1    **Heterogeneity of Transcription Factor binding specificity models within and across cell lines**

2    Mahfuza Sharmin[1,2], Héctor Corrada Bravo[1,2], Sridhar Hannenhalli[2,3*]

3    Center for Bioinformatics and Computational Biology

4    [1]Department of Computer Science

5    [2]Center for Bioinformatics and Computational Biology

6    [3]Department of Cell and Molecular Biology

7    University of Maryland, College park, MD

8    [*]Corresponding author

9

10    **Corresponding author**

11    Sridhar Hannenhalli

12    3104G Biomolecular Sciences Building (#296)

13    University of Maryland, College Park, MD 20742, USA

14    301 405 8219 (v) 301 314 1341 (f)

15    sridhar@umiacs.umd.edu

16

17

18    **Key words:** Heterogeneity, Sequence Pattern, TF binding specificity, Ensemble models, Clustering, Co-
19    factors, Cell specificity

20

21     **Abstract**

22     Complex gene expression patterns are mediated by binding of transcription factors (TF) to specific
23     genomic loci. The *in vivo* occupancy of a TF is, in large part, determined by the TF's DNA binding
24     interaction partners, motivating genomic context based models of TF occupancy. However, the
25     approaches thus far have assumed a uniform binding model to explain genome wide bound sites for a TF
26     in a cell-type and as such heterogeneity of TF occupancy models, and the extent to which binding rules
27     underlying a TF's occupancy are shared across cell types, has not been investigated. Here, we develop an
28     ensemble based approach (*TRISECT*) to identify heterogeneous binding rules of cell-type specific TF
29     occupancy and analyze the inter-cell-type sharing of such rules. Comprehensive analysis of 23 TFs, each
30     with ChIP-Seq data in 4-12 cell-types, shows that by explicitly capturing the heterogeneity of binding
31     rules, *TRISECT* accurately identifies *in vivo* TF occupancy (93%) substantially improving upon previous
32     methods. Importantly, many of the binding rules derived from individual cell-types are shared across
33     cell-types and reveal distinct yet functionally coherent putative target genes in different cell-types.
34     Closer inspection of the predicted cell-type-specific interaction partners provides insights into context-
35     specific functional landscape of a TF. Together, our novel ensemble-based approach reveals, for the first
36     time, a widespread heterogeneity of binding rules, comprising interaction partners within a cell-type,
37     many of which nevertheless transcend cell-types. Notably, the putative targets of shared binding rules in
38     different cell-types, while distinct, exhibit significant functional coherence.

39      **Introduction**

40      Transcriptional regulation is critically mediated by the binding of transcription factors (TF) to specific
41      DNA elements in the genome (JACOB & MONOD 1961; Busby & Ebright 1994). While the *in vitro* binding
42      specificity of many human TFs has been determined, it is well recognized that the *in vitro* binding
43      specificity of a TF does not explain its condition-specific *in vivo* binding specificity (Zinzen et al. 2009;
44      Yáñez-Cuna et al. 2012). This recognition has spurred investigations of additional determinants of *in vivo*
45      binding, such as heterogeneity of TF's binding motif (Hannenhalli & Levy 2002), homotypic clusters of
46      binding sites (Dror et al. 2015),  cooperative binding of the TF with its partners (Wang et al. 2006),
47      condition-specific chromatin context (Heintzman et al. 2009), local DNA properties (Dror et al. 2015),
48      epigenomic context (Gheldof et al. 2010) etc. While overall, both local genomic and epigenomic features
49      have been deemed important in determining in *in vivo* occupancy of a TF, recent reports suggest that *in*
50      *vivo* binding of a TF can be accurately predicted based only on the genomic signatures near the binding
51      site (BS) without relying on the epigenomic context (Arvey et al. 2012; Dror et al. 2015); this is consistent
52      with very recent reports showing that the epigenome itself is encoded by the genomic context
53      (Whitaker et al. 2015; Benveniste et al. 2014). Taken together, these results strongly suggest that
54      proximal genomic elements are the primary driver of *in vivo* TF binding. Prior sequence-based models of
55      *in vivo* TF binding have shown that, somewhat counter-intuitively, the genomic context of a BS, which is
56      the property of the genome, effectively encodes the condition-specific *in vivo* binding specificity (Arvey
57      et al. 2012). This can be explained by the substantial plasticity of a TF's interaction with other TFs' and
58      the modular nature of TF binding cooperatively with other TFs (Frietze & Farnham 2011), such that
59      availability of specific combination of interacting TFs can guide *in vivo* binding to specific loci where the
60      BS of the interacting TF are present in close proximity to each other, along with the availability of
61      corresponding TFs (Hannenhalli & Levy 2002).

62      Previous sequence-based modeling of *in vivo* TF binding was done in a cell type-specific fashion. These
63      cell type-specific models exhibit substantial *inter*-cell type heterogeneity, as expected, given variation in
64      the availability of the potentially interacting TFs. However, these previous approaches build a single
65      model for a cell type, thus implicitly assuming a homogeneous cell type-specific model, and as such have
66      not investigated *intra*-cell-type model heterogeneity. Such heterogeneity of TF binding 'rule' across the
67      genome can be expected for the same reason as for the inter-cell type heterogeneity. Moreover, in
68      many instances, a binding specificity model trained in one cell type can predict a subset of *in vivo*
69      binding in a different cell type (Arvey et al. 2012), suggesting that models of binding, or parts thereof,
70      may be shared across cell types. Overall, the heterogeneity of sequence-based models of cell type-
71      specific *in vivo* TF binding, and the extent to which a subset of binding rules (*sub-models*) are shared
72      across cell types, is not known, motivating the present study.

73      To this end, we have developed an ensemble model based approach (***TRISECT***) to reveal both cell-
74      specific and cell-independent rules for *in vivo* TF binding. We applied *TRISECT* to 23 TFs, each with
75      genome-wide *in vivo* binding data in 4 – 12 cell types (a total of 135 TF-cell type combinations). For each
76      TF, for each cell type, we built ensemble models of *in vivo* TF binding (***EMT***), then decomposed each
77      *EMT* model into sub-models and clustered the pooled set of sub-models across all cell types using
78      feature selection. Our comprehensive analyses strongly suggest that the cell type-specific binding rule
79      for a TF consists of multiple sub-models, supported by our result showing that *EMT* captures the binding
80      specificity better than previous non-ensemble models (Arvey et al. 2012). Moreover, for many TFs, the
81      sub-models are shared across cell-types, and interestingly, we found that the putative target genes for

82    similar sub-models across cell types exhibit a high degree of expression and functional coherence,
83    suggesting that the *in vivo* binding rules are related to function of the gene targets, much more so than
84    the cell type they are derived from.

85    In further probing the superior performance of *EMT*, we demonstrate that while a model based only on
86    the known motifs of the reference TF, *i.e.* without incorporating additional potential TF interaction
87    partners (*NonInteraction* model), can predict *in vivo* binding with ~78% accuracy, when motifs for other
88    TFs are used in the model (*Interaction model*), the prediction accuracy is substantially increased to over
89    90%. Moreover, we found that the improvement in prediction accuracy by the *Interaction* model
90    strongly correlates with the increase in the number of interaction partners, i.e., with model complexity,
91    suggesting that the *Interaction* model effectively captures the heterogeneity of the binding rule. We
92    identified and validated, based on literature, the potential interaction partners (we will refer to these as
93    co-factors) that mediate context-specific binding and function of a TF. Finally, we show that certain TFs
94    with multiple distinct binding motifs prefer binding to different motifs in different cell types, which may
95    in part be associated with their inter-cell type variability of co-factors (Slattery et al. 2011).

96    In sum, our analysis reveals distinct sub-models of *in vivo* TF binding within a cell type that are
97    nevertheless shared across cell types, and the shared sub-models across cell types target distinct yet co-
98    functional genes in different cell types. A refined understanding of the genomic context of *in vivo*
99    binding specificity can facilitate future investigations of transcriptional regulation and understanding of
100   its genetic determinants.

