

## Supplementary Materials

Methods.....	1
1 Motif Extraction algorithm (MEX).....	1
2 Expression Coherence (EC) Method .....	6
3 Expression Data .....	8
4 Clustering Algorithm .....	9
Results.....	15
5 Clusters .....	15
6 EC patterns of clusters .....	16
7 Localization of motifs along promoters.....	20
8 Fisher distances between clusters .....	23
9 Intersections between sets of genes of couples of clusters .....	24
10 Comparing MEX to alternative approaches.....	26

## Methods

### 1 MOTIF EXTRACTION ALGORITHM (MEX)

MEX is a motif extraction algorithm that extracts statistically significant motifs from sequential data. MEX is a data driven unsupervised algorithm, hence does not need any preprocessing of the data or additional information apart from the data set itself. Furthermore, MEX finds motifs that are not necessarily over-represented in the data.

MEX was originally developed in a linguistic context, as a distillation tool for extracting words from corpora of natural language. As more intuitive, let us first describe the algorithm in its original context.

Consider a corpus of sentences, whose word delimiters have been removed (such as spaces, capital letters, punctuations, etc.). The problem at hand is to uncover the words that have originally constructed the sentences. MEX receives as an input such corpus, consisting of many sequences of a given finite alphabet of size  $N$  (e.g.  $N=26$  letters in the English alphabet,  $N=20$  amino acids in proteins and  $N=4$  nucleic acids in the case of DNA). The algorithm uses a directed graph, whose vertices,  $V$ , are composed of the letters of the given alphabet, in addition to a 'begin' and an 'end' vertices. A set of ordered pairs of vertices (directed edges) represent the order in which the letters appear in the corpus. For example, the edge  $e(t,h)$ , represents a connection from the vertex 't' to the vertex 'h', which means that the letter 'h' appears at some point along the corpus after the letter 't'. MEX loads the given corpus onto a directed graph, one sentence after the other. The edges representing each sentence are built, starting with the 'begin' vertex, followed by the letters composing the sentence, one after the other, and ending with the 'end' vertex. This way, ordered paths are created in the graph, such that each sentence is represented by a path. Each path is saved by MEX and will be used as a search path for patterns. This procedure is demonstrated in figure 1.1.

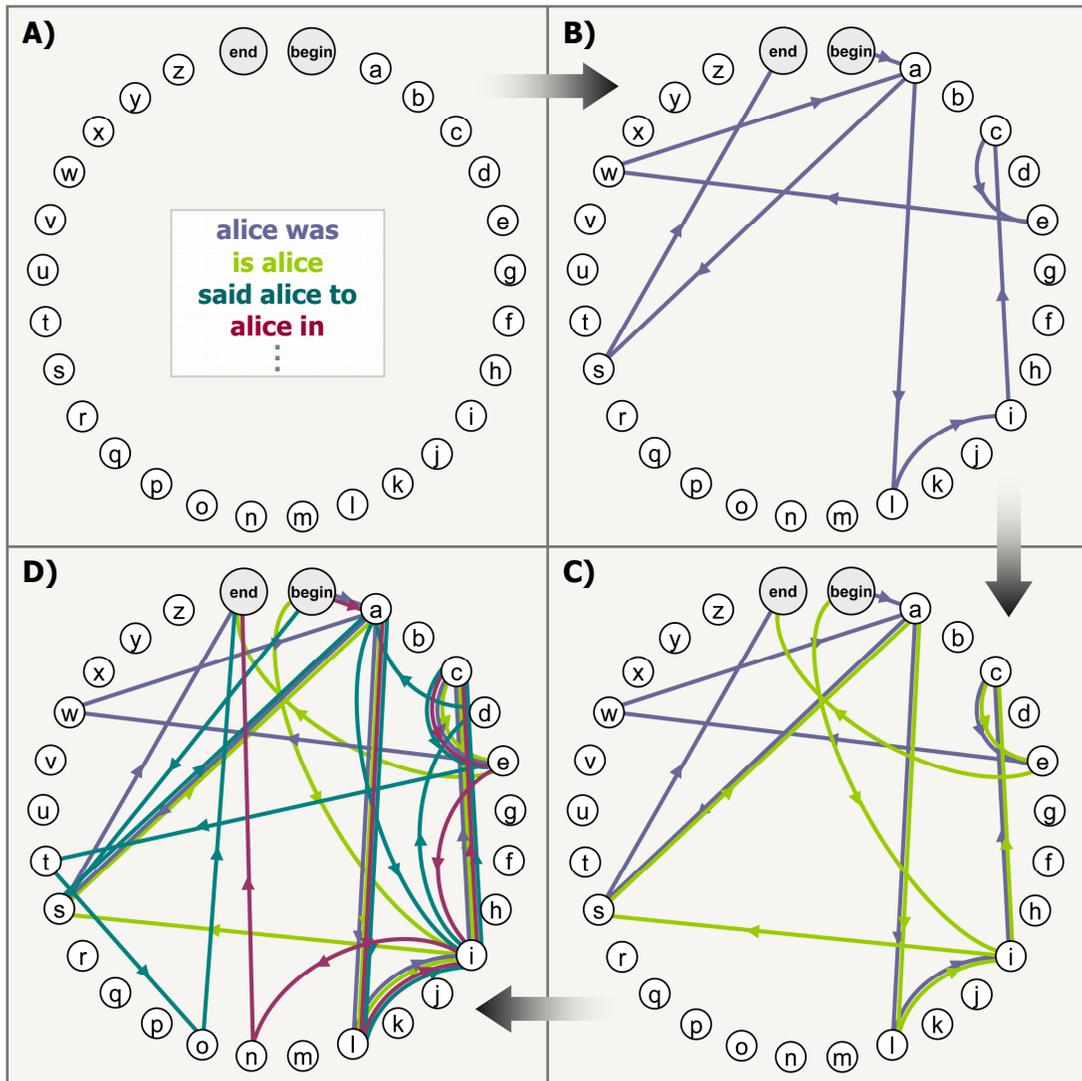


Figure 1.1 MEX loads the corpus onto a directed graph, one sentence after the other. The graph is composed of vertices representing the letters of the given alphabet, in addition to two vertices representing the beginnings and the endings of sentences (A). One sentence at a time, directed edges are added to the graph, representing the order in which letters appear in each sentence (B-D). The ordered edges composing a sentence are considered a path along the graph. In this example, four paths are loaded onto the graph: ‘alicewas’ (blue path), ‘isalice’ (light green path), ‘saidalice’ (turquoise path) and ‘alicein’ (red path), one after the other.

Once the entire corpus has been represented as search paths on a directed graph, the algorithm starts searching for statistically significant patterns. Intuitively, for each search path MEX looks for sub-paths that may be considered as candidates for being significant patterns. A sub-path that represents a significant pattern is expected to be shared by other paths throughout the graph, such that these paths will converge into the sub-path at its first vertex, form a bundle along the sub-path and scatter after the sub-path's last vertex. This follows the assumption that at different instances of a given word throughout the corpus, after the word ends, it is likely to find many different possible words following it. In such a case many paths will form a bundle along the sub-path representing the word and scatter immediately after it ends. The vertex after which such a divergence occurs may be considered as the last vertex of the pattern. A similar notion underlies the way MEX searches the start points of patterns, by looking for a divergence of a bundle while going leftwards through a search path. Figure 1.2 demonstrates this idea. The four paths in figure 1.2 converge and form a bundle along the sub-path 'a→l→i→c→e', after which they diverge.

This can be rephrased into a probabilistic language; for each search path (sentence) that is to be explored for patterns, two probability functions are defined, based on information inheres in the complete graph. The first one,  $P_{Right}$ , is the right moving ratio of the through-going flux of paths to the incoming flux of paths, which varies along the search path. Starting at the vertex  $e_1$  we define  $P_{Right}$  at  $e_2$  as:

$$P_{Right}(e_1, e_2) = p(e_2 | e_1) = \frac{\text{total no. of paths passing from } e_1 \text{ to } e_2}{\text{total no. of paths entering } e_1}$$

At  $e_3$   $P_{Right}$  becomes:

$$P_{Right}(e_1, e_3) = p(e_3 | e_1 e_2) = \frac{\text{total no. of paths passing from } e_1 \text{ through } e_2 \text{ to } e_3}{\text{total no. of paths passing from } e_1 \text{ to } e_2}$$

And generally:

$$P_{Right}(e_i, e_j) = p(e_j | e_i e_{i+1} e_{i+2} \dots e_{j-1}) = \frac{\text{total no. of paths passing from } e_i \text{ up to } e_{j-1} \text{ and continue to } e_j}{\text{total no. of paths passing from } e_i \text{ up to } e_{j-1}}$$

Similarly, a second function,  $P_{Left}$ , is defined as we proceed leftward from some vertex  $e_j$  down the search path towards the vertex  $e_i$  and examine the left-going ratio of the through-going flux of paths to the incoming flux of paths:

$$P_{Left}(e_j, e_i) = p(e_i | e_{i+1} e_{i+2} \dots e_{j-1} e_j) = \frac{\text{total no. of paths passing from } e_i \text{ to } e_j}{\text{total no. of paths passing from } e_{i+1} \text{ to } e_j}$$

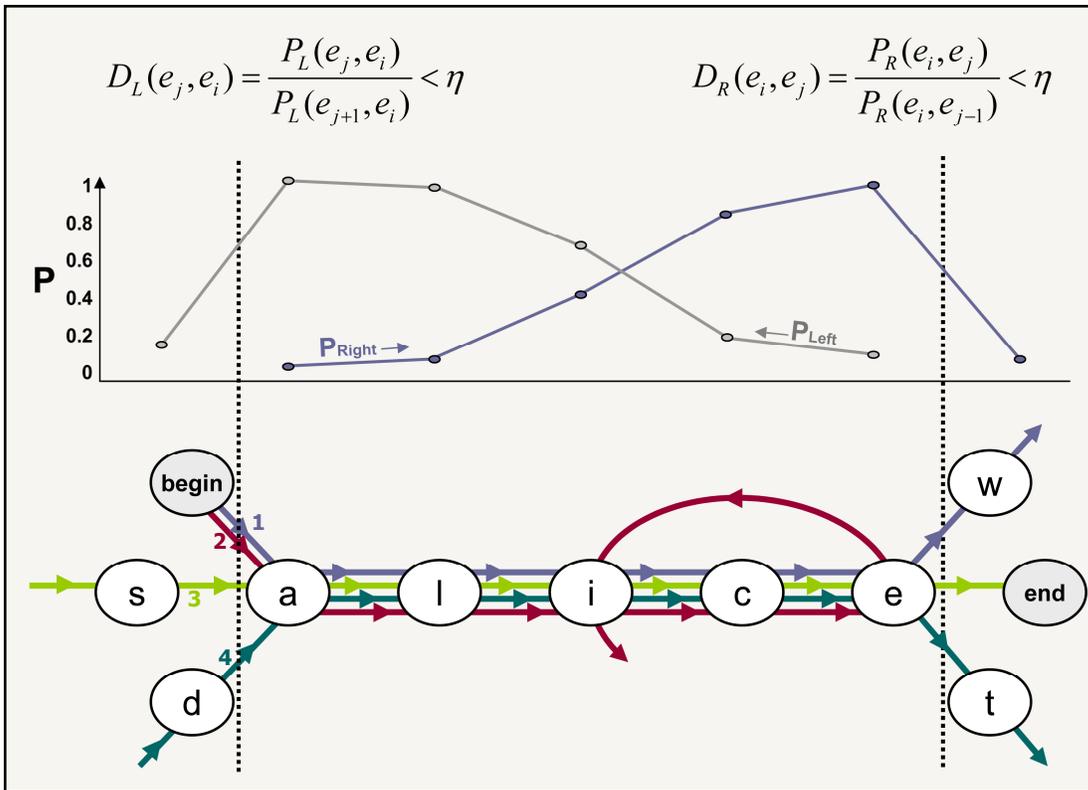


Figure 1.2 A partial view of the graph used by MEX. The search path no. 1, ‘alicewas’ (blue line), shares the sub-path ‘a→l→i→c→e’ with three other paths: ‘isalice’ (2), ‘saidalice’ (3) and ‘alicein’ (4). The four paths form a bundle that may constitute a significant pattern. The conditional probabilities  $P_{Right}$  and  $P_{Left}$ , originating at the vertices ‘a’ and ‘e’, respectively, are illustrated for the example shown here. A sharp drop in the right moving probability,  $P_{Right}$ , indicates that the paths constructing the bundle have scattered, thus may denote the end of the pattern. Similarly, a sharp drop in  $P_{Left}$  may indicate the beginning of the pattern, hence reveal the pattern ‘alice’.

