

Phylogenetically and Spatially Conserved Word Pairs Associated with Gene Expression Changes in Yeasts

Derek Y. Chiang

Department of Molecular and Cell Biology, University of California, Berkeley

Alan M. Moses

Graduate Group in Biophysics, University of California, Berkeley

Manolis Kamvysselis

Laboratory of Computer Science, Massachusetts Institute of Technology

MIT / Whitehead Institute Center for Genome Research

Eric S. Lander

MIT / Whitehead Institute Center for Genome Research

Department of Biology, Massachusetts Institute of Technology

Michael B. Eisen

Department of Molecular and Cellular Biology, University of California, Berkeley and

Department of Genome Sciences, Ernest Orlando Lawrence Berkeley National Lab

Corresponding author:

Derek Y. Chiang

1 Cyclotron Road, Mailstop 84-355

Berkeley, CA 94703

dchiang@ocf.berkeley.edu

Running head: Conserved word pair templates

ABSTRACT

Background

In most eukaryotes, transcriptional regulation is often multifactorial, involving multiple factors binding to the same transcriptional control region (*i.e.*, promoters and enhancers). Phylogenetic footprinting among related species is a powerful approach to detect putative binding sites for individual transcription factors. We extend this concept of functional conservation to higher-order features of transcriptional control regions.

Results

We present a method to discover simple conserved word templates, defined as sequence pairs that are jointly conserved in the intergenic regions of four *Saccharomyces sensu stricto* genomes and are closely spaced within the upstream regions of open reading frames. Our method identified 1117 pairs of 6-bp long words that met the conserved word pair template criteria. Groups of genes that shared a template in their transcriptional control regions were over ten times more likely to be associated with gene expression changes. Many sequences comprising templates matched known transcription factor binding sites. In addition, these associations with gene expression were often significant in conditions where the corresponding transcription factors are known to be active.

Conclusions

The incorporation of both joint conservation and spacing constraints of sequence pairs yields groups of target genes that were specific for common patterns of gene expression. Our work suggests that positional information, especially the relative spacing between transcription factor binding sites, represent a common organizing principle of TCR's.

BACKGROUND

At a cellular level, all organisms have evolved intricate signaling networks that sense and respond to the environment. The activation of one or more signaling networks often leads to coordinated changes in gene expression, via the regulated binding of transcription factors to transcriptional control regions (TCR's) of genes, *e.g.* enhancers and upstream activating sequences. In yeast and higher organisms, the transcriptional regulation of individual genes is often multifactorial: multiple transcription factors may bind to a single TCR [1] [2] [3]. In some cases, transcription factors bind to TCR's independently of one another in order to alter gene expression in response to distinct cellular cues [4]; in other examples, multiple factors bind and act cooperatively to modulate gene expression via direct or indirect physical interaction with each other [5] [6] [7].

Due to multifactorial control, the regulation of a gene under a particular environmental condition may involve various assortments of binding sites for transcription factors that are active under that condition. The challenges in understanding transcriptional regulation include both the identification of regulatory sequences in TCR's and the recognition of some design principles underlying multifactorial regulation. One core principle seems to be the encoding of positional information within transcriptional control regions. Many experimental results have noted the importance of spatial constraints among individual binding sites [8] [9] [10].

Previous computational work to explore multifactorial regulation usually began with the identification of putative transcription factor binding sites. A plethora of computational methods has been developed to find over-represented sequences in a subset of genes believed to contain a common transcription factor binding site (reviewed in [11] [12]). The rapid pace of genome sequencing has enabled a complementary approach – phylogenetic footprinting (reviewed in [13]

[14]) – that recognizes that the conservation of sequences across related organisms often reflects evolutionary selection for their presence in TCR's. To perform phylogenetic footprinting analyses systematically, several algorithms endeavor to detect putative binding sites using orthologous TCR's among related genome sequences [15] [16] [17].

After compiling a collection of putative binding sites, associations can be made between various binding site assortments and gene expression. Some recent studies include Boolean logic [18], regression methods [19] [20] [21], spatial clustering [22], and multiple binding site matrix classifiers [23] [24] [25]. Spatial information on the relative locations of binding sites is ignored in all but the last two classes of approaches. Yet even these methods, which often use fixed orderings among the individual binding sites, may miss permutations of binding sites within TCR's that may still be bound by their cognate transcription factors.

The primary aim of this work was to evaluate the role of spacing between pairs of jointly conserved sequences in the regulation of gene expression. In doing so, we expanded the focus of phylogenetic footprinting from the conservation of contiguous sequences to higher-order features of TCR's, namely the spatial organization of individual binding sites. Because transcription factors participating in multifactorial regulation may require physical proximity among their binding sites, we searched for groups of conserved sequences that were more closely spaced than expected in TCR's. We refer to these spatially organized sequences as conserved word templates. As a proof of principle, we started with the simplest example of such templates: pairs of conserved words. Conservation was assessed using the genome sequences of three additional *Saccharomyces sensu stricto* species, which were chosen to be sequenced in order elucidate regulatory sequences [26]. To exploit this comparative genome data, we have devised a method that systematically tested sequence pairs for joint conservation and close spacing. To

evaluate the regulatory information contained in these conserved word pair templates, we adopted a group-by-sequence approach to test for significant associations between sequence rules and gene expression changes [27]. Since genes regulated by the same set of transcription factors often display similar gene expression patterns in certain experimental conditions; we searched for TCR templates whose gene targets were associated with common changes in gene expression. Significant associations with gene expression changes, as well as the prevalence of known transcription factor binding sites, suggest that conserved word pair templates comprise sequences important for multifactorial regulation in yeast.

RESULTS

Discovery of conserved word pair templates

We present a method to find conserved higher-order sequence templates from related *Saccharomyces sensu stricto* genomes (Figure 1). Our method incorporates sequential statistical tests, with each step focusing on a distinct property of conserved sequence templates. The simplest instances of sequence templates involve word pairs and their relative spacing. As described in detail below, pairs of words that were conserved in the same intergenic regions of the *Saccharomyces sensu stricto* genomes were identified using a chi-square test for independence. Next, the relative spacing of co-conserved word pairs was assessed for proximity using a permutation test. Finally, to evaluate the transcriptional information contained in conserved word pairs with close spacing, the expression of genes containing TCR templates were compared to the rest of the genome.

By pooling instances of either a sequence or its reverse complement, we initialized our word list using all 2080 words of length six. For each TCR (up to 600 nucleotides of an

intergenic region found upstream of an open reading frame), a word was labeled conserved if all six bases were identical among three or more *Saccharomyces sensu stricto* genomes in the CLUSTALW alignment of that TCR. To systematically test whether words were conserved more often in the same intergenic regions of the *Saccharomyces sensu stricto* genomes than expected by independent conservation, a chi-square test was performed on all possible pairwise combinations of words (see Materials and Methods). Pairs of words that overlapped each other by more than three nucleotides were excluded. A significant proportion of word pairs showed dependent conservation: among the 2.18 million word pairs tested, 8452 of them (~4%) had conservation χ^2 scores greater than 31.1. This threshold corresponds to a probability of 0.05 for obtaining one or more false positives after a Bonferroni correction for multiple testing.

Next, we selected word pairs that displayed closer physical spacing in intergenic regions than expected by chance. As a metric for the closeness in relative spacing between word pairs, the average minimum distance in *S. cerevisiae* (\bar{D}), was calculated among the genes whose TCR's conserved both words. If two non-overlapping words were closely spaced in all TCR's, we should find \bar{D} to be smaller than expected by chance, as assessed using a permutation test of the TCR labels (see Materials and Methods). After accounting for multiple testing by controlling the False Discovery Rate, a total of 1117 out of 8452 word pairs (~13%) had significantly small values for \bar{D} (FDR $q < 0.05$). A list of these closely spaced and jointly conserved word pairs can be obtained from <word_pair_templates.txt>. As a negative control, we also assayed a sample of word pairs that did not show dependent conservation ($\chi^2 < 1$), yet were jointly conserved in at least 10 TCR's. Only 161 out of 42801 (~0.4%) random word pairs with non-dependent conservation ($\chi^2 < 1$) showed significantly small values for \bar{D} . Figure 2 illustrates the distributions of \bar{D} for conserved word pair templates, jointly conserved word pairs, and

randomly conserved word pairs. The medians of these distance distributions were 100 nucleotides, 116.5 nucleotides and 132 nucleotides, respectively. Notably, the median \bar{D} for template pairs was significantly smaller ($p < 0.05$) than the median \bar{D} for randomly conserved pairs, as assessed by the statistical package R [28]. These results indicate that many of the word pairs that were conserved in the same intergenic regions of multiple *Saccharomyces sensu stricto* genomes also exhibited closer spacing in TCR's.

