

1 **A framework for supervised enhancer prediction with epigenetic pattern**
2 **recognition and targeted validation across organisms**

3
4 Anurag Sethi^{1,2,†}, Mengting Gu^{1,†}, Emrah Gumusgoz⁶, Landon Chan³, Koon-Kiu Yan^{1,2},
5 Kevin Yip⁴, Joel Rozowsky^{1,2}, Iros Barozzi⁷, Veena Afzal⁷, Jennifer Akiyama⁷, Ingrid
6 Plajzer-Frick⁷, Catherine Pickle⁷, Momoe Kato⁷, Tyler Garvin⁷, Quan Pham⁷, Anne
7 Harrington⁷, Brandon Mannion⁷, Elizabeth Lee⁷, Yoko Fukuda-Yuzawa⁷, Axel Visel⁷,
8 Diane E. Dickel⁷, Richard Sutton⁶, Len A. Pennacchio⁷ and Mark Gerstein^{1,2,5}

9
10
11
12 ¹ Program in Computational Biology and Bioinformatics, Yale University, New Haven,
13 Connecticut, United States of America

14 ² Molecular Biophysics and Biochemistry, Yale University, New Haven, Connecticut,
15 United States of America

16 ³ School of Medicine, The Chinese University Hong Kong, China

17 ⁴ Department of Computer Science, The Chinese University Hong Kong, China

18 ⁵ Department of Computer Science, Yale University, New Haven, Connecticut, United
19 States of America

20 ⁶ Department of Internal Medicine, Section of Infectious Diseases, Yale University School
21 of Medicine, New Haven, Connecticut, United States of America

22 ⁷ Functional Genomics Department, Lawrence Berkeley National Laboratory, Berkeley,
23 California, United States of America

24
25
26

27 **Abstract**

28

29 Enhancers are important noncoding elements, but they have been traditionally hard to
30 characterize experimentally. Only a few mammalian enhancers have been validated,
31 making it difficult to train statistical models for their identification properly. Instead,
32 postulated patterns of genomic features were used heuristically for identification. The
33 development of massively parallel assay allows the characterization of large numbers of
34 enhancers for the first time. Here, we develop a framework that uses them to create
35 shape-matching filters based on enhancer-associated meta-profiles of epigenetic
36 features. These features are combined with supervised machine learning algorithms (i.e.,
37 SVMs) to predict enhancers. We demonstrated that our model can be applied to predict
38 enhancers in mammalian species (eg, mouse and human). The predictions are
39 comprehensively validated using a combination of *in vivo* and *in vitro* assays (133
40 mouse transgenic enhancer assays in 6 different tissues and 25 human H1 hESC
41 transduction-based reporter assays). The validation results confirm that our model can
42 accurately predict enhancers in different species without re-parameterization. Finally, we
43 predict enhancers in cell lines with many transcription-factor binding sites. This highlights
44 distinct differences between the type of binding at enhancers and promoters, enabling
45 the construction of a secondary model discriminating between these two.

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78
79
80
81

82 Introduction

83

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

95 Traditionally, the regulatory activity of enhancers and promoters were experimentally
96 validated in a non-native context using low throughput heterologous reporter constructs
97 leading to a small number of validated enhancers that function in the same mammalian
98 cell-type [15, 16]. In addition to the small numbers, the validated enhancers were
99 typically selected based on conserved noncoding regions [17] with particular patterns of
100 chromatin [18], transcription-factor binding, [19] or noncoding transcription [20]. The
101 small number and biases within the validated enhancers make them inappropriate for
102 parameterizing tissue-specific enhancer prediction models [16]. As a result, most
103 theoretical methods to predict enhancers could not optimally parameterize their models
104 using a gold-standard set of functional elements. Instead, most of these models were
105 parameterized based on certain heuristic features associated with enhancers, which
106 were then utilized to predict enhancers [19, 21-30]. For example, two widely used
107 methods for predicting enhancers were based on the fact that these elements are
108 expected to contain a cluster of transcription factor binding sites [24] and their activity is
109 often correlated with an enrichment of particular post-translational modifications on
110 histone proteins [27, 30]. These predictions could not be comprehensively assessed as
111 few putative enhancers could be validated experimentally due to the low throughput of
112 validation assays and it remains challenging to assess the performance of different
113 methods for enhancer prediction.

114

115 In recent times, due to the advent of next-generation sequencing, a number of
116 transfection and transduction-based assays were developed to experimentally test the
117 regulatory activity of thousands of regions simultaneously in a massively parallel fashion
118 [31-37]. In these experiments, several plasmids that each contains a single core
119 promoter upstream of a luciferase or GFP gene are transfected or transduced into cells.
120 These plasmids are used to test the regulatory activity of different regions by placing one
121 region within the screening vector in each plasmid as differences in the gene's
122 expression occur due to the differences in the activity of the tested region. STARR-seq
123 was one such massively parallel reporter assay (MPRA) that was used to test the
124 regulatory activity of the fly genome by inserting candidate fragments from the genome
125 within the 3' untranslated region of the luciferase gene. STARR-seq identified thousands
126 of cell-type specific enhancers and promoters within the fly genome [31, 38]. MPRA

127 have confirmed that active enhancers and promoters tend to be depleted of histone
128 proteins and contain accessible DNA on which various transcription factors and
129 cofactors bind [39, 40]. These regulatory regions also tend to be flanked by
130 nucleosomes that contain histone proteins with certain characteristic post-translational
131 modifications. These attributes lead to an enriched peak-trough-peak (“double peak”)
132 signal in different ChIP-Seq experiments for various histone modifications such as
133 acetylation on H3K27 and methylations on H3K4. The troughs in the double peak ChIP-
134 seq signal represent the accessible DNA that leads to a peak in the DNase-I
135 hypersensitivity (DHS) at the enhancers [41]. However, the optimal method to combine
136 information from multiple epigenetic marks to make cell-type specific regulatory
137 predictions remains unknown. For the first time, using data from several MPRA, we
138 have the ability to properly train our models based on a large number of experimentally
139 validated enhancers and test the performance of different models for enhancer
140 prediction using cross validation.

141
142 Our goal in this paper is to develop a framework for making supervised enhancer
143 prediction models using MPRA datasets. We make use of all published data resources
144 to provide a comprehensive model for enhancer prediction that can be applied across
145 different contexts (i.e., different species and tissue types); we validate our model in a
146 variety of different contexts. In particular, we utilized extensive datasets from STARR-
147 seq experiments performed on fly cell lines to create and parameterize our model. Unlike
148 previous prediction methods that focused on the enrichment (or signal) of different
149 epigenetic datasets, we developed a method to also take into account the enhancer-
150 associated pattern within different epigenetic signals. As the epigenetic signal around
151 each enhancer is noisy, we aggregated the signal around thousands of enhancers
152 identified using MPRA to increase signal-to-noise ratio, and identified the shape
153 associated with active regulatory regions. Previous ENCODE and modENCODE efforts
154 showed that the chromatin modifications on active promoters and enhancers were
155 conserved across higher eukaryotes [42-48]. The signal of different chromatin
156 modifications upstream of a gene have been used to create a universal model for
157 predicting its expression and the parameters of the model were transferable across
158 humans, flies, and worm. Here, we further explored this conservation of epigenetic
159 signal shapes for constructing simple-to-use transferrable statistical models with six
160 parameters that were used to predict enhancers and promoters in diverse eukaryotic
161 species including fly, mouse, and human. We showed that the enhancer predictions from
162 our transferrable model was comparable to the prediction accuracy of species-specific
163 models.

164
165 Working across organisms also allowed us to take advantage of different assays to
166 validate our predictions in a robust fashion using multiple experimental approaches. In
167 the first stage, we predicted enhancers in six different embryonic mouse tissues and
168 tested the activity of these predictions *in vivo* with transgenic mouse assays. Due to the
169 obvious ethical considerations of performing such transgenic assays in human embryos,
170 we then proceeded to test the activity of these elements in a human cell-line *in vitro*.

171
172 H1-hESC is a highly studied human cell-line in which a comprehensive set of
173 transcription factor (TF) binding experiments are available. After validating our
174 predictions, the many TFs provided us with the opportunity to differentiate between the
175 enhancers and promoters. The pattern of TF and co-TF binding at active enhancers is
176 much more heterogeneous than the corresponding patterns on promoters, which can be
177 used to distinguish enhancers from promoters with high accuracy. Thus, our methods

178 provide a framework that utilizes different epigenetic genomics datasets to predict active
179 regulatory regions in a cell-type specific manner. Further functional genomics datasets
180 can be utilized to identify key TFs associated with active regulatory regions within these
181 cell types.

