

Supplementary Information: Network analysis identifies
chromosome intermingling regions as regulatory hotspots for
transcription

Anastasiya Belyaeva Saradha Venkatachalapathy Mallika Nagarajan
G.V. Shivashankar Caroline Uhler

November 30, 2017

SI Appendix

Obtaining Hi-C matrices

The Hi-C matrices were obtained from [12] at 250kb resolution. Matrices were corrected for bias using interchromosomal matrix balancing based on the Knight-Ruiz algorithm using software in [3]. Centromeric regions, as well as peri-centromeric regions within 2Mb of the centromere were filtered out. Repeat regions, outliers based on row and column sums (outside of $1.5 \times$ interquartile range interval) in the Hi-C contact matrix, and regions already masked in [12], were removed from the analysis. The final Hi-C matrix was $\log(1+x)$ transformed and normalized by mean contact frequency and standard deviation, computed over all interchromosomal contact pairs that were not filtered out.

LAS Algorithm for identifying highly interacting regions

The LAS algorithm [13] takes a real-valued data matrix $X(m \times n)$ as input and outputs contiguous submatrices $U(k \times l)$ that have a high average, τ . This is done via the following iterative algorithm:

Repeat until $\tau\sqrt{kl} < threshold$:

1. Search: greedily, by updating one row and column set at a time, find a submatrix U that maximizes the submatrix score

$$S(U) = -\log[(m-k+1)(n-l+1)\Phi(-\tau\sqrt{kl})] \quad (1)$$

2. Remove: identify rows and columns corresponding to U in X and subtract the submatrix average τ from this set of rows and columns.

The LAS algorithm search space was limited to contiguous submatrices of at most $10\text{Mb} \times 10\text{Mb}$ in size, i.e. 40×40 submatrices (at 250kb resolution). For each chromosome pair each iteration of the search procedure was initialized at a random contiguous $k \times l$ submatrix in the interchromosomal Hi-C map. The threshold for the algorithm was chosen based on a Gaussian approximation such that $P(\tau\sqrt{kl} > threshold) = 1E-15$. This stringent cutoff guarantees that highly interacting regions in the whole-genome Hi-C map are identified with FDR controlled at 4.16×10^{-8} (see the following paragraph). Returning submatrices U , as determined by LAS, for each interchromosomal contact matrix results in a list of highly interacting pairs of 250kb regions.

FWER and FDR computation for Large Average Submatrix (LAS) algorithm

The LAS algorithm takes a real-valued matrix $X(m \times n)$ as input and outputs contiguous submatrices $U(k \times l)$ of high average [13]. The null hypothesis is that the interchromosomal Hi-C matrix is a standard Gaussian random matrix, and the alternative hypothesis is that the interchromosomal Hi-C matrix is a sum of K constant (> 0) submatrices plus a standard Gaussian random matrix, i.e., that the Hi-C contact matrix contains substructure [9]. More precisely, each entry in the alternative model can be expressed as

$$x_{i,j} = \sum_{k=1}^K \alpha_k I(i \in A_k, j \in B_k) + \epsilon_{ij}, \quad (2)$$

where $A_k \subseteq [m]$ and $B_k \subseteq [n]$ are the row and column sets of the k th submatrix, α_k is the constant corresponding to the k th submatrix, ϵ_{ij} are independent noise variables sampled from $\mathcal{N}(0,1)$, and $I(\cdot)$ is the indicator function. Note that $K=0$ corresponds to the null model. The individual entries in the $\log(1+x)$ transformed Hi-C matrices have an R^2 of 0.972 with the standard normal distribution, which justifies using a standard Gaussian random matrix as null hypothesis.

Let $\tau := \text{Avg}(U)$, i.e., the average of the submatrix U . Under the null hypothesis, $\tau\sqrt{kl} \sim \mathcal{N}(0,1)$ and thus the probability of observing a $k \times l$ submatrix V with an average of τ or greater is $P(\text{Avg}(V) \geq \tau) = \Phi(-\tau\sqrt{kl})$, where Φ is the standard normal cdf. Let A denote the event that there exists a $k \times l$ submatrix V with average greater than or equal to τ in an $m \times n$ matrix. Note that this event is bounded as follows: $P(A) \leq \sum P(\text{Avg}(V) \geq \tau)$, where the sum is over all $k \times l$ submatrices in the $m \times n$ matrix. Hence, under the null hypothesis, $P(A) \leq N\Phi(-\tau\sqrt{kl})$, where $N = (m-k+1) \times (n-l+1)$, the total number of contiguous submatrices of size $k \times l$ in an $m \times n$ matrix.

The search space of the LAS algorithm was limited to contiguous submatrices of at most $10\text{Mb} \times 10\text{Mb}$ in size, which corresponds to 40×40 submatrices (at 250kb resolution). In order to calculate

the total number of hypotheses for each interchromosomal matrix, we summed the number of possible contiguous submatrices for all combinations of k and l within the $[1,40]$ range. Considering all pairs of interchromosomal matrices, the total number of hypotheses was 9.33×10^{10} . In our procedure, we applied a p-value threshold, namely $P(\tau\sqrt{kl} > threshold) = 1 \times 10^{-15}$, for the discovery of significant submatrices. Using a formulation based on the Bonferroni correction, we can estimate the familywise error rate (FWER), which is the probability of making at least one type I error. Let p be the p-value threshold, b the number of hypotheses and α the FWER level. The Bonferroni correction rejects the null hypothesis when p-value $\leq \frac{\alpha}{b}$, thereby controlling the FWER at $\leq \alpha$. With our p-value threshold of 1×10^{-15} , the FWER is ≤ 0.0000933 .

We can also calculate the false discovery rate (FDR), i.e. the fraction of false discoveries among all discoveries, using the Benjamini-Hochberg procedure. Let a be the number of discoveries, b the number of hypotheses, p the p-value threshold, and α the FDR level. The Benjamini-Hochberg procedure rejects the null hypothesis when p-value $\leq \frac{a}{b}\alpha$, thereby controlling the FDR at $\leq \alpha$. With our p-value threshold of 1×10^{-15} and the resulting number of discoveries $a = 2244$, the FDR is $\leq 4.16 \times 10^{-8}$.

Genomic features

Pre-processed data for 48 features including histone modifications, transcription factor ChIP-seq, DNase-seq and RNA-seq, were retrieved from ENCODE (32), Roadmap Epigenomics (33), the GEO database, and previous studies (49) (SI Appendix, Table S2) for the IMR90 cell line. In order to obtain the genomic profile for a 250kb region, matching the resolution of Hi-C data, the number of peaks overlapping the 250kb region was calculated for each feature. For each feature, the feature matrix was $\log(1+x)$ transformed and z-scored by computing the mean and standard deviation of the feature across all regions in the genome that were not removed by the Hi-C filtering step.

Weighted Correlation Clustering

The weighted network of 250kb regions was partitioned into clusters using weighted correlation clustering on networks [4]. This method determines clusters by drawing cluster boundaries across edges with low weights but not across edges with high weights by solving a non-convex minimization problem. Weighted correlation clustering was run in 25 replicates and the clustering with the lowest value of the objective function was chosen for further analysis. The resulting clusterings were robust across replicate runs, as evidenced by the high adjusted mutual information between cluster labels across runs (Fig. S12).

Classification into Intermingling and Non-intermingling Domains

In order to identify features that may be important for chromosome intermingling, a binary classification task was performed. The training and test data consisted of genomic feature profiles (SI Appendix, Table S2) for intermingling versus non-intermingling regions, weighted by the number of samples in each class. Classification was done using eXtreme gradient boosting trees with $n_estimators = 1000$, $learning_rate = 0.1$, $max_depth = 5$, $min_child_weight = 1$ with 10-fold cross-validation. Feature importances were computed by the relative rank of a feature in the decision tree, calculated via `feature_importances_` function in scikit-learn [11] in python. Additionally, features were evaluated using iterative feature elimination by removing one feature at a time and optimizing the AUC.

Fold Enrichment of Genomic Features

Fold enrichment for the intermingling regions as well as for specific clusters was calculated as follows:

$$\frac{\# \text{ bases in cluster and having feature}}{\# \text{ bases in genome}} \bigg/ \left(\frac{\# \text{ bases in cluster}}{\# \text{ bases in genome}} \right) \left(\frac{\# \text{ bases having feature}}{\# \text{ bases in genome}} \right) \quad (3)$$

A fold enrichment of 1 indicates that the two events - belonging to the intermingling regions or a particular cluster and belonging to a particular feature - are independent events.

Comparison to a random network - Stochastic Block Model

To analyze the importance of the spatial interactions for the function and properties of the determined clusters, we performed a comparison based on a "similar" network in which the spatial interactions have been randomized. To be more precise, we generated a network from a stochastic block model, where each chromosome is a community and the edge probabilities within and between communities are computed from the number of interactions in the Hi-C matrix as determined by the LAS algorithm. In order to obtain similar cluster sizes as in the original network, we sampled the edge weights from the observed distribution of edge weights.

Using a 2-sided χ^2 -test we tested whether the proportion of intermingling regions in the observed network was equal (H_A : not equal) to the proportion of intermingling regions in a random network. Generating 50 networks from the stochastic block model and using the average proportion of intermingling regions as test statistic, the null hypothesis was rejected with a p-value $< 2.2 \times 10^{-16}$.

Further, in order to test the functional relevance of the determined clusters, we tested using a 1-sided χ^2 -test whether the proportion of clusters enriched for all five active marks (RNAPII, H3K9ac, H3K36me3, H3K4me3, H3K4me1) in the observed network was equal (H_A : larger than) in a random network. As in the previous test, generating 50 networks from the stochastic block model and using the average proportion as test statistic, the null hypothesis was rejected with a p-value of 1.398×10^{-5} .

Finally, we tested the regulatory event that active and inactive clusters are spatially separated. To do this, we tested using a 1-sided χ^2 -test whether the proportion of clusters enriched for all five active marks and the inactive mark (H3K9me3) in the observed network was equal (H_A : smaller than) in a random network. Following the same procedure as in the previous test, the null hypothesis was rejected with a p-value of 4.699×10^{-4} .

Venn diagrams

Venn diagrams were constructed using the software provided by [1].

Gene ontology

Expressed genes for IMR90 with reads per kilobase of transcript per million mapped reads (RPKM) > 0 were obtained from ENCODE [5]. For each cluster, we identified the genes that resided in the cluster and were expressed. For each cluster, we then performed gene ontology (GO) term analysis on these genes using DAVID [6, 7].

Cluster ranking

In order to select clusters for experimental validation, we ranked each cluster based on the number of TFBS present in each region in the cluster. Several methods and databases were used for ranking the clusters in order to choose a robust set of clusters for experimental validation. The whole genome was scanned for TFBS using position frequency matrices (PFMs) from the JASPAR2016 database for humans. MOODS software [8] was used to identify motif matches. TFBS were further filtered by ChIP-seq from ENCODE (32), resulting in 52 TFs or robust CAGE peaks, resulting in 386 TFs. The CAGE peaks, which indicate transcription start sites, were obtained from the FANTOM5 project [2], which pooled CAGE analysis over 573 human cell samples. These peaks were flanked by 400bp upstream and 50bp downstream as suggested by [10] and overlaid with the TFBS data to obtain the final set of TFBS.

In the following, we explain how we used the determined TFBS to rank the clusters based on a permutation test. First, we constructed a score function to compare observed and randomized matrices. For each cluster we construct a matrix Z of size $m \times n$ consisting of the TFBS counts for each of the m 250kb domains that are clustered together for each of the n transcription factors that were analyzed. The number m may change from cluster to cluster, while the number of considered transcription factors n is the same for all clusters. Let A be the set of TFs that have TFBS on multiple chromosomes in the considered cluster. Then the score function for each cluster is computed as follows:

$$Score(Z) = \sum_{j \in A} \sum_{i=1}^m Z_{ij} \quad (4)$$

For each matrix Z , a set of 1000 random matrices is generated to compute the background score distribution. Assuming that the number of TFBS for a specific transcription factor is independent of the other

transcription factors, a random matrix for a particular cluster with corresponding matrix Z is generated by the following procedure:

1. Let k denote the number of nonzero entries in Z . The probability of having a nonzero entry for each of the n transcription factors is defined by p_j , where $p_j = \frac{\# \text{ nonzero entries for TF}_j}{\text{total \# of nonzero entries}}$. The number of nonzero entries for each transcription factor, x_j , is drawn from a multinomial distribution, $(x_1, \dots, x_n) \sim \text{Mult}(k, p_1, \dots, p_n)$.
2. After determining the number of nonzero entries for each transcription factor, these nonzero entries must be distributed across the m clustered 250kb regions. Let q_i be the probability of assigning a nonzero entry to that specific region, where $q_i = \frac{\# \text{ nonzero entries in region}_i}{k}$. For each transcription factor j , the number of nonzero entries for each region, (y_{1j}, \dots, y_{mj}) is drawn from a multinomial, $(y_{1j} \dots y_{mj}) \sim \text{Mult}(x_j, q_1, \dots, q_m)$.
3. By now the positions of nonzero entries within the randomly generated matrix have been chosen, and only the number of TFBS (counts) remain to be assigned to each of the k nonzero entries. For each transcription factor, samples are drawn from the observed count distribution for that transcription factor over active clusters.

Finally, the p-value of the observed score was computed using the background score distribution that we obtained by calculating a score for each of the 1000 randomly generated matrices described above. In order to ensure stability of the ranking procedure, the background distribution was computed in 10 replicates, resulting in 10 different p-values. We observed that the p-values across different runs were consistent. For each replicate, we obtained the cluster rankings based on their p-values. The final rank of each cluster was computed from the median rank across the 10 replicate runs.

Negative controls - chromosomes that do not intermingle

As negative controls, we identified by a whole-genome analysis pairs of chromosomes that do not intermingle (in Hi-C) and are anti-correlated in terms of genomic features (Fig. S8). First, we determined chromosome pairs for which the LAS analysis did not result in any intermingling regions. These chromosomes formed the nodes of a network with edges drawn between pairs of chromosomes with no intermingling regions. The weight of the edges was calculated as $1 - |\rho|$, where ρ is the correlation between the genomic features averaged over the whole chromosome. Chromosomes 3 and 20 were chosen as a negative control pair, since they were representative of the network of non-intermingling chromosomes.

