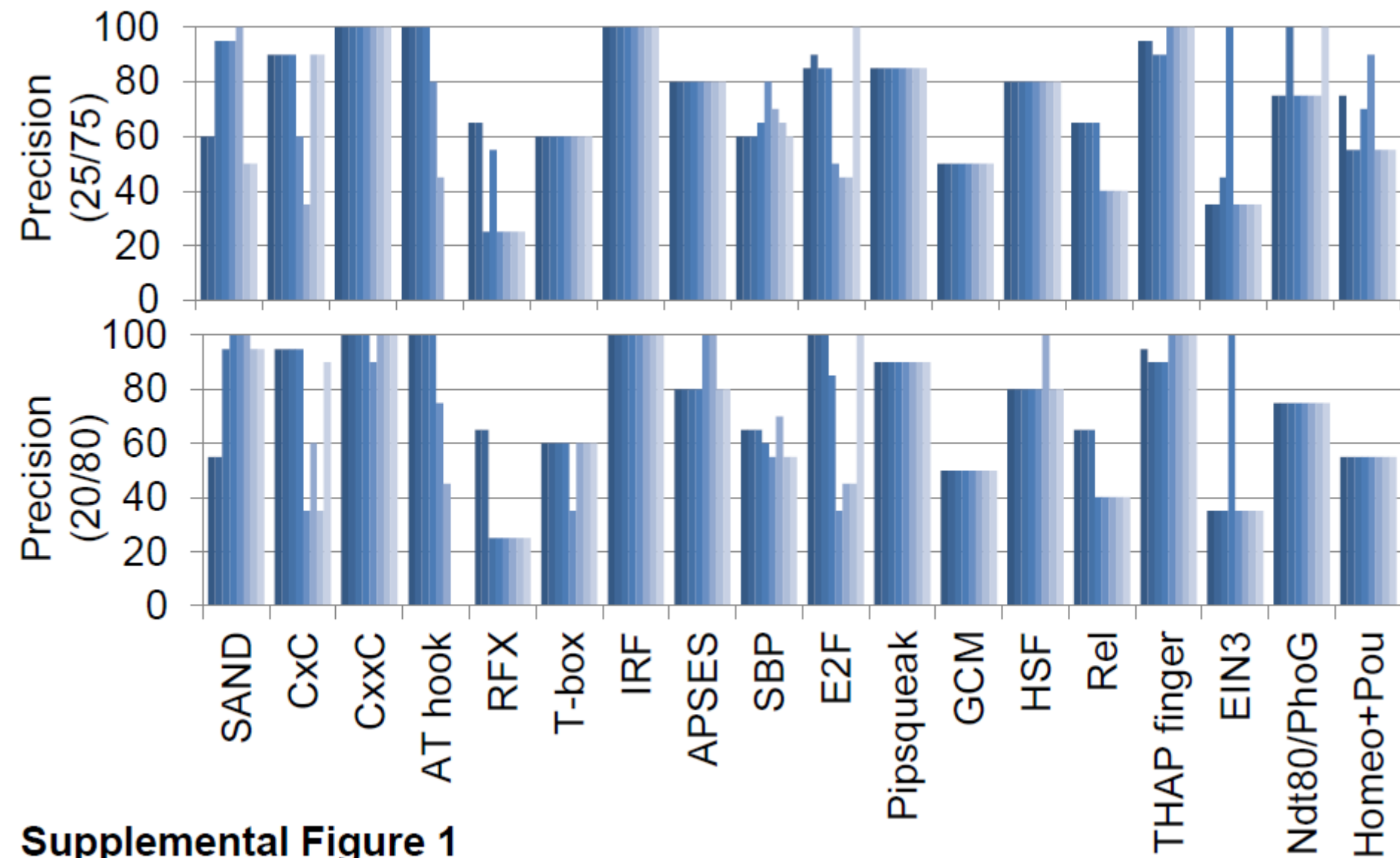


Escore48

# Cutoffs across families and parameters

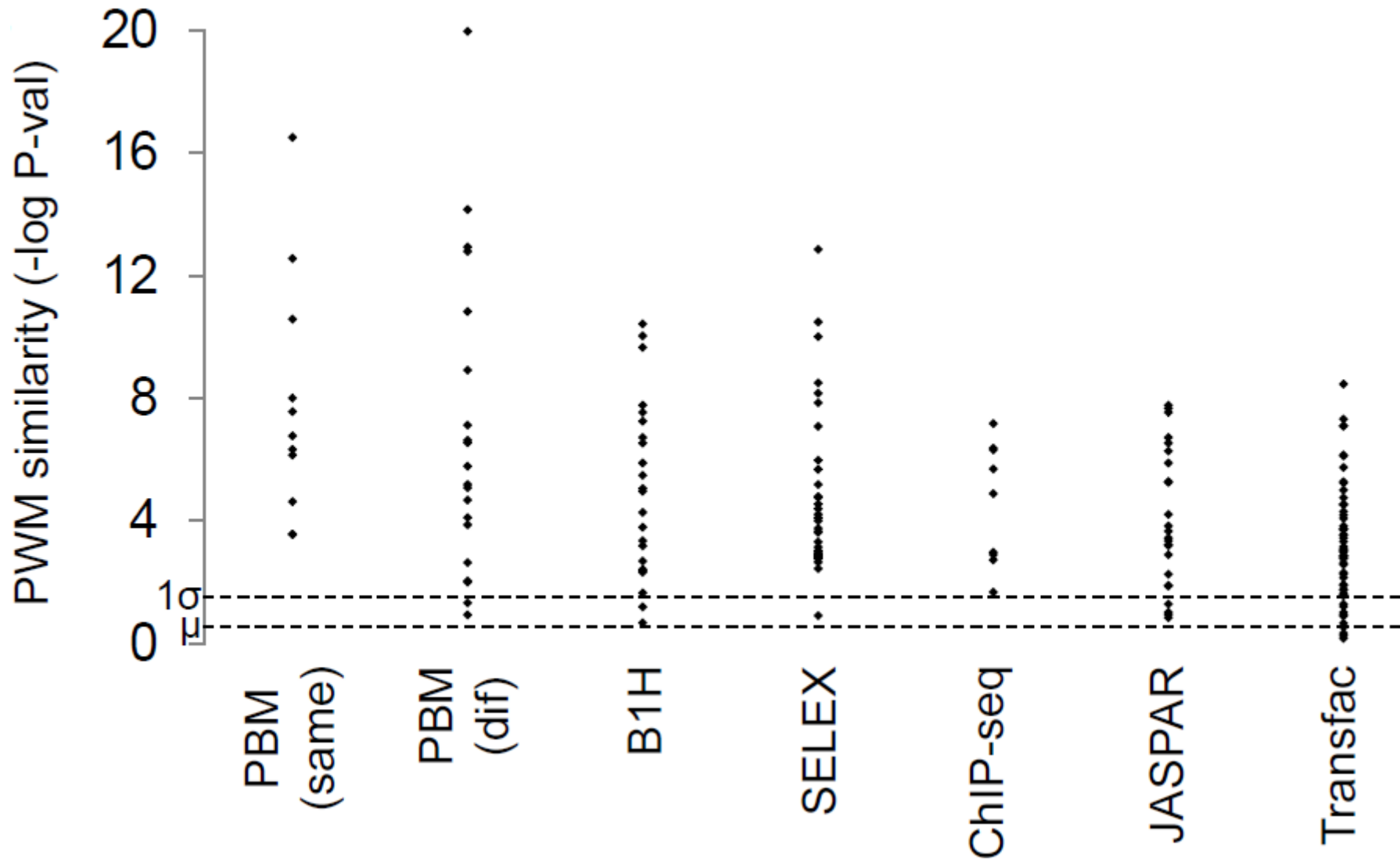


