Supplemental Figures for:

## Modular Combinatorial Binding among Human *Trans*-acting Factors Reveals Direct and Indirect Factor Binding

Yuchun Guo[1], David K Gifford[1]

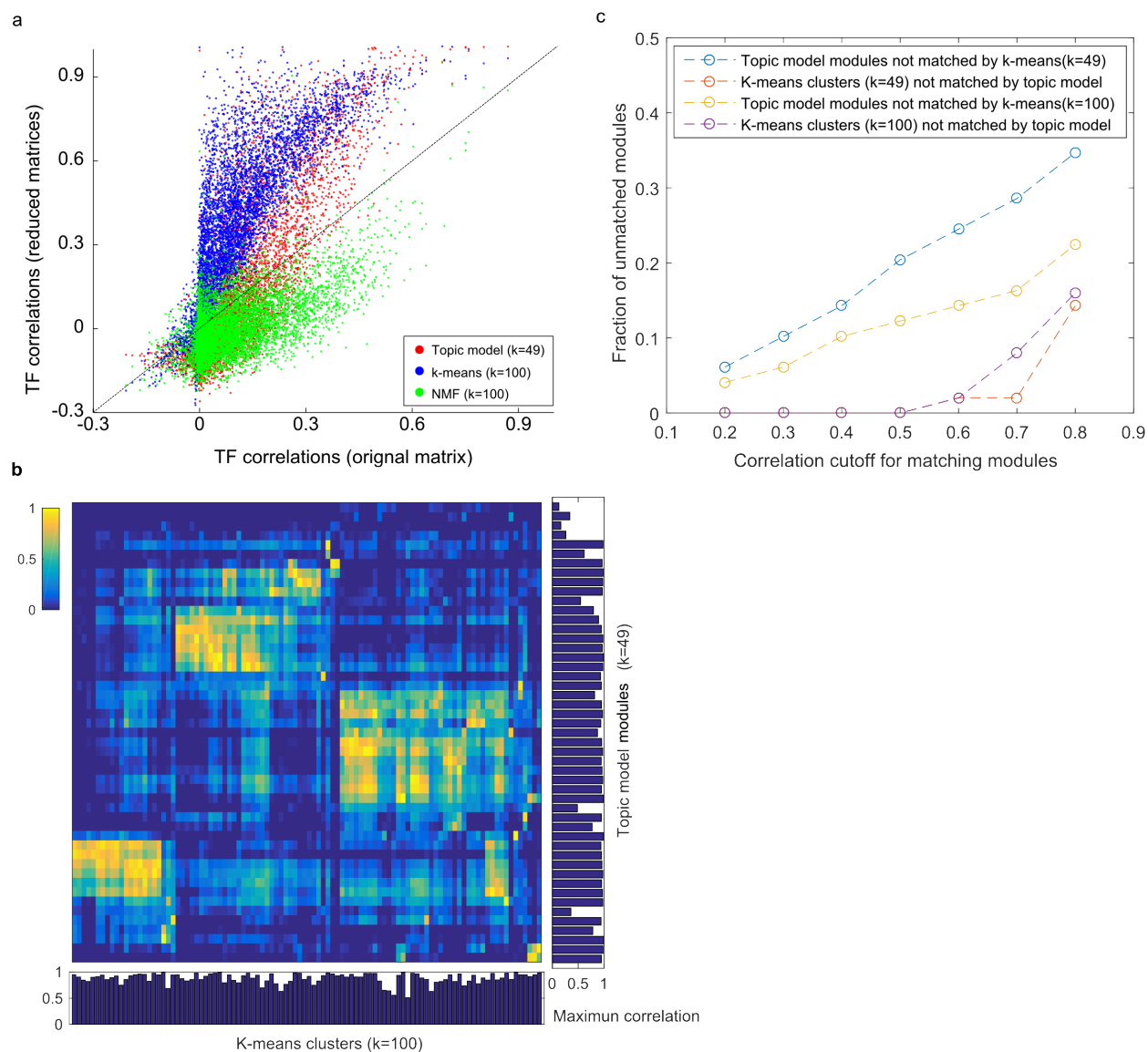[1]MIT, Computer Science and Artificial Intelligence Laboratory, Cambridge, MA 02139, USA