Cell culture and chromosome FISH

BJ fibroblast cells were cultured in Low Glucose DMEM (Life Technologies, USA) supplemented with 10% (vol/vol) FBS (GIBCO, Life Technologies, USA) at 37°C in 5% CO₂. BJ fibroblast cells were cultured overnight on fibronectin-coated cleaned glass slides. Cells were then washed with 1× PBS to remove cell culture medium followed by incubation on ice for 5-8 minutes, with 0.25% Triton in CSK buffer (100 mM NaCl, 300 mM Sucrose, 3 mM MgCl₂, 10 mM PIPES with pH 6.8). Cultured cells were fixed with 4% PFA (Paraformaldehyde) for 10 minutes, briefly rinsed with 0.1 M Tris-HCl and washed with 1× PBS wash. This was followed by permeabilization with 0.5% Triton for 10-15 minutes. The cells were then incubated overnight in 20% glycerol at 4°C and subjected to 5-6 freeze-thaw cycles in liquid nitrogen. After this, cells were washed with 1× PBS a few times, before and after treatment with 0.01% HCl for 5-10 minutes, followed by protein digestion with 0.002% porcine pepsin (Sigma Aldrich, USA) in 0.01N HCl at 37°C for 4 minutes. Cells were then fixed with 1% PFA for 4 minutes, briefly rinsed in 1× PBS before being treated with RNase (Promega, USA, 200 microgram/ml made in 2× SSC-0.3M sodium chloride and 30mM trisodium citrate) at 37°C for 15-20 minutes to digest RNA. The cells were then washed with 2× SSC and equilibrated in 50% Formamide / 2× SSC (pH 7.4) overnight at 4°C. Hybridization was set up the following day. Chromosome fish probes (Chrombios, Germany) tagged with different fluorophores were thawed to room temperature and mixed with hybridization buffer provided by the supplier. The DNA was denatured in 50% Formamide / 2× SSC at 85°C for 2-3 minutes and then incubated with the fluorescently labeled human chromosome FISH probe mix. The slides were then sealed with a Sigmacote (Sigma Aldrich, USA) coated hydrophobic coverslip and rubber cement to incubate for 18-48 hours in a moist chamber at 37°C with shaking. At the end of the incubation period, slides were washed three times in 50% Formamide / 2× SSC at 45°C and 0.1× SSC at 60°C. After the last stringent wash with 50% Formamide made in 0.1× SSC at 45°C, the nuclei were blocked in 5% BSA

solution made in $2\times$ SSC and then subjected to primary and secondary antibodies diluted in 5% BSA solution made in $2\times$ SSC. In case indirect labels such as chromosome probes conjugated with biotin and digoxigenin (DIG), were used during hybridization detection step, the procedure also involved the use of fluorophore labeled streptavidin/avidin and anti-DIG. The primary antibodies used here were: RNAPII CTD repeat YSPTSPS (phospho S5) (Abcam - ab5131, 1:500 dilution), mouse monoclonal (21H8) to DIG (Abcam-ab420; 1:500 dilution). Finally, the nuclei were stained with Hoechst 33342 (Sigma Aldrich, USA) for 10 minutes and then mounted with Prolong Gold antifade mounting medium (Life Technologies, USA), sealed with a coverslip, and imaged.

Confocal Imaging and Image Analysis

Slides for chromosome FISH were scanned using a Nikon A1 Confocal microscope (Nikon, USA) with a $100\times$, 1.4 NA oil objective. Stacks of 12-bit gray scale two-dimensional images were obtained with a pixel size of 130 nm in XY direction and 500 nm in the Z direction and used for the quantitative evaluation. The image analysis was performed using a custom code in ImageJ2. The code first identified the nuclear boundary using Otsu 3D thresholding method. This was followed by the identification of chromosome territories in the nuclear region using ReyniEntropy 3D thresholding. The threshold for identifying signal and background in each image was determined using the intensity histogram from the 3D image stack. The thresholded image was binarized. The overlapping region between two chromosomes, i.e. the intermingling region (IMR), was identified by performing the AND function over the 3D binary stacks of both chromosomes. The chromosome and IMR volumes were computed by summing up the volumes of the non-zero voxels in the respective binary images. The intermingling degree was calculated by dividing the volume of the IMR between two chromosomes by the total volume of the two chromosomes. Similarly, the amount of active RNAPII in the nucleus and the IMR was obtained by passing the RNAPII image and the binary images of the nucleus and IMR through the AND filter, respectively. The enrichment of active RNAPII in the intermingling regions was obtained by dividing the mean intensity of active RNAPII in the IMR by the mean intensity of the active RNAPII in the entire nucleus. R was used for testing statistical significance and for data visualization.

References

- [1] P. Bardou, J. Mariette, F. Escudié, C. Djemiel, and C. Klopp. jvenn: an interactive Venn diagram viewer. *BMC Bioinformatics*, 15(1):293, 2014.
- [2] T. F. Consortium, the RIKEN PMI, and C. (DGT). A promoter-level mammalian expression atlas. *Nature*, 507(7493):462–470, 2014.
- [3] N. C. Durand, J. T. Robinson, M. S. Shamim, I. Machol, J. P. Mesirov, E. S. Lander, and E. L. Aiden. Juicebox Provides a Visualization System for Hi-C Contact Maps with Unlimited Zoom. *Cell Systems*, 3(1):99–101, 2016.
- [4] M. Elsner and W. Schudy. Bounding and comparing methods for correlation clustering beyond ILP. *Proceedings of the Workshop on Integer Linear Programming for Natural Language Processing*, (June):19–27, 2009.
- [5] ENCODE Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature*, 489(7414):57–74, 2012.
- [6] D. W. Huang, R. a. Lempicki, and B. T. Sherman. Systematic and integrative analysis of large gene lists using DAVID bioinformatics resources. *Nature Protocols*, 4(1):44–57, 2009.
- [7] D. W. Huang, B. T. Sherman, and R. A. Lempicki. Bioinformatics enrichment tools: Paths toward the comprehensive functional analysis of large gene lists. *Nucleic Acids Research*, 37(1):1–13, 2009.
- [8] J. H. Korhonen, K. Palin, J. Taipale, and E. Ukkonen. Fast motif matching revisited: high-order PWMs, SNPs and indels. *Bioinformatics*, 33(4):514–521, 2016.
- [9] E. Lieberman-aiden, N. L. V. Berkum, L. Williams, M. Imakaev, T. Ragoczy, A. Telling, I. Amit, B. R. Lajoie, P. J. Sabo, M. O. Dorschner, R. Sandstrom, B. Bernstein, M. A. Bender, M. Groudine, A. Gnirke, J. Stamatoyannopoulos, and L. A. Mirny. Comprehensive Mapping of Long-Range

- Interactions Reveals Folding Principles of the Human Genome. *Science*, 326(October):289–293, 2009.
- [10] D. Marbach, D. Lamparter, G. Quon, M. Kellis, Z. Kutalik, and S. Bergmann. Tissue-specific regulatory circuits reveal variable modular perturbations across complex diseases. *Nature Methods*, 13(4):366–370, 2016.
- [11] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and É. Duchesnay. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2012.
- [12] S. S. P. Rao, M. H. Huntley, N. C. Durand, E. K. Stamenova, I. D. Bochkov, J. T. Robinson, A. L. Sanborn, I. Machol, A. D. Omer, E. S. Lander, and E. L. Aiden. A 3D map of the human genome at kilobase resolution reveals principles of chromatin looping. *Cell*, 159(7):1665–1680, 2014.
- [13] A. A. Shabalín, V. J. Weigman, C. M. Perou, and A. B. Nobel. Finding large average submatrices in high dimensional data. *Annals of Applied Statistics*, 3(3):985–1012, 2009.

SI Tables

Table S1: Total sizes of intermingling regions

Region	Size
Intermingling domains from LAS	903.25 Mb
Intermingling regions after clustering	459.5 Mb
Intermingling regions in active clusters	179.75 Mb

Table S2: Datasets and accessions used to obtain the genomic features. The symbol * indicates that peaks were retrieved from a previous study (49) that had already pre-processed the data from the GEO database.

Name	Accession	Category
RNA-seq	GSE24565	active
RNAPII	GSE31477	active
H3K4me1	GSE16256	active
H3K4me2	GSE16256	active
H3K4me3	GSE16256	active
H3K36me3	GSE16256	active
H3K9ac	GSE16256	active
H3K27me3	GSE16256	repressive
H3K9me3	GSE16256	repressive
YAP1*	GSE61852	other
RFX5	GSE31477	other
RELA*	GSE43070	other
RCOR1	GSE31477	other
RBL2*	GSE19899	other
RB1*	GSE19899	other
RAD21	GSE31477	other
MXI1	GSE31477	other
MECP2*	GSE47678	other
MAZ	GSE31477	other
MAFK	GSE31477	other
LMNB1*	GSE53332	other
H4K91ac	GSE16256	other
H4K8ac	GSE16256	other
H4K5ac	GSE16256	other
H4K20me1	GSE16256	other
H3K9me1	GSE16256	other
H3K79me2	GSE16256	other
H3K79me1	GSE16256	other
H3K56ac	GSE16256	other
H3K4ac	GSE16256	other
H3K27ac	GSE16256	other
H3K23ac	GSE16256	other
H3K18ac	GSE16256	other
H3K14ac	GSE16256	other
H2BK5ac	GSE16256	other
H2BK20ac	GSE16256	other
H2BK15ac	GSE16256	other
H2BK12ac	GSE16256	other
H2BK120ac	GSE16256	other
H2A.Z	GSE16256	other
MacroH2A1.1	GSE54847	other
H2AK9ac	GSE16256	other
H2AK5ac	GSE16256	other

Table S3 : Clusters identified via network analysis, enrichment in genomic features, and cluster type.

Clusters (kb)	H3K36me3	RNAPII	H3K9ac	H3K4me3	H3K4me1	H3K9me3	Label
chr15:64000, chr10:32250	1.699	1.183	1.174	0.914	1.443	0.499	active
chr2:211750, chr4:179750	0.014	0.000	0.056	0.000	0.004	2.129	inactive
chr21:43250, chr5:250	4.284	7.366	3.647	3.850	3.270	0.512	active
chr1:26750, chr7:99750, chr22:39250, chr17:76000, chr1:26500, chr11:66250	2.601	2.151	2.782	3.094	1.654	0.497	active
chr7:100250, chr7:100750, chr11:65250, chr17:42750	2.942	10.116	4.162	6.824	2.254	0.542	active
chr12:125000, chr2:217500, chr3:197000	1.204	1.938	3.385	2.155	2.996	0.851	active
chr2:224250, chr1:241750, chr2:224500, chr2:223750	0.875	1.432	1.254	0.830	1.844	0.410	active
chr21:34750, chr20:3750	4.801	4.828	3.022	4.542	2.644	0.161	active
chr16:4750, chr19:11000	4.777	3.499	3.601	3.363	1.879	0.266	active
chr14:102250, chr12:123750	4.778	1.552	2.580	2.851	1.959	0.127	active
chr15:47000, chr9:69250	0.002	0.000	0.005	0.103	0.011	1.181	inactive
chr21:26250, chr20:7000, chr20:7250, chr20:6750	0.034	0.000	0.129	0.081	0.048	1.327	inactive
chr6:163500, chr1:241000, chr1:241250	0.020	0.000	0.180	0.078	0.071	2.818	inactive
chr22:19500, chr17:4250, chr22:18750	0.716	1.405	1.946	1.725	1.378	0.748	active
chr12:50000, chr6:35250	3.305	5.211	4.071	3.695	2.959	0.311	active
chr22:40000, chr10:72250, chr15:73500	0.890	0.575	1.361	0.327	1.460	1.203	inactive
chr20:33500, chr14:75250, chr15:68500	3.449	1.319	2.876	2.120	2.818	0.319	active
chr4:96750, chr9:44250	0.000	0.000	0.009	0.051	0.007	1.133	inactive
chr5:172000, chr2:239750	1.693	4.227	2.929	2.526	2.330	0.867	active
chr20:12500, chr21:27750, chr20:12250	0.008	0.000	0.083	0.098	0.029	2.230	inactive
chr7:101250, chr7:101750, chr17:37250, chr9:134500	2.647	1.083	2.300	1.976	2.017	0.560	active
chr9:24750, chr15:23250	0.155	0.000	0.037	0.142	0.000	1.112	inactive
chr14:74500, chr15:64250, chr17:43500	2.590	1.404	1.752	1.954	1.235	0.435	active
chr2:226500, chr3:190000	0.076	1.441	0.690	0.922	0.923	0.321	active
chr17:80750, chr14:103250, chr14:103000	4.074	1.744	1.953	1.522	2.308	0.583	active
chr6:159500, chr1:236750, chr6:159750	0.512	0.466	0.649	0.562	0.344	1.650	inactive
chr4:167250, chr11:89500	0.004	0.000	0.013	0.000	0.009	1.149	inactive
chr21:43500, chr22:28000	2.108	0.431	2.572	1.525	2.187	1.250	inactive
chr11:2250, chr4:3500, chr11:2000	0.902	0.648	1.810	1.867	1.079	1.698	inactive
chr20:59500, chr20:58750, chr20:59000, chr20:59250, chr19:32500	0.048	0.173	0.346	0.117	0.240	2.953	inactive
chr6:167000, chr1:15250	2.255	3.144	2.708	1.945	2.622	0.928	active
chr14:75000, chr20:33750	3.305	3.009	3.358	2.495	3.306	0.230	active
chr3:126500, chr1:26250, chr3:126250	2.018	1.263	2.060	1.802	1.826	0.652	active
chr17:18000, chr22:43500, chr21:34500	2.844	2.199	3.090	2.922	2.016	0.643	active
chr5:171000, chr2:239500	0.442	0.000	0.953	0.167	0.972	1.879	inactive
chr22:26250, chr20:22000, chr20:21750	0.011	0.000	0.101	0.000	0.000	3.180	inactive
chr22:31000, chr17:41750	1.296	3.118	2.121	2.411	1.755	0.510	active
chr16:18500, chr7:85750	0.153	0.000	0.089	0.255	0.033	1.084	inactive
chr7:77250, chr4:129000	3.458	1.922	1.859	1.990	2.510	0.510	active
chr9:67250, chr3:96000	0.038	0.000	0.147	0.226	0.123	1.328	inactive
chr22:50750, chr17:7250	4.009	8.212	4.114	7.026	2.330	0.379	active
chr1:8000, chr1:8250, chr2:201500	1.912	8.022	3.473	2.939	3.169	0.337	active
chr21:25750, chr18:2000, chr21:25250, chr14:26000	0.004	0.000	0.063	0.101	0.083	2.095	inactive
chr1:13000, chr2:156500, chr2:156750	0.001	0.000	0.028	0.016	0.001	1.749	inactive
chr1:148750, chr7:122500	0.258	0.431	0.275	0.816	0.223	1.363	inactive
chr11:74500, chr14:50500	3.703	2.628	1.569	1.047	1.424	0.485	active
chr17:48750, chr12:56250	2.341	3.686	3.359	4.556	2.131	0.355	active
chr9:43500, chr4:35000	0.006	0.000	0.005	0.050	0.000	1.248	inactive
chr17:74250, chr22:45500	3.395	9.272	3.875	3.819	3.048	0.237	active
chr18:50250, chr6:94750	0.008	0.000	0.013	0.000	0.000	2.185	inactive
chr19:46500, chr15:74000	1.264	0.885	1.669	1.725	1.270	1.580	inactive
chr5:138750, chr1:9750	3.071	2.643	2.611	1.631	2.488	0.278	active
chr4:132250, chr6:123500	0.004	0.000	0.071	0.000	0.003	2.125	inactive
chr13:42000, chr15:67000	2.075	1.369	3.628	1.317	3.524	0.611	active
chr12:750, chr15:40250, chr11:73750	3.137	1.712	1.928	1.926	2.219	0.364	active
chr2:236500, chr4:186000	2.826	1.503	1.840	1.479	2.496	0.233	active
chr3:192500, chr2:227500	1.484	1.191	2.374	3.144	2.752	0.285	active
chr7:49000, chr16:32750	0.000	0.000	0.045	0.111	0.002	1.675	inactive
chr12:27000, chr21:38250	2.675	2.585	1.401	2.292	1.615	0.464	active
chr22:24500, chr21:40500	3.886	2.481	2.730	2.319	2.185	0.331	active
chr18:51250, chr9:44000	0.011	0.000	0.036	0.102	0.007	1.661	inactive
chr6:36000, chr12:49750	3.422	1.485	1.614	1.645	2.239	0.529	active
chr1:1250, chr12:124750	4.612	5.699	4.692	5.677	3.030	0.619	active
chr16:33000, chr14:48500	0.008	0.000	0.037	0.013	0.000	1.104	inactive
chr17:62250, chr15:63750	2.392	3.266	1.613	1.707	2.339	0.282	active
chr10:2750, chr8:4750	0.029	0.000	0.179	0.034	0.130	2.943	inactive
chr21:43750, chr4:6000, chr4:6250	0.451	0.288	1.150	0.919	0.448	1.762	inactive
chr1:20750, chr12:122500	2.355	2.797	2.540	3.067	1.713	0.587	active
chr1:25750, chr8:124000, chr15:70750, chr11:74000	2.119	1.210	2.020	1.927	2.134	0.727	active
chr9:68250, chr1:143000	0.135	0.000	0.123	0.592	0.133	1.173	inactive
chr9:134750, chr21:38000, chr22:29250, chr12:2000	2.299	1.597	1.721	2.063	2.350	0.737	active
chr2:211500, chr4:179250	0.055	0.000	0.032	0.000	0.002	2.138	inactive
chr4:65000, chr9:69000	0.082	0.000	0.174	0.538	0.097	1.085	inactive
chr4:97500, chr1:143250	0.027	0.000	0.045	0.113	0.000	2.069	inactive
chr11:63750, chr16:67750	4.255	4.689	4.485	4.824	2.781	0.098	active
chr2:13750, chr7:74500	0.069	0.000	0.086	0.125	0.056	1.168	inactive
chr12:53750, chr20:34250	4.557	5.981	3.342	3.659	3.046	0.175	active
chr3:47250, chr12:122250	3.474	2.642	2.998	2.097	2.140	0.265	active
chr17:45500, chr16:70000	3.437	1.180	1.407	1.455	1.193	0.082	active
chr15:75500, chr11:0	2.093	3.742	1.518	2.394	0.810	0.197	active
chr1:239250, chr6:163000	0.008	0.000	0.082	0.150	0.026	2.992	inactive
chr4:178750, chr2:214250, chr2:214500	0.004	0.000	0.128	0.049	0.178	1.321	inactive
chr1:27750, chr3:46750	2.212	1.008	3.330	2.640	2.765	0.589	active
chr16:86250, chr12:22250	1.691	4.227	4.145	1.433	4.180	0.598	active
chr7:102000, chr12:121750	3.615	3.167	2.451	2.288	1.243	0.105	active
chr10:75250, chr15:43250	3.421	1.987	1.505	1.929	1.410	0.260	active
chr16:28750, chr19:11250	3.369	4.231	3.245	4.878	1.669	0.301	active
chr16:30750, chr19:17250, chr16:31000	3.641	5.382	3.843	6.132	1.980	0.325	active
chr6:161000, chr1:239500	0.017	0.000	0.091	0.144	0.012	2.839	inactive
chr10:1500, chr8:2750	0.025	0.000	0.192	0.000	0.076	3.523	inactive
chr6:164750, chr5:166000	0.013	0.000	0.037	0.000	0.000	3.116	inactive