182

183 **Results**

184

185 **Aggregation of epigenetic signal (in fly) to create metaprofile:**

186

187 We developed a framework to predict active regulatory elements using the epigenetic
188 signal patterns associated with experimentally validated promoters and enhancers [31].
189 We aggregated the signal of histone modifications on MPRA peaks to remove noise in
190 the signal and created a metaprofile of the double peak signals of histone modifications
191 flanking enhancers and promoters. MPRA peaks typically consist of a mixture of
192 enhancers and promoters, and at this stage, we do not differentiate between the two
193 sets of regulatory elements. These metaprofiles were then utilized in a pattern
194 recognition algorithm for predicting active promoters and enhancers in a cell-type
195 specific manner.

196

197 The STARR-seq studies on fly cell-lines provide the most comprehensive MPRA
198 datasets as the whole genome was tested for regulatory activity within these assays and
199 these assays were performed with multiple core promoters (cite31, 50). Hence, we
200 chose to create metaprofiles using the histone modification H3K27ac at active STARR-
201 seq peaks (see Figure 1 and Methods) identified within the S2 cell-line of the fly.
202 Approximately 70% of the active STARR-seq peaks contain an easily identifiable double
203 peak pattern even though there is a lot of variability in the distance between the two
204 maxima of the double peak in the ChIP-chip signal (Figure S1). While the minimum
205 tends to occur in the center of these two maxima on average, the distance between the
206 two maxima in the double peaks can vary between 300 and 1100 base pairs. During
207 aggregation, we aligned the two maxima in the H3K27ac signal across different STARR-
208 seq peaks, followed by interpolation and smoothing the signal before calculating the
209 average metaprofile. In addition, an optional flipping step was performed to maintain the
210 asymmetry in the underlying H3K27ac double peak because it may be associated with
211 the directionality of transcription [49]. We also calculated the dependent metaprofiles for
212 thirty other histone marks and DHS signal by applying the same set of transformations to
213 these datasets. The metaprofile for the histone marks associated with active regulatory
214 regions were also double peak signals, and the maxima across different histone
215 modification signals tended to align with each other on average (Figure S2). This
216 indicates that a large number of histone modifications tend to simultaneously co-occur
217 on the nucleosomes flanking an active enhancer or promoter. In contrast, as expected,
218 the DHS signal displayed a single peak at the center of the H3K27ac double peak
219 (Figure 1). In addition, repressive marks such as H3K27me3 were depleted in these
220 regions, and the metaprofile for these regions did not contain a double peak signal
221 (Figure S2).

222

223 **Match of a metaprofile is predictive of regulatory activity:**

224

225 We evaluated whether these metaprofiles can be utilized to predict active promoters and
226 enhancers using matched filters, a well-established algorithm in template recognition. A
227 matched filter is the optimal pattern recognition algorithm that uses a shape-matching
228 filter to recognize the occurrence of a template in the presence of stochastic noise [50].

229 We evaluated whether the occurrence of the epigenetic metaprofiles identified for the
230 histone marks and DHS can be used to predict active enhancers and promoters using
231 receiver operating characteristic (ROC) and precision-recall (PR) curves. PR curves are
232 particularly useful to assess the performance of classifiers in skewed or imbalanced data
233 sets in which one of the classes is observed much more frequently compared to the
234 other class, as it plots the fraction of true positives among all predicted positives. If the
235 area under a PR curve is higher, the corresponding model has a low false discovery rate
236 and can easily distinguish between the positives from the negatives. On the other hand,
237 in skewed datasets, the area under ROC curves could be high even when the FDR is
238 high even. This is because, in these cases, even if a small fraction of negatives are
239 predicted to be positive by the model, the false discovery rate can be high as the total
240 number of true positives are much smaller than the total number of true negatives [51].
241 The matched filter score is higher in genomic regions where the template pattern occurs
242 in the corresponding signal track while it is low when only noise is present in the signal
243 (Figure 1). Due to the aforementioned variability in the double peak pattern, the
244 H3K27ac signal track is scanned with multiple matched filters with templates that vary in
245 width between the two maxima in the double peak and the highest matched filter score
246 with these matched filters is used to rate the regulatory potential of this region (see
247 Methods). The dependent profiles are then used on the same region with the matched
248 filter to score the corresponding genomic tracks.

249
250 We used 10-fold cross validation to assess the performance of matched filters for
251 individual histone marks to predict active STARR-seq peaks. In Figure 2, we observe
252 that the H3K27ac matched filter is the single most accurate feature for predicting active
253 regulatory regions (AUROC=0.92, AUPR=0.72) identified using STARR-seq. This is
254 consistent with the literature as H3K27ac enriched peaks are often used to predict active
255 promoters and enhancers [23, 52, 53]. In general, several histone acetylations (H3K27ac,
256 H3K9ac, H4K12ac, H2BK5ac, H4K8ac, H4K5ac, H3K18ac) marks as well as the H1,
257 H3K4me2, and DHS are the most accurate prediction features (Table S1) because the
258 matched filter scores for these features are higher on the STARR-seq peaks. The
259 degree to which the matched filter scores for promoters and enhancers are higher than
260 the matched filter scores for the rest of the genome is a measure of the signal to noise
261 ratio for regulatory region prediction in the corresponding feature's genomic track. The
262 larger the separation between positives and negatives, the greater the accuracy of the
263 corresponding matched filter for predicting active regulatory regions. Interestingly, the
264 distribution of matched filter scores for STARR-seq peaks are unimodal for each histone
265 mark except for H3K4me1, H3K4me3, and H2Av, which are bimodal (Figure S3). We
266 also show that the matched filter scores are more accurate for predicting active STARR-
267 seq peaks than the enrichment of signal alone as they outperform histone peak calling
268 on ROC and PR curves (Figure S4).

269
270 While a single STARR-seq experiment identifies thousands of active regulatory regions,
271 these regions display core-promoter specificity, and different sets of enhancers are
272 identified when different core promoters are used in the same cell-type [54-58]. As we
273 wanted to create a framework to predict all the enhancers and promoters active in a
274 particular cell type, we combined the peaks identified from multiple STARR-seq
275 experiments in the S2 cell-type and reassessed the performance of the matched filters at
276 predicting these regulatory regions. Merging the STARR-seq peaks from multiple core
277 promoters in the S2 cell-type leads to higher AUROC and AUPR for the matched filters
278 from most histone marks (Figure 2 and Table S2).

279

280 **Machine learning can combine matched filter scores from different epigenetic**
281 **features**

282
283 We built an integrated model with combined matched filter scores of the most
284 informative epigenetics marks (H3K27ac, H3K4me1, H3K4me2, H3K4me3, H3K9ac,
285 and DHS) associated with active regulatory regions using a linear SVM [54]. The
286 selection of six features ensures that the integrated model can be applied to a variety of
287 cell lines and tissues, as many relevant ChIP-seq and DNase experiments have been
288 performed by the Roadmap Epigenomics Mapping [59] and the ENCODE [60] Consortia
289 in a wide variety of samples. We also assessed the performance of other statistical
290 approaches including a nonlinear SVM for combining the features. While all these
291 approaches performed similarly (Figure S5), a linear SVM is used in our framework for
292 its better interpretability.

293
294 During integration, the normalized matched filter score for each epigenetic feature in a
295 particular region is scaled by its optimized weight and added together to form a
296 discriminant function. The sign of the discriminant function is then used to predict
297 whether the region is regulatory. The features with large positive and negative weights
298 are predicted to be important for discriminating regulatory from non-regulatory regions.
299 The optimized weights can also be used to measure the amount of non-redundant
300 information added by each feature in the integrated model. According to the model, the
301 acetylations (H3K27ac and H3K9ac) are the most important feature for predicting active
302 regulatory regions. The DHS matched filter performed well as an individual feature
303 (AUPR in Figure 2) to predict enhancers and can be highly predictive of regulatory
304 activity in combination with other marks such as H3K27ac (Moore et al., in review).
305 However, in the integrated model, the information in DHS is redundant with the
306 information contained within the five histone marks as indicated by the fact that it has the
307 lowest weight among the six features in the integrated model. The integrated model, as
308 expected, achieved a higher accuracy than the individual matched filter scores (Figure 2),
309 as they can leverage information from multiple epigenetic marks. We also trained a 6-
310 parameter SVM model using STARR-seq data in BG3 cell-line. The model is highly
311 accurate at predicting active enhancers and promoters in the S2-cell line (Figure S6),
312 indicating our framework of combining epigenetic features with a linear SVM model to
313 predict enhancers is applicable across species of great evolutionary distance.

314
315
316 To assess the information contained in other epigenetic marks, we combined the
317 matched filters from all 30 measured histone marks along with the DHS matched filter in
318 separate statistical models (Figure S7) and these models displayed higher accuracy
319 (AUROC=0.97, AUPR=0.93 for SVM model with multiple core promoters) than the 6
320 feature model presented in Figure 2. The feature weights in this model indicated that
321 H3K27ac contains the most information regarding the activity of regulatory regions.
322 However, we found that a few other acetylations such as H2BK5ac, H4ac, and H4K12ac
323 contain additional non-redundant information regarding the activity of these regulatory
324 regions and might improve the accuracy of promoter and enhancer prediction from
325 machine learning models.

