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

Anurag Sethi[1,2], Mengting Gu[1], Emrah Gumusgoz[6], Landon Chan[3], Koon-Kiu Yan[1,2], Kevin Yip[4], Joel Rozowsky[1,2], Richard Sutton[6], and Mark Gerstein[1,2,5]


[1] Program in Computational Biology and Bioinformatics, Yale University, New Haven, Connecticut, United States of America
[2] Molecular Biophysics and Biochemistry, Yale University, New Haven, Connecticut, United States of America
[3] School of Medicine, The Chinese University Hong Kong, China
[4] Department of Computer Science, The Chinese University Hong Kong, China
[5] Department of Computer Science, Yale University, New Haven, Connecticut, United States of America
[6] Department of Internal Medicine, Section of Infectious Diseases, Yale University School of Medicine, New Haven, Connecticut, United States of America

**Abstract**

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

**Significance Statement**

Enhancers are import regulatory elements in the genome. The distance between the enhancer and its regulating genes varies between several kilobytes to megabytes, making it hard annotate enhancer region both experimentally and computationally. Here we demonstrate that by integrating epigenetic features with supervised machine learning models, we can achieve high accuracy of enhancer prediction. The match filter tool providing a general framework to identify enhancers across cell lines.

**\body**

## Introduction

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

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

In recent times, due to the advent of next generation sequencing, a number of transfection and transduction-based assays were developed to experimentally test the regulatory activity of thousands of regions simultaneously in a massively parallel fashion [31-37]. In these experiments, several plasmids that each contains a single core promoter upstream of a luciferase or GFP gene are transfected or transduced into cells. These plasmids are used to test the regulatory activity of different regions by placing one region near the core promoter in each plasmid as differences in the gene's expression occur due to the differences in the activity of the tested region. STARR-seq was one such massively parallel reporter assay (MPRA) that was used to test the regulatory activity of the fly genome in several cell-types [31, 38] and was used to identify thousands of cell-type specific enhancers and promoters. MPRAs have confirmed that active enhancers and promoters tend to be depleted of histone proteins and contain accessible DNA on which various transcription factors and cofactors bind [39, 40]. These regulatory regions also tend to be flanked by nucleosomes that contain histone proteins with certain characteristic post-translational modifications. These attributes lead to an

enriched peak-trough-peak ("double peak") signal in different ChIP-Seq experiments for various histone modifications such as acetylation on H3K27 and methylations on H3K4. The troughs in the double peak ChIP-seq signal represent the accessible DNA that leads to a peak in the DNase-I hypersensitivity (DHS) at the enhancer [41]. However, the optimal method to combine information from multiple epigenetic marks to make cell-type specific regulatory predictions remains unknown. For the first time, using data from several MPRAs, we have the ability to properly train our models based on a large number of experimentally validated enhancers and test the performance of different models for enhancer prediction using cross validation.

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

**Results**

**Aggregation of epigenetic signal to create metaprofile:**

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

These metaprofiles were initially created using the histone modification H3K27ac at active STARR-seq peaks (see Figure 1 and Methods) identified in the S2 cell-line of fly.

Approximately 70% of the active STARR-seq peaks contain an easily identifiable double peak pattern even though there is a lot of variability in the distance between the two maxima of the double peak in the ChIP-chip signal (Figure S1). Even though the minimum tends to occur in the center of these two maxima on average, the distance between the two maxima in the double peaks can vary between 300 and 1100 base pairs. During aggregation, we aligned the two maxima in the H3K27ac signal across different STARR-seq peaks, followed by interpolation and smoothening the signal before calculating the average metaprofile. In addition, an optional flipping step was performed to maintain the asymmetry in the underlying H3K27ac double peak because it may be associated with the directionality of transcription [42]. For the first time, we also calculated the dependent metaprofiles for thirty other histone marks and DHS signal by applying the same set of transformations to these datasets. The metaprofile for the histone marks associated with active regulatory regions were also double peak signals and the maxima across different histone modification signals tended to align with each other on average (Figure S2). This indicates that a large number of histone modifications tend to simultaneously co-occur on the nucleosomes flanking an active enhancer or promoter. In contrast, as expected, the DHS signal displayed a single peak at the center of the H3K27ac double peak (Figure 1). In addition, repressive marks such as H3K27me3 were depleted in these regions and the metaprofile for these regions did not contain a double peak signal (Figure S2).