Table S4: Fold enrichment of high-occupancy target (HOT) regions as compared to low-occupancy target regions (LOT) in active clusters.

Feature	Enrichment
HOT:LOT	2.94

Table S5: Top 15 inactive clusters, ranked by enrichment for H3K9me3. Clusters are given by chromosome number and start position in kb; each region in the cluster is 250kb in length. Only clusters with less than seven 250kb regions in the cluster are included.

Cluster #	Clusters (kb)	H3K9me3	RNAPII
443	chr10:1500, chr8:2750	3.523	0.000
114	chr10:125000, chr1:5750, chr10:125250	3.377	0.000
204	chr8:137000, chr6:164500	3.353	0.000
213	chr17:11500, chr17:10750, chr5:6250, chr17:11000, chr17:11250	3.264	0.079
260	chr2:0, chr5:178250	3.225	0.431
370	chr22:26250, chr20:22000, chr20:21750	3.180	0.000
63	chr6:170250, chr1:13750	3.152	0.000
316	chr1:238250, chr6:162250	3.127	0.000
445	chr6:164750, chr5:166000	3.116	0.000
49	chr1:34750, chr1:34000, chr1:34250, chr8:142500, chr1:35000, chr1:34500	3.054	0.000
265	chr19:30750, chr20:19250	3.048	0.000
251	chr6:165250, chr8:138500	3.019	0.000
235	chr21:32250, chr20:15250	3.001	0.431
230	chr14:104500, chr14:104750, chr13:112000	2.994	0.000
434	chr1:239250, chr6:163000	2.992	0.000

Table S6: Top 15 active clusters, ranked by p-value of permutation test based on TFBS (JASPAR 2016, threshold = 0.000001, CAGE). Clusters are given by chromosome number and start position in kb; each region in the cluster is 250kb in length.

Cluster #	Clusters (kb)	P-value	RNAPII	H3K9me3
137	chr7:100250, chr7:100750, chr11:65250, chr17:42750	0.0000	10.1162	0.5417
111	chr1:22250, chr19:45250	0.0000	4.3250	0.2820
61	chr16:30000, chr19:10750, chr16:29750, chr19:10500	0.0002	4.9110	0.3573
57	chr19:56000, chr19:55500, chr22:50250, chr19:55750	0.0001	8.1535	0.9214
157	chr22:50750, chr17:7250	0.0005	8.2122	0.3790
27	chr20:48750, chr20:48500, chr21:47250, chr22:30500	0.0021	13.5989	0.4189
29	chr12:53500, chr4:1000, chr12:53250	0.0024	5.0042	0.5440
144	chr12:50000, chr6:35250	0.0031	5.2105	0.3111
92	chr1:1000, chr12:123250, chr1:750, chr9:133500, chr8:145000, chr13:114750	0.0043	9.0570	0.5125
180	chr1:27750, chr3:46750	0.0068	1.0079	0.5889
84	chr17:73000, chr8:144250	0.0068	3.5725	0.5660
68	chr2:220250, chr1:16750, chr2:220000	0.0070	5.4663	0.4099
185	chr16:30750, chr19:17250, chr16:31000	0.0073	5.3819	0.3246
17	chr12:49250, chr17:38250, chr12:49500	0.0075	6.0798	0.2533
78	chr1:203250, chr17:74500	0.0103	7.9689	0.5743

SI Figures

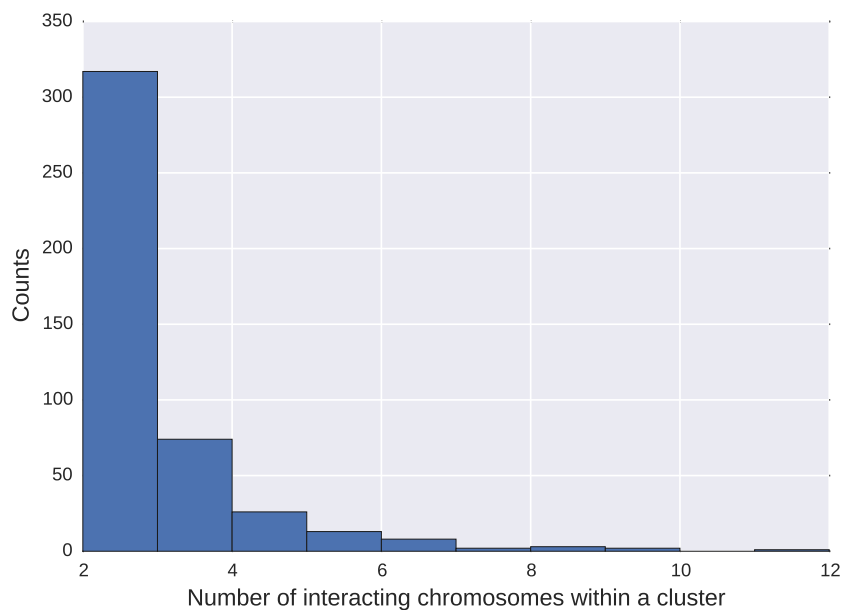


Figure S1: Number of chromosomes within an intermingling cluster.

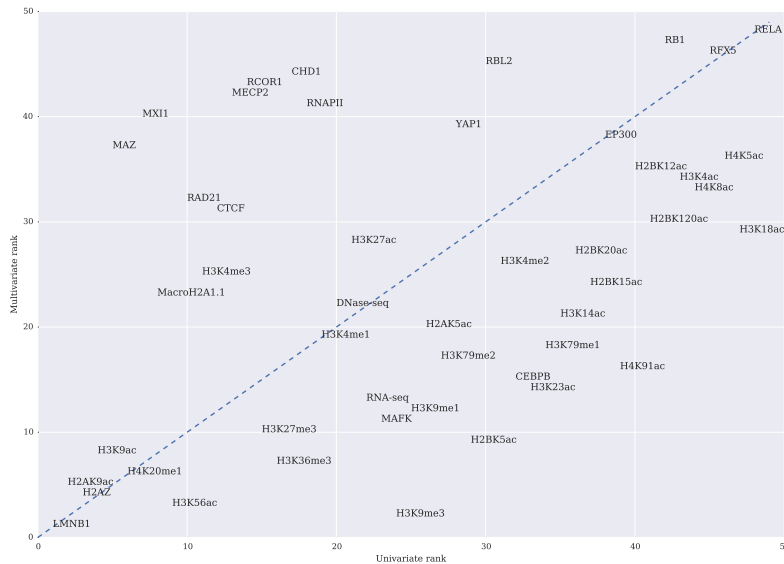


Figure S2: Scatterplot of univariate rank of a feature (two-sample Kolmogorov–Smirnov test) versus predictive rank of a feature when it is combined with all other features (relative importance in decision trees) for classification of intermingling versus non-intermingling regions. A similar analysis has been performed in (49) to explore context-dependency of features for classification.