Conserved word pair templates were significantly associated with gene expression

Our method identified conserved word pair templates that were statistically significant with respect to both co-conservation in multiple genomes and close spacing in *S. cerevisiae* TCR's. To evaluate the regulatory information in these templates, we assessed the statistical association between gene groups that shared a template and changes in gene expression. Similar to other group-by-sequence approaches for finding regulatory sequences, we expect that gene subsets defined by common sequence rules should have similar and significantly different gene expression patterns under conditions where the transcription factors are active [27].

To assess the association between conserved word pair templates and differentially expressed genes, we identified gene subsets that contain both conserved words in the template within their TCR's and observed their expression patterns in publicly available datasets ([29] [30] [31] [32] [33] [34] [35] [36], see Materials and Methods for details). We then conducted Kolmogorov-Smirnov (K-S) tests to evaluate for differential gene expression between each gene subset and the whole genome. K-S tests provide a nonparametric, sensitive and robust way to compare two distributions. Similar results were obtained using other statistical tests, such as t -tests and likelihood ratio tests (A. M. M., unpublished data). A $P \times C$ output matrix was

computed: each conserved word pair in P was assigned a K-S p -value for each experimental condition observed in C . (see Materials and Methods). A gene subset associated with a word pair template was only deemed significant, and thus was included in the output matrix, if its K-S p -value: (1) met the threshold for multiple testing; and (2) was at least 10 times more significant than the K-S p -value for a gene subset associated with either word alone (see Materials and Methods). The latter criterion discounts gene expression changes that are due predominantly to the action of a single transcription factor.

Figure 3A displays the number of conserved word pair templates that were associated with significantly gene expression changes, for varying significance levels of the K-S test, which have been corrected for multiple testing (see Materials and Methods). Each line indicates the number of gene subsets that were significant in a different minimum number of experimental conditions. Several hundred closely spaced word pairs were significantly associated with differential gene expression. For example, 293 word pairs met an FDR-corrected significance threshold of $p < 10^{-3}$ for 3 or more experimental conditions, which represented 26% of all closely spaced word pairs. By contrast, using 721 randomly conserved word pairs with modest spatial constraints ($p < 0.1$), only 15 word pairs (3%) had significant expression at an FDR-corrected threshold of $p < 10^{-3}$ for 5 or more experimental conditions. For these particular thresholds, gene groups that contain a common conserved word pair template in their TCR's are about 12.5 times more likely to be associated with significant gene expression changes, compared with gene groups selected using randomly conserved word pairs. As shown in Figure 3B, the odds ratios for association with gene expression changes in multiple conditions varies between 5 and 25. Thus, gene groups that contain jointly conserved and spatially close word pairs are more significantly associated with gene expression changes than expected by chance.

Many identified sequences represent known transcription factor binding sites

Conserved word pair templates that were most strongly associated with gene expression changes also agreed with biological knowledge about transcription factors. For visualization purposes, we organized the $P \times C$ output matrix by hierarchically clustering the K-S p -values for the 339 word pairs that had significant associations with gene expression changes at a multiple testing threshold of 0.005 for 10 or more experiments (see Materials and Methods for details).

Hierarchical clustering of this output matrix identified groups of word pairs with similar K-S p -values in specific subsets of experimental conditions (Figure 4). In many cases, the grouped word pairs also comprised overlapping hexamer sequences, suggesting that some of the hexamers in different pairs may represent a larger, somewhat variable sequence. For example, group #13 in Figure 4 includes seven word pairs. In each of these word pairs, one of the component words – such as GCACGT or CACGTG – matches part of the Cbf1 or Pho4 binding sites. The other component word in each pair – such as AACTGT, ACTGTG, CTGTGG, TGTGGC or GTGGCT – represents part of the known Met31/32 binding site (AAACTGTGG). Therefore, genes whose TCR's contained any word pair within this group likely contain a conserved Cbf1 or Pho4 binding site, along with a conserved Met31/32 binding site. The distances between the conserved sites in these genes were also smaller than expected by chance.

A partial list of the 13 most significant groups of word pairs is shown in Table 1. Adjacent word pairs in the clustered output matrix with overlapping sequences were joined into consensus sequences (Table 2). Many of these consensus sequences matched transcription factor binding sites that had been biochemically verified.

Templates are associated with conditions when transcription factors are active

Further support that templates contain transcription information was obtained from a key observation: the experimental conditions with significant gene expression changes often corresponded to conditions in which the cognate transcription factors are known to be active (Figure 4). For example, gene groups with both a conserved Hap1 binding site and a conserved Hap2/3/4/5 binding site in their TCR's were significantly associated with induction in late stationary phase (Figure 5A). In addition, these genes were significantly associated with induction in ethanol – as well as with repression in fructose, glucose and sucrose – when compared to a common reference pool (Figure 5A). Among these 16 genes, 7 are involved in aerobic respiration. Indeed, both the Hap1 transcription factor and the Hap2/3/4/5 transcription factor complex may be active in the presence of nonfermentable carbon sources [37] [38]. By contrast, gene groups with both a conserved Hap1 binding site and a conserved Ecm22 or Upc2 binding site in their TCR's were only significantly associated with induction in the presence of a drug that inhibited ergosterol biosynthesis (Figure 5B). This group of 30 genes contained 8 ergosterol biosynthesis genes; this proportion represented an enrichment compared to the rest of the genome (hypergeometric p -value $< 2 \times 10^{-11}$). The transcription factors Ecm22 and Upc2 have been shown to induce the expression of ergosterol biosynthesis genes in response to low intracellular concentrations of ergosterol [39]. Note that the gene groups shown in Figures 5A and 5B showed gene expression changes in different sets of environmental conditions.

We briefly discuss a few examples of gene expression changes associated with templates; gene expression profiles for genes that contain TCR templates highlighted in Figure 4 and Table 1 can be found in the Supplementary Figures. Genes that contained conserved binding sites for Cbf1 and Met31/32 in their TCR's were strongly induced in cadmium, amino acid starvation and

early nitrogen starvation (Supplementary Figure 1). Groups 1 and 2 both share a word that corresponds to the Swi6/Swi4 binding site. In the former group, the other word matches the Swi6/Mbp1 binding site. In group 2, the other word represents a novel binding site (CGCCAA) that corresponds to part of a sequence enriched among the upstream regions of Stb1 chromatin immunoprecipitation targets [40]. Both of these gene groups peak in gene expression at G₁ (Supplementary Figure 2), which corresponds to the cell cycle phase when Swi6/Swi4, Mbp1/Swi6 and Swi6/Stb1 are all active.

DISCUSSION

This work describes two innovations for analyzing combinations of regulatory sequences. First, sequence conservation among closely related yeast species was used to find sequences that were more likely to be functionally important. Secondly, a template approach that considered joint positional distributions of word pairs increased the specificity of gene expression predictions using sequence-based rules. We have demonstrated that higher-order sequence features within TCR's were conserved across multiple *Saccharomyces sensu stricto* genomes. Closely spaced and jointly conserved word pairs were also more likely to be associated with gene expression changes. In many cases, associations between templates and gene expression changes were significant in conditions when the corresponding transcription factors are known to be active. In addition, a large proportion of words contained in templates matched known transcription factor binding sites. These results suggest that conserved word pair templates, which were discovered strictly based on higher-order properties of sequence conservation, also have biological relevance.

Conserved word pair templates may be grouped under several distinct classes of regulatory elements in TCR's. One possible interpretation of templates is that closely spaced sequence pairs may augment direct or indirect interactions between transcription factors by increasing the local concentration of the individual factors. For example, the proximity of Cbf1 and Met31/32 binding sites may promote their interaction. Experimental studies on the TCR's of MET3 and MET28 have demonstrated that the binding of Cbf1 enhances the DNA binding affinity of Met31/32 [41]. Since both Cbf1 and Met31/32 recruit the transcriptional activators Met4 and Met28 to induce gene expression, it is likely that all of these proteins interact at the TCR's of some sulfur assimilation genes.

By contrast, individual sequences found in templates may correspond to binding sites for transcription factors that bind independently under the same or separate conditions. The Hap1 and Hap2/3/4/5 transcription factors, whose binding sites were identified in a template, represent one of the earliest examples of multifactorial control under different conditions [42]. In some cases, templates could discern genes that shared binding sites for one transcription factor, but were differentially expressed under certain sets of conditions. Genes that conserved both Hap1p and Hap2/3/4/5p binding sites in their TCR's included mitochondrial enzymes and showed significant induction during growth in stationary phase and ethanol (Figure 5A). By contrast, genes that conserved both the Hap1p and Upc2/Ecm22p binding sites in their TCR's were enriched for ergosterol biosynthesis enzymes. These genes were significantly induced under treatment with drugs that inhibit ergosterol biosynthesis – itraconazole, lovastatin and terbinafine – but showed no change under different carbon sources (Figure 5B). The common link between these two enzyme classes is that they both require heme as a cofactor, which is sensed by Hap1 [43] [44] [45]. These results suggest that the differential expression of Hap1 target genes under

nonfermentable carbon sources versus ergosterol inhibition may be mediated by transcription factors that bind independently of Hap1. Further biochemical studies are required to distinguish between the mechanisms of independent binding and transcription factor interactions.

