1

# Identifying Transcriptional Regulatory Modules among Different Chromatin States in Mouse Neural Stem Cells

**Sharmi Banerjee** [1,2]**, Hongxiao Zhu** [3]**, Man Tang** [3]**, Wu-chun Feng** [4]**, Xiaowei Wu** [3]**, Hehuang Xie**[2,5,6,7,*]

[1]*Bradley Department of Electrical and Computer Engineering,Virginia Tech, Blacksburg, 24061, USA*
[2]*Biocomplexity Institute of Virginia Tech, Blacksburg, 24061, USA*
[3]*Department of Statistics,Virginia Tech, Blacksburg, 24061, USA*
[4]*Department of Computer Science,Virginia Tech, Blacksburg, 24061, USA*
[5]*Department of Biomedical Sciences and Pathobiology, Virginia-Maryland College of Veterinary Medicine,Blacksburg, 24061, USA*
[6]*Department of Biological Sciences,Virginia Tech, Blacksburg, 24061, USA*
[7]*School of Neuroscience,Virginia Tech, Blacksburg, 24061, USA*

Correspondence*:
Hehuang Xie, Biocomplexity Institute of Virginia Tech, Department of Biomedical Sciences and Pathobiology, Virginia-Maryland College of Veterinary Medicine, Department of Biological Sciences, School of Neuroscience,Virginia Tech, Blacksburg, 24060, USA
davidxie@vt.edu

2 **ABSTRACT**

3   Gene expression regulation is a complex process involving the interplay between transcription
4 factors and chromatin states. Significant progress has been made towards understanding the
5 impact of chromatin states on gene expression. Nevertheless, the mechanism of transcription
6 factors binding combinatorially in different chromatin states to enable selective regulation of
7 gene expression remains an interesting research area. We introduce a nonparametric Bayesian
8 clustering method for inhomogeneous Poisson processes to detect heterogeneous binding
9 patterns of multiple proteins including transcription factors to form regulatory modules in different
10 chromatin states. We applied this approach on ChIP-seq data for mouse neural stem cells
11 containing 21 proteins and observed different groups or modules of proteins clustered within
12 different chromatin states. These chromatin-state-specific regulatory modules were found to have
13 significant influence on gene expression. We also observed different motif preferences for certain
14 TFs between different chromatin states. Our results reveal a degree of interdependency between
15 chromatin states and combinatorial binding of proteins in the complex transcriptional regulatory
16 process. The software package is available on Github at - https://github.com/BSharmi/DPM-LGCP.

17 **Keywords: transcription factor, regulatory network, Poisson process, chromatin states, neural stem cell**

# 1 INTRODUCTION

Transcription factors (TFs) and other proteins that bind to specific DNA sequences play key roles in the regulation of gene expression. Binding locations of a protein of interest can be determined with chromatin immunoprecipitation followed by sequencing (ChIP-seq). This produces millions of short reads covering the protein-DNA binding sites across the genome. Several computational tools have been developed to identify these binding locations from ChIP-seq data. Widely used among these is MACS2 (Feng et al., 2012) which can identify transcription factor binding regions or 'peaks'. Recently, efforts have been devoted to integrate multiple ChIP-seq datasets to uncover protein-protein interactions. SignalSpider (Wong et al., 2015) uses Gaussian mixture model to reveal regions co-regulated by multiple TFs. Sharmin et. al. identified cell-type specific TF binding events (Sharmin et al., 2016) using ensemble model. Cha and Zhou developed a method based on inhomogeneous Poisson processes and Ripley's K-function that detects pairwise TF clustering and ordering patterns (Cha and Zhou, 2014).

Recent studies have also revealed new insights into the interplay between proteins, specifically TFs and histone marks that define chromatin states. Most TFs bind to open chromatin regions that are highly accessible and nucleosome-depleted. Such chromatin regions are often enriched with specific histone modifications in promoters and enhancers, such as H3K4me1 and H3K27ac marks. It has been found that histone-modification-dependent TF binding is protein family specific (Xin and Rohs, 2018; Sugathan and Waxman, 2013; Liu et al., 2015, 2016). In addition, a small number of TFs act as pioneers with the ability to reach inaccessible chromatin regions and shape the chromatin landscape to facilitate the binding of other TFs. ChIP-seq data from histone modifications have been used to partition the genome into different chromatin states using semi-automated genome annotation (SAGA) tools (Libbrecht et al., 2015). Early examples of the SAGA tools are HMMSeg (Day et al., 2007) and ChromHMM (Ernst and Kellis, 2012). Since then more sophisticated chromatin segmentation tools, Segway (Hoffman et al., 2012) and diHMM (Marco et al., 2017), were developed providing refined genome-wide map of the chromatin states. ChromHMM and diHMM use hidden Markov models while Segway applies a dynamic Bayesian network to segment the genome and identify distinct chromatin states. Segway and ChromHMM perform genome segmentation and classification at a single length scale while diHMM segments the genome at multiple length scales (narrow or broad corresponding to nucleosome-level states and domain-level states respectively). We studied protein bindings through ChIP-seq data among different chromatin states in mice neural stem cells (detailed description of datasets provided in Supplementary document section 3.1). Our results showed several known co-binding rules such as NFIC-bHLH-SOX in Upstream Enhancer state and Poised Enhancer state (Mateo et al., 2015) and JMJD3-SMAD3 in all chromatin states (Estarás et al., 2012). We also showed that the regulatory effects of the predicted modules on proximal genes vary across chromatin states. Also, for certain classes of DNA binding proteins, the de-novo binding sequences compiled from ChIP-seq peaks were dependent on the chromatin states.

# 2 MATERIALS AND METHODS

In this paper we propose a two-step process (Figure 1) to investigate how chromatin configurations may affect the binding affinity of proteins. In the first step, uniquely aligned BAM files containing genomic regions of histone marks and TFs are used along with the diHMM software to segment the genome and identify distinct chromatin states (illustrated by chromatin state examples X and Y). In the second step, using the identified chromatin states from the previous step and protein binding regions obtained from ChIP-seq (data used in this study were obtained from ChIP-Atlas (http://chip-atlas.org)), a nonparametric Bayesian clustering method DPM-LGCP is applied to identify transcriptional regulatory modules within

