

Supplementary materials to the following article:

Deciphering Principles of Transcription Regulation in Eukaryotic Genomes

Dat H. Nguyen and Patrik D'haeseleer

Department of Genetics, Harvard Medical School, Boston, MA 02115

All correspondences should be sent to Dat H. Nguyen
(dnguyen@genetics.med.harvard.edu)

Supplementary 1:

Figures 3a and 3b in the main text show that the expression correlation between genes that contain the PAC motif quickly drops to zero the further the PAC motif is located from the start codon, whereas MED predicts a significant influence for the PAC motif over a much wider range. Applying cluster analysis on the genes containing PAC spaced between 600 and 750 bp upstream, we find that their expression profiles follow two somewhat anti-correlated patterns (Fig. SF1a). Similarly, the genes containing PAC spaced between 750 and 1000 bp upstream cluster into three distinct expression patterns (Fig. SF1b), two that are anti-correlated and a third one with low expression levels. This clearly demonstrates that lack of co-expression does not necessarily imply lack of co-regulation, because average pairwise correlation is a poor measure to infer motif strength accurately as mentioned in the main text. In contrast, for a very tightly co-expressed gene ensemble (Fig SF1c), average correlation does accurately reflect the degree of co-regulation caused by the set of motifs in the ensemble.

Supplementary 2: Apparent PAC/RRPE synergetic behavior on gene expression

Expression data for PAC-only, RRPE-only, and PAC/RRPE-only gene ensembles are presented in Fig SF2. In this figure, we show both an average gene expression for each gene ensemble for each condition derived from experimental data (SF2a) and the corresponding one derived from MED using Eq. 1 (SF2b). These figures show a consistent feature that when a gene contains both PAC and RRPE motifs, its expression signal is larger in magnitude compared to a gene that contains either motif alone. Note that the term “only” here refers to the fact that the gene ensemble of one particular motif contains only such motif and nothing else. This choice of gene ensembles (instead of PAC, RRPE and PAC/RRPE ensembles) eliminates side effects of other motifs on gene expression pattern of the corresponding gene ensemble, thereby allowing us to pinpoint the root cause of apparent enhancement in expression level of PAC/RRPE-only containing gene ensemble compared to PAC-only and RRPE-only containing gene ensembles.

In Fig. SF2c, the strengths of the PAC and RRPR motifs in the gene ensemble containing both of them are plotted as a function of RRPE motif's position averaged over

all PAC positions. The two curves in this Figure SF2-c are complement to the two black diamond curves shown in Figure 3a in the main text.

Supplementary 3: PAC/RRPE motif data

Figure SF3-A shows the position distribution of PAC and RRPE in associated promoters of genes containing each. In Figure SF3-B, the distribution of expression correlation coefficients for a set of instances of PAC- and RRPE-containing gene ensembles is shown.

Supplementary 4: Proof of uniqueness (see supplementary 11 for definition of mathematical notations)

As shown in Eq. 2 in the main text, the first part of the MED algorithm decomposes the gene expression data matrix \mathbf{E} (m-genes by n-conditions) in order to obtain motif matrix \mathbf{M} (m-genes by k-motifs) and the matrix of proxy regulator activity \mathbf{A} (k-motifs by n-conditions) according to the following equation:

$$E \approx M \bullet A \tag{2}$$

For any invertible $k \times k$ matrix \mathbf{X} , Eq. 2 above can always be written as follows:

$$E \approx M \bullet A = M \bullet X \bullet X^{-1} \bullet A = M^* \bullet A^* \tag{S1}$$

In order to prove that matrices \mathbf{M} and \mathbf{A} are unique, it is sufficient to prove two things. First we prove that the MED formalism restricts \mathbf{X} to be an identity matrix, provided that (1) \mathbf{M} and \mathbf{A} are full rank, and (2) the matrix \mathbf{M} is sparse and rectangular (i.e. there must be at least $k-1$ zero-matrix elements in each column in general, and the number of rows is larger than the number of columns). Secondly, when rows of matrix \mathbf{A} are linearly dependent or near-linearly dependent, we prove that MED formalism can still produce unique motif matrix \mathbf{M} regardless of rank deficiency of matrix \mathbf{A} , due to the second term of Eq. 3 or 4.

Proof

1. Matrix X is an identity matrix as an outcome of the MED formalism:

In Eq. 1 in the main text:

$$E_{gc} \approx \sum_{j \in \Omega_g} M_{gj} A_{jc}$$

if $j \notin \Omega_g$ (i.e. motif j does not exist in the promoter of gene \mathbf{g}) then the corresponding matrix element M_{gj} of the motif matrix \mathbf{M} must be zero and remain so and as it is always the case in MED. As a result, for each matrix element M_{gj} of matrix \mathbf{M} whose value is zero, the corresponding matrix element M_{gj}^* of matrix \mathbf{M}^* must also be zero after the

linear transformation with matrix X . Therefore, for a zero element M_{gj} in the column \mathbf{j}^{th} of the motif matrix M , we have:

$$M_{gj}^* = \sum_{l=1}^{j-1} M_{gl} X_{lj} + \overset{0}{\cancel{M_{gj} X_{jj}}} + \sum_{l=j+1}^k M_{gl} X_{lj} = \sum_{l=1}^{j-1} M_{gl} X_{lj} + \sum_{l=j+1}^k M_{gl} X_{lj} = 0 \quad [\text{S2}]$$

For all genes with zero-entry in the column \mathbf{j}^{th} , Eq. S2 can be written in the matrix form:

$$\begin{bmatrix} M_{g_1 j}^* \\ M_{g_2 j}^* \\ \dots \\ M_{g_p j}^* \end{bmatrix} = \begin{bmatrix} M_{g_1 1} & M_{g_1 2} & \dots & M_{g_1 j-1} & M_{g_1 j+1} & \dots & M_{g_1 k} \\ M_{g_2 1} & M_{g_2 2} & \dots & M_{g_2 j-1} & M_{g_2 j+1} & \dots & M_{g_2 k} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ M_{g_p 1} & M_{g_p 2} & \dots & M_{g_p j-1} & M_{g_p j+1} & \dots & M_{g_p k} \end{bmatrix} \begin{bmatrix} X_{1j} \\ X_{2j} \\ \dots \\ X_{j-1,j} \\ X_{j+1,j} \\ \dots \\ X_{kj} \end{bmatrix} = \vec{0} \quad [\text{S3}]$$

if

$$M = \begin{bmatrix} M_{g_1 1} & M_{g_1 2} & \dots & M_{g_1 j-1} & M_{g_1 j+1} & \dots & M_{g_1 k} \\ M_{g_2 1} & M_{g_2 2} & \dots & M_{g_2 j-1} & M_{g_2 j+1} & \dots & M_{g_2 k} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ M_{g_p 1} & M_{g_p 2} & \dots & M_{g_p j-1} & M_{g_p j+1} & \dots & M_{g_p k} \end{bmatrix} \text{ and } \vec{X}_{\bullet j} = \begin{bmatrix} X_{1j} \\ X_{2j} \\ \dots \\ X_{j-1,j} \\ X_{j+1,j} \\ \dots \\ X_{kj} \end{bmatrix} \quad [\text{S4}]$$

then the Eq. S3 becomes:

$$M \bullet \vec{X}_{\bullet j} = \vec{0} \quad [\text{S5}]$$

Since the number of zero elements \mathbf{p} (number of genes that do not contain motif indexed as \mathbf{j}) in each column \mathbf{j} of matrix M is much larger than $k-1$ (number of motifs in a motif set) for eukaryotic genome, and the way we estimate the initial non-zero values for matrix M in Eq. 2, the matrix M in Eq. S5 is full-rank. One possible exception is when a motif only occurs in the presence of a second motif (as may be the case with two variants of the same motif). However, in this case it is easy to see that the use of the

lambda parameter still guarantees a single optimum, similar to the proof for rank-deficient A below. Accordingly, the equality in Eq. S5 is true if and only if:

$$\vec{X}_{\cdot j} = \vec{0} \quad [S6]$$

Therefore, all entries in column j of matrix X must be zero except X_{jj} , which can be of any value, for all $j=1..k$. As a result, matrix X is a diagonal matrix. Since the step (c) of the MED formalism requires that each row of matrix A has the unit norm, each diagonal element of matrix X must be 1. Therefore X is an identity matrix. This completes the proof of the first part.