**Figure S1. RMD is able to discover modules that are not discovered by k-means clustering.**

**a** The topic model re-capitulates the original binding data more accurately than k-means clustering and NMF. Each point in the scatter plot represents the Pearson correlation coefficient between a pair of TFs calculated using the original binding data (x-axis) or calculated using the reduced data matrix (y-axis) by topic model (k=49), k-means clustering (k=100), or NMF (k=100). **b** A heatmap shows the Pearson correlations between topic model regulatory modules (k=49) and k-means clusters (k=100). The modules are ordered as in Fig. 2. Bottom bar chart shows the maximum correlation values for each k-means cluster. All k-means clusters are matched by at least one modules, with a maximum correlation value larger than 0.5. Right bar chart shows the maximum correlation values for each module. Six modules cannot be matched by any k-means clusters. **c** The fractions of unmatched modules at various cutoffs for correlation values. A larger fraction of topic model modules are not matched by the k-means clusters (k=49 or k-100) than the fraction of the unmatched k-means clusters.

*Trans-acting Factors*
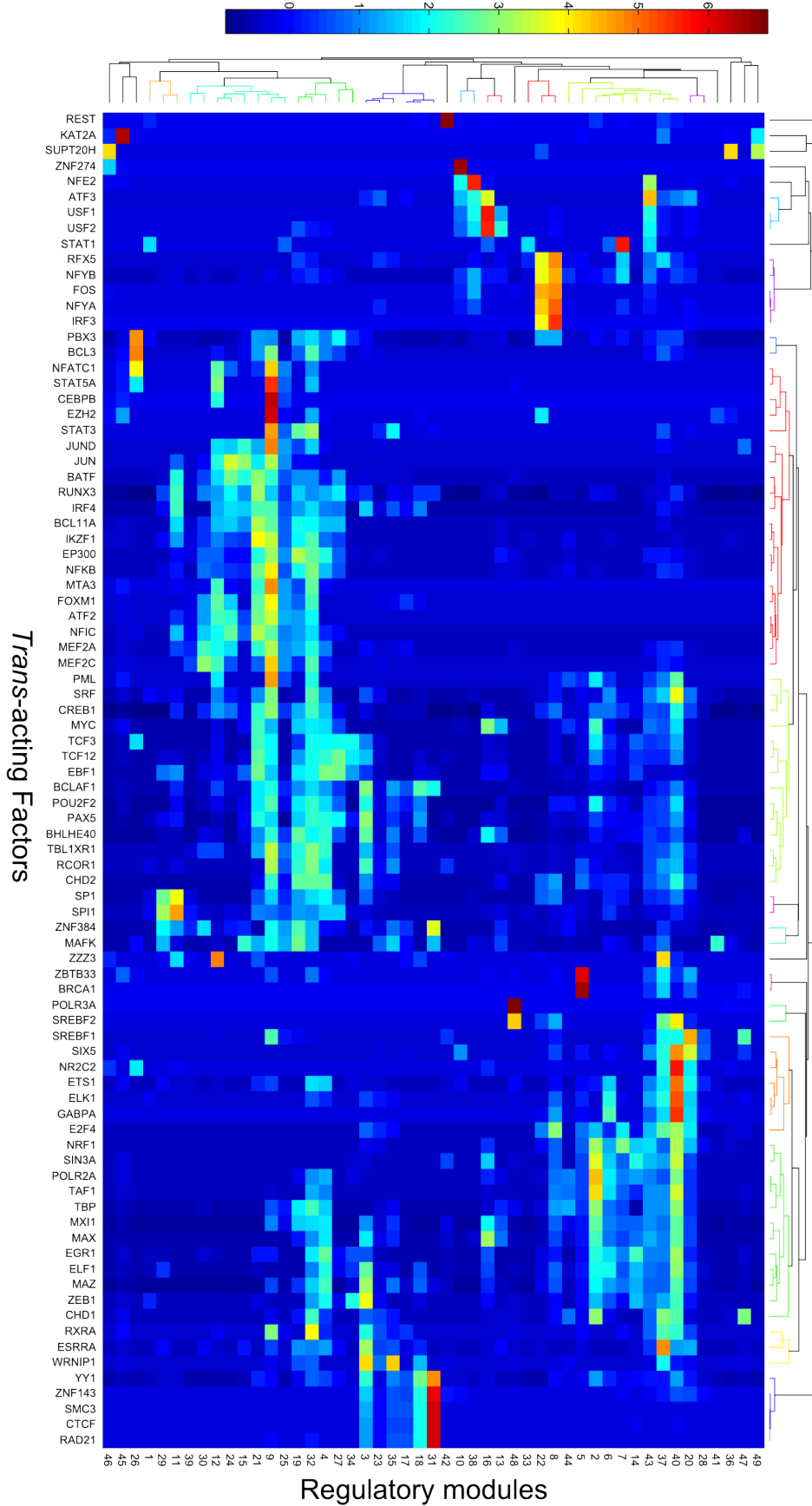
Regulatory modules