59  each chromatin state. In downstream analyses, proximal (+/- 2kb from transcription start site genes are used
60  to compare the Transcripts Per Kilobase Million or TPM expression level when regulated by individual
61  proteins to that when regulated combinatorially by the predicted regulatory modules in step 2. Finally, using
62  de-novo motif enrichment analysis, the binding sequences of the proteins are compared across different
63  chromatin states to study the effect of histone marks and co-factors on motif preferences. Details of the
64  datasets used in the study can be found in Supplementary Table S2.
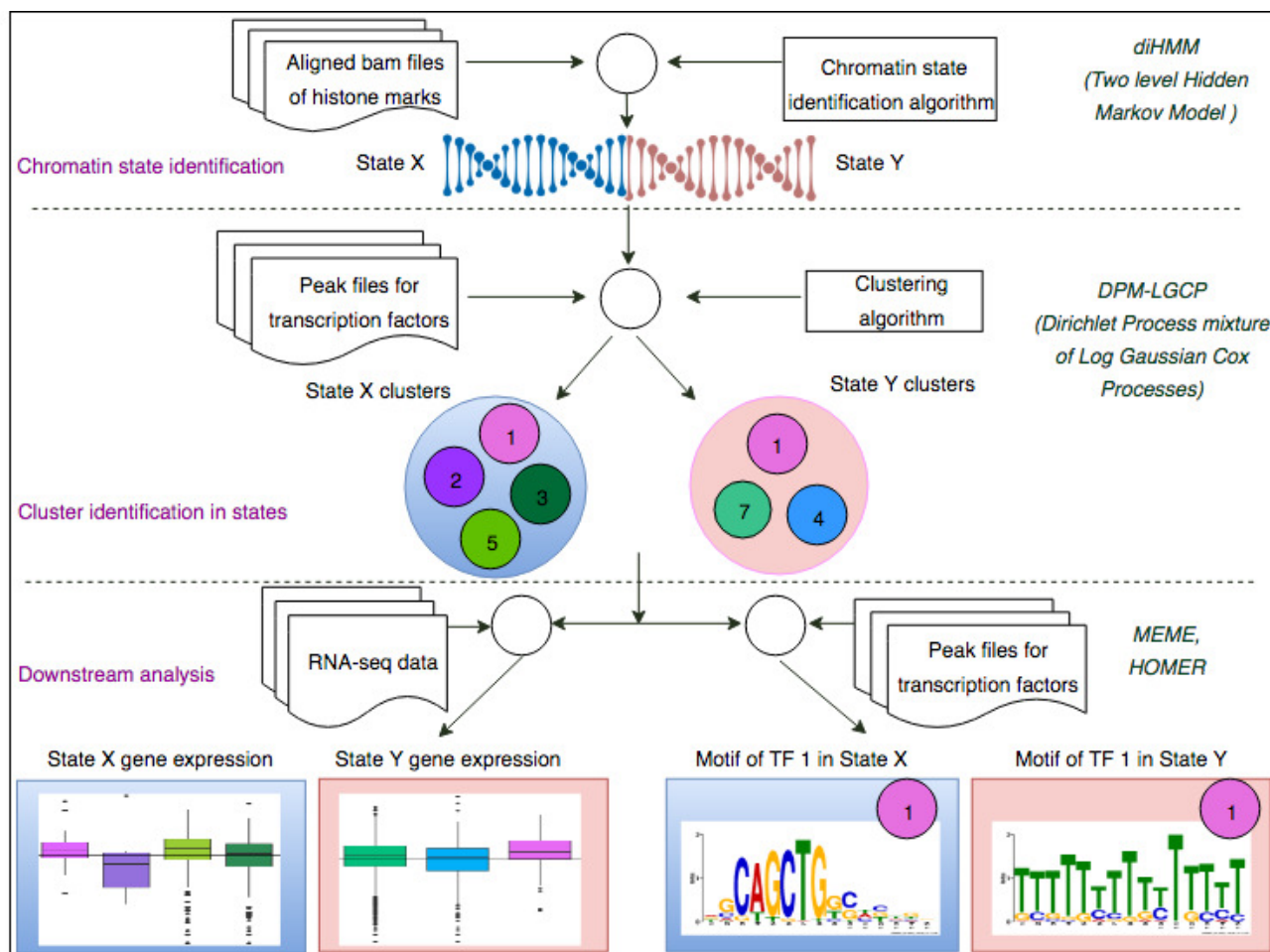


**Figure 1.** A two-step process to identify chromatin-state-specific transcriptional regulatory modules. In the first step, uniquely aligned bam files of histone marks are used along with the diHMM software to segment the genome and identify distinct chromatin states (illustrated by State X and State Y). In the second step, using the identified chromatin states from the previous step and ChIP-seq peak files for different TFs, the proposed Bayesian clustering method is applied to identify transcriptional regulatory modules within each chromatin state. In downstream analyses, proximal (+/- 2kb from TSS) genes are used to compare the TPM expression level when regulated by individual TFs to that when regulated combinatorially by the predicted regulatory modules in step 2. Finally, using de-novo motif enrichment analysis, the binding sequences of the TFs are compared across different chromatin stats to study the effect of histone marks and co-factors on TF binding sequences.

## 2.1  Chromatin state identification through genome segmentation

66  diHMM (Marco et al., 2017) is a tool based on hidden Markov model that models the presence or
67  absence of a histone mark to a high degree of accuracy. It segments and annotates the genome into different
68  chromatin states at multiple length scales by modeling the genome wide distribution of histone marks. By

69 default, diHMM has two scales of classification: (a) nucleosome level, with finer resolution chromatin state
70 windows of around 200 base-pair (bp) length and (b) domain level, formed by stitching together similar
71 nucleosome-level windows and having broader chromatin state windows extending over 100kbp-long
72 regions. The domain-level states identified by diHMM are able to recapitulate known patterns in the
73 chromatin literature and capture functional differences among diverse regulatory elements (Marco et al.,
74 2017). The first step in identifying chromatin states is to binarize uniquely aligned BAM files. This is
75 implemented in ChromHMM (Ernst and Kellis, 2012), a predecessor of diHMM. The diHMM software
76 provides several nucleosome- and domain-level statistics including nucleosome-level emissions, combined
77 nucleosome-level fold enrichments for each domain, fractional genome coverage of each nucleosome- and
78 domain-level state, and nucleosome and domain state lengths. These statistics, together with the relative
79 distance information of nucleosome- and domain-level states from transcription start site (TSS) and the
80 enrichment of nucleosome-level states in genomic regions, were jointly analyzed to annotate each state to a
81 biologically relevant functional category (details provided in RESULTS section).

## 2.2 Protein binding intensity estimation using Dirichlet Process Mixture of Log Gaussian Cox Processes (DPM-LGCP)

Binding regions of the proteins were obtained using MACS2 acting as inputs to our proposed clustering algorithm. Treating the center of each region as a binary binding event, we modeled binding events of each protein along the genome by an inhomogeneous Poisson process (*IP*). We chose this modeling strategy for the following reasons: (i) the event of each binding site falling into a minuscule interval is a rare event, independent of the events in other non-overlapping intervals, and (ii) the non-uniform distribution of the peaks at different genomic locations can be well characterized by the intensity function of the inhomogeneous Poisson process. For a protein with $n$ binding site locations, we map these locations to points in a closed interval $D$ on the real line, denoted by $S = \{s_1, \ldots, s_n\}$. Following the inhomogeneous Poisson process model setting, the likelihood of observing $S$ can be written as (Simpson et al., 2016)

$$f(S|\lambda(s)) = \exp\left\{|D| - \int_D \lambda(s)ds\right\} \prod_{j=1}^n \lambda(s_j),\tag{1}$$

where $|D|$ is the interval length and $\lambda(s), s \in D$ is the intensity function. The Poisson process likelihood (1) provides the basis for nonparametric clustering of proteins based on their binding patterns, resulting in identification of modules of co-binding proteins that share similar regulatory functions. For a given ChIP-seq dataset of $N$ proteins coming from $K$ clusters (with $K$ unknown), we assume that proteins in the same cluster share a common intensity function, distinct from those in other clusters.
Under this assumption, we implement a Dirichlet process mixture of log Gaussian Cox process (DPM-LGCP) model that employs a Dirichlet process (DP) prior to the latent log intensity functions to facilitate clustering of the intensity functions. Let $S_i$ denote the binding site locations of the $i$th protein, the DPM-LGCP model can be described as follows:

$$S_i|\lambda_i(s) \sim IP(\lambda_i(s)), s \in D, \quad i = 1, ..., N,$$
$$\log(\lambda_i(s)) = z_i(s), \quad z_i(s) \sim G,\tag{2}$$
$$G \sim DP(m, G_0), \quad G_0 = GP(0, C_\theta),$$

84 where $G$ is a random distribution with a DP prior. The DP prior is characterized by two parameters $m$ and
85 $G_0$, where $m$ is the precision parameter, and $G_0$ is the base measure. The base measure $G_0$ is assumed to

86  be a Gaussian process with mean 0 and covariance kernel $C_\theta(,)$, and $\theta$ contains parameters that control
87  the shape of the covariance kernel. The introduction of this DP prior to the latent log intensity functions
88  naturally facilitates clustering of the $N$ point processes based on their intensity functions. With this model,
89  neither the number of clusters nor ad-hoc distance measure between two point processes needs to be
90  specified.