2. In MED, the motif matrix M is unique even when matrix A is rank deficient:

When two or more distinct motifs are regulated by the same regulator, two or more corresponding rows of the global proxy regulator activity matrix A become linearly dependent (i.e. A is rank deficient), an event that cannot be controlled *a priori* mathematically. The same is also true when two or more regulators work together within a finite number of experimental conditions under consideration. For example, let $\vec{A}_{i_1 \cdot} = \alpha_2 \vec{A}_{i_2 \cdot} = \dots = \alpha_q \vec{A}_{i_q \cdot}$. In that case, the contribution of these motifs to the expression of gene g can be written as:

$$\begin{bmatrix} \vec{A}_{i_1 \cdot}^T & \vec{A}_{i_2 \cdot}^T & \dots & \vec{A}_{i_q \cdot}^T \end{bmatrix} \bullet \begin{bmatrix} M_{g i_1} \\ M_{g i_2} \\ \vdots \\ M_{g i_q} \end{bmatrix} = \begin{bmatrix} \vec{A}_{i_1 \cdot}^T & \alpha_2^{-1} \vec{A}_{i_2 \cdot}^T & \dots & \alpha_q^{-1} \vec{A}_{i_q \cdot}^T \end{bmatrix} \bullet \begin{bmatrix} M_{g i_1} \\ M_{g i_2} \\ \vdots \\ M_{g i_q} \end{bmatrix} \approx \begin{bmatrix} E_g^T \end{bmatrix} \quad [S7]$$

Obviously a set of infinite solutions can result from Eq. S7, because the rank of the matrix in Eq. S7 is one and hence values can be shifted freely among the q motif strength variables $M_{g i_l}$, $l=1..q$.

However, in the MED formalism, we require that the sum of square of motif strengths in each gene be minimal or as close as possible to some user predefined values (if they are known *a priori*) as the second term of Eq. 4 or Eq. 5 represents, respectively. Mathematically, this can be written in the matrix form as follows:

$$\begin{bmatrix} \vec{A}_{i_1 \cdot}^T & \alpha_2^{-1} \bullet \vec{A}_{i_2 \cdot}^T & \dots & \alpha_q^{-1} \bullet \vec{A}_{i_q \cdot}^T \\ & \sqrt{\lambda} \bullet I & & \end{bmatrix} \bullet \begin{bmatrix} M_{g i_1} \\ M_{g i_2} \\ \dots \\ M_{g i_q} \end{bmatrix} \approx \begin{bmatrix} E_g^T \\ \vec{0} \end{bmatrix} \quad \text{for Eq. 4} \quad [S8]$$

or

$$\begin{bmatrix} \bar{A}_{i_1}^T & \alpha_2^{-1} \bar{A}_{i_2}^T & \cdots & \alpha_q^{-1} \bar{A}_{i_q}^T \\ & \sqrt{\lambda} \bullet I & & \end{bmatrix} \bullet \begin{bmatrix} M_{g^{i_1}} \\ M_{g^{i_2}} \\ \vdots \\ M_{g^{i_q}} \end{bmatrix} \approx \begin{bmatrix} E_g^T \\ \sqrt{\lambda} \bullet (\bar{M}_g^*)^T \end{bmatrix} \quad \text{for Eq. 5} \quad [\text{S9}]$$

In the Eq. S9, \bar{M}_g^* is a vector of user defined priors for the motif strength for gene \mathbf{g} . Since the matrix in Eqs. S8-S9 is always full-rank, due to the block represented by the identity matrix I regardless of data in matrix \mathbf{A} , Eqs. S8 or S9 has a unique least square solution. This completes the proof of the second part.

Note that in Eq. S1, the linear transformations of matrix M to a matrix represented by $M \bullet X$ and of matrix A to a matrix represented by $X^{-1} \bullet A$ are known as (oblique) rotation (Paatero et al., 2002). Historically, in the proof to part 1, the idea of using an *a priori* known pattern of zero elements in matrix factorization problem to reduce the rotational degrees of freedom (i.e. number of off-diagonal non-zero elements of matrix X) is a well-known one. If there are enough pre-determined known zero elements in one of factored matrices, the rotation becomes prohibited as discussed by Anderson (Anderson, 1984); Paatero and co-workers (Paatero et al., 2002); and an alternative proof to part 1 on the structure of matrix X^{-1} instead of matrix X was previously provided by Liao and co-workers (Liao et al., 2003). In such case, the matrix X and hence X^{-1} become diagonal matrices, otherwise such pre-determined pattern of zero is violated. However, the use of the pre-determined pattern of zeros alone is not enough, since the rotation cannot deal with, and therefore guarantee, the uniqueness of motif matrix M when rows of matrix A are linearly dependent as discussed in the proof to part 2 above. The main mathematical step that makes our MED formalism work is the proof in step 2, which guarantees the uniqueness of the motif matrix M regardless of the data pattern in matrix A .

Supplementary 5: Procedure for initializing non-zero matrix elements of the motif matrix M

If the regulatory motifs were defined by a consensus sequence, the motif matrix M could simply be initialized by counting how many times each specific consensus sequence occurs in each promoter (Bussemaker et al., 2001). However, since we are dealing in this case with motifs which are specified by a position specific scoring matrix (PSSM), we can weight this count by the score of each motif instance.

The initial weight M_{gm} for motif m and gene g is given by:

$$M_{gm} = \sum_i e^{S_i(g,m)} / \max_{i,g} (e^{S_i(g,m)}) \quad [\text{S10}]$$

where $S_i(g,m)$ is the ScanAce score (Hughes et al., 2000) of the i^{th} instance of motif m in the promoter of gene g , which is approximately proportional to the log

probability of binding to that sequence (Stormo, 2000). Note that for the best-fit sequences to the PSSM for the motif, the initial value for M_{gm} reverts to a count of the number of motif instances.

Supplementary 6: Standard least square formalism used in MED

Given a matrix $M \in \mathfrak{R}^{m \times n}$, $m \geq n$ and a vector of observable data $\vec{b} \in \mathfrak{R}^m$, the least square problem is to find a vector $\vec{x} \in \mathfrak{R}^n$ such that $\|M\vec{x} - \vec{b}\|$ is minimized. Geometrically, this can be represented in Fig SF4 where the vector solution \vec{x} can be calculated by projecting the vector of observable data \vec{b} orthogonally onto the range of the matrix M spanned by its columns vectors. Mathematically we can write:

$$M\vec{x} \approx \vec{b} \quad [\text{S11}]$$

$$\Rightarrow \vec{x} = (M^T M)^{-1} M^T \vec{b} \quad [\text{S12}]$$

Using singular value decomposition for the matrix M , we have:

$$M = U \bullet \Sigma \bullet V^T \quad [\text{S13}]$$

Substitute S13 into S12, the solution to the least square problem can be computed using the following formula:

$$\vec{x} = V \bullet \Sigma^{-1} \bullet U^T \bullet \vec{b} \quad [\text{S14}]$$

If the vector of observable data \vec{b} in Eq. S14 is a matrix, then Eq. S14 becomes:

$$X = V \bullet \Sigma^{-1} \bullet U^T \bullet B \quad [\text{S15}]$$

where X is now a matrix instead of a vector.