101

102   **Results**

103   **TRISECT – *Ensemble model of TF binding and the Clustering of sub-models across cell types***

104   ***Overview.*** The full analysis pipeline, *TRISECT*, is illustrated by Fig 1A. As the first step, , we developed an
105   ensemble model (*EMT*) to discriminate a TF's *in vivo* bound genomic loci from the background, balancing
106   model complexity (number of sub-models in the ensemble) against the cross-validation classification
107   accuracy. Given a set of genome-wide loci bound by a specific TF, we first construct a foreground set of
108   sequences (100 bps) centered at the ChIP-Seq peak. As a stringent background control, as done
109   previously (Arvey et al. 2012), we use 100 bps regions ~200 bps away from the peak location (M&M).
110   We considered a variety of feature sets for discrimination (see below). The *EMT* model was trained using
111   Adaboost method where each sub-model is a decision tree built from a bootstrap sample (Friedman et
112   al. 2000; Friedman 2002; Freidman 2008). Next, for a TF, given EMT models for all cell types, we
113   represented each cell type-specific sub-model as a point in a *d*-dimensional space corresponding to *d*
114   selected features (M&M). We clustered the data points, representing all sub-models in all cell types
115   considered for a TF, using XY-fused network (*XYF*) (Melssen et al. 2006) such that sub-models within a
116   cluster represent similar binding rules, either within a cell-type or across cell types.

117   ***EMT Feature sets.*** We considered three types of feature sets for a 100 bps sequence – (i) Kmer:
118   frequency of occurrence for all 4096 6-mers, (ii) KmerRC: frequency of occurrence for all 2080 6-mers
119   where a k-mer and its reverse complement were unified, and (iii) aggregate binding scores for 981
120   vertebrate TF motifs from TRANSFAC database (we used four stringencies for motif match) (M&M); we
121   refer to these as *pwm* models. We applied *TRISECT* to 23 TFs, each with ChIP-Seq data in 4 to 12 cell

122  types (a total of 135 TF-cell pair *EMT*s), listed in Supplementary Table 1. A TF was included in this study if
123  (i) TF has narrow-peak data for at least 4 cell lines with at least 4k sites in each cell line, and (ii) TF has
124  established PWM in TRANSFAC 2011 database. The performance assessment of *EMT*s was conducted
125  based on 25% held-out dataset.  The overall performance is summarized in Fig 1B and details are
126  provided in Supplementary Table 2.

127  ***EMT performance.*** Fig 1B shows the overall accuracy distribution (over 135 TF-cell type pairs) for the 6
128  types of models, where the accuracy is quantified using ROCAUC on the test set. We compared the
129  performances, using Wilcoxon test, among 6 sets of *EMT*s (kmer, kmerRC, and PWM at 4 stringencies)
130  containing 135 TF-cell type pairs in each set (Fig 1C). We found that kmerRC significantly outperforms
131  kmer model (Wilcoxon p-value 2.65E-20), consistent with the fact that TF binding occurs on double-
132  stranded DNA and as such does not have directionality (except in relation with other interacting TFs)
133  and therefore unifying each kmer with its reverse complement provides a better abstraction of
134  biological determinants of TF binding. Following this line of reasoning, PWMs provide an even better
135  abstraction of DNA binding specificity and as expected, the PWM-based models outperform kmer-based
136  models (p-value, 2.29E-06 comparing kmerRC and pwm1k). Based on relative performances we selected
137  pwm1k-based EMT for feature selection and clustering of sub-models and all subsequent analyses.

138  ***Comparison with previous model.*** Next, we compared *EMT* model (using kmerRC and pwm1k) with
139  previously published model based on string kernel SVM (*SVM-kmer*) (Arvey et al. 2012). Supplementary
140  Table 3 lists 17 TFs for which ROCAUC was reported in (Arvey et al. 2012), where the mean accuracy
141  across multiple cell lines was reported for each TF. We therefore compared the published accuracy with
142  the mean *EMT* performance of the TF across only the cell types that were considered previously. As
143  shown in Fig 2, in most cases, *EMT* outperforms *SVM-kmer*. DNAse hypersensitive (DHS) of a region
144  represents its accessibility by DNA-binding proteins and previous studies have shown that integrating
145  DHS with *in vitro* binding specificity can substantially enhance *in vivo* binding prediction (Arvey et al.
146  2012; Pique-Regi et al. 2011). Surprisingly, using pwm1k features 6 cases *EMT* outperforms even the
147  model that integrated DHS with the kmer frequencies in the SVM (green). In a few cases (blue), SVM-
148  kmer yields either comparable or improved predictability. Overall, the *EMT* models predict *in vivo*
149  binding with a greater accuracy than a non-ensemble SVM approach represented by *SVM-kmer* (Arvey
150  et al. 2012).

151  In sum, we have described a novel ensemble-based approach to in vivo binding modeling and
152  established its superiority relative to SVM-kmer across a wide variety of TFs and cell types.

153

154  ***TRISECT reveals intra-cell type heterogeneity and inter-cell type sharing of binding rules across cell***
155  ***types***

156  The architectural difference and performance advantage of *EMT* relative to *SVM-kmer* suggests that
157  *EMT* might be better able to exploit heterogeneous binding rules across the genome dictated by
158  different combinations of interacting TFs. For each TF, we clustered the sub-models obtained from
159  different cell types. As an illustrative example, Fig 3A-B show the cluster-membership matrix for TF *ATF3*
160  for number of clusters *k* = 16 and 20. Fig S1 includes such mapping for all other TFs for *k* = 16. We found
161  both cell type-specific (Fig 3B, cluster #6) and ubiquitous (Fig 3C, cluster #20) clusters. Examining the
162  cluster mapping for all TFs (Fig S1), a wide range of patterns emerge: for certain TFs most clusters map

163 to single cell type, suggesting cell type-specific binding modalities of these TFs (*EP300, JUN*), while
164 certain other TFs have ubiquitously applicable binding rules, such as *YY1* and *TBP*, suggesting cell type
165 independent binding rules and, presumably, function. Importantly, many clusters consist of sub-models
166 from multiple, but not all, cell types. We ensured that inter-cell type sharing of *in vivo* binding rule is not
167 simply due to shared binding loci across cell types (Supplementary Notes & Fig S2). Subsequent analyses
168 are based on $k$ = 16; reasons for this choice are discussed in Supplementary Notes & Fig S3).

169 It is possible that *EMT* can falsely yield multiple sub-models, even in absence of heterogeneity, and
170 those sub-models can be falsely clustered. We ascertained heterogeneity across sub-models for a TF
171 from multiple cell types using a *Dudahart test* (Duda et al. 2001) and assessed the clustering tendency of
172 the sub-models in the $d$-dimensional feature space using *Hopkins statistics* (Jain & Dubes 1988). The
173 *Dudahart* test verifies whether or not a set of data points should be split into two clusters from the
174 estimate of within-cluster sum of squares for all pairs of clusters versus overall sum of squares; the ratio
175 of the two sum of squares is quantified as the *dh-ratio*. On the other hand, the *Hopkins statistic (H)*
176 compares the nearest neighbor distribution for a random set of points to the same distribution for the
177 clustered sub-models (M&M). A value close to 0.5 indicates the sub-models are random set of points
178 with no clustering, a value close to 1 indicates that they form a cluster. Fig 3C-D summarize the *dh-ratio*
179 and *Hopkins statistic* respectively for 135 TF-cell pairs based on sub-models of TF-cell type pair, and for
180 each TF after gathering all sub-models under a TF. We found that in all cases the *dh-ratio* is lower than 1
181 rejecting homogeneity (Fig 3C) and the set of sub-models form clusters (Fig 3D). All tests done for the
182 analysis are significant (p-value <0.001) (M&M). Together, the *Dudahart test* and *Hopkins statistic*
183 strongly suggest that the sub-models are distinct and clusterable, i.e., TF binding rules are
184 heterogeneous and partly shared across cell types.

185 Next we assessed the functional underpinning of shared binding rules across cell types. Specifically, we
186 assessed whether two co-clustered loci from different cell types (i.e., obeying similar binding rule) are
187 functionally associated relative to loci from the same cell type but belonging to different clusters, i.e.,
188 obeying different binding rules. We devised a cluster-specific scoring of each binding sequence and
189 assigned each binding site in each cell type to one or more clusters (M&M). As per convention, we
190 assigned each binding site to the nearest gene as a potential transcriptional target; 95% of the target
191 genes were within 100 kb from the binding site (median distance 4.5 kbp) (Fig S4). To assess functional
192 coherence of a cluster, we determined the fraction of gene-pairs in the cluster (regardless of cell type)
193 that participate in the same pathway as compared to all pairs of target genes within each cell type, and
194 assessed the significance of enrichment using Fisher test. Likewise, we also estimated the expression
195 coherence of genes within a cluster (M&M). As shown in Fig 4 and S5: ~40% (respectively, ~18%) multi-
196 cell type clusters show significantly higher (p-value <= 0.05) expression-coherence (respectively,
197 pathway-coherence) than the background (expectation is 5%). Moreover, the pathway and expression
198 coherence are highly correlated across clusters (spearman correlation=0.56, p-value=0.02). We
199 conducted the same set of tests for random clusters of same size as real clusters. In both cases, the
200 coherence was no greater than the null expectation (Fig 4A-B). In Supplementary Tables 5a-b, we
201 catalogue all the clusters with mapped target genes and their enriched GO terms.