91      To overcome the difficulty of calculating the marginal likelihood of the point process $S_i$, we employed an
92  approximate but efficient posterior inference using the Integrated Nested Laplace Approximations (INLA)
93  package (Simpson et al., 2016; Rue et al., 2009).

94      The INLA approximation of the LGCP transforms the continuous covariance kernel of $z_i(s)$ into a
95  discrete precision matrix of the B-spline basis coefficients on a regular grid, which enables very fast
96  covariance computation (Rue and Held, 2005). Finally, posterior inference on the assignment of proteins
97  into clusters is performed through a Markov chain Monte Carlo (MCMC) algorithm using Neal's Gibbs
98  sampler (Neal, 2000) (detailed description provided in the Supplementary document).

## 3 RESULTS

### 99  3.1  Genome segmentation and chromatin state identification

100     As described in the methods section, diHMM segments a genome into distinct chromatin states and
101  outputs the states as regions within two bed files labeled by nucleosome and domain indexes (e.g. N1,
102  N2... and D1, D2... respectively). For the nucleosome level states, annotation of the chromatin states
103  to functionally relevant categories was performed by using information from the emission probabilities
104  of the nucleosome states (Figure 2(a)), fractional genome coverage (Figure 2(b)), relative enrichment
105  in different genomic regions (Supplementary Figure S3), and distribution of nucleosome states around
106  TSS (Supplementary Figure S4(A)). Similarly, by comparing the nucleosome-level fold enrichments in
107  each domain level state and the distribution of the domain level states around TSS (Supplementary Figure
108  S4(B)), the domain-level states were further grouped into different broader functional categories as shown
109  in Figure 2(c). Details of functional annotation of the nucleosome and domain-level states are presented in
110  Section 3 of the Supplementary Document.

### 111  3.2  Chromatin state preference of individual protein binding and gene expression
### 112       regulation

113     To analyze the distribution of protein-DNA binding sites in each chromatin state, we integrated ChIP-seq
114  data with the chromatin state map of mouse neural stem cells (NSCs) (Figure 3(a)). For most proteins,
115  the binding events occur in open chromatin regions, although some pioneer transcription factors have the
116  ability to bind directly to condensed chromatin and recruit co-factors (Zaret and Carroll, 2011; Soufi et al.,
117  2015; Cuesta et al., 2007). We observed, in both active and repressed states, enrichment of pioneer TFs as
118  well as other proteins (that might have been recruited by the former). BMI1, which is known to bind to
119  regions marked by both H3K27me3 and H3K4me3 (Bhattacharya et al., 2015), was found to be highly
120  enriched in the Bivalent Promoter and Poised Enhancer states (Figure 3(a)). In addition, most TFs were
121  found to be enriched in the Super Enhancer states except for RAD21, BMI1, SMCHD1 and NUP153. A
122  similar observation was made by the authors in Mateo et al. (2015) where they showed that OLIG2, NFI
123  family, SOX2, SOX9, TCF3, FOXO3, ASCL1, SOX21, and MAX were associated with active enhancer
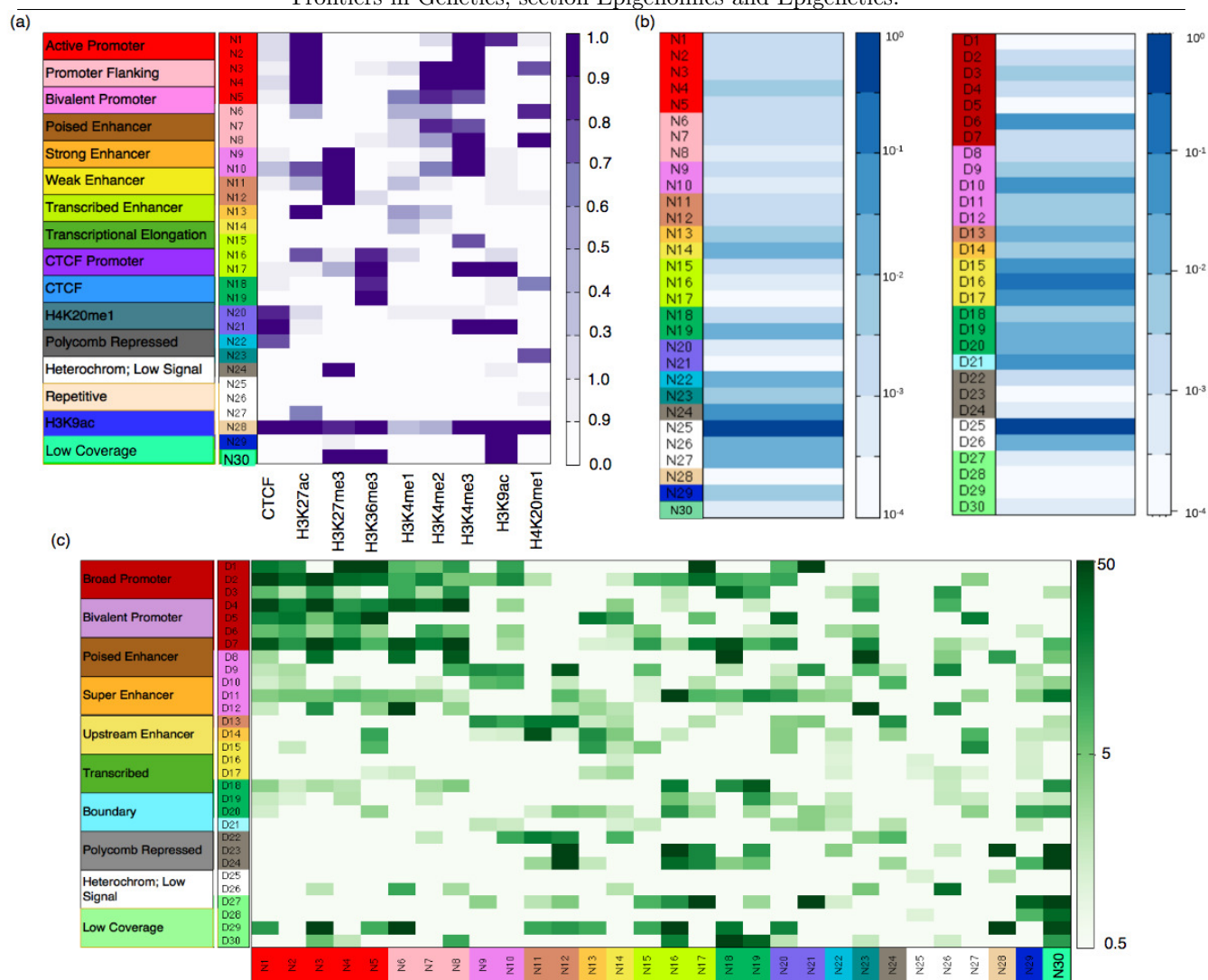124  regions.

**Figure 2.** (a) Nucleosome level emission matrix generated by diHMM. Functional annotations of the nucleosome level states are shown in the color bar on the left. Scale varies linearly between 0 and 1. (b) Fractional genome coverage for nucleosome and domain level states. Scale varies logarithmically between $10^{-4}$ and 1. (c) Combined nucleosome-domain fold change obtained by diHMM. Functional annotation of the states are shown in the color bar on the left. Scale varies logarithmically between 0.5 and 50.