326
327 **Distinct epigenetic signals associated with promoters and enhancers**

328
329 We proceeded to create individual metaprofiles and machine learning models for the two
330 classes of regulatory activators – promoters (or proximal) and enhancers (or distal). We

331 divided all the active STARR-seq peaks into promoters or enhancers based on their
332 distance to the closest transcription start site (TSS) to delineate their likely function in the
333 native context. Due to the conservative distance metric used in this study (1kb upstream
334 and downstream of TSS in fly), the enhancers are regulatory elements that are not close
335 to any known TSS and could be considered to enhance gene transcription from a
336 distance. However, a few of the promoters may also regulate distal genes in addition to
337 their promoter activity. We then created metaprofiles of the different epigenetic marks on
338 the promoters and enhancers and assessed the performance of the matched filters for
339 predicting active regulatory regions within each category (Figure 3). The highest
340 matched filter scores are typically observed on promoters, and the matched filters for
341 each of the six features tended to perform better for promoter prediction. The H3K27ac
342 matched filter continues to outperform other epigenetic marks for predicting active
343 promoters and enhancers. In addition, the DHS, H3K9ac, and H3K4me2 matched filters
344 also performed reasonably for promoter and enhancer prediction. Similar to previous
345 studies [61, 62], we observed that the H3K4me1 metaprofile performs better for
346 predicting enhancers while it is close to random for predicting promoters. In contrast, the
347 H3K4me3 metaprofile can be utilized to predict promoters and not enhancers. The
348 histogram for matched filter scores shows that H3K4me1 matched filter score is higher
349 near enhancers while the H3K4me3 matched filter score tends to be higher near
350 promoters (Figure S8). The mixture of these two populations lead to bimodal
351 distributions for H3K4me1 and H3K4me3 matched filter scores when calculated over all
352 regulatory regions (Figure S3).

353
354 We created different integrated models to learn the combination of features associated
355 with promoters and enhancers respectively. These integrated models outperformed the
356 individual matched filters at predicting active enhancers and promoters (Figures 3 and
357 S9). In addition, the weights of the individual features identified the difference in roles of
358 the H3K4me1 and H3K4me3 matched filter scores at discriminating active promoters
359 and enhancers from inactive regions in the genome. The promoter-based (enhancer-
360 based) model performed much more poorly at predicting enhancers (promoters)
361 indicating the unique properties of these regions (Figures S10 and S11). We also
362 created two integrated models utilizing matched filter scores of all thirty histone marks as
363 features for predicting enhancers and promoters. The additional histone marks provided
364 independent information regarding the activity of promoters and enhancers as these
365 features increased the accuracy of these models (Figure S12). The weights of different
366 features indicate that H2BK5ac again displays the most independent information for
367 accurately predicting active enhancers and promoters. We observe similar trends and
368 accuracy with several different machine learning methods (Figures S9 and S12).

369

370

371 **Application of STARR-seq model to predict enhancers in mammalian species**

372

373 One of the important findings of previous ENCODE and model organism ENCODE
374 efforts is the conservation of chromatin marks close to regulatory elements across
375 hundreds of millions of years of evolution [42-48]. The relationship of chromatin marks to
376 gene expression was very similar, for instance, in worms, flies, mice and human, so
377 much that one could build a statistical model relating chromatin modification to gene
378 expression that would work without re-parameterization across different organisms. This
379 motivated us to apply our well-parameterized model based on the STARR-seq data from
380 flies to mammalian systems -- eg. mouse and human -- and test our model performance.

381

382 We started with genome-wide predictions of regulatory regions in mouse. Tissue-specific
383 epigenetic signals were processed and applied to our model to account for the tissue
384 specificity of enhancers. Predictions are made in six different tissues (forebrain, midbrain,
385 hindbrain, limb, heart and neural tube) at mouse e11.5 stage (Data available through our
386 website at <https://github.com/gersteinlab/MatchedFilter>). These tissues are selected as
387 their epigenetic signals are highly studied in mouse ENCODE, providing us with a rich
388 source of raw data that can be utilized for making enhancer and promoter predictions. In
389 addition, the VISTA database contains close to 100 validated enhancers that can be
390 used for test for each of these tissues. Using our model, we predicted 31K to 39K
391 regulatory regions in individual tissues in mouse, with each region ranging from 300bp to
392 1100bp. Notably, a consistent proportion of two-thirds (66%~70%) of these predicted
393 regulatory regions are distal regulatory elements for all six tissues, with the other one-
394 third (30%~34%) being proximal regulators (Table S3). These numbers agree with a
395 previous enhancer evolution study [8], and suggest that the amount of enhancers and
396 promoters are likely comparable in different tissues.

397

398

399 Similarly, we did genome wide prediction of regulatory regions in ENCODE top tier
400 human cell lines, including H1-hESC, GM12878, K562, HepG2 and MCF-7 (all available
401 through our website). For each cell line, we utilized the 6-parameter integrated model to
402 predict active enhancers and promoters based on the epigenetic datasets measured by
403 the ENCODE consortium [60]. In H1-hESC, for example, we predicted 43463 active
404 regulatory regions, of which 22828 (52.5%) are within 2kb of the TSS and are labeled as
405 promoters. A large proportion of the predicted enhancers are found in the introns
406 (30.41%) and intergenic regions (13.93%) (Figure S13). The predicted promoters and
407 enhancers are significantly closer to active genes than might be expected randomly
408 (Figure S14).

409

410 **Comparison of STARR-seq model to mammalian models for enhancer prediction**

411

412 We next tried to evaluate how well the STARR-seq model did on predicting mammalian
413 enhancers. Particularly, we want to compare the current mouse enhancer predictions
414 with predictions from models directly trained on mouse data. The relatively large number
415 of known mouse enhancers from VISTA database enabled us to parameterize a model
416 in a same way as what we did with the fly STARR-seq data. However, the VISTA
417 database is not nearly at the same scale as the fly STARR-seq dataset. In total, we
418 pulled together 1253 tissue specific positive regions and 8631 tissue specific negative
419 regions from the assays.

420

421

422 With VISTA database, we trained four models based on four sets of available E11.5
423 mouse tissue-specific enhancers (hindbrain, limb, midbrain and neural tube), and
424 assessed them using 10-fold cross-validation respectively. (There are no DHS data
425 available for E11.5 forebrain and heart thus these two tissues are excluded for fair
426 comparison). The average AUROC value is compared to the AUROC of testing STARR-
427 seq trained model on the same VISTA enhancer data. Despite the significantly
428 unbalanced negative to positive ratios of mouse enhancers in the database, the 6-
429 parameter integrative SVM models learned using balanced fly STARR-seq data were
430 highly accurate at predicting active enhancers and promoters in mouse (Figure S15 A).
431 The cross-validated mouse model, while it did well, performed no better on predicting
432 mouse tissue specific enhancers. We found that the best performing one among the

433 mouse models is for tissue midbrain, likely due to the fact that the number of validated
434 midbrain enhancers is the largest. To construct a larger training sample for mouse, we
435 pooled together the normalized z-scores of matched filter scores for six epigenetic
436 signals of all four tissues, and parameterized a model using this larger set of data. Again,
437 we observed that the original model trained with fly STARR-seq data performed equally
438 well on predicting mouse enhancers and much better in predicting fly enhancers (Figure
439 S15 B). Overall, the result suggests that using the larger and more comprehensive
440 STARR-seq data set for parameter tuning was superior to using the smaller mouse data
441 set, even on mouse.

442
443 In human we did not have an extensive amount of validated enhancer data to allow us to
444 re-parameterize our model and compare to the STARR-seq model. Instead, we
445 compared our predicted enhancers to the enhancer predictions from popular
446 segmentation-based algorithms in human cells, eg, chromHMM [63] and SegWay [27].
447 We observe that a majority of the predicted enhancers and promoters are also predicted
448 to be enhancers and promoters by chromHMM and SegWay respectively (Figures S16
449 to S19).

450
451 Given the above overall statistical and computational evaluations, we are confident in
452 the STARR-seq parameterized model. We then set out to do targeted unbiased
453 validations of the mammalian enhancers predicted, which is described in the next two
454 sections.

455

456

457 **Validation in vivo in Mouse**

458

459 To test the activity of predicted mouse enhancers in vivo, we performed transgenic
460 mouse enhancer assay in e11.5 mice for 133 regions in heart and forebrain, including
461 102 regions selected based on the H3K27ac signals rank of corresponding mouse
462 tissues, and 31 regions selected by an ensemble approach from human homolog
463 sequences (See Methods and Supplement Table S4, S5). In addition, we obtained
464 another set of transgenic mouse enhancer assay results from ENCODE Phase III
465 Encyclopedia (Moore et al., in review), which assessed 151 regions in mouse e11.5
466 hindbrain, midbrain and limb. The combined results from these two large sets of
467 validations, as well as any previously tested tissue-specific e11.5 enhancers from VISTA
468 database, allow us to comprehensively evaluate our enhancer predictions in all six e11.5
469 mouse tissues.