202 Taken together, these analyses support existence of heterogeneous sets of rule governing *in vivo* TF
203 binding and that subset of rules are shared across cell types with functional implication.

204

205 ***The role of interaction partners in a TF's binding occupancy across cell types***

206 By using 981 PWMs for a comprehensive set of vertebrate TFs as the basis for features, *EMT* implicitly
207 incorporates the contributions of interaction partners in predicting *in vivo* binding of the reference TF.
208 To quantify the contribution of interacting motifs, we repeated the *EMT* training and testing using only
209 the PWMs corresponding to the reference TF. Individual TFs have multiple motifs reported in the
210 literature (ranging from 1 to 8, with a median of 3; Supplementary Table 6), which can differ
211 substantially from each other with potential functional implications (Bulyk et al. 2002; Hannenhalli
212 2008); we refer to these motifs as the *reference motifs*, and the *EMT* model utilizing only the reference
213 motifs as the *NonInteraction* model and to contrast we refer pwm1k model as *Interaction* model.
214 Supplementary Table 7 shows the prediction accuracies for the *Interaction* and the *NonInteraction*
215 *models*; the diagonal elements represent the cross-validation accuracies within a cell type, while the off-
216 diagonal elements represent the accuracy when *EMT* is trained on one cell type (row) and tested on
217 another (column).  Comparing the diagonal elements for the two models (summarized in Fig 5A), it is
218 evident that *Interaction models* have higher predictive accuracy than *NonInteraction* models, which is
219 consistent with the expectation that *in vivo* binding of a TF relies on interactions among several TFs.

220 Next, we conjectured that in the *Interaction* model, allowing for greater numbers of partners allows
221 learning of more complex binding rules and increase binding prediction accuracy. We therefore assessed
222 the effect of the length of the region flanking the binding site on prediction accuracy (M&M). We note
223 that beyond 100bp, due to narrowing of the gap between the foreground and the background region,
224 the discrimination accuracy is expected to decrease. Despite this, in some cases (Fig 5B & S6), the
225 increase in ROCAUC beyond 100bp suggests that a larger context may be necessary in these cases to
226 capture the binding rules. Nevertheless, we chose a sequence context of 100bp to make our model
227 comparable to the previously published *SVM-kmer* (Arvey et al. 2012).

228 For a given TF, we also quantified the variability of the model accuracy in different cell types (M&M). We
229 expect a model that relies on cell type-specific interaction partners to be more variable in its
230 performance accuracy than the one that relies only on the reference motifs. This expectation is borne
231 out in our analysis (Fig 5C). This suggests that part of the sequence information for *in vivo* binding is
232 encoded by the TF's own motifs and this does not vary substantially across cell types, while the
233 additional context- and interaction-dependent part does. However, the small variability in cross-cell type
234 prediction accuracy when using *NonInteraction* model is likely to come from the heterogeneity of
235 binding motifs for a TF. We quantified the inter-motif divergence for each TF as either the number of
236 motifs annotated for the TF, or motif-divergence defined over all motifs-pairs) (M&M). We found that
237 the *NonInteraction* model performance variability is positively correlated with both measures of motif
238 divergences (Spearman correlation=0.63, 0.67; p-value=1.2e-3, 6.3e-4 respectively).

239 For the *Interaction* model, the off-diagonal elements in Supplementary Table 7 show relatively high
240 cross-cell type performance accuracy, suggesting that the binding 'rules' are shared between cell types.
241 We ensured that the high cross-cell type prediction accuracy is not simply due to shared sequence
242 information, i.e., the genomic loci on which the model was trained in one cell type does not substantially
243 overlap with the loci tested in another cell type. Overall, across all TFs and all pairs of cell types, the
244 fractional overlap in genomic loci ranges from 0 to 10%, with a mean and median of ~4% (Fig 5D).  This
245 suggests that it is the binding rule, independent of specific sequence instances, that is shared across cell
246 types.

247 Furthermore, we found that when using the *Interaction* models, the cross-cell type accuracy is
248 symmetric (Spearman correlation of upper and lower triangle in Supplementary Table 7 is 0.68, p-value
249 9.5e-53). In other words, a high (respectively, low) accuracy in cell type *Y* using *EMT* trained on cell type
250 *X* implies a respectively high (respectively, low) accuracy in cell type *X* using the model learnt from cell
251 type *Y*. This further supports that the interaction-dependent (therefore genomic-context dependent)
252 binding rules are shared across cell types. In stark contrast, there is a lack of symmetry in cross-cell
253 prediction accuracy when *NonInteraction* model is used (Spearman correlation = 0.04, p-value 0.4).

254 In sum, our analyses suggest that the cell type-specific TF interactions play critical role in determining
255 cell type-specific *in vivo* binding. In addition to that, these revealed by *EMT* might be responsible for cell
256 specific binding of the reference motifs.

257

258 **_TRISECT reveals putative co-factors providing insights into cell-specific biological roles of a TF_**

259 Our results so far suggest that cell type-specific co-factors of a TF are a major driver of variability in the
260 *in vivo* binding rules across cell types. To further probe into the functional implications of cell type-
261 specific co-factors, for each reference TF, we identified its cell type-specific co-factors using the feature
262 importance of the corresponding motif as estimated by the model. To minimize redundancy, we
263 excluded motifs with substantially high co-occurrence frequency with at least one of the reference
264 motifs (M&M). To further minimize false positives, we assessed the enrichment of motif occurrence
265 near the cell-specific ChIP-Seq peaks of the reference TF relative to background and retained only those
266 putative co-factor motifs that were significantly enriched (odds ratio > 1.2 and p-value < 0.05, M&M).
267 The choice of enrichment odds ratio threshold is rationalized in Fig S7, which shows that increasing the
268 threshold would result in a loss of information for some TFs e.g. REST.

269 Several lines of evidence support the cell type-specific co-factors for a TF identified by *TRISECT*. First, we
270 found that for ~70% of the models, the putative co-factors are enriched for either heterodimerizing TFs
271 or for the TF family that the reference TF belongs to (M&M & Supplementary Table 8). The enrichment
272 of same family as that of reference TF is consistent with the fact that TFs forms dimer with other TFs
273 preferably from same family (Amoutzias et al. 2008; Dror et al. 2015). We also performed protein
274 domain enrichment analysis (Supplementary Table 9) using DAVID tool (Huang, Brad T. Sherman, et al.
275 2009; Huang, Brad T Sherman, et al. 2009), and found that more than 80% of enriched domains are
276 involved in homo- or hetero-dimerization consistent with Supplementary Table 8.

277 Second, we expect putative co-factors to be expressed at higher level in the specific cell types where
278 they are deemed as co-factors. For each co-factor (excluding ubiquitous co-factors), we determined the
279 log-fold difference in expression between the cell types where it is identified as co-factor relative to cell
280 types where it is not (M&M). The distribution of log fold changes of the co-factors are compared with a
281 control set of fold ratios as presented in Fig 6A. For most TFs, the co-factors show significantly higher
282 expression in the relevant cells. This is not true only in 5 cases. Among these, *CTCF* is known as cell type-
283 independent TF and for two of them (*GABPA* and *NRF1*) we show below, via an independence test, that
284 they show higher cell independence than other TFs.

285 Third, for each TF's cell type-specific co-factors, we performed biological processes GO term enrichment
286 analysis using the Gorilla tool (Eden et al. 2009) relative to all 981 motifs as the background. We found

287  significant differences in function among co-factors for a TF in different cell types. Remarkably, the
288  biological processes can vary across cell types while still being functionally related to the reference TF.
289  As an illustrative example, Fig 6B shows the enriched BP (false discovery rate <= 10%) for ATF3 in 4 cell
290  types. ATF3 is a stress-inducible TF involved in homeostasis (Allen-Jennings et al. 2001; Tanaka et al.
291  2011), specifically regulating cell-cycle, apoptosis, cell adhesion and signaling (Tanaka et al. 2011). We
292  found that ATF3 co-factors are enriched for functions related to cell cycle and proliferation in 3 out of 4
293  cell lines. In stem cell, the identified co-factors are involved in liver regeneration and inflammatory
294  response, consistent with previous studies showing direct link between ATF3 induction and liver injury
295  and regeneration in mice (Chen et al. 1996; Su et al. 2002). Furthermore, enrichment of NOTCH and
296  apoptotic signaling among co-factors in Hepg2 cell line is consistent with role of ATF3 in glucose
297  homeostasis and other primary functions of the liver (Allen-Jennings et al. 2001). Surprisingly, we find
298  enrichment of cognition, learning and memory among the co-factors in leukemia cell line. Since
299  leukemia is a cancerous cell line, non-native gene expression is not unexpected (Lotem et al. 2004;
300  Lotem et al. 2005). However, even though ATF3 is not known to play a direct role in neuronal function, a
301  closely functionally and structurally related protein CREB has well documented role in neuronal activity
302  and long-term memory formation in brain (Mayr & Montminy 2001), raising the possibility that either
303  ATF3 has a hitherto unknown role in cognition or, alternatively, the same set of co-factors are involved
304  in memory formation in conjunction with other TFs.