Next, to study the regulatory effect of histone marks on proximal genes, we compared the expression levels of genes (Transcripts Per Kilobase Million or TPM) with promoters located in different chromatin states. We observed that proximal genes in the Broad Promoter state had a higher median expression than proximal genes in the Polycomb Repressed or Low Coverage states (Figure 3(b)). To understand the influence of chromatin states on transcriptional regulation, we further grouped genes in each state based on the presence of binding sites of different proteins surrounding their TSSs. We observed that, for most proteins, the median expression of the genes in active states was higher than those in repressed states (Figure 3(c),(d), Supplementary Figure S8). Also, fewer proteins had binding sites in repressed states as compared to active states (In Figure 3(c), there are 16 proteins whereas in Figure 3(d), there are 14 proteins). Additional gene expression analysis for individual proteins is shown in Supplementary Figure S8.
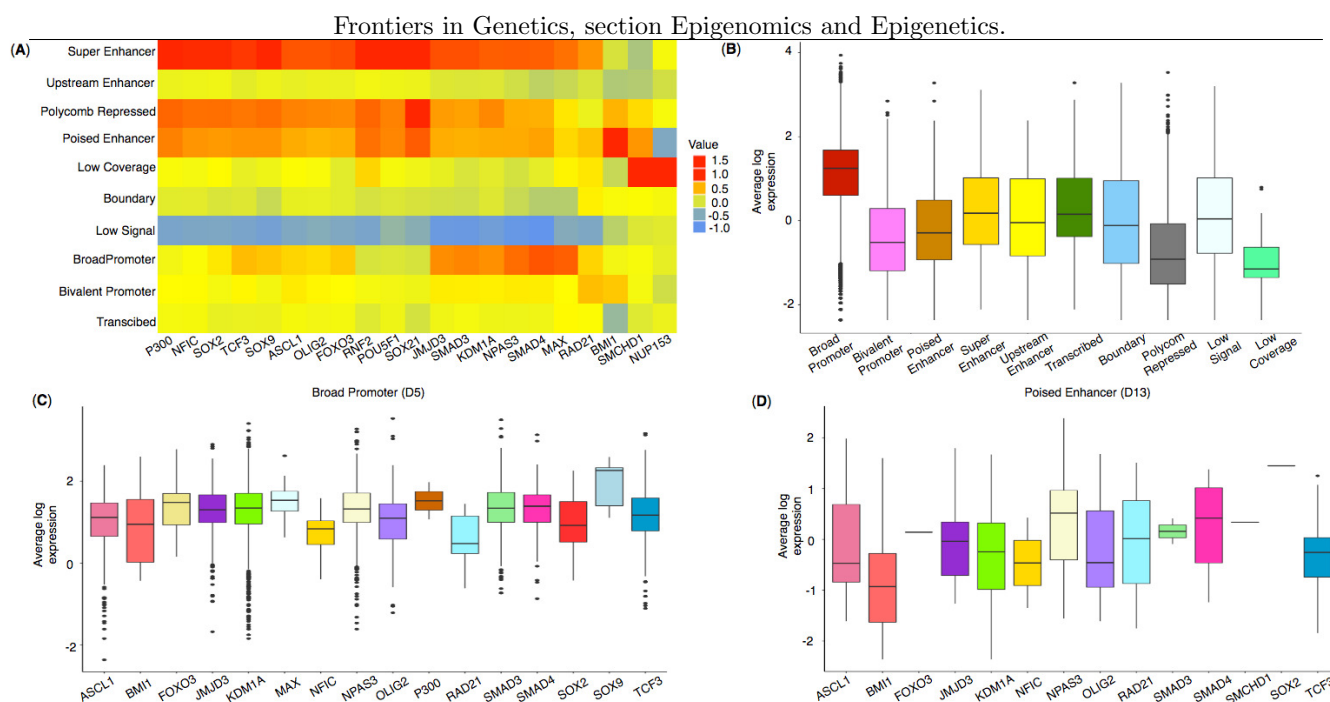
**Figure 3.** (a) Enrichment (in log scale) of TF peaks in different chromatin states showing binding preference of individual TFs. (b) Comparison of average TPM expression (in log scale) of proximal genes (+/- 2kb from TSS) in different domain level chromatin states. Genes were mapped to the nucleosome-level states for the corresponding domain-level states. (c) Comparison of average TPM expression (in log scale) of proximal genes (+/- 2kb from TSS) mapped to individual TFs in the Broad Promoter state and in (d) the Poised Enhancer state.

## 3.3  Chromatin state and preferential clustering of proteins

The distributions of ChIP-seq peaks across distinct chromatin states indicate that functionally relevant proteins may have similar binding patterns (Supplementary Figure S2). We determined the co-occupancy of proteins in a specific chromatin state through a nonparametric Bayesian clustering approach that identifies the combinatorial binding patterns of proteins (detailed description available in Supplementary Document). Each state at the domain level had multiple windows over different chromosomes across the genome. We observed that most windows are with very few peaks although the average domain-level window length ranged from 3.8 kb to over 450 kb. This prevented prediction of modules within a single domain window. To ensure that the unique properties of the domain-level states were preserved during clustering, we merged all windows of a single domain-level state (e.g. D1) across the entire genome and mapped the genome positions to a common interval [0, 50] on an imaginary real line. Adopting this approach for all domain level states eliminated the problem that different domains may have different sizes. Next, for each domain level state, the proposed algorithm used these mapped binding locations, computed individual binding intensity of each protein and clustered proteins having similar intensity patterns together to construct transcriptional regulatory modules. This process was repeated for each domain level state.