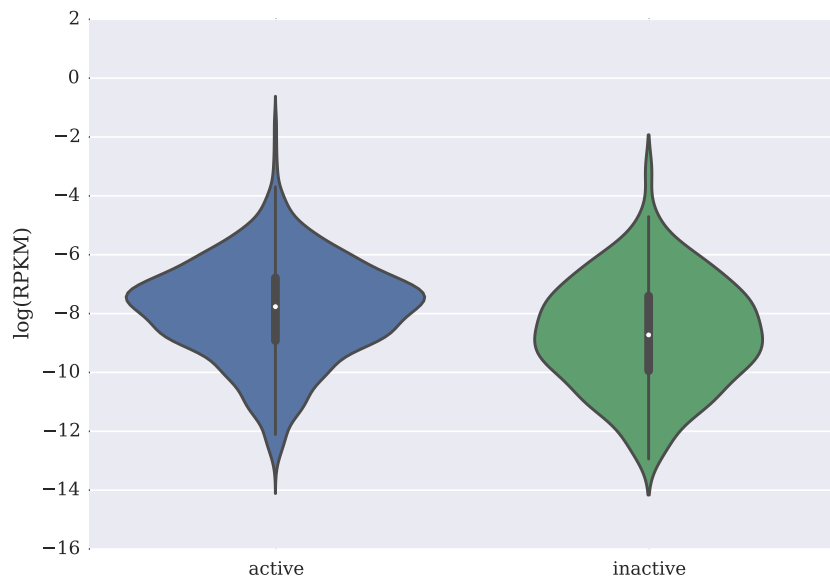
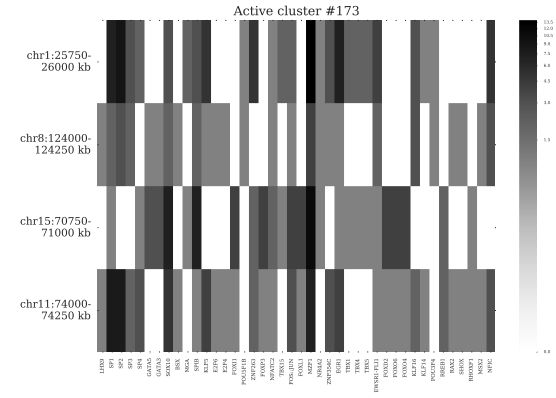
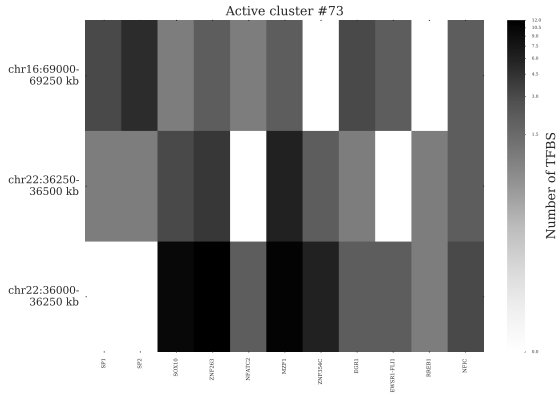
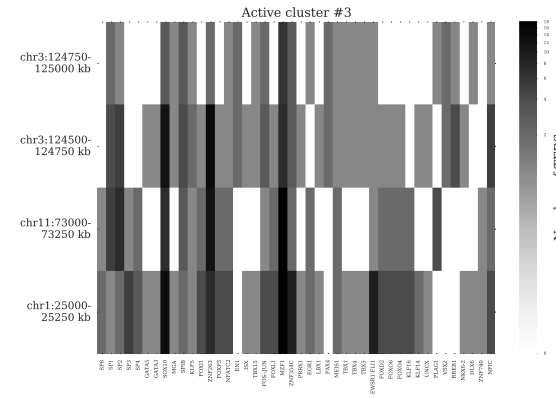
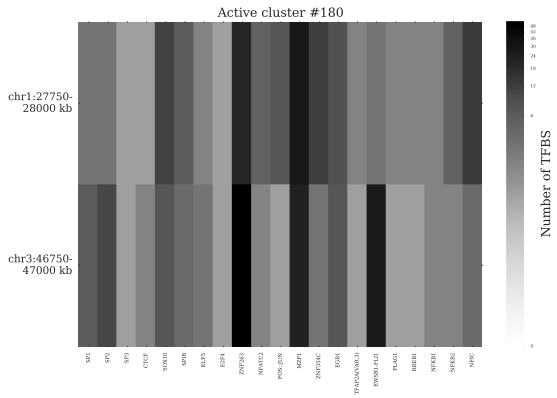
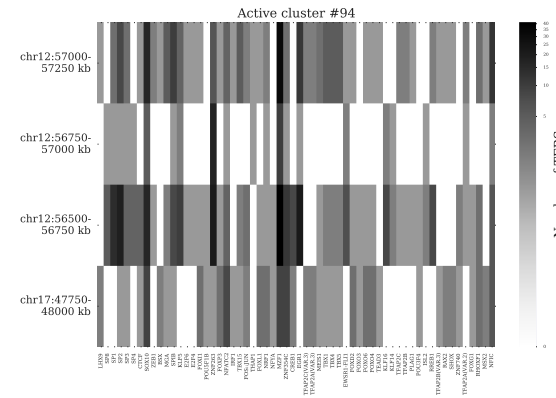
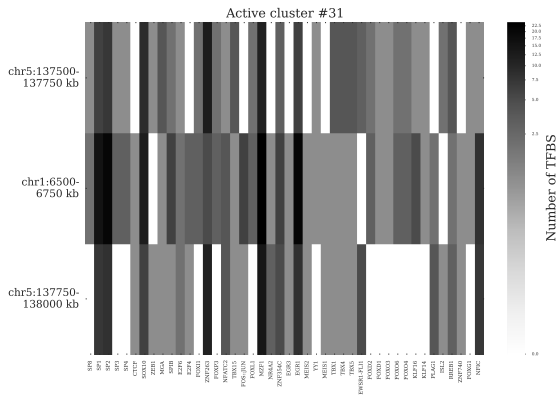
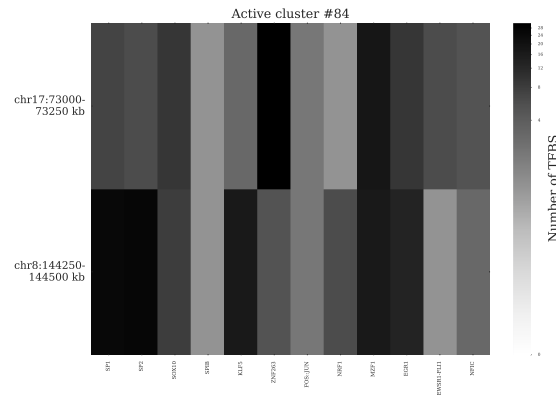
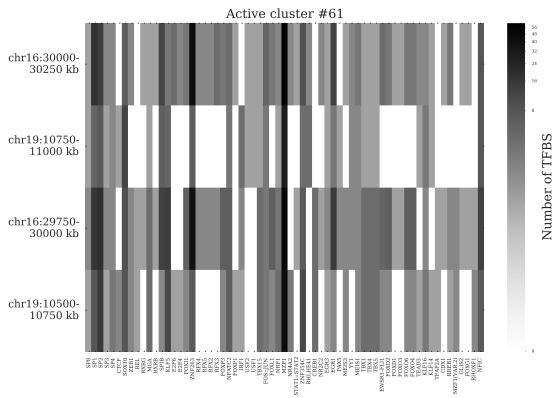
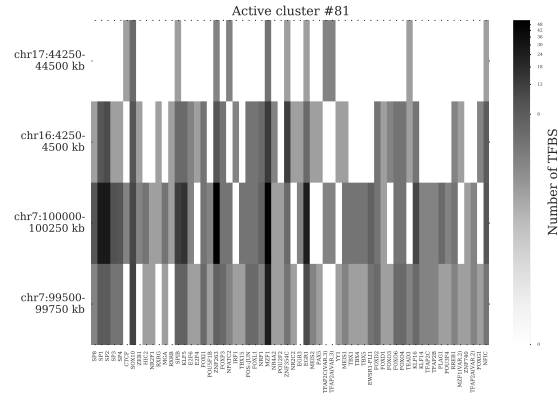
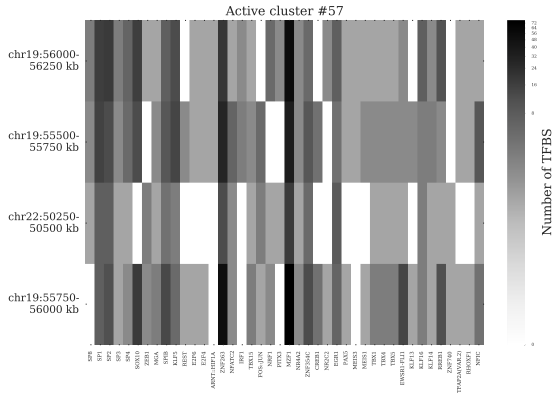
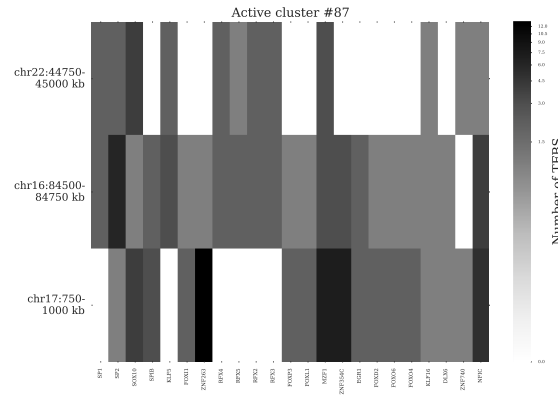
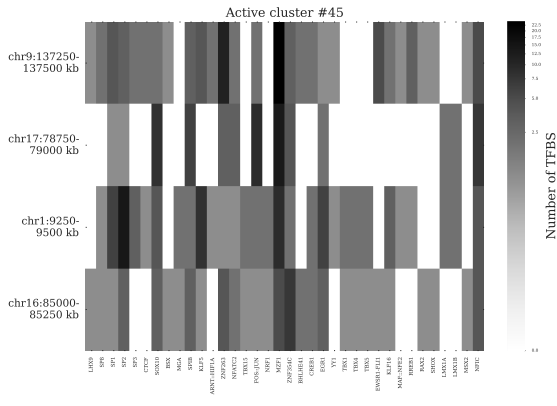
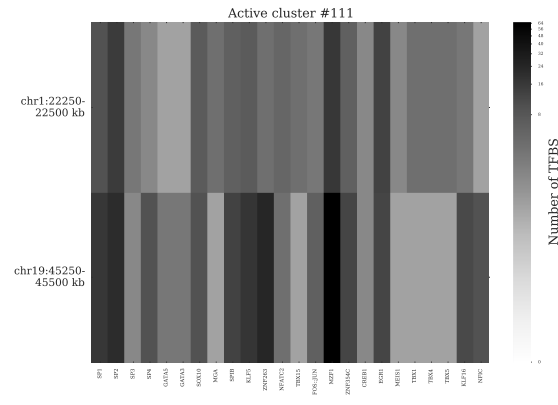
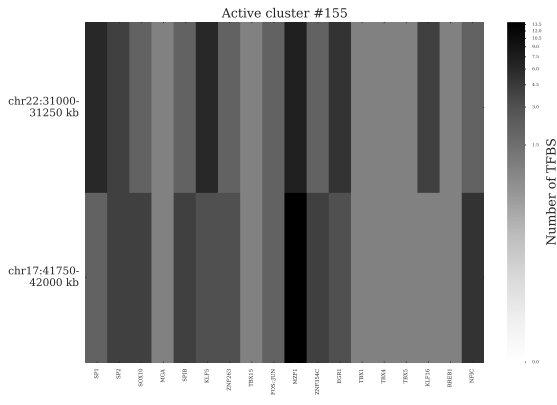
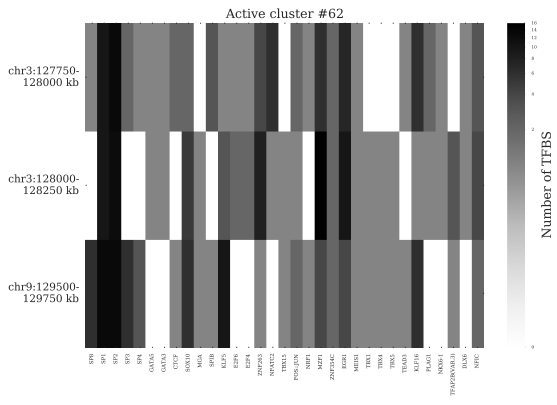


Figure S3: Comparison of gene expression in reads per kilobase of transcript per million mapped reads (RPKM) on log scale between active and inactive clusters. Active clusters show significantly higher gene expression (p-value = 0.004 under t-test).