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

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

We used 10-fold cross validation to assess the performance of matched filters for individual histone marks to predict active STARR-seq peaks. In Figure 2, we observe that the H3K27ac matched filter is the single most accurate feature for predicting active regulatory regions (AUROC=0.92, AUPR=0.72) identified using STARR-seq. This is consistent with the literature as H3K27ac enriched peaks are often used to predict active promoters and enhancers [23, 45, 46]. In general, several histone acetylation (H3K27ac,

H3K9ac, H4K12ac, H2BK5ac, H4K8ac, H4K5ac, H3K18ac) marks as well as the H1, H3K4me2, and DHS matched filters are the most accurate marks (see Figure 2 and Table S1) because the matched filter scores for these regions on these marks are higher for STARR-seq peaks (Figure S3). The degree to which the matched filter scores for promoters and enhancers are higher than the matched filter scores for the rest of the genome is a measure of the signal to noise ratio for regulatory region prediction in the corresponding feature's genomic track and the larger the separation between positives and negatives, the greater the accuracy of the corresponding matched filter for predicting active regulatory regions. Interestingly, the distribution of matched filter scores for STARR-seq peaks are unimodal for each histone mark except for H3K4me1, H3K4me3, and H2Av, which are bimodal (Figure S3). We also show that the matched filter scores are more accurate for predicting active STARR-seq peaks than enrichment of signal alone as they outperform the histone peaks on ROC and PR curves (Figure S4).

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

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

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

feature in the integrated model. According to the model, the acetylations (H3K27ac and H3K9ac) are the most important feature for predicting active regulatory regions from inactive regions. While the DHS matched filter performed well as an individual feature (AUPR in Figure 2), the information in DHS is redundant with the information in the histone marks as indicated by the fact that it has the lowest weight among the six features in the integrated model. We compared several other machine learning algorithms including nonlinear SVM (results not shown) to combine the machine learning models and found that they all displayed nearly similar accuracy and similar features were more important across these different models (Figure S5).

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

**Distinct epigenetic signals associated with promoters and enhancers:**

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

We created two different integrated models to learn the combination of features associated with promoters and enhancers. These integrated models outperformed the individual matched filters at predicting active enhancers and promoters (Figures 3 and S8). In addition, the weights of the individual features identified the difference in roles of

the H3K4me1 and H3K4me3 matched filter scores at discriminating active promoters and enhancers from inactive regions in the genome. The promoter-based (enhancer-based) model performed much more poorly at predicting enhancers (promoters) indicating the unique properties of these regions (Figures S10 and S11). We also created two integrated models utilizing matched filter scores for all thirty histone marks as features for predicting enhancers and promoters. The additional histone marks provided independent information regarding the activity of promoters and enhancers as these features increased the accuracy of these models (Figure S9). The weights of different features indicate that H2BK5ac again displays the most independent information for accurately predicting active enhancers and promoters (Figures S9). We observe similar trends and accuracy with several different machine learning models (Figures S8 and S9).