To visualize the predicted regulatory modules in different chromatin states, we have shown the estimated binding intensities of the proteins and the corresponding clusters in Figure 4(a), (b) and in Supplementary Figures S6, S7. We took a closer look at the clustering results in two contrasting states—Broad Promoter (Figure 4(a)) and Poised Enhancer (Figure 4(b)), and found noticeable differences in the binding intensity shape of both individual proteins and the predicted clusters between the two states. In addition, the set
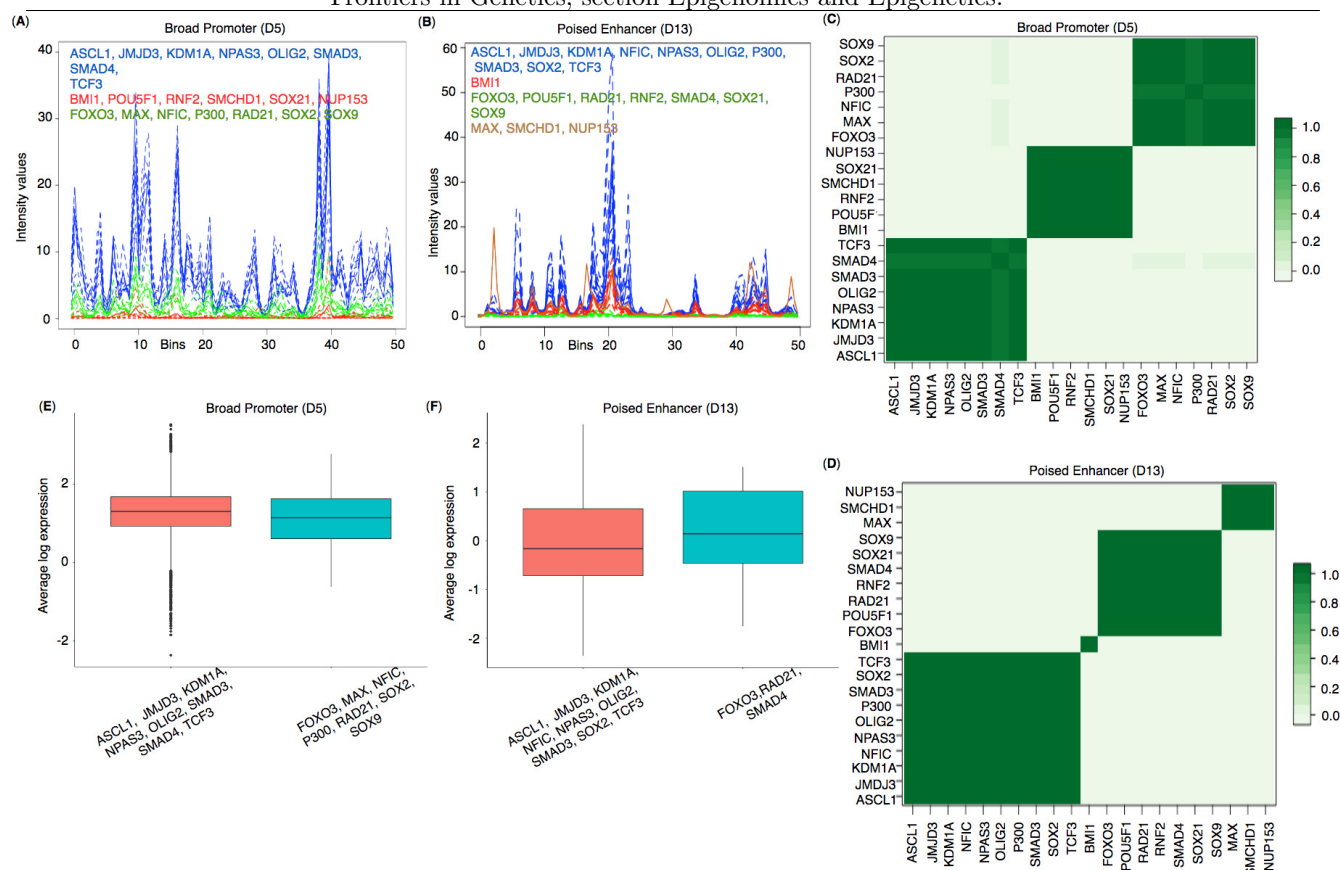
**Figure 4.** (a), (b) Estimated cluster binding intensities along with the individual TF binding intensities in the Broad Promoter and Poised Enhancer states, respectively. In each figure, the estimated binding intensities of the individual TFs are shown in dotted lines and the estimated binding intensities of the clusters are shown in solid line. TFs in each cluster are shown in the same color as that of the cluster. The X axis represents the genomic locations mapped on the real line between 0 and 50. The Y axis represents the estimated binding intensities, both for the individual TFs and for the identified clusters. (c), (d) Pairwise protein co-binding probabilities corresponding to (a) and (b) respectively. (e), (f) Comparison of proximal gene expressions (TPM) regulated by the clusters in (a) and (b) respectively. Only those clusters having (1) multiple TFs and (2) proximal genes for at least two TFs are shown in the figure to explain the combinatorial regulation of gene expressions by multiple TFs.

156  of co-factors for different proteins varied between the two states. BMI1 is known to bind to repressed
157  and poised states (Bhattacharya et al., 2015) and was predicted as a single-protein cluster in the Poised
158  Enhancer (Figure 4(b)) and Bivalent Promoter states (Supplementary Figure S6). In other states such as
159  Broad Promoter, Super Enhancer, and Upstream Enhancer, BMI1 was predicted with RNF2, RAD21, or
160  SMCHD1 (Supplementary Figures S6, S7). It is worth noting that both BMI1 and RNF2 are components of
161  the Polycomb group multi-protein, whereas SMCHD1, a non-canonical member of the SMC super-family,
162  is also known to be associated with transcriptional repression (Chen et al., 2015) and polycomb recruitment
163  mechanisms (Gendrel et al., 2012). The proposed approach was able to cluster several other functionally
164  relevant proteins that shared similar binding patterns, for example, JMJD3-SMAD3 (Figure 4) in most
165  chromatin states (in Estarás et al., 2012, the authors found that JMJD3 is recruited to gene promoters by
166  SMAD3 in neural stem cells and is essential to activate TGF-$\beta$ -responsive genes), FOXO3-NFIC-SOX-
167  TCF3 (Supplementary Figures S6) in Upstream Enhancer states (in Mateo et al., 2015, the authors showed

168 interactions among NFI family, TCF3, SOX2, SOX9, and FOXO3. We have shown additional predicted
169 protein-protein interactions in Supplementary Table S1.

170    To assess the strength of association between two co-binding proteins, we calculated a pairwise protein
171 co-binding probability matrix from the posterior samples of the MCMC procedure (Figure 4(c), (d)). Each
172 value in Figure 4(c), (d) indicates the frequency of observing the corresponding two proteins in the same
173 cluster out of the total 200 MCMC iterations. A high protein co-binding probability (indicated by darker
174 color) provides stronger evidence of the existence of the protein pair in a cluster. We further performed a
175 three-fold assessment on the robustness of the clustering algorithm explained in Supplementary document
176 section 5.

177    We next examined the expression levels of proximal genes (Transcripts Per Kilobase Million or TPM)
178 regulated by the predicted clusters in each state to understand transcriptional regulation by combinatorial
179 binding of proteins in different chromatin states. We observed that the median expression level of the
180 genes regulated by distinct clusters are close to each other in the Broad Promoter state (Figure 4(e)).
181 On the contrary, the median expression level of the proximal genes combinatorially regulated by the
182 FOXO3-RAD21-SMAD4 cluster in Poised Enhancer was higher than that of the genes combinatorially
183 regulated by the other cluster (Figure 4(f)) (Similar behavior was observed in Bivalent Promoter, Upstream
184 Enhancer and Boundary states shown in Supplementary Figure S9). These results show that gene expression
185 could change due to combinatorial binding of proteins in different chromatin states.

## 3.4    Comparison of results with other clustering methods

187    We compared the clustering results of the proposed algorithm with K-means and CLARANS (Ng and
188 Han, 2002). Instead of applying these two clustering methods directly on the binding locations of the
189 proteins, we first estimated individual protein binding intensities and used these intensity matrices as inputs
190 for clustering (we assumed each protein was in its own cluster). For both methods, we first obtained the
191 optimal number of clusters using the NBclust package (Charrad et al., 2014). From the results in Table
192 1, we observe that for both methods, the number of optimal clusters was 2 for the two chromatin states.
193 However, the cluster compositions that contain the regulatory TF modules are very similar to that of the
194 proposed approach. Furthers comparisons are provided in Supplementary Table S5.