470

471 Among the first 102 tested regions, 62 are selected based on forebrain H3K27ac signal
472 rank, with 20, 22, 20 regions being in the top, middle and bottom rank respectively.
473 Another 40 regions are selected by heart H3K27ac signal rank with half of them coming
474 from the top rank and the other half coming from the middle rank. The bottom ranked
475 regions were skipped because the activity of middle ranked regions dropped off so much.
476 Consistently, the observed active rate of assessed regions decreases from top tier to
477 bottom tier. The validation result suggested a great prediction accuracy of our model: 61%
478 predicted active rate versus 70% observed active rate for top tier, 45% predicted active
479 rate versus 32% observed active rate for middle tier, and 34% predicted active rate
480 versus 35% observed active rate for bottom tier in forebrain, etc. For the other 31 human
481 homolog sequences, 12.9% and 9.7% of the assessed regions are active in heart and
482 forebrain respectively. The lower active rate is likely due to the fact that these human

483 sequences are less well behaved in mouse tissues compared to their original native
484 environment.

485
486

487 For systematic comparison, we evaluate the predictability of our matched filter model for
488 each individual histone marks and DHS, as well as the integrated SVM model (Figure 4).
489 Consistent with previous result from STARR-seq data, H3K27ac signal is the single best
490 performed histone marks for predicting enhancers, while DHS signal performs well as an
491 independent source. The integrated model, as expected, out-performs the individual
492 histone mark models. We then did similar evaluation using the regulatory elements
493 identified by the transduction-based FIREWACH assay in mouse embryonic stem cells
494 (mESC) [36]. With the same metaprofiles, the predictions are based on epigenetic
495 signals of mESC available from ENCODE website. Again, we observe similar results for
496 individual histone marks and combined SVM model (Figure S20). As the *in vivo* and
497 FIREWACH assays utilized a single core promoter to validate regulatory regions, the
498 performance of the different models in Figures 4 and S20 are probably underestimated.

499
500
501

500 **Validation in human cell lines**

501

502 We proceeded to validate our STARR-seq based model for predicting human enhancers
503 using an in vitro transduction assay. A third generation, self-inactivating HIV-1 based
504 vector system in which the eGFP reporter was driven by the DNA element of interest
505 was used to validate putative enhancers after stable transduction of various cell lines,
506 including H1 hESC (Figure 5). The predicted enhancers, ranging from 650 to 2500 bp,
507 were PCR amplified from human genomic DNA and inserted just upstream of a basal
508 Oct-4 promoter of 142 bp (a housekeeping promoter is used so that the activity of the
509 putative enhancers should be similar across different cell lines). VSV G-pseudotyped
510 vector supernatants from each were prepared by co-transfection of 293T cells, and
511 these were used to transduce the various cell lines, with empty vector and FG12 vector
512 serving as negative and positive controls, respectively. Putative enhancer activity was
513 assessed by flow cytometric readout of eGFP expression 48-72 h post-transduction,
514 normalized to the negative control.

515

516 A total of 25 predicted intergenic enhancers were randomly selected for validation
517 (Supplementary Table S6). These predictions were chosen randomly to ensure that
518 these truly represented the whole spectrum of predicted enhancers and not just the top
519 tier of predicted enhancers. Of these 25 putative enhancers, 23 were successfully
520 amplified and cloned into the HIV vector. To measure the distribution of gene
521 expression in the absence of enhancer, we also amplified and cloned 25 non-repetitive
522 elements with similar length distribution that were predicted to be inactive using the
523 same HIV vector. All positive and negative DNA elements were transduced and tested
524 for activity in both forward and reverse strand orientations since enhancers are thought
525 to function in an orientation-independent manner. Functional testing was performed in
526 HOS, TZMBL, and A549 cell lines in addition to H1-hESCs.

527

528 Insertion of twelve of the 23 putative enhancers into the HIV vector resulted in a
529 significant increase in eGFP expression (P-value < 0.05 over the distribution of gene
530 expression for negative elements) in the H1-hESCs (Supplementary Table S7). While
531 most of the positive enhancers displayed a significant increase in gene expression
532 irrespective of their orientation, a few elements showed significantly higher levels of
533 gene expression in one of the orientations. In contrast, the negatives displayed much

534 lower levels of gene expression typically (Figure 5 and Supplementary Figure S21). In
535 addition, most of these elements increased gene expression of GFP in the four different
536 cell lines even though some of the elements were preferentially active in one of the cell
537 lines. Overall, 16 of the 23 tested predictions displayed a statistically significant increase
538 in gene expression of the reporter gene in at least one of the cell lines (Supplementary
539 Table S7 and Supplementary Figure S21). Given the promoter specificity of enhancers
540 in such assays, we would anticipate that some of the elements that could not be
541 validated in this particular vector would function as enhancers in a more natural
542 biological context.

543
544

545 **Integrative analysis in human cell-lines: Different Transcription Factors bind to** 546 **enhancers and promoters**

547

548 We further studied the differences in TF binding at promoters and enhancers (Figure 6
549 and Figure S22). We focused on the human H1-hESC cell line as there is large amount
550 of functional genomic assays from the ENCODE [60] and Roadmap Epigenomics
551 Mapping Consortium [59] within these cell lines. Together, the consortia have generated
552 ChIP-Seq data for 60 transcription related factors in H1-hESC cell line, including a few
553 chromatin remodelers and histone modification enzymes. Collectively we call all these
554 transcription related factors "TF"s for simplicity.

555
556

557 We show that the patterns of TF binding within regulatory regions can be utilized in a
558 logistic regression model to distinguish active enhancers from promoters with high
559 accuracy (AUPR = 0.89, AUROC = 0.87) (Figure 6). We were also able to identify the
560 most important features that distinguish promoters from enhancers. In addition to TATA-
561 box associated factors such as TAF1, TAF7, and TBP, the RNA polymerase-II binding
562 patterns as well as chromatin remodelers such as KDM5A and PHF8 are some of the
563 most important factors that distinguish promoters from enhancers in H1-hESC. This
564 provides a framework that can be utilized to identify the most important TFs associated
565 with active enhancers and promoters in each cell-type.

566
567

568 We found that while most promoters and enhancers contain multiple TF binding sites,
569 the pattern of TF binding at promoters is different from that at enhancers and that TF-
570 binding at enhancers displays more heterogeneity: more than 70% of the promoters bind
571 to the same set of 2-3 sequence-specific TFs, which is not observed for enhancers
572 (Figure 6C and S23). The majority of the promoters also contain peaks for several
573 TATA-associated factors (TAF1, TAF7, and TBP). These TF co-associations could lead
574 to mechanistic insights of cooperativity between TFs. For example, similar to a previous
575 study [64], CTCF and ZNF143 may function cooperatively as they are observed to co-
576 occur frequently at distal regulatory regions in this study. Overall, the high heterogeneity
577 associated with enhancer TF-binding is consistent with the absence of a sequence code
578 (or grammar) which can be utilized to easily identify active enhancers on a genome-wide
579 fashion.

580
581

581 **Discussion**

582

583 In this paper, we have developed a framework using transferable supervised machine
584 learning models trained on regulatory regions identified by MPRA to accurately predict

585 active enhancers in a cell-type specific manner. Current, most existing methods were
586 parameterized (not properly “trained”) on regions that had various features associated
587 with promoters and enhancers and only a small number of these regions were typically
588 tested for regulatory activity experimentally in an *ad hoc* manner [19, 21-30]. The rich
589 amount of whole genome STARR-seq experiments [31] can now establish the
590 characteristic pattern flanking active regulatory regions within certain histone
591 modifications. This motivated us to train a shape-matching and filtering model that can
592 be used to identify these patterns within the shape of the ChIP-seq signals. As the
593 chromatin marks and epigenetic profiles associated with active regulatory regions are
594 highly conserved among organisms [42-48], we showed that a well parameterized model
595 in one model organism can be transferred to another with high prediction accuracy.
596

597 In the model, we compared close to 30 epigenetic signals for their ability to predict
598 regulatory elements individually. The H3K27ac matched filter remains the single most
599 important feature for predicting active regions while H3K4me1 and H3K4me3 are shown
600 to distinguish promoters and enhancers. We characterized the amount of redundant
601 information within the metaprofile of different epigenetic features and showed that the
602 ChIP-seq signals of H2BK5ac, H4ac and H2A provide independent information that
603 helps to improve the accuracy of promoter and enhancer predictions. In addition to these
604 30-feature models, we also provide a simple to use six-parameter SVM model for
605 combining H3K27ac, H3K9ac, H3K4me1, H3K4me2, H3K4me3, and DHS to predict
606 active promoters and enhancers in a cell-type specific manner. These six histone marks
607 have been measured for a number of different tissues and cell-types by the Roadmap
608 Epigenomics Mapping [39], the ENCODE [60], and the modENCODE Consortia [65].
609 Based on these signals, our model could be applied in a tissue and cell-type specific
610 fashion in other organisms like mouse and human. We trained our models with datasets
611 from different species and demonstrated that the high-quality STARR-seq data from fly
612 is sufficient to train a well transferable model. We also compared our result with
613 chromHMM [63] and SegWay [27] predictions and observed the majority of them overlap
614 (Figure S17 to S20).
615