305  For other TFs, the enriched GO-terms at false discovery rate cutoff of 10% (enrichment scores ranges
306  from 1.22 to 93.75 with a median of 7.44) are listed in Supplementary Table 10 with corresponding
307  discussion based on literature survey is provided as Supplementary Notes. This can serve as a resource
308  for further investigation into cell type-specific binding and function of a broad array of TFs. In
309  Supplementary Tables 5a-b, we catalogue all the clusters with their specific TF interactions (M&M), and
310  their enriched GO terms.

311  We noted substantial variability in the number of detected co-factors across cell types for a TF.
312  Interestingly, a literature survey suggests that the cell types where the reference TF has specific
313  function, the number of co-factors in that cell type is comparatively higher. For example, REST has well-
314  known neuronal functions and its binding sites in neurons exhibit lack of cognate RE1 motifs (Rockowitz
315  et al. 2014), suggestive of dependence on co-factors. Consistently, Sknsh (brain cancer cell line) has
316  highest co-factor cardinality for REST. Similarly, JUN plays specific role in hematopoetic differentiation
317  and we found that Gm12878 (normal blood cell line) has the largest number of co-factors (Liebermann
318  et al. 1998). We reasoned that TF with greater cell type-specific roles would exhibit greater variability in
319  co-factor cardinality. For each TF we measured the variability of its co-factor cardinality across cell types.
320  As shown in Fig 7A, interestingly, TFs with ubiquitous and invariant roles such as TBP and CTCF have the
321  least variable co-factor cardinality.

322  We also assessed whether the difference in prediction accuracy achieved by *Interaction* model and the
323  *NonInteraction* model for a particular TF-cell type pair may reflect the TF's dependence on co-factors.
324  We measured the normalized distance between the performance (*performance distance*) of *Interaction*
325  and *NonInteraction* model (M&M) and compared it with co-factor cardinality. As shown in Fig 7C, we
326  found that the *performance distance* is positively correlated with co-factor cardinality (Spearman
327  correlation = 0.65, p-value = 2.7E-17).

328    Previous studies have found that the DNA sequence specificity of a TF can be influenced by interaction
329    with co-factors (Siggers et al. 2011; Slattery et al. 2011). Interestingly, a close inspection of the feature
330    importance estimated by the *NonInteraction EMT* model shows that in different cell types different
331    compositions of the reference motifs are utilized. Fig S8 presents all cell type-specific usage of a TF's
332    motifs; the cells where the motif usage is significantly different from expected usage are marked with
333    asterisk (M&M). Notably, such diverse usage is observed using *NonInteraction* models, suggesting cell
334    type-specific motif preference even without any modulation by the co-factors.

335    Taken together, the cell type-specific co-factors revealed by TRISECT are consistent with their cell type-
336    specific expression and function and may be critical in modulating a TF's cell type-specific biological
337    function.

338

339    **Discussion**

340    In this study, we have presented a novel ensemble-based framework –*TRISECT*, to investigate intra-cell
341    type heterogeneity of *in vivo* TF binding rules and inter-cell type commonality thereof. To the best of our
342    knowledge, this is the first study to show, based on a comprehensive analysis, that *in vivo* binding
343    specificity rule is composed of multiple components, or sub-models, many of which are shared across
344    multiple cell types. Tellingly, non-orthologous targets of binding sites across cell types governed by a
345    shared binding sub-model exhibit a greater functional and expression coherence than targets of binding
346    sites in the same cell type that are governed by different binding rules. For each TF, *TRISECT* identified
347    cell type-specific co-factors that are supported by gene expression data and literature studies supporting
348    their cell type-specific function. As a useful functional resource, for 23 TFs included in this study, we
349    provide a catalogue of clusters of shared sub-models, along with their putative cell type-specific targets,
350    the co-factors characterizing the cluster and their function.

351    Our ensemble model not only outperformed the previously reported sequence-based discriminative
352    model (*SVM-kmer*), but in several cases it outperformed the model that utilizes the chromatin
353    accessibility in addition to the sequence flanking the binding site (Arvey et al. 2012); paradoxically, some
354    of the TFs (e.g., JUND) whose *in vivo* binding were deemed to depend less on the sequence context and
355    more on the chromatin accessibility by the previous SVM approach were found to be adequately
356    modeled by sequence alone when using the EMT approach. Taken together with our observation that
357    these TFs depend on a large number of cell-type exclusive co-factors for their *in vivo* binding, these
358    results suggest that cell type-specific chromatin accessibility is captured, to some extent, by binding sites
359    for cell type-specific co-factors, shown independently by recent work (Whitaker et al. 2015; Benveniste
360    et al. 2014). Apart from the modeling approach of a TF's *in vivo* binding specificity, our study differs from
361    Arvey et al (Arvey et al. 2012) in several other aspects. In discussing cell type-specificity, the previous
362    study compared the models only in two cell types – GM12878 and K562, while we have investigated in-
363    depth the cell type-specificity of *TRISECT* across 4-12 cell types. While the previous work primarily
364    discusses cell type-specificity and ubiquity of their models, by clustering the cell type-specific sub-
365    models, our work investigates the extent of shared binding rules; cell type-specificity and ubiquity are
366    extreme cases thereof. In addition to cell type-specific variability in proximal co-factors, we investigated
367    in much greater depth than the previous work the cross-cell type variability in the preferred motif for
368    the reference TF. Together, these novel aspects of our study adds to the knowledge of sequence
369    information that specify a TF's *in vivo* binding in various cell types.

370    Another recent study (Dror et al. 2015) aiming to decipher the determinants of *in vivo* occupancy of a TF
371    showed that TF binding specificity is influenced by nearby homotypic sites (for the reference TF), the
372    local nucleotide composition, and certain DNA physical properties. Moreover, a preferred *in vivo* binding
373    in a homotypic cluster was shown to be related to a preferred nucleotide composition (GC-rich for zinc
374    finger TFs and AT-rich for homeodomain reference TFs) in the flanking region of the binding site. These
375    previous findings are consistent with the fact that the co-factors identified by *TRISECT* are enriched for
376    same family of TFs as the reference TF and thus have similar preference for nucleotide composition as
377    the reference TF. In the previous work (Dror et al. 2015), the accuracy in discriminating bound vs.
378    unbound sequences after controlling for the presence of a putative site for the reference TF was modest
379    (ROCAUC ~ 0.6). Whereas, we have shown that the motifs for the reference TF alone can discriminate
380    bound from the unbound control sites with ROCAUC ~ 0.78, suggesting that the reference TF are most
381    informative in determining in vivo binding, as also observed in Pique-Regi et al (Pique-Regi et al. 2011),
382    and the additional power of discriminations comes from the presence of co-factor motifs, as suggested
383    before (Arvey et al. 2012; Hannenhalli & Levy 2002), or from nucleotide composition and various DNA
384    physical properties (Dror et al. 2015). Interestingly, DNA flexibility measured by propeller twist (el
385    Hassan & Calladine 1996) is highly dependent on GC-content (Hancock et al. 2013), which in turn is
386    related to motif composition, as we have noted. Overall, these seemingly independent properties
387    (nucleotide composition and DNA physical properties on one hand and motif composition on the other)
388    may be related. Specific advantage of an ensemble model based on motif composition is that apart from
389    being highly accurate, it is functionally interpretable and provides insights into a TF's cell type-specific
390    functions.

391    Context-dependent function of a *cis* regulatory region requires binding of a specific combination of TFs.
392    This modularity contributes to morphological evolution through changes in cis elements controlling
393    transcription, while avoiding the pleiotropic effects of TF gene's expression change (Prud'homme et al.
394    2007). Shared sub-models of TF binding rules across cell types, as revealed by *TRISECT*, may suggest
395    shared history of cell types.

