

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17

**Using pattern recognition of epigenetic signals for supervised enhancer prediction**

Anurag Sethi<sup>1,2</sup>, Mengting Gu<sup>1</sup>, Landon Chan<sup>3</sup>, Koon-Kiu Yan<sup>1,2</sup>, Kevin Yip<sup>4</sup>, Joel Rozowsky<sup>1,2</sup>, and Mark Gerstein<sup>1,2,5</sup>

---

<sup>1</sup> Program in Computational Biology and Bioinformatics, Yale University, New Haven, Connecticut, United States of America

<sup>2</sup> Molecular Biophysics and Biochemistry, Yale University, New Haven, Connecticut, United States of America

<sup>3</sup> School of Medicine, The Chinese University Hong Kong, China

<sup>4</sup> Department of Computer Science, The Chinese University Hong Kong, China

<sup>5</sup> Department of Computer Science, Yale University, New Haven, Connecticut, United States of America

18 **Abstract**

19

20 Enhancers are important noncoding elements. Unfortunately, until recently, they were  
21 difficult to characterize experimentally, and only a few mammalian enhancers were  
22 validated, making it difficult to train statistical models for their identification. Instead,  
23 postulated patterns of genomic features were used heuristically for identification.  
24 Recently, a large number of massively parallel assays for characterizing enhancers have  
25 been developed. Here, we use them to create shape-matching filters based on  
26 enhancer-associated metaprofiles in epigenetic features. We then combine different  
27 features with simple, linear models and predict enhancers in a supervised fashion. By  
28 cross-validating and testing our models, we show that they can be transferred without re-  
29 parameterization between cell lines and even between organisms. Finally, we predict  
30 enhancers in cell lines with many transcription-factor binding sites. In turn, this highlights  
31 distinct differences between the type of binding at enhancers and promoters, enabling  
32 the construction of a secondary model discriminating between these two.

33

34

## 35 Introduction

36

37 Enhancers are gene regulatory elements that activate expression of target genes from a  
38 distance [1]. Enhancers are turned on in a space and time-dependent manner  
39 contributing to the formation of a large assortment of cell-types with different  
40 morphologies and functions even though each cell in an organism contains a nearly  
41 identical genome [2-4]. Moreover, changes in the sequences of regulatory elements are  
42 thought to play a significant role in the evolution of species[5-9]. Understanding  
43 enhancer function and evolution is currently an area of great interest because variants  
44 within distal regulatory elements are also associated with various traits and diseases  
45 during genome-wide association studies [10-12]. However, the vast majority of  
46 enhancers and their spatiotemporal activities remain unknown because it is not easy to  
47 predict their activity based on DNA sequence or chromatin state [13, 14].

48 Traditionally, the regulatory activity of enhancers and promoters were experimentally  
49 validated in a non-native context using low throughput heterologous reporter constructs  
50 leading to a small number of validated enhancers that function in the same mammalian  
51 cell-type [15, 16]. In addition to the small numbers, the validated enhancers were  
52 typically selected based on conserved noncoding regions [17] with particular patterns of  
53 chromatin [18], transcription-factor binding, [19] or noncoding transcription [20]. The  
54 small number and biases within the validated enhancers make them inappropriate for  
55 parameterizing tissue-specific enhancer prediction models [16]. As a result, most  
56 theoretical methods to predict enhancers could not optimally parameterize their models  
57 using a gold standard set of functional elements. Instead, most of these models were  
58 parameterized based on certain heuristic features associated with enhancers, which  
59 were then utilized to predict enhancers [19, 21-30]. For example, two of the widest used  
60 methods for predicting enhancers were based on the fact that these elements are  
61 expected to contain a cluster of transcription factor binding sites [24] and their activity is  
62 often correlated with an enrichment of certain post-translational modifications on histone  
63 proteins [27, 30]. These predictions were not rigorously assessed as very few putative  
64 enhancers could be validated experimentally and it remains challenging to assess the  
65 performance of different methods for enhancer prediction.

66

67 In recent times, due to the advent of next generation sequencing, a number of  
68 transfection and transduction-based assays were developed to experimentally test the  
69 regulatory activity of thousands of regions simultaneously in a massively parallel fashion  
70 [31-37]. In these experiments, several plasmids that each contains a single core  
71 promoter upstream of a luciferase or GFP gene are transfected or transduced into cells.  
72 These plasmids are used to test the regulatory activity of different regions by placing one  
73 region near the core promoter in each plasmid as differences in the gene's expression  
74 occur due to the differences in the activity of the tested region. STARR-seq was one  
75 such MPRA that was used to test the regulatory activity of the fly genome in several cell-  
76 types [31, 38] and was used to identify thousands of cell-type specific enhancers and  
77 promoters. MPRA have confirmed that active enhancers and promoters tend to be  
78 depleted of histone proteins and contain accessible DNA on which various transcription  
79 factors and cofactors bind [39, 40]. These regulatory regions also tend to be flanked by  
80 nucleosomes that contain histone proteins with certain characteristic post-translational  
81 modifications. These attributes lead to an enriched peak-trough-peak ("double peak")  
82 signal in different ChIP-Seq experiments for various histone modifications such as  
83 acetylation on H3K27 and methylations on H3K4. The troughs in the double peak ChIP-

84 seq signal represent the accessible DNA that leads to a peak in the DNase-I  
85 hypersensitivity (DHS) at the enhancer [41]. However, the optimal method to combine  
86 information from multiple epigenetic marks to make cell-type specific regulatory  
87 predictions remains unknown. For the first time, using data from several MPRA, we  
88 have the ability to properly train our models based on a large number of experimentally  
89 validated enhancers and test the performance of different models for enhancer  
90 prediction using cross validation.

91  
92 We developed a new supervised machine-learning method that was trained and tested  
93 on large number of experimentally active regulatory regions identified in MPRA to  
94 accurately predict active enhancers and promoters in a cell-type specific manner. Unlike  
95 previous prediction methods that focused on the enrichment (or signal) of different  
96 epigenetic datasets, we developed a method to also take into account the enhancer-  
97 associated pattern within different epigenetic signals. As the epigenetic signal around  
98 each enhancer is noisy, we aggregated the signal around thousands of enhancers  
99 identified using MPRA to increase the signal-to-noise ratio and identified the shape  
100 associated with active regulatory regions. The epigenetic signal shapes associated with  
101 promoters and enhancers are conserved across millions of years of evolution and these  
102 models can be used to predict enhancers and promoters in different cell-types and  
103 tissues and across diverse eukaryotic species. We further created simple to use  
104 transferrable statistical models with six parameters that can be used to predict  
105 enhancers and promoters in several eukaryotic species including fly, mouse, and  
106 human. We applied these models to predict active enhancers and promoters in the H1-  
107 human embryonic stem cell (H1-hESC), a highly studied human cell-line in the ENCODE  
108 datasets. These analyses show that the pattern of transcription factor (TF) binding and  
109 co-binding varies between enhancers and promoters. The pattern of TF and co-TF  
110 binding at active enhancers is much more heterogeneous than the corresponding  
111 patterns on promoters. The pattern of TF binding can be used to distinguish enhancers  
112 from promoters with high accuracy. Thus, our methods provide a framework that utilizes  
113 different epigenetic genomics datasets to predict active regulatory regions in a cell-type  
114 specific manner and then utilizes further functional genomics datasets to identify key TFs  
115 associated with active regulatory regions within these cell-types.

## 116 117 **Results**

### 118 119 **Aggregation of epigenetic signal to create metaprofile:**

120  
121 We developed a framework to predict activating regulatory elements utilizing the  
122 epigenetic signal patterns associated with experimentally validated promoters and  
123 enhancers [31]. We aggregated the signal of histone modifications on MPRA peaks to  
124 remove noise in the signal and created a metaprofile of the double peak signals of  
125 histone modifications flanking enhancers and promoters. MPRA peaks typically consist  
126 of a mixture of enhancers and promoters, and at this stage, we do not differentiate  
127 between the two sets of regulatory elements. These metaprofiles were then utilized in a  
128 pattern recognition algorithm for predicting active promoters and enhancers in a cell-type  
129 specific manner.