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

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

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

**Validation of Enhancer Prediction Models**

The ENCODE consortium has ChIP-Seq data for 60 transcription related factors in H1-hESC cell line, including a few chromatin remodelers and histone modification enzymes. Collectively we call all these transcription related factors "TF"s for simplicity. We utilized the 6 parameter integrated model to predict active enhancers and promoters in the hESC cell-line based on the epigenetic datasets measured by the ENCODE consortium. This provides us with a system to validate our enhancer prediction model as well as to study the patterns of TF binding within enhancers and promoters. Using these models, we predicted 43463 active regulatory regions, of which 22828 (52.5%) are within 2kb of

the TSS and are labeled as promoters. A large proportion of the predicted enhancers are found in the introns (30.41%) and intergenic regions (13.93%) (Figure S13). The predicted promoters and enhancers are significantly closer to active genes than might be expected randomly (Figure S14). By comparing the matched filter predicted enhancers and promoters with chromatin states predicted by chromHMM [30] and SegWay [27], we observe that a majority of the predicted enhancers and promoters are also predicted to be enhancers and promoters by chromHMM and SegWay respectively (Figures S15 to S18).

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

A total of 25 predicted intergenic enhancers were randomly selected for validation (Supplementary Table S3). These predictions were chosen randomly to ensure that these truly represented the whole spectrum of predicted enhancers and not just the top tier of predicted enhancers. Of these 25 putative enhancers, 23 were successfully amplified and cloned into the HIV vector. To measure the distribution of gene expression in the absence of enhancer, we also amplified and cloned 25 non-repetitive elements with similar length distribution that were predicted to be inactive using the same HIV vector. All positive and negative DNA elements were transduced and tested for activity in both forward and reverse strand orientations since enhancers are thought to function in an orientation-independent manner. Functional testing was performed in HOS, TZMBL, and A549 cell lines in addition to H1-hESCs.

Insertion of twelve of the 23 putative enhancers into the HIV vector resulted in a significant increase in eGFP expression (P-value < 0.05 over distribution of gene expression for negative elements) in the H1-hESCs (Supplementary Table S3). While most of the positive enhancers displayed a significant increase in gene expression irrespective of their orientation during orientation, a few elements showed significantly higher levels of gene expression in one of the orientations (Supplementary Table S4). In contrast, the negatives displayed much lower levels of gene expression typically (Figure 5 and Supplementary Figure S19). In addition, most of these elements increased gene expression of GFP in the four different cell lines even though some of the elements were preferentially active in one of the cell lines. Overall, 16 of the 23 tested predictions displayed statistically significant increase in gene expression of the reporter gene in at least one of the cell lines (Supplementary Table S3 and Supplementary Figure S19). Given the promoter specificity of enhancers in such assays, we would anticipate that some of the elements that could not be validated in this particular vector would function as enhancers in a more natural biological context.

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

We further studied the differences in TF binding at promoters and enhancers (Figure 6 and Figure S20). Most promoters and enhancers contain multiple TF-binding sites. However, the TF-binding of enhancers is more heterogeneous than promoters: in particular, more than 70% of the promoters bind to the same set of 2-3 sequence-specific TFs, which is not observed for enhancers. The majority of the promoters also contain peaks for several TATA-associated factors (TAF1, TAF7, and TBP). Overall, the high heterogeneity associated with enhancer TF-binding is consistent with the absence of a sequence code (or grammar) which can be utilized to easily identify active enhancers on a genome-wide fashion.

In Figure 6, we show that the patterns of TF binding within regulatory regions can be utilized in a logistic regression model to distinguish active enhancers from promoters with high accuracy (AUPR = 0.89, AUROC = 0.87). We were also able to identify the most important features that distinguish promoters from enhancers. In addition to TATA-box associated factors such as TAF1, TAF7, and TBP, the RNA polymerase-II binding patterns as well as chromatin remodelers such as KDM4A and PHF8 are some of the most important factors that distinguish promoters from enhancers in H1-hESC. This provides a framework that can be utilized to identify the most important TFs associated with active enhancers and promoters in each cell-type.