396    The ability of a TF to bind to diverse reference motifs and in conjunction, interact with diverse
397    combinations of co-factors serves to enhance its functional repertoire across contexts (Meijsing et al.
398    2009; Arvey et al. 2012). Our analyses indeed reveal cell type-specific preference for the reference motif
399    as well as the cell type-specific interaction partners of a TF. We found that the expression of cell type-
400    specific interaction partners to be higher in the cell types where they are expected to interact with the
401    TF and their function are consistent with the context based on the literature. Thus our study provides
402    further support for a TF's cell type-specific functions, and more importantly, enables further
403    investigation into the mechanisms underlying a TF's diverse cell-specific functions.

404

405    **Methods**

406    *Data Processing*

407    We downloaded the ChIP-Seq peaks or 23 TFs from ENCODE (Supplementary Table 1). For each TF we
408    selected only those cell lines for which narrow-peak data was available. We chose the more stringent of
409    the two criteria – top 5000 most significant peaks, or FDR q-values<0.2 to select binding sites (Arvey et
410    al. 2012). Relative to the center of ChIP-Seq peaks, the DNA regions of length 100bp were identified as

411  the foreground. As negative control, we sampled flanking regions of 100bp from 200bp away from the
412  positive sequences. Moreover, control sequences overlapping with any peak were excluded. Due to the
413  proximity of the negative examples, both foreground and background are expected to have similar GC-
414  composition (Arvey et al. 2012) and chromatin accessibility. However, we explicitly controlled for the GC
415  composition using sequence set balancing technique when comparing the foreground and the
416  background (Whitaker et al. 2015). We discarded any cell line resulting in fewer than 4000 sites.

417  ### *Learning EMT*

418  We considered three types of feature set for the sequence specificity model: (a) kmers - frequencies of
419  4096 6-mers in the 100bp sequence, (b) kmerRC - frequencies of 2080 6-kmer groups equating a k-mer
420  and its reverse complement, and (c) pwm*lk* – we take all the positional weight matrices (pwm) from
421  TRANSFAC 2011 as the features and get the motif hits using PWMSCAN (Levy & Hannenhalli 2002). The
422  feature value is the sum of pwm-score (-log10(hit score)) obtained from the PWMSCAN; we took the log
423  of feature values to compensate for the skewed distribution of the number of binding sites. Here, *lk*
424  refers to the PWM hit threshold (hit expected every *l* kb on average in the genome); we used *l* =
425  1/2/5/10kb.

426  We chose Adaptive boosting (Freidman 2008; Friedman 2002) as our composite model where each sub-
427  model within the ensemble is a decision tree and each decision tree is constructed based on a bootstrap
428  sample. We used the Adaboost framework implemented in R gbm package (Ridgeway 2015). In the
429  framework, Huber loss function is selected to reduce over-fitting. We estimated the classification
430  accuracy of the model based on 25% held out data set, while 75% data were being used to build each
431  tissue-specific model.

432  ### *Model conversion, Dudahart test and Hopkins statistics*

433  Each sub-model is represented by a point in a $d$-dimensional space. Each dimension denotes a feature
434  and the value along the dimension indicates the importance of the feature for the sub-model.
435  Therefore, each model (consisting of multiple sub-models) can be represented as a set of points in an $n$-
436  dimensional space where $n \leq 981$. For a model, the feature importance was measured based on the
437  prediction performance improvement by evaluating predictions on an out-of-bag samples. We modified
438  the gbm package (Ridgeway 2015) implementation of feature-importance to accommodate the
439  calculation for single tree or the sub-model in question. In other words, we determined the contribution
440  of a single tree (sub-model) in prediction performance improvement using the same out-of-bag samples.
441  We disregard the features which do not contribute to any sub-model. We measured dh-ratio (ratio of
442  within-cluster sum of clusters and overall sum of squares) for all cluster pairs, based on either cell type-
443  specific set of sub-models, or the pooled set of sub-models across all cell types for a TF. While
444  calculating dh-ratio, K-nearest neighborhood (KNN) approach was used for clustering. Since the final
445  output of KNN depends on initial random set of centers, the dh-ratio calculation was repeated 1000
446  times to ascertain robustness. We noted that all test results were significant (p-value < 0.01).

447  To measure Hopkins statistics (H) the sub-models are again represented as a set of points. H is defined
448  by the following.

$$H = \frac{\sum_{j=1...m} U_j^d}{\sum_{j=1...m} U_j^d + \sum_{j=1...m} W_j^d}$$

449   $W_j$ are the nearest-neighbor distances of $m$ randomly chosen points (sub-models), which demarcate the
450   sampling window. $U_j$ are the minimum distances of the sub-models from $m$ random points in the
451   sampling window. To define the sampling window, we either took 25 to 75 percentile of the feature
452   values or from δ to max.value-δ along each dimension, where δ denotes the standard deviation of the
453   feature value (Dubes & Zeng 1987; Zeng & Richard C Dubes 1985; Zeng & Richard C. Dubes 1985). To
454   estimate p-value, we repeat the above procedure 1000 times and measured the H value. The p-values
455   ranges from 0.026 to less than 0.001.

456   ***Clustering sub-models***

457   For a TF, we obtained sub-models in all cell types, and then clustered all sub-models using K-nearest
458   neighbor (KNN), where each sub-model is an instance and the features of the instances are individual
459   feature-importance obtained in the context of respective tissue-specific model. Before feeding into the
460   KNN, we remove all the features whose cumulative importance over all sub-models is zero. The sub-
461   models are also clustered using XY-fused version of self-organizing map (Melssen et al. 2006) from
462   kohonen R package (Wehrens 2015). To make it comparable to KNN, we assumed 100% weight for X
463   map, i.e. sub-models will be clustered without preexisting label of which sub-models belonged to which
464   cell.

465   ***Assignment of sequences and target genes to the clusters***

466   A cluster of sub-models can be viewed as a new ensemble. We scored each binding site sequence
467   against each cluster, and a sequence is assigned to a cluster when it is scored above a threshold (of 1) by
468   the cluster. The choice of the threshold was based on the rationale that the *intercept* (Ridgeway 2015)
469   of tissue-specific models are ~1, and for a high-confidence positive sequence, the model-score should be
470   greater than the intercept. Each bound sequence (from all cell lines) is mapped to a set of clusters. For
471   each bound sequence, the nearest gene on the genome is considered to be its putative target, as per
472   convention (Zhu et al. 2010). Hence, each cluster corresponds to a set of target genes coming from
473   different tissues. We arranged the target genes into an $M$-by-$N$ array, where $M$ is the number of cell
474   lines and $N$ is the number of clusters. The enriched pathway among the target genes of each cluster was
475   determined using clusterProfiler R package (Yu et al. 2012).

476   ***Measuring functional and expression coherence using Fisher test***

477   We downloaded the KEGG pathways (www.genome.jp/kegg). We use the following contingency table to
478   determine whether the target genes from different cell lines that are assigned to the same cluster are
479   more functionally related than the target genes coming from the same tissue but from different clusters.

| Gene pair across | | Cluster (Foreground) | Cell line (Background) |
|---|---|---|---|
| In same | Yes | $a$ | $c$ |
| Pathway? | No | $b$ | $d$ |

480

481   In the $M$-by-$N$ target gene array, we compared all gene-pairs along columns from different rows (same
482   cluster, different tissues) and the gene-pairs along rows from different columns (same tissue, different

483    cluster) as the background. Then we apply the Fisher exact test in a cluster-centric fashion by comparing
484    the fraction of foreground gene-pairs in the same pathway relative to the background.

485    Expression coherence tests were designed similarly, based on the following contingency table.

| Gene pair across | | Cluster (Foreground) | Tissue (Background) |
|---|---|---|---|
| Co-expressed? | Yes | a | c |
| | No | b | d |

486

487    A gene-pair is considered co-expressed if both of the genes are turned on (RNA-seq log2CPM > 1) in
488    their respective tissues; CPM stands for Counts per Million. CPM, instead of the standard FPKM measure
489    to quantify gene expression suffices for our purpose as we only compare a gene's expression across
490    samples, and not with other genes in the same sample. We showed similar trend of expression
491    coherence with different expression threshold (log2CPM>=5) (Fig S5).

492    ***Model variability, and Motif-divergenece***

493    Model variability is defined by its normalized-predictability across cell lines. For each model, n ROCAUC
494    values are obtained on held-out dataset of n cell-lines. Cross-ROCAUC values are normalized by self-
495    ROCAUC value. Mathematically, $var_{model_i} = \frac{\sum_{j \neq i, j \in cells} rocauc_j}{rocauc_i}$ .