130  
131 These metaprofiles were initially created using the histone modification H3K27ac at  
132 active STARR-seq peaks (see Figure 1 and Methods) identified in the S2 cell-line of fly.  
133 Approximately 70% of the active STARR-seq peaks contain an easily identifiable double  
134 peak pattern even though there is a lot of variability in the distance between the two



135 maxima of the double peak in the ChIP-chip signal (Figure S1). Even though the  
136 minimum tends to occur in the center of these two maxima on average, the distance  
137 between the two maxima in the double peaks can vary between 300 and 1100 base  
138 pairs. During aggregation, we aligned the two maxima in the H3K27ac signal across  
139 different STARR-seq peaks, followed by interpolation and smoothing the signal before  
140 calculating the average metaprofile. In addition, an optional flipping step was performed  
141 to maintain the asymmetry in the underlying H3K27ac double peak because it may be  
142 associated with the directionality of transcription [42]. For the first time, we also  
143 calculated the dependent metaprofiles for thirty other histone marks and DHS signal by  
144 applying the same set of transformations to these datasets. The metaprofile for the  
145 histone marks associated with active regulatory regions were also double peak signals  
146 and the maxima across different histone modification signals tended to align with each  
147 other on average (Figure S2). This indicates that a large number of histone modifications  
148 tend to simultaneously co-occur on the nucleosomes flanking an active enhancer or  
149 promoter. In contrast, as expected, the DHS signal displayed a single peak at the center  
150 of the H3K27ac double peak (Figure 1). In addition, repressive marks such as  
151 H3K27me3 were depleted in these regions and the metaprofile for these regions did not  
152 contain a double peak signal (Figure S2).

153

#### 154 **Occurrence of metaprofile is predictive of regulatory activity:**

155

156 We evaluated whether these metaprofiles can be utilized to predict active promoters and  
157 enhancers using matched filters, a well-established algorithm in template recognition. A  
158 matched filter is the optimal pattern recognition algorithm that uses a shape-matching  
159 filter to recognize the occurrence of a template in the presence of stochastic noise [43].  
160 We evaluated whether the occurrence of the epigenetic metaprofiles identified for the  
161 histone marks and DHS can be used to predict active enhancers and promoters using  
162 receiver operating characteristic (ROC) and precision-recall (PR) curves. The PR curves  
163 are particularly useful to assess the performance of classifiers in skewed or imbalanced  
164 data sets in which one of the classes is observed much more frequently as compared to  
165 the other. On these imbalanced data sets, PR curves are useful alternative to ROC  
166 curves as the precision is directly related to the false detection ratio at different  
167 thresholds. The PR curve highlights differences in performance of different models even  
168 when their ROC curves remain comparable [44]. The matched filter score is higher in  
169 genomic regions where the template pattern occurs in the corresponding signal track  
170 while it is low when only noise is present in the signal (Figure 1). Due to the  
171 aforementioned variability in the double peak pattern, the H3K27ac signal track is  
172 scanned with multiple matched filters with templates that vary in width between the two  
173 maxima in the double peak and the highest matched filter score with these matched  
174 filters is used to rate the regulatory potential of this region (see Methods). The  
175 dependent profiles are then used on the same region with the matched filter to score the  
176 corresponding genomic tracks.

177

178 We used 10-fold cross validation to assess the performance of matched filters for  
179 individual histone marks to predict active STARR-seq peaks. In Figure 2, we observe  
180 that the H3K27ac matched filter is the single most accurate feature for predicting active  
181 regulatory regions (AUROC=0.92, AUPR=0.72) identified using STARR-seq. This is  
182 consistent with the literature as H3K27ac enriched peaks are often used to predict active  
183 promoters and enhancers [23, 45, 46]. In general, several histone acetylation (H3K27ac,  
184 H3K9ac, H4K12ac, H2BK5ac, H4K8ac, H4K5ac, H3K18ac) marks as well as the H1,  
185 H3K4me2, and DHS matched filters are the most accurate marks (see Figure 2 and

186 Table S1) because the matched filter scores for these regions on these marks are higher  
187 for STARR-seq peaks (Figure S3). The degree to which the matched filter scores for  
188 promoters and enhancers are higher than the matched filter scores for the rest of the  
189 genome is a measure of the signal to noise ratio for regulatory region prediction in the  
190 corresponding feature's genomic track and the larger the separation between positives  
191 and negatives, the greater the accuracy of the corresponding matched filter for  
192 predicting active regulatory regions. Interestingly, the distribution of matched filter scores  
193 for STARR-seq peaks are unimodal for each histone mark except for H3K4me1,  
194 H3K4me3, and H2Av, which are bimodal (Figure S3). We also show that the matched  
195 filter scores are more accurate for predicting active STARR-seq peaks than enrichment  
196 of signal alone as they outperform the histone peaks on ROC and PR curves (Figure  
197 S4).

198  
199 While a single STARR-seq experiment identifies thousands of active regulatory regions,  
200 these regions display core-promoter specificity and different sets of enhancers are  
201 identified when different core promoters are used in the same cell-type [47-51]. As we  
202 wanted to create a framework to predict all the enhancers and promoters active in a  
203 particular cell-type, we combined the peaks identified from multiple STARR-seq  
204 experiments in the S2 cell-type and reassessed the performance of the matched filters at  
205 predicting these regulatory regions. Merging the STARR-seq peaks from multiple core  
206 promoters in the S2 cell-type leads to higher AUROC and AUPR for the matched filters  
207 from most histone marks (Figure 2).

### 208 209 **Machine learning can combine matched filter scores from different epigenetic** 210 **features:**

211  
212 We combined the normalized matched filter scores (see Methods) from six different  
213 epigenetic marks (H3K27ac, H3K4me1, H3K4me2, H3K4me3, H3K9ac, and DHS)  
214 associated with active regulatory regions by the Roadmap Epigenomics Mapping [52]  
215 and the ENCODE [53] Consortia using a linear SVM [54] and the integrated model  
216 achieved a higher accuracy than the individual matched filter scores (Figure 2). We also  
217 assessed the performance of other statistical approaches for combining the features  
218 (including non-linear models) in Figure S6 and all these models performed similarly. By  
219 using only six features, we ensure that our model is capable of being applied to many  
220 cell-lines and tissues on which the relevant experiments have been performed. These  
221 models are trained to learn the patterns in the matched filter scores for different  
222 epigenetic marks within experimentally verified regulatory regions and we chose these  
223 marks as we wanted to assess the applicability of these machine learning models to  
224 predict active enhancers and promoters across different cell-types and species. As  
225 expected, the integrated models outperformed the individual matched filter scores, as  
226 they are able to leverage information from multiple epigenetic marks. In addition, the six-  
227 parameter integrated model displayed higher accuracy after combining the peaks  
228 identified using different core promoters. In the integrated model, the normalized  
229 matched filter score for each epigenetic feature in a particular region is scaled by its  
230 optimized weight and added together to form the discriminant function. The sign of the  
231 discriminant function is then used to predict whether the region is regulatory. The  
232 features with large positive and negative weights are predicted to be important for  
233 discriminating regulatory regions from non-regulatory regions in such models. They can  
234 also be used to measure the amount of non-redundant information added by each  
235 feature in the integrated model. According to the model, the acetylations (H3K27ac and  
236 H3K9ac) are the most important feature for predicting active regulatory regions from

237 inactive regions. While the DHS matched filter performed well as an individual feature  
238 (AUPR in Figure 2), the information in DHS is redundant with the information in the  
239 histone marks as indicated by the fact that it has the lowest weight among the six  
240 features in the integrated model. We compared several other machine learning  
241 algorithms including nonlinear SVM (results not shown) to combine the machine learning  
242 models and found that they all displayed nearly similar accuracy and similar features  
243 were more important across these different models (Figure S5).

244  
245 To assess the information contained in other epigenetic marks, we combined the  
246 matched filters from all 30 measured histone marks along with the DHS matched filter in  
247 separate statistical models (Figure S6) and these model displayed higher accuracy  
248 (AUROC=0.97, AUPR=0.93 for SVM model with multiple core promoters) than the 6  
249 feature model presented in Figure 2. The feature weights in this model indicated that  
250 H3K27ac contains the most information regarding the activity of regulatory regions.  
251 However, we found that a few other acetylations such as H2BK5ac, H4ac, and H4K12ac  
252 contain additional non-redundant information regarding the activity of these regulatory  
253 regions and might improve the accuracy of promoter and enhancer prediction from  
254 machine learning models (Figure S6).

### 255 256 **Distinct epigenetic signals associated with promoters and enhancers:**

257  
258 We proceeded to create individual metaprofiles and machine learning models for the two  
259 classes of regulatory activators – promoters (or proximal) and enhancers (or distal). We  
260 divided all the active STARR-seq peaks into promoters or enhancers based on their  
261 distance to the closest transcription start site (TSS) to delineate their likely function in the  
262 native context. Due to the conservative distance metric used in this study (1kb upstream  
263 and downstream of TSS in fly), the enhancers are regulatory elements that are not close  
264 to any known TSS even though a few of the promoters may actually function as  
265 enhancers. We then created metaprofiles of the different epigenetic marks on the  
266 promoters and enhancers and assessed the performance of the matched filters for  
267 predicting active regulatory regions within each category (Figure 3). The highest  
268 matched filter scores are typically observed on promoters and the matched filters for  
269 each of the six features tended to perform better for promoter prediction. The H3K27ac  
270 matched filter continues to outperform other epigenetic marks for predicting active  
271 promoters and enhancers (Figure 3). In addition, the DHS, H3K9ac, and H3K4me2  
272 matched filters also performed reasonably for promoter and enhancer prediction. Similar  
273 to previous studies [55, 56], we observed that the H3K4me1 metaprofile performs better  
274 for predicting enhancers while it is close to random for predicting promoters. In contrast,  
275 the H3K4me3 metaprofile can be utilized to predict promoters and not enhancers. The  
276 histogram for matched filter scores show that H3K4me1 matched filter score is higher  
277 near enhancers while the H3K4me3 matched filter score tends to be higher near  
278 promoters (Figure S7). The mixture of these two populations lead to bimodal  
279 distributions for H3K4me1 and H3K4me3 matched filter scores when calculated over all  
280 regulatory regions (Figure S3).