Close spacing between word pairs may be important for reasons other than the promotion of transcription factor interactions. Rather, different regions of TCR's at varying windows away from translation start sites may be more competent at recruiting or inhibiting RNA polymerase. These differences may be influenced by nucleosome accessibility, chromatin structure, or DNA physical properties, which can be correlated with A/T content (see [46] for references). Notably, we have also found that the relative proportions of A and T nucleotides vary considerably within the 200 bp closest to translation start sites (A. M. M. and Audrey Gasch, unpublished results). Low-complexity words that contained 4 or more A's or T's could be found in many templates (denoted by T_nC in Figure 4 and Table 1); these words may serve as surrogates for distance from translation start. Transcription factor binding sites that are closely spaced to these low-complexity words may be found in more transcriptionally competent regions of TCR's.

Much recent research has been devoted to finding dyad motifs with rigid spacing constraints among the individual sequences [47] [48] [49] [50]. These dyad search algorithms perform best with spacing ranges under 50 nucleotides, due to the computational complexity of scanning all intermediate positions in a fixed window size. Since instances of transcription factor binding sites that occur in the same TCR may be spaced further apart than this search window allows (note the large values of the average minimum distance statistic, \bar{D} in Table 1), our permutation test for close spacing is more flexible than many current dyad search algorithms. A key limitation of our approach is the use of exact words. Transcription factor binding sites often contain degenerate positions that reflect specificity [11]. The known binding sites listed in

Table 1 correspond to transcription factors with high sequence specificities. Since other known binding sites are poorly modeled by exact words, our method has failed to include them in conserved word pair templates. In addition, our method currently requires sequence identity for a word to be labeled as conserved. This strict requirement omits binding sites that may retain their function, despite mutations in degenerate positions.

The consideration of joint conservation and close spacing has provided insights into how TCR organization may influence the multifactorial regulation of gene expression. These criteria were motivated by experimental studies on the positional organization of individual binding sites within TCR's, with the hypothesis that this underlying architecture would be functionally conserved. Even more complicated higher-order sequence rules are apparent in the organization of metazoan enhancer sequences [51]. However, a common theme of organization in the TCR's of both of these organisms is the importance of relative spacing between transcription factor binding sites. The discovery of further principles for TCR organization would further advance our understanding of how regulatory information is encoded in genome sequences.

MATERIALS AND METHODS

Datasets

Whole-genome shotgun sequencing of *Saccharomyces bayanus*, *Saccharomyces mikatae*, and *Saccharomyces paradoxus* has been previously described [26]. All of these organisms are highly related to *Saccharomyces cerevisiae*, as they are grouped within the *sensu stricto* branch of the *Saccharomyces* genus [52]. Intergenic regions were aligned using CLUSTALW as described [26] and are available from the Saccharomyces Genome Database [53]. A total of 4101 CLUSTALW alignments were analyzed. These alignments correspond to transcriptional control regions and were filtered for orthologs in at least 3 genomes.

Gene expression measurements were obtained from the Stanford Microarray Database [54] and Rosetta [36]. The main experimental types among the 342 conditions examined include diauxic shift [29], cell cycle [31] [32], environmental stress response [30], DNA damage [33] [34], low phosphate [35], cadmium (N. Ogawa and P. O. Brown, unpublished data), and inhibition of ergosterol biosynthesis [36]. This data has been log-transformed (base 2), and each experimental condition has been median normalized.

Dependent Conservation of Word Pairs

To assess whether two words were co-conserved in the same intergenic regions, a chi-square test of independence was systematically conducted for all possible words of length six. We defined a word to include a 6-bp sequence and its reverse complement. Define a transcriptional control region (TCR) for a gene as the 600 base pairs upstream of its translation start site. TCR's shared between divergently transcribed genes less than 600 bp long were only counted once. A word was labeled conserved in a TCR if all six bases were identical among

three or more genomes in the CLUSTALW alignment. For each word pair (W , V) whose overlap was less than 4, a contingency table C_{wv} was constructed. In this table, $C_{wv} = \# \text{TCR}(I_w \cap I_v)$, where I_w, I_v are indicator variables for the presence of each conserved word in a TCR. TCR's shared between divergently transcribed genes less than 600 bp long were only counted once. The expected counts E_{wv} were obtained from an independence assumption, *i.e.* the product of the individual word conservation probabilities, multiplied by the total number of TCR's. Thus the chi-square statistic with Yates continuity correction was computed according to the definition:

$$\chi_{wv}^2 = \sum_{I_w=0}^1 \sum_{I_v=0}^1 \frac{(|C_{wv} - E_{wv}| - \frac{1}{2})^2}{E_{wv}} \quad (1)$$

Spatial Proximity of Constrained Word Pairs

The second requirement for a conserved sequence template involved constraints on spatial arrangements between individual words. Any method that evaluates spacing distributions between word pairs must take into account positional biases that may be present for individual words (A. M. M., unpublished results). We used a permutation test to evaluate the significance of the average minimum distance, excluding overlaps, between conserved word pairs. By permuting the TCR labels for one of the words, but not the word positions themselves, we retained the positional biases of individual words within intergenic regions. Within any given TCR t , define $p_t(W) = \{p_t^1(W), \dots, p_t^j(W)\}$ as a vector of positions in *S. cerevisiae* where word W is conserved. Suppose that words W and V were jointly conserved in TCR's $T_1 \dots T_N$. Then the average minimum distance, \bar{D} , can be computed as:

$$\bar{D}_{wv} = \frac{1}{T} \sum_{t=1}^T \min_{j,k} |p_t^j(W) - p_t^k(V)| \quad (2)$$

We used a permutation test to generate an empirical null distribution of $\bar{\mathbf{D}}$ for all word pairs with $N \geq 10$. After randomly permuting the labels \mathbf{t} for the position vectors of word V , a permutation test statistic, $\bar{\mathbf{D}}^*$, can be calculated as above. By repeating this resampling procedure R times, an empirical null distribution $\bar{\mathbf{D}}_{null} = \{\bar{\mathbf{D}}^{*1}, \dots, \bar{\mathbf{D}}^{*R}\}$ can be obtained. The significance of the observed average minimum distance, $\bar{\mathbf{D}}$, in the N promoters was calculated as its quantile in the empirical null distribution $\bar{\mathbf{D}}_{null}$. We set an upper bound of $R = 10^6$, but stopped permutations early if 20 or more values in $\bar{\mathbf{D}}_{null}$ were found less than $\bar{\mathbf{D}}$.

Correction for multiple testing involved control of the proportion of false positives using a False Discovery Rate method [55]. This method has increased power over Bonferroni-type methods. Permutation quantiles for all N word pairs tested were sorted in non-decreasing order:

$q_1 \leq \dots \leq q_N$. Let $k = \max\left(i : q_i < \frac{0.05i}{N}\right)$. Then the first k word pairs in the ordering had a

corrected significance level of $q < 0.05$, *i.e.* the rate of false positives is approximately 5%.

Association between Template-Specified Gene Groups and Gene Expression Changes

So far we have identified word pairs with two properties: dependent conservation and spatial proximity among all TCR's in the whole genome. These word pairs can be viewed as sequence-based rules for selecting a subset of genes based on the conservation of an element of TCR architecture. In this stage, we would like to evaluate the transcriptional information associated with these rules by assaying for gene expression changes among genes that match these sequence constraints.

For each gene expression condition c in our dataset, $c \in \{1, \dots, 342\}$, we tested the null hypothesis that a gene subset $\mathbf{G}_{wv} \subseteq \mathbf{G}$ selected by a conserved word pair (w, v) had the same

distribution of gene expression ratios (E_{wv}^c) as the entire genome (E^c). The alternate hypothesis stated that the two gene expression distributions were significantly different. Any gene was an element of G_o if its corresponding TCR conserved both sequences in the word pair. Since the size N_o of gene subsets may be small and the distributions may not be normally distributed, we used the nonparametric Kolmogorov-Smirnov (K-S) test. The test statistic K compares the cumulative distribution functions F_{wv}^c and F^c corresponding to E_{wv}^c and E^c by the formula $K = \max_x |F_{wv}^c(x) - F^c(x)|$. The significance level of an observed value K^* can be obtained using a numerical approximation [56].