Supplementary 7: Definition of the total residual variance, convergence and selection of λ

From Eq. S11, the residual matrix R is defined as follows:

$$R = M \bullet X - B \quad [\text{S16}]$$

therefore total residual variance R is then defined as:

$$R = \sum_{i=1}^m \text{var}(\vec{R}_i) \quad [\text{S17}]$$

Since the second term of the Eq. 4 imposes additional constraint on the factorization for any $\lambda > 0$, test error will always be larger than the corresponding with $\lambda = 0$. Therefore, the value of λ must be chosen in such a way that it does not affect the test error significantly and Figs. SF5a and SF5b show both training and test errors, respectively, for $\lambda = 0$ and $\lambda = 10^{-4}$, which is the value we used in this work. Given this λ , we also show the trajectory of MED algorithm for both training and test error in factoring the expression matrix in Figs. SF4c and SF4d, respectively.

Supplementary 8: Extended form of Eq. 1 to account for non-linear motif-motif interactions

In the Eq. 1 in the main text, the log expression ratio of gene g under environmental condition c , E_{gc} , can be approximated as follows:

$$E_{gc} \approx \sum_{j \in \Omega_g} M_{gj} A_{jc} \quad [1]$$

Although motif-motif interactions are not explicitly included in Eq. 1, such interactions can still be captured to some degree by assigning a different weight to a motif depending on the presence of other motifs it interacts with. Nevertheless, one can also choose to include interaction terms explicitly in Eq. 1, according to the following equation:

$$E_{gc} \approx \sum_{j \in \Omega_g} M_{gj} A_{jc} + \sum_{jl \in \Omega_g \times \Omega_g} M_{gjl} A_{jlc} \quad [S18]$$

where M_{gjl} represents the influence of interaction between motifs j and l present in the promoter of gene g ; and A_{jlc} represents the global proxy activity of the virtual regulator for such interaction. With this explicit form, one can investigate combinatorial regulation principles in terms of motif-motif interactions, which is an important feature in eukaryotic genomes (Kellis et al., 2003). The differences between Eqs. 1 and S18 are extra non-zero entries for each gene (extra columns) in the motif matrix M and the MED formalism can still be applied transparently.

Supplementary 9: Data

Motif data and expression data sets along with other essential data used in this work are available from our website at <http://arep.med.harvard.edu/MED>

Supplementary 10: p-Value calculations for motif strengths derived from the random shuffling of expression data

In the random shuffling test, we recalculated the strength of each motif (random motif strength) in the same manner as being done in Fig. 3 (in the main text) for each randomly shuffled expression matrix, for 100 times. The supplementary Figures SF6a-c show the distribution of such random motif strengths for the PAC, RRPE, MCB, and RAP1 motifs in the promoter location and with orientation, where they are most active,

along with the corresponding ones (actual motif strength) derived from the actual expression data. Since the actual motif strength is far away from the mean of the distribution of the corresponding random motif strengths by at least 28 standard deviations, the p-value for these actual motif strengths can be much smaller than 0.01, due to the Chebyshev inequality (Abramowitz and Stegun, 1972). Note that for the MCB and RAP1 motifs, we also show their distribution and actual motif strength for both motif's orientations with " \leftarrow " denotes $5' \leftarrow$ and " \rightarrow " denotes $5' \rightarrow$ orientations.

Supplementary 11: Mathematical notations

Throughout the text, we employ a set of mathematical notations and symbols. Here are their definitions.

\vec{X}_j : denotes the j^{th} row vector of a given matrix X

\vec{X}_j : denotes the j^{th} column vector of a given matrix X

X^T : denotes the transpose of a give matrix X

\bullet : represents the multiplication operation

\vec{x} : denotes a column vector.

Reference

- Abramowitz, M. and Stegun, I.A. (1972) *Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables*. Dover, New York.
- Anderson, T.W. (1984) *An Introduction to Multivariate Statistical Analysis*. Wiley, Chichester ; New York.
- Bussemaker, H.J., Li, H. and Siggia, E.D. (2001) Regulatory element detection using correlation with expression. *Nat Genet*, **27**, 167-171.
- Hughes, J.D., Estep, P.W., Tavazoie, S. and Church, G.M. (2000) Computational identification of Cis-regulatory elements associated with groups of functionally related genes in *Saccharomyces cerevisiae*. *Journal of Molecular Biology*, **296**, 1205-1214.
- Kellis, M., Patterson, N., Endrizzi, M., Birren, B. and Lander, E.S. (2003) Sequencing and comparison of yeast species to identify genes and regulatory elements. *Nature*, **423**, 241-254.
- Liao, J.C., Boscolo, R., Yang, Y.L., Tran, L.M., Sabatti, C. and Roychowdhury, V.P. (2003) Network component analysis: reconstruction of regulatory signals in biological systems. *Proc Natl Acad Sci U S A*, **100**, 15522-15527.
- Paatero, P., Hopke, P.K., Song, X. and Ramadan, Z. (2002) Understanding and controlling rotations in factor analytic models. *Chemometrics and Intelligent Laboratory Systems*, **60**, 253-264.
- Stormo, G.D. (2000) DNA binding sites: representation and discovery. *Bioinformatics*, **16**, 16-23.