616
617 To avoid potential biases, we chose to validate our model using multiple regulatory
618 assays including *in vivo* transgenic assays and *in vitro* transductions assays, in which
619 the predicted region is tested for regulatory activity in the native chromatin environment.
620 The transgenic assays are performed in E11.5 mice for 133 regions of three rank tiers
621 predicted active in mouse heart and forebrain. The experiment is supplemented by
622 another set of 151 assayed regions predicted active in mouse hindbrain, midbrain and
623 limb in ENCODE Phase III Encyclopedia (Moore et al., in review). Together with other
624 validated regulatory regions from VISTA database, we were able to comprehensively
625 validate our tissue-specific predictions in six different tissues in mouse. As we show in
626 figure 4, the H3K27ac and DHS signals continue to be the highest predictive signals in
627 mouse. We also did a similar evaluation with publicly available FIREWACH assay data
628 [36] in mouse, and the results are consistent. Taken together, we showed that the
629 matched filter model is transferable with high accuracy in predicting active enhancers in
630 mouse tissues.
631

632
633 The human cell-line specific regulatory elements predictions are validated through *in*
634 *vitro* transduction assays in human H1-hESC cells. The majority of the predicted
635 elements displayed a significant increase in expression of the reporter gene, further

636 confirming the predictability of our model in mammalian organisms. H1-hESC is a highly
637 studied cell line, allowing us to analyze the differences in the patterns of TF binding at
638 proximal and distal regulatory regions. The TF binding and co-binding patterns at
639 enhancers are much more heterogeneous than that at promoters. This heterogeneity in
640 TF binding patterns makes it more difficult to predict enhancers due to the absence of
641 obvious sequence patterns in distal regulatory regions. However, we were able to create
642 accurate machine learning models that can distinguish proximal promoter regions from
643 distal enhancers based on the patterns of TF ChIP-seq peaks within these regulatory
644 regions. The conservation of the epigenetic underpinnings underlying active regulatory
645 regions sets the stage for our method to study the evolution of tissue-specific enhancers
646 and their genomic properties across different eukaryotic species.

647
648
649 Our results echo to the previous findings that the epigenetic profiles associated with
650 active enhancers and promoters are highly conserved in evolution [42-48]. Therefore,
651 our model of integrating shape-matching epigenetic scores using fly STARR-seq
652 enhancers can be applied to predict on a variety of tissues and cell lines in other species.
653 In the cross-comparison, we show that the six-parameter integrated model trained in
654 STARR-seq data performs equally well at predicting mouse tissue enhancers with a
655 model trained in VISTA mouse enhancer data. This highlights the advantage of modeling
656 based on a comprehensive genome-wide experimental assay. In the future, we expect
657 that more extensive whole-genome STARR-seq dataset will become available on
658 mammalian systems. It could thus be advantageous to re-train the matched filter model
659 on the state-of-art datasets. With the set up of our framework, re-training the model with
660 newly generated datasets should be straightforward. We envision that our framework
661 would benefit from these datasets and generate more comprehensive regulatory
662 elements annotations across different eukaryotic species.

663
664
665
666
667
668

669 **Acknowledgement:**

670
671
672
673
674
675
676
677
678

M.G. were supported by the NIH grant HG009446-01. A.V. and L.A.P. were supported
by NHLBI grant R24HL123879, and NHGRI grants R01HG003988, and U54HG006997,
and UM1HG009421 where research was conducted at the E.O. Lawrence Berkeley
National Laboratory and performed under Department of Energy Contract DE-AC02-
05CH11231, University of California.

679 **References:**

680

- 681 1. Banerji, J., S. Rusconi, and W. Schaffner, *Expression of a beta-globin gene is*
682 *enhanced by remote SV40 DNA sequences.* Cell, 1981. **27**(2 Pt 1): p. 299-308.
- 683 2. Ong, C.T. and V.G. Corces, *Enhancer function: new insights into the regulation*
684 *of tissue-specific gene expression.* Nat Rev Genet, 2011. **12**(4): p. 283-93.
- 685 3. Woolfe, A., et al., *Highly conserved non-coding sequences are associated with*
686 *vertebrate development.* PLoS Biol, 2005. **3**(1): p. e7.
- 687 4. Spitz, F. and E.E. Furlong, *Transcription factors: from enhancer binding to*
688 *developmental control.* Nat Rev Genet, 2012. **13**(9): p. 613-26.
- 689 5. Cotney, J., et al., *The evolution of lineage-specific regulatory activities in the*
690 *human embryonic limb.* Cell, 2013. **154**(1): p. 185-96.
- 691 6. Degner, J.F., et al., *DNase I sensitivity QTLs are a major determinant of human*
692 *expression variation.* Nature, 2012. **482**(7385): p. 390-4.
- 693 7. Shibata, Y., et al., *Extensive evolutionary changes in regulatory element activity*
694 *during human origins are associated with altered gene expression and positive*
695 *selection.* PLoS Genet, 2012. **8**(6): p. e1002789.
- 696 8. Villar, D., et al., *Enhancer evolution across 20 mammalian species.* Cell, 2015.
697 **160**(3): p. 554-66.
- 698 9. Xiao, S., et al., *Comparative epigenomic annotation of regulatory DNA.* Cell,
699 2012. **149**(6): p. 1381-92.
- 700 10. Wray, G.A., *The evolutionary significance of cis-regulatory mutations.* Nat Rev
701 Genet, 2007. **8**(3): p. 206-16.
- 702 11. Corradin, O. and P.C. Scacheri, *Enhancer variants: evaluating functions in*
703 *common disease.* Genome Med, 2014. **6**(10): p. 85.
- 704 12. Gusev, A., et al., *Partitioning heritability of regulatory and cell-type-specific*
705 *variants across 11 common diseases.* Am J Hum Genet, 2014. **95**(5): p. 535-52.
- 706 13. Slattery, M., et al., *Absence of a simple code: how transcription factors read the*
707 *genome.* Trends Biochem Sci, 2014. **39**(9): p. 381-99.
- 708 14. Levo, M., et al., *Unraveling determinants of transcription factor binding outside*
709 *the core binding site.* Genome Res, 2015. **25**(7): p. 1018-29.
- 710 15. Pennacchio, L.A., et al., *Enhancers: five essential questions.* Nat Rev Genet,
711 2013. **14**(4): p. 288-95.
- 712 16. Erwin, G.D., et al., *Integrating diverse datasets improves developmental*
713 *enhancer prediction.* PLoS Comput Biol, 2014. **10**(6): p. e1003677.
- 714 17. Pennacchio, L.A., et al., *In vivo enhancer analysis of human conserved non-*
715 *coding sequences.* Nature, 2006. **444**(7118): p. 499-502.
- 716 18. Nord, A.S., et al., *Rapid and pervasive changes in genome-wide enhancer usage*
717 *during mammalian development.* Cell, 2013. **155**(7): p. 1521-31.
- 718 19. Visel, A., et al., *ChIP-seq accurately predicts tissue-specific activity of enhancers.*
719 Nature, 2009. **457**(7231): p. 854-8.
- 720 20. Andersson, R., et al., *An atlas of active enhancers across human cell types and*
721 *tissues.* Nature, 2014. **507**(7493): p. 455-61.
- 722 21. Narlikar, L., et al., *Genome-wide discovery of human heart enhancers.* Genome
723 Res, 2010. **20**(3): p. 381-92.