A gene subset determined by a word pair was deemed to have significantly different expression if its K-S p -value was less than a certain threshold. To correct for multiple testing, this threshold was established by controlling the False Discovery Rate. The significance levels p_i from each K-S test were ordered in ascending order. Let N represent the total number of K-S tests performed, i.e. the number of jointly conserved, closely spaced word pairs times the number of gene expression experiments). If k was the largest i such that $p_i < i\alpha / N$, then the first k word pairs in the ordering were deemed to have a significance level of $p < \alpha$.

We ensured that the K-S p -value for the conserved word pair subset G_o was more significant than subsets G_w or G_v , comprised of only one conserved word by computing K for E_w^c vs. E_v^c , as well as for E_w^c vs. E^c . The marginal improvement of the joint word pair was defined as: $K(F_o^c \text{ vs. } F^c) - \max(K(F_w^c \text{ vs. } F^c), K(F_v^c \text{ vs. } F^c))$.

Hierarchical Clustering of Word Pair Associations

From the K-S tests for associations with gene expression changes described above, a $P \times C$ output matrix was obtained: each conserved word pair in P was assigned a K-S p -value for

each experimental condition observed in C . This output matrix was log-transformed (base 10). Rows in the output matrix were reordered using average linkage hierarchical clustering of the output matrix, performed by the program Cluster [57]. Since the log-transformed K-S p -values were all negative, a centered Pearson correlation was used as the similarity metric.

ACKNOWLEDGEMENTS

We thank Justin Fay and Audrey Gasch for their critical reading of the manuscript and insightful suggestions during the development of this work, as well as Peter Bickel and John Storey for their statistical advice. D.Y.C. is a Howard Hughes Medical Institute Predoctoral Fellow, and M.B.E. is a Pew Scholar in the Biomedical Sciences. This work was conducted under the US Department of Energy contract No. ED-AC03-76SF00098.

REFERENCES

1. **Composite regulatory elements: structure, function and classification** [<http://www.gene-regulation.com/pub/databases/transcompel/compel.html>]
2. Wolberger C: **Multiprotein-DNA complexes in transcriptional regulation.** *Annu Rev Biophys Biomol Struct* 1999, **28**:29-56
3. Kel OV, Romaschenko AG, Kel AE, Wingender EKolchanov NA: **A compilation of composite regulatory elements affecting gene-transcription in vertebrates.** *Nucleic Acids Res* 1995, **23**:4097-4103
4. Gasch AP: **The environmental stress response: a common yeast response to diverse environmental stresses.** In *Yeast Stress Responses*. Edited by S. Hohmann and W. H. Mager. Berlin: Springer; 2003: 11-70.
5. Mead J, Bruning AR, Gill MK, Steiner AM, Acton TBVershon AK: **Interactions of the Mcm1 MADS box protein with cofactors that regulate mating in yeast.** *Mol Cell Biol* 2002, **22**:4607-4621
6. Bhoite LT, Allen JM, Garcia E, Thomas LR, Gregory ID, Voth WP, Whelihan K, Rolfes RJStillman DJ: **Mutations in the Pho2 (Bas2) transcription factor that differentially affect activation with its partner proteins Bas1, Pho4, and Swi5.** *J Biol Chem* 2002, **277**:37612-37618
7. Verger ADuterque-Coquillaud M: **When Ets transcription factors meet their partners.** *Bioessays* 2002, **24**:362-370
8. Ambrosetti DC, Basilico C Dailey L: **Synergistic activation of the fibroblast growth factor 4 enhancer by Sox2 and Oct-3 depends on protein-protein interactions facilitated by a specific spatial arrangement of factor binding sites.** *Mol Cell Biol* 1997, **17**:6321-6329

9. Ludwig MZ, Patel NH, Kreitman M: **Functional analysis of eve stripe 2 enhancer evolution in *Drosophila*: rules governing conservation and change.** *Development* 1998, **125**:949-958
10. Liu Z, Little JW: **The spacing between binding sites controls the mode of cooperative DNA-protein interactions: implications for evolution of regulatory circuitry.** *J Mol Biol* 1998, **278**:331-338
11. Stormo GD: **DNA binding sites: representation and discovery.** *Bioinformatics* 2000, **16**:16-23
12. Ohler U, Niemann H: **Identification and analysis of eukaryotic promoters: recent computational approaches.** *Trends Genet* 2001, **17**:56-60
13. Duret L, Bucher P: **Searching for regulatory elements in human noncoding sequences.** *Curr Opin Struct Biol* 1997, **7**:399-406
14. Pennacchio LA, Rubin EM: **Genomic strategies to identify mammalian regulatory sequences.** *Nat Rev Genet* 2001, **2**:100-109
15. Blanchette M, Tompa M: **Discovery of regulatory elements by a computational method for phylogenetic footprinting.** *Genome Res* 2002, **12**:739-748
16. Schwartz S, Zhang Z, Frazer KA, Smit A, Riemer C, Bouck J, Gibbs R, Hardison R, Miller W: **PipMaker--a web server for aligning two genomic DNA sequences.** *Genome Res* 2000, **10**:577-586
17. Loots GG, Ovcharenko I, Pachter L, Dubchak I, Rubin EM: **rVista for comparative sequence-based discovery of functional transcription factor binding sites.** *Genome Res* 2002, **12**:832-839
18. Pilpel Y, Sudarsanam P, Church GM: **Identifying regulatory networks by combinatorial analysis of promoter elements.** *Nat Genet* 2001, **29**:153-159

19. Bussemaker HJ, Li H, Siggia ED: **Regulatory element detection using correlation with expression.** *Nat Genet* 2001, **27**:167-171
20. Keles S, van der Laan ME, Eisen MB: **Identification of regulatory elements using a feature selection method.** *Bioinformatics* 2002, **18**:1167-1175
21. Wang W, Cherry JM, Botstein D, Li H: **A systematic approach to reconstructing transcription networks in *Saccharomyces cerevisiae*.** *Proc Natl Acad Sci USA* 2002, **99**:16893-16898
22. Wagner A: **Genes regulated cooperatively by one or more transcription factors and their identification in whole eukaryotic genomes.** *Bioinformatics* 1999, **15**:776-784
23. Klingenhoff A, Frech K, Quandt K, Werner T: **Functional promoter modules can be detected by formal models independent of overall nucleotide sequence similarity.** *Bioinformatics* 1999, **15**:180-186
24. Pavlidis P, Furey TS, Liberto M, Haussler D, Grundy WN: **Promoter region-based classification of genes.** *Proceedings of the Pacific Symposium on Biocomputing* 2001, **6**:151-164
25. Kel-Margoulis OV, Ivanova TG, Wingender EK, Kel AE: **Automatic annotation of genomic regulatory sequences by searching for composite clusters.** *Pac Symp Biocomput* 2002, **7**:187-198
26. Kamvysselis M, Patterson N, Birren B, Berger BL, Lander ES: **Whole-genome comparative annotation and regulatory motif discovery in multiple yeast species.** *Proceedings of the 7th International Conference on Research in Computational Molecular Biology* 2003, **7**:

27. Chiang DY, Brown PO, Eisen MB: **Visualizing associations between genome sequences and gene expression data using genome-mean expression profiles.** *Bioinformatics* 2001, **17** Suppl 1:S49-S55
28. **The R Project for Statistical Computing** [<http://www.r-project.org>]
29. DeRisi JL, Iyer VR, Brown PO: **Exploring the metabolic and genetic control of gene expression on a genomic scale.** *Science* 1997, **278**:680-686
30. Gasch AP, Spellman PT, Kao CM, Carmel-Harel O, Eisen MB, Storz G, Botstein D, Brown PO: **Genomic expression programs in the response of yeast cells to environmental changes.** *Mol Biol Cell* 2000, **11**:4241-4257
31. Spellman PT, Sherlock G, Zhang MQ, Iyer VR, Anders K, Eisen MB, Brown PO, Botstein D, Futcher B: **Comprehensive identification of cell cycle-regulated genes of the yeast *Saccharomyces cerevisiae* by microarray hybridization.** *Mol Biol Cell* 1998, **9**:3273-3297
32. Cho RJ, Campbell MJ, Winzler EA, Steinmetz L, Conway A, Wodicka L, Wolfsberg TG, Gabrielian AE, Landsman D, Lockhart DJ, *et al.*: **A genome-wide transcriptional analysis of the mitotic cell cycle.** *Mol Cell* 1998, **2**:65-73
33. Gasch AP, Huang M, Metzner S, Botstein D, Elledge SJ, Brown PO: **Genomic expression responses to DNA-damaging agents and the regulatory role of the yeast ATR homolog Mec1p.** *Mol Biol Cell* 2001, **12**:2987-3003
34. Lee SE, Pellecioli A, Demeter J, Vaze MP, Gasch AP, Malkova A, Brown PO, Botstein D, Stearns T, Foiani M, *et al.*: In *Biological Responses to DNA Damage*. Edited by Cold Spring Harbor, NY: Cold Spring Harbor Laboratory Press; 2000: 303-314.