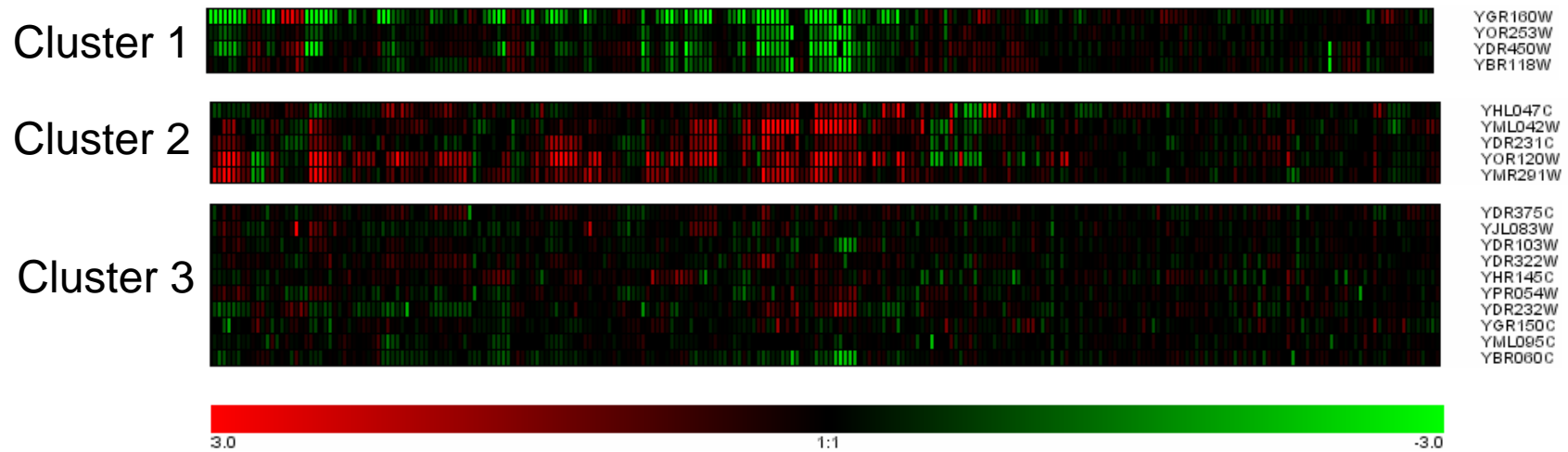
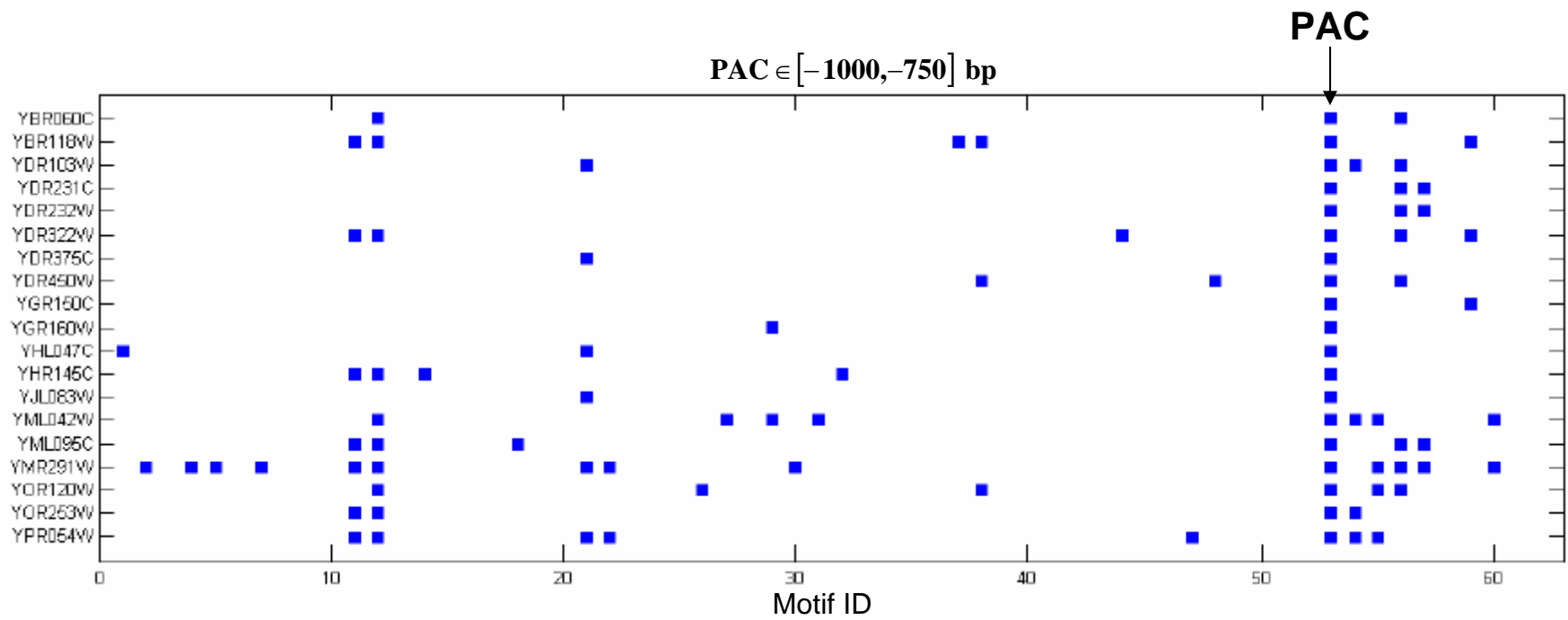
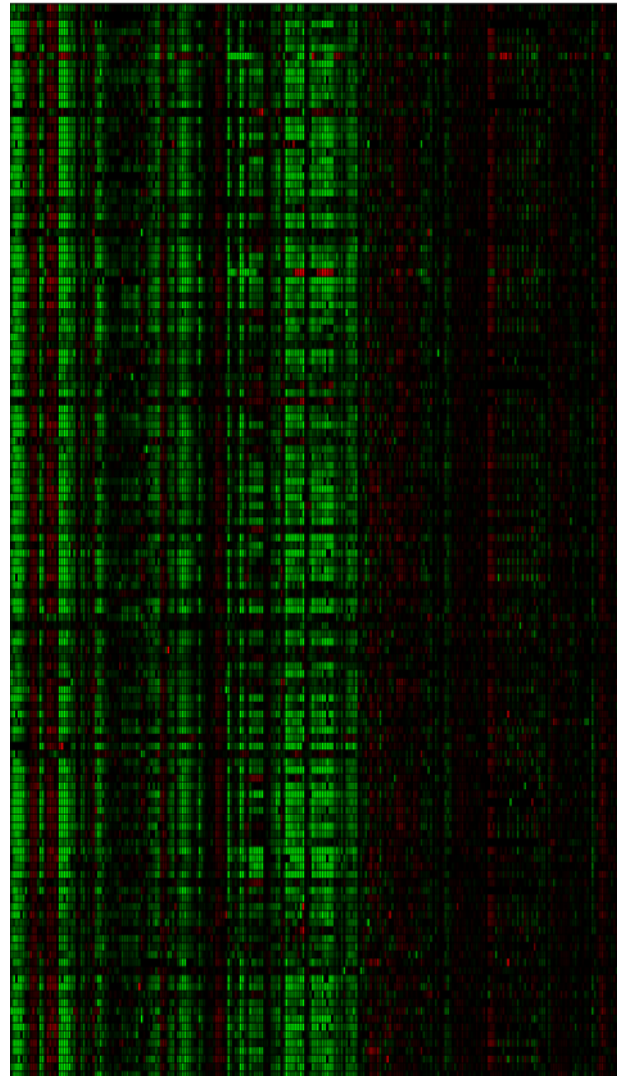


Figure SF1b

PAC $\in [-150, \text{ATG}]$ bp



VAL009W
VBL014C
VER034C
VER207W
VER271W
VCL056C
VCL056C
VCR015W
VCR051C
VCR072C
VCR072C
VLD011W
VLD011W
VLD030C
VDR000W
VDR007C
VDR011C
VDR250W
VDR251C
VDR305C
VDR305W
VDR440C
VDR440C
VDR460C
VDR460C
VDR477W
VER002W
VER002W
VER050W
VER120C
VER127C
VFL029W
VFL029W
VFL079C
VFL171W
VFL355W
VFR103W
VFR140W
VFR175W
VFR242C
VFR242C
VFR272C
VFR272C
VFR282C
VFR282C
VFR002W
VFR002W
VFR002C
VFR070W
VFR091W
VFR091W
VFR140W
VFR140W
VFR165W
VFR170W
VLD109W
VLD109C
VLD127C
VLD239W
VJL020W
VJL106C
VJL106C
VJL122W
VJL122W
VJL281C
VJL281C
VJL281C
VJL281C
VLD009W
VLD078W
VLD078C
VLD082C
VLD082C
VLD099C
VLD117W
VLD117W
VLD101W
VDR003C
VDR003C
VLD022C
VLD022C
VLD022W
VLD022W
VLD034C
VLD034C
VLD014C
VLD015W
VLD017W
VLD000W
VLD008W
VLD074C
VLD222C
VLD336C
VLD400C
VLD000W
VLD022W
VLD125C
VDR040C
VDR040W
VDR128W
VDR131C
VDR185W
VDR230C
VDR200C
VDR200C
VDR091W
VDR075W
VDR110C
VDR113W
VDR124W
VDR124W
VDR205W
VDR247W
VDR248C
VDR248C
VDR305C
VDR003C
VDR039W
VOL041C
VOL080C
VOL140W
VOL144W
VOR075W
VOR115C
VOR207C
VOR207C
VOR211C
VOR272W
VFL047W
VFL068C
VFL068C
VFL089C
VFL090W
VFL120W
VFL211W
VFL217C
VFR200W
VFR010C
VFR112C
VFR144C
VFR190C