281  
282 We created two different integrated models to learn the combination of features  
283 associated with promoters and enhancers. These integrated models outperformed the  
284 individual matched filters at predicting active enhancers and promoters (Figures 3 and  
285 S8). In addition, the weights of the individual features identified the difference in roles of  
286 the H3K4me1 and H3K4me3 matched filter scores at discriminating active promoters  
287 and enhancers from inactive regions in the genome. The promoter-based (enhancer-

288 based) model performed much more poorly at predicting enhancers (promoters)  
289 indicating the unique properties of these regions (Figures S10 and S11). We also  
290 created two integrated models utilizing matched filter scores for all thirty histone marks  
291 as features for predicting enhancers and promoters. The additional histone marks  
292 provided independent information regarding the activity of promoters and enhancers as  
293 these features increased the accuracy of these models (Figure S9). The weights of  
294 different features indicate that H2BK5ac again displays the most independent  
295 information for accurately predicting active enhancers and promoters (Figures S9). We  
296 observe similar trends and accuracy with several different machine learning models  
297 (Figures S8 and S9).

### 299 **The epigenetic underpinnings of active regulatory regions are highly conserved in** 300 **evolution:**

302 In order to assess the transferability of these metaprofiles and machine learning models  
303 for predicting regulatory regions in other tissues and cell-types, we assessed the  
304 accuracy of these models for predicting regulatory elements identified using the  
305 transduction-based FIREWACH assay in mouse embryonic stem cells (mESC) [36]. The  
306 metaprofiles for individual histone marks learned using active promoters and enhancers  
307 identified with the STARR-seq assay in the S2 cell-line were used with matched filters to  
308 predict the regulatory activity of different regions in mESC based on the epigenetic  
309 signals in mESC (Figure 4). The matched filters for individual histone marks displayed  
310 similar accuracy for predicting enhancers and promoters in mESC as in the original S2  
311 cell-line. In addition, the 6-parameter SVM models learned using STARR-seq data in S2  
312 cell-line were also highly accurate at predicting active enhancers and promoters in  
313 mouse (Figure 4).

315 This indicates that the epigenetic profiles associated with active enhancers and  
316 promoters are conserved over 600 million years of evolution underscoring the  
317 importance of such epigenetic modifications in maintaining the regulatory role of  
318 enhancers and promoters across different cell-types and species. As these regulatory  
319 regions were identified using a single core promoter in FIREWACH, the performance of  
320 the different models in Figure 4 is probably underestimated. The accuracy of these  
321 models enables us to use the metaprofiles and statistical models learned using STARR-  
322 seq data in fly to predict enhancers in different cell-lines and eukaryotic species.  
323 Consistent with this, the metaprofile and machine learning models learned using  
324 STARR-seq experiment in BG3 cell-line (fly) can be utilized to predict active promoters  
325 and enhancers in the S2 cell-line (Figure S12).

### 327 **Different Transcription Factors bind to enhancers and promoters**

329 The ENCODE consortium has ChIP-Seq data for 60 transcription related factors in H1-  
330 hESC cell line, including a few chromatin remodelers and histone modification enzymes.  
331 Collectively we call all these transcription related factors “TF”s for simplicity. We utilized  
332 the 6 parameter integrated model to predict active enhancers and promoters in the  
333 hESC cell-line based on the epigenetic datasets measured by the ENCODE consortium  
334 to study the patterns of TF binding within enhancers and promoters. Using these  
335 models, we predicted 43463 active regulatory regions, of which 22828 (52.5%) are  
336 within 2kb of the TSS and are labeled as promoters. A large proportion of the predicted  
337 enhancers are found in the introns (30.41%) and intergenic regions (13.93%) (Figure

338 S13). The predicted promoters and enhancers are significantly closer to active genes  
339 than might be expected randomly (Figure S14).

340  
341 We further studied the differences in TF binding at promoters and enhancers (Figure 5  
342 and Figure S15). Most promoters and enhancers contain multiple TF-binding sites.  
343 However, the TF-binding of enhancers is more heterogeneous than promoters: in  
344 particular, more than 70% of the promoters bind to the same set of 2-3 sequence-  
345 specific TFs, which is not observed for enhancers. The majority of the promoters also  
346 contain peaks for several TATA-associated factors (TAF1, TAF7, and TBP). Overall, the  
347 high heterogeneity associated with enhancer TF-binding is consistent with the absence  
348 of a sequence code (or grammar) which can be utilized to easily identify active  
349 enhancers on a genome-wide fashion.

350  
351 In Figure 5, we show that the patterns of TF binding within regulatory regions can be  
352 utilized in a logistic regression model to distinguish active enhancers from promoters  
353 with high accuracy (AUPR = 0.89, AUROC = 0.87). We were also able to identify the  
354 most important features that distinguish promoters from enhancers. In addition to TATA-  
355 box associated factors such as TAF1, TAF7, and TBP, the RNA polymerase-II binding  
356 patterns as well as chromatin remodelers such as KDM4A and PHF8 are some of the  
357 most important factors that distinguish promoters from enhancers in H1-hESC. This  
358 provides a framework that can be utilized to identify the most important TFs associated  
359 with active enhancers and promoters in each cell-type.

360  
361 In Figure 5A, we show that the pattern of TF binding at promoters is different from that at  
362 enhancers and TF-binding at enhancers displaying more heterogeneity. As the set of  
363 TFs binding promoters is fairly uniform, the same pairs of TF also tend to bind together  
364 on promoters. In contrast, for enhancers, the patterns of TF co-binding is much more  
365 heterogeneous and different enhancers tend to contain different TF-pairs. This can be  
366 observed in the patterns of TF co-binding in Figures 5C and S16. These TF co-  
367 associations could lead to mechanistic insights of cooperativity between TFs. For  
368 example, similar to a previous study [57], CTCF and ZNF143 may function cooperatively  
369 as they are observed to co-occur frequently at distal regulatory regions in this study.

## 370 371 **Discussion**

372  
373 Our ability to accurately predict active enhancers in a cell-type specific manner using  
374 transferable supervised machine learning models that were trained based on regulatory  
375 regions identified using new NGS-enabled MPRA distinguishes our method from  
376 previous enhancer prediction methods. Currently, most existing methods were  
377 parameterized (not properly “trained”) with regions that had various features associated  
378 with promoters and enhancers and only a small number of these regions were typically  
379 tested for regulatory activity experimentally in an *ad hoc* manner. The MPRA were able  
380 to firmly establish that certain histone modifications occur on nucleosomes flanking  
381 active regulatory regions leading to the formation characteristic double peak pattern  
382 within the ChIP-signal [39]. This motivated us to create matched filter models that were  
383 able to identify these patterns within the shape of the ChIP-signal in the presence of  
384 stochastic noise with the highest signal to noise ratio. Furthermore, we were able to  
385 combine the matched filter scores from different epigenetic features using simple  
386 transferrable linear SVM models and learned the most informative epigenetic features  
387 for regulatory region predictions.

388

389 The sensitivity and selectivity of various MPRAs is currently a matter of debate. A  
390 majority of these MPRAs test the regulatory activity of different regions by assessing  
391 their ability to induce gene expression in a plasmid after transfecting it into a cell-type of  
392 interest [31]. Such assays may not recapitulate the native chromatin environment found  
393 in chromosomes, which may be necessary for assessing whether the regulatory region  
394 is active in its genomic environment.

395  
396 Here, we show for the first time, that the patterns in the epigenetic signals associated  
397 with active enhancers identified using a transfection-based assay (STARR-seq) can be  
398 utilized to predict the activity of enhancers in a transduction-based assay (FIREWACH).  
399 During the FIREWACH assay, random nucleosome-free regions in mESC were captured  
400 and assayed for regulatory activity of the GFP gene by utilizing a lentiviral plasmid vector  
401 and inserted (or transduced) these vectors into the chromosome in mESC cells. As the  
402 FIREWACH assay tests the regulatory activity of enhancers after transduction, we  
403 assume that these regions were tested in their native chromatin environment and  
404 transduction-based assays form a more stringent test for regulatory activity. However,  
405 due to the shorter length of the tested region (< 300 bp) and the single core promoter  
406 used in the FIREWACH assay, we think that the accuracy of the statistical models in  
407 Figure 4 is underestimated.

