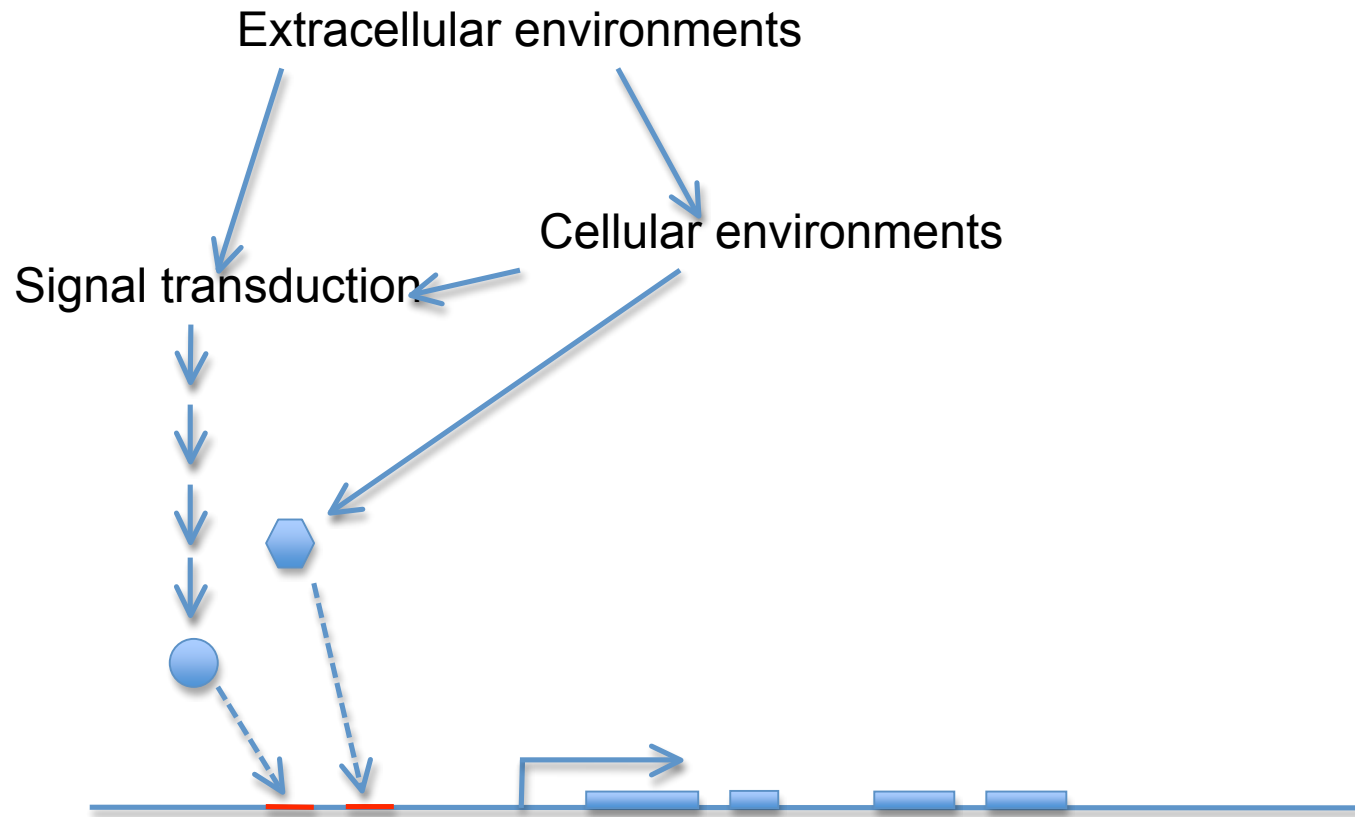


Computational Identification of Cis-
regulatory Elements Associated with
Groups of Functionally Related Genes in
Saccharomyces cerevisiae