Vertex	Conditional Probability Expression	$P_{Right}$
a	$P(a) = 8770 / 109625$	0.08
l	$P(l   a) = 1046 / 8770$	0.12
i	$P(i   al) = 486 / 1046$	0.45
c	$P(c   ali) = 397 / 486$	0.85
e	$P(e   alic) = 397 / 397$	1
w	$P(w   alicew) = 48 / 397$	0.12
a	$P(a   alicewa) = 21 / 48$	0.44
s	$P(s   alicewas) = 17 / 21$	0.81
b	$P(b   alicewasb) = 2 / 17$	0.12
e	$P(e   alicewasbe) = 2 / 2$	1
g	$P(g   alicewasbeg) = 2 / 2$	1
.	.	.
.	.	.
.	.	.

Table 1.1 Calculating right-going conditional probabilities for the search path ‘alicewasbeginning...’. Probabilities are calculated for a given search path, based on information in the entire graph. The corpus used in this example was the sentences from Alice in wonderland, by Lewis Carroll.

MEX calculates  $P_{\text{Right}}$  from different starting points to each vertex down the search path. Going rightwards through a sub-path that represents a significant pattern, it is expected that  $P_{\text{Right}}$  will first increase since other paths join the search path to form a coherent bundle, and then decrease as many paths leave the search path.

In order to demonstrate this, let us examine as a toy problem the corpus of Alice in wonderland, by Lewis Carroll. MEX has received as an input the sentences within Alice in wonderland, after all word delimiters have been removed. Going through the first search path ‘alicewasbeginningtogetverytired...’ MEX calculates the rightward-going probabilities,  $P_{\text{Right}}$ , along the path, as demonstrated at table 1.1. MEX starts at the first vertex ‘a’ and calculates the probability of its appearance in the corpus,  $P_{\text{Right}}(a)$ ; as ‘a’ appears in 8770 cases out of the total of 109625 letters in the corpus,  $P_{\text{Right}}(a) = \frac{8770}{109625} = 0.08$ .

MEX continues to the next vertex ‘l’, calculating the probability of its appearance after the previous vertex, i.e.  $P_{\text{Right}}(al|a)$ ; in this case, ‘l’ appears 1046 times after the 8770 instances of ‘a’, hence  $P_{\text{Right}}(al|a) = \frac{1046}{8770} = 0.12$ . MEX continues calculating the rightward-going probabilities  $P_{\text{Right}}(ali|al)$ ,  $P_{\text{Right}}(alici|ali)$  and so on, up to the end of the search path. As can be seen in table 1.1, the rightward-going probabilities initially rise and then drop sharply. Such a dramatic drop may occur owing to the sudden divergence of a coherent bundle, and will be considered as a candidate for terminating a pattern.

We will define the end of a motif as the vertex after which a dramatic drop in the right-moving probabilities is apparent (expressing the divergence of edges from that vertex), and the beginning of a motif as a dramatic drop in the left moving probabilities (expressing the convergence of edges to that vertex).

Formally, let us define a “decrease ratio”:

$$D_{\text{Right}}(e_i, e_j) = \frac{P_{\text{Right}}(e_i, e_j)}{P_{\text{Right}}(e_i, e_{j-1})}$$

$$D_{\text{Left}}(e_j, e_i) = \frac{P_{\text{Left}}(e_j, e_i)}{P_{\text{Left}}(e_j, e_{i+1})}$$

We will declare  $e_{j-1}$  as a candidate end point of the pattern if  $D_{\text{Right}}(e_i, e_j)$  is smaller than a preset cutoff parameter  $\eta < 1$ . Similarly,  $e_{i+1}$  will be declared as candidate start point of a pattern if  $D_{\text{Left}}(e_j, e_i) < \eta$ .

The statistical significance of the decreases in  $P_{\text{Right}}$  and  $P_{\text{Left}}$  must be evaluated.  $P_{\text{Right}}$  and  $P_{\text{Left}}$  can be regarded as variable-order Markov probability functions. We can define their significance in terms of a null hypothesis stating that  $P_{\text{Right}}(e_i, e_j) \geq \eta P_{\text{Right}}(e_i, e_{j-1})$  and  $P_{\text{Left}}(e_j, e_i) \geq \eta P_{\text{Left}}(e_j, e_{i+1})$ , and require that the p-values of both  $D_{\text{Right}}(e_i, e_j) < \eta$  and  $D_{\text{Left}}(e_j, e_i) < \eta$  be, on average smaller than a preset threshold parameter  $\alpha < 1$ .

A bundle of coinciding paths whose end-points obey these significance conditions is declared as a possibly significant pattern. Given a search path, we calculate both  $P_{\text{Right}}$  and  $P_{\text{Left}}$  from all of the possible starting points, traversing each path leftward and rightward, correspondingly. This technique defines many search-sections, which may be candidates for significant patterns. The most significant ones of these candidates are returned as the outcome patterns for the search path in question.

## 2 EXPRESSION COHERENCE (EC) METHOD

### Rational

The Expression coherence (EC) score is a measure of how clustered a set of genes is in expression space. It may be defined for any gene set for which expression profiles are available. Given an expression profile of N time points for M genes, each gene can be thought of as a point in an N dimensional space, where the ith dimension has the expression level of the gene at the ith time point (Figure 2.1). Given a set of genes, one wishes to determine whether they are tightly clustered, or rather spread "all over the place". One way to accomplish this could be to calculate the "center of mass" of the cloud of genes and then sum over distances of each gene from it (other variations may be to sum over squares of such distances, take standard deviation around that mean etc). Yet, this measure has a clear shortcoming - in cases where the gene cluster is split, say to two, equally sized very tightly clustered subsets, that are yet remote from one another, any deviation-from-mean score will be low. This fails to capture the unique substructure of this gene set, which is composed of two tight sub-clusters. The intuitive reason why this gene set is 'impressive' is that out of  $P=M*(M-1)*0.5$  gene pairs in it  $p=(M/2)*((M/2)-1)$  pairs are close (defined below) to each other. So the ratio  $p/P$  is a good measure for how tight the cluster is. This is the expression coherence score (see figure 2.1 for illustration).

### Biological relevance

Applying the EC score to a set of genes that share a given motif in their promoters, gives a measure of the extent to which the motif may influence expression. Moreover it allows to functionally annotate the motif, by describing the biological conditions in which it governs coherent expression, along with the regulatory effect it exerts (e.g. increased expression in response to a particular stress, or peak in expression at a specific pint during cell cycle ). Such analyses can be performed online for the *S. cerevisiae* genome, via the Motif Analysis Workbench (Lapidot and Pilpel 2003) at <http://longitude.weizmann.ac.il/services.html>

### Algorithm

Given a gene set  $S$ , we compute the Euclidean distances between the centered and variance-normalized expression profiles of each of its  $P=|S|*(|S|-1)*0.5$  gene pairs. The EC score is defined as  $p/P$  where  $p$  is the number of gene pairs whose Euclidian distance is smaller than a threshold distance ( $D$ ).  $D$  is determined based on the distribution of pair-wise distances between expression profiles of all genes in the genome (or more precisely of all genes for which expression level was measured). The original definition of the EC score (Pilpel et al Nat. Genet. 2001) used the 5th percentile as the cutoff for defining "close" expression profiles, but other cutoffs may be applied.

$$EC(S) = \frac{|\{g_i, g_{j \neq i} \in S\}: ExpDist(g_i, g_j) < D|}{|S| * (|S| - 1) \div 2}$$

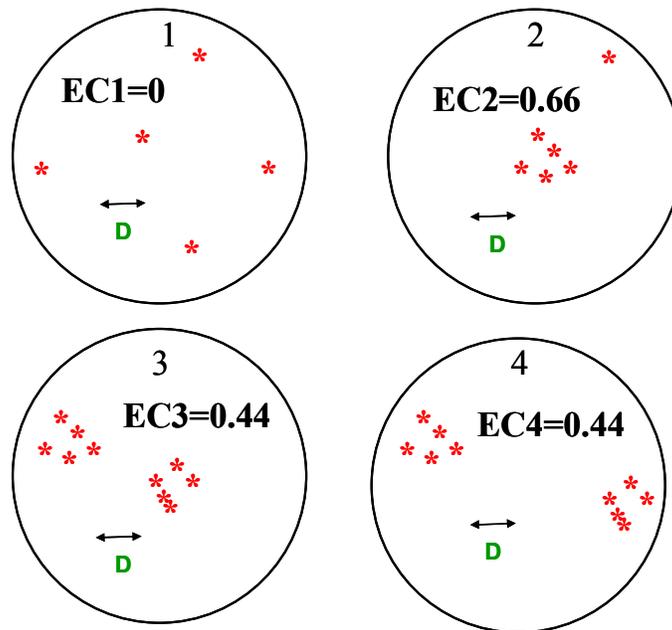


Figure 2.1 EC scores of different scenarios. To illustrate how the EC score is used to measure the extent to which genes are clustered in expression space, we show four scenarios; each disc displays gene sets embedded in a putative expression space. In the first scenario the genes are evenly spread in expression space, no two genes are closer than the threshold distance  $D$ , and thus the EC score is 0. In the second scenario, all genes, but one are tightly clustered and thus the EC score is high, 10/15 possible gene pairs are closer than the threshold  $D$ . In scenarios 3 and 4, two tight clusters are observed. The EC score is the same for both cases, despite the fact that the clusters in scenario 3 are closer to one another. This is because in both cases only genes within each of the two subsets are closer than the threshold (20/45 gene pairs) whereas any two genes belonging to different subsets are further away than the threshold. Measures based on distance from the center of mass would not capture the sub-cluster structure depicted in scenarios 3 and 4.

