

Table 1: Motif comparison for 30 TFs with six different programs. Table shows the motifs learned by various algorithms used by Harbison *et al.* and those learned by Our algorithm. For comparison, we use the motifs with the top MAP score for AlignACE [6], MEME [1], MDscan [4], MEME.c (modified MEME with conservation information), method by [3], and CONVERGE (a method by Harbison *et al.*), the latter three being based on conservation information) In the ninth column we report the top motif according to our score. We also report the predicted class and the percentage of entries in *c* contributing to that class. The last column is the literature consensus as used by Harbison *et al.* collected from YPD, SCPD, and TRANSFAC databases at the time the paper was published. The bold sections in the motifs indicate either a match with the literature consensus in the final column or to a motif we found in the literature search we conducted. In cases where the match is not obvious, it is probably because the reverse complement of the sequence matches the literature consensus. Lower case letters in the motifs indicate a weaker preference (less information content at that position). Ambiguity codes: S=C/G, W=A/T, R=A/G, Y=C/T, M=A/C, K=G/T, and ‘.=’= A/C/G/T.

TF	Environ cond.#	Harbison <i>et al.</i>							PRIORITY		Literature
		AlignACE Top MAP	MEME Top MAP	MDscan Top MAP	MEME.c Top MAP	Kellis Top score	CONVERGE Top enrichment ^{&}	Top MAP	Predicted class		
Basic Leucine Zipper TFs											
Arr1	YPD	R.AmA.a.A.A.AmA.A	cAmAcACmcAmAmayrCA	CACACACAC	TkTGtgTkttTstwT	GGA.....GTG	AAACAmAsAmACAMA	YAAACaca	fork	78%	TTACTAA
Cad1	SM YPD	AArarARAAA..A	AAGrAArArAgR	ATTA.TmA	.ATTAGTMAkCa	TGC..GGG	sAAAArKTAtTCt.A	cTKACWAA	bZip	76%	TTACTAA
		GtGTGTGkGTGTG	GCTKACTAAT.	GKGTGTGK	GCTkACTAAT	GCC.....TAT	t..yGm Tta.tAA.	GCTTACTA	bZip	73%	
Cin5	H2O2Hi H2O2Lo	mR.RAAARAAR	YYYYT.ctTTTykysT.	TTTttTTTT	TkTtTt.TyTtktkyyTT	CTGCTTA	gcTTA..TmAk....	CGGTG...	bHLH	83%	TTAC[r]TAA*
		AARAAAAA..AA.A	TTTYyytytTy.ytyYYK	.GSGssgG	AAAa..AAraaraAAA	TAC...AGC	.gmTTAyrtAA....	TTAyGTAA	bZip	94%	
	YPD	RRARAARAAA	GGC.cgGCCGS	TTAYrTAA	—	GGG.....CTC	aTka.rTaAgc....	.gcTGG..	bHLH	92%	
Cst6	YPD	A.A..rAAA.A..A..a.A	rmAtk.mAwRCRAAaA	AgTY.AsT	rmAtk.mAwRCRAAaA	TTG....ACA	TCCyAkTAYTCTkCA	.ACTGGAC	bZip	80%	—
Gcn4	RAPA SM	RAAAAARAAr	TGAGTCA	.TGAsTCA	.TGAGTCAY	rTGAsTCATgAsTCaY	ATGACTCA	bZip	95%	ArTGACTCw