- 724 22. Visel, A., et al., *Ultraconservation identifies a small subset of extremely*
725 *constrained developmental enhancers*. Nat Genet, 2008. **40**(2): p. 158-60.
- 726 23. Bonn, S., et al., *Tissue-specific analysis of chromatin state identifies temporal*
727 *signatures of enhancer activity during embryonic development*. Nat Genet,
728 2012. **44**(2): p. 148-56.
- 729 24. Yip, K.Y., et al., *Classification of human genomic regions based on*
730 *experimentally determined binding sites of more than 100 transcription-*
731 *related factors*. Genome Biol, 2012. **13**(9): p. R48.
- 732 25. Ghandi, M., et al., *Enhanced regulatory sequence prediction using gapped k-*
733 *mer features*. PLoS Comput Biol, 2014. **10**(7): p. e1003711.
- 734 26. Heintzman, N.D., et al., *Distinct and predictive chromatin signatures of*
735 *transcriptional promoters and enhancers in the human genome*. Nat Genet,
736 2007. **39**(3): p. 311-8.
- 737 27. Hoffman, M.M., et al., *Unsupervised pattern discovery in human chromatin*
738 *structure through genomic segmentation*. Nat Methods, 2012. **9**(5): p. 473-6.
- 739 28. Kharchenko, P.V., et al., *Comprehensive analysis of the chromatin landscape in*
740 *Drosophila melanogaster*. Nature, 2011. **471**(7339): p. 480-5.
- 741 29. He, H.H., et al., *Nucleosome dynamics define transcriptional enhancers*. Nat
742 Genet, 2010. **42**(4): p. 343-7.
- 743 30. Ernst, J., et al., *Mapping and analysis of chromatin state dynamics in nine*
744 *human cell types*. Nature, 2011. **473**(7345): p. 43-9.
- 745 31. Arnold, C.D., et al., *Genome-wide quantitative enhancer activity maps identified*
746 *by STARR-seq*. Science, 2013. **339**(6123): p. 1074-7.
- 747 32. Dickel, D.E., et al., *Function-based identification of mammalian enhancers*
748 *using site-specific integration*. Nat Methods, 2014. **11**(5): p. 566-71.
- 749 33. Gisselbrecht, S.S., et al., *Highly parallel assays of tissue-specific enhancers in*
750 *whole Drosophila embryos*. Nat Methods, 2013. **10**(8): p. 774-80.
- 751 34. Kwasniewski, J.C., et al., *High-throughput functional testing of ENCODE*
752 *segmentation predictions*. Genome Res, 2014. **24**(10): p. 1595-602.
- 753 35. Melnikov, A., et al., *Systematic dissection and optimization of inducible*
754 *enhancers in human cells using a massively parallel reporter assay*. Nat
755 Biotechnol, 2012. **30**(3): p. 271-7.
- 756 36. Murtha, M., et al., *FIREWACH: high-throughput functional detection of*
757 *transcriptional regulatory modules in mammalian cells*. Nat Methods, 2014.
758 **11**(5): p. 559-65.
- 759 37. Patwardhan, R.P., et al., *Massively parallel functional dissection of mammalian*
760 *enhancers in vivo*. Nat Biotechnol, 2012. **30**(3): p. 265-70.
- 761 38. Yanez-Cuna, J.O., et al., *Dissection of thousands of cell type-specific enhancers*
762 *identifies dinucleotide repeat motifs as general enhancer features*. Genome
763 Res, 2014. **24**(7): p. 1147-56.
- 764 39. Shlyueva, D., G. Stampfel, and A. Stark, *Transcriptional enhancers: from*
765 *properties to genome-wide predictions*. Nat Rev Genet, 2014. **15**(4): p. 272-86.
- 766 40. Maston, G.A., et al., *Characterization of enhancer function from genome-wide*
767 *analyses*. Annu Rev Genomics Hum Genet, 2012. **13**: p. 29-57.
- 768 41. Thurman, R.E., et al., *The accessible chromatin landscape of the human*
769 *genome*. Nature, 2012. **489**(7414): p. 75-82.

- 770 42. Yue, F., et al., *A comparative encyclopedia of DNA elements in the mouse*
771 *genome*. Nature, 2014. **515**(7527): p. 355-64.
- 772 43. Gerstein, M.B., et al., *Comparative analysis of the transcriptome across distant*
773 *species*. Nature, 2014. **512**(7515): p. 445-8.
- 774 44. Dong, X., et al., *Modeling gene expression using chromatin features in various*
775 *cellular contexts*. Genome Biol, 2012. **13**(9): p. R53.
- 776 45. Cheng, C. and M. Gerstein, *Modeling the relative relationship of transcription*
777 *factor binding and histone modifications to gene expression levels in mouse*
778 *embryonic stem cells*. Nucleic Acids Research, 2012. **40**(2): p. 553-568.
- 779 46. Cheng, Y., et al., *Principles of regulatory information conservation between*
780 *mouse and human*. Nature, 2014. **515**(7527): p. 371-+.
- 781 47. Boyle, A.P., et al., *Comparative analysis of regulatory information and circuits*
782 *across distant species*. Nature, 2014. **512**(7515): p. 453-6.
- 783 48. Gjoneska, E., et al., *Conserved epigenomic signals in mice and humans reveal*
784 *immune basis of Alzheimer's disease*. Nature, 2015. **518**(7539): p. 365-9.
- 785 49. Kundaje, A., et al., *Ubiquitous heterogeneity and asymmetry of the chromatin*
786 *environment at regulatory elements*. Genome Res, 2012. **22**(9): p. 1735-47.
- 787 50. Kumar, V.B.V.K., A. Mahalanobis, and R.D. Juday, *Correlation Pattern*
788 *Recognition*. 2005.
- 789 51. Davis, J. and M. Goadrich, *The Relationship Between Precision-Recall and ROC*
790 *Curves*. Proceedings of the 23rd international conference on Machine
791 Learning, 2006: p. 233-240.
- 792 52. Creighton, M.P., et al., *Histone H3K27ac separates active from poised*
793 *enhancers and predicts developmental state*. Proc Natl Acad Sci U S A, 2010.
794 **107**(50): p. 21931-6.
- 795 53. Rada-Iglesias, A., et al., *A unique chromatin signature uncovers early*
796 *developmental enhancers in humans*. Nature, 2011. **470**(7333): p. 279-83.
- 797 54. Butler, J.E. and J.T. Kadonaga, *Enhancer-promoter specificity mediated by DPE*
798 *or TATA core promoter motifs*. Genes Dev, 2001. **15**(19): p. 2515-9.
- 799 55. Li, X. and M. Noll, *Compatibility between enhancers and promoters determines*
800 *the transcriptional specificity of gooseberry and gooseberry neuro in the*
801 *Drosophila embryo*. EMBO J, 1994. **13**(2): p. 400-6.
- 802 56. Merli, C., et al., *Promoter specificity mediates the independent regulation of*
803 *neighboring genes*. Genes Dev, 1996. **10**(10): p. 1260-70.
- 804 57. Ohtsuki, S., M. Levine, and H.N. Cai, *Different core promoters possess distinct*
805 *regulatory activities in the Drosophila embryo*. Genes Dev, 1998. **12**(4): p.
806 547-56.
- 807 58. Zabidi, M.A., et al., *Enhancer-core-promoter specificity separates*
808 *developmental and housekeeping gene regulation*. Nature, 2015. **518**(7540):
809 p. 556-9.
- 810 59. Roadmap Epigenomics, C., et al., *Integrative analysis of 111 reference human*
811 *epigenomes*. Nature, 2015. **518**(7539): p. 317-30.
- 812 60. Consortium, E.P., *An integrated encyclopedia of DNA elements in the human*
813 *genome*. Nature, 2012. **489**(7414): p. 57-74.

814 61. Rajagopal, N., et al., *RFECs: a random-forest based algorithm for enhancer*
815 *identification from chromatin state*. PLoS Comput Biol, 2013. **9**(3): p.
816 e1002968.

817 62. Koch, C.M., et al., *The landscape of histone modifications across 1% of the*
818 *human genome in five human cell lines*. Genome Res, 2007. **17**(6): p. 691-707.

819 63. Ernst, J. and M. Kellis, *ChromHMM: automating chromatin-state discovery and*
820 *characterization*. Nat Methods, 2012. **9**(3): p. 215-6.

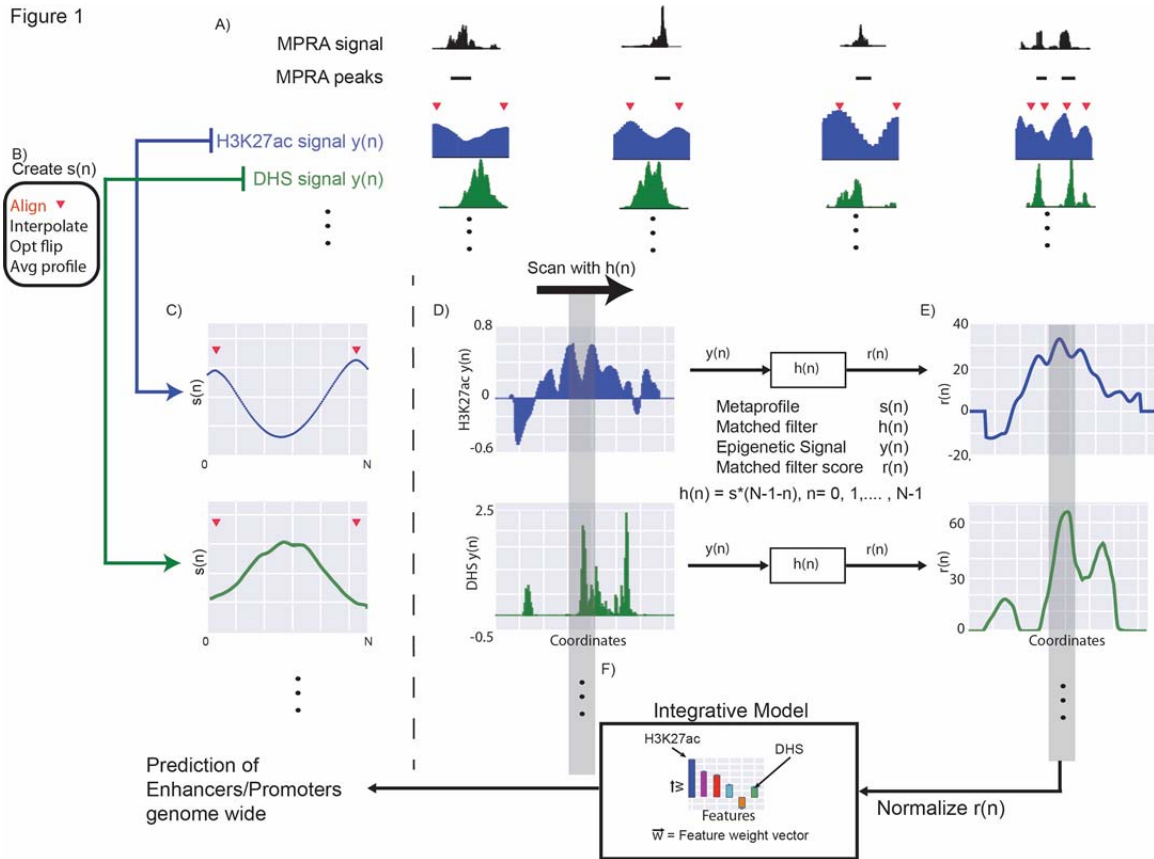
821 64. Bailey, S.D., et al., *ZNF143 provides sequence specificity to secure chromatin*
822 *interactions at gene promoters*. Nat Commun, 2015. **2**: p. 6186.