35. Ogawa N, DeRisi JB, Brown PO: **New components of a system for phosphate accumulation and polyphosphate metabolism in *Saccharomyces cerevisiae* revealed by genomic expression analysis.** *Mol Biol Cell* 2000, **11**:4309-4321
36. Hughes TR, Marton MJ, Jones AR, Roberts CJ, Stoughton R, Armour CD, Bennett HA, Coffey E, Dai HY, He YDD, *et al.*: **Functional discovery via a compendium of expression profiles.** *Cell* 2000, **102**:109-126
37. Gancedo JM: **Yeast carbon catabolite repression.** *Microbiol Mol Biol Rev* 1998, **62**:334-361
38. Kwast KE, Burke P, Poyton RO: **Oxygen sensing and the transcriptional regulation of oxygen-responsive genes in yeast.** *J Exp Biol* 1998, **201**:1177-1195
39. Vik AR, Rine J: **Upc2p and Ecm22p, dual regulators of sterol biosynthesis in *Saccharomyces cerevisiae*.** *Mol Cell Biol* 2001, **21**:6395-6405
40. Lee TI, Rinaldi NJ, Robert F, Odom DT, Bar-Joseph Z, Gerber GK, Hannett NM, Harbison CT, Thompson CM, Simon I, *et al.*: **Transcriptional regulatory networks in *Saccharomyces cerevisiae*.** *Science* 2002, **298**:799-804
41. Blaiseau PL, Thomas D: **Multiple transcriptional activation complexes tether the yeast activator Met4 to DNA.** *EMBO J* 1998, **17**:6327-6336
42. Guarente L, Lalonde B, Gifford P, Alani E: **Distinctly regulated tandem upstream activation sites mediate catabolite repression of the *CYCI* gene of *S. cerevisiae*.** *Cell* 1984, **36**:503-511
43. Burke P, Poyton RO: **Structure/function of oxygen-regulated isoforms in cytochrome c oxidase.** *J Exp Biol* 1998, **201**:1163-1175

44. Lees ND, Skaggs B, Kirsch DR, Bard M: **Cloning of the late genes in the ergosterol biosynthetic pathway of *Saccharomyces cerevisiae*--a review.** *Lipids* 1995, **30**:221-226
45. Kennedy MA, Barbuch R, Bard M: **Transcriptional regulation of the squalene synthase gene (ERG9) in the yeast *Saccharomyces cerevisiae*.** *Biochim Biophys Acta* 1999, **1445**:110-122
46. Liao GC, Rehm EJ, Rubin GM: **Insertion site preferences of the P transposable element in *Drosophila melanogaster*.** *Proc Natl Acad Sci USA* 2000, **97**:3347-3351
47. van Helden J, Rios AF, Collado-Vides J: **Discovering regulatory elements in non-coding sequences by analysis of spaced dyads.** *Nucl Acid Res* 2000, **28**:1808-1818
48. GuhaThakurta D, Stormo GD: **Identifying target sites for cooperatively binding factors.** *Bioinformatics* 2001, **17**:608-621
49. Liu X, Brutlag DL, Liu JS: **BioProspector: discovering conserved DNA motifs in upstream regulatory regions of co-expressed genes.** 2001, 127-38
50. Eskin E, Pevzner PA: **Finding composite regulatory patterns in DNA sequences.** 2002, **18 Suppl 1**:S354-S363
51. Berman BP, Nibu Y, Pfeiffer BD, Tomancek P, Celniker SE, Levine M, Rubin GM, Eisen MB: **Exploiting transcription factor binding site clustering to identify cis-regulatory modules involved in pattern formation in the *Drosophila* genome.** *Proc Natl Acad Sci USA* 2002, **99**:757-762
52. Cliften PF, Hillier LW, Fulton L, Graves T, Miner T, Gish WR, Waterston RH, Johnston M: **Surveying *Saccharomyces* genomes to identify functional elements by comparative DNA sequence analysis.** *Genome Res* 2001, **11**:1175-1186
53. **Saccharomyces Genome Database** [<http://genome-www.stanford.edu/Saccharomyces/>]

54. Sherlock G, Hernandez-Boussard T, Kasarskis A, Binkley G, Matese JC, Dwight SS, Kaloper M, Weng S, Jin H, Ball CA, *et al.*: **The Stanford Microarray Database.** *Nucleic Acids Res* 2001, **29**:152-155
55. Benjamini Y Hochberg Y: **Controlling the false discovery rate: a practical and powerful approach to multiple testing.** *J R Statist Soc B* 1995, **57**:289-300
56. Press WH, Teukolsky SA, Vetterling WT Flannery BP: *Numerical Recipes in C, Second Edition.* Cambridge: Cambridge University Press;1992.
57. Eisen MB, Spellman PT, Brown PO Botstein D: **Cluster analysis and display of genome-wide expression patterns.** *Proc Natl Acad Sci USA* 1998, **95**:14863-14868

FIGURE LEGENDS

Figure 1) Overview of method to discover conserved word pair templates in yeast

Figure 2) Word pairs in conserved word pair templates are closely spaced in *S. cerevisiae*

A comparison is shown among three categories of word pairs. “Template” denotes closely spaced and jointly conserved word pairs ($\chi^2 > 31.1$, spacing $q < 0.05$, $N = 1117$). “Depend” denotes dependently conserved word pairs ($\chi^2 > 31.1$, $N = 8452$) and includes all of the word pairs in the “close” category. “Random” denotes a sample of randomly conserved word pairs ($\chi^2 < 1$, $N = 4667$). For each category, the distribution of average minimum distances is represented by a box-and-whisker plot, which was generated using the statistical software package R. In a box-and-whisker plot, the box extends from the 25th percentile to the 75th percentile. The vertical line within the box denotes the median of the distribution. Dashed lines extend for 1.5 times the range of the box. Circles indicate extreme values.

Figure 3) Conserved word pair templates are associated with significant changes in gene expression

(A) *Total number of conserved word pair template associations at different K-S significance values.* The horizontal axis shows different False Discovery Rate-corrected significance levels for the Kolmogorov-Smirnov test (see Materials and Methods). The number of closely spaced word pairs meeting this cutoff for different minimum numbers of expression conditions is shown on the vertical axis. Word pairs were also filtered for an improvement of 10× over the K-S significance from any single word.

(B) *Relative enrichment for significant gene expression associations.* Relative enrichment was computed as an odds ratio: the fraction of gene groups selected by conserved word pair templates associated with significant gene expression changes, divided by the fraction of gene groups selected by randomly conserved word pairs associated with significant gene expression changes. The odds ratio is shown on the vertical axis; various FDR-corrected significance thresholds for gene expression associations are shown on the horizontal axis. Once again, word pairs were filtered for an improvement of 10× over the K-S significance from any single word.

Figure 4) Specific patterns of gene expression changes associated with conserved word pair templates

The output matrix O shown above was clustered by rows and visualized with TreeView (<http://rana.lbl.gov>). Each row corresponds to a conserved word pair template, and each column represents a single gene expression condition. The value in each cell O_{pc} corresponds to the K-S p -value of gene expression changes in condition c for a group of genes that contain the conserved word pair template p in their TCR's. Brighter intensities of orange denote more significant p -values from the Kolmogorov-Smirnov test. K-S p -values that failed to meet a False Discovery Rate critical value of 0.005 for multiple testing are displayed in grey. Word pairs that failed to meet a False Discovery Rate critical value of 0.005 for multiple testing in 3 or more experiments are not shown. Some of the most significant conserved word pair associations are labeled and annotated in Table 1. Abbreviations for significant conditions include: ND (nitrogen depletion), SP (stationary phase).

Figure 5) Multifactorial regulation of Hap1p target genes

(A) Genes with Hap1p (CCGATA) and Ecm22p/Upc2 (TCGTTT) sites conserved in their TCR's; (B) Genes with Hap1p (CCGATA) and Hap2/3/4/5p (CCAATC) sites conserved in their TCR's. Rows denote single genes and columns denote experiments. Genes are listed in ascending order of minimum distance between the two conserved words in the corresponding TCR of *S. cerevisiae*. Gene names highlighted in orange (A) or in purple (B) correspond to genes with similar biological roles. Arrows above the columns indicate conditions in which the displayed gene groups show significant gene expression changes according to the Kolmogorov-Smirnov test, after False Discovery Rate correction for multiple testing at a p -value of 0.005 (blue) or 0.01 (grey). Related sets of experiments are grouped between light grey vertical lines. The first set of experiments is a stationary phase timecourse, with samples at 2 h, 4 h, 8 h, 12 h, 1 day, 2 days, 3 days, 5 days, 7 days, 13 days, 22 days, 28 days. The second set of experiments includes samples taken from growth on different carbon sources, normalized against a common reference. The corresponding carbon sources in the growth media are: ethanol (E), fructose (F), galactose (C), glucose (G), mannose (M), raffinose (R) and sucrose (S). The last set of experiments, performed by Rosetta, corresponds to growth in the presence of drugs that inhibit ergosterol biosynthesis: itraconazole (I), lovastatin (L) and terbinafine (T) [36].