# Cutoffs across families and parameters

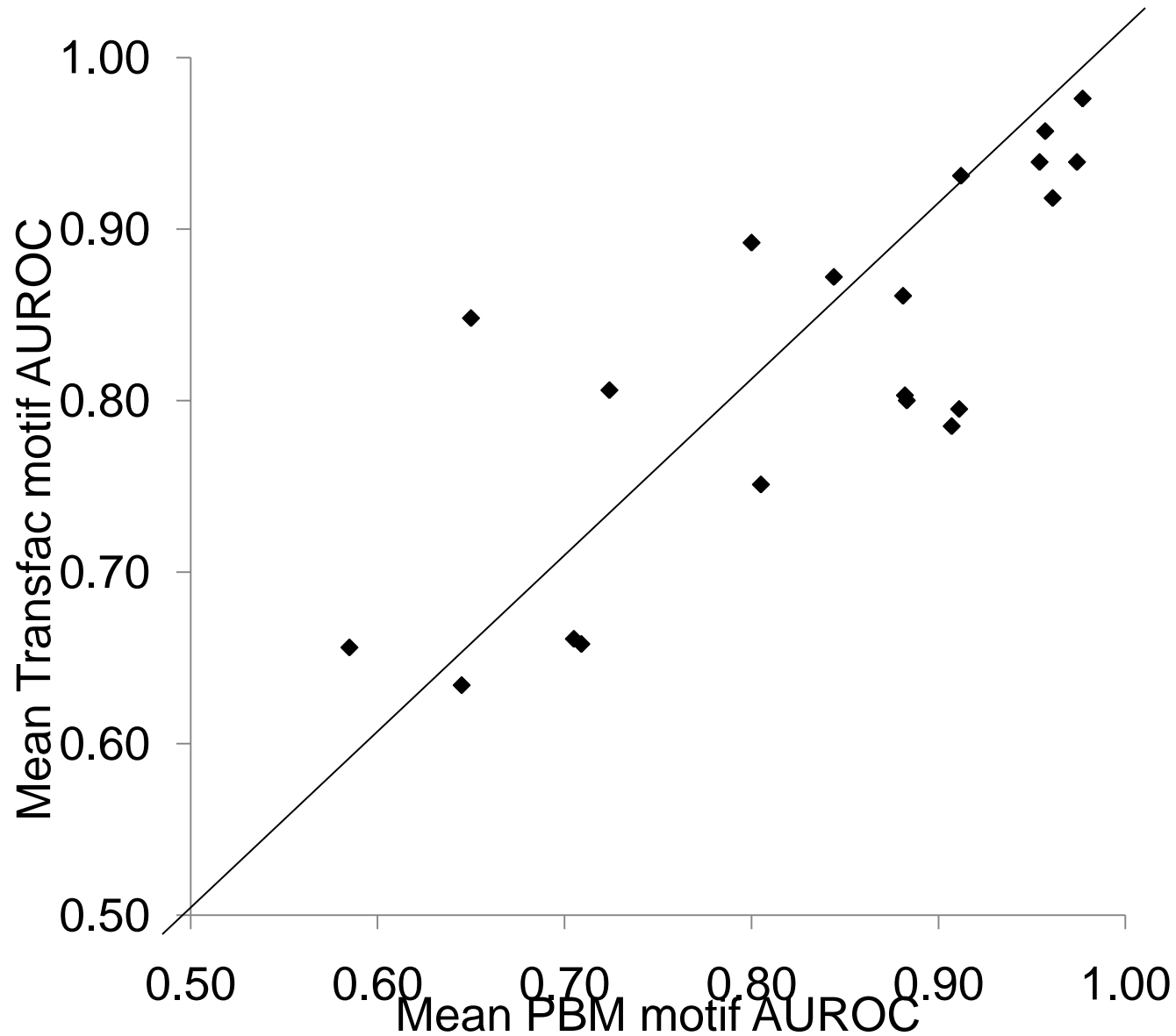


Supplemental Figure 1