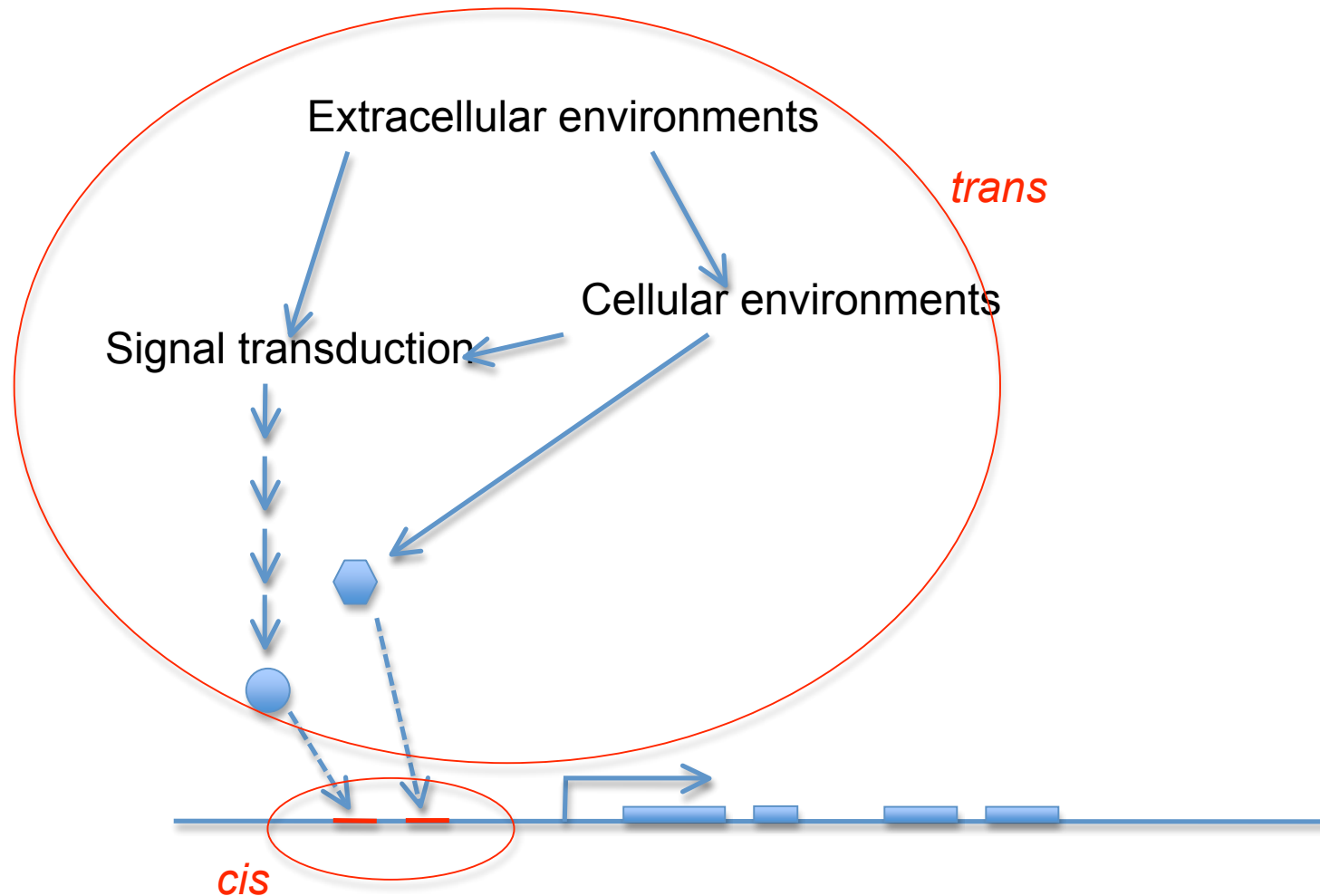
Jason D. Hugher, Preston W. Estep,
Saeed Tavazoie, and George M. Church

Why study transcription factor binding sites (TFBSs)?

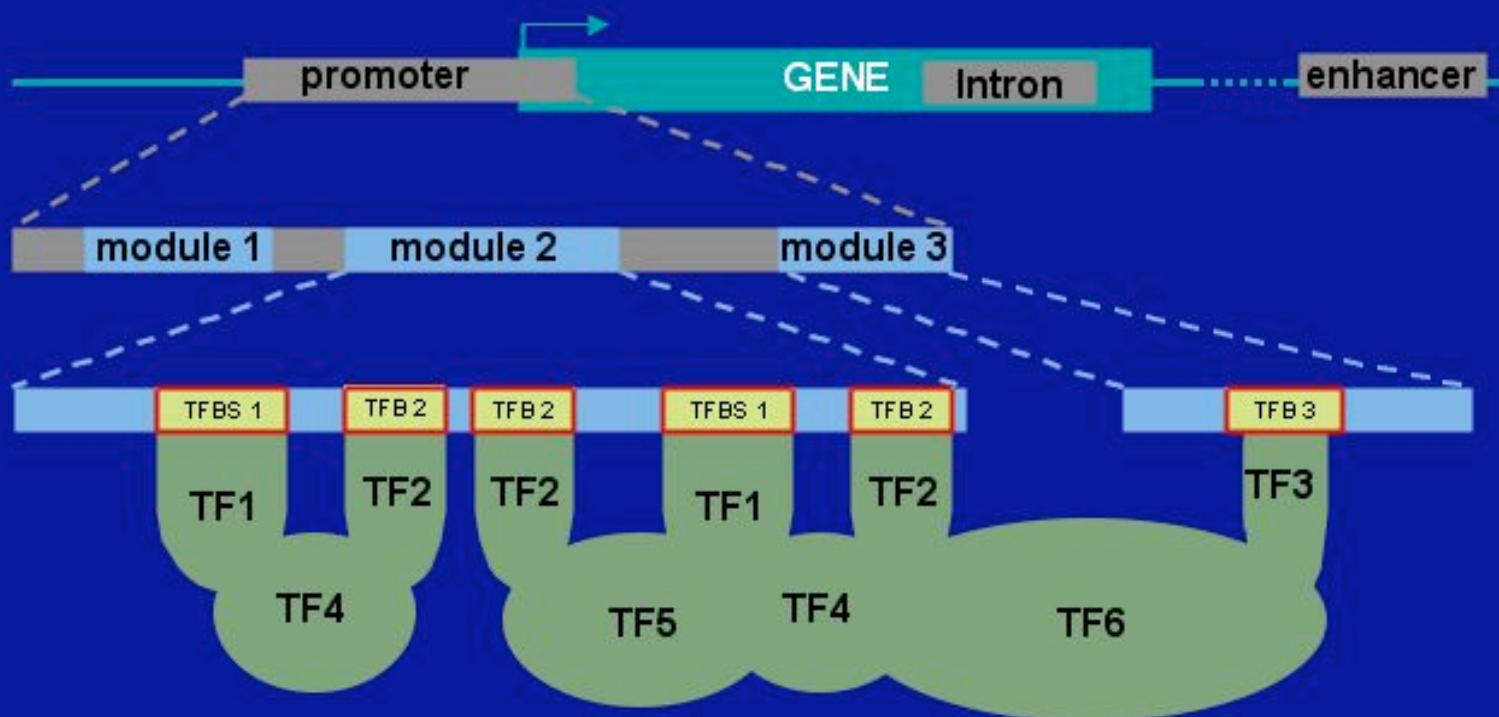
- All protein coding genes are regulated by transcription factors.