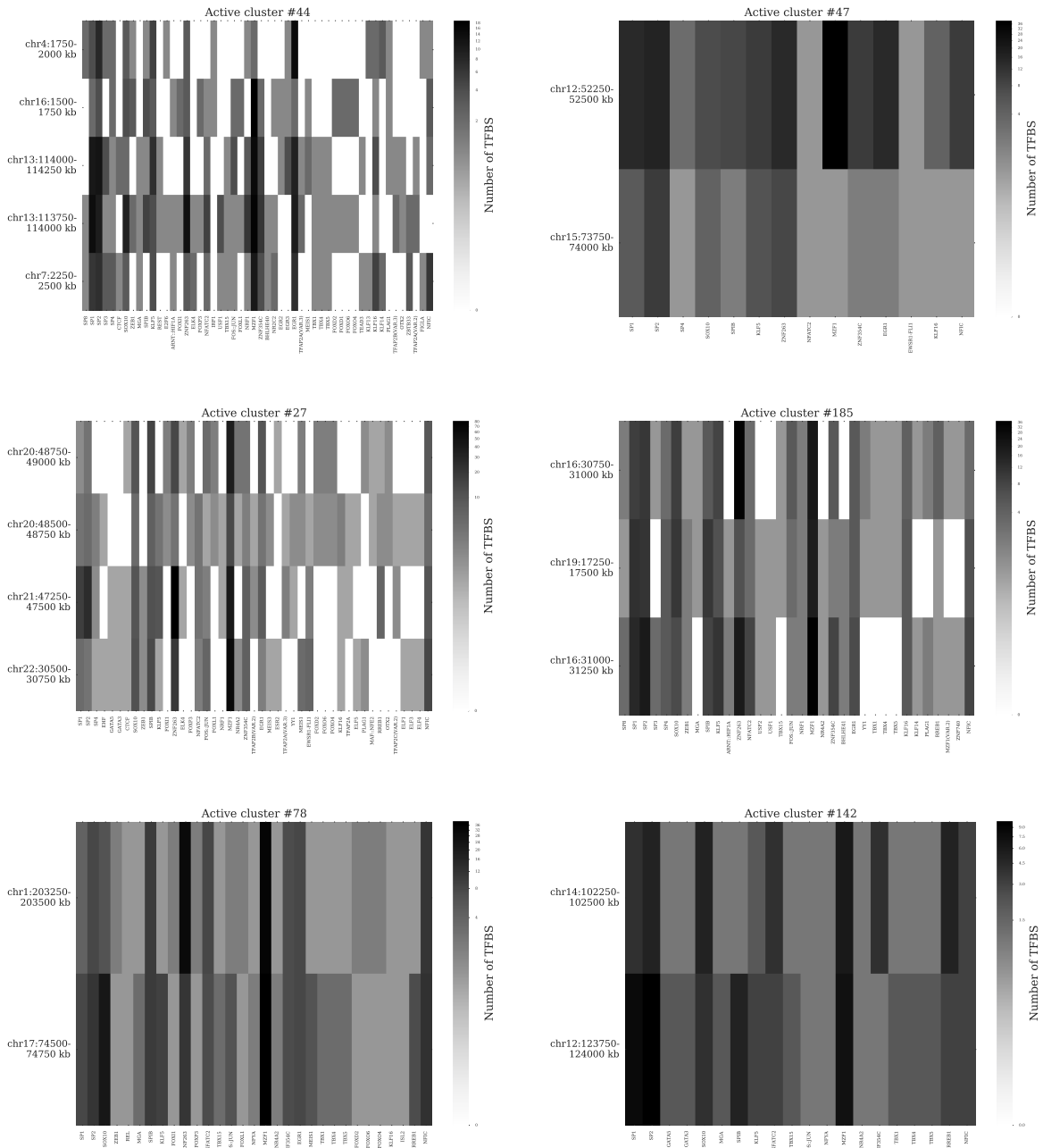
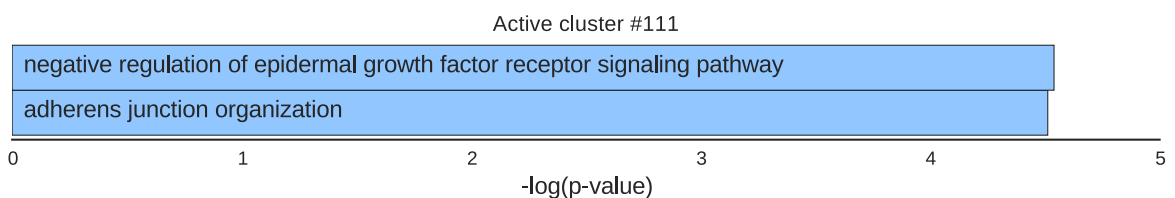
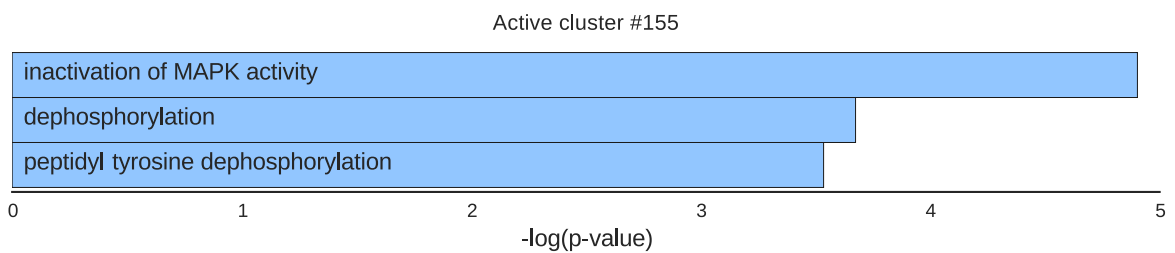
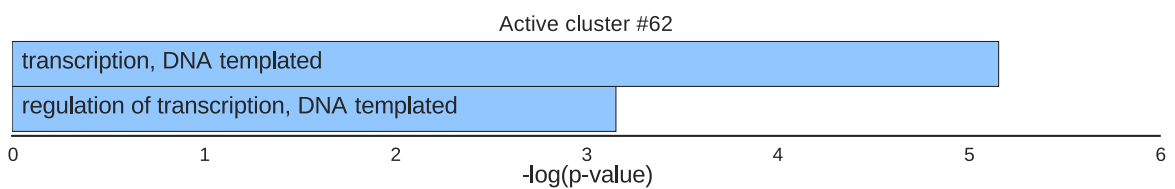
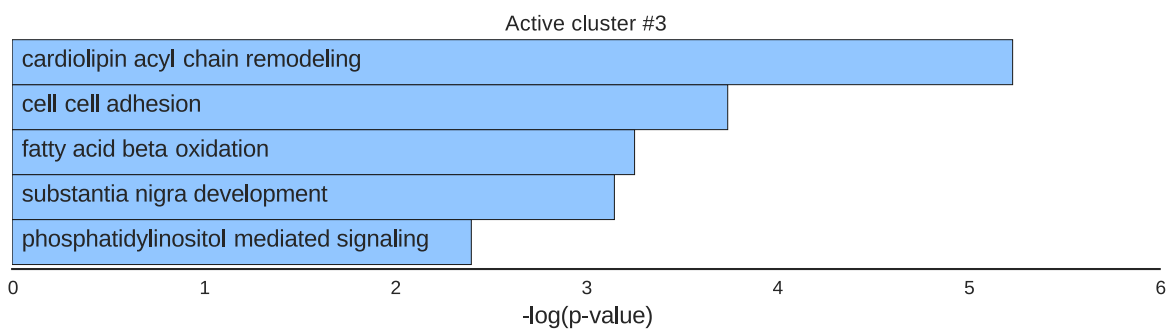
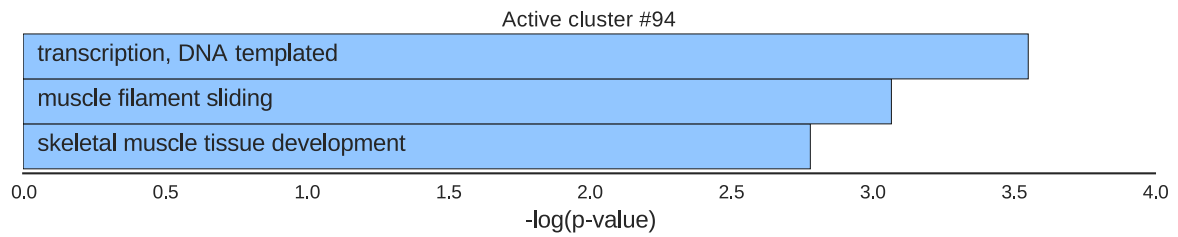
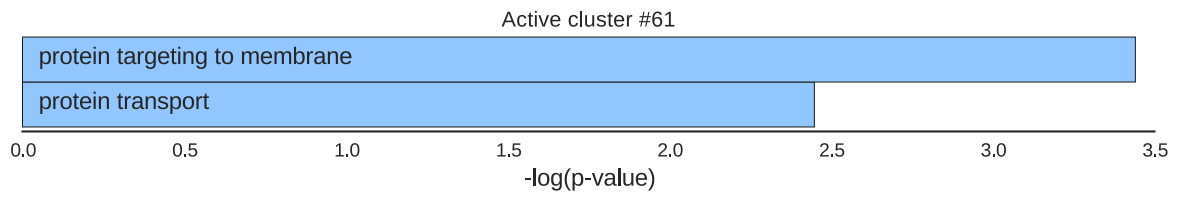
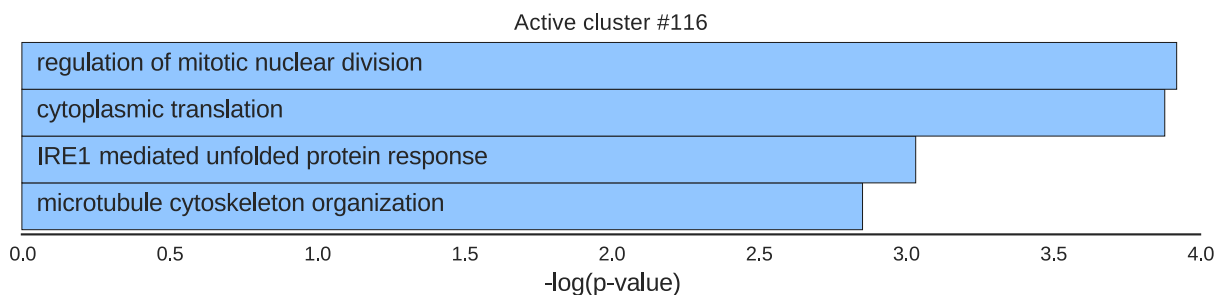
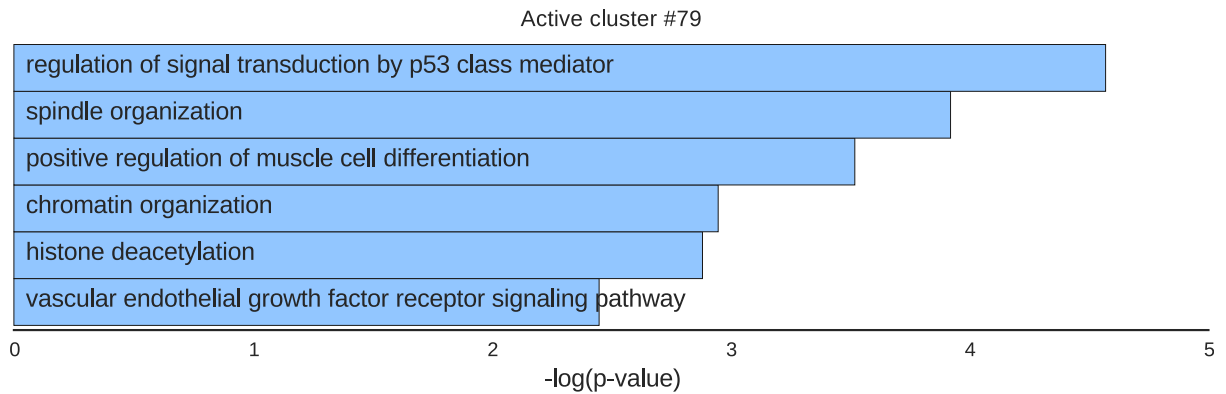
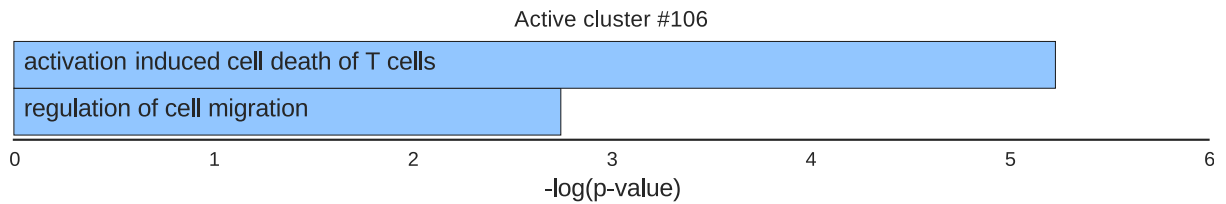
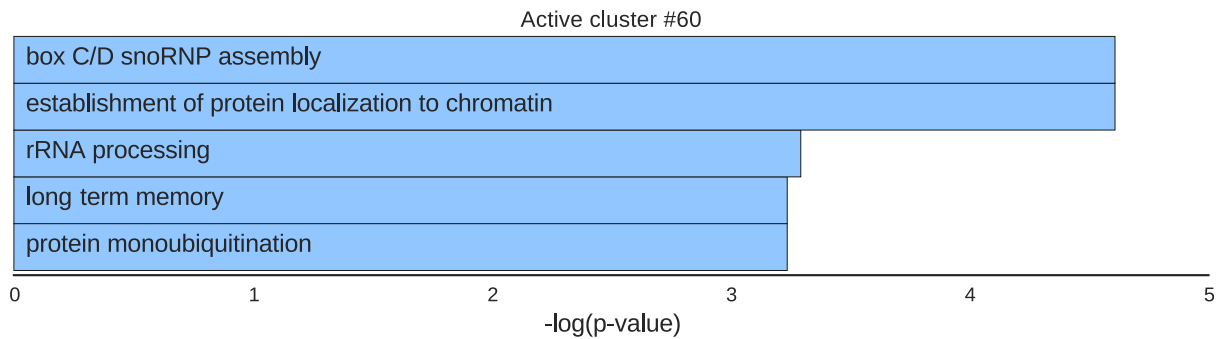
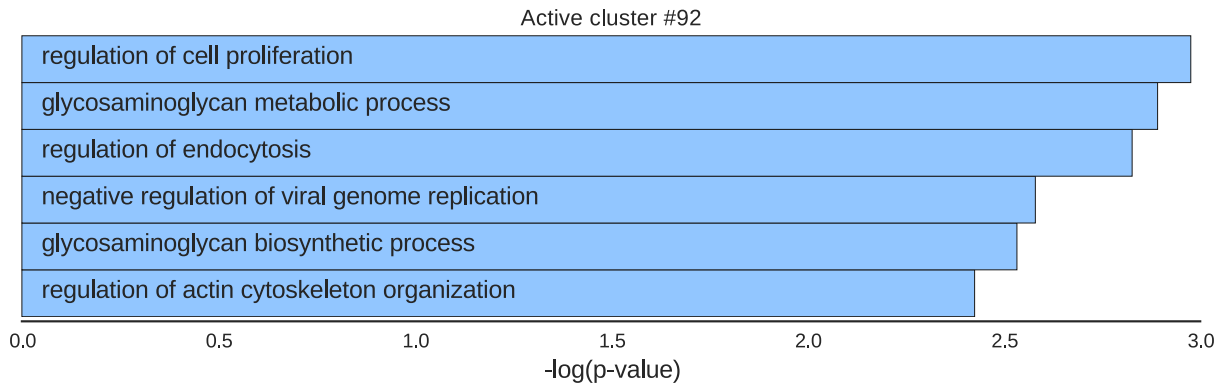
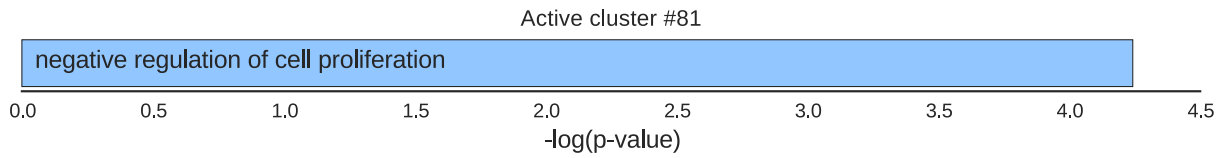
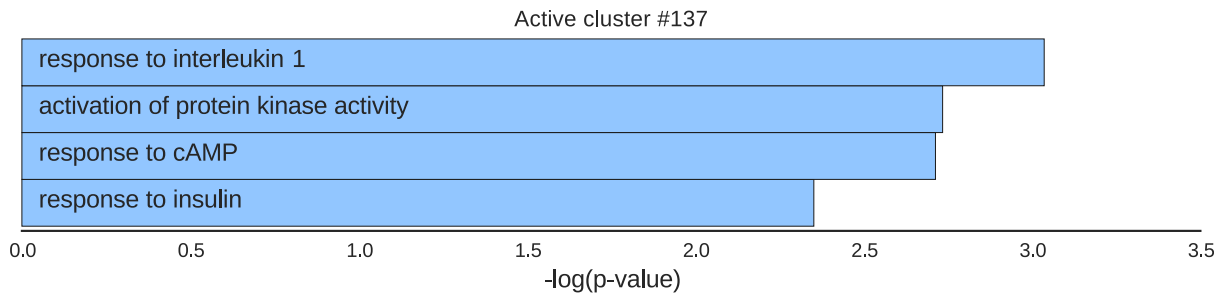
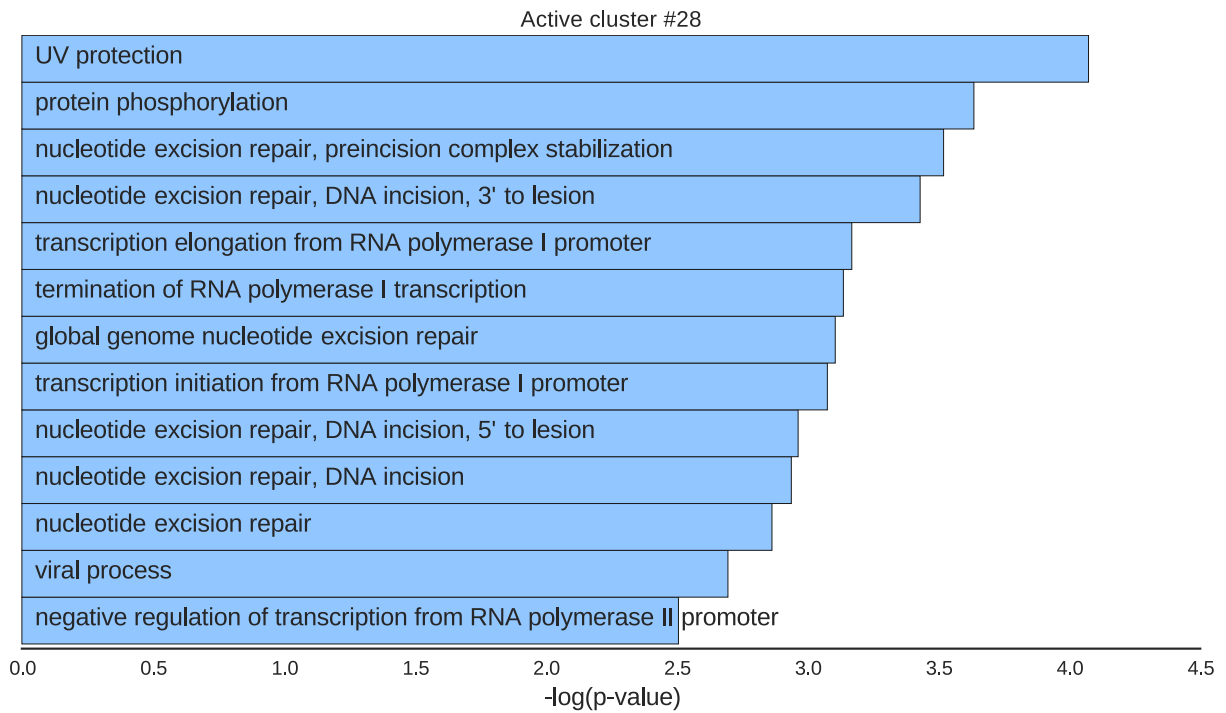
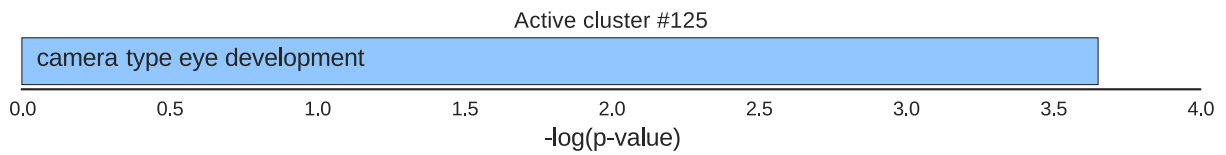
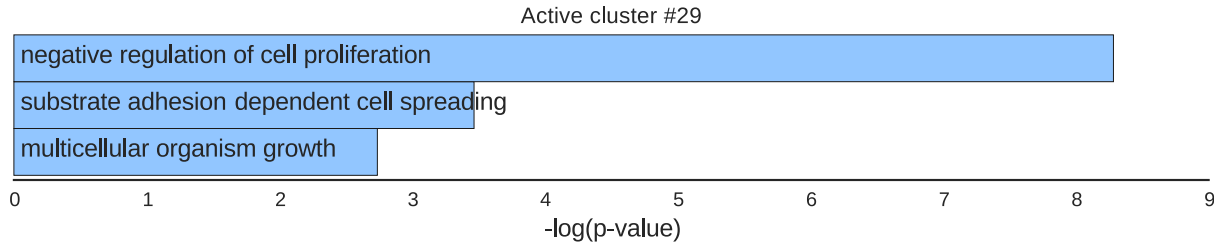
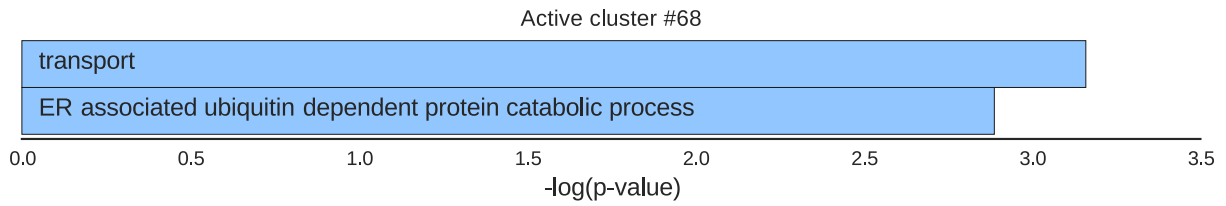