Figure SF1c

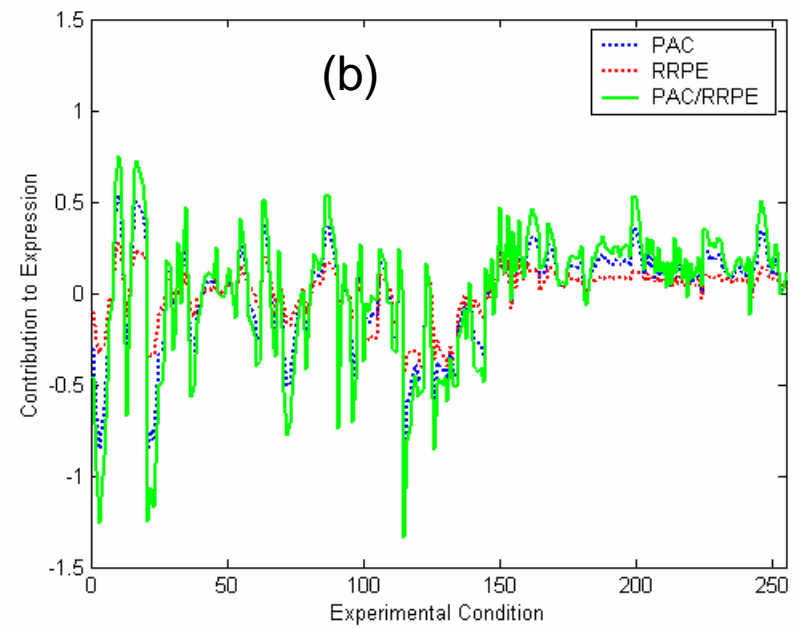
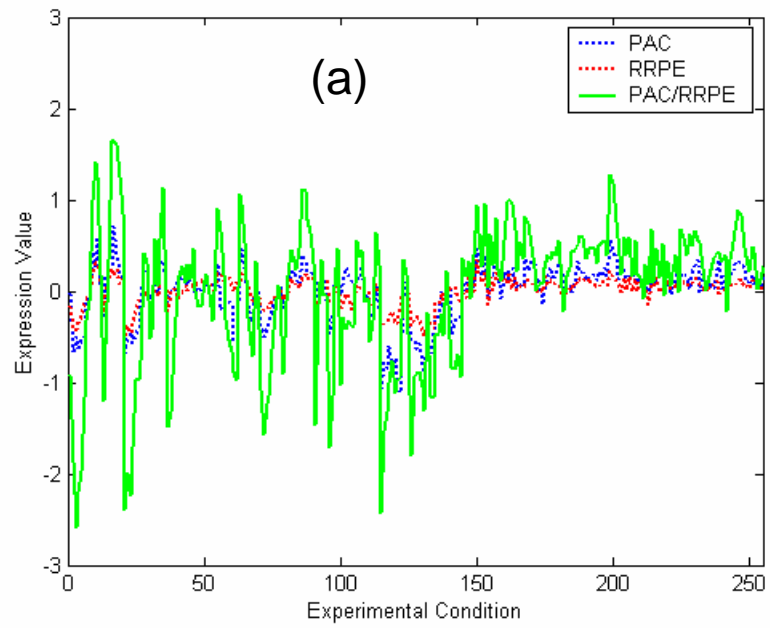


Figure SF2

(c)