## Figure S2. RMD discovered 49 regulatory modules using GM12878 TF binding data

RMD was applied to a compendium of ChIP-seq binding sites of 86 TFs in human GM12878 cells and discovered 49 regulatory modules. Each cell in the heatmap represents the z-score of the TF binding site count (standardized along the columns) of a TF (column) in a regulatory module (row). The top and left dendrograms were computed by applying hierarchical clustering on the regulatory module matrix with Pearson correlation distance and average linkage.
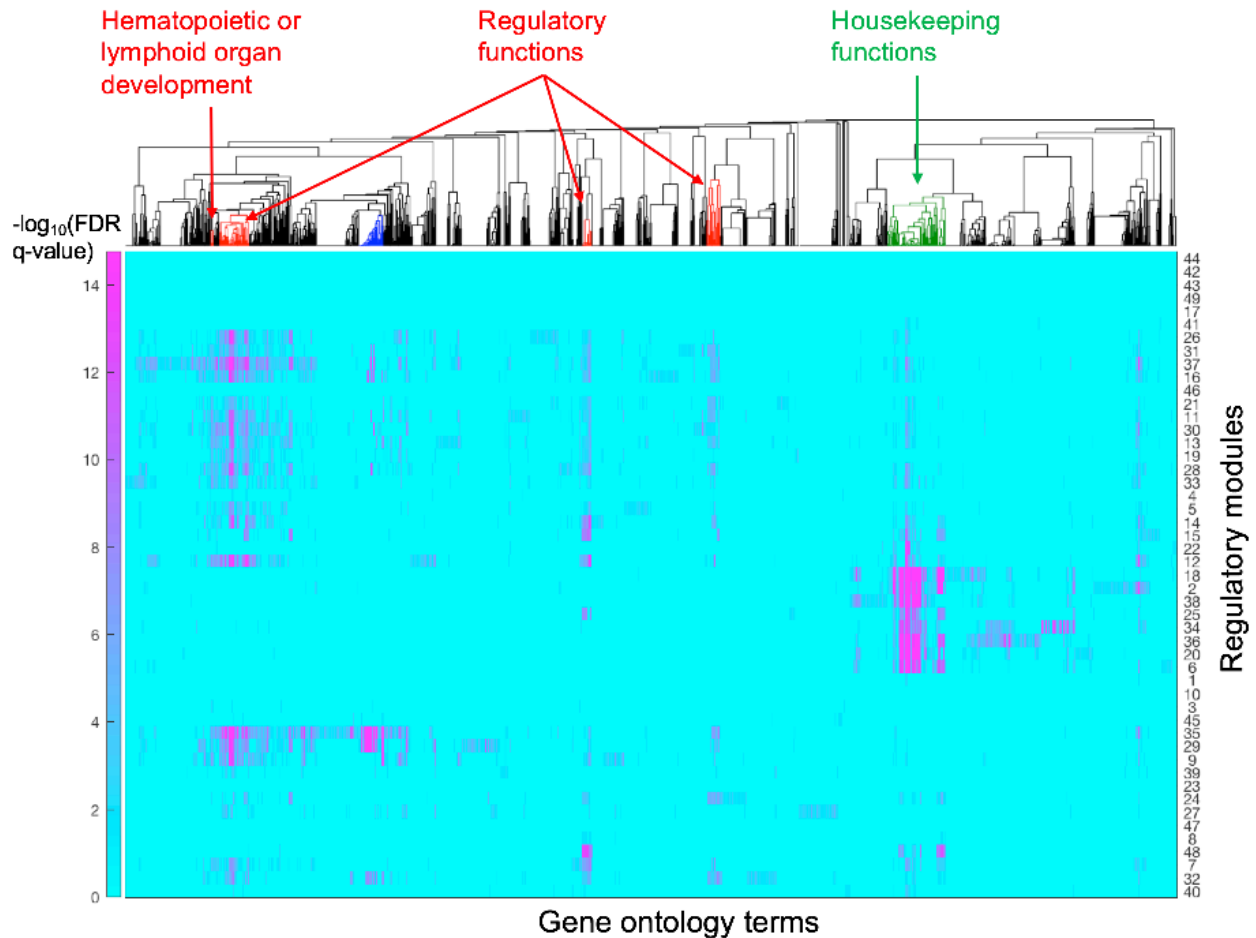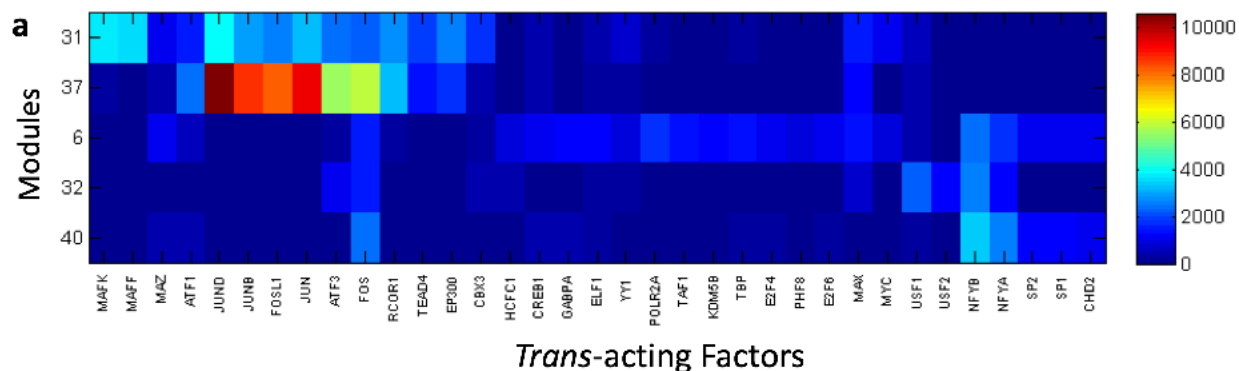
**Figure S3. The GO term enrichment in regulatory modules.**
The enriched GO terms are obtained from GREAT analysis on the regions that use the regulatory modules. The top dendrogram was computed by applying hierarchical clustering on the -log$_{10}$(FDR q-value) of GO terms with Pearson correlation distance and average linkage. The regulatory modules (rows) are ordered as in Figure 2.