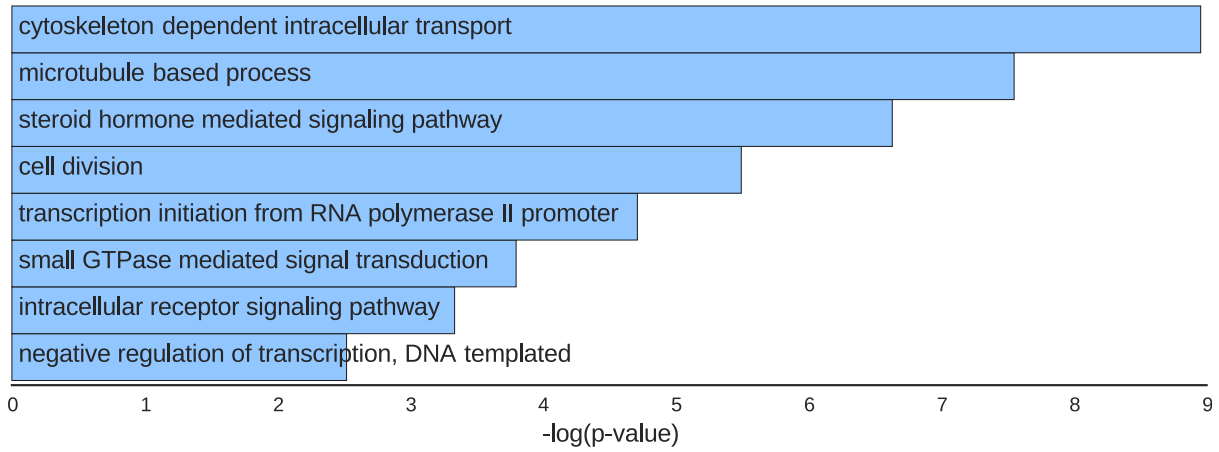
Figure S4: TFBS for top 15 active clusters given in Fig. S7 (JASPAR 2016, threshold = 0.000001, CAGE). Only TFs that span at least 2 regions in the cluster are shown.