408  
409 We were able to assess the accuracy of different epigenetic metaprofiles for predicting  
410 regulatory activity using our statistical models. While different acetylation modifications  
411 are associated with active regions of the genome, we were able to compare close to 30  
412 histone marks for enhancer and promoter predictions. The H3K27ac matched filter  
413 remains the single most important feature for predicting active regulatory regions while  
414 H3K4me1 and H3K4me3 are known to distinguish promoters from enhancers. However,  
415 our analysis characterizes the amount of redundancy in information within the  
416 metaprofile of different epigenetic features for predicting active regulatory regions and  
417 shows that ChIP-experiments of H2BK5ac, H4ac, and H2A variants could also produce  
418 independent information that can improve the accuracy of promoter and enhancer  
419 predictions. In addition to these 30-feature models, we also provide a simple to use six-  
420 parameter SVM model for combining H3K27ac, H3K9ac, H3K4me1, H3K4me2,  
421 H3K4me3, and DHS to predict active promoters and enhancers in a cell-type specific  
422 manner. We also showed that the metaprofiles and the combination of epigenetic marks  
423 associated with active regulatory regions are highly conserved in evolution making these  
424 models highly transferable. These six histone marks have been measured for a number  
425 of different tissues and cell-types by the Roadmap Epigenomics Mapping Consortium  
426 [39], the ENCODE [53], and the modENCODE Consortium [58].

427  
428 One aspect that is discussed less frequently is the effect of core promoter on enhancer  
429 and promoter prediction. MPRAs show that the regulatory activity of enhancers and  
430 promoters in a regulatory assay depends on the core promoter used during the  
431 experiment [51]. As the transcription factors that bind to each regulatory region are  
432 thought to play a key role in core-promoter specificity [47, 51], we suspect that machine  
433 learning models that contain sequence or motif-based features may be biased towards  
434 certain transcription factor binding sites when trained with regulatory regions identified  
435 using a single-core promoter. To avoid such biases, it would be more appropriate to train  
436 models with sequence-based features when the validation experiments are performed  
437 with multiple core promoters. In the absence of validation data with multiple core  
438 promoters, it may be more suitable to train models using epigenetic features as such  
439 models contain no sequence-based information. In comparing the predictions from such

440 models with experiments using a single core promoter, some of the strongest predictions  
441 may be mislabeled as negatives even though they contain some regulatory activity  
442 leading to a lower accuracy estimate as shown in Figure 2.

443

444 As the epigenetic profiles and statistical models learned in this study are transferable  
445 across different cell-lines and species, we are able to apply these models to predict  
446 active enhancers and promoters in different cell-types. We applied these models to  
447 predict enhancers and promoters in H1-hESC, a highly studied ENCODE cell-line. This  
448 allowed us to analyze the differences in the patterns of TF binding at proximal and distal  
449 regulatory regions. The TF binding and co-binding patterns at enhancers is much more  
450 heterogeneous than that at promoters. We think that this heterogeneity in TF binding  
451 patterns makes it much more difficult to predict enhancers due to the absence of obvious  
452 sequence patterns in distal regulatory regions. However, we were also able to create  
453 highly accurate machine learning models that are able to distinguish proximal promoter  
454 regions from distal enhancers based on the patterns of TF ChIP-seq peaks within these  
455 regulatory regions. The conservation of the epigenetic underpinnings underlying active  
456 regulatory regions sets the stage for our method to study the evolution of tissue-specific  
457 enhancers and their genomic properties across different eukaryotic species.

458

459

460

## Figure Captions

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

**Figure 1: Creation of metaprofile.** A) We identified the “double peak” pattern in the H3K27ac signal close to STARR-seq peaks. The red triangles denote the position of the two maxima in the double peak. B) We aggregated the H3K27ac signal around these regions after aligning the flanking maxima, using interpolation and smoothing on the H3K27ac signal, and averaged the signal across different MPRA peaks to create the metaprofile in C). The exact same operations can be performed on other histone signals and DHS to create metaprofiles in other dependent epigenetic signals. D) Matched filters can be used to scan the histone and/or DHS datasets to identify the occurrence of the corresponding pattern in the genome. E) The matched filter scores are high in regions where the profile occurs (grey region shows an example) and it is low when only noise is present in the data. The individual matched filter scores from different epigenetic datasets can be combined using integrated model in F) to predict active promoters and enhancers in a genome wide fashion.

**Figure 2: Performance of matched filters and integrated models for predicting MPRA peaks.** The performance of the matched filters of different epigenetic marks and the integrated model for predicting all STARR-seq peaks is compared here using 10-fold cross validation. A) The area under the receiver-operating characteristic (AUROC) and the precision-recall (AUPR) curves are used to measure the accuracy of different matched filters and the integrated model. B) The weights of the different features in the integrated model are shown and these weights may be used as a proxy for the importance of each feature in the integrated model. C) The individual ROC and PR curves for each matched filter and the integrated model are shown. The performance of these features and the integrated model for predicting the STARR-seq peaks using multiple core promoters and single core promoter are compared. The numbers within the parentheses in A) refer to the AUROC and AUPR for predicting the peaks using a single STARR-seq core promoter while the numbers outside the parentheses refers to the performance of the model for predicting peaks from multiple core promoters.

**Figure 3: Performance of matched filters and integrated models for predicting promoters and enhancers.** The performance of the matched filters of different epigenetic marks and the integrated model for predicting active promoters and enhancers are compared here using 10-fold cross validation. A) The numbers within parentheses refer to the AUROC and AUPR for predicting promoters while the numbers outside parentheses refer the performance of the models for predicting enhancers. B) The weights of the different features in the integrated models for promoter and enhancer prediction are shown. C) The individual ROC and PR curves for each matched filter and the integrated model are shown. The performance of these features and the integrated model for predicting the active promoters and enhancers using multiple core promoters are compared.

**Figure 4: Conservation of epigenetic features.** The performance of the fly-based matched filters and the integrated model for predicting active promoters and enhancers in mouse embryonic stem cells identified using FIREWACH. A Similar to Figure 3, the numbers within parentheses refer to the AUROC and AUPR for predicting promoters while the numbers outside parentheses refer the performance of the models for predicting enhancers. B) The weights of the different features in the integrated models for promoter and enhancer prediction are shown. C) The individual ROC and PR curves



510 for each matched filter and the integrated model are shown. The performance of these  
511 features and the integrated model for predicting the active promoters and enhancers  
512 identified using FIREWACH are shown.

513

514 **Figure 5: Differences in TF binding patterns at enhancers and promoters.** A) The  
515 fraction of predicted promoters and enhancers that overlap with ENCODE ChIP-seq  
516 peaks for different TFs in H1-hESC are shown. The names of all TFs in the figure can be  
517 viewed in Figure S15. B) The AUROC and AUPR for a logistic regression model created  
518 from the pattern of TF binding at each regulatory region to distinguish enhancers from  
519 promoters are shown. The weight of each feature in the logistic regression model can be  
520 used to identify the most important TFs that distinguish enhancers from promoters. C)  
521 The patterns of TF co-binding at active promoters and enhancers are shown. The names  
522 of all the TFs in this graph can be viewed in Figure S16.  
523