### Estimation of EC Score Significance

The significance of an EC score calculated for a set of genes, relies on the set size and on the analyzed condition. Thus for each of the expression time series experiments and for each gene set sizes (varying from 3-100 genes), we selected 100,000 random gene sets and computed an EC score for each such set at each cutoff definition. We define the p-value of a given EC score as the fraction of random sets (of the same size and in the same condition) that scored similarly or higher (Note that this sets a lower bound of  $10^{-5}$  on the significance that can be assigned to a given EC score). Since we assume that for a given EC score the probability to get the same score for random sets of genes drops with the set size, gene sets larger than 100 are assigned an upper bound approximated p-value, using the randomly sampled sets of size 100.

For large sets of genes even a small deviation from an EC score of 0.05 (the mean value for random sets) can be statistically significant, whereas for very small sets large deviations from an EC score of 0.05 can be expected purely by chance, as demonstrated in (<http://bioportal.weizmann.ac.il/~lapidotm/rMotif/html/doc/ECscore.html>).

### 3 EXPRESSION DATA

Table 3.1

Experiment short name	Experiment name <sup>1</sup>
Cell cycle (1)	ExpressDB Cho - cell cycle
Cell cycle (2)	ExpressDB Spellman - cell-cycle alpha
Cell cycle (3)	ExpressDB Spellman - cell-cycle cdc15
Cell cycle (4)	ExpressDB Spellman - cell-cycle cdc28
Cell cycle (5)	ExpressDB Spellman - cell-cycle eluteration
Sporulation	ExpressDB Chu - sporulation
MapK	ExpressDB - MapK
Diaux shift	ExpressDB Gasch environmental response - diaux shift
YPD (1)	ExpressDB Gasch environmental response - YPD1
YPD (2)	ExpressDB Gasch environmental response - YPD2
X media vs. car1	ExpressDB Gasch environmental response - x media vrs car1
YPx media vs. car2	ExpressDB Gasch environmental response - YPx media vrs car2
Nitrogen depletion	ExpressDB Gasch environmental response - Nitrogen Deplation
Amino acid starvation	ExpressDB Gasch environmental response - Amino Acid starv
Acid	ExpressDB Environmental response - Acid
Alkali	ExpressDB Environmental response - Alkali
Diamide	ExpressDB Gasch environmental response - diamide
Hydrogen Peroxide (H2O2)	ExpressDB Environmental response - Peroxide
Constant H2O2	ExpressDB Gasch environmental response - constatat h2o2
Menadione	ExpressDB Gasch environmental response - Menadione
NaCl	ExpressDB Environmental response - NaCl
Hypo-osmotic	ExpressDB Gasch environmental response - Hypo-osmotic
DTT (1)	Eisen - dtt
DTT (2)	ExpressDB Gasch environmental response - DTT1
DTT (3)	ExpressDB Gasch environmental response - DTT2
Sorbitol (1)	ExpressDB Environmental response - Sorbitol
Sorbitol (2)	ExpressDB Gasch environmental response - sorbitol
DNA damage	Jelinsky - DNA Damage
Cold	Eisen - cold
Heat shock (1)	ExpressDB Environmental response - Heat
Heat shock (2)	Eisen - heat
Heat shock (3)	ExpressDB Gasch environmental response - 37-25 shock
Heat shock (4)	ExpressDB Gasch environmental response - Heat Shock 1
Heat shock (5) & sorbitol	ExpressDB Gasch environmental response - hs 29-33 1m sorbitol
Heat shock (6)	ExpressDB Gasch environmental response - hs 29-33
Heat shock (7)	ExpressDB Gasch environmental response - hs 29-33 No sorbitol
Heat shock (8)	ExpressDB Gasch environmental response - Heat Shock2 (3 time zero)
Heat shock (9)	ExpressDB Gasch environmental response - hs various temp to 37c
Various temp growth	ExpressDB Gasch environmental response - various temp growth
Various temp steady state	ExpressDB Gasch environmental response - var temp steady state

<sup>1</sup> Whole-genome mRNA expression data of 40 time series experiments in *S. cerevisiae*, were obtained from ExpressDB (<http://arep.med.harvard.edu/cgi-bin/ExpressDByeast/EXDStart>). These time series represent a wide range of natural (e.g. cell cycle) and perturbed conditions. This set of conditions was utilized by us before (Garten et al. Nucleic Acids Res 2005) and a complete list of the conditions is available at [http://longitude.weizmann.ac.il/TFLocation/conditions\\_explist.html](http://longitude.weizmann.ac.il/TFLocation/conditions_explist.html)

## 4 CLUSTERING ALGORITHM

We have formulated an iterative method for clustering motifs, according to their sequences and EC scores information. We first initiate clusters by gathering motifs that share some building blocks, or 'seeds'. Then, a series of iterations improves the clusters, using various procedures detailed below. The clusters refinement steps include the addition and removal of motifs from existing clusters and splitting and merging of clusters. We have quantified the quality of clusters using several criteria associated with sequential patterns and EC score patterns of the motifs. The metrics as well as the refinement steps are listed below.

The clustering algorithm may be used to cluster any type of sequential data that are linked to numerical data. The input to the algorithm is a set of sequences of a given alphabet (e.g. motifs) and a complementary set of vectors (e.g. EC vectors), holding an additional information that needs to be taken into account in the clustering process.

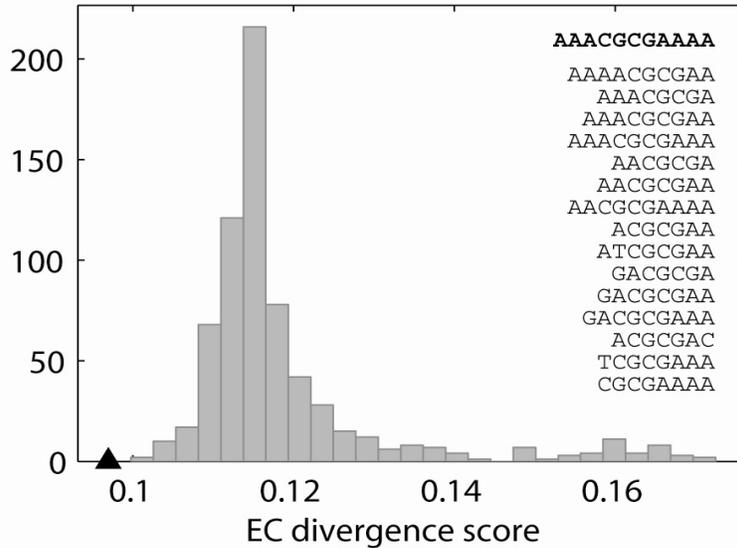
Our clustering method may be considered 'fuzzy' in the sense that single motifs may belong to several clusters. Additionally, not all motifs must be clustered and may be left as singletons.

### Initiating clusters by seeds

Our set of motifs was scanned to find short strings of nucleotides (of length 6) that appear within at least three motifs, to be called 'seeds'. Selecting all motifs that contain a given seed defines a preliminary cluster.

### Pruning clusters to increase EC tightness

For each motif one defines an EC vector of length 40 whose entries specify the p-values of significantly successful EC experiments (that had passed the FDR criterion). Let us define the space of all these vectors as EC space and define an *EC divergence* measure for a cluster of motifs as the average distance of all pairs of its EC vectors. In order to decide whether to eliminate a motif from a given cluster, we ask whether its presence increases the divergence of the cluster. To decide whether a motif  $m$  should be eliminated from a cluster  $M_C$ , we compare the EC divergence of  $M_C$  with the empirical distribution of EC divergence scores resulting from replacing  $m$  with every one of the motifs that lie outside the cluster  $M_C$  (that is, with a background sample  $M_B$ ). The motif will be pruned from the cluster if it does not significantly reduce the cluster's EC divergence, in comparison to motifs from the random background. The deletion of motifs from a cluster occurs after all motifs have been tested, thus the order of tested motifs does not affect their chances of remaining in the cluster. A pseudo code describing this procedure is available in box 4.1. An example is shown in Figure 4.1.



**Figure 4.1** An example for testing the contribution of a specific motif to the cluster’s tightness. The EC-divergence score of the cluster including the motif AAACGCGAAAA (black triangle) is compared to the empirical distribution of EC-divergence of clusters, in which the motif in question has been replaced with random motifs (histogram). Our null hypothesis claims that the motif does not reduce EC-divergence of the group (which is equivalent to saying that the motif harms the tightness of the cluster). In this example, however, the divergence-score of the cluster with the motif included in it is very small. Hence, we can reject the null hypothesis with a probability value of 0.001 and include the motif in the cluster.

#### Box 4.1: Pruning clusters to increase EC tightness

Given a set  $M$  of motifs, their EC scores vectors,  $EC_{i \in M} \in \mathfrak{R}^{40}$ , and two disjoint subsets,  $M_C, M_B \subset M$  (the cluster in question and a background subset of motifs, respectively), we wish to eliminate from cluster  $M_C$  motifs that increase its EC divergence. We will test the contribution of  $M_C$ ’s motifs to its EC divergence, by comparing them to  $M_B$ ’s motifs contribution to  $M_C$ ’s EC divergence.

##### **Pseudo code:**

1. Calculate the EC distance between every pair of EC vectors in  $M_C \cup M_B$ :
 
$$ECdist_{ij} = avg(|EC_i - EC_j|) ; i, j \in M_C \cup M_B$$
2. Calculate  $M_C$ ’s divergence score:  $DivScore_{M_C} = avg(ECdist_{ij}) ; i < j \in M_C$ .
3. Create a new subgroup  $M_{ij}$  by replacing the  $i$ ’th motif of  $M_C$  with the  $j$ ’th motif of  $M_B$ .
4. Calculate  $M_{ij}$ ’s divergence score.
5. Repeat steps 3, 4 for all  $i \in M_C, j \in M_B$  in order to examine the effect of each motif  $i$  on  $M_C$ ’s tightness.
6. For each motif  $i \in M_C$ , generate the empirical distribution of divergence scores, as found in the replacement of motif  $i$  with every motif  $j \in M_B$ .
7. For each motif  $i \in M_C$ , calculate the p-value of getting the divergence score of  $M_C$  by chance.
8. Compare each p-value to a preset significance value  $\alpha$ .
9. Eliminate motifs that are not significantly reducing the divergence of the group in comparison to the randomly sampled motifs.

## Expanding clusters

We search for new motifs to be added to the cluster without increasing its EC divergence. To decide whether a motif  $m$  should be added to a cluster  $M_C$  we compare the EC divergence resulting from its addition ( $M_C+m$ ) with the empirical EC divergence distribution resulting from additions of each of the motifs lying outside  $M_C$  (that is, in a background sample  $M_B$ ), one at a time.

At the same time we also require sequential similarity of the new motif to the ones that belong to the cluster. The sequential distance between motifs is defined as the edit distance of their best alignment, not allowing gaps. The sequential distance score,  $D$ , is normalized between 0 and 1, such that  $D=0$  if the short motif is fully contained in the long one and  $D=1$  if the motifs have no match at all.

A cluster will be expended by motifs that keep its tightness, as well as being strongly similar to the cluster by sequence. The addition of motifs to a cluster occurs after all motifs in  $M_B$  have been tested, thus the order of tested motifs does not affect their chances of being added to the cluster. A pseudo code is available in box 4.2.