496    Motif-divergence is defined by the following equation. $motif.div._{pwms} = \sum_{i,j \in pwms} \frac{dist_{i,j}}{IC_i + IC_j}$. Here,
497    $dist_{i,j} = 1/similarity_{i,j}$ and $IC_i$ is the information content of ith motif. Similarity between two pwms
498    is calculated following the normalized version of the sum of column correlations (Pietrokovski 1996).

499    ***Identification of co-factors***

500    EMT provides importance of all features in discriminating the foreground from the background. We
501    retained all features with nonzero importance. From the initial set, we removed any motif that has 60%
502    pwm-similarity (consensus overlap) for at least 50% of the binding site locations with any of the
503    reference motifs. Next, we calculated enrichment of the motif in the foreground binding sites relative to
504    control sites. We retained the motifs with greater than 1.2-fold enrichment and p-value <= 0.05. The
505    resulting motifs were considered as cofactor. For further analysis, we considered tissue specific
506    cofactors by removing common motifs across tissue. For *unique-relaxed* set we excluded co-factors that
507    are common across all cell-lines, and for *unique-strict* set co-factors common to any two cell lines were
508    excluded. The functional tissue-specificity measure for a TF is determined using the cardinality-
509    variability of unique-strict co-factors.

510    ***Gene expression and differential gene expression***

511    For gene expression, we used RNA-seq data downloaded from ENCODE (Supplementary Table 4). For
512    each tissue, we obtained between 2 and 4 RNA-seq samples depending on the availability and obtained
513    the number of reads aligned to the gene. We corrected for batch effect using ComBat tool (Leek &
514    Storey 2007). To estimate differential expression between two set of cell lines (those in which a TF is

515     deemed a co-factor, and those where it is not), we used linear model from R package, limma (Smyth
516     2005).

### Enrichment of same family TFs and heterodimerizing TFs

518     We collected the family name of each PWM and the list of heterodimerizing PWMs based on semi-
519     automated inspection of TRANSFAC 2011 annotations, based on keywords and further reading of the
520     description. For hyper-geometric test of family-enrichment, we compared how many co-factors belong
521     to the family of reference motifs relative to the 981 motifs. Heterodimer enrichment was tested
522     similarly.

### Cluster specific TF-interactions mapping

524     Cluster-specific co-factors are identified by treating a cluster as a new ensemble of sub-models. We
525     computed an aggregated relative importance of the features, considering the decision trees of the new
526     ensemble corresponding to a cluster. Since the set of decision tree has been changed from the original
527     set of trees from the EMT, some of the detected co-factors may be false positives. We took the
528     intersection of the features (with non-zero importance) with the 'enriched-nonoverlapped' (or 'distinct-
529     relaxed' or 'distinct-strict') co-factors of the original EMT. The corresponding enriched GO terms are
530     determined using a R package called clusterProfiler (Yu et al. 2012).

### Tissue-specific pwm for the reference TF

532     We obtained relative feature importance of the reference motifs from the *Noninteraction* models and
533     compared them with random expectation. To calculate the random expectation, 1000 *Noninteraction*
534     models are learned based on randomly sampled 4k sites from among all binding sites across cell-lines.
535     From 1000 models 1000 relative feature importance is calculated. Each set of relative importance is
536     assumed a point in p-dimensional space where p is the number of reference motifs. We considered the
537     relative importance vectors as data points from multivariate normal distribution and for each vector we
538     calculated the Mahalanobis distances from the centroid which follows a chi-square distribution (Slotani
539     1964). The degrees of freedom (d) for the chi-squared distribution is determined using maximum
540     likelihood estimate and a P-value is generated from a chi-square distribution function of d degrees of
541     freedom.

542

### Figures

544     **Figure 1**: (A). Schematic of *TRISECT* pipeline. Different colors represent different binding rules or sub-
545     models. Rows (a, b, c) represent cell types. Green, pink and yellow colors indicate cell type-specific sub-
546     models. Each EMT is represented by a bucket of sub-models (top right). Star denotes sub-models and
547     diamond denotes the corresponding data point after transformation into reduced feature space. The
548     sub-models across all cell types are clustered. Cyan is common between cell types *a* and *b*, light-brown is
549     common between cell types *b* and *c,* and purple is common across all three cell types. (B). Accuracy
550     (ROCAUC) distribution for 6 choices of feature sets for EMT. (C). Comparison of accuracy between all
551     pairs of 6 feature-set choices. Nodes are labeled with feature type and mean accuracy. Directional edges
552     are labeled with Wilcoxon p-value.

553    **Figure 2**: Prediction accuracy comparison of EMT against svm-kmer and svm trained using both kmer
554    and DNase (kmer+DNase), where (A) EMT is trained using kmerRC features, and (B) EMT is trained using
555    pwm hits with 1kb stringency (*pwm1k*). Each point represent a TF. Except for 3~4 TFs (blue), EMT
556    outperform svm in all other cases. For some TFs (green), sequence based EMT outperforms
557    sequence+chromatin based model as well.

558    **Figure 3**: Assessing the existence of sub-models shared across cell types. (A&B). Cluster membership
559    matrix using k-nearest neighbor clustering. Each row represents a cluster and column represents a cell
560    type. Each element in the matrix denotes the number of sub-models in the cluster coming from each cell
561    type. Some clusters consist of sub-models from multiple cells (cluster#20 in B), while some other consist
562    of sub-models from a single cell type (cluster#6 in A). (C&D). Boxplot of *dh-ratio* and *Hopkins statistic* for
563    135 TF-cell pairs based on sub-models of TF-cell type pair, and pooling all sub-models for each TF.

564    **Figure 4**: Functional and Expression coherence of sub-model clusters. (A&B) Fraction of multi-cell
565    clusters found to be coherent using k-nearest neighbor (KNN) and XY-Fused (XYF) self-organizing map
566    respectively. Mapped.targets denotes when genes are assigned to cluster based on *TRISECT* pipeline,
567    random.targets indicates the clusters consisting of random genes among all targets and random.genes
568    indicates the cluster consisting of random genes.

569    **Figure 5**: Association between number of interaction partners and model-accuracy. (A) The trend of
570    model accuracy with increasing sequence size for TF ZNF143 (selected arbitarily for illustration). (B).
571    Comparison of cross-validation prediction accuracy for *Interaction* and *Noninteraction* models. (C).
572    Comparison of model variability  in log scale (cross-cell type performance variability) for *Interaction* and
573    *Noninteraction* models. (D). Distribution of the fraction of test sequences that fall in one of the four
574    categories: Overlapped_true (respectively, overlapped_false) denotes the  correctly (respectively,
575    incorrectly) classified sequences having at least 50% overlap between the training sequences in one cell
576    type and the test sequences in another cell type.  Nonoverlapped_true (respectively,
577    nonoverlapped_false) denotes correctly (respectively, incorrectly) classified sequences that do not
578    overlap with any sequence in the training set.

579    **Figure 6**: Functional validation of putative co-factors. (A). Identified co-factors have higher expression in
580    the cell lines they are detected in. For a TF motif detected as a co-factor in *n* cell lines, and not in
581    another *m* cell lines, we calculated fold difference in the TF's expression between the two sets of cell
582    lines. Each boxplot corresponds to all co-factors of a TF in X-axis. (B). As an example, for ATF3, GO
583    enrichment analysis of co-factors in four cell types recapitulate the known cell type-specific biological
584    roles.

585    **Figure 7**: EMT model heterogeneity is associated with cell type-specificity of co-factors. (A) The plot
586    shows for each TF the variability of co-factor cardinality across cell types. Each point is further labeled
587    with cell type where the relevant TF has specific usage, based on literature and has largest number of
588    co-factors. TBP and CTCF are the most ubiquitous TFs. (B) Normalized ROCAUC difference of *Interaction*
589    and *NonInteraction* models for a specific TF-cell type pair correlates with co-factor cardinality. (C-D)
590    Cross-cell type variability in motif usage for the reference TF in the *NonInteraction* model, for JUN and
591    TBP as two extreme examples. JUN shows different binding specificity in different cell types, while TBP
592    does not.

593

594    **Supplemental information**

595    Supplemental Figures, S1-S8

596    Supplementary Tables, 1-10

597    Supplementary Notes, 1-3

598

599    **Disclosure Declaration**

600    None

601

602    **Authors' contributions**

603    S.H. conceived the project. S.H. and M.S. designed the analyses in consultation with H.C.B. M.S.
604    performed the analyses. S.H. and M.S. wrote the manuscript with help from H.C.B.