Figure 1

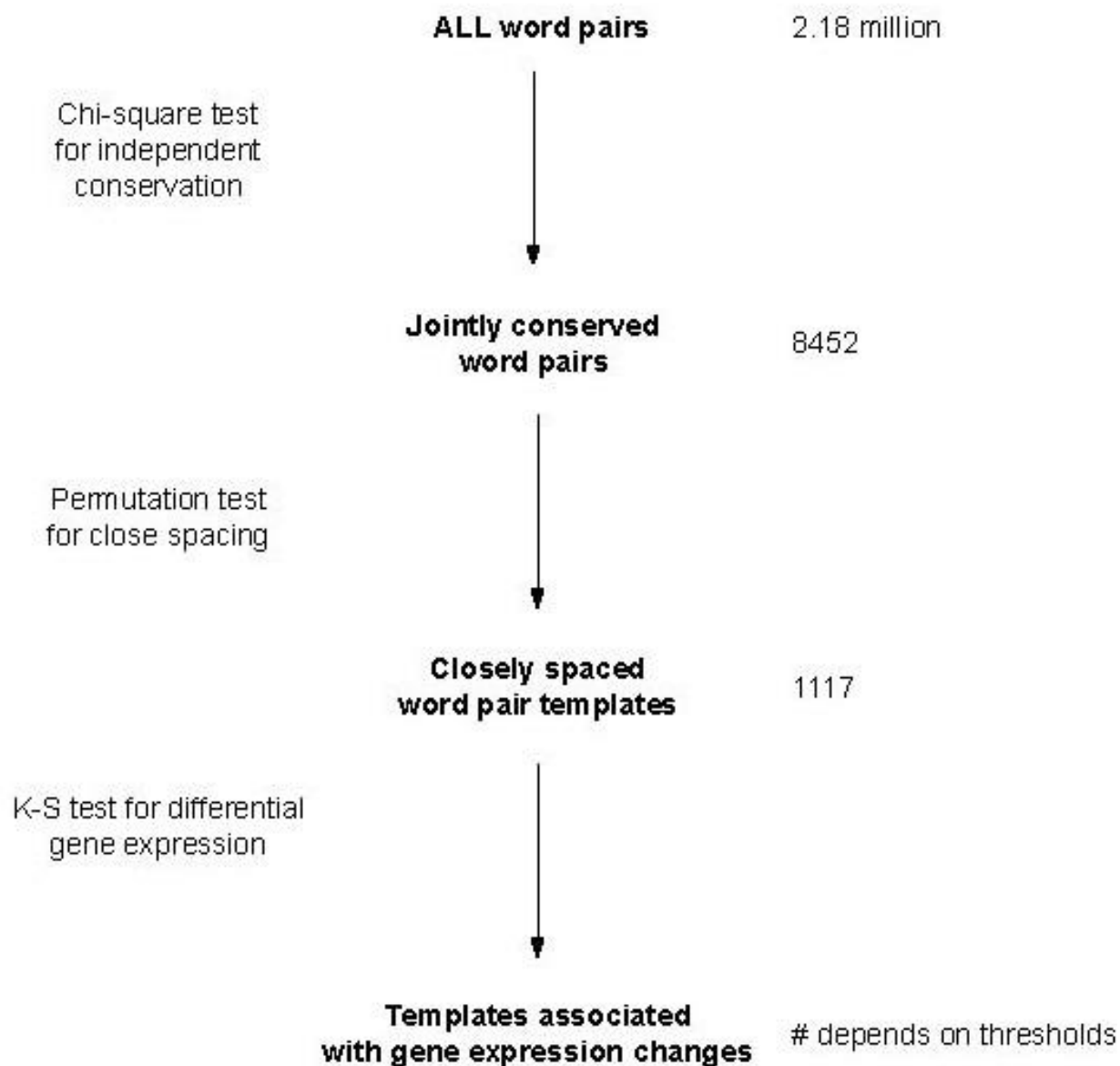


Figure 2

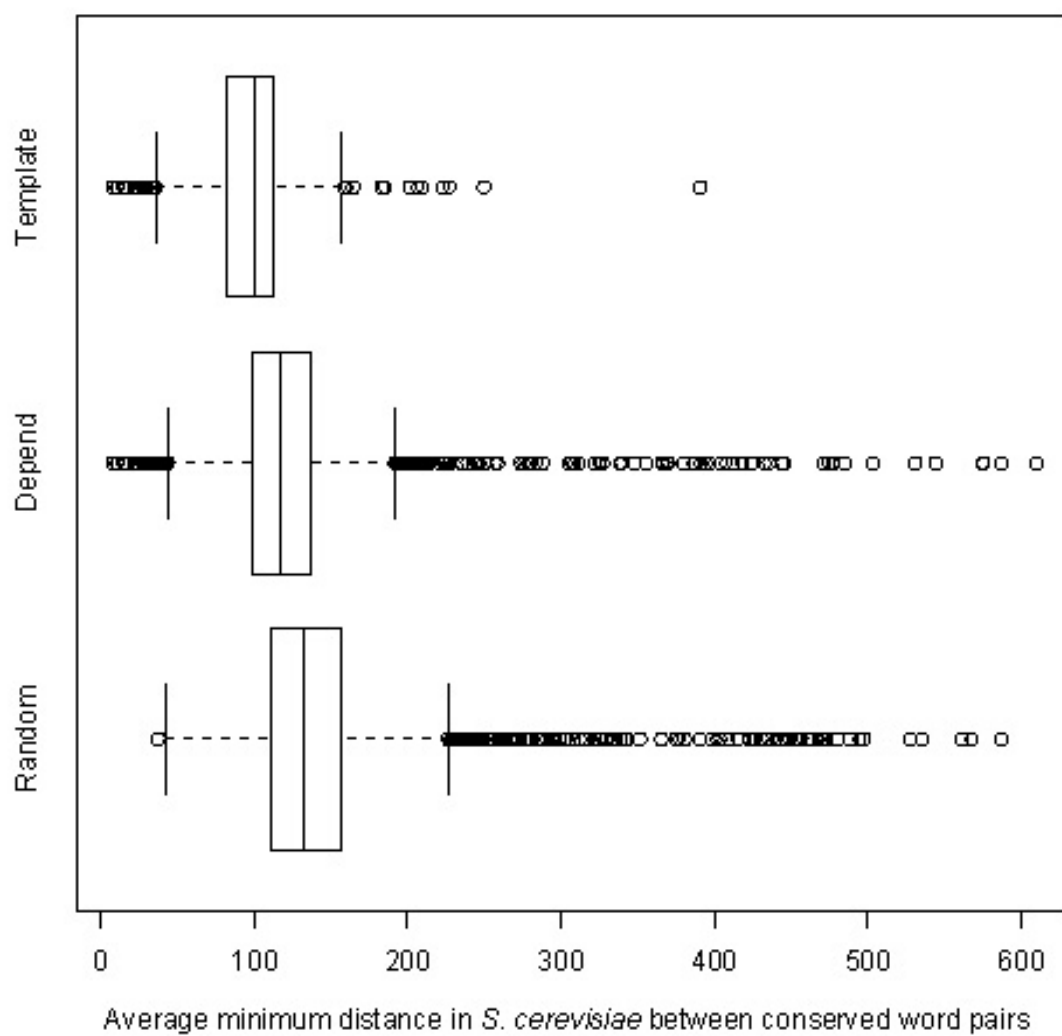
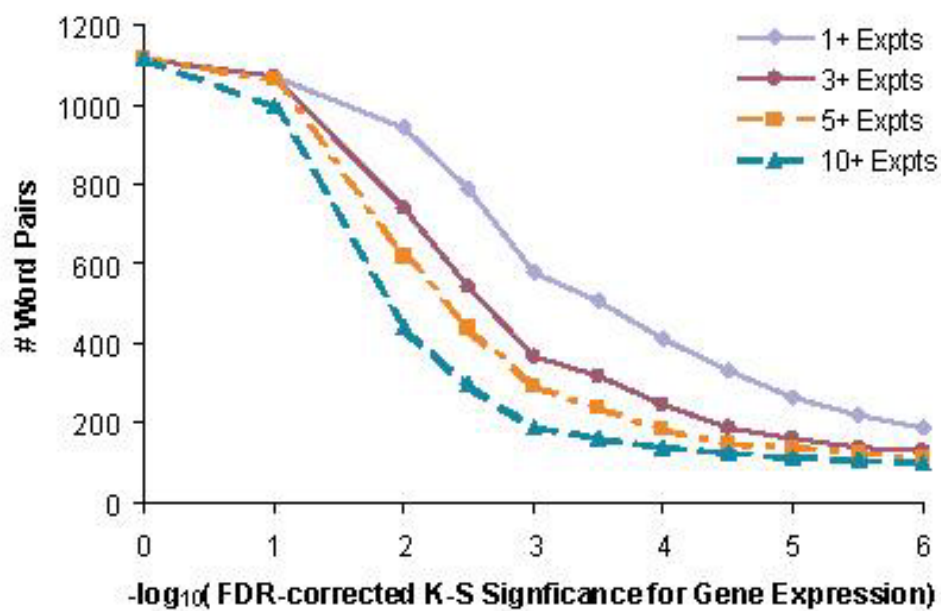