**a** A subset of the regulatory module matrix that involves FOS binding are shown. Each cell in the heatmap represents the TF binding site count of a TF (column) in a regulatory module (row).

| Modules | Xie et. al categories |
|---------|----------------------|
| 31 | EP300 mediated distal |
| 37 | Canonical AP1 |
| 6 | Proximal - HOT |
| 32 | FOS+NFYB |
| 40 | FOS+NFYB |
| 37 + 6 | AP1 - HOT |

**Figure S4. The regulatory modules reveal context-dependent FOS co-binding as reported in previous work.**

**a** A subset of the regulatory module matrix that involves FOS binding are shown. Each cell in the heatmap represents the TF binding site count of a TF (column) in a regulatory module (row). **b** The modules recapitulate five categories of FOS co-localization patterns reported previously (Xie et al., 2013).
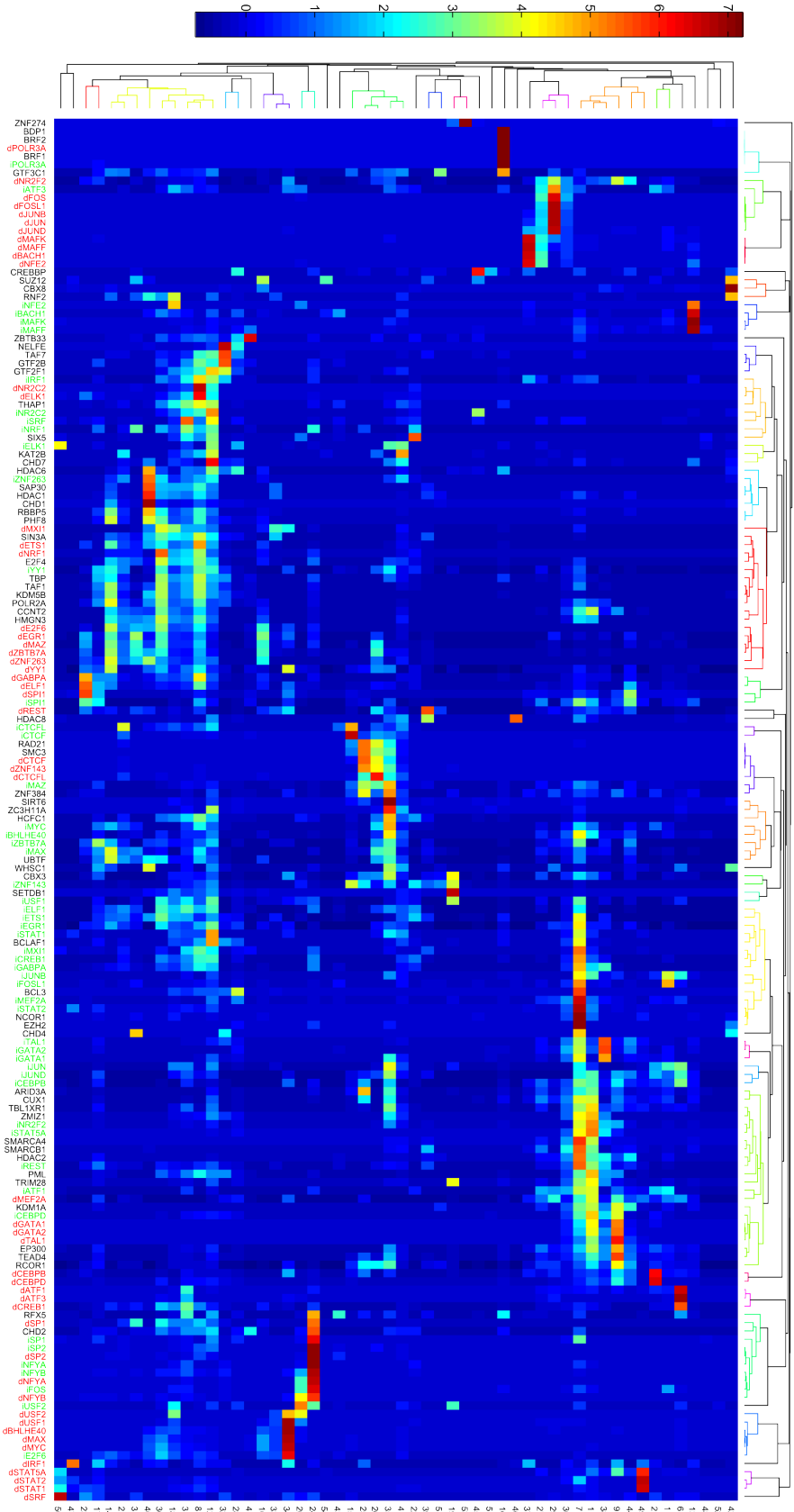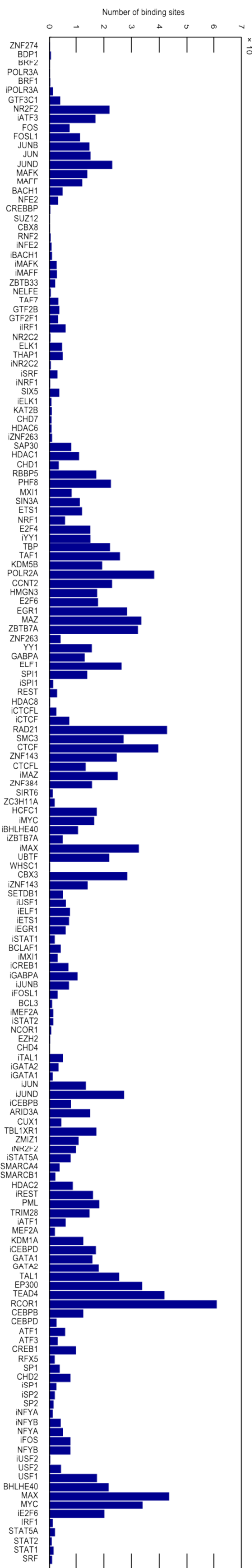
## Figure S5. The direct/indirect binding regulatory modules in K562 cells