605

606    **Acknowledgements**

610

611    **References**

612    Allen-Jennings, A.E. et al., 2001. The roles of ATF3 in glucose homeostasis. A transgenic mouse
613          model with liver dysfunction and defects in endocrine pancreas. *The Journal of biological*
614          *chemistry*, 276(31), pp.29507–29514.

615    Amoutzias, G.D. et al., 2008. Choose your partners: dimerization in eukaryotic transcription
616          factors. *Trends in Biochemical Sciences*, 33(5), pp.220–229.

617    Arvey, A. et al., 2012. Sequence and chromatin determinants of cell-type-specific transcription
618          factor binding. *Genome Research*, 22(9), pp.1723–1734.

619    Benveniste, D. et al., 2014. Transcription factor binding predicts histone modifications in human
620          cell lines. *Proceedings of the National Academy of Sciences of the United States of America*,
621          111(37), pp.13367–13372.

622    Bulyk, M.L., Johnson, P.L.F. & Church, G.M., 2002. Nucleotides of transcription factor binding
623          sites exert interdependent effects on the binding affinities of transcription factors. *Nucleic*
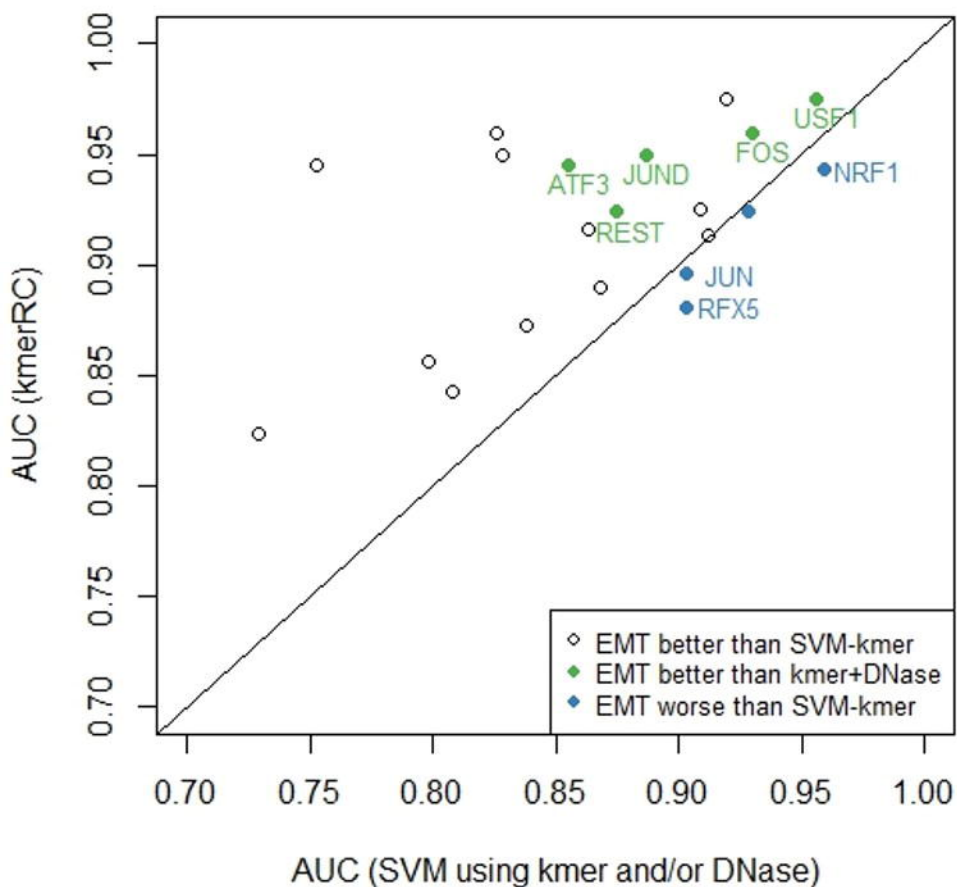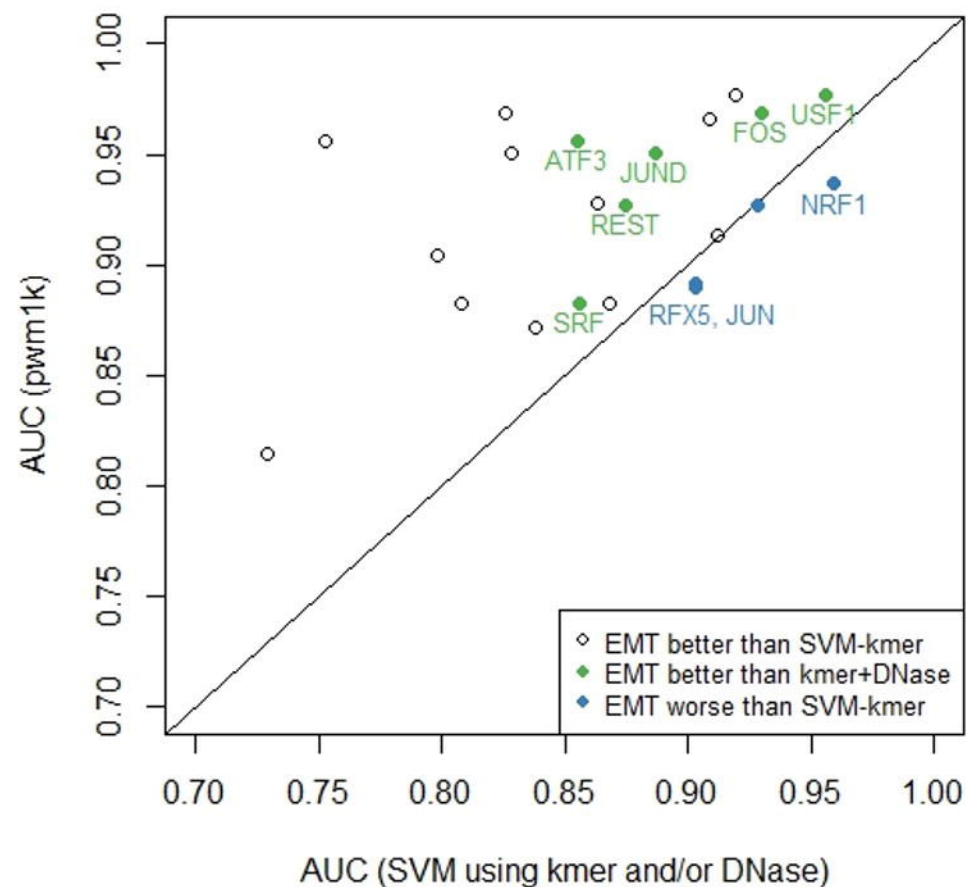624          *acids research*, 30(5), pp.1255–1261.

625    Busby, S. & Ebright, R.H., 1994. Promoter structure, promoter recognition, and transcription
626         activation in prokaryotes. *Cell*, 79(5), pp.743–746.

627    Chen, B.P., Wolfgang, C.D. & Hai, T., 1996. Analysis of ATF3, a transcription factor induced
628         by physiological stresses and modulated by gadd153/Chop10. *Molecular and cellular*
629         *biology*, 16(3), pp.1157–1168.

630    Dror, I. et al., 2015. A widespread role of the motif environment in transcription factor binding
631         across diverse protein families. *Genome research*.

632    Dubes, R.C. & Zeng, G., 1987. A test for spatial homogeneity in cluster analysis. *Journal of*
633         *Classification*, 4(1), pp.33–56.

634    Duda, R., Hart, P. & Stork, D., 2001. Pattern Classification. *New York: John Wiley, Section*,
635         p.680.

636    Eden, E. et al., 2009. GOrilla: a tool for discovery and visualization of enriched GO terms in
637         ranked gene lists. *BMC bioinformatics*, 10, p.48.

638    Freidman, J.H., 2008. Greedy Function Approximation□: A Gradient Boosting Machine Author
639         ( s ): Jerome H . Friedman Source□: The Annals of Statistics , Vol . 29 , No . 5 ( Oct ., 2001
640         ), pp . 1189-1232 Published by□: Institute of Mathematical Statistics Stable URL□:
641         http://www. *Institue of Mathematical Statistics*, 29(5), pp.1189–1232.

642    Friedman, J., Hastie, T. & Tibshirani, R., 2000. Additive Logistic Regression: a Statistical View
643         of Boosting. *The Annals of Statistics*, 28(2), pp.337–407.

644    Friedman, J.H., 2002. Stochastic gradient boosting. *Computational Statistics and Data Analysis*,
645         38(4), pp.367–378.

646    Frietze, S. & Farnham, P.J., 2011. Transcription factor effector domains. *Sub-cellular*
647         *biochemistry*, 52, pp.261–277.