- All protein coding genes are regulated by transcription factors.



Gene expression regulatory elements

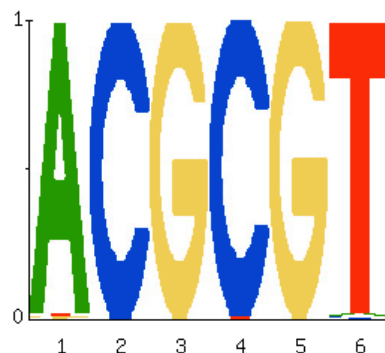


Transcription factors (TF) bind to regulatory motifs

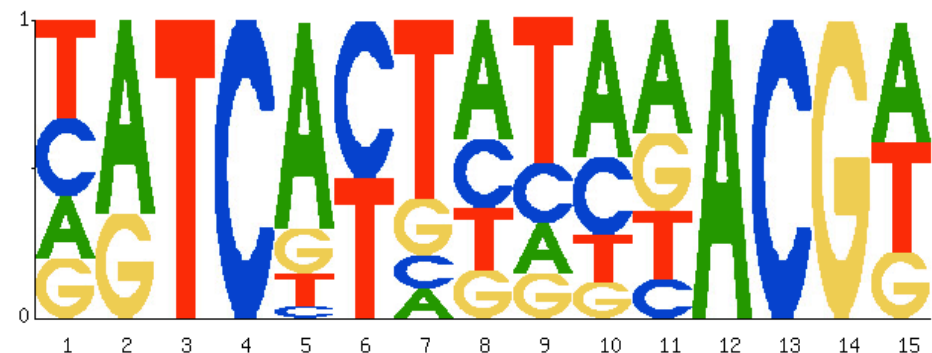
binding sites = binding motifs = regulatory elements = regulatory motifs

The difficulties

- Transcription factor binding sites (TFBSs) are **short** and **degenerated**.
 - Usually 6 ~ 15 bp.
- TFBSs can be on **either** strand of the DNA