### Box 4.2: Expanding clusters

Given a set  $M$  of motifs, their EC scores vectors,  $EC_{i \in M} \in \mathfrak{R}^{40}$  and two disjoint subsets,  $M_C, M_B \subset M$ , we wish to expend  $M_C$  by similar motifs from a background set  $M_B$  that do not increase its EC divergence.

#### *Pseudo code:*

1. Find candidate motifs for addition,  $M_{cand} \subset M_B$ , that show strong similarity by sequence to at least one motif in  $M_C$ .
2. Calculate the EC distance between every pair of EC vectors in  $M_C \cup M_B$ :  
$$ECdist_{ij} = avg(|EC_i - EC_j|) ; i, j \in M_C \cup M_B.$$
3. Create a new candidate subgroup  $M_{C,cand}$  by adding  $M_C$  a single motif from  $M_{cand}$ .
4. Calculate  $M_{C,cand}$ 's divergence score:  $DivScore_{M_{C,cand}} = avg(ECdist_{ij}) ; i < j \in M_{C,cand}$ .
5. Create a new test group  $M_{C,j}$  by adding  $M_C$  a single motif from  $M_B \neq M_{cand}$ .
6. Calculate  $M_{C,j}$ 's divergence score.
7. Repeat steps 5, 6 for all  $j \in M_B \neq M_{cand}$  to generate the empirical distribution of divergence scores of the test groups.
8. Calculate the probability value for getting the divergence score of  $M_{C,cand}$  by chance.
9. Expend the group by the current motif candidate if it produces a significantly low divergence score (lower than some preset significance value,  $\alpha$ ).
10. Repeat steps 3-9 for every motif candidate in respect to the original cluster  $M_C$ .

### Fusion of clusters

Clusters will be merged if they share a minimum percentage of motifs and are also found to be similar in EC. EC distance between two clusters A and B is defined by a *Fisher* criterion, as the distance between the centers of the clusters, divided by the sum of their standard deviations:

$$F_{A,B} = \frac{\|\mu_A - \mu_B\|}{\|\sigma_A\| + \|\sigma_B\|}$$

$\mu_A$  and  $\mu_B$  are the mean EC vectors of the two EC matrices (the center of each cluster). For each cluster we define  $\sigma$  as the vector of the 40 standard deviations corresponding to the 40 EC experiments. Clusters will be merged if their *Fisher distance*,  $F$ , is smaller than some threshold, as long as they also obey the sequential similarity criterion.

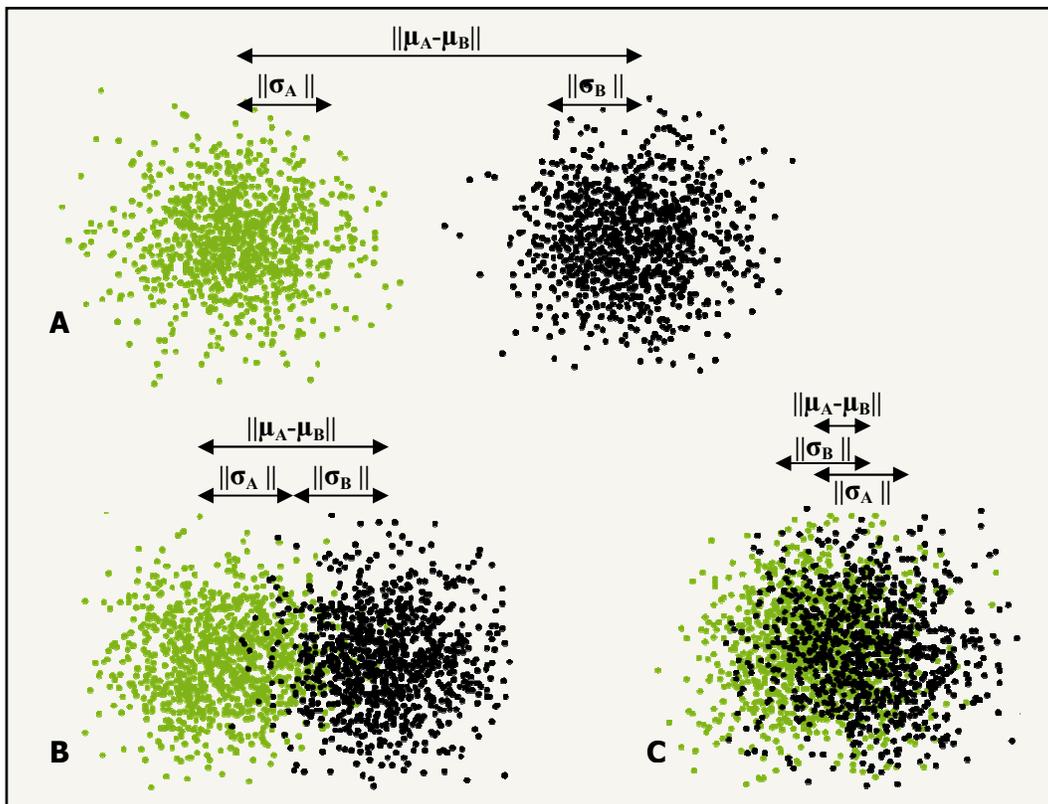


Figure 4.2 A demonstration of the Fisher criterion. The fisher distance for distant clusters exceeds the value 1 (A). The smaller the fisher distance is, the more difficult it gets to distinguish between the clusters (B, C).

### Splitting of clusters

Clusters will be split into K smaller clusters if they exceed a given size. Splitting is done using the K-means algorithm on the EC space of the cluster. After applying this indiscriminative step, however, a fusion step is applied, so that unnecessary splitting will be reversed. In this work we have used K=3.

### Fine refinement of clusters

The former procedures are applied iteratively in a preset order, to generate clusters that are rather tight in EC and in sequence and differ from each other in sizes, EC patterns and motif sequences. In a final pruning step finer parameters are used. Then the improvement of each cluster is tested with respect to a *cluster score*, assessing the quality of the cluster, and the pruning is accepted or rejected accordingly. The clusters are given a *cluster score*, a heuristic function encapsulating the various measures used in the analysis:

$$ClusterScore = \frac{\sqrt{\frac{SeqDivScore_{cluster}}{SeqDivScore_{others}} \cdot \frac{DivScore_{cluster}}{DivScore_{others}}}}{(MC \cdot F_{cluster,others})^2}$$

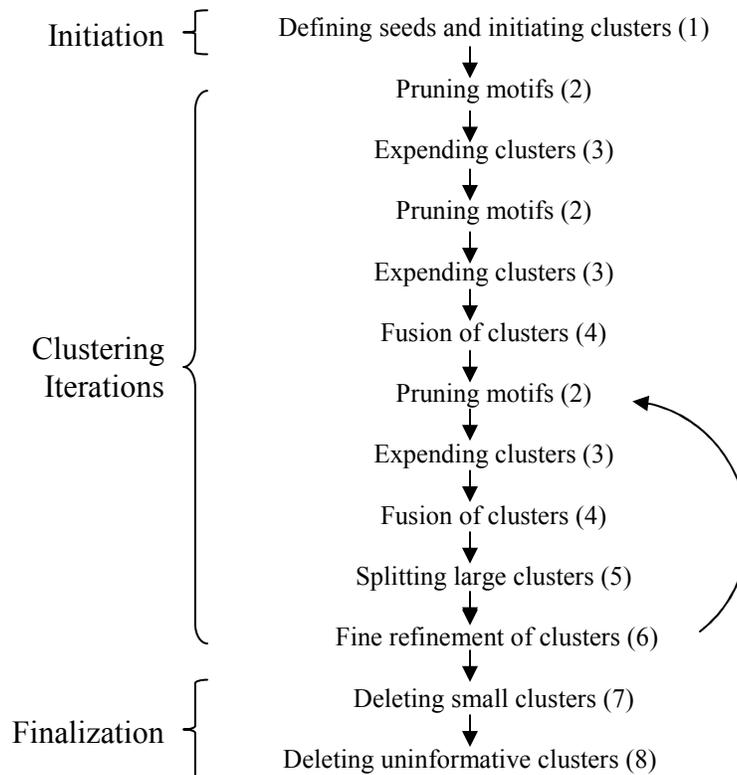
$SeqDivScore_{cluster}$ ,  $DivScore_{cluster}$  are the cluster's sequential and EC divergence scores, respectively (the former is defined similarly to the later, as the average sequential distance of all pairs of motifs within the cluster).  $SeqDivScore_{others}$ ,  $DivScore_{others}$  are the sequential and EC divergence scores of all the motifs outside the cluster, respectively.  $F_{cluster,others}$  is the EC fisher distance between the cluster and the rest of the motifs, and  $MC$  is the number of motifs within the cluster. The smaller the cluster score is, the better the quality of the cluster is considered.

The cluster score quantifies the quality of a cluster in terms of its internal tightness relatively to the background. As affected by many different factors, the cluster score is sensitive to noise. Hence it is only used at a late stage along the algorithm, when clusters are already coherent to a great extent.

### Flow of the algorithm

After initiation, cycles of the various iterations occur, gradually improving the clusters with respect to their sequences and EC patterns. The algorithm stops when the rate of change of the clusters falls below a certain cutoff (a stopping criterion) or if no clusters are found. Clusters that are too small (below a preset threshold) are disregarded.

The order of iterations used for clustering our motifs is described below:



## Results

### 5 CLUSTERS

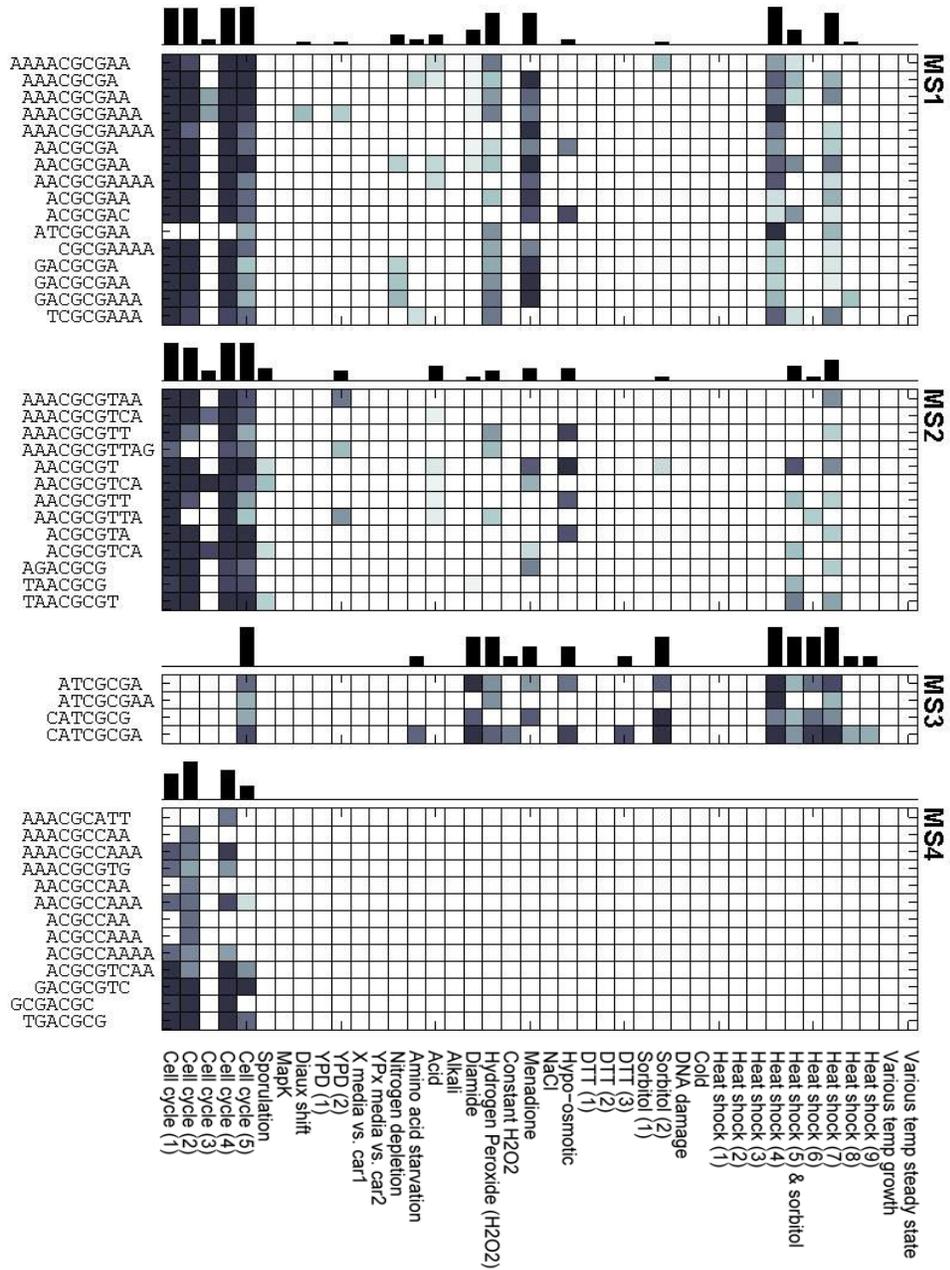
We have clustered a distilled set of motifs, consisting of 694 motifs that have passed the FDR criterion and have also had at least one EC success with a p-value of 0.001 or lower. Our algorithm finds 20 clusters, covering a total of 182 motifs. 14 of our clusters have large overlaps with known clusters.