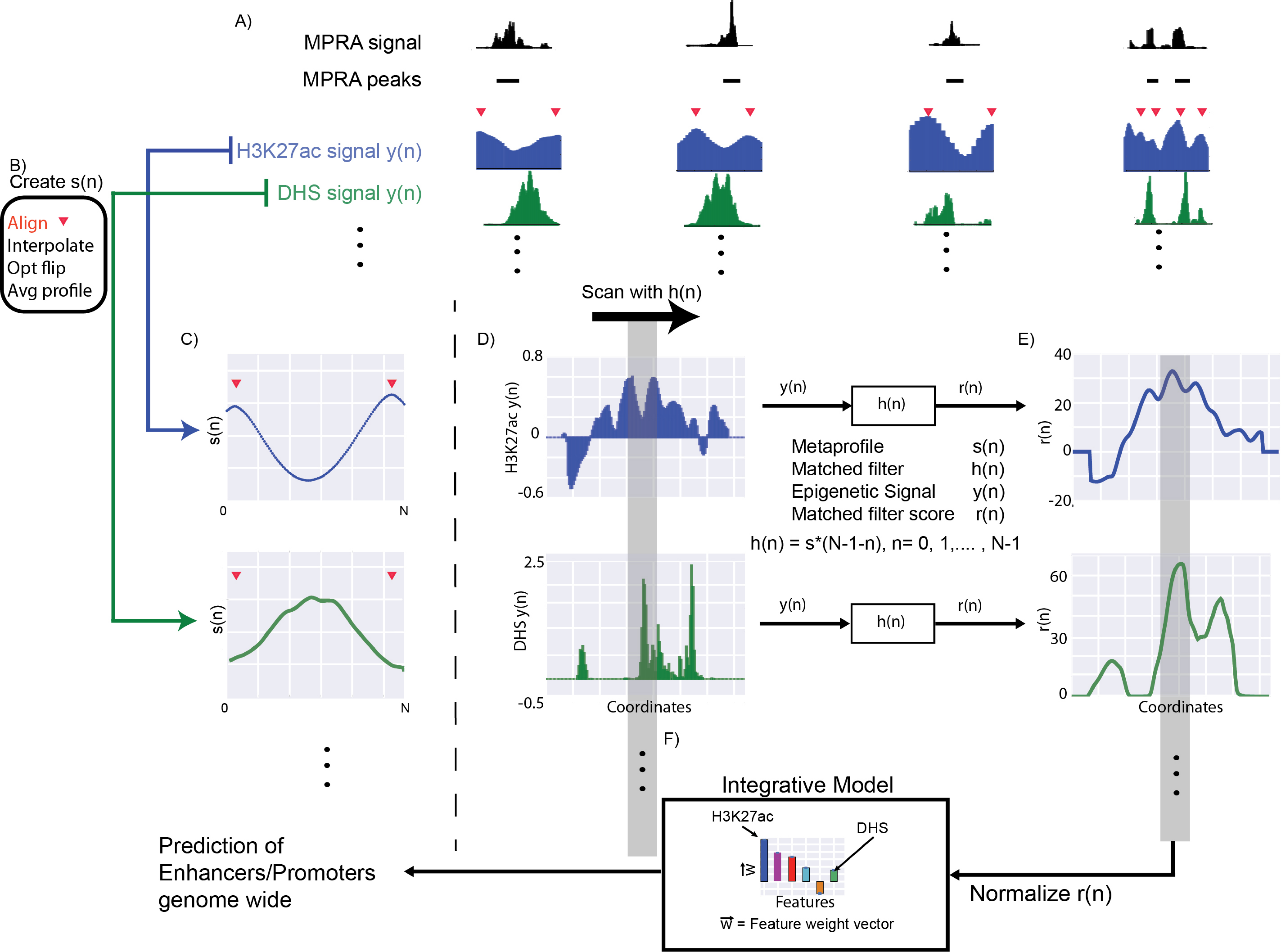
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568

## References:

1. Banerji, J., S. Rusconi, and W. Schaffner, *Expression of a beta-globin gene is enhanced by remote SV40 DNA sequences*. Cell, 1981. **27**(2 Pt 1): p. 299-308.
2. Ong, C.T. and V.G. Corces, *Enhancer function: new insights into the regulation of tissue-specific gene expression*. Nat Rev Genet, 2011. **12**(4): p. 283-93.
3. Woolfe, A., et al., *Highly conserved non-coding sequences are associated with vertebrate development*. PLoS Biol, 2005. **3**(1): p. e7.
4. Spitz, F. and E.E. Furlong, *Transcription factors: from enhancer binding to developmental control*. Nat Rev Genet, 2012. **13**(9): p. 613-26.
5. Cotney, J., et al., *The evolution of lineage-specific regulatory activities in the human embryonic limb*. Cell, 2013. **154**(1): p. 185-96.
6. Degner, J.F., et al., *DNase I sensitivity QTLs are a major determinant of human expression variation*. Nature, 2012. **482**(7385): p. 390-4.
7. Shibata, Y., et al., *Extensive evolutionary changes in regulatory element activity during human origins are associated with altered gene expression and positive selection*. PLoS Genet, 2012. **8**(6): p. e1002789.
8. Villar, D., et al., *Enhancer evolution across 20 mammalian species*. Cell, 2015. **160**(3): p. 554-66.
9. Xiao, S., et al., *Comparative epigenomic annotation of regulatory DNA*. Cell, 2012. **149**(6): p. 1381-92.
10. Wray, G.A., *The evolutionary significance of cis-regulatory mutations*. Nat Rev Genet, 2007. **8**(3): p. 206-16.
11. Corradin, O. and P.C. Scacheri, *Enhancer variants: evaluating functions in common disease*. Genome Med, 2014. **6**(10): p. 85.
12. Gusev, A., et al., *Partitioning heritability of regulatory and cell-type-specific variants across 11 common diseases*. Am J Hum Genet, 2014. **95**(5): p. 535-52.
13. Slattery, M., et al., *Absence of a simple code: how transcription factors read the genome*. Trends Biochem Sci, 2014. **39**(9): p. 381-99.
14. Levo, M., et al., *Unraveling determinants of transcription factor binding outside the core binding site*. Genome Res, 2015. **25**(7): p. 1018-29.
15. Pennacchio, L.A., et al., *Enhancers: five essential questions*. Nat Rev Genet, 2013. **14**(4): p. 288-95.
16. Erwin, G.D., et al., *Integrating diverse datasets improves developmental enhancer prediction*. PLoS Comput Biol, 2014. **10**(6): p. e1003677.
17. Pennacchio, L.A., et al., *In vivo enhancer analysis of human conserved non-coding sequences*. Nature, 2006. **444**(7118): p. 499-502.
18. Nord, A.S., et al., *Rapid and pervasive changes in genome-wide enhancer usage during mammalian development*. Cell, 2013. **155**(7): p. 1521-31.
19. Visel, A., et al., *ChIP-seq accurately predicts tissue-specific activity of enhancers*. Nature, 2009. **457**(7231): p. 854-8.
20. Andersson, R., et al., *An atlas of active enhancers across human cell types and tissues*. Nature, 2014. **507**(7493): p. 455-61.
21. Narlikar, L., et al., *Genome-wide discovery of human heart enhancers*. Genome Res, 2010. **20**(3): p. 381-92.

- 569 22. Visel, A., et al., *Ultraconservation identifies a small subset of extremely*  
570 *constrained developmental enhancers*. Nat Genet, 2008. **40**(2): p. 158-60.
- 571 23. Bonn, S., et al., *Tissue-specific analysis of chromatin state identifies temporal*  
572 *signatures of enhancer activity during embryonic development*. Nat Genet,  
573 2012. **44**(2): p. 148-56.
- 574 24. Yip, K.Y., et al., *Classification of human genomic regions based on*  
575 *experimentally determined binding sites of more than 100 transcription-*  
576 *related factors*. Genome Biol, 2012. **13**(9): p. R48.
- 577 25. Ghandi, M., et al., *Enhanced regulatory sequence prediction using gapped k-*  
578 *mer features*. PLoS Comput Biol, 2014. **10**(7): p. e1003711.
- 579 26. Heintzman, N.D., et al., *Distinct and predictive chromatin signatures of*  
580 *transcriptional promoters and enhancers in the human genome*. Nat Genet,  
581 2007. **39**(3): p. 311-8.
- 582 27. Hoffman, M.M., et al., *Unsupervised pattern discovery in human chromatin*  
583 *structure through genomic segmentation*. Nat Methods, 2012. **9**(5): p. 473-6.
- 584 28. Kharchenko, P.V., et al., *Comprehensive analysis of the chromatin landscape in*  
585 *Drosophila melanogaster*. Nature, 2011. **471**(7339): p. 480-5.
- 586 29. He, H.H., et al., *Nucleosome dynamics define transcriptional enhancers*. Nat  
587 Genet, 2010. **42**(4): p. 343-7.
- 588 30. Ernst, J., et al., *Mapping and analysis of chromatin state dynamics in nine*  
589 *human cell types*. Nature, 2011. **473**(7345): p. 43-9.
- 590 31. Arnold, C.D., et al., *Genome-wide quantitative enhancer activity maps identified*  
591 *by STARR-seq*. Science, 2013. **339**(6123): p. 1074-7.
- 592 32. Dickel, D.E., et al., *Function-based identification of mammalian enhancers*  
593 *using site-specific integration*. Nat Methods, 2014. **11**(5): p. 566-71.
- 594 33. Gisselbrecht, S.S., et al., *Highly parallel assays of tissue-specific enhancers in*  
595 *whole Drosophila embryos*. Nat Methods, 2013. **10**(8): p. 774-80.
- 596 34. Kwasnieski, J.C., et al., *High-throughput functional testing of ENCODE*  
597 *segmentation predictions*. Genome Res, 2014. **24**(10): p. 1595-602.
- 598 35. Melnikov, A., et al., *Systematic dissection and optimization of inducible*  
599 *enhancers in human cells using a massively parallel reporter assay*. Nat  
600 Biotechnol, 2012. **30**(3): p. 271-7.
- 601 36. Murtha, M., et al., *FIREWACH: high-throughput functional detection of*  
602 *transcriptional regulatory modules in mammalian cells*. Nat Methods, 2014.  
603 **11**(5): p. 559-65.
- 604 37. Patwardhan, R.P., et al., *Massively parallel functional dissection of mammalian*  
605 *enhancers in vivo*. Nat Biotechnol, 2012. **30**(3): p. 265-70.
- 606 38. Yanez-Cuna, J.O., et al., *Dissection of thousands of cell type-specific enhancers*  
607 *identifies dinucleotide repeat motifs as general enhancer features*. Genome  
608 Res, 2014. **24**(7): p. 1147-56.
- 609 39. Shlyueva, D., G. Stampfel, and A. Stark, *Transcriptional enhancers: from*  
610 *properties to genome-wide predictions*. Nat Rev Genet, 2014. **15**(4): p. 272-86.
- 611 40. Maston, G.A., et al., *Characterization of enhancer function from genome-wide*  
612 *analyses*. Annu Rev Genomics Hum Genet, 2012. **13**: p. 29-57.
- 613 41. Thurman, R.E., et al., *The accessible chromatin landscape of the human*  
614 *genome*. Nature, 2012. **489**(7414): p. 75-82.