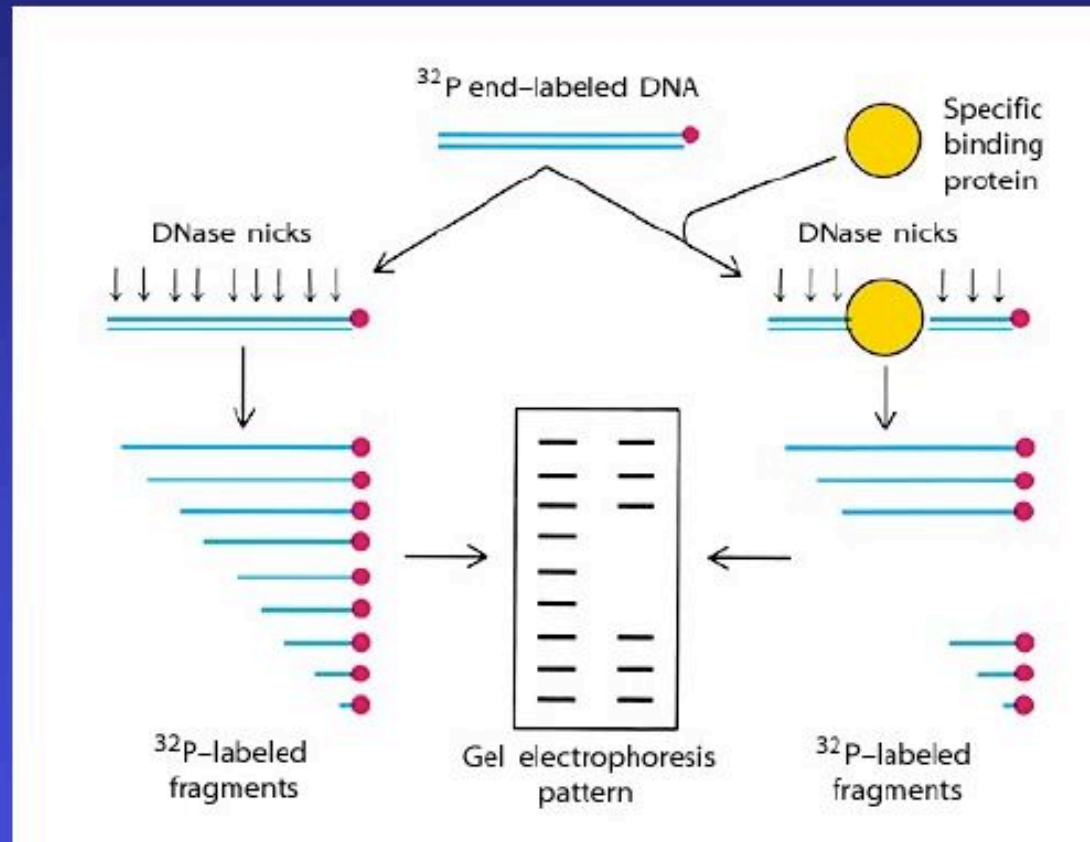


SWI6



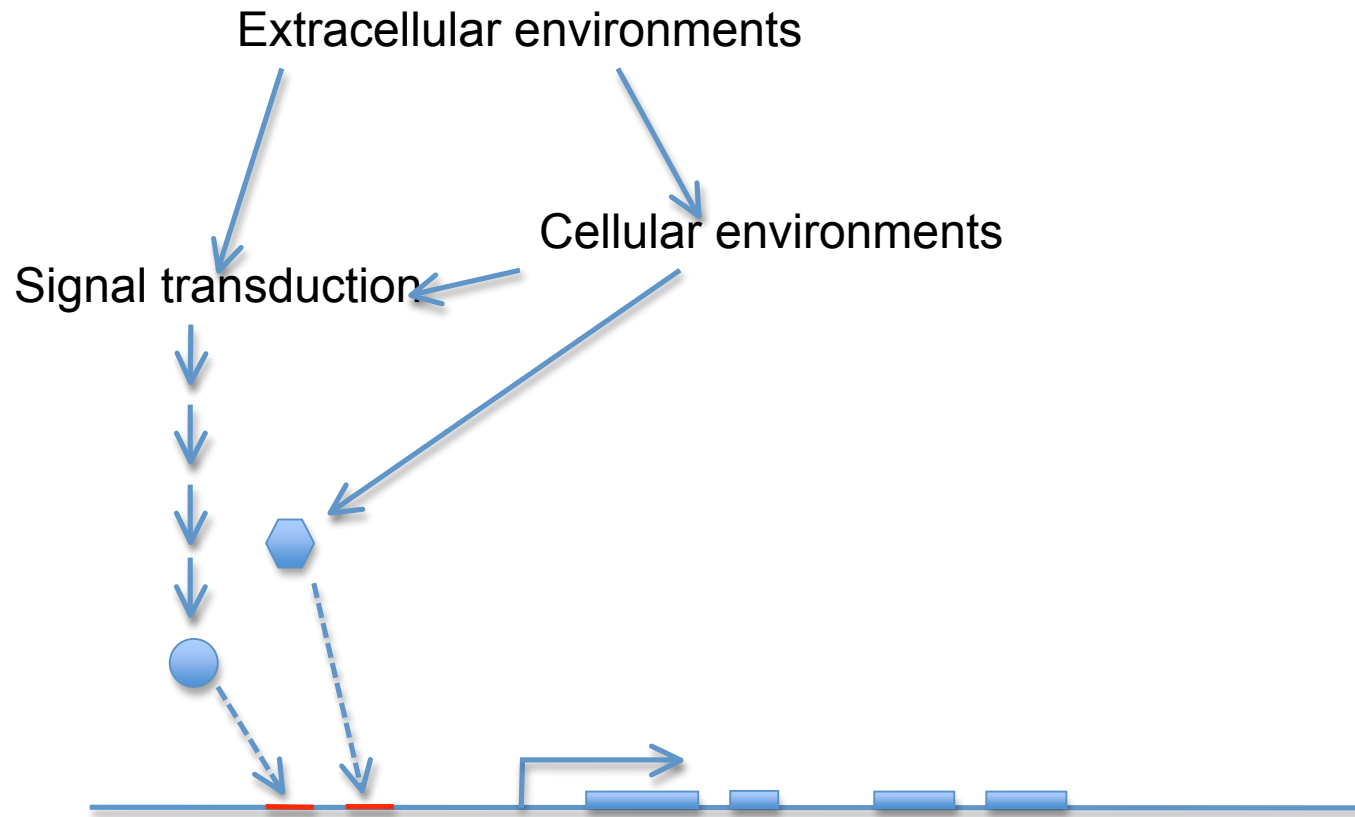
ABF1

Principle of DNA Footprinting



- In yeast, *Saccharomyces cerevisiae*:
 - ~6000 genes
 - >200 TFs

- All protein coding genes are regulated by transcription factors.



Basic idea

- Genes with the same (or similar) functions could be regulated by the same TF, therefore, have the same TFBSs at their promoter regions.