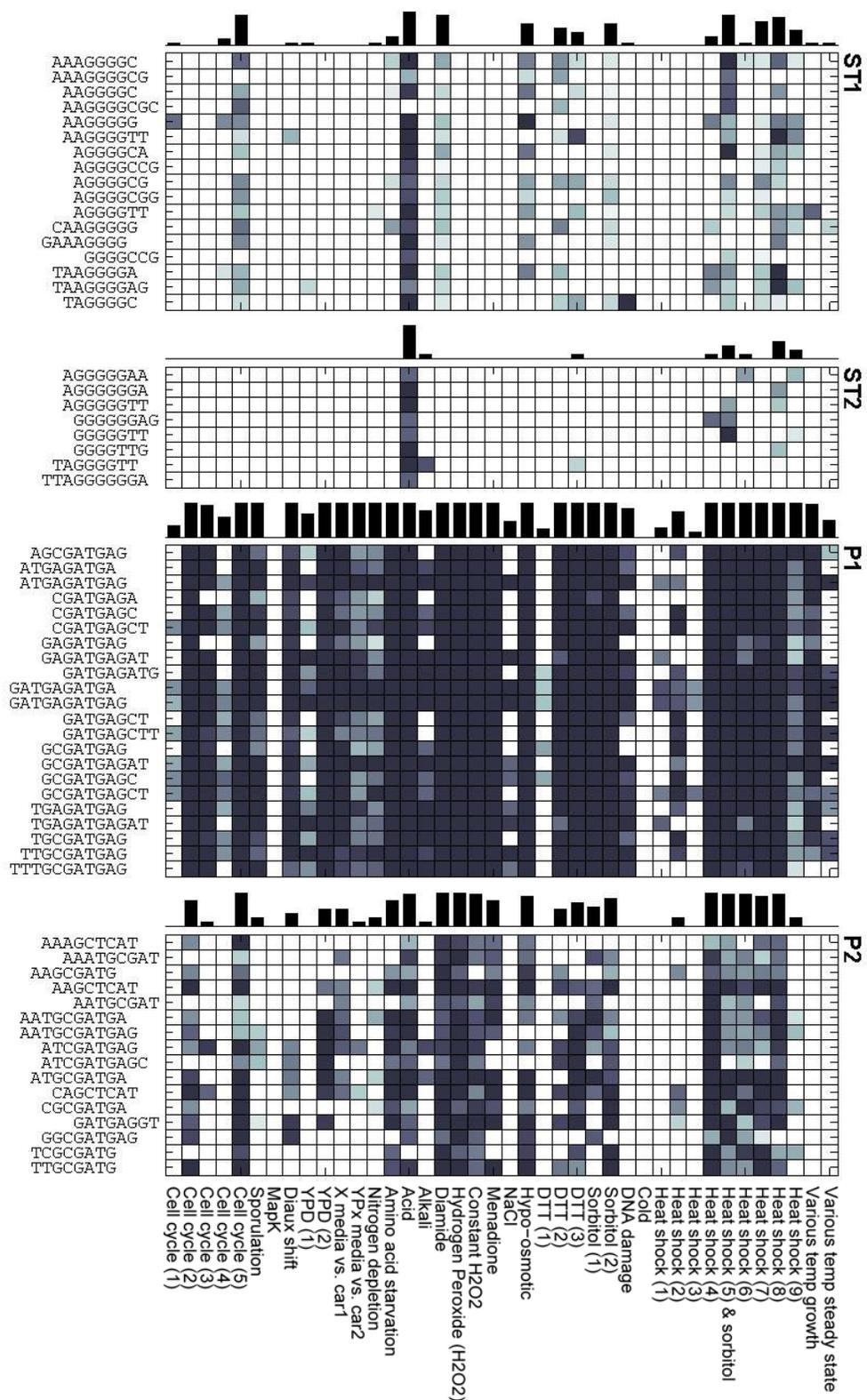
Table 5.1

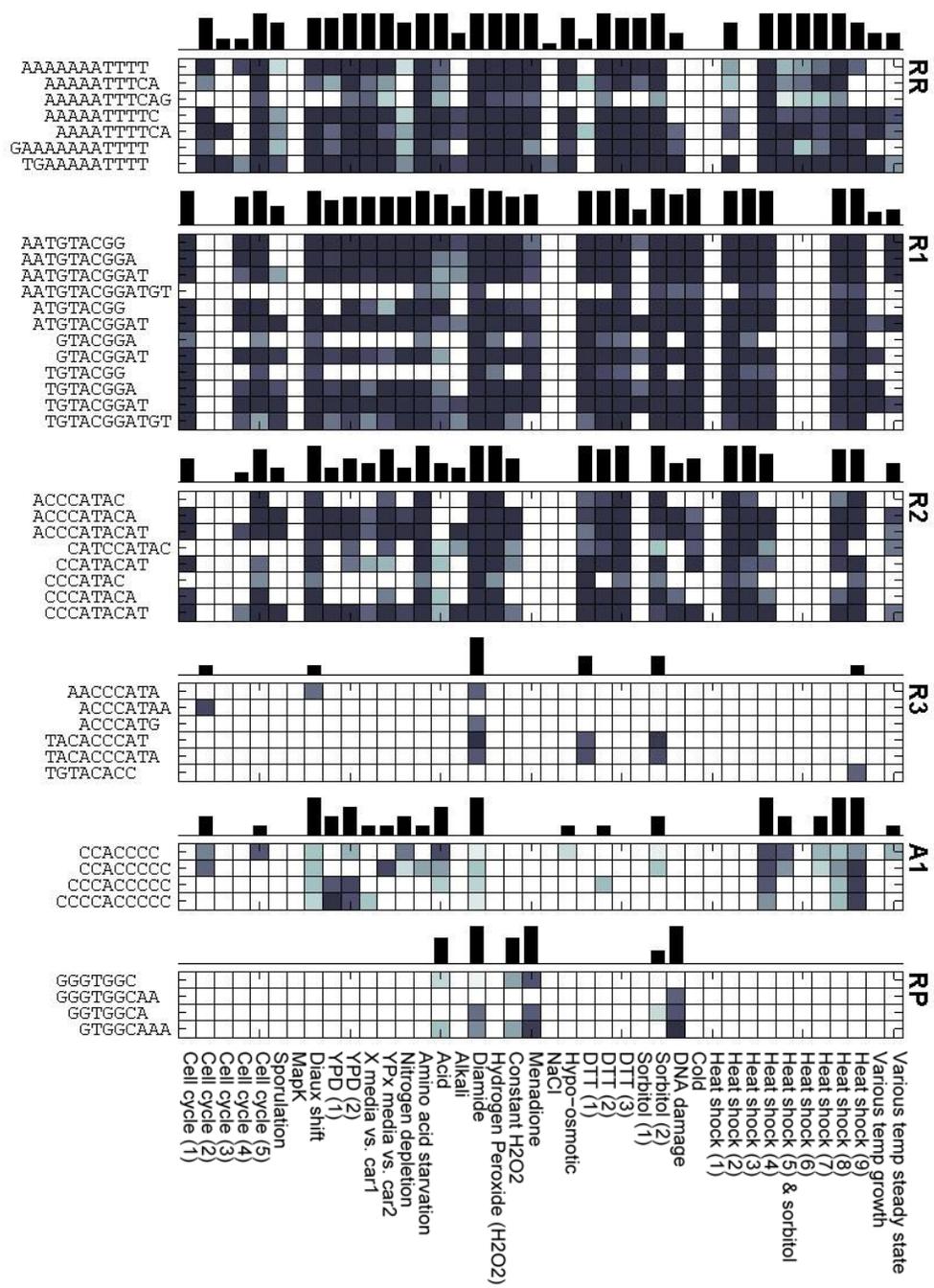
Cluster name	Identified as
MS1	MBF / SBF
MS2	MBF / SBF
MS3	MBF / SBF
MS4	MBF / SBF
ST1	STRE
ST2	STRE
P1	PAC
P2	PAC
RR	RRPE
R1	RAP1
R2	RAP1
R3	RAP1
A1	ADR1 / STRE
RP	RPN4
C15	Unknown
C16	Unknown
C17	Unknown
C18	Unknown
C19	Unknown
C20	Unknown