## 3.5    Protein-DNA binding motif preferences in chromatin states

196    It is known that local epigenetic states affect bindings of proteins to targets and protein-DNA binding
197 may prevent or facilitate epigenetic changes on their binding sites (Blattler and Farnham, 2013; Xin and
198 Rohs, 2018). A protein is known to bind to the DNA with different motifs depending on the presence of
199 its co-binding partners (Bais et al., 2011). To examine the influence of chromatin states and co-binding
200 partners on the binding sequences of a protein, we grouped ChIP-seq peaks for each protein overlapped
201 with each chromatin state and analyzed the binding motifs of the protein in an active (Broad Promoter/Super
202 Enhancer) and a repressed state (Poised Enhancer/Polycomb Repressed) (Figure 5(a), (b)). We used the
203 MEME suite (Bailey et al., 2009) to identify de-novo motif sequences and from the results we selected the
204 motif that matched with the candidate protein's consensus motif or was known as a secondary motif. In
205 both the HOMER (Heinz et al., 2010) or JASPAR (Mathelier et al., 2016) databases, no reference motif is
206 documented for BMI1, KDM1A, JMJD3, NPAS3, NUP153, RNF2, RAD21, P300, and SMCHD1. For
207 the remaining proteins with known motifs, we extracted genomic sequences from two different subsets of
208 peaks overlapped with two contrasting chromatin states as mentioned before and determined the de-novo
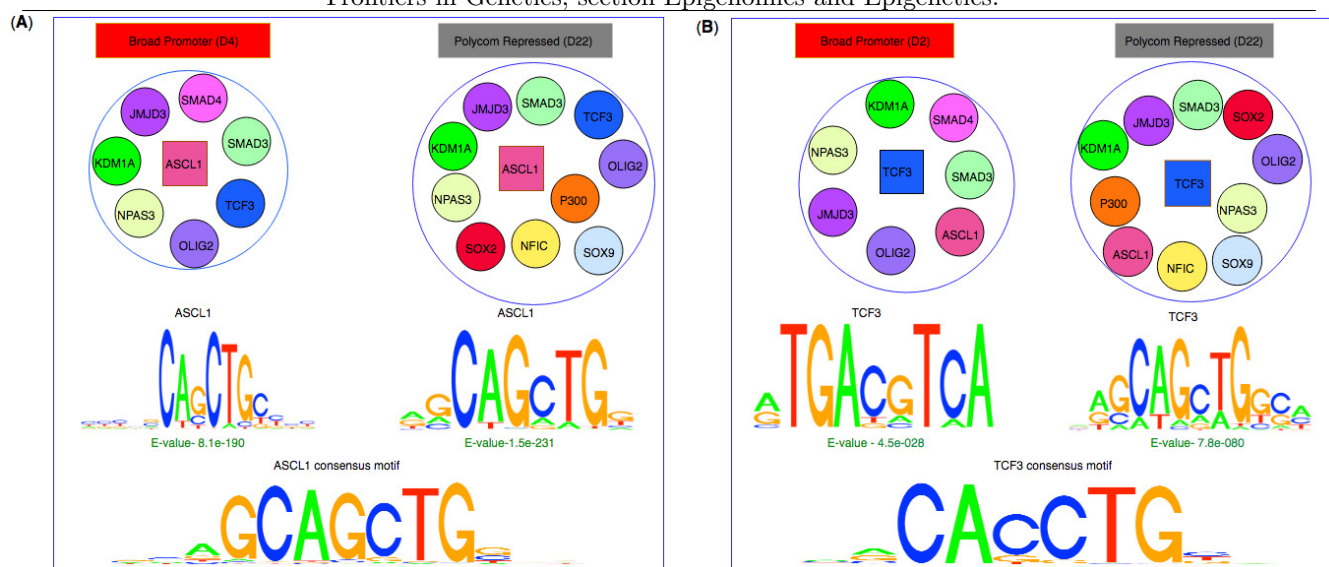209 motifs.

**Figure 5.** Effect of chromatin states and co-binding partner on binding motifs. (a) De-novo motifs obtained using MEME for ASCL1 are similar to the consensus motif in both Broad Promoter and Polycomb Repressed states although the co-factors of ASCL1 are different in the two states. (b) De-novo motifs obtained using MEME for TCF3 show differences in motifs between the two states with different co-factors. The motifs in active state resemble the $\beta$-catenin/TCF/LEF motif whereas the motifs in repressed state resemble the E-Box consensus motif.

210    Based on the MEME results, a protein's binding preferences may be broadly categorized into one of the
211  three types: (1) De-novo sequences that closely matched the protein's consensus motif such as ASCL1
212  (Figure 5(a)), MAX, NFIC, FOXO3, and TFs from the SOX family. (2) De-novo sequences that either
213  did not match with the consensus/secondary motifs or matched the consensus motif but were weakly
214  enriched. It has been observed that the ATF/CREB motifs ('TGAYRTCA') are often enriched in genes
215  targeted by $\beta$-catenin/TCF/LEF (Taniue et al., 2016; Lien et al., 2014). For TCF3, we observed highly
216  enriched de-novo sequences resembling its consensus motif in the repressed state (Figure 5(b)). However,
217  in the active state we observed that the 'TGACGTCA' pattern was highly enriched. This could imply
218  that TCF3 might have been recruited by other co-factors resulting in indirect binding in that particular
219  state. For OLIG2, both active and repressed chromatin states contained de-novo sequences resembling
220  its consensus motif. However, these sequences were highly enriched in the repressed state and weakly
221  enriched in the active state. The fact that the E-value of the de-novo sequences of OLIG2 was not significant
222  in the active state might suggest indirect binding in the state, probably being governed by other factors. (3)
223  De-novo sequences resembling the secondary motifs such the SMAD family. For SMAD4, we observed
224  that sequences with 'GCCGC' pattern were highly enriched in both active and repressed chromatin states,
225  as reported previously in (Hu et al., 2013) where the authors found that SMAD4 can bind to both methylated
226  and un-methylated motifs of distinct sequences. Similarly, for SMAD3, we observed highly enriched
227  sequences rich in 'GC' content in both chromatin states, which have been reported as secondary SMAD3
228  motifs, often associated with known SMAD binding partners in TGF-$\beta$ responses (Vidakovic et al., 2015).
229  Interestingly, for POU5F1, we observed that the E-Box element 'CANNTG' was significantly enriched
230  in both active and repressed chromatin states. In Yin et al. (2017), the authors had also observed that the
231  E-Box motif was significantly enriched with a p-value of 1e-6 in a POU5F1 ChIP-seq experiment of ES
232  cell with Dnmt1, Dnmt3A and Dnmt3B triple knockout, whereas the consensus POU5F1 motif was weakly
233  enriched with a p-value of 0.1. Detailed results are provided in Supplementary Table S3.

## 4   DISCUSSION

234   Development of the semi-automated genome annotation tools has enabled genome segmentation and
235   identification of distinct chromatin states at fine resolutions. In this study, we designed a two-step process to
236   identify transcriptional regulatory modules within distinct chromatin states. First, we segmented the genome
237   using the diHMM software. Second, we designed a novel nonparametric Bayesian clustering algorithm to
238   identify clusters of co-binding proteins on the segmented genome. Existing work have adopted distance
239   thresholds and empirical tests to define pairwise co-bound regions and context-dependent co-regulators (Ji
240   et al., 2006; Chen et al., 2008; Orlov et al., 2009; Lee and Zhou, 2013). The statistically principled approach
241   we proposed models protein-DNA binding site locations through inhomogeneous Poisson processes. It
242   also employs a Dirichlet process prior to the random distribution of the latent log-intensity functions to
243   facilitate clustering of the binding patterns. Such a nonparametric Bayesian clustering procedure is based
244   on joint likelihood rather than pairwise protein-protein relationship and is flexible in capturing the intricate
245   protein-DNA binding patterns in ChIP-seq data. This approach does not require pre-specified parameters
246   such as window size, distance threshold, and number of clusters, and hence achieves improved robustness.

247   We applied the approach on ChIP-seq data for neural stem cells obtained from ChIP-Atlas, an
248   integrated and comprehensive database rapidly gaining importance in cell replacement therapy. Despite
249   the methodological advantages, this approach may have limitations in practical use. First, ChIP-seq can
250   produce millions of short reads, which may result in varying strengths of signal intensities along the
251   genome. In the current study, we did not consider the peak-height for different proteins but treated the
252   center of each peak as a binary binding event along the genome. The overlook of the signal intensity effects
253   may impact the modeling of protein binding patterns. Another possible limitation of our approach lies in
254   handling the three dimensional structural information of the histone marks. This restricted our downstream
255   gene expression analysis to gene promoters present in the Enhancer states. While not in scope of the current
256   study, including such information may improve the accuracy of the model and enable the prediction of
257   long distance Enhancer activity.