To identify where are TFBSs, we need to know:

- 1. What sequences they are
 - Find motifs: AlignACE
- 2. Which TFs will bind to them
 - Assign AlignACE motifs to known *cis*-regulatory elements: CompareACE

- AlignACE
 - Gibbs sampling
- CompareACE
 - Pearson correlation coefficient between the nucleotide base frequencies of two motif alignments

Maximum a priori log likelihood (MAP) score:

$$N \log R$$

N: number of aligned sites

R: degree of over-representation

- Group specificity

$$S = \sum_{i=x}^{\min(s_1, s_2)} \frac{\binom{s_1}{i} \binom{N-s_1}{s_2-i}}{\binom{N}{s_2}}$$

N: the total number of ORFs

s1, s2: the numbers of ORFs in the group used to find the motif and in the list of target gene, respectively

X: the number of ORFs in the intersection of the two lists

- Position bias

The probability of observing m or more sites out of possible t in a 50bp window (w) of a 600bp region (s)

$$P = \sum_{i=m}^t \binom{t}{i} \left(\frac{w}{s}\right)^i \left(1 - \frac{w}{s}\right)^{t-i}$$

Result 1

- Motifs ranked by specificity score

Cluster	MAP	Spec	PosBias	PrefPos	Logo	Notes
S1	123	8.6×10^{-48}	4.03×10^{-6}	279	ATGTATGG	Rap1
S2	128	1.27×10^{-32}	3.02×10^{-10}	125	GGTGGCAAA	Rpn4
S3	31.1	5.29×10^{-23}	0.0163	270	AAATGAGTCA	Gcn4
S4	55	8.73×10^{-21}	0.0015	169	GAA TTG GAA	HSE
S5	26.8	3.3×10^{-18}	0.00305	557	TG GGGGTA	Mig1
S6	30.9	5.36×10^{-18}	0.00159	246	AAT TCACGTG	Cbf1
S7	12.4	9.98×10^{-14}	0.000546	241	TG CTGG TG	RPS genes
S8	15.9	2.32×10^{-13}	0.0113	558	TGATTGG	Hap2,3,4
S9	25.3	6.13×10^{-13}	1.49×10^{-5}	191	ATTC GGAA	Lys14
S10	39.3	1.43×10^{-12}	4.85×10^{-8}	124	ACGCGT	MCB
S11	11.1	1.44×10^{-12}	0.00527	168	GCC G GAG	RPL genes
S12	102	2.34×10^{-12}	2.04×10^{-43}	123	TGAAAAATTT	RRPE
S13	18.7	2.37×10^{-12}	0.0447	316	GGG GGG A	STRE
S14	20.6	2.6×10^{-12}	0.00213	171	AACTGTGG	Met31,32
S15	19.4	3.08×10^{-12}	0.00269	153	CCG ACGG	Leu3
S16	17.2	3.08×10^{-12}	0.0204	105	TCGG GTTA	Oaf1
S17	14.1	4.17×10^{-12}	0.000374	201	GG CGGG G	mito. transp.
S18	21	1.14×10^{-11}	0.000413	344	G GGG G GG	stress
S19	24.6	1.57×10^{-11}	9.29×10^{-6}	308	GA GGA	CSRE
S20	16	2.07×10^{-11}	8.66×10^{-5}	388	CGGGGG	TCA cycle
S21	21.2	2.2×10^{-11}	0.000205	172	AA G TAAACA	cytoskel. transp.
S22	20.3	2.27×10^{-11}	0.00241	374	CGGAG TCGG	TCA cycle
S23	10.5	2.49×10^{-11}	0.0221	349	SCA GTGGA	Pho4
S24	10.7	4.27×10^{-11}	0.000847	161	TGAAACA IT	Ste12
S25	44.8	8.64×10^{-11}	0.00647	382	TCGGGA G	Pdr3

Cluster	MAP	Spec	PosBias	PrefPos	Logo	Notes
S1	123	8.6×10^{-48}	4.03×10^{-6}	279	▲ ATGT▲I▲GG▲	Rap1
S2	128	1.27×10^{-32}	3.02×10^{-10}	125	GGTGGCAAA▲	Rpn4
S3	31.1	5.29×10^{-23}	0.0163	270	▲▲▲I▲GA▲TCA	Gcn4
S4	55	8.73×10^{-21}	0.0015	169	G▲A_ TTc_▲GAA	HSE
S5	26.8	3.3×10^{-18}	0.00305	557	▲G_ GGGG_▲_▲▲	Mig1
S6	30.9	5.36×10^{-18}	0.00159	246	▲▲T_ TCACGTG	Cbf1
S7	12.4	9.98×10^{-14}	0.000546	241	▲G_ ▲CTGG_ ▲G	RPS genes
S8	15.9	2.32×10^{-13}	0.0113	558	TGATTGG_▲▲▲	Hap2,3,4
S9	25.3	6.13×10^{-13}	1.49×10^{-5}	191	▲TTC_ _ GGAA▲	Lys14
S10	39.3	1.43×10^{-12}	4.85×10^{-8}	124	▲▲ACGCGT_▲	MCB
S11	11.1	1.44×10^{-12}	0.00527	168	▲G_ GCG_ G_ _▲▲GA	RPL genes
S12	102	2.34×10^{-12}	2.04×10^{-43}	123	TGAAAA▲T TT	RRPE
S13	18.7	2.37×10^{-12}	0.0447	316	GcGc_▲G_ _A	STRE
S14	20.6	2.6×10^{-12}	0.00213	171	AAcTGTGGc_▲	Met31,32
S15	19.4	3.08×10^{-12}	0.00269	153	cC_▲_▲cCGGc	Leu3
S16	17.2	3.08×10^{-12}	0.0204	105	TcGGc_▲IT_▲A	Oaf1
S17	14.1	4.17×10^{-12}	0.000374	201	▲▲_ CGGG_ _▲_ G▲	mito. transp.
S18	21	1.14×10^{-11}	0.000413	344	G_ ▲G_ G_ ▲G_ GG	stress
S19	24.6	1.57×10^{-11}	9.29×10^{-6}	308	▲▲G_▲_▲_ cGG▲	CSRE
S20	16	2.07×10^{-11}	8.66×10^{-5}	388	c_ C_ G_ G_ ▲_ G▲	TCA cycle
S21	21.2	2.2×10^{-11}	0.000205	172	A_ AA_ G_ TAAACA	cytoskel. transp.
S22	20.3	2.27×10^{-11}	0.00241	374	CcGAG_ ITcGG	TCA cycle
S23	10.5	2.49×10^{-11}	0.0221	349	▲c▲_ GTGcGA	Pho4
S24	10.7	4.27×10^{-11}	0.000847	161	T▲AAACA_ _T_ _▲	Ste12
S25	44.8	8.64×10^{-11}	0.00647	382	▲TCCGcG_▲_ G	Pdr3