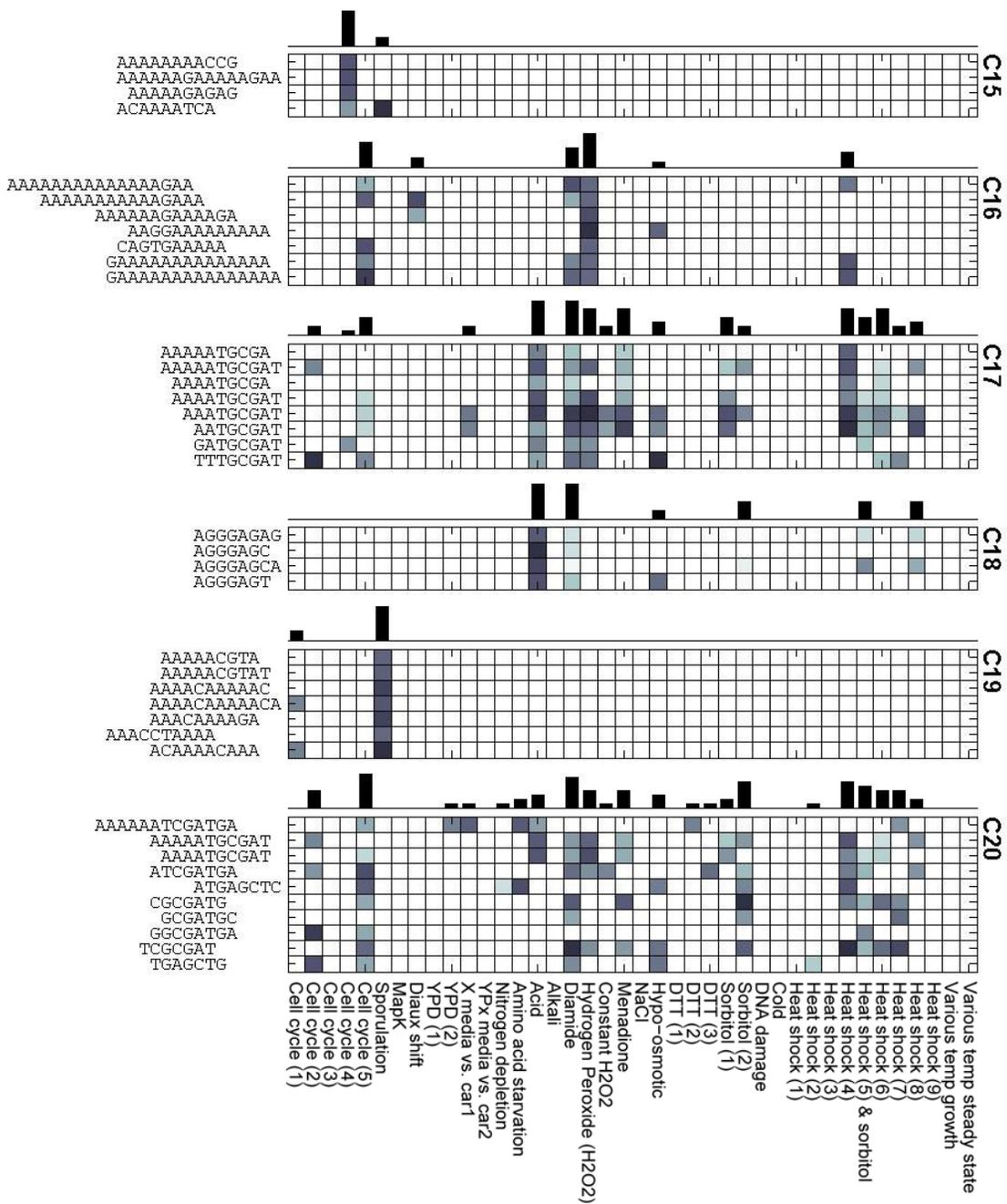
## 6 EC PATTERNS OF CLUSTERS

Figure 6.1





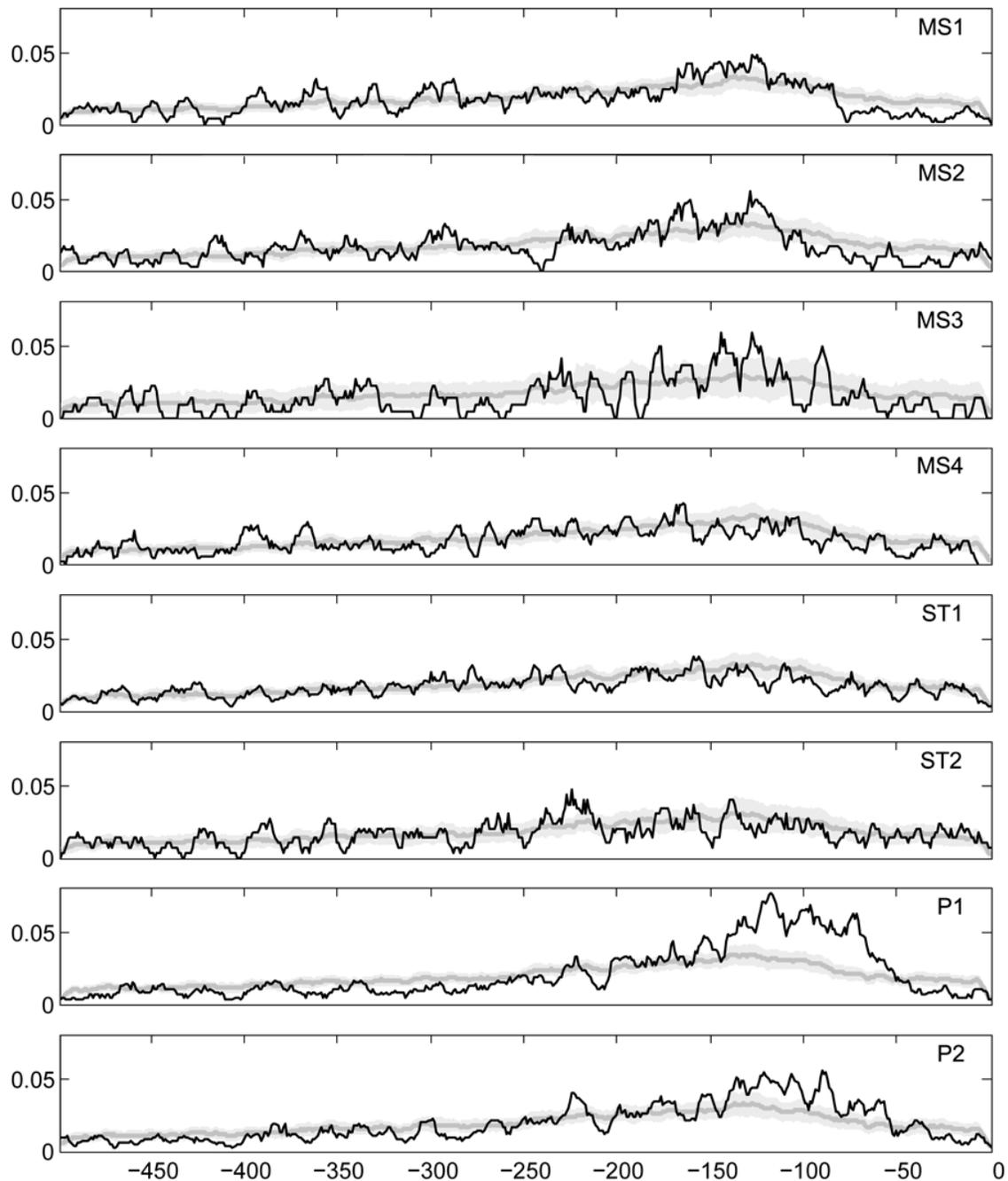


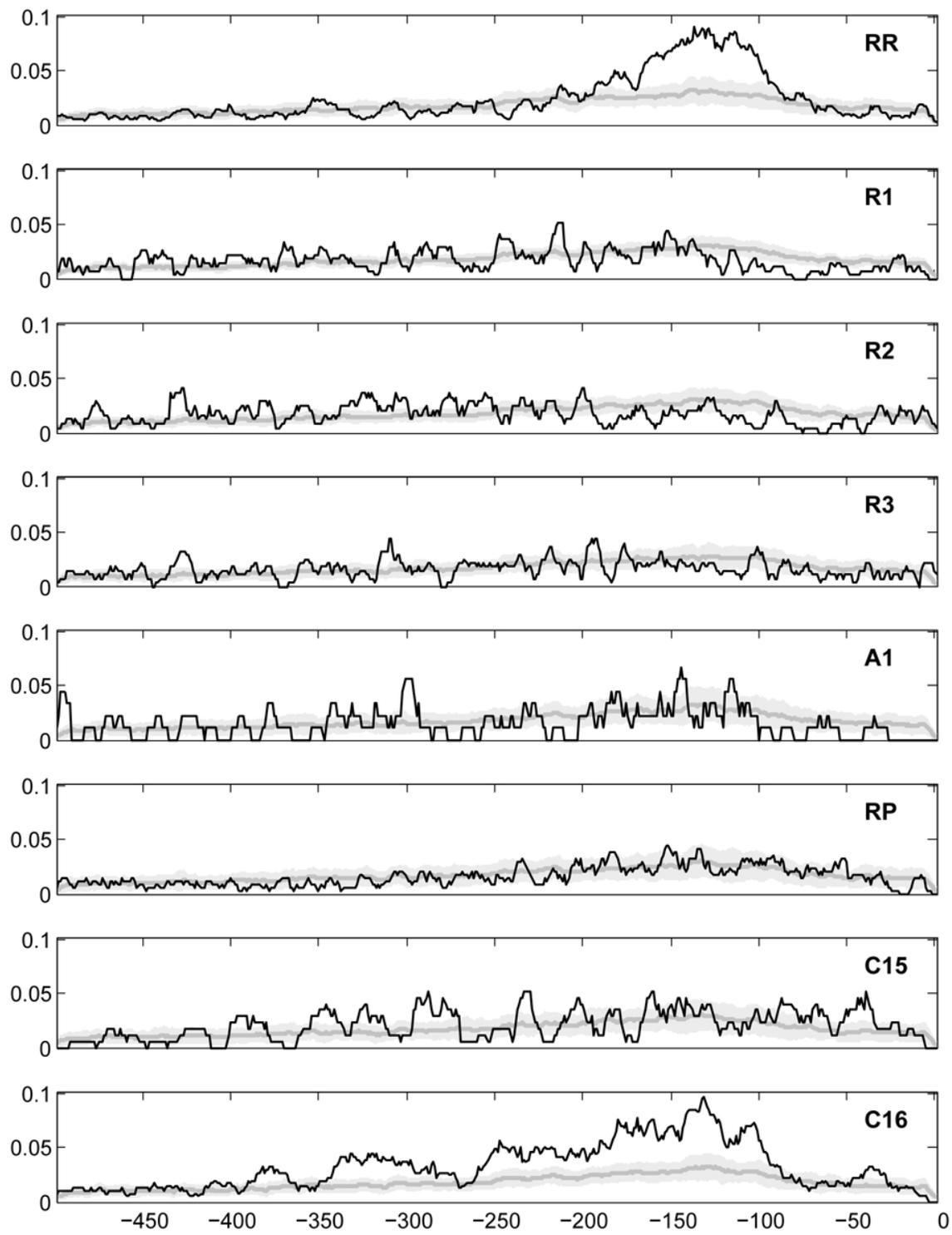


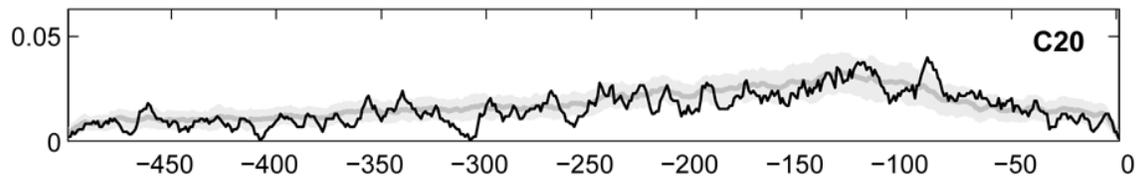
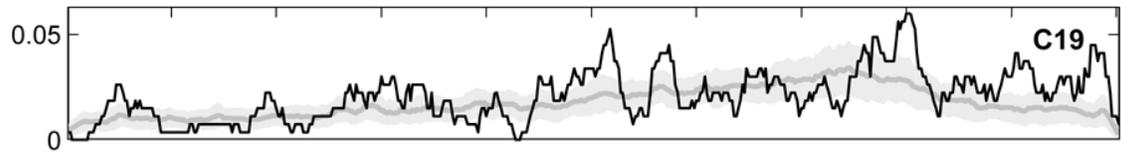
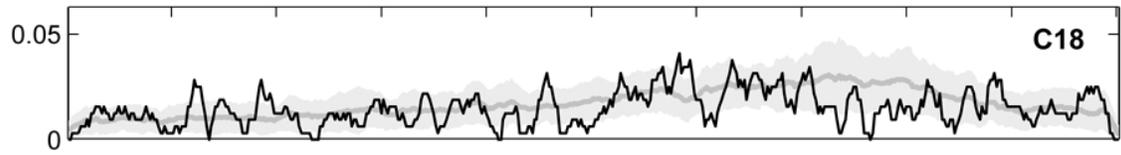
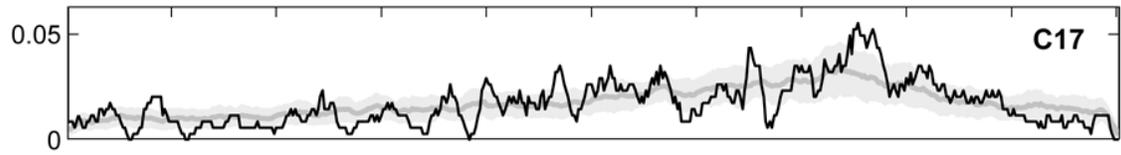
## 7 LOCALIZATION OF MOTIFS ALONG PROMOTERS

Black lines indicate, for each position upstream to the genes (up to -500bp), the percentage of promoters on which the cluster's sequences have been found. This can be compared to the localization of sampled groups of motifs (of the same sizes as those of the clusters in question) out of the set of 694 motifs. For each cluster, the dark gray line shows the mean motif occurrence per position over 1000 such sampled groups, while the light gray area represents the samples' standard deviation of occurrences per position.

Figure 7.1





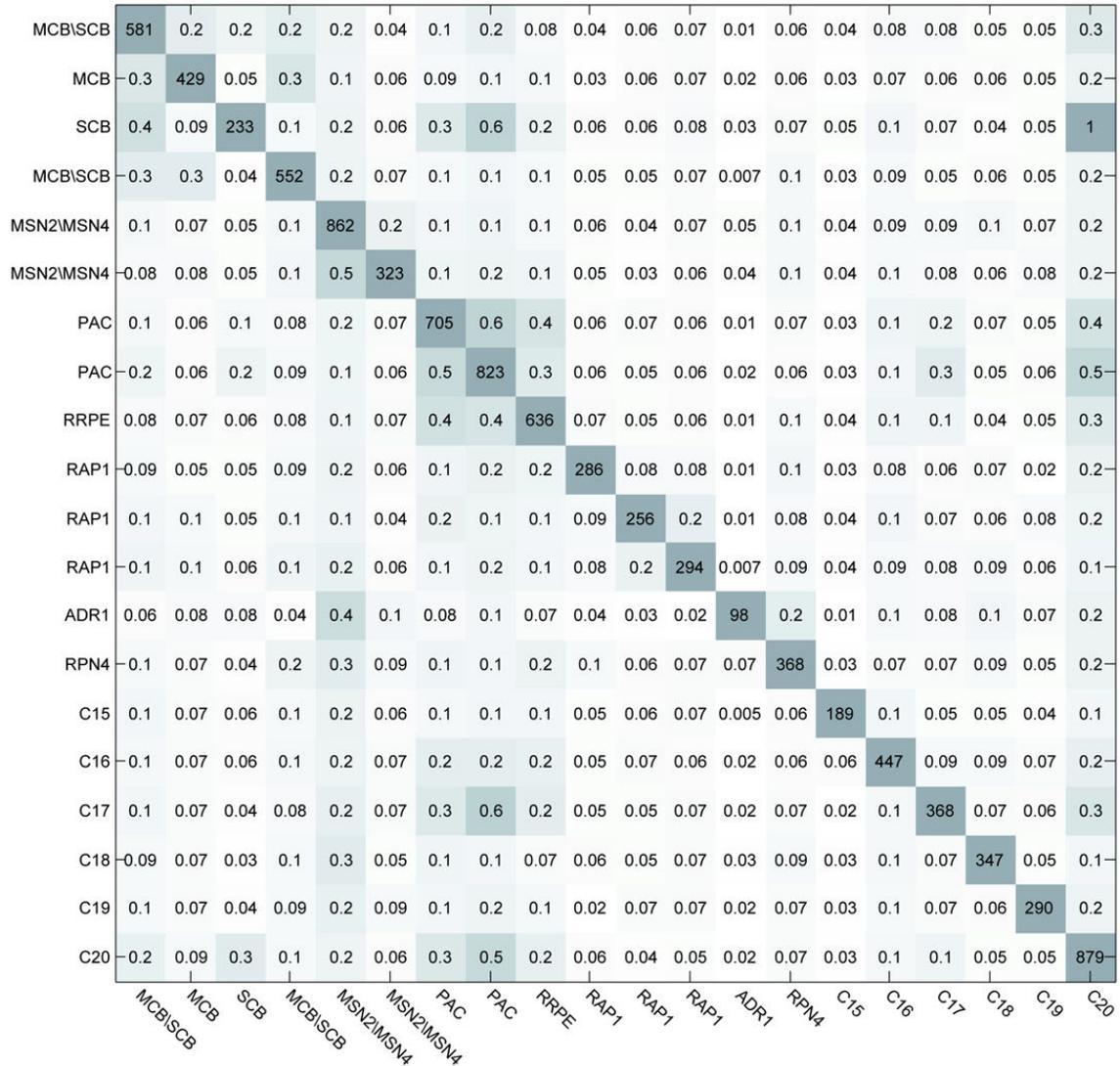


## 8 FISHER DISTANCES BETWEEN CLUSTERS

MS1	0.25	0.48	0.69	0.87	0.84	1.1	1.6	0.9	1.3	1.2	1.1	1	0.85	1	1.3	0.88	0.7	1.2	1.5	0.55
MS2	0.48	0.27	0.78	0.55	0.73	0.9	1.6	0.97	1.3	1.2	1.1	0.86	0.85	0.89	0.96	0.85	0.75	0.97	1.1	0.55
MS3	0.69	0.78	0.51	1	0.55	0.9	1.3	0.58	0.98	1.1	0.94	0.82	0.68	0.83	1.2	0.69	0.47	0.9	1.2	0.31
MS4	0.87	0.55	1	0.26	0.94	0.91	2.2	1.3	1.7	1.6	1.4	0.74	0.98	0.92	0.88	0.95	0.93	1.1	1.2	0.67
ST1	0.84	0.73	0.55	0.94	0.24	0.6	1.4	0.69	1	1.1	0.91	0.71	0.54	0.73	1.1	0.76	0.52	0.52	1.1	0.39
ST2	1.1	0.9	0.9	0.91	0.6	0.35	2	1.1	1.5	1.4	1.2	0.65	0.75	0.68	1	0.85	0.67	0.42	1	0.58
P1	1.6	1.6	1.3	2.2	1.4	2	0.21	0.83	0.37	0.9	1	2	1.2	1.8	2.6	1.9	1.3	2	2.6	1.2
P2	0.9	0.97	0.58	1.3	0.69	1.1	0.83	0.25	0.54	0.88	0.81	1.1	0.7	1.1	1.6	1.1	0.57	1.1	1.6	0.57
RR	1.3	1.3	0.98	1.7	1	1.5	0.37	0.54	0.38	0.76	0.79	1.5	0.89	1.4	2	1.5	0.96	1.5	2	0.93