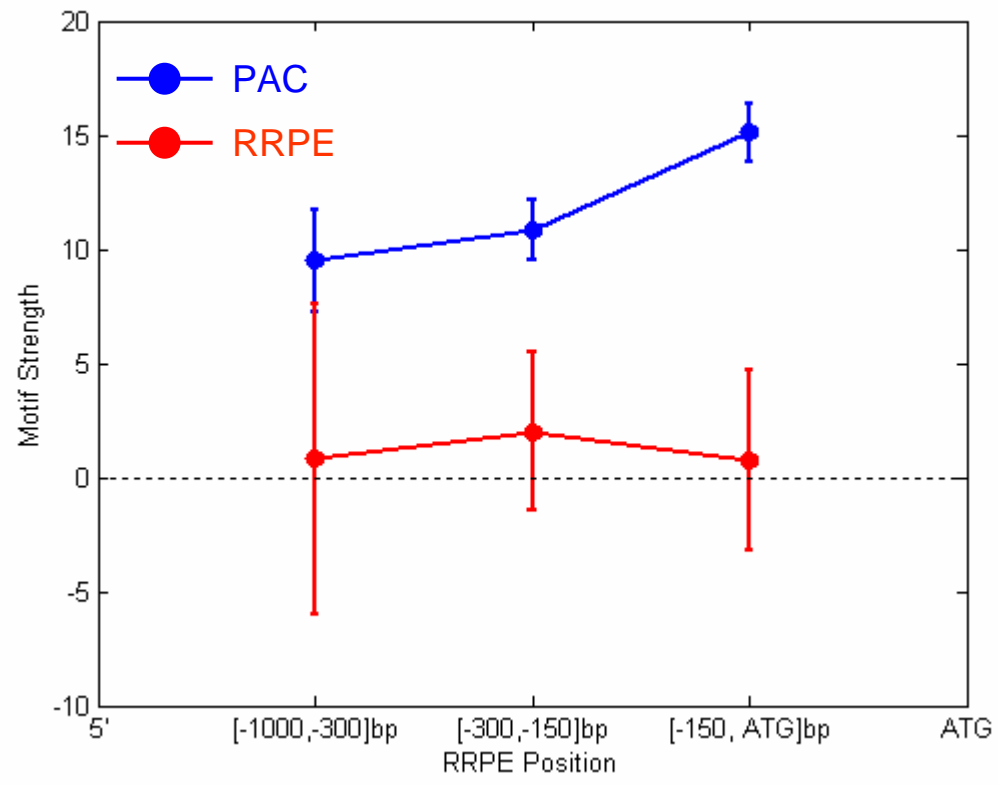


Figure SF2C

A

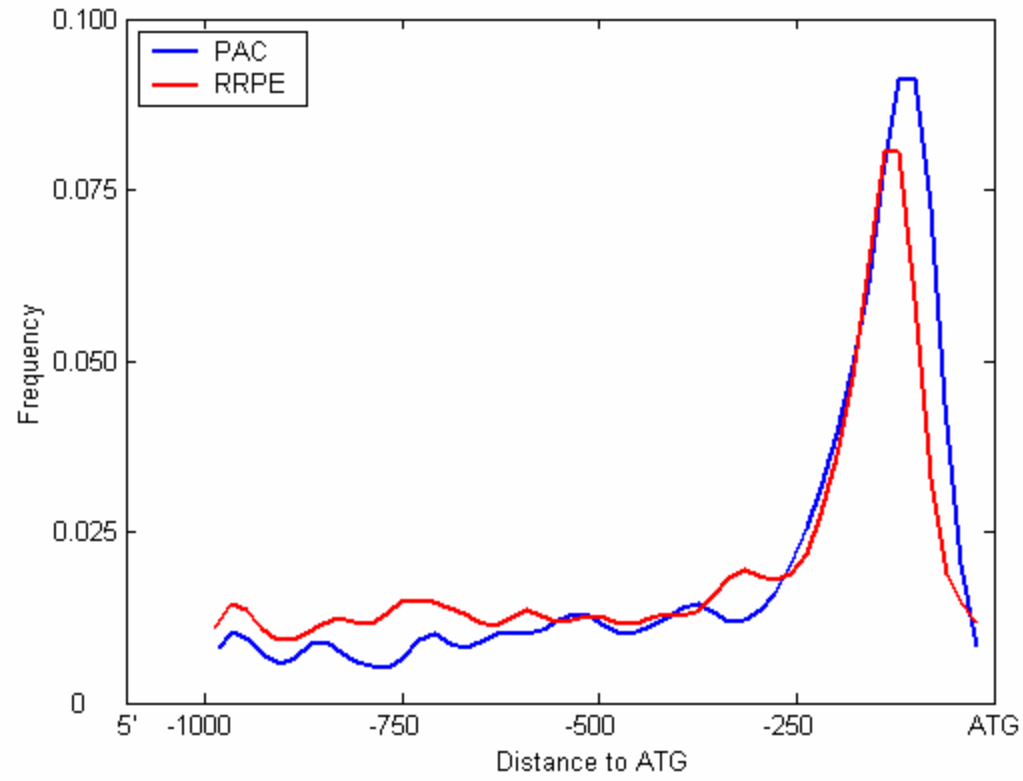


Figure SF3A

B

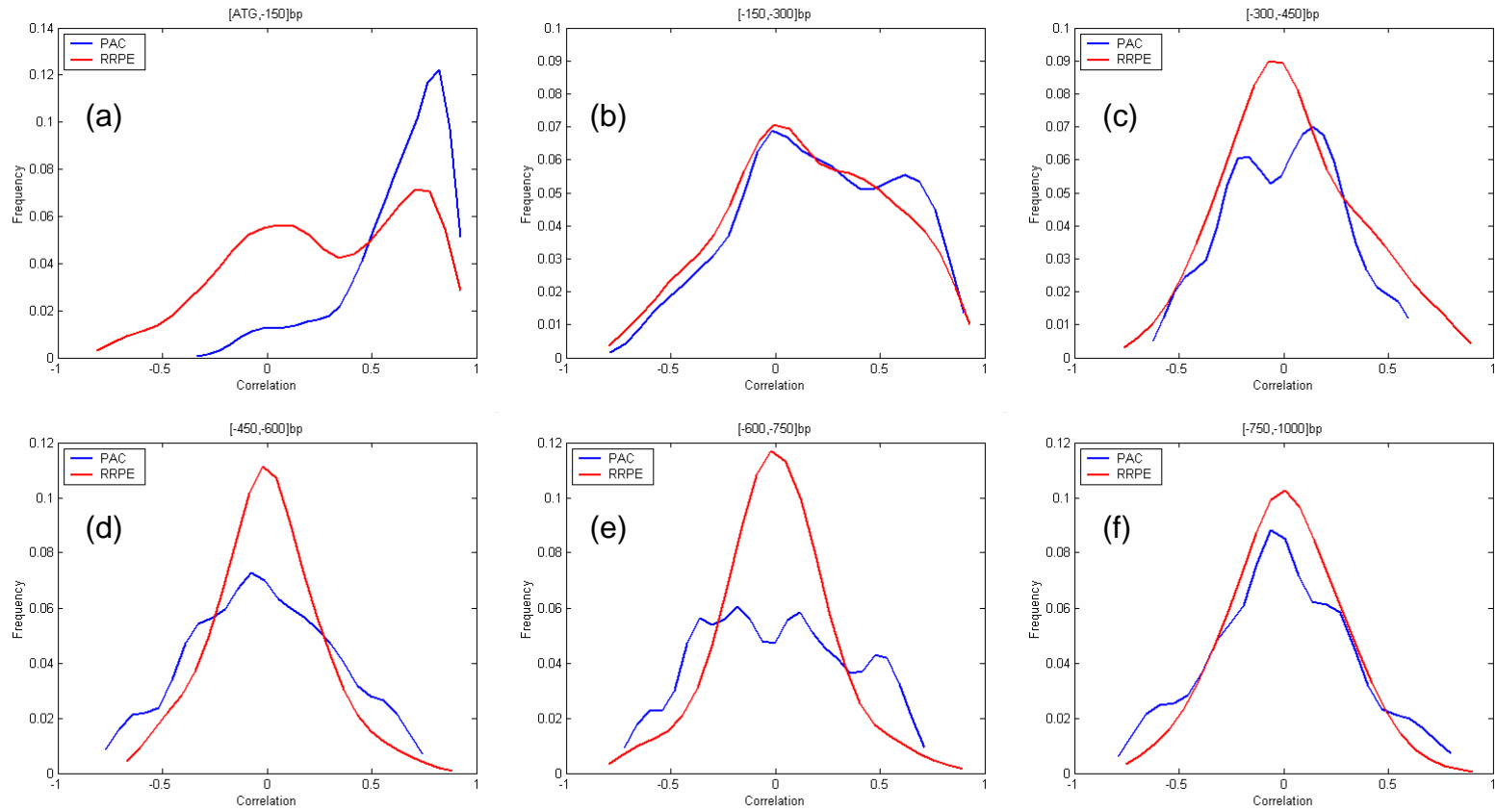


Figure SF3B

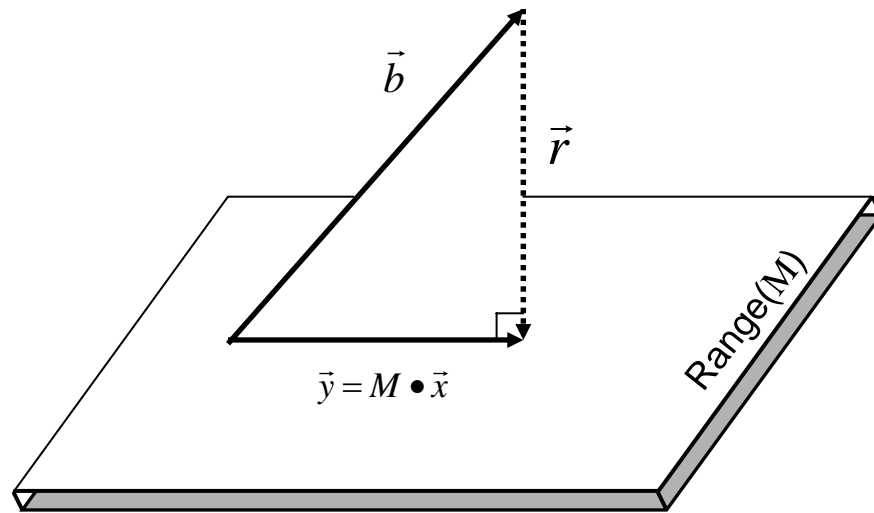