Figure 3

A



B

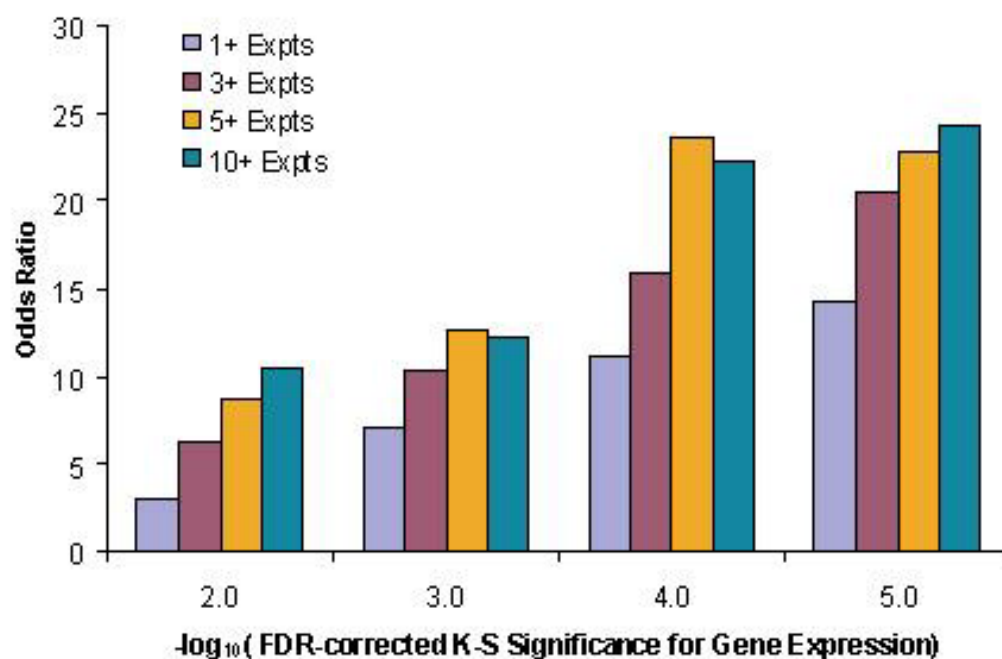
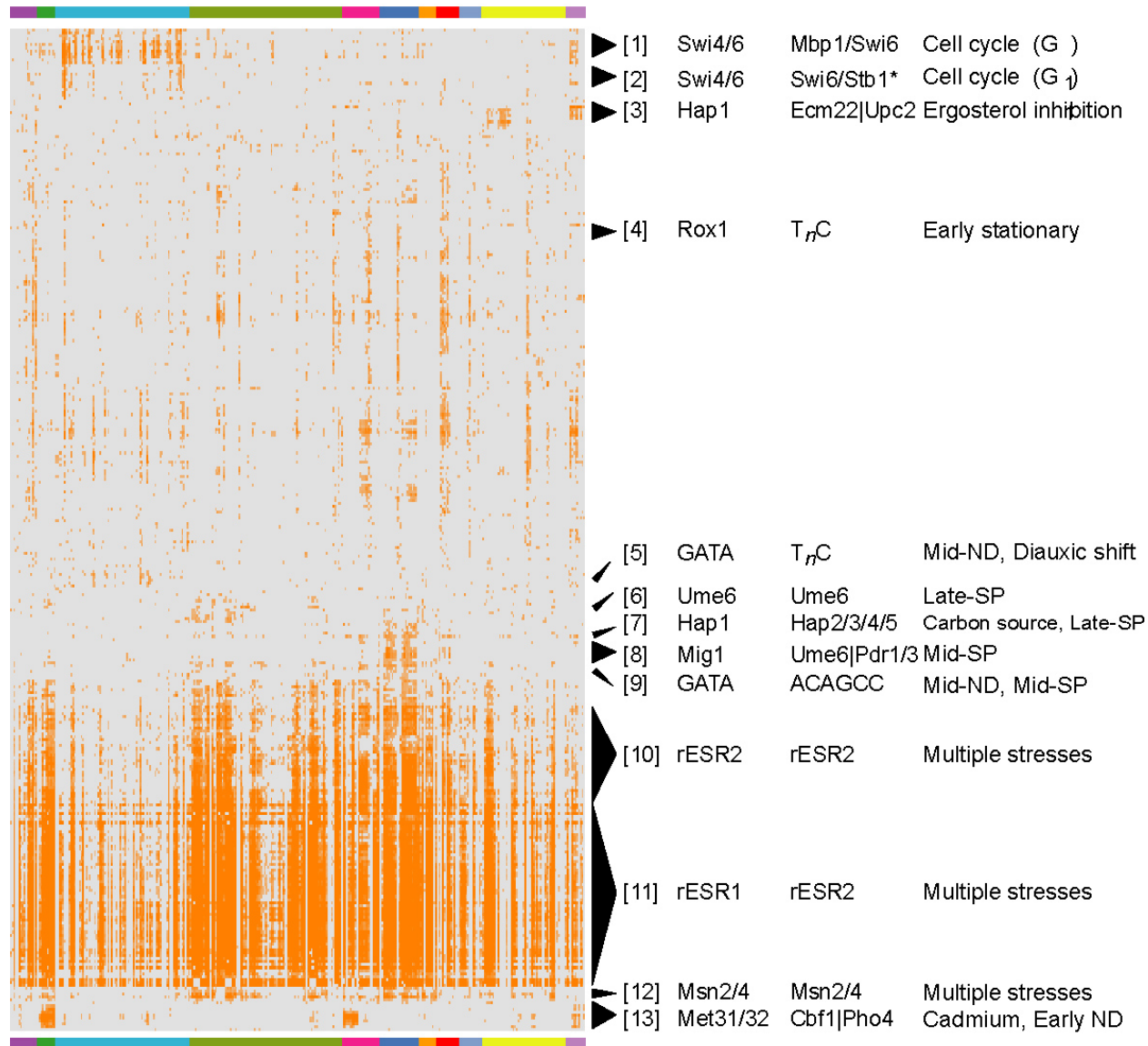