R1	1.2	1.2	1.1	1.6	1.1	1.4	0.9	0.88	0.76	0.28	0.36	1.4	0.9	1.3	1.8	1.4	1	1.5	1.8	0.97
R2	1.1	1.1	0.94	1.4	0.91	1.2	1	0.81	0.79	0.36	0.34	1.2	0.78	1.2	1.6	1.1	0.95	1.3	1.6	0.84
R3	1	0.86	0.82	0.74	0.71	0.65	2	1.1	1.5	1.4	1.2	0.41	0.78	0.56	0.82	0.64	0.73	0.62	0.83	0.5
A1	0.85	0.85	0.68	0.98	0.54	0.75	1.2	0.7	0.89	0.9	0.78	0.78	0.53	0.84	1.2	0.84	0.63	0.76	1.2	0.54
RP	1	0.89	0.83	0.92	0.73	0.68	1.8	1.1	1.4	1.3	1.2	0.56	0.84	0.53	1	0.79	0.64	0.66	1	0.54
C15	1.3	0.96	1.2	0.88	1.1	1	2.6	1.6	2	1.8	1.6	0.82	1.2	1	0.71	1.1	1.1	1.2	1.3	0.81
C16	0.88	0.85	0.69	0.95	0.76	0.85	1.9	1.1	1.5	1.4	1.1	0.64	0.84	0.79	1.1	0.4	0.64	0.88	1.1	0.45
C17	0.7	0.75	0.47	0.93	0.52	0.67	1.3	0.57	0.96	1	0.95	0.73	0.63	0.64	1.1	0.64	0.35	0.62	1.1	0.31
C18	1.2	0.97	0.9	1.1	0.52	0.42	2	1.1	1.5	1.5	1.3	0.62	0.76	0.66	1.2	0.88	0.62	0.58	1.2	0.54
C19	1.5	1.1	1.2	1.2	1.1	1	2.6	1.6	2	1.8	1.6	0.83	1.2	1	1.3	1.1	1.1	1.2	0.46	0.82
C20	0.55	0.55	0.31	0.67	0.39	0.58	1.2	0.57	0.93	0.97	0.84	0.5	0.54	0.54	0.81	0.45	0.31	0.54	0.82	0.31
	MS1	MS2	MS3	MS4	ST1	ST2	P1	P2	RR	R1	R2	R3	A1	RP	C15	C16	C17	C18	C19	C20

Figure 8.1 Fisher distances between clusters. On the diagonal (where  $F=0$ ) we have added the mean  $F$ -values obtained by randomly dividing each of the clusters into two arbitrary ones (mean over 100 random divisions for each cluster).

## 9 INTERSECTIONS BETWEEN SETS OF GENES OF COUPLES OF CLUSTERS



**Figure 9.1** Intersections between sets of genes of couples of clusters. The percentage of common genes between cluster  $i$  and cluster  $j$ , out of cluster  $i$  ( $i$  in rows,  $j$  in columns).

$i=j \rightarrow$  The number of genes on the promoter of which the clusters motifs are found (the color on the diagonal is set as 100%, in respect to the rest of the matrix).