RMD was applied to a compendium of ChIP-seq binding sites of 167 direct and indirect binding "factors" in human K562 cells and discovered 54 regulatory modules. Each cell in the heatmap represents the z-score of the TF binding site count (standardized along the columns) of a TF (column) in a regulatory module (row). The factor names are colored according to their binding types: direct binding factors dTF (red), indirect binding factors iTF (green) and factors with unknown binding type (black). The top and left dendrograms were computed by applying hierarchical clustering on the regulatory module matrix with Pearson correlation distance and average linkage. The bottom bar plot shows the total number of binding sites of the "factors".
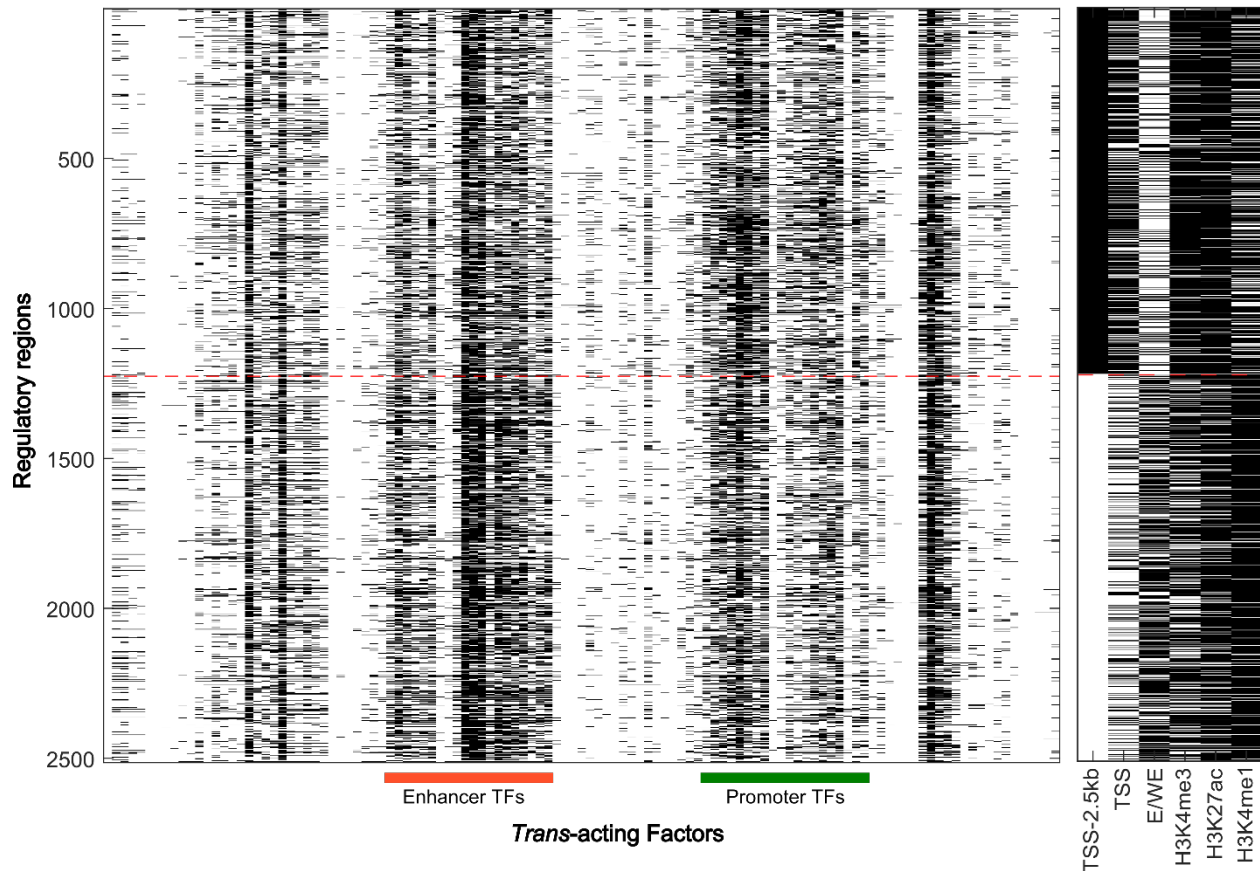
**Figure S6. Co-occurrence of enhancer and promoter modules.**

Left panel: A heatmap showing the region-TF binding matrix of the 2514 regions co-bound by enhancer and promoter modules. There are 30477 regions bound by enhancer modules and 27349 regions bound by promoter modules (not shown). The regions are divided into those that are within 2.5kb of the TSSs and those that are not within 2.5kb of the TSSs. The TFs are in the same order as in Figure 2. Right panel: A heatmap showing the chromatin state and histone mark annotations of the same regions. E/WE: enhancer/weak enhancer chromatin state; TSS: TSS/promoter chromatin state.