258   Nevertheless, we were able to establish several interesting findings. It has been known that protein-
259   DNA binding sites are not randomly distributed but rather clustered together at enhancer or promoter
260   regions. Hence, some specific proteins may team up to have a significant epigenetic impact on gene
261   expression. In our study, transcriptional regulatory modules identified in different chromatin states revealed
262   several known protein-protein interactions in neural stem cells, for example, SOX family and NF1 in the
263   Enhancer states (Webb et al., 2013), MAX-FOXO3-OLIG2 in Upstream Enhancer (Mateo et al., 2015),
264   and JMJD3-SMAD3 in most chromatin states (Estarás et al., 2012). These results suggest chromatin-
265   state-specific protein-protein co-occupancy. In addition, diverse gene expression levels were observed
266   through combinatorial regulation by the predicted transcriptional regulatory modules in different states.
267   The uncovered links between gene expression and protein binding patterns on a genome-wide scale will
268   enhance our understanding on how chromatin-state-specific regulatory network is assembled to coordinate
269   tissue differentiation and cell specification.

270   An important issue in transcription regulation is to understand the binding specificity and affinity of a
271   protein. A TF may have several thousands of DNA binding sites along the genome, which collectively
272   can be represented as a motif—a consensus sequence demonstrating the nucleotide preferences at each
273   position of the binding site. In this study, we observed that chromatin state can have an impact on the
274   binding preferences of transcription factors and their co-activators (Jolma et al., 2015). For example, the
275   de-novo sequences predicted for the some proteins resembled the consensus PWM across distinct chromatin
276   states whereas for certain proteins such as SMAD family the sequences resembled secondary motifs in

277 specific chromatin states. Further, we also noticed that the prediction of binding preferences might help the
278 identification of indirect protein bindings when the de-novo sequences do not match the consensus PWM
279 (Yin et al., 2017). In conclusion, we expect that our work will help understand the causality of chromatin
280 state and combinatorial protein-DNA binding in regulating gene expression in neural stem cells.

## AUTHOR CONTRIBUTIONS

281 H.X. conceived and designed the study; H.Z. and X.W. designed and implemented the clustering model;
282 S.B and M.T. designed computational experiments and performed data analyses; H.X., W.F., X.W., and
283 S.B. wrote the original draft. All authors read and approved the final manuscript.

## FUNDING

## REFERENCES

286 Bailey, T. L., Boden, M., Buske, F. A., Frith, M., Grant, C. E., Clementi, L., et al. (2009). Meme suite:
287     tools for motif discovery and searching. *Nucleic acids research* 37, W202–W208
288 Bais, A. S., Kaminski, N., and Benos, P. V. (2011). Finding subtypes of transcription factor motif pairs
289     with distinct regulatory roles. *Nucleic Acids Research* 39, e76–e76
290 Bhattacharya, R., Mustafi, S. B., Street, M., Dey, A., and Dwivedi, S. K. D. (2015). Bmi-1: At the
291     crossroads of physiological and pathological biology. *Genes & Diseases* 2, 225–239
292 Blattler, A. and Farnham, P. J. (2013). Cross-talk between site-specific transcription factors and DNA
293     methylation states. *Journal of Biological Chemistry* 288, 34287–34294
294 Cha, M. and Zhou, Q. (2014). Detecting clustering and ordering binding patterns among transcription
295     factors via point process models. *Bioinformatics* 30, 2263–2271
296 Charrad, M., Ghazzali, N., Boiteau, V., Niknafs, A., and Charrad, M. M. (2014). Package 'nbclust'.
297     *Journal of Statistical Software* 61, 1–36
298 Chen, K., Hu, J., Moore, D. L., Liu, R., Kessans, S. A., Breslin, K., et al. (2015). Genome-wide binding and
299     mechanistic analyses of smchd1-mediated epigenetic regulation. *Proceedings of the National Academy*
300     *of Sciences* 112, E3535–E3544
301 Chen, X., Xu, H., Yuan, P., Fang, F., Huss, M., Vega, V. B., et al. (2008). Integration of external signaling
302     pathways with the core transcriptional network in embryonic stem cells. *Cell* 133, 1106–1117
303 Cuesta, I., Zaret, K. S., and Santisteban, P. (2007). The forkhead factor foxe1 binds to the thyroperoxidase
304     promoter during thyroid cell differentiation and modifies compacted chromatin structure. *Molecular and*
305     *Cellular Biology* 27, 7302–7314
306 Day, N., Hemmaplardh, A., Thurman, R. E., Stamatoyannopoulos, J. A., and Noble, W. S. (2007).
307     Unsupervised segmentation of continuous genomic data. *Bioinformatics* 23, 1424–1426
308 Ernst, J. and Kellis, M. (2012). Chromhmm: automating chromatin-state discovery and characterization.
309     *Nature Methods* 9, 215–216
310 Estarás, C., Akizu, N., García, A., Beltrán, S., de la Cruz, X., and Martínez-Balbás, M. A. (2012).
311     Genome-wide analysis reveals that smad3 and jmjd3 hdm co-activate the neural developmental program.
312     *Development* 139, 2681–2691

313 Feng, J., Liu, T., Qin, B., Zhang, Y., and Liu, X. S. (2012). Identifying ChIP-seq enrichment using MACS.
314     *Nature Protocols* 7, 1728–1740

315 Gendrel, A.-V., Apedaile, A., Coker, H., Termanis, A., Zvetkova, I., Godwin, J., et al. (2012). Smchd1-
316     dependent and-independent pathways determine developmental dynamics of CpG island methylation on
317     the inactive x chromosome. *Developmental cell* 23, 265–279

318 Heinz, S., Benner, C., Spann, N., Bertolino, E., Lin, Y. C., Laslo, P., et al. (2010). Simple combinations of
319     lineage-determining transcription factors prime cis-regulatory elements required for macrophage and b
320     cell identities. *Molecular Cell* 38, 576–589

321 Hoffman, M. M., Buske, O. J., Wang, J., Weng, Z., Bilmes, J. A., and Noble, W. S. (2012). Unsupervised
322     pattern discovery in human chromatin structure through genomic segmentation. *Nature Methods* 9,
323     473–476

324 Hu, S., Wan, J., Su, Y., Song, Q., Zeng, Y., Nguyen, H. N., et al. (2013). DNA methylation presents distinct
325     binding sites for human transcription factors. *Elife* 2

326 Ji, H., Vokes, S. A., and Wong, W. H. (2006). A comparative analysis of genome-wide chromatin
327     immunoprecipitation data for mammalian transcription factors. *Nucleic Acids Research* 34, e146–e146

328 Jolma, A., Yin, Y., Nitta, K. R., Dave, K., Popov, A., Taipale, M., et al. (2015). DNA-dependent formation
329     of transcription factor pairs alters their binding specificity. *Nature* 527, 384–388

330 Lee, Y. and Zhou, Q. (2013). Co-regulation in embryonic stem cells via context-dependent binding of
331     transcription factors. *Bioinformatics* 29, 2162–2168

332 Libbrecht, M. W., Ay, F., Hoffman, M. M., Gilbert, D. M., Bilmes, J. A., and Noble, W. S. (2015). Joint
333     annotation of chromatin state and chromatin conformation reveals relationships among domain types
334     and identifies domains of cell-type-specific expression. *Genome Research* 25, 544–557