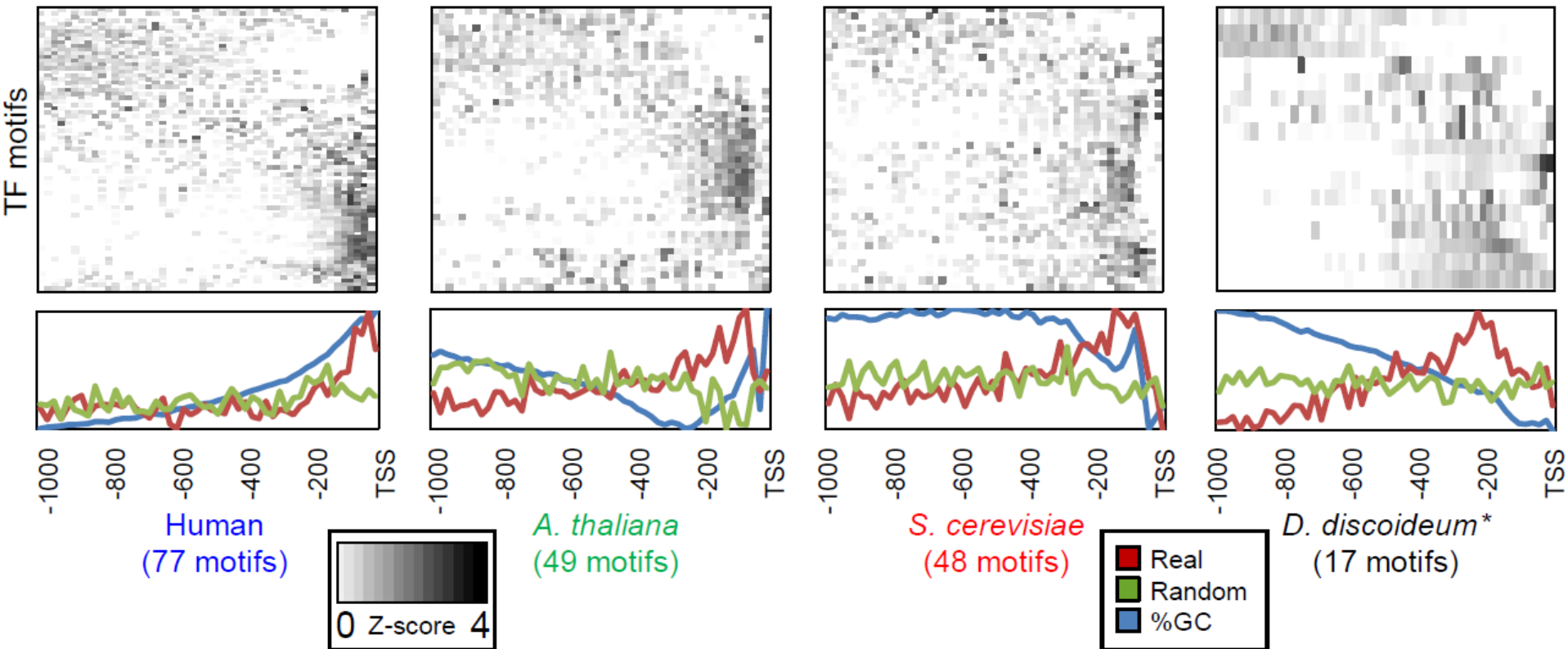
# Our motifs vs. other studies



# PBM-derived motifs predict ChIP peaks as accurately as