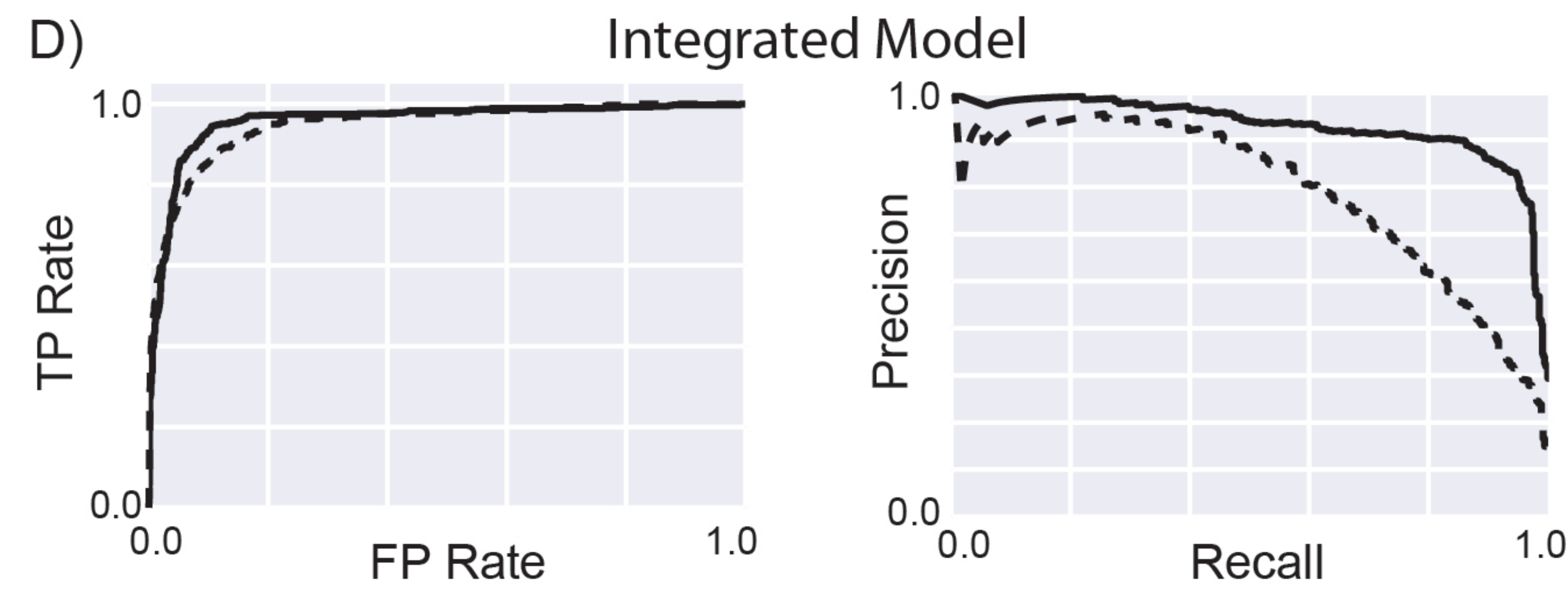
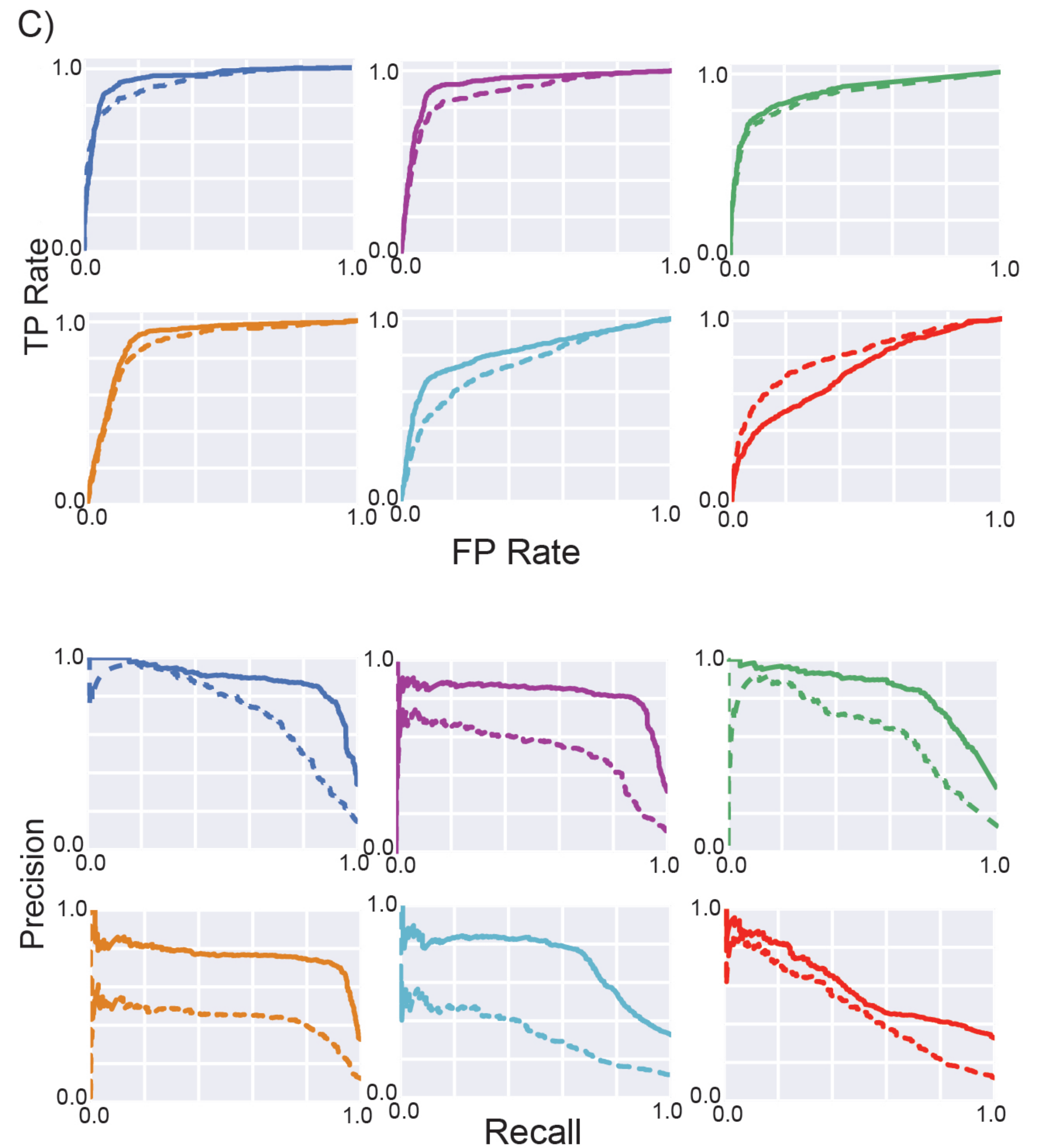
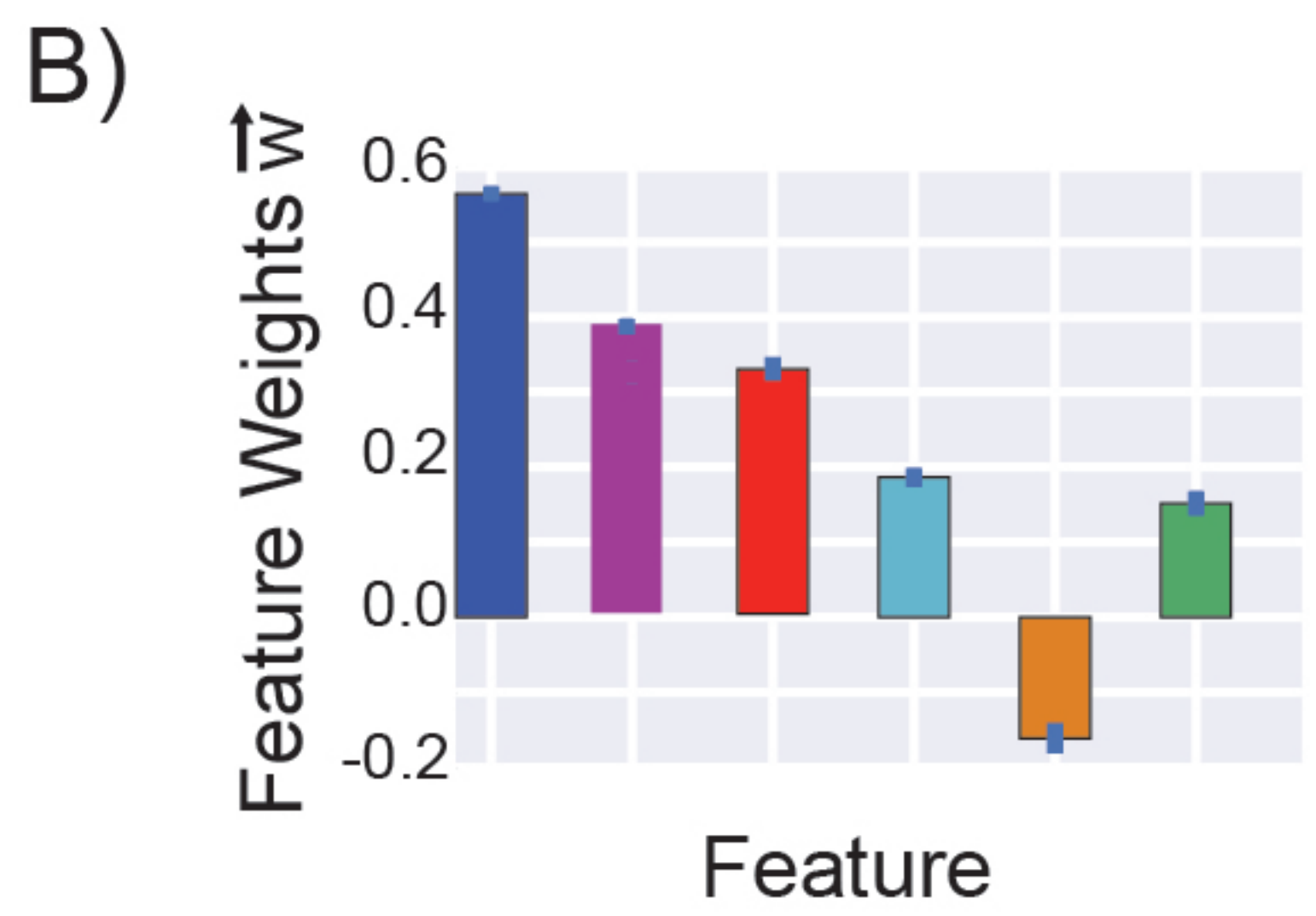
- 615 42. Kundaje, A., et al., *Ubiquitous heterogeneity and asymmetry of the chromatin*  
616 *environment at regulatory elements*. Genome Res, 2012. **22**(9): p. 1735-47.
- 617 43. Kumar, V.B.V.K., A. Mahalanobis, and R.D. Juday, *Correlation Pattern*  
618 *Recognition*. 2005.
- 619 44. Davis, J. and M. Goadrich, *The Relationship Between Precision-Recall and ROC*  
620 *Curves*. Proceedings of the 23rd international conference on Machine  
621 Learning, 2006: p. 233-240.
- 622 45. Creyghton, M.P., et al., *Histone H3K27ac separates active from poised*  
623 *enhancers and predicts developmental state*. Proc Natl Acad Sci U S A, 2010.  
624 **107**(50): p. 21931-6.
- 625 46. Rada-Iglesias, A., et al., *A unique chromatin signature uncovers early*  
626 *developmental enhancers in humans*. Nature, 2011. **470**(7333): p. 279-83.
- 627 47. Butler, J.E. and J.T. Kadonaga, *Enhancer-promoter specificity mediated by DPE*  
628 *or TATA core promoter motifs*. Genes Dev, 2001. **15**(19): p. 2515-9.
- 629 48. Li, X. and M. Noll, *Compatibility between enhancers and promoters determines*  
630 *the transcriptional specificity of gooseberry and gooseberry neuro in the*  
631 *Drosophila embryo*. EMBO J, 1994. **13**(2): p. 400-6.
- 632 49. Merli, C., et al., *Promoter specificity mediates the independent regulation of*  
633 *neighboring genes*. Genes Dev, 1996. **10**(10): p. 1260-70.