335 Lien, W.-H., Polak, L., Lin, M., Lay, K., Zheng, D., and Fuchs, E. (2014). In vivo transcriptional
336     governance of hair follicle stem cells by canonical wnt regulators. *Nature Cell Biology* 16, 179

337 Liu, L., Jin, G., and Zhou, X. (2015). Modeling the relationship of epigenetic modifications to transcription
338     factor binding. *Nucleic Acids Research* 43, 3873–3885

339 Liu, L., Zhao, W., and Zhou, X. (2016). Modeling co-occupancy of transcription factors using chromatin
340     features. *Nucleic Acids Research* 44, e49–e49

341 Marco, E., Meuleman, W., Huang, J., Glass, K., Pinello, L., Wang, J., et al. (2017). Multi-scale chromatin
342     state annotation using a hierarchical hidden Markov model. *Nature Communications* 8

343 Mateo, J. L., van den Berg, D. L., Haeussler, M., Drechsel, D., Gaber, Z. B., Castro, D. S., et al. (2015).
344     Characterization of the neural stem cell gene regulatory network identifies olig2 as a multifunctional
345     regulator of self-renewal. *Genome Research* 25, 41–56

346 Mathelier, A., Fornes, O., Arenillas, D. J., Chen, C.-y., Denay, G., Lee, J., et al. (2016). Jaspar 2016: a
347     major expansion and update of the open-access database of transcription factor binding profiles. *Nucleic*
348     *Acids Research* 44, D110–D115

349 Neal, R. M. (2000). Markov chain sampling methods for Dirichlet process mixture models. *Journal of*
350     *Computational and Graphical Statistics* 9, 249–265

351 Ng, R. T. and Han, J. (2002). Clarans: A method for clustering objects for spatial data mining. *IEEE*
352     *transactions on knowledge and data engineering* 14, 1003–1016

353 Orlov, Y. L., Huss, M. E., Joseph, R., Xu, H., Vega, V. B., Lee, Y. K., et al. (2009). Genome-wide statistical
354     analysis of multiple transcription factor binding sites obtained by ChIP-seq technologies. In *Proceedings*
355     *of the 1st ACM Workshop on Breaking Frontiers of Computational Biology* (ACM), 11–18

356 Rue, H. and Held, L. (2005). *Gaussian Markov Random Fields: Theory and Applications* (CRC press)

357   Rue, H., Martino, S., and Chopin, N. (2009). Approximate Bayesian inference for latent Gaussian models
358     by using integrated nested Laplace approximations. *Journal of the Royal Statistical Society: Series B*
359     *(Statistical Methodology)* 71, 319–392

360   Sharmin, M., Bravo, H. C., and Hannenhalli, S. (2016). Heterogeneity of transcription factor binding
361     specificity models within and across cell lines. *Genome Research* 26, 1110–1123

362   Simpson, D., Illian, J. B., Lindgren, F., Sørbye, S. H., and Rue, H. (2016). Going off grid: Computationally
363     efficient inference for log-Gaussian Cox processes. *Biometrika* 103, 49–70

364   Soufi, A., Garcia, M. F., Jaroszewicz, A., Osman, N., Pellegrini, M., and Zaret, K. S. (2015). Pioneer
365     transcription factors target partial DNA motifs on nucleosomes to initiate reprogramming. *Cell* 161,
366     555–568

367   Sugathan, A. and Waxman, D. J. (2013). Genome-wide analysis of chromatin states reveals distinct
368     mechanisms of sex-dependent gene regulation in male and female mouse liver. *Molecular and cellular*
369     *biology* 33, 3594–3610

370   Taniue, K., Kurimoto, A., Takeda, Y., Nagashima, T., Okada-Hatakeyama, M., Katou, Y., et al. (2016).
371     Asbel-tcf3 complex is required for the tumorigenicity of colorectal cancer cells. *Proceedings of the*
372     *National Academy of Sciences of the United States of America* 113, 12739–12744

373   Vidakovic, A. T., Rueda, O. M., Vervoort, S. J., Batra, A. S., Goldgraben, M. A., Uribe-Lewis, S., et al.
374     (2015). Context-specific effects of tgf-$\beta$/smad3 in cancer are modulated by the epigenome. *Cell Reports*
375     13, 2480–2490

376   Webb, A. E., Pollina, E. A., Vierbuchen, T., Urbán, N., Ucar, D., Leeman, D. S., et al. (2013). Foxo3
377     shares common targets with ascl1 genome-wide and inhibits ascl1-dependent neurogenesis. *Cell Reports*
378     4, 477–491

379   Wong, K.-C., Li, Y., Peng, C., and Zhang, Z. (2015). Signalspider: probabilistic pattern discovery on
380     multiple normalized ChIP-seq signal profiles. *Bioinformatics* 31, 17–24

381   Xin, B. and Rohs, R. (2018). Relationship between histone modifications and transcription factor binding
382     is protein family specific. *Genome Research* 1, gr–220079

383   Yin, Y., Morgunova, E., Jolma, A., Kaasinen, E., Sahu, B., Khund-Sayeed, S., et al. (2017). Impact
384     of cytosine methylation on DNA binding specificities of human transcription factors. *Science* 356,
385     eaaj2239

386   Zaret, K. S. and Carroll, J. S. (2011). Pioneer transcription factors: establishing competence for gene
387     expression. *Genes & Development* 25, 2227–2241

**Table 1.** Comparison of clustering results with other methods

| Chromatin state | DPM-LGCP | K-means | CLARANS |
|---|---|---|---|
| Broad Promoter (D5) | (1) ASCL1, JMJD3, KDM1A, NPAS3, OLIG2, SMAD3, SMAD4, TCF3; (2) BMI1, POU5F1, RNF2, SMCHD1, SOX21, NUP153; (3) FOXO3, MAX, NFIC, P300, RAD21, SOX2, SOX9 | (1) ASLC1, JMJD3, KDM1A, NFIC, NPAS3, OLIG2, SMAD3, SMAD4, TCF3; (2) BMI1, FOXO3, MAX, P300, POU5F1, RAD21, RNF2, SMCHD1, SOX2, SOX21, SOX9, NUP153 | (1) ASCL1, FOXO3, JMJD3, KDM1A, NFIC, NPAS3, OLIG2, RAD21, SMAD3, SMAD4, SOX2, SOX9; (2) BMI1, MAX, P300, POU5F1, RNF2, SMCHD1, SOX21, NUP153 |
| Poised Enhancer (D13) | (1) ASCL1, JMJD3, KDM1A, NFIC, NPAS3, OLIG2, P300, SMAD3, SOX2, TCF3; (2) BMI1; (3) FOXO3, POU5F1, RAD21, RNF2, SMAD4, SOX21, SOX9, TCF3; (4) MAX, SMCHD1, NUP153 | (1) ASCL1, JMJD3, KDM1A, NFIC, NPAS3, OLIG2, P300, SMAD3, SOX2, SOX9, TCF3; (2) BMI1, FOXO3, MAX, POU5F1, RAD21, RNF2, SMAD4, SMCHD1, SOX21, NUP153 | (1) ASCL1, FOXO3, JMJD3, KDM1A, NFIC, NPAS3, OLIG2, P300, POU5F1, SMAD3, SMAD4, SOX2, SOX9, TCF3; (2) BMI1, MAX, RAD21, RNF2, SMCHD1, SOX21, NUP153 |

For each method the clusters are preceded by the cluster number within parentheses. Further comparisons are shown in Supplementary Table S5.