Literature: STRE binding site: **AGGGG**

Result 2

- Positionally biased motifs

Cluster	MAP	Spec	PosBias	PrefPos	Logo	Notes
P1	637	0.000114	1.55×10^{-166}	140	AA₂AA₂AAAA	
P2	15.6	0.0559	2.8×10^{-48}	114	A₂T₂TATATA	AT repeat
P3	16.4	0.109	8.09×10^{-48}	133	AA₂ A A₂ AA₂ AA₂	
P4	102	2.34×10^{-12}	2.04×10^{-43}	123	TGAAAA₂TT	RRPE
P5	14.3	0.025	2.29×10^{-41}	148	ATCA₂ A₂CG₂	Abf1
P6	23.8	0.0906	1.58×10^{-35}	130	A₂ A A₂ AA₂ A A₂ AA₂A	
P7	29.6	0.00888	2.79×10^{-33}	101	G₂GATGAG₂T	PAC
P8	17.6	0.00102	3.53×10^{-31}	148	CGGGTAA₂	Reb1
P9	10.2	0.0111	5.66×10^{-21}	35	GTGTG₂GTGT	GT repeat
P10	32.6	1.98×10^{-14}	5.66×10^{-21}	35	GT TGGGT₂	Rap1
P11	125	4.08×10^{-28}	1.13×10^{-14}	112	GGTGGCAAA₂	Rpn4
P12	12.5	0.00818	6.6×10^{-11}	123	AAA₂ T₂ A₂ AAA₂	
P13	13.3	7.5×10^{-6}	7.01×10^{-10}	249	AA₂ A TAA₂ ATA₂ A₂	
P14	19.2	0.028	7.45×10^{-10}	114	AA₂GC₂AAAA	
P15	10.5	9.72×10^{-5}	5×10^{-9}	141	G₂CGACGCG₂A	MCB
P16	11.8	0.000216	5×10^{-9}	37	GAGAAA₂AA	
P17	20.2	0.00495	9.18×10^{-9}	127	T₂TTGAAAA	

Cluster	MAP	Spec	PosBias	PrefPos	Logo	Notes
P1	637	0.000114	1.55×10^{-166}	140	AA _Δ AA _Δ AAAA	
P2	15.6	0.0559	2.8×10^{-48}	114	<u>A</u> _Δ <u>T</u> _Δ TATATA	AT repeat
P3	16.4	0.109	8.09×10^{-48}	133	AA _Δ A A _Δ AA _Δ AA _Δ	
P4	102	2.34×10^{-12}	2.04×10^{-43}	123	<u>T</u>GAAAA<u>A</u><u>T</u><u>T</u>	RRPE
P5	14.3	0.025	2.29×10^{-41}	148	ATCA _Δ <u>T</u> A _Δ G _Δ	Abf1
P6	23.8	0.0906	1.58×10^{-35}	130	A _Δ A _Δ A _Δ AA _Δ A A _Δ AA _Δ A	
P7	29.6	0.00888	2.79×10^{-33}	101	<u>G</u><u>A</u>GATGA<u>G</u><u>T</u>	PAC
P8	17.6	0.00102	3.53×10^{-31}	148	CGGGTAA _Δ _Δ _Δ	Reb1
P9	10.2	0.0111	5.66×10^{-21}	35	GTGTG _Δ GTGT	GT repeat
P10	32.6	1.98×10^{-14}	5.66×10^{-21}	35	GT TGGGT _Δ _Δ	Rap1
P11	125	4.08×10^{-28}	1.13×10^{-14}	112	GGTGGCAAA _Δ	Rpn4
P12	12.5	0.00818	6.6×10^{-11}	123	AAA _Δ T _Δ A _Δ AAA	
P13	13.3	7.5×10^{-6}	7.01×10^{-10}	249	AA _Δ A TAA _Δ AT _Δ _Δ A	
P14	19.2	0.028	7.45×10^{-10}	114	AA _Δ GC _Δ AAAA	
P15	10.5	9.72×10^{-5}	5×10^{-9}	141	G _Δ _Δ ACGCG _Δ _Δ A	MCB
P16	11.8	0.000216	5×10^{-9}	37	GAGAAA _Δ _Δ AA	
P17	20.2	0.00495	9.18×10^{-9}	127	T _Δ TTGAAAAA	