- 634 50. Ohtsuki, S., M. Levine, and H.N. Cai, *Different core promoters possess distinct*  
635 *regulatory activities in the Drosophila embryo*. Genes Dev, 1998. **12**(4): p.  
636 547-56.
- 637 51. Zabidi, M.A., et al., *Enhancer-core-promoter specificity separates*  
638 *developmental and housekeeping gene regulation*. Nature, 2015. **518**(7540):  
639 p. 556-9.
- 640 52. Roadmap Epigenomics, C., et al., *Integrative analysis of 111 reference human*  
641 *epigenomes*. Nature, 2015. **518**(7539): p. 317-30.
- 642 53. Consortium, E.P., *An integrated encyclopedia of DNA elements in the human*  
643 *genome*. Nature, 2012. **489**(7414): p. 57-74.
- 644 54. Burges, C.J.C., *A Tutorial on Support Vector Machines for Pattern Recognition*.  
645 Data Mining and Knowledge Discovery, 1998. **2**: p. 121--167.
- 646 55. Rajagopal, N., et al., *RFECs: a random-forest based algorithm for enhancer*  
647 *identification from chromatin state*. PLoS Comput Biol, 2013. **9**(3): p.  
648 e1002968.
- 649 56. Koch, C.M., et al., *The landscape of histone modifications across 1% of the*  
650 *human genome in five human cell lines*. Genome Res, 2007. **17**(6): p. 691-707.
- 651 57. Bailey, S.D., et al., *ZNF143 provides sequence specificity to secure chromatin*  
652 *interactions at gene promoters*. Nat Commun, 2015. **2**: p. 6186.
- 653 58. mod, E.C., et al., *Identification of functional elements and regulatory circuits by*  
654 *Drosophila modENCODE*. Science, 2010. **330**(6012): p. 1787-97.
- 655