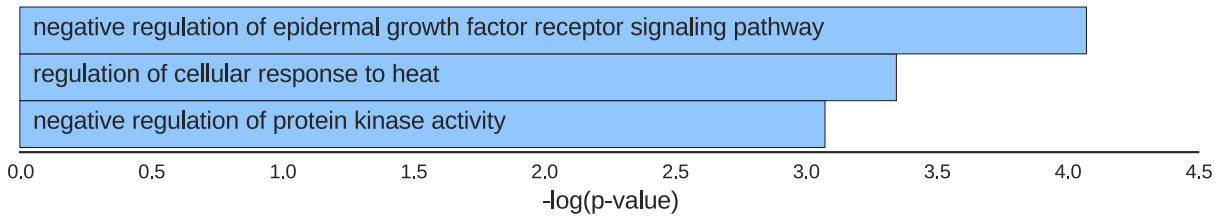




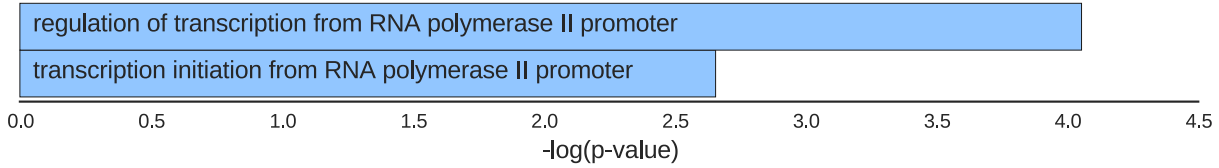
Active cluster #17



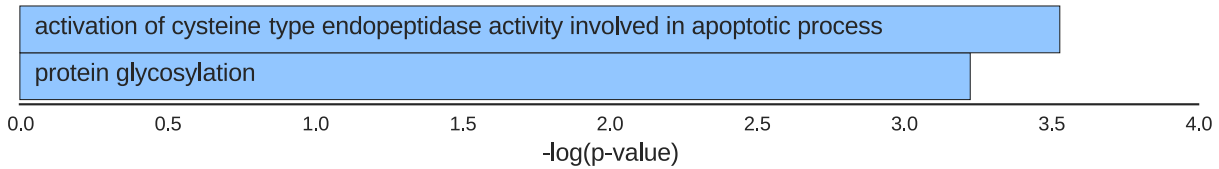
Active cluster #12



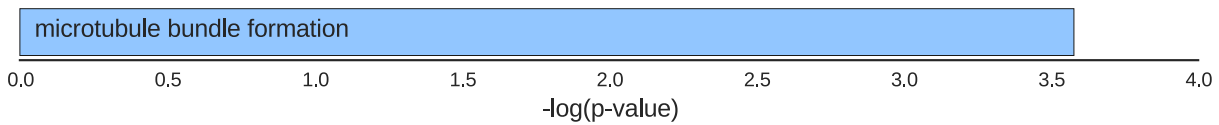
Active cluster #144



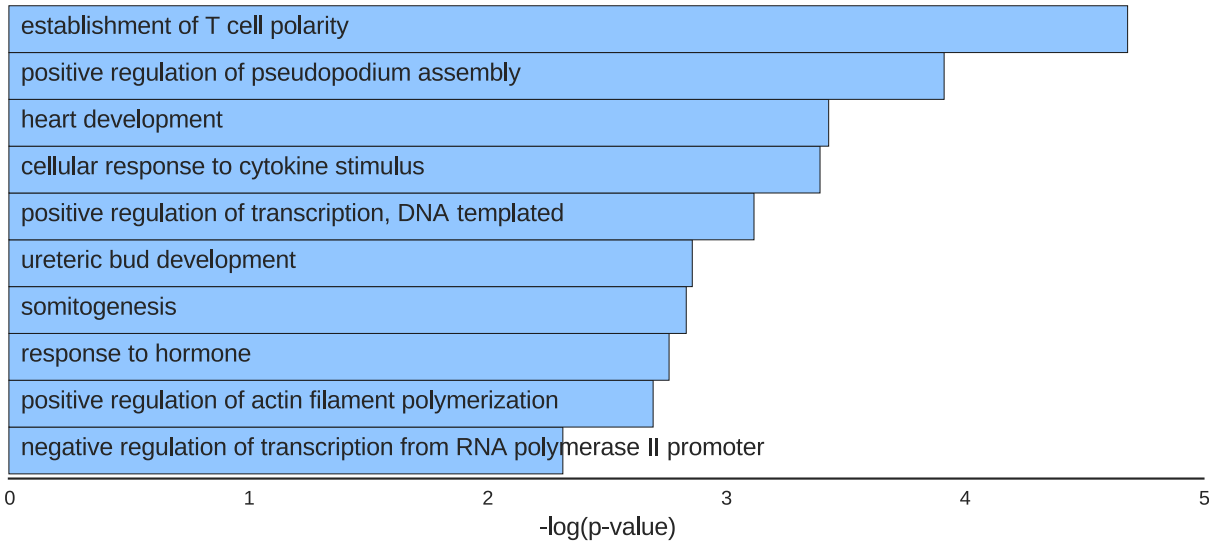
Active cluster #172



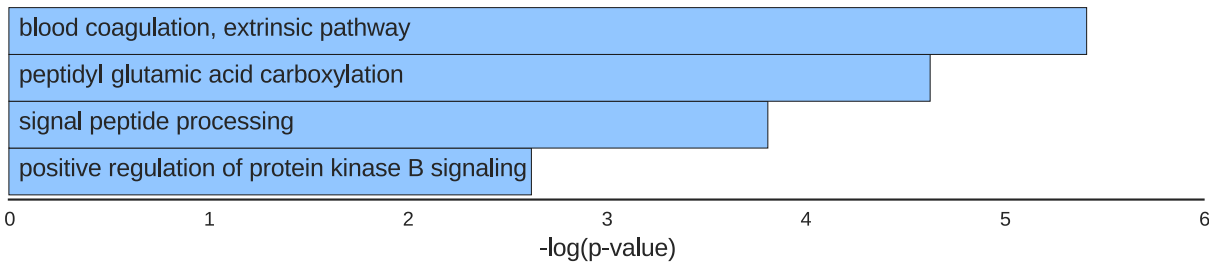
Active cluster #4



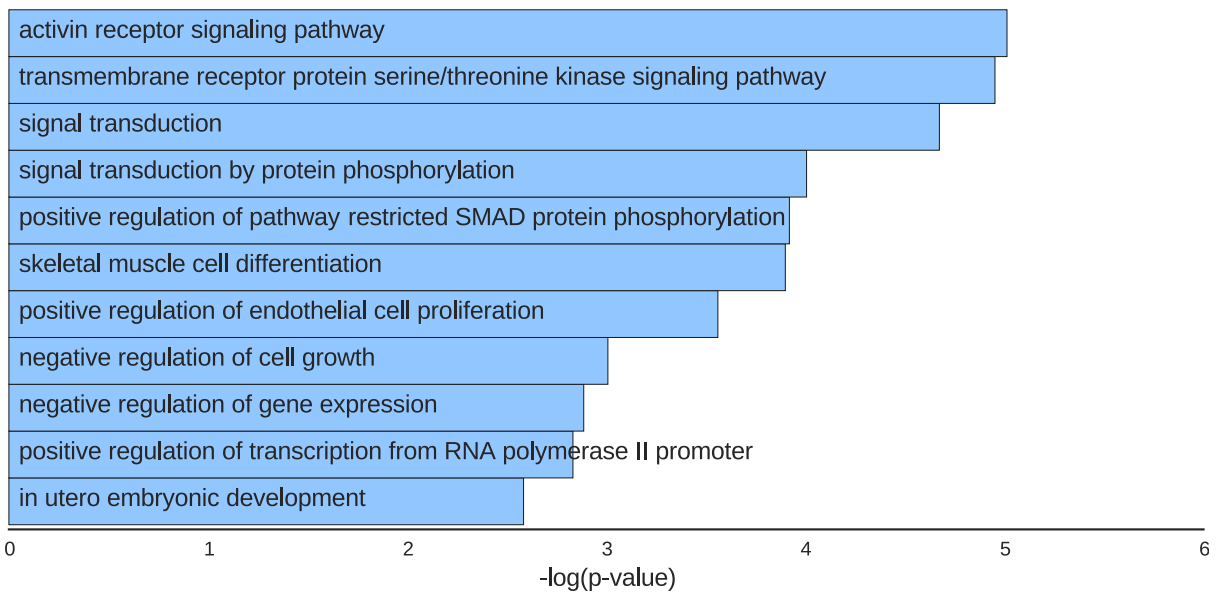
Active cluster #69



Active cluster #44



Active cluster #47



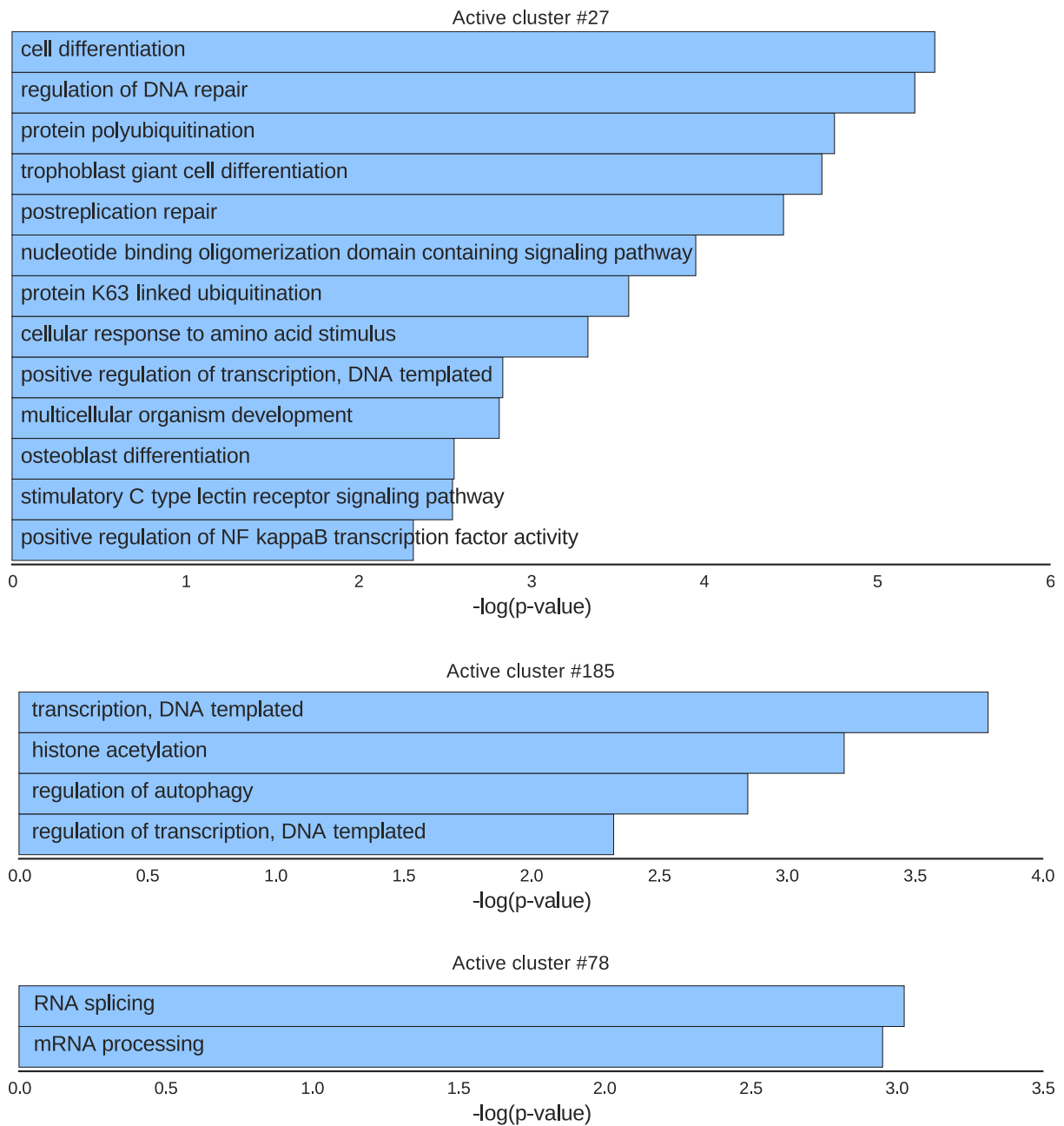
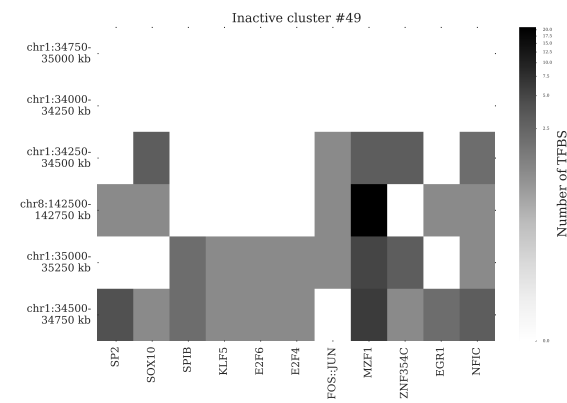
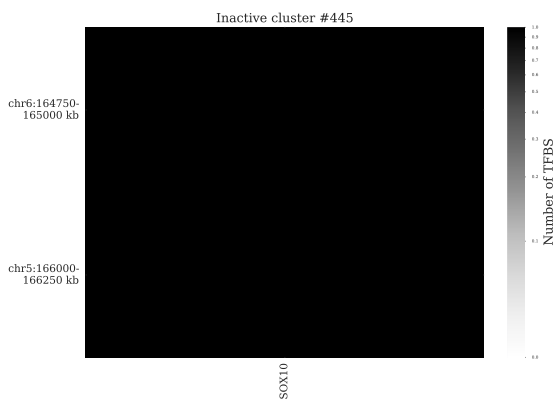
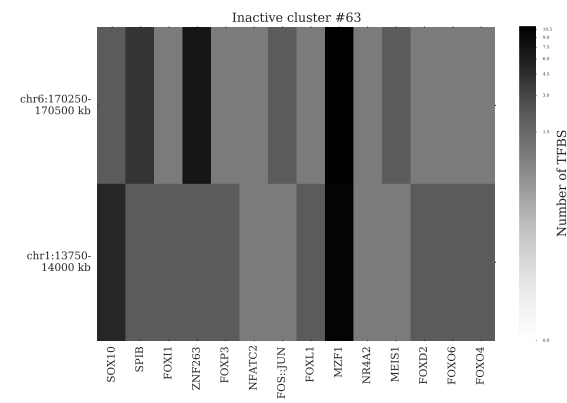
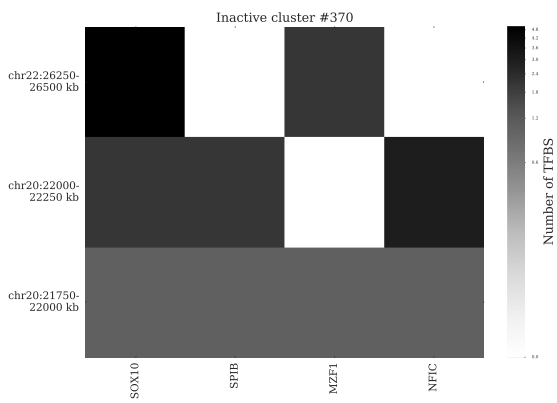
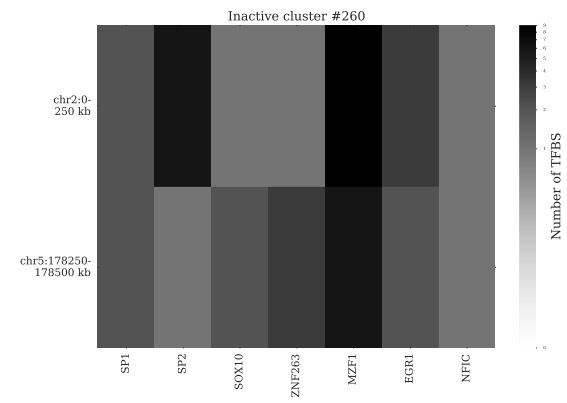
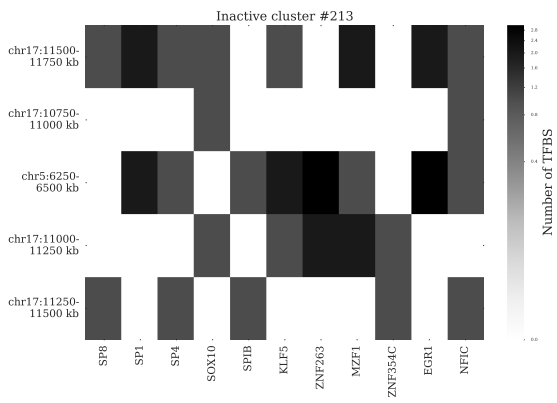
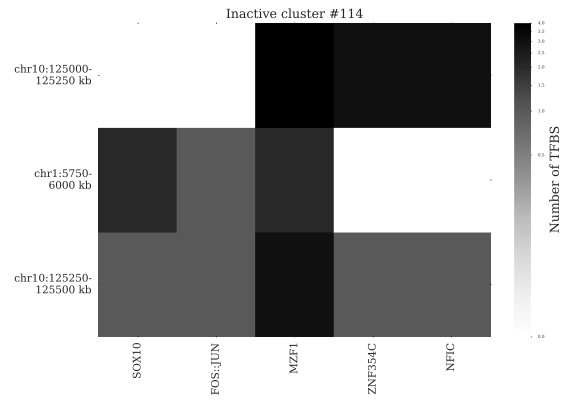
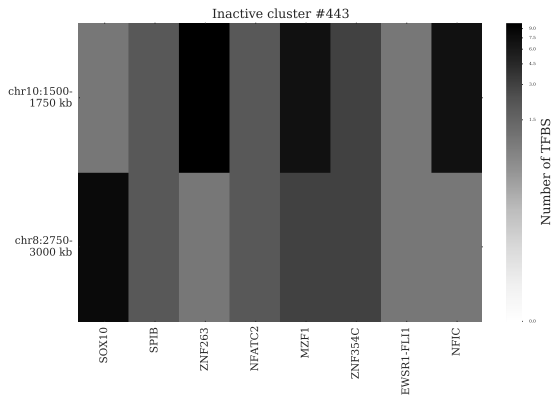


Figure S5: GO terms associated with top active clusters in Fig. S7. Only clusters that have less than 10 chromosomes in them and that had GO terms associated with the cluster are displayed.



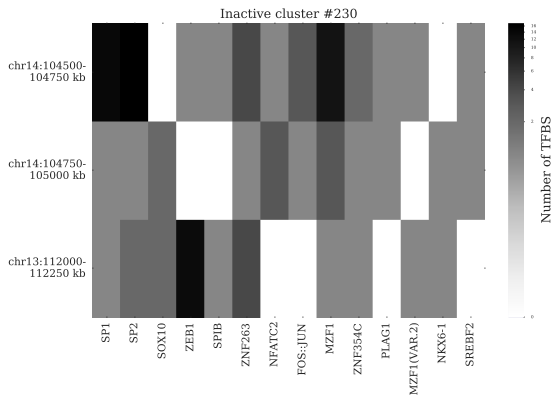


Figure S6: TFBS for top 15 inactive clusters given in Table S5 (with TFBS based on JASPAR 2016, threshold = 0.000001, CAGE). Only TFs that span at least 2 regions in the cluster are shown. Clusters with zero TFs spanning at least 2 regions are not shown.

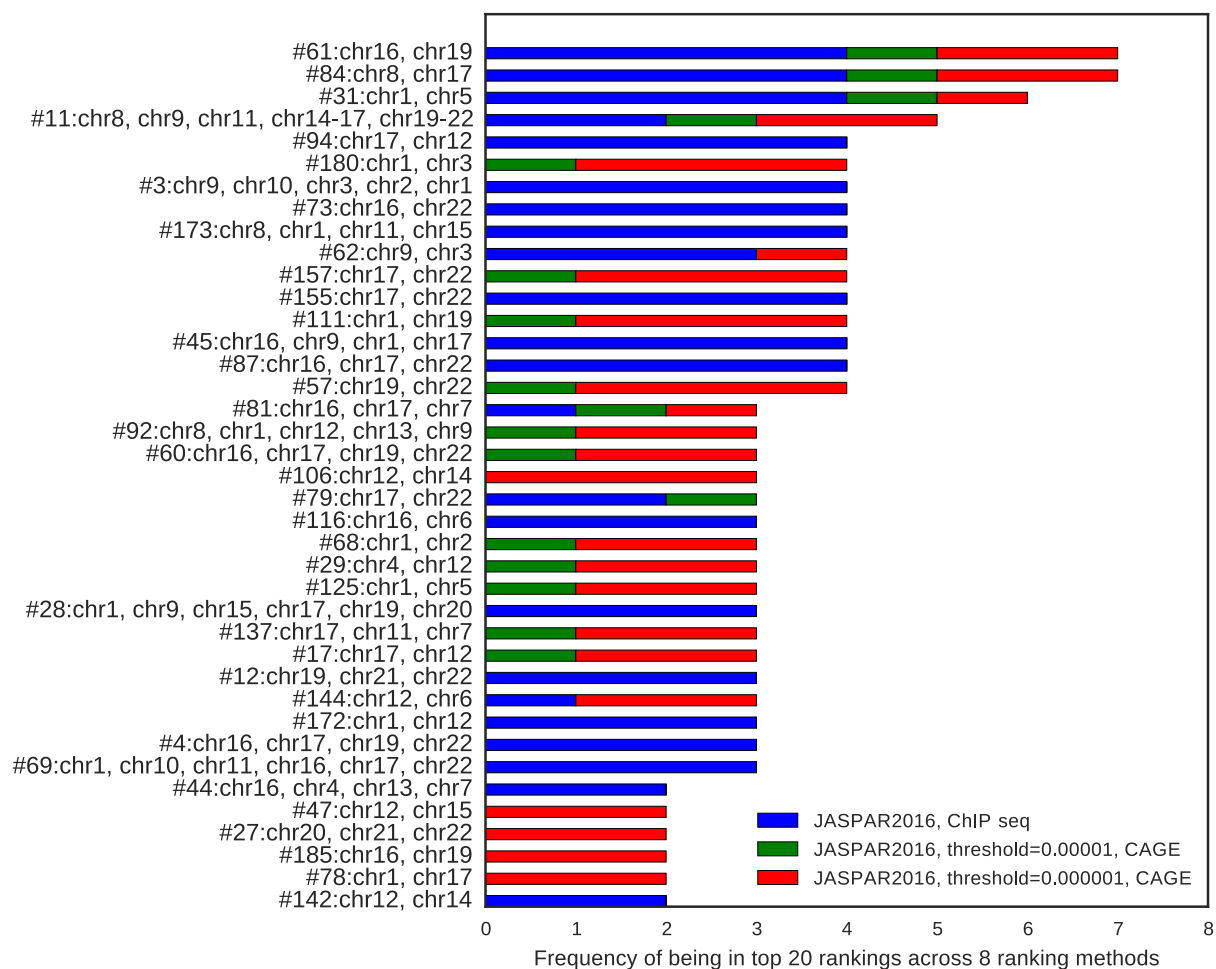


Figure S7: Top ranked clusters that appeared in the top 20 clusters across 8 different methods of evaluation. The eight rankings are grouped by the different choices for filtering the JASPAR 2016 database: The JASPAR2016 database was filtered by ChIP-seq (blue) and the rankings were obtained i) using active clusters as background and by generating random matrices from the observed counts for the permutation test, ii) using active clusters as background and by generating random matrices based on the Dirichlet distribution for the permutation test, iii) using all intermingling regions as background and by generating random matrices from the observed counts for the permutation test, iv) using the whole genome as background and by generating random matrices from the observed counts for the permutation test. The JASPAR2016 database TFBS were obtained with a threshold of 0.00001 and filtered by CAGE (green) and the rankings were obtained v) using active clusters as background and by generating random matrices from the observed counts for the permutation test. The JASPAR2016 database TFBS were obtained with a threshold of 0.000001 and filtered by CAGE (red) and the rankings were obtained vi) using active clusters as background and by generating random matrices from the observed counts for the permutation test, vii) using active clusters as background and by generating random matrices from the Dirichlet distribution for the permutation test, and viii) using intermingling regions as background and by generating random matrices from the observed counts for the permutation test.