Figure SF4

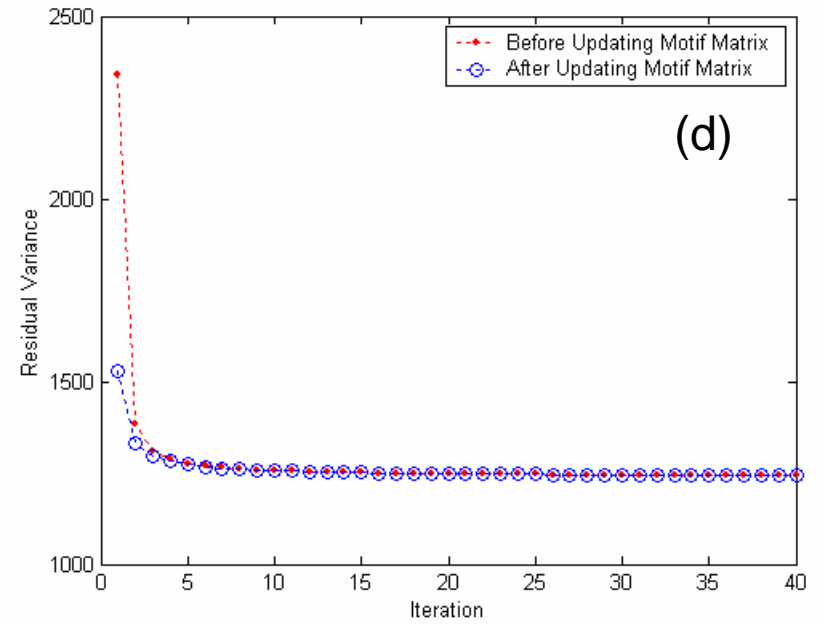
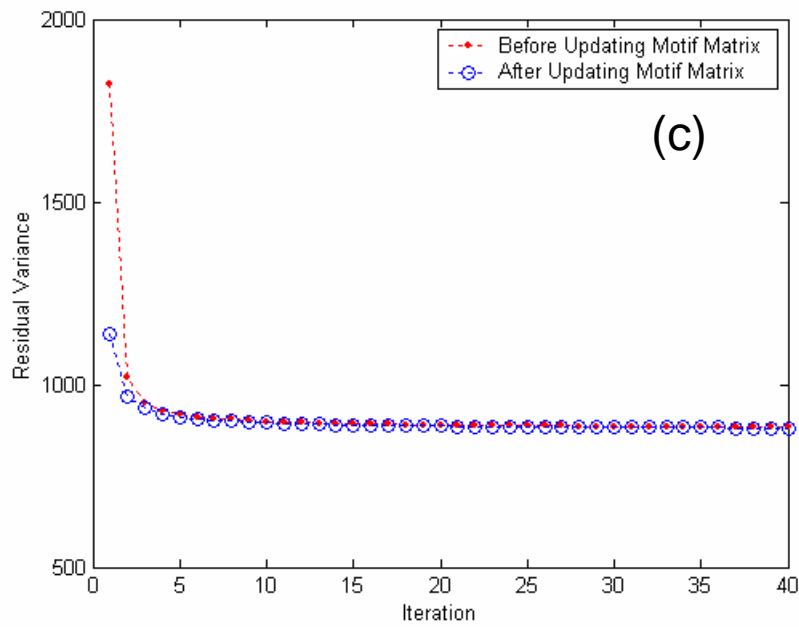
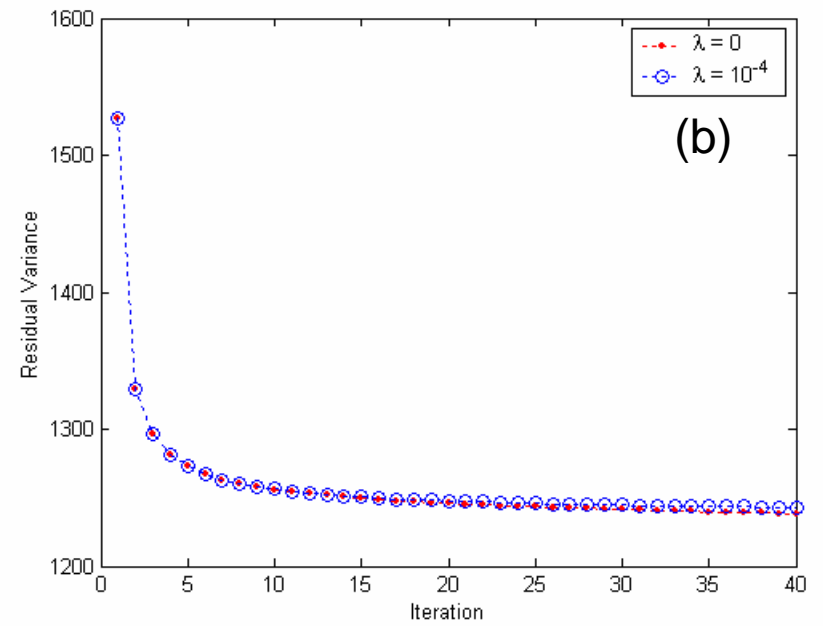
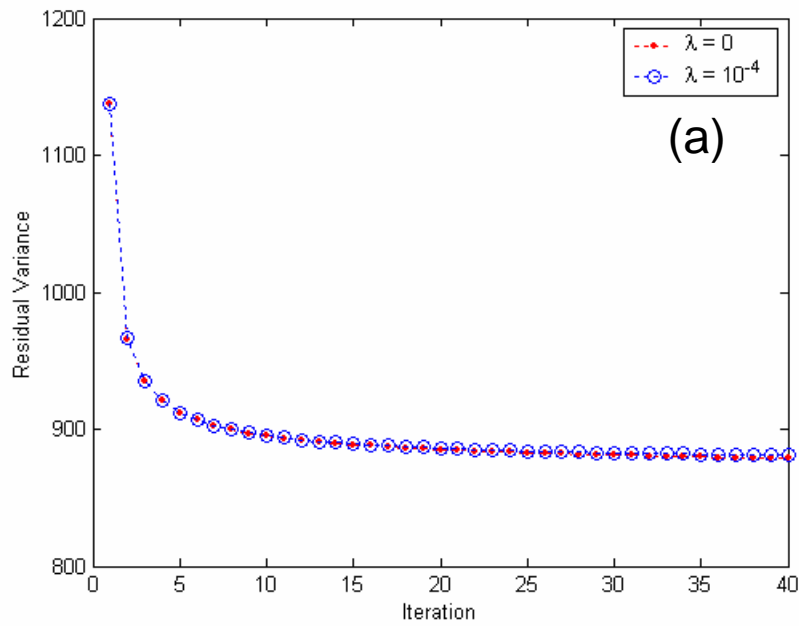