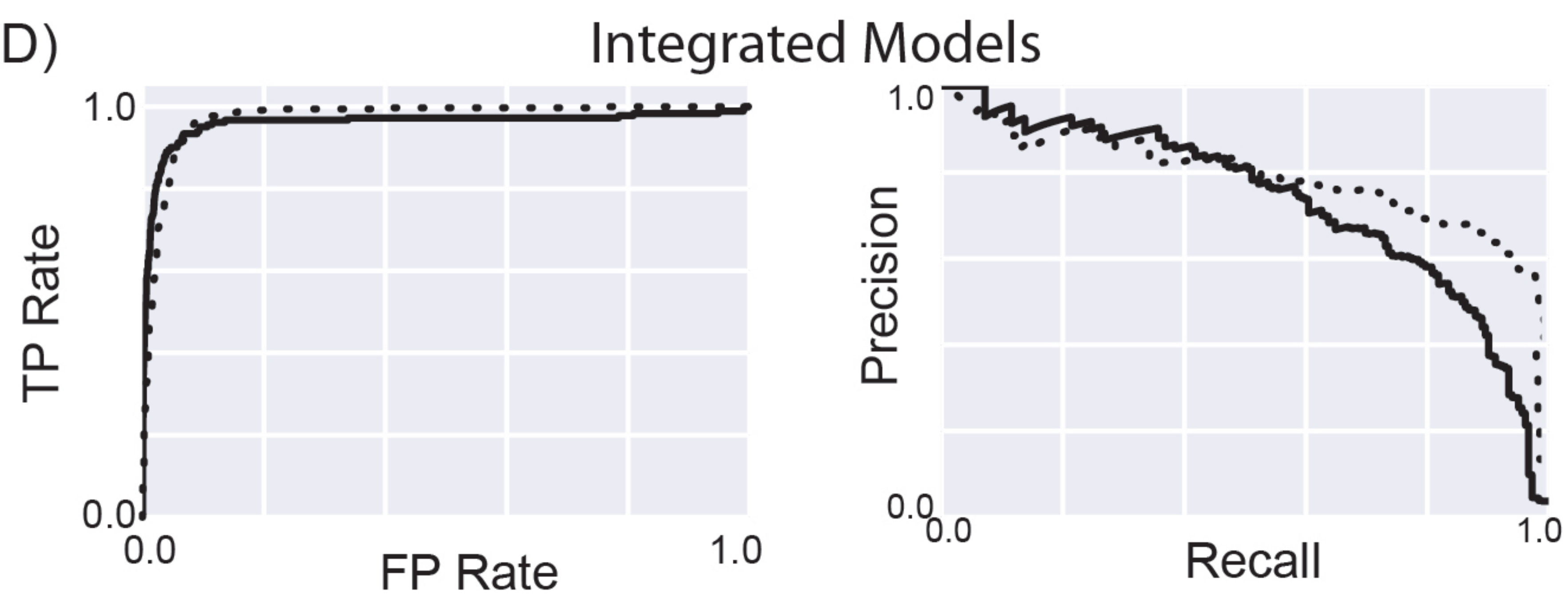
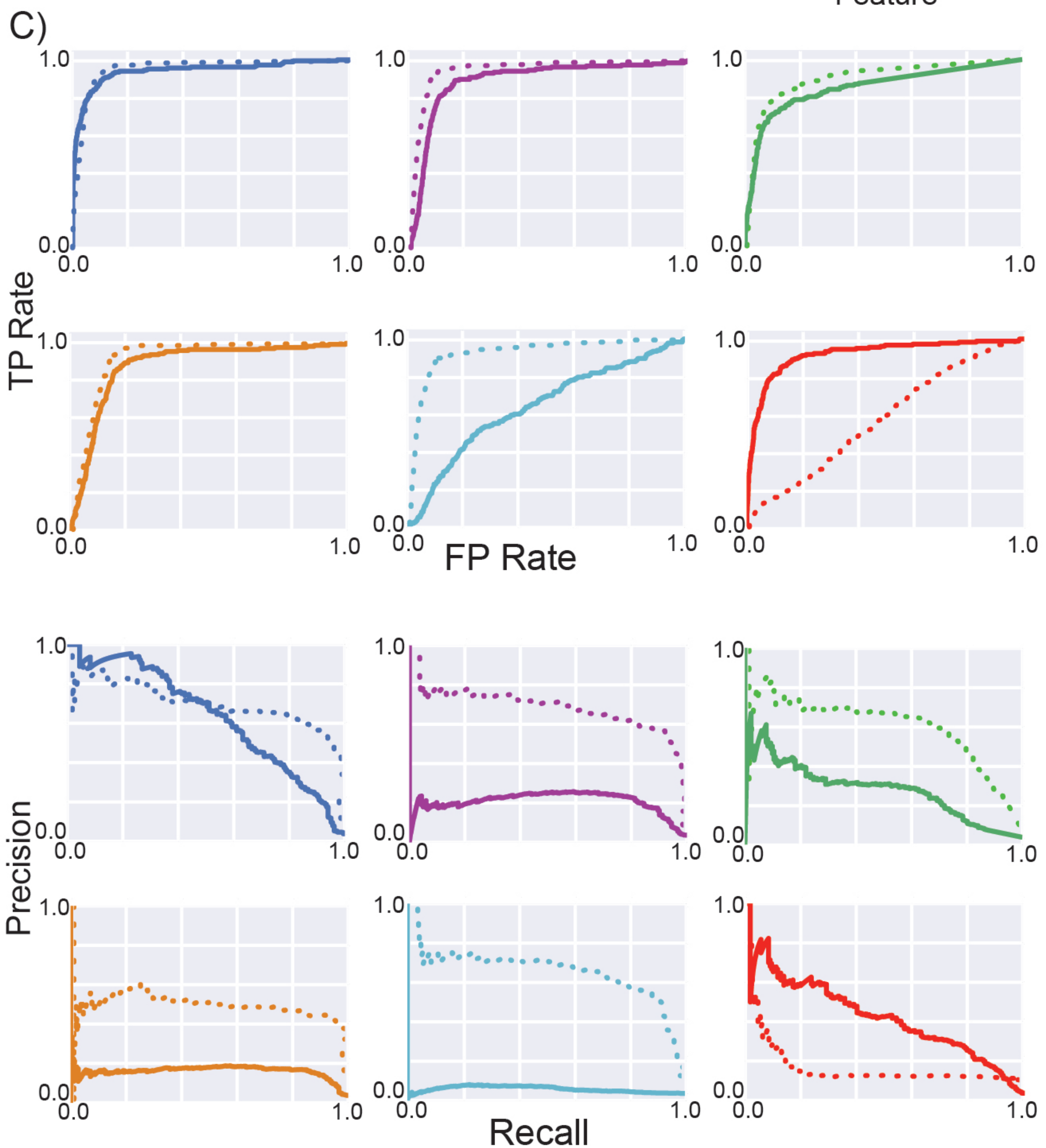
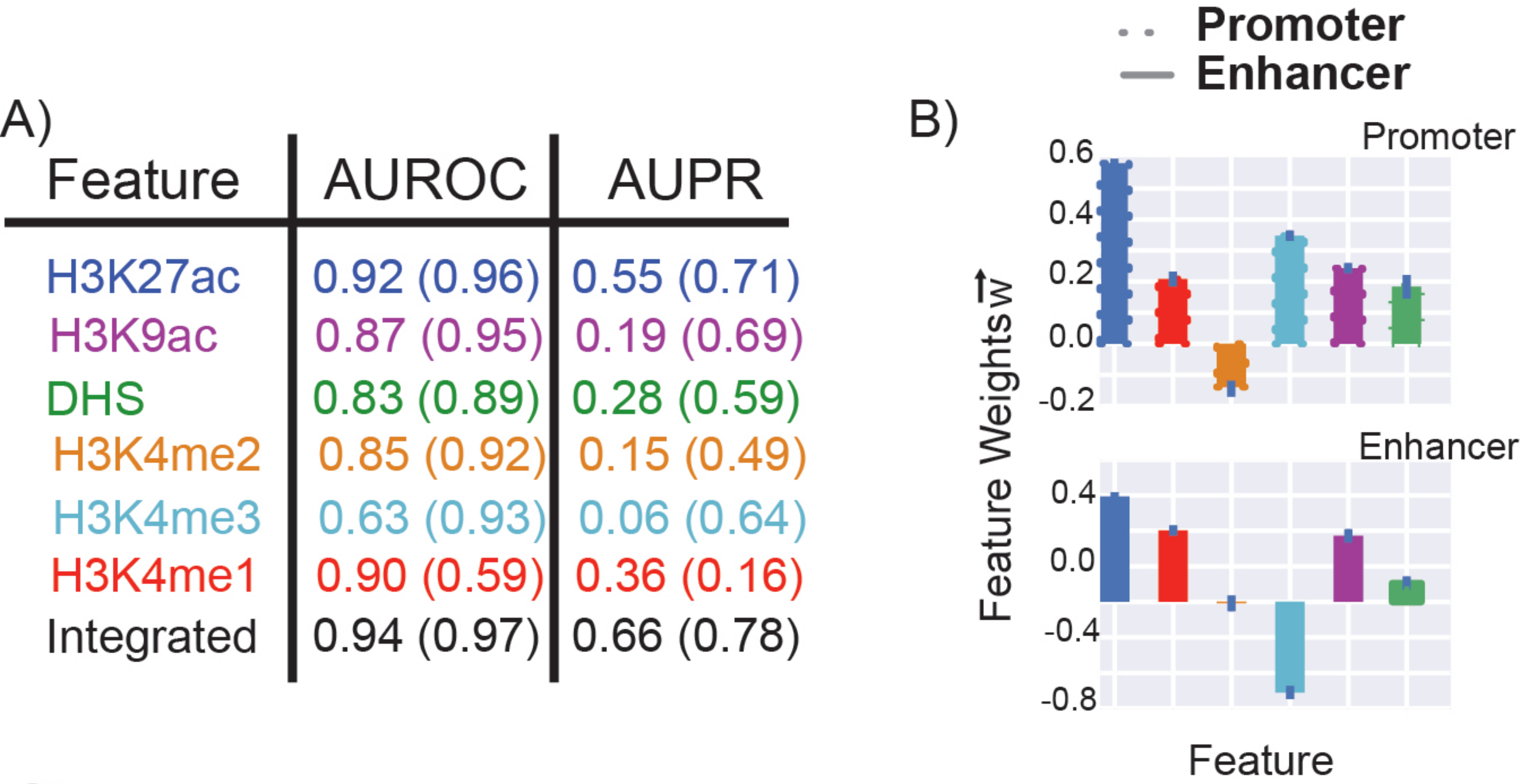


-- Single Core Promoter  
 — Multiple Core Promoters

A) Feature	AUROC	AUPR
H3K27ac	0.95 (0.92)	0.80 (0.72)
H3K9ac	0.92 (0.89)	0.82 (0.52)
DHS	0.88 (0.86)	0.79 (0.58)
H3K4me2	0.90 (0.87)	0.73 (0.41)
H3K4me3	0.82 (0.73)	0.71 (0.32)
H3K4me1	0.70 (0.80)	0.56 (0.46)
Integrated	0.96 (0.95)	0.91 (0.76)









# Fly-based models on mouse

·· Promoter  
— Enhancer

A)

Feature	AUROC	AUPR
H3K27ac	0.86 (0.95)	0.38 (0.71)
H3K9ac	0.80 (0.97)	0.23 (0.83)
DHS	0.90 (0.96)	0.34 (0.70)
H3K4me3	0.74 (0.97)	0.21 (0.82)
H3K4me1	0.83 (0.66)	0.27 (0.17)
Integrated	0.87 (0.98)	0.40 (0.83)