*i*

MCB\SCB	581	138	98	142	91	26	74	124	48	25	32	40	6	35	25	49	48	31	29	174
MCB	138	429	22	145	61	27	40	52	45	13	26	32	8	25	13	30	24	24	21	79
SCB	98	22	233	23	39	15	69	130	41	14	14	19	8	16	12	28	16	9	12	233
MCB\SCB	142	145	23	552	92	39	56	73	53	25	29	40	4	58	19	50	28	33	25	92
MSN2\MSN4	91	61	39	92	862	167	106	119	86	49	35	59	44	94	33	77	77	94	59	138
MSN2\MSN4	26	27	15	39	167	323	48	50	45	16	10	18	12	34	12	32	26	19	26	50
PAC	74	40	69	56	106	48	705	396	252	41	48	39	8	51	20	101	110	49	38	269
PAC	124	52	130	73	119	50	396	823	234	53	38	47	13	48	23	102	218	44	47	436
RRPE	48	45	41	53	86	45	252	234	636	44	33	35	7	66	23	91	88	25	31	183
RAP1	25	13	14	25	49	16	41	53	44	286	22	23	4	36	10	23	17	21	5	54
RAP1	32	26	14	29	35	10	48	38	33	22	256	55	3	21	11	31	19	16	20	39
RAP1	40	32	19	40	59	18	39	47	35	23	55	294	2	25	13	26	24	25	19	42
ADR1	6	8	8	4	44	12	8	13	7	4	3	2	98	24	1	10	8	12	7	18
RPN4	35	25	16	58	94	34	51	48	66	36	21	25	24	368	12	25	26	32	20	65
C15	25	13	12	19	33	12	20	23	23	10	11	13	1	12	189	25	9	9	8	25
C16	49	30	28	50	77	32	101	102	91	23	31	26	10	25	25	447	40	38	30	95
C17	48	24	16	28	77	26	110	218	88	17	19	24	8	26	9	40	368	24	21	123
C18	31	24	9	33	94	19	49	44	25	21	16	25	12	32	9	38	24	347	17	45
C19	29	21	12	25	59	26	38	47	31	5	20	19	7	20	8	30	21	17	290	45
C20	174	79	233	92	138	50	269	436	183	54	39	42	18	65	25	95	123	45	45	879
		MCB	SCB	MCB\SCB	MSN2\MSN4	MSN2\MSN4	PAC	PAC	RRPE	RAP1	RAP1	RAP1	ADR1	RPN4	C15	C16	C17	C18	C19	C20

*j*

**Figure 9.2 Intersections between sets of genes of couples of clusters.**

$i \neq j \rightarrow$  The number of common genes between cluster  $i$  and cluster  $j$ .

$i=j \rightarrow$  The number of genes on the promoter of which the clusters motifs are found.

## 10 COMPARING MEX TO ALTERNATIVE APPROACHES

### **Comparison of applying the MEX algorithm, followed by EC to an exhaustive k-mer enumeration followed by EC**

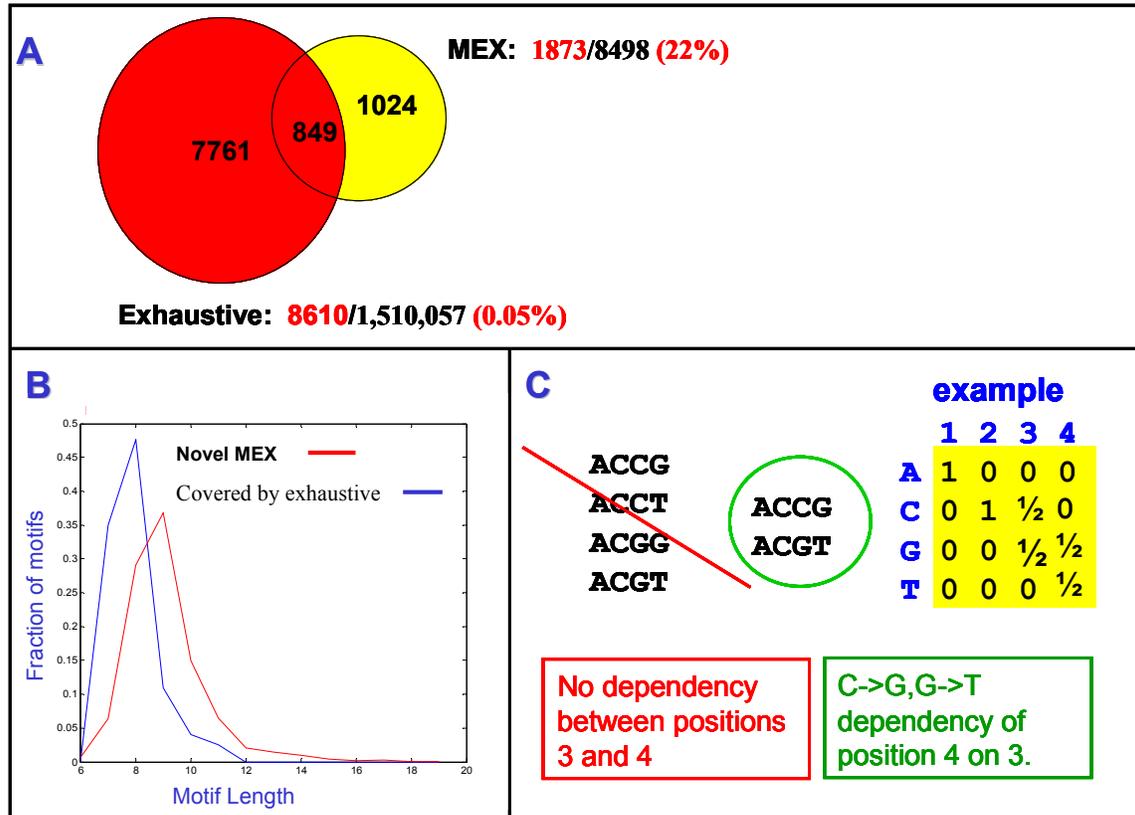
1,873/8,498 (22 %) of MEX's extracted motifs had a significant EC score (passed FDR of 0.1) in at least one of the examined biological conditions. For comparison, in an exhaustive enumeration of all k-mers of length 7-11 residing in yeast promoters (Shalgi et al 2005, Lapidot & Pilpel in prep), we produced 1,510,057 hypotheses, only 8,610(0.6%) of which scored significantly under the same FDR cutoff. In other words we see a striking enhancement in the probability of a k-mer to pass an EC test if it was pre-selected by MEX as a motif that obeys some syntactic rules.

There was an overlap of 849 motifs between the motif sets obtained by the two approaches. 1024 (55%) of MEX's motifs were not discovered by the exhaustive approach (Figure 10.1A). These are mostly weaker motifs, that could not be identified within a very noisy background; MEX provides an enrichment in signal which relaxes the p-value thresholds set by FDR, allowing for weaker motifs to be detected as significant. In addition, MEX extracted sequence motifs of length up to 19 nucleotides. 57 of the unique MEX motifs were longer than 11 bases, and thus were not even scanned by the exhaustive approach (Figure 10.1B). Expanding the exhaustive enumeration to larger sequence lengths is extremely expensive computationally, whereas MEX is easily scalable to longer sequences and to larger genomes. Motifs that were detected by the exhaustive approach, but not by MEX most likely do not obey the inherent position dependencies, selected for by MEX (figure 10.1C). It has been reported that some, but not all functional TFBS display such position dependencies (Tomovic et al. Bioinformatics 2007). The relative success of MEX in identifying high scoring motifs suggests however that there are some syntactic rules that characterize functional binding sites. This thought is intriguing because it implies that we may be able to identify at least some of the TF binding sites based on their sequence context alone.

### **Coverage of the Harbison PSSM set**

To estimate the comprehensiveness of our MEX-extracted motif set, we examined its coverage of the well accepted reference PSSM set published by Harbison et al. We scanned each of the 1,873 motifs of the MEX-extracted set, as well as each of the 8,610 motifs of the exhaustive enumeration set against all 102 Harbison PSSMs. We applied a scoring method that assesses how likely a given k-mer is to be generated by a given PSSM. For each pair of known PSSM and k-mers we summed up the frequencies corresponding to the nucleotides observed in the k-mer, over all PSSM relevant positions. This score was then scaled to the range [0-100] by subtracting the minimal possible score that may be obtained from the PSSM (that is, the score obtained for a k-mer that corresponds to the least frequent position in each column of the PSSM) and dividing by the range of possible scores (obtained after additionally identifying the maximal possible scoring k-mer from that PSSM). Applying a cutoff of 99% identity, 16% of the exhaustive set provide a coverage of 91% of Harbison's PSSMs, and 45% of the MEX set, cover 66% of the Harbison set. Namely MEX is not as comprehensive as the exhaustive set, but it is enriched in signal and contains less false positives. Table 10.1

summarizes the coverage of the Harbison set by the MEX-extracted motifs when applying different cutoffs.



**Figure 10.1** Comparison between motif sets obtained by the exhaustive approach versus the syntax based approach (MEX). **A.** Overlap between motif sets obtained by exhaustive enumeration and by MEX. **B.** Length distribution of MEX motifs that are covered by the exhaustive search versus novel MEX motifs. MEX has a clear advantage in identifying longer motifs, for which the exhaustive search is computationally too demanding. There are also motifs in the length range of 7-11 which scored significantly when pre-selected by MEX, and not in the exhaustive dictionary. These motifs are weaker and when embedded in a very noisy background (of all possible k-mers), their score is not high enough to pass the threshold set by FDR. **C.** Motifs in which there are no clear inter-position dependencies will be missed by MEX. MEX learns simple syntax rules from the promoter sequences and searches for motifs that obey these rules.

Comparison score cutoff	Coverage of MEX motif set	Coverage of known	Unique known clusters
99%	837/1873=45%	63/102=61%	51/77=66%
98%	866/1873=46%	68/102=67%	56/77=72%
97%	904/1873=48%	73/102=72%	60/77=78%
95%	996/1873=53%	79/102=77%	65/77=84%

**Table 10.1 Coverage of the Harbison PSSM set by the significant motifs extracted by ME. A scoring method was devised to assess how likely a given string is to be generated from a given PSSM. The score is on a scale of 0 to 100. It is computed by summing up the frequencies corresponding to the observed nucleotides over all motif positions, and normalizing this score to a scale of 0-100. The scaling is done by subtracting the minimal possible score and dividing by the range of possible scores. For example for the PSSM [A: 0.0191 0.0191 0.9733 0.9733 0.0120, C:0.9500 0.9500 0.0074 0.0074 0.0074, G: 0.0117 0.0117 0.0074 0.0074 0.0074 T:0.0191 0.0191 0.0120 0.0120 0.9733] the lowest possible score 0.0455 is obtained for the string GG(C/G)(C/G)(C/G), the highest possible score 4.8198 is obtained for the string CCAAT. After scaling GGCCC will score 0%, CCAAT will score 100% and CCATT will score 79.9% ( (3.8585-0.0455)/(4.8198-0.0455) ). We computed this score for each of the significant (based on the EC score analysis) motifs extracted by MEX versus all 102 Harbison PSSMs. The coverage of Harbison's motif set was assessed for several different score cutoffs. Note that one motif may match more than one Harbison PSSM, because of redundancy in Harbison's dataset. For comparison 16% to 25% of the exhaustive set provide a coverage of 91% to 97% of the Harbison set, depending on the cutoff employed. There are two very long (17 and 18 positions) gapped motif in Harbison's set, for which we have no match, because the dictionary only covers motifs of length 7-11.**