	YPD	rAAAAAArAAa	yTyTyyTyyTyyTTc	TGAGTCA	.TGAsTCat	.rTGACTCwT	rTGACTca	ATGACTCA	bZip	96%	
		rAAAAAArAAA	.TGAGTCAY	.TGAsTCA	.TGAGTCAY	rTGACTCA.	.rTGAsTCA.....	.TGACTCA	bZip	91%	
Hac1	YPD	A..rAA..MAAARA	TrCSTSkkccwywtmM	TAcGTGkC	TrCSTSkkccwywtmM	CTC.....CAT	amAAyGaksAgCAA	ACGTGTCA	bZip	76%	kGmCA[G]CGTGTc **
Met28	YPD	a.A.A....AAAA.AAA	TkyTTTkskskscTTw	ATrTayAT	CGGCCCTC	ATA.CTT	mAmASsrYGkCTrTr	SKAACY	bZip	75%	—
Met4	SM	AA.AR.RARAAA	AArAAMmmRmAA	TATATATAT	ASTGTGgySGyS	TCT.....TAC	MMAAsTGTGgySsCr	TGTACGT	bHLH	80%	—
Sko1	YPD	mAAA.RARr.AA	TkTTkyyykTTTkyKKCk	sSgtacSs	kwTyTTkTyyTyktTTkT	CTG.....AGG	ATGACGTr	.tACGTCA	bZip	72%	ACGTCA
Yap1	H2O2Lo	AwMArrAAr..A	ssTTTyCrT	TTA.TAAk	TTAGTmAGC	CCG..TTT	cTkACtAA	TKACTAAw	bZip	87%	TTAsTAA
	YPD	ArYAAAArAR	yTyy.TkC.TKkyyytmt	CAC.CasA	TTT..tTtt.Tt.tktTT	TAGTCAGC	aTTAGTmA	.YAAAC.	fork	87%	
Yap3	YPD	A.A...A.A.A.A.A.Amr	T.kyttcTT.mTTkTT	CACACACAC	gkkkkCkykmTtrT..	TTC.....AAC	rAAArscATGTATyT	.CTAAaTS	bZip	65%	TTACTAA
Yap5	YPD	GTGTG.GTGTG	GTGCAgsAgAAcsAgGAt	CACACmCAC	ATGTAyggrtg	TTG.CGG	yGmrsAGmAcAsAgRa	CGTCKGYG	bHLH	93%	TTACTAA
Yap6	H2O2Hi H2O2Lo YPD	ARRRAAA.AAR	s.TtTyytykYc.yt.YT	cySggcCG	ssGgCsGAgcss.s.m	TTT.TCC	wAtmTGKwCCTrYGs	CGTGGG.g	bHLH	87%	TTACTAA
		RARARRA.AAA	A.r.ArAArrrrrAAA	scyGCRgs	yTTTtCtTyTTtytCsTt	ATC.CTT	ARartTtA	..CGTGG.	bHLH	91%	
		RArAArAAAAA	magAAA..rrrArArR	smYGCAks				..CGTGG.	bHLH	91%	
Yap7	H2O2Hi	AAAARRAAAR	mTTAsTmAkc..	TKAsTMAk	raA.aArrAAAarGAAAA	GAC.AAG	AAcyYYra	TKAsTAAk	bZip	82%	TTACTAA
	H2O2Lo	RRAAA.ARAAR	TtT..yyttTctyTTyyT	TTAsTMAk	Gc TKAsTAA	TTA...AGC	.y..TkAsTaAk...	ATKAsTAA	bZip	92%	
Forkhead TFs											
Fhl1	RAPA SM YPD	RTGTayGGrTg	.t.tacAyCCrTACAYyy	yCCrTACA	rrrrTGtayGGrTgtA.a	TGT.TGGrT	rrTGYayGGrTgtw.	TGTAYGG	other	91%	—
		RTGTayGGrTg	.trcAyCCrTACAYyy	GTayGGrT	rrrrTGtayGGrTgtA.a	TGT.TGGrT	CAyCCrTACAYyy..	TGTAYGG	other	94%	
		RTGTayGGrTG	.t.taCAyCCrTACAYyy	TGyryGGr	rrrrTGtayGGrTgtA.a	TGT.TGGrT	rrRTGTayGGrTgtw	TGTAYGG	other	94%	
Fkh1	YPD	RAR.ARA.RAA.A	aaa.rtAAACaa..r.a	tTGTtTAC	.k...tTGTtTAC.tTtk	.TTGTTTAC	.raa.gTAAACaa..	GTAAACAA	fork	95%	GGTAAACAA
Fkh2	H2O2Hi	RRaARR.AAA.R	TTtTyc.TTYyyt.y.TTT	T.TKTmCm	gg.aAAa. GTAAACAr	TTGTTTAC	gg.aAAa.GTAAACa	GTAAACAA	fork	92%	GGTAAACAA
Fkh2	H2O2Lo YPD	ArMArr.AAR.A	gg.AAAA. GTMAACA.ya	S.AgAgAG	gG.rAaw. GTMAACA	ACG.....ATT	gG.AAA.GTMAACA	GTAAACAA	fork	80%	
Fkh2		rRaAR.AAA.R	AArrr.rAAAa.r.AAA	GtAAACAA	g.rAaa.gTAAACaa.r.	TTGTTTAC..TT	.raaa.gTAAACaa.	GTAAACAA	fork	94%	