823 65. mod, E.C., et al., *Identification of functional elements and regulatory circuits by*
824 *Drosophila modENCODE*. Science, 2010. **330**(6012): p. 1787-97.

825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862

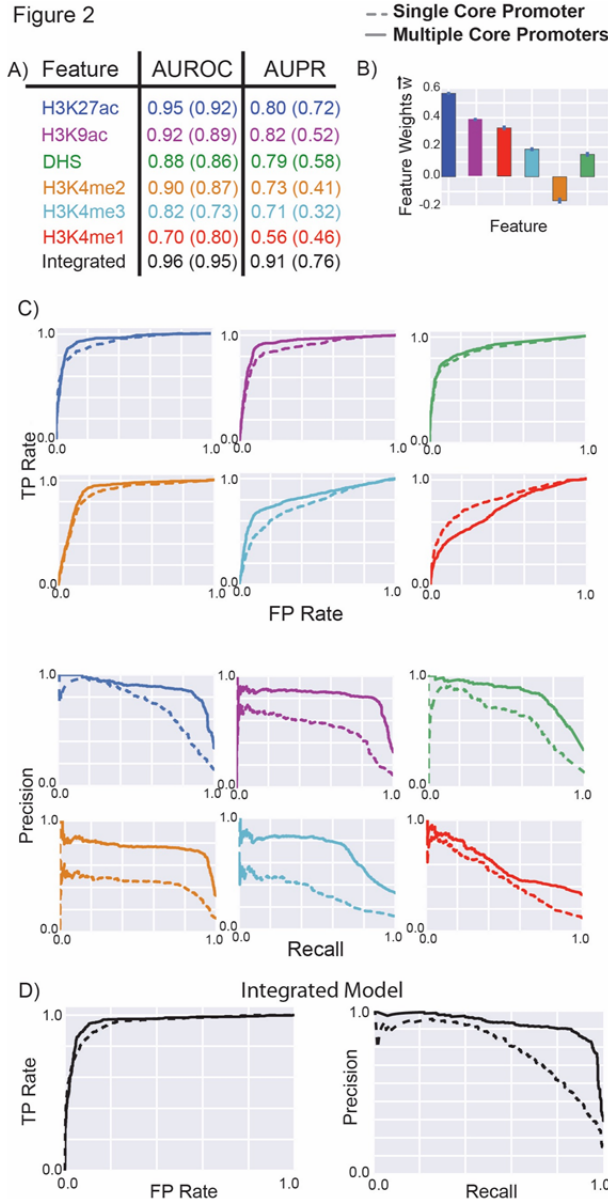
863
864

Figures and Captions



865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881

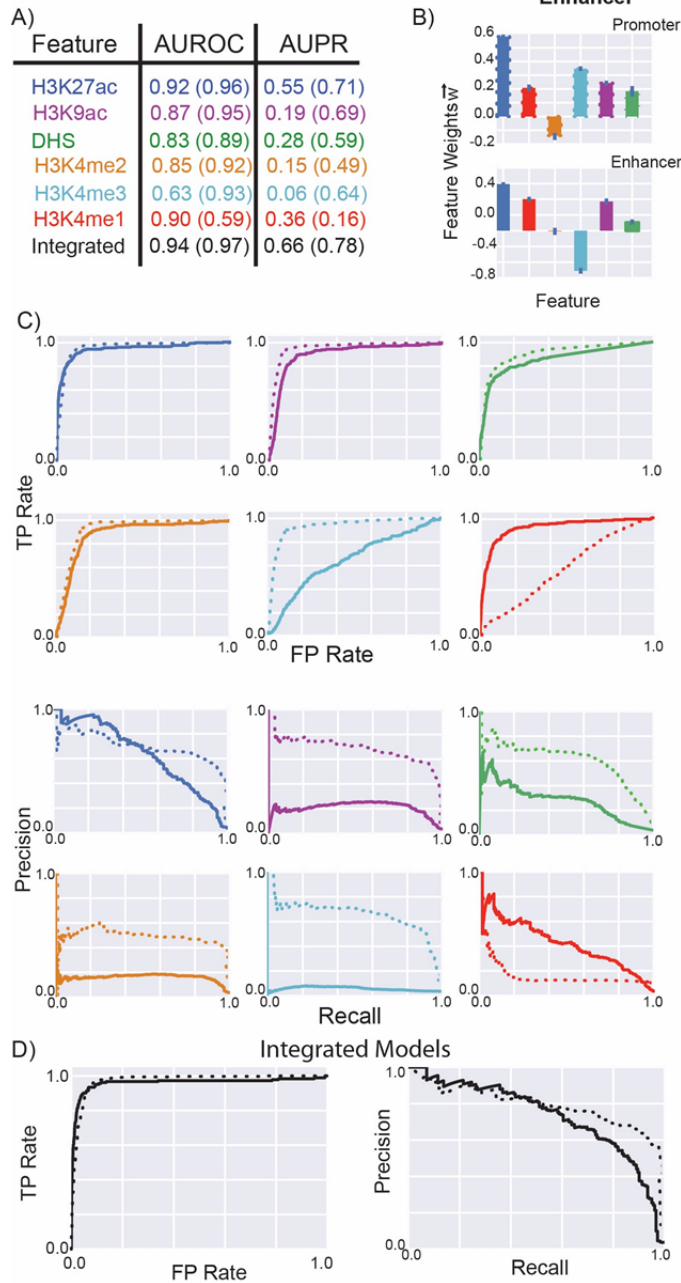
Figure 1: Creation of metaprofile. A) We identified the “double peak” pattern in the H3K27ac signal close to STARR-seq peaks. The red triangles denote the position of the two maxima in the double peak. B) We aggregated the H3K27ac signal around these regions after aligning the flanking maxima, using 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.



882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897

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.