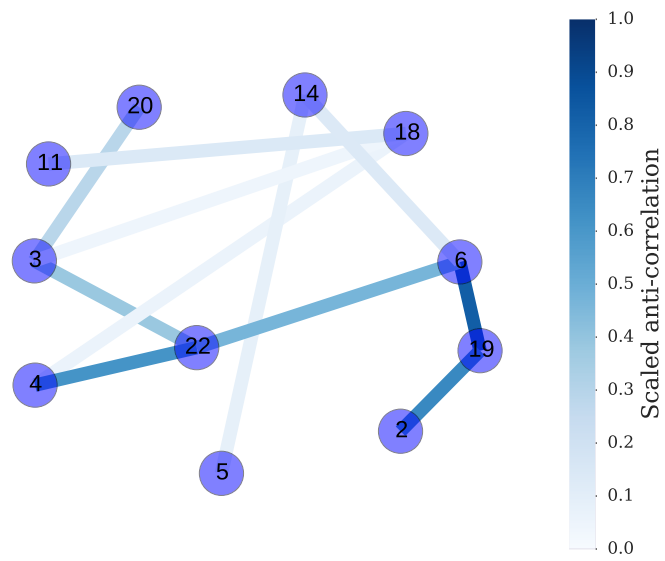


Figure S8: Determining chromosomes that do not intermingle for predicting negative controls. Edges link nodes, i.e., pairs of chromosomes, that showed no intermingling regions in the LAS analysis. The edge weights are given by the absolute value of the anti-correlation between the genomic features of the adjacent chromosomes at a whole chromosome level.

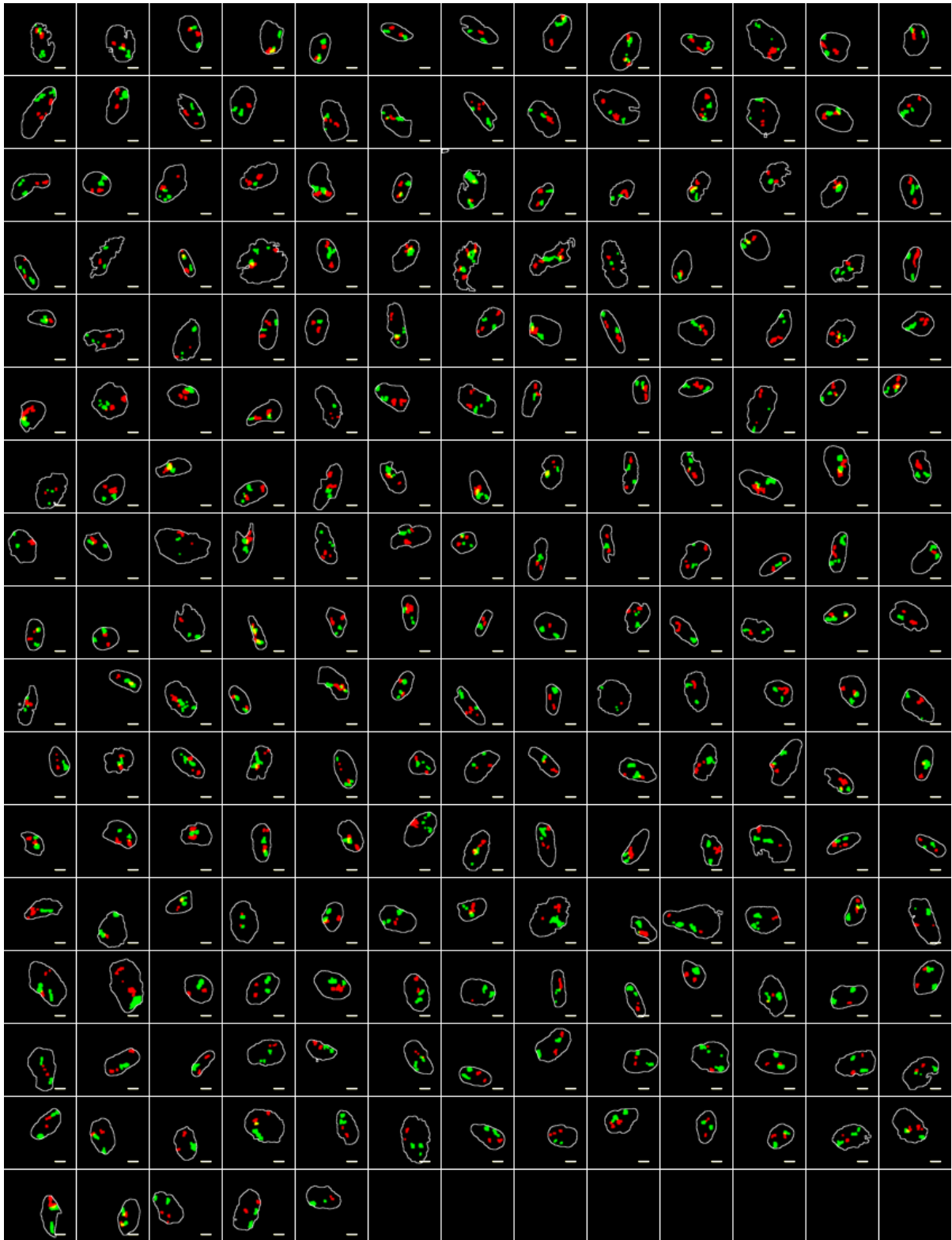


Figure S9: Experimental validation using FISH for predicted active cluster on chromosomes 12 and 17. Segmented images of nuclei from a population of cells with nuclear boundary shown in white, chromosome 17 in red, and chromosome 12 in green. The scale bar has a length of 5 μm .

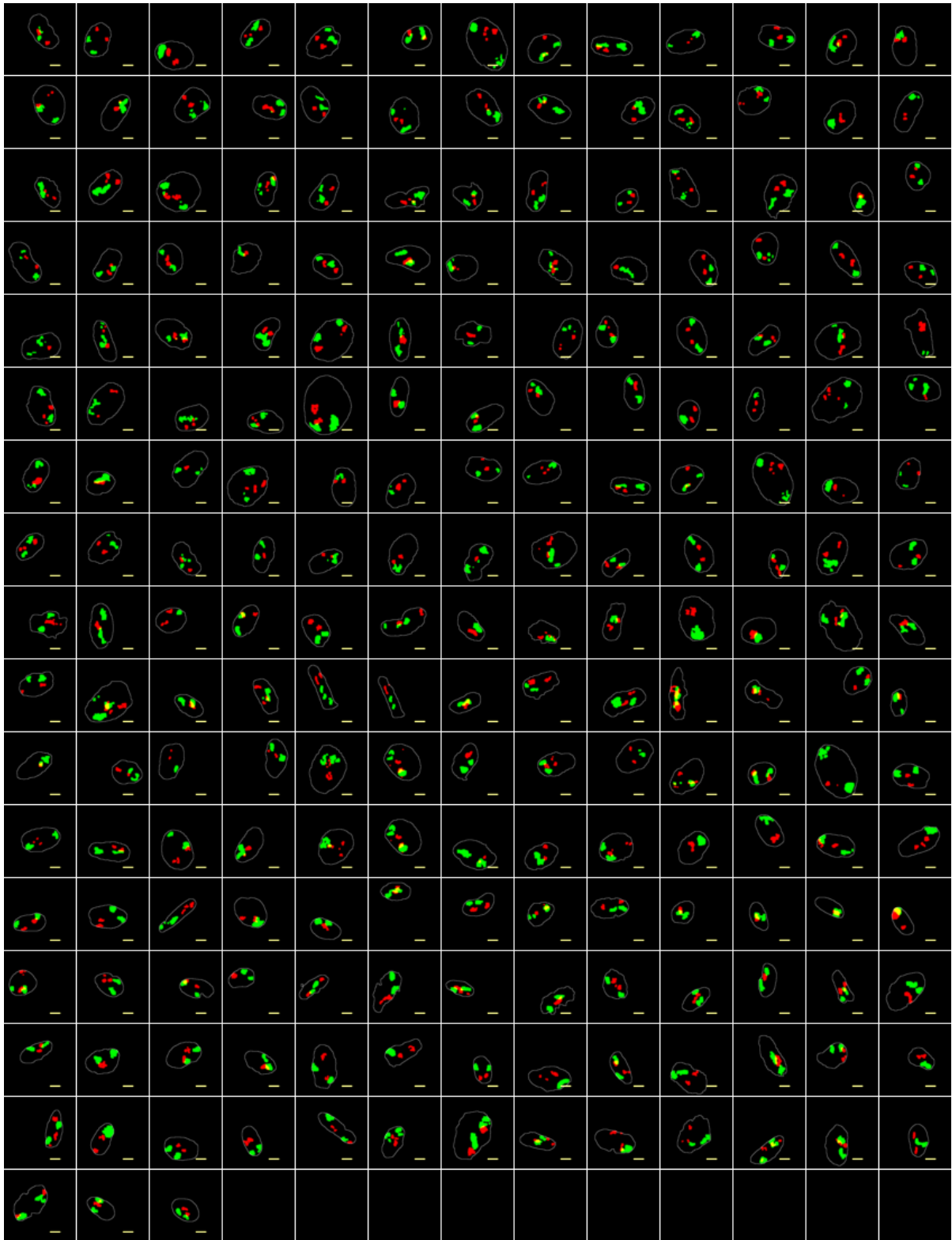


Figure S10: Experimental validation using FISH for predicted negative control, i.e. chromosomes 3 and 20 that are predicted to not intermingle. Segmented images of nuclei from a population of cells with nuclear boundary shown in white, chromosome 20 in red, and chromosome 3 in green. The scale bar has a length of 5 μm .

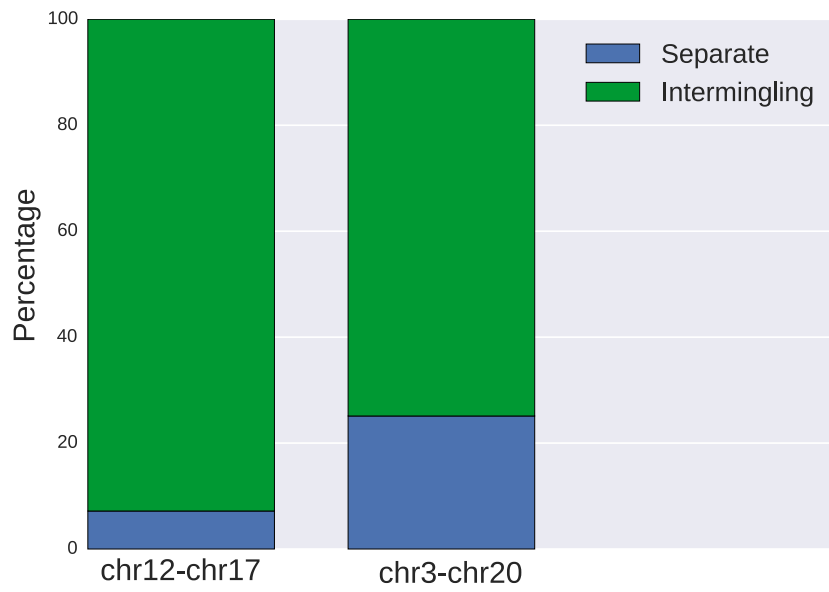


Figure S11: Breakdown of the percentage of nuclei that intermingle (intermingling degree > 0) versus don't intermingle for the chromosome pairs 12-17 and 3-20, as found by FISH experiments.

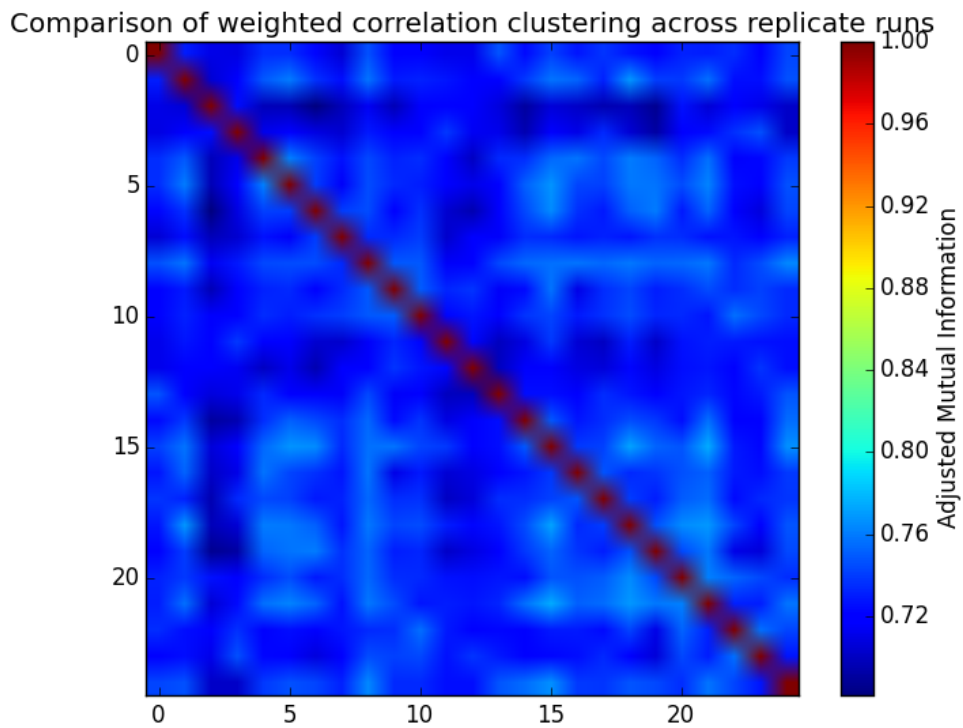


Figure S12: Adjusted mutual information between replicate clusterings from weighted correlation clustering.