literature-derived motifs



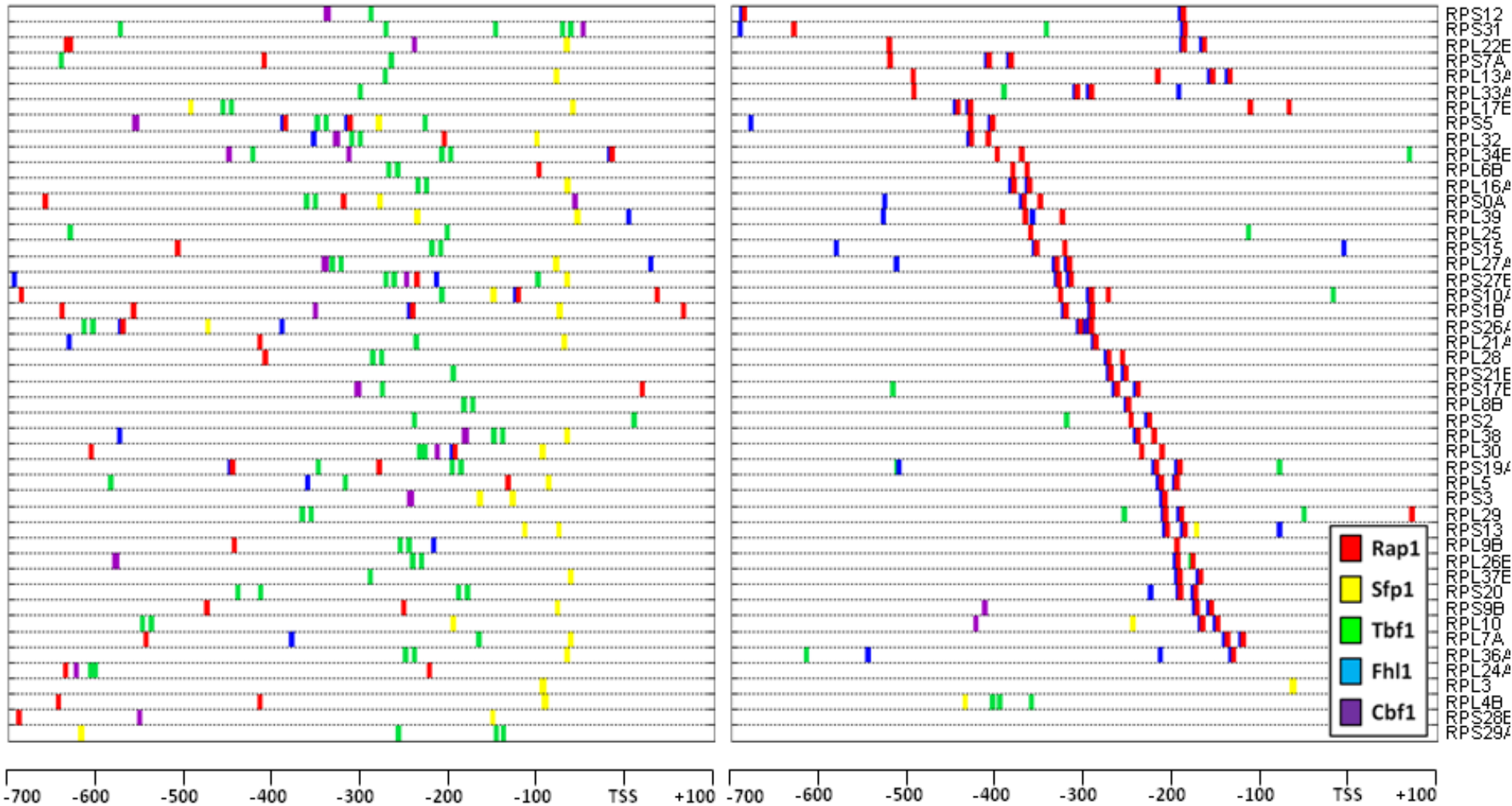
# Positional bias of (non-redundant) motifs in eukaryotic promoters



# Completely different in another yeast

*C. albicans*

*S. cerevisiae*



**Major TF controller switched from Tbf1 to Rap1**