Figure 4



Experimental Conditions

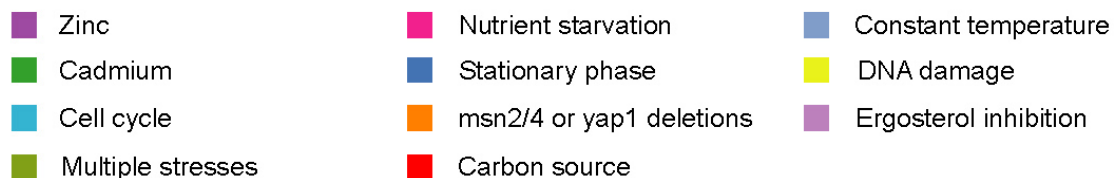


Figure 5

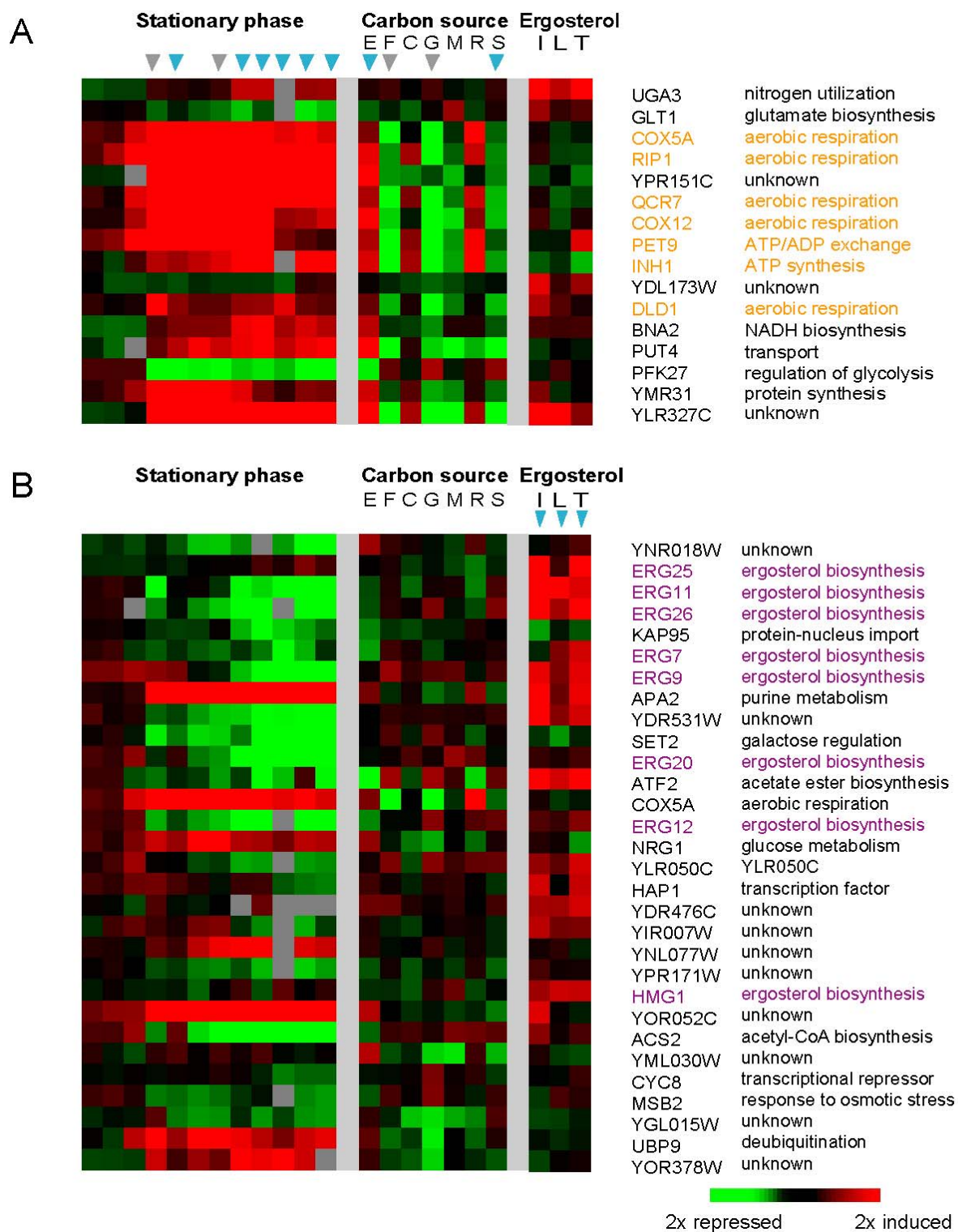


Table 1) Most significant groups of conserved word pairs

Statistics are listed for one representative word pair in clusters of overlapping sequences shown in Figure 4. Pairs in bold denote word pairs for which at least one cognate transcription factor is known. Multiple transcription factors that may bind the same sequence motif are separated by “|”.

	Conserved Word Pairs (Consensus of overlapping words)	Known transcription factors or motifs	Most significant word pair in consensus group			
			Conservation (χ^2 , p -val via Bonferroni)	Avg. min dist \bar{D}	# TCR	Expression conditions with significant gene subsets (FDR significance)
1	RCGAAA, RACGCG,	Swi4/6, Swi6/Mbp1	83.0 (7×10^{-13})	68.1	36	Cell cycle, G1 phase (10^{-6})
2	CACGAAA, CGCCAA	Swi4/6, Stb1 (putative)	55.6 (2×10^{-7})	78.2	25	Cell cycle, G1 phase (10^{-5})
3	CCGATA, TC[GT]TTT	Hap1, Ecm22 Upc2	36.2 (0.004)	81.7	30	MMS (DNA damage) (10^{-4}) Ergosterol inhibition (10^{-3})
4	TCTATT, CTTTTT	Rox1, T_nC	49.0 (6×10^{-6})	87.7	102	Early stationary phase (10^{-5})
5	GATAAG, TTCTTT	GATA, T _n C	36.0 (0.005)	100.5	88	Nitrogen depletion 8h (10^{-5})
6	CGGCTA, CGGCGG	Ume6, Ume6	179.2 (2×10^{-34})	81.9	15	Late stationary phase (10^{-4})
7	CCGATA, CCAATC	Hap1, Hap2/3/4/5	35.0 (0.007)	88.6	16	Stationary phase (10^{-5}) Diauxic shift (10^{-4}) Ethanol (10^{-4})
8	ACCCCA, CCGCCG	Mig1, Ume6 Pdr1/3	66.7 (7×10^{-10})	70.5	16	Stationary phase (10^{-6}) Ethanol (10^{-5})
9	GATAAG, ACAGCC	GATA, Novel	39.5 (0.004)	75.5	21	Nitrogen depletion 8h-2d (10^{-4}) Stationary phase 10h-2d (10^{-3})
10	ANTGAAA, GAAAAWT	rESR2 (Overlap)	96.9 (2×10^{-16})	96.8	68	Repressed in multiple environmental stresses (10^{-6})
11	G[AC]GATGAG, TGAAAATTTT	rESR1 motif, rESR2 motif	240.6 (10^{-49})	41.6	183	Repressed in multiple environmental stresses (10^{-6})
12	AWAAGG, AGGGG	Msn2/4 (Overlap)	94.7 (5×10^{-16})	99.0	29	Multiple stresses (10^{-3})
13	ACTGTGGC, [GT]CACGTG	Met31/32, Cbf1 Pho4	47.5 (2×10^{-5})	43.5	22	Amino acid starv. (10^{-6}) Nitrogen depletion (10^{-6}) Cadmium (10^{-6})

Table 2) Consensus sequences for significant groups of word pairs

The output $P \times C$ matrix of word pairs (P) that were significantly associated with at least 5 or more environmental conditions (C) was ordered using hierarchical clustering. Consensus sequences were assembled from groups of word pairs that were found in adjacent rows in the ordering. Residues are shown in bold if it is contained in at least two hexamers. Numbers denote the groups that are indicated in Figure 4. IUPAC codes used: K (G or T); M (A or C); R (A or G); S (C or G); W (A or T).

	Hexamer list for word 1	Consensus sequence 1	TF for consensus 1	Hexamer list for word 2	Consensus sequence 2	TF for consensus 2
1	CG CGAA AC GAAA GC GAAA C GAAAC CG AAAA	CRC GAAAM	Swi4/6	AA ACGC G ACGCG A ACGCG A CGCGT A CGCGA	AR ACGCGW	Swi6/Mbp1
2	AC CACGA C ACGAA A CGAAA CG AAAC	AC CAGAAAC	Swi4/6	CGCCAA	CGCCAA	Stb1 (putative)
3	CC GATA C GATAA	CC GATAA	Hap1	TCGTTT TCTTTT C GTTTA	TCKTTTA	Ecm22 Upc2
4	TCTATT	TCTATT	Rox1	CTTTTT TTTTTT	T _n C	Unknown
5	GATAAG	GATAAG	GATA motif	TTCTTT	T _n C	Unknown
6	GGCTAA CG GCTA	CG GCTAA	Ume6	CG GCGG	CGGCGG	Ume6
7	CCGATA	CCGATA	Hap1	CCAATC	CCAATC	Hap2/3/4/5
8	T ACCCC A CCCCA C CCCCA	T ACCCCAA	Mig1	CC GCCG C GCCGA	CC GCCGA	Ume6 Pdr1/3
9	GATAAG	GATAAG	GATA motif	ACAGCC	ACAGCC	Unknown
10	AAGTGA AATGAA AGTGAA AT GAAA CT GAAA T GAAAA	ANTGAAAAA	Unknown (rESR2 motif)	GAAAAA GAAAAAT AAAAAT AAATTT	GAAAAWTT	Unknown (rESR2 motif)

11	GAGATG GCGATG AGATGA CGATGA GATGAG ATGAGA ATGAGC TGAGAT TGAGCT GAGATG AGATGA AGCTCA	GMATGAGMTSA	Unknown (rESR1 motif)	TGAAAA GAAAAA AAAAAT AAAATT AAATTT	TGAAAATTT	Unknown (rESR2 motif)
12	ATAAGG AAAAGG	AWAAGG	Msn2/4	AGGGGT	AGGGG	Msn2/4
13	AACTGT ACTGTG CTGTGG TGTGGC GTGGCT	AACTGTGGCT	Met31/32	TCACGT GCACGT CACGTG	KCACGTG	Cbf1 Pho4