Figure SF5

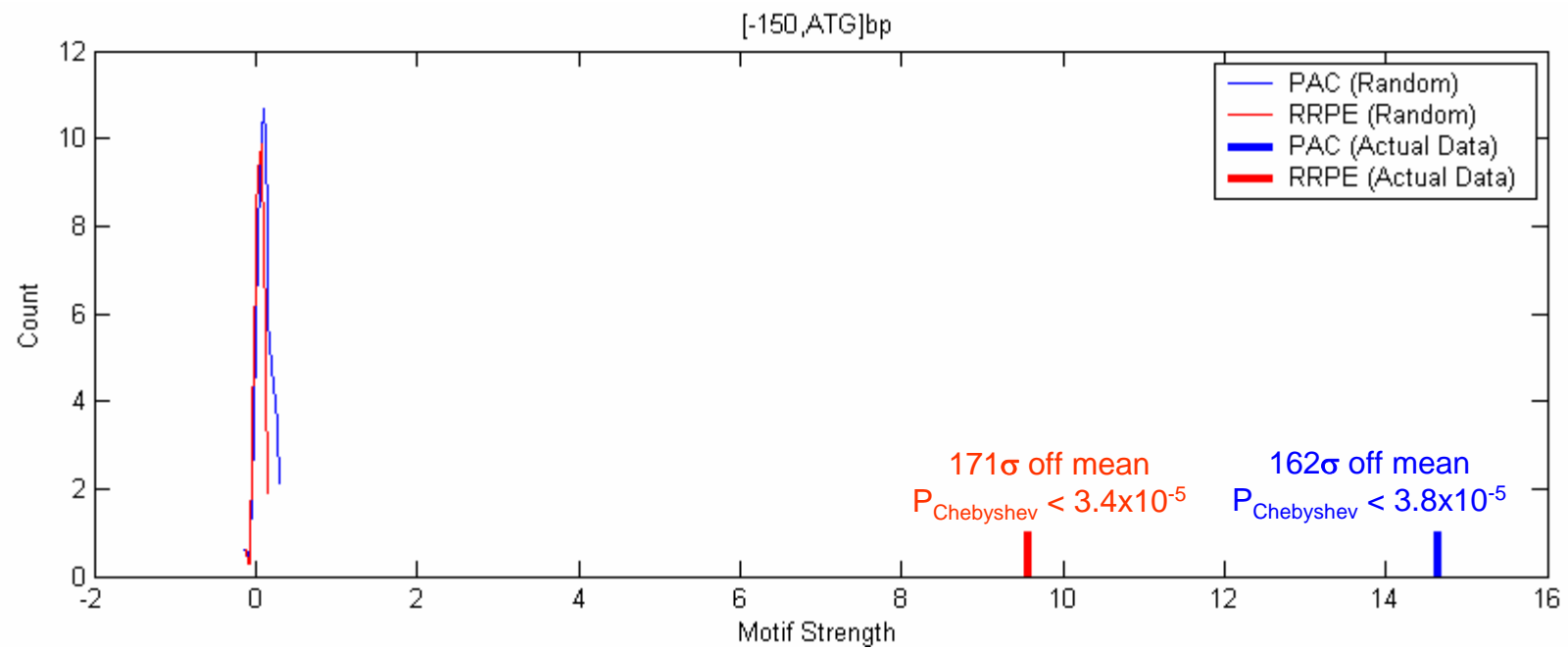


Figure SF6a

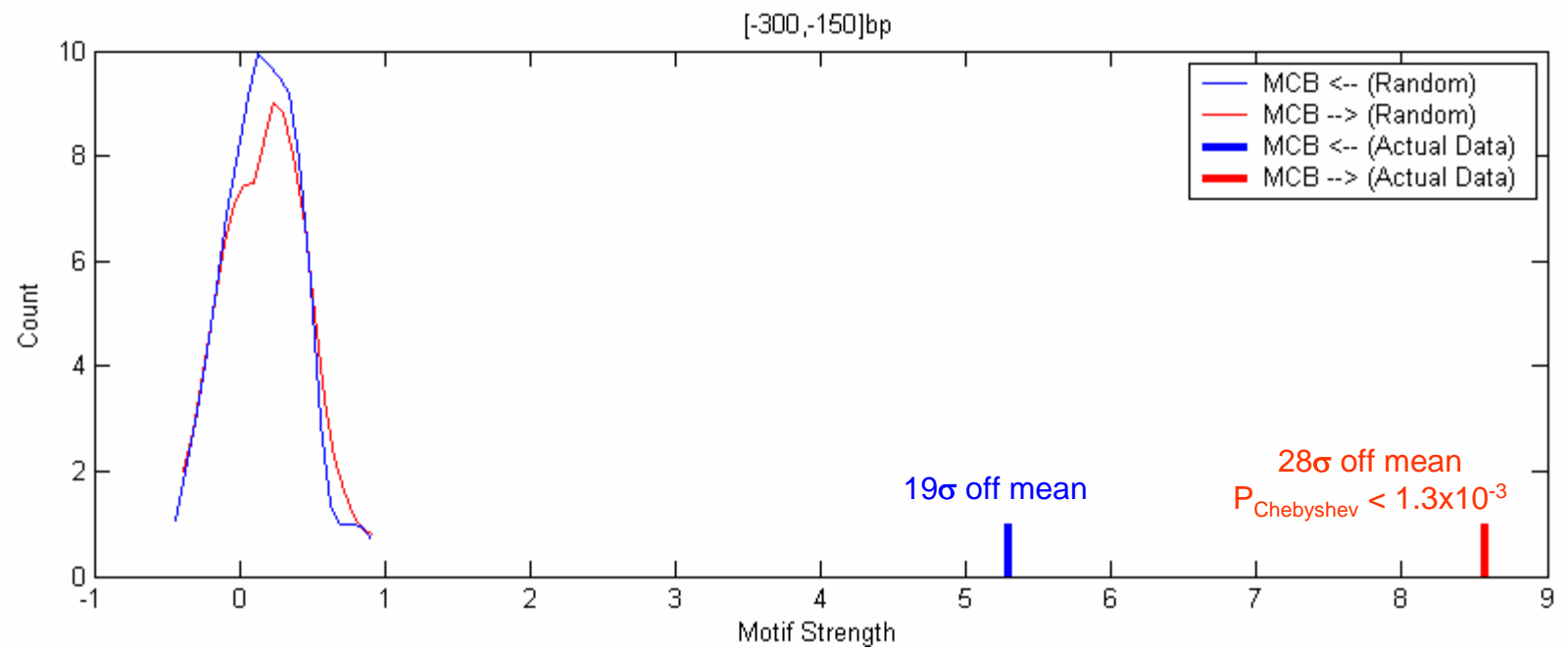


Figure SF6b

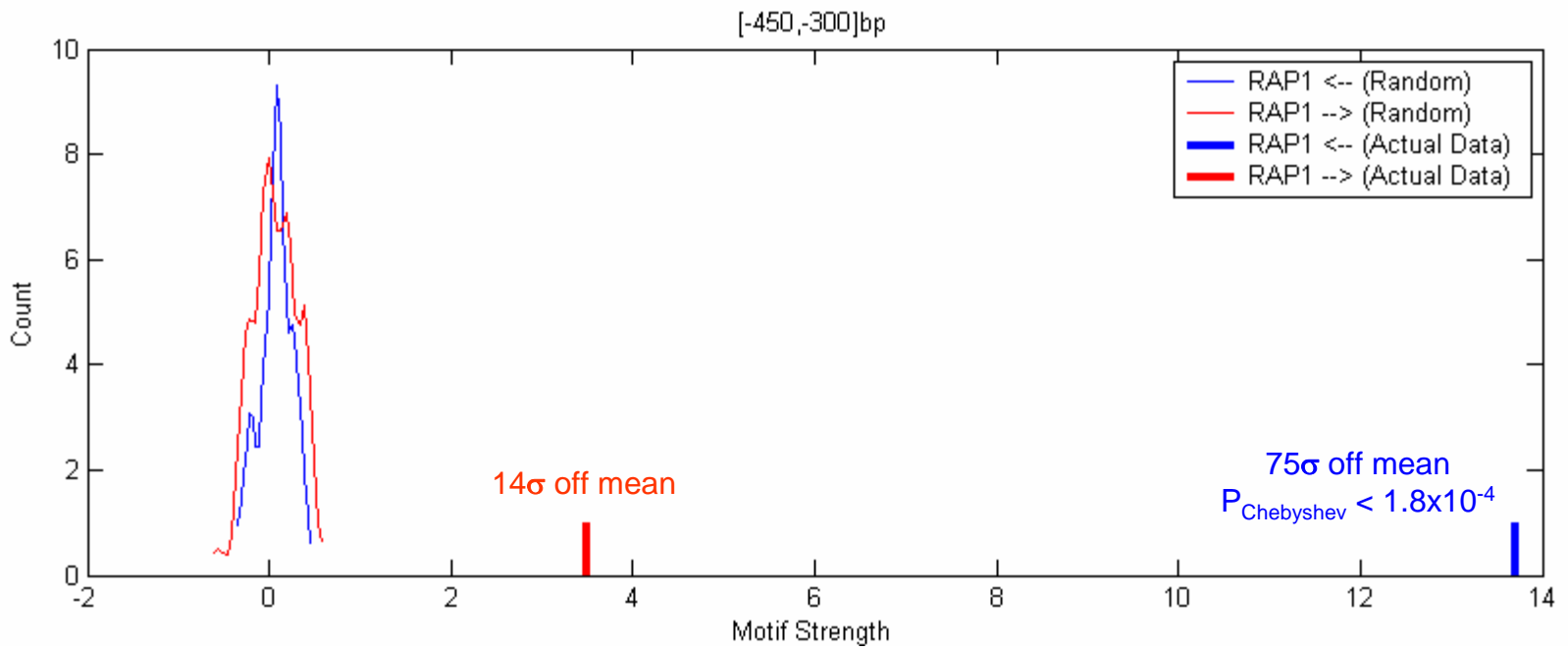


Figure SF6c