Negative controls

Spec. score cutoff	Random groupings motifs		Functional category motifs	
	MAP > 0	MAP > 10	MAP > 0	MAP > 10
1	3692(1063)	1792(205)	3311(1324)	1234(208)
10 ⁻¹	2766(1038)	1047(202)	2713(1284)	815(194)
10 ⁻²	2026(978)	553(181)	2198(1201)	530(179)
10 ⁻³	1416(831)	285(149)	1622(1016)	337(153)
10 ⁻⁴	935(641)	151(104)	1109(753)	226(121)
10 ⁻⁵	554(425)	72(56)	750(543)	160(90)
10 ⁻⁶	329(290)	31(29)	446(329)	122(67)
10 ⁻⁷	151(143)	15(15)	270(199)	91(47)
10 ⁻⁸	60(59)	9(9)	164(118)	73(35)
10 ⁻⁹	37(36)	6(6)	97(62)	60(28)
10 ⁻¹⁰	14(14)	5(5)	69(38)	54(25)

Positive controls

- Search 29 literature reported motifs
 - 21 of 29 were found

Conclusion

The method (AlignACE + group specificity + position bias) can be applicable to groups of genes other than functional categories.

After that ...

- Microarray
- Chromatin Immunoprecipitation chip
(ChIP-chip)

ChIP-chip

Transcriptional regulatory code of a eukaryotic genome

**Christopher T. Harbison^{1,2*}, D. Benjamin Gordon^{1*}, Tong Ihn Lee¹,
Nicola J. Rinaldi^{1,2}, Kenzie D. Macisaac³, Timothy W. Danford³,
Nancy M. Hannett¹, Jean-Bosco Tagne¹, David B. Reynolds¹, Jane Yoo¹,
Ezra G. Jennings¹, Julia Zeitlinger¹, Dmitry K. Pokholok¹,
Manolis Kellis^{1,3,4}, P. Alex Rolfe³, Ken T. Takusagawa³, Eric S. Lander^{1,2,4},
David K. Gifford^{3,4}, Ernest Fraenkel^{1,3} & Richard A. Young^{1,2,4}**

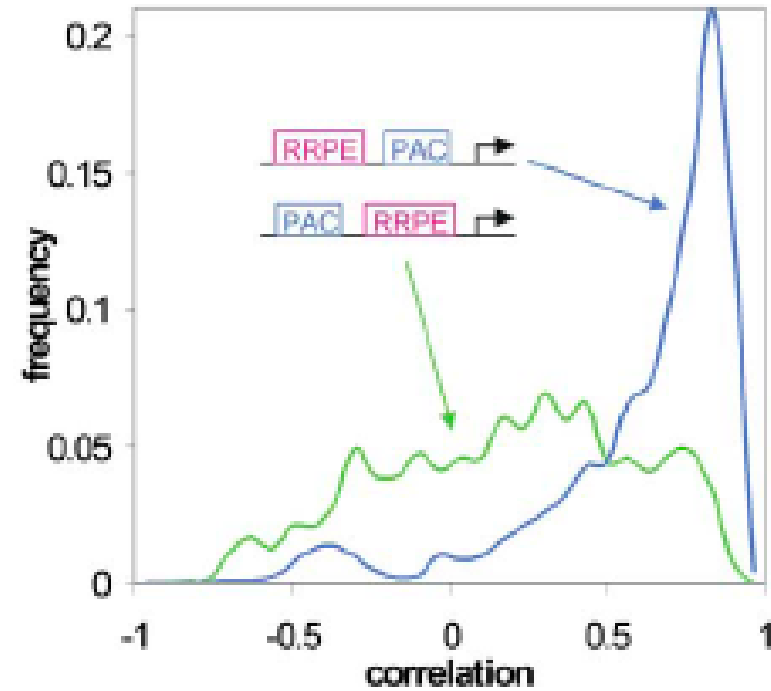
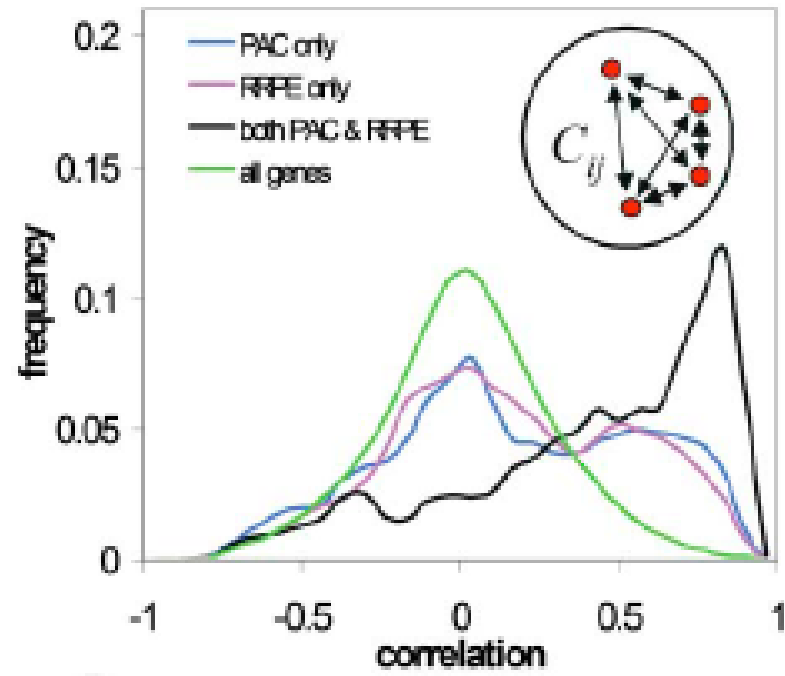
Position bias