TF	Environ cond. [#]	Harbison <i>et al.</i>							PRIORITY		Literature
		AlignACE Top MAP	MEME Top MAP	MDscan Top MAP	MEME_c Top MAP	Kellis Top score	CONVERGE Top enrichment&	Top MAP	Predicted class		
Basic Helix Loop Helix TFs											
Cbf1	SM YPD	rAAAAARAAR CACGTGacc.	RTCACGTGm CACGTGACCm	kCACGTGm CACGTGay	— gGTCACGTG	CACGTGry TCT....AGC	CACGTGAY CACGTGACcm..wk.	rTCACGTG GTCACGTG	bHLH 97% bHLH 67%	rTCACrTGA	
Ino2	YPD	rARARARR.AA	sCAYSTGMw.a	kCAsrTGc	.sCATGTGma.A	TTCACATGC	gCATGTGa	TCACATGC	bHLH 86%	ATTTCACATC	
Ino4	YPD	A...A.AArArAR	tTTYCACATGs	CAygtGma	tTTTCACATGs	T.TTCACATG	.CATGTGa	CATGTGAA	bHLH 85%	CATGTGAAAT	
Phd1	BUT90 YPD	rRAAArRRAA RARARR.RAA	rAaA.grAaA.RrRaA yyTtk...yyycykTTkyC	.SSSsSSS sctGCags	TT.TTTtT.TTkTyTT GMAGGCACg	CGT.....TTT	RAA.rrrm.armrAA GCA..GAT	kCGTgsc. .mAGGCAC	bHLH 95% bHLH 92%	—	
Pho4	Pi- YPD	AAAAArAAA..A r.AAR.RRAAAA	sCACGTg yTgyTkTtcYTtwT.Y	sCACGTGs GTGTGtgT	g.CACGTGs AaAakAgAaAAaggawAr	CCACGTG CCT.CCT	cAcGTGss TSCkYGAGYAsTWAR	CACGTGcs YRTGKG.G	bHLH 81% bHLH 89%	wcacgtk.g	
Sok2	BUT14	ARAArRAAA.R	ArrM..AAAmr.RrAA	SSss.sSG	AaaAr.AAAaAgAaAA	CGT.....TTT	yAGGsAm.	.sCGTg..	bHLH 93%	—	
Tyc7	YPD	rRAARArAAAYs	rYCACstGAYg	TCACGTGA	rYCACstGAYg	CACGTGA	TCACGTGA	CasGTGAT	bHLH 92%	CA..TG	
TFs belonging to structural classes other than the three modeled explicitly											
Leu3	SM YPD	aA.AAAaaa...A rARArAArAA	gCCsGtacCGSwc m.GmaAsaaaArkaAArm	CGgtacCG CGgtacCG	gCSgGTacCGs m.GmaAsaaaArkaAArm	CCGGT.CCGG CCGGT.CCGG	cSGgTwcC rATkyrTcAkTMATA	GgtACyGG CGRTwCG	other 80% other 86%	yGCCGGTACCGGyk	
Nrg1	H2O2Hi H2O2Lo YPD	rAAAAArAAR AArAARAArAAR rAARArAAr	srAarmSrAAA rrAArakArAamsraAAr rAAA.rrrAA.RaAAR	gGACCCTk ggsGgsGG CACACACA	GGACCCTk rAaAr.gAAAAAsaAAr CssAr.CsCmcS.s.m	TAAACGA TTC.....GGC CTT.....ATG	.aGGGtcc mAGGcAc AAGgAAAA	.GGACCT mAGGcAc G.GTGTGt	other 89% bHLH 86% bHLH 90%	CCCT	
Rap1	YPD	rTGYayGGrTg	..grTGYayGGrTgyr	yCCCrtrCM	rrTGYayGGrTgya	TGT.TGGGT	.wrCAyCCrtrCAYc	ACCCRTAC	other 90%	wrmACCCATACAYY	
Reb1	H2O2Hi YPD	ksCGGGTAAy ksCGGGTAAy	.rTTACCCG.myt ...ksCGGGTAAy.	rTTACCCG rTTACCCG	rTTACCCG.m .rTTACCCGsm.	TTACCCG GTTACCCG	SCGGGTAA rTTACCCGsm.....	mTTACCCG mTTACCCG	other 86% other 93%	TTACCCGG	
Ste12	Alpha BUT14 BUT90 YPD	AAAArRAAA..R RAARAA.RAAA rAr.AARrAAA AArAaaaaAA	gAaAca..t.TgAaAca tTTtyyTtyt.k.YTTY smaaA..A.r.Ar.Rrra AmAaRAAgsAArA.aaaa	tGTTCA. tGTTCA. .TgAAACA .TrAAACA	wTTtwAmwwkwwTta GCaraarM.rcAagAmms CCsrAw..GGaA aAACa.t.T. TGAAACm	TGAAACr TGAAACA TGAAACA TGAAAC	Aaacr...yrAaAc rcATTcyy ...r..tGTTCA.m. AAacr..t.yrAAC	ttTGAAAC AaGAATCT .TGAAACA .TgAAACA	fork 95% other 91% other 94% other 78%	ATGAAAC	
Ume6	H2O2Hi YPD	TsgCGGGCTA TsGGCGGCTA	wwTAGCCGCcsa.. wwTAGCCGCCsA.s	TAGCCGCC TAGCCGCC	awwTAGCCGCcsa.. wwTAGCCGCCsA.s	TAGCCGCCs.. .wTAGCCGCCsA..wtTAGCCGCC s.tsgGCGGCTAww.	TAGCCGCC AGCCGCCs	other 92% other 86%	wGCCGCCGw	