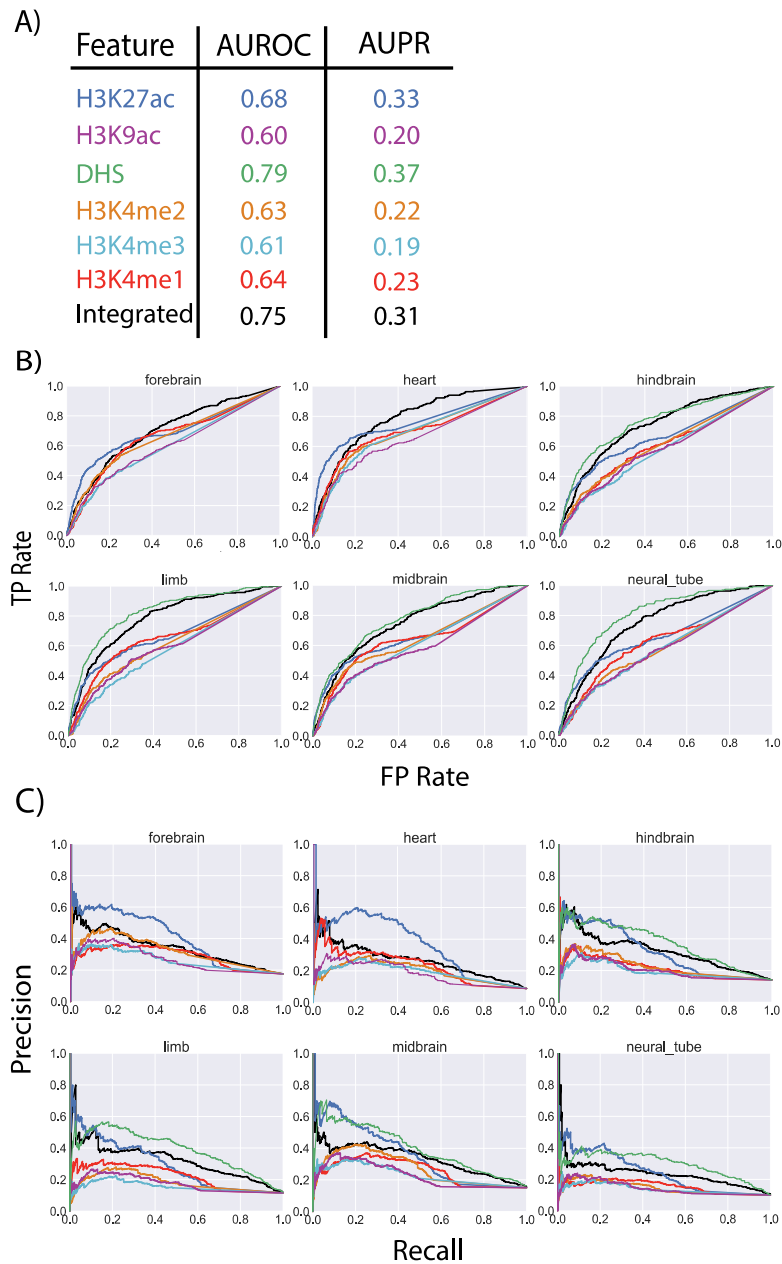
Figure 3



898
899
900
901
902
903
904
905
906
907
908
909
910

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.

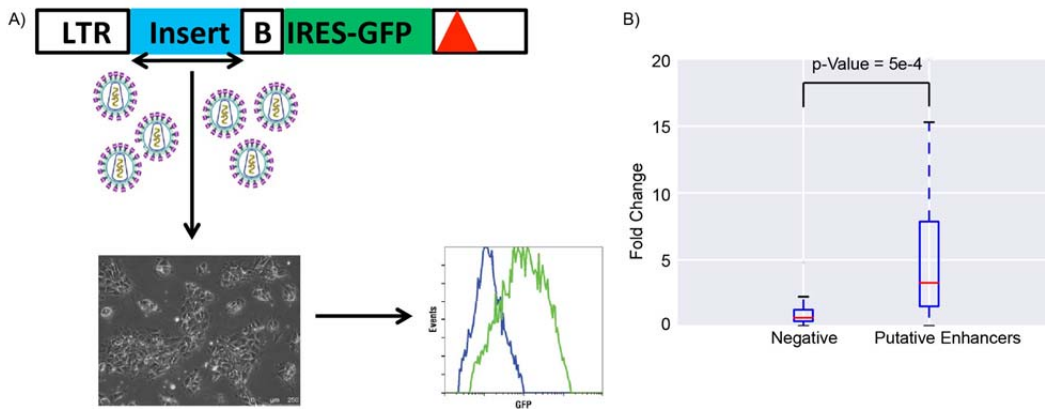
Figure 4



911
912
913
914
915
916
917
918
919
920
921
922

Figure 4: Conservation of epigenetic features. The performance of the fly-based matched filters and the integrated model for predicting active enhancers identified by transgenic mouse enhancer assays at 6 different tissues in E11.5 mice. A) Average AUROC and AUPR for predicting enhancers by different features and by the integrated model. The weights of the different features in the integrated model is the same as the weights shown in Figure 3 for enhancers. B) The individual ROC curves of each feature and the integrated model for each tissue are shown. C) The individual PR curves of each feature and the integrated model for each tissue are shown.

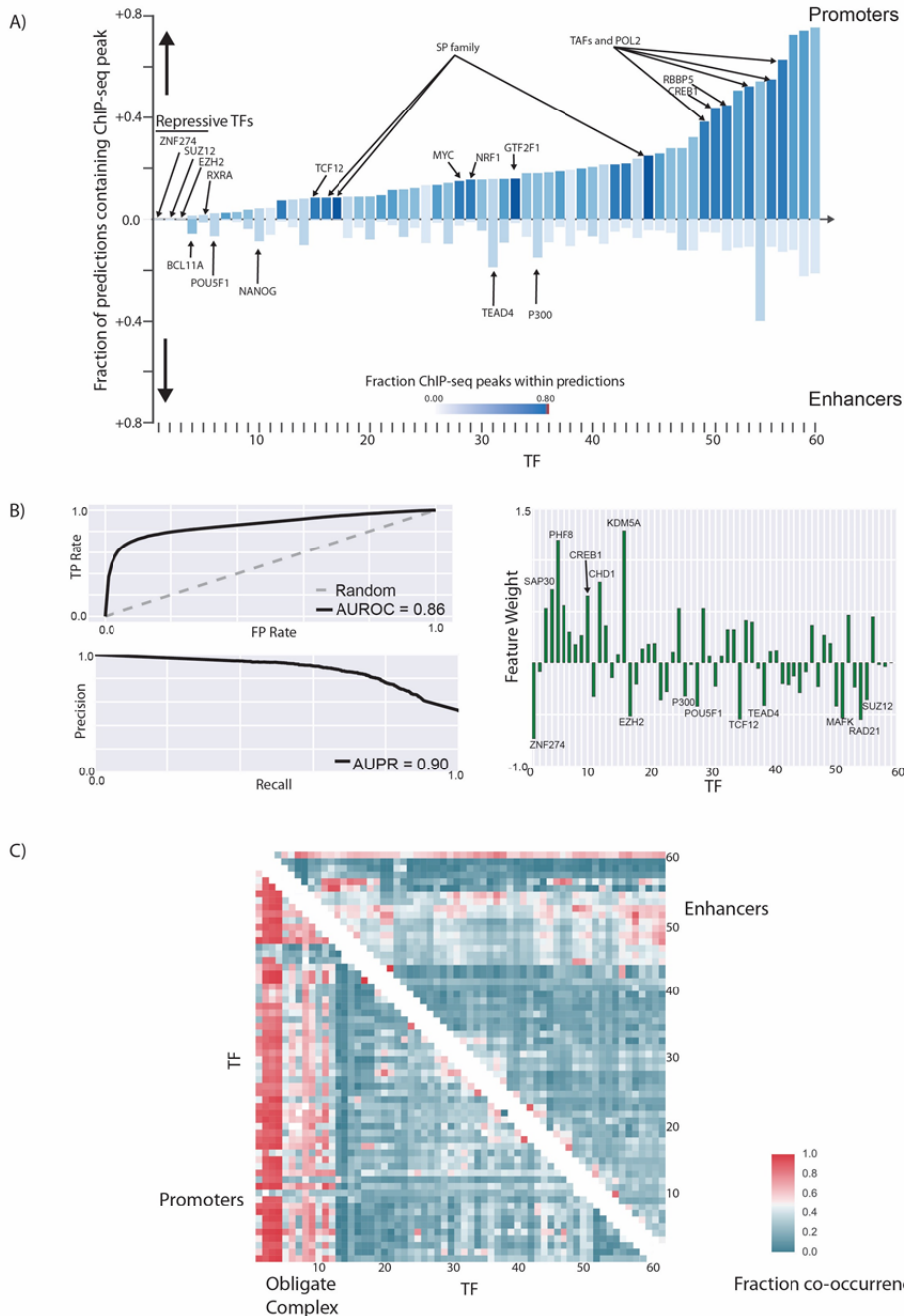
Figure 5



923
924
925
926
927
928
929
930
931
932
933
934
935
936
937

Figure 5: Enhancer Validation Experiments. A) A schematic of the enhancer validation scheme is shown. At top is third generation HIV-based self-inactivating vector (deletion in 3' LTR indicated by red triangle), with PCR-amplified test DNA (blue, two-headed arrow indicates fragment was cloned in both orientations), inserted just 5' of a basal (B) Oct4 promoter driving IRES-eGFP (green). Vector supernatant was prepared by plasmid co-transfection of 293T cells and used to transduce cellular targets and analyzed by flow cytometry a few days later. B) The fold change of gene expression of eGFP is compared between negative elements and putative enhancers chosen for experiments. The p-Value of the difference in activity is measured using a Wilcoxon signed-rank test.

Figure 6



938
939
940
941
942
943
944
945
946
947
948
949

Figure 6: Differences in TF binding patterns at enhancers and promoters. A) The fraction of predicted promoters and enhancers that overlap with ENCODE ChIP-seq peaks for different TFs in H1-hESC are shown. The names of all TFs in the figure can be viewed in Figure S20. B) The AUROC and AUPR for a logistic regression model created using the pattern of TF binding at each regulatory region to distinguish enhancers from promoters are shown. The weight of each feature in the logistic regression model can be used to identify the most important TFs that distinguish enhancers from promoters. C) The patterns of TF co-binding at active promoters and enhancers are shown. The names of all the TFs in this graph can be viewed in Figure S21.