Cell, Vol. 117, 185–198, April 16, 2004, Copyright ©2004 by Cell Press

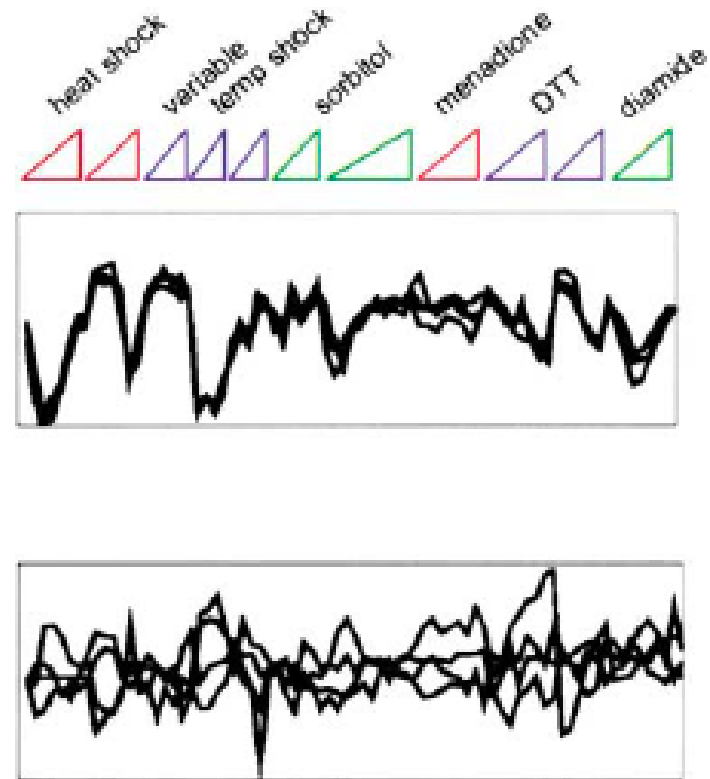
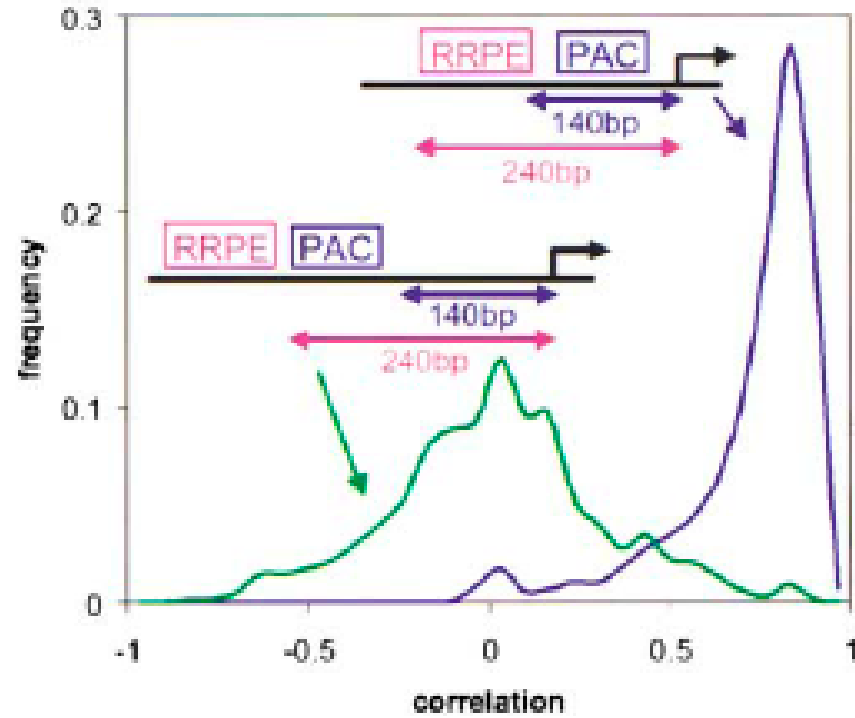
Predicting Gene Expression from Sequence

Michael A. Beer and **Saeed Tavazoie***
Lewis-Sigler Institute for Integrative Genomics
and Department of Molecular Biology
Princeton University
Princeton, New Jersey 08544

Cis-rule



Cis-rule



Dealing with gapped motifs

Discovering gapped binding sites of yeast transcription factors

Chien-Yu Chen*, Huai-Kuang Tsai†, Chen-Ming Hsu‡, Mei-Ju May Chen§, Hao-Geng Hung¶, Grace Tzu-Wei Huang||, and Wen-Hsiung Li|***††

- Questions