648    Gheldof, N. et al., 2010. Cell-type-specific long-range looping interactions identify distant
649         regulatory elements of the CFTR gene. *Nucleic Acids Research*, 38(13), pp.4325–4336.

650    Hancock, S.P. et al., 2013. Control of DNA minor groove width and Fis protein binding by the
651         purine 2-amino group. *Nucleic acids research*, 41(13), pp.6750–6760.

652    Hannenhalli, S., 2008. Eukaryotic transcription factor binding sites--modeling and integrative
653         search methods. *Bioinformatics (Oxford, England)*, 24(11), pp.1325–1331.

654    Hannenhalli, S. & Levy, S., 2002. Predicting transcription factor synergism. *Nucleic acids*
655         *research*, 30(19), pp.4278–4284.

656    El Hassan, M.A. & Calladine, C.R., 1996. Propeller-twisting of base-pairs and the
657        conformational mobility of dinucleotide steps in DNA. *Journal of molecular biology*,
658        259(1), pp.95–103.

659    Heintzman, N.D. et al., 2009. Histone modifications at human enhancers reflect global cell-type-
660        specific gene expression. *Nature*, 459(7243), pp.108–112.

661    Huang, D.W., Sherman, B.T. & Lempicki, R.A., 2009. Bioinformatics enrichment tools: Paths
662        toward the comprehensive functional analysis of large gene lists. *Nucleic Acids Research*,
663        37(1), pp.1–13.

664    Huang, D.W., Sherman, B.T. & Lempicki, R.A., 2009. Systematic and integrative analysis of
665        large gene lists using DAVID bioinformatics resources. *Nature protocols*, 4(1), pp.44–57.

666    JACOB, F. & MONOD, J., 1961. Genetic regulatory mechanisms in the synthesis of proteins.
667        *Journal of molecular biology*, 3, pp.318–356.

668    Jain, A.K. & Dubes, R.C., 1988. *Algorithms for Clustering Data*,

669    Leek, J.T. & Storey, J.D., 2007. Capturing heterogeneity in gene expression studies by surrogate
670        variable analysis. *PLoS Genetics*, 3(9), pp.1724–1735.

671    Levy, S. & Hannenhalli, S., 2002. Identification of transcription factor binding sites in the human
672        genome sequence. *Mammalian genome : official journal of the International Mammalian
673        Genome Society*, 13(9), pp.510–514.

674    Liebermann, D.A., Gregory, B. & Huffman, B., 1998. AP-1 (Fos/Jun) transcription factors in
675        hematopoietic differentiation and apoptosis (Review). *International Journal of Oncology*,
676        12(3), pp.685–700.

677    Lotem, J. et al., 2005. Human cancers overexpress genes that are specific to a variety of normal
678        human tissues. *Proceedings of the National Academy of Sciences of the United States of
679        America*, 102(51), pp.18556–18561.

680    Lotem, J. et al., 2004. Induction in myeloid leukemic cells of genes that are expressed in
681        different normal tissues. *Proceedings of the National Academy of Sciences of the United
682        States of America*, 101(45), pp.16022–16027.

683    Mayr, B. & Montminy, M., 2001. Transcriptional regulation by the phosphorylation-dependent
684        factor CREB. *Nature reviews. Molecular cell biology*, 2(8), pp.599–609.

685    Meijsing, S.H. et al., 2009. DNA binding site sequence directs glucocorticoid receptor structure
686        and activity. *Science (New York, N.Y.)*, 324(5925), pp.407–410.

687    Melssen, W., Wehrens, R. & Buydens, L., 2006. Supervised Kohonen networks for classification
688        problems. *Chemometrics and Intelligent Laboratory Systems*, 83(2), pp.99–113.

689  Pietrokovski, S., 1996. Searching databases of conserved sequence regions by aligning protein
690      multiple-alignments. *Nucleic Acids Research*, 24(19), pp.3836–3845.

691  Pique-Regi, R. et al., 2011. Accurate inference of transcription factor binding from DNA
692      sequence and chromatin accessibility data. *Genome Research*, 21(3), pp.447–455.

693  Prud'homme, B., Gompel, N. & Carroll, S.B., 2007. Emerging principles of regulatory
694      evolution. *Proceedings of the National Academy of Sciences of the United States of*
695      *America*, 104 Suppl , pp.8605–8612.

696  Ridgeway, G., 2015. Generalized Boosted Regression Models.

697  Rockowitz, S. et al., 2014. Comparison of REST Cistromes across Human Cell Types Reveals
698      Common and Context-Specific Functions. *PLoS Computational Biology*, 10(6).

699  Siggers, T. et al., 2011. Non-DNA-binding cofactors enhance DNA-binding specificity of a
700      transcriptional regulatory complex. *Molecular Systems Biology*, 7.

701  Slattery, M. et al., 2011. Cofactor binding evokes latent differences in DNA binding specificity
702      between hox proteins. *Cell*, 147(6), pp.1270–1282.

703  Slotani, M., 1964. Tolerance regions for a multivariate normal population. *Annals of the Institute*
704      *of Statistical Mathematics*, 16(1), pp.135–153.

705  Smyth, G., 2005. limma: Linear Models for Microarray Data. In R. Gentleman et al., eds.
706      *Bioinformatics and Computational Biology Solutions Using R and Bioconductor*. Springer-
707      Verlag, pp. 397–420. Available at: http://dx.doi.org/10.1007/0-387-29362-0_23.

708  Su, A.I. et al., 2002. Gene expression during the priming phase of liver regeneration after partial
709      hepatectomy in mice. *Proceedings of the National Academy of Sciences of the United States*
710      *of America*, 99(17), pp.11181–11186.

711  Tanaka, Y. et al., 2011. Systems analysis of ATF3 in stress response and cancer reveals opposing
712      effects on pro-apoptotic genes in p53 pathway. *PLoS ONE*, 6(10).

713  Wang, L., Jensen, S. & Hannenhalli, S., 2006. An interaction-dependent model for transcription
714      factor binding. *Systems Biology and Regulatory Genomics*, pp.225–234.

715  Wehrens, R., 2015. kohonen: Supervised and Unsupervised Self-Organising Maps.

716  Whitaker, J.W., Chen, Z. & Wang, W., 2015. Predicting the human epigenome from DNA
717      motifs. *Nature methods*, 12(3), pp.265–72, 7 p following 272.

718  Yáñez-Cuna, J.O. et al., 2012. Uncovering cis-regulatory sequence requirements for context-
719      specific transcription factor binding. *Genome Research*, 22(10), pp.2018–2030.

720  Yu, G. et al., 2012. clusterProfiler: an R Package for Comparing Biological Themes Among
721      Gene Clusters. *OMICS: A Journal of Integrative Biology*, 16(5), pp.284–287.

722  Zeng, G. & Dubes, R.C., 1985. A comparison of tests for randomness. *Pattern Recognition*,
723      18(2), pp.191–198. Available at:
724      http://www.sciencedirect.com/science/article/pii/0031320385900433 [Accessed August 17,
725      2015].

726  Zeng, G. & Dubes, R.C., 1985. A test for spatial randomness based on k-NN distances. *Pattern*
727      *Recognition Letters*, 3(2), pp.85–91. Available at:
728      http://www.sciencedirect.com/science/article/pii/0167865585900133 [Accessed August 17,
729      2015].

730  Zhu, L.J. et al., 2010. ChIPpeakAnno: a Bioconductor package to annotate ChIP-seq and ChIP-
731      chip data. *BMC bioinformatics*, 11, p.237.

732  Zinzen, R.P. et al., 2009. Combinatorial binding predicts spatio-temporal cis-regulatory activity.
733      *Nature*, 462(7269), pp.65–70.

734

**A**

AUC (kmerRC) vs AUC (SVM using kmer and/or DNase)

Labeled points: USF1, NRF1, FOS, JUND, ATF3, REST, JUN, RFX5

Legend:
- EMT better than SVM-kmer
- EMT better than kmer+DNase
- EMT worse than SVM-kmer

**B**

AUC (pwm1k) vs AUC (SVM using kmer and/or DNase)

Labeled points: USF1, FOS, NRF1, JUND, ATF3, REST, SRF, RFX5, JUN

Legend:
- EMT better than SVM-kmer
- EMT better than kmer+DNase
- EMT worse than SVM-kmer

A

Color Key

ATF3 (k=16)

B

Color Key

ATF3 (k=20)

C

D

**A**

## cofactor-cardinality - cell specific behavior of the TFs

**B**

## spearman corr. 0.65 with p-value 2.7e-17



**C**

## JUN

**D**

## TBP