* The motif with the inserted r was experimentally confirmed by Harbison *et al.* after they conducted a gel-shift assay to verify the authenticity of the motif they obtained by their *in silico* analysis for Cin5.

** Harbison *et al.* report the longer motif with the central G as literature consensus for HAC1, but in a literature search we conducted, we found that a new binding site TYACGTGTYM without the central G has been experimentally confirmed by [5] using gel-retardation assays.

Regulators were profiled in one or more specific environments by Harbison *et al.* The condition in bold (one per TF) is the condition used by Harbison *et al.* in their final motif, which we have used in the main text of the paper. In cases where they do not report a final motif, we arbitrarily select for the main text of the paper the condition which yields the largest number of probes. Here, we list all. Environmental conditions in brief are Alpha: mating, BUT14/BUT90: filamentation, H2O2Hi: highly peroxide, H2O2Lo: mildly peroxide, Pi-: phosphate deprived, RAPA: nutrient deprived, SM: amino acid starved, YPD: rich medium.

& Although all other programs mentioned in this table report a score for the discovered motif, we did not find a score attached to CONVERGE from the results reported by Harbison *et al.* Here, we use the reported motif with the top enrichment score. It is thus important to note that this is a post-processed motif, giving CONVERGE a slight benefit over all other programs in the table.

Table 1 reports results for all TFs considered in the main text of the paper, and in all environmental conditions used by Harbison *et al.* Sequence sets contained at least 10 probes, each bound with a p -value < 0.001 by the respective TF. There are 54 total sequence sets covering the 30 TFs. Known binding sites are reported in the literature for effectively 45 of these sequence sets (assuming the same binding mechanism for a TF in all conditions, which may not be necessarily true). Out of these 45 sequence sets, PRIORITY finds the correct binding site motif in **33**, AlignACE finds it in **6**, MEME finds it in **19**, MDscan finds it in **29**, MEME.c finds it in **24**, Kellis finds it in **22**, and CONVERGE finds it in **29**.

In 10 of the 12 cases that PRIORITY does not find a motif matching the literature, it predicts a class other than the class of the TF being profiled. This is helpful in diagnosing that the reported motif may not be the correct motif for the profiled TF.

Among the 9 sequence sets with no known binding site in the literature, PRIORITY predicts a motif of the correct class in the sequence sets of bZips Cst6 and Met4, so these might be good candidates for experimental validation. The cases of Met4, Fhl1 (where PRIORITY consistently predicts a class of “other” in all sequence sets), Phd1 (where PRIORITY finds a similar bHLH motif in both sequence sets), and Sok2 have all been discussed in detail in the main text of the paper.

References

- [1] Bailey,T. and Elkan,C. (1994) Fitting a mixture model by expectation maximization to discover motifs in biopolymers, *ISMB '94*, AAAI Press, Menlo Park, California, pp. 28–36.
- [2] Harbison,C., *et al.* (2004) Transcriptional regulatory code of a eukaryotic genome, *Nature*, 431:99–104.
- [3] Kellis,M., Patterson,N., Endrizzi,M., Birren,B., and Lander,E. (2003) Sequencing and comparison of yeast species to identify genes and regulatory elements, *Nature*, 432:241–254.
- [4] Liu,X., Brutlag,D., and Liu,J. (2002) An algorithm for finding protein-DNA binding sites with applications to chromatin immunoprecipitation microarray experiments, *Nature Biotech.*, 20:835–839.
- [5] Patil,C., Li,H., and Walter P. (2004) Gcn4p and novel upstream activating sequences regulate targets of the unfolded protein response, (*PLoS Biology*), 2(8):E246.
- [6] Roth,F., Hughes,J., Estep,P., and Church,G. (1998) Finding DNA regulatory motifs within unaligned non-coding sequences clustered by whole-genome mRNA quantitation, *Nature Biotech.*, 16:939–945.