In Figure 6A, we show that the pattern of TF binding at promoters is different from that at enhancers and TF-binding at enhancers displaying more heterogeneity. As the set of TFs binding promoters is fairly uniform, the same pairs of TF also tend to bind together on promoters. In contrast, for enhancers, the patterns of TF co-binding is much more heterogeneous and different enhancers tend to contain different TF-pairs. This can be observed in the patterns of TF co-binding in Figures 6C and S21. These TF co-associations could lead to mechanistic insights of cooperativity between TFs. For example, similar to a previous study [57], CTCF and ZNF143 may function cooperatively as they are observed to co-occur frequently at distal regulatory regions in this study.

**Discussion**

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

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

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

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

One aspect that is discussed less frequently is the effect of core promoter on enhancer and promoter prediction. MPRAs show that the regulatory activity of enhancers and promoters in a regulatory assay depends on the core promoter used during the experiment [51]. As the transcription factors that bind to each regulatory region are thought to play a key role in core-promoter specificity [47, 51], we suspect that machine learning models that contain sequence or motif-based features may be biased towards certain transcription factor binding sites when trained with regulatory regions identified using a single-core promoter. To avoid such biases, it would be more appropriate to train models with sequence-based features when the validation experiments are performed with multiple core promoters. In the absence of validation data with multiple core

promoters, it may be more suitable to train models using epigenetic features as such models contain no sequence-based information. In comparing the predictions from such models with experiments using a single core promoter, some of the strongest predictions may be mislabeled as negatives even though they contain some regulatory activity leading to a lower accuracy estimate as shown in Figure 2.

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


/body

**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.

22.     Visel, A., et al., *Ultraconservation identifies a small subset of extremely constrained developmental enhancers.* Nat Genet, 2008. **40**(2): p. 158-60.

23.     Bonn, S., et al., *Tissue-specific analysis of chromatin state identifies temporal signatures of enhancer activity during embryonic development.* Nat Genet, 2012. **44**(2): p. 148-56.

24.     Yip, K.Y., et al., *Classification of human genomic regions based on experimentally determined binding sites of more than 100 transcription-related factors.* Genome Biol, 2012. **13**(9): p. R48.

25.     Ghandi, M., et al., *Enhanced regulatory sequence prediction using gapped k-mer features.* PLoS Comput Biol, 2014. **10**(7): p. e1003711.

26.     Heintzman, N.D., et al., *Distinct and predictive chromatin signatures of transcriptional promoters and enhancers in the human genome.* Nat Genet, 2007. **39**(3): p. 311-8.

27.     Hoffman, M.M., et al., *Unsupervised pattern discovery in human chromatin structure through genomic segmentation.* Nat Methods, 2012. **9**(5): p. 473-6.

28.     Kharchenko, P.V., et al., *Comprehensive analysis of the chromatin landscape in Drosophila melanogaster.* Nature, 2011. **471**(7339): p. 480-5.

29.     He, H.H., et al., *Nucleosome dynamics define transcriptional enhancers.* Nat Genet, 2010. **42**(4): p. 343-7.

30.     Ernst, J., et al., *Mapping and analysis of chromatin state dynamics in nine human cell types.* Nature, 2011. **473**(7345): p. 43-9.

31.     Arnold, C.D., et al., *Genome-wide quantitative enhancer activity maps identified by STARR-seq.* Science, 2013. **339**(6123): p. 1074-7.

32.     Dickel, D.E., et al., *Function-based identification of mammalian enhancers using site-specific integration.* Nat Methods, 2014. **11**(5): p. 566-71.

33.     Gisselbrecht, S.S., et al., *Highly parallel assays of tissue-specific enhancers in whole Drosophila embryos.* Nat Methods, 2013. **10**(8): p. 774-80.

34.     Kwasnieski, J.C., et al., *High-throughput functional testing of ENCODE segmentation predictions.* Genome Res, 2014. **24**(10): p. 1595-602.

35.     Melnikov, A., et al., *Systematic dissection and optimization of inducible enhancers in human cells using a massively parallel reporter assay.* Nat Biotechnol, 2012. **30**(3): p. 271-7.

36.     Murtha, M., et al., *FIREWACh: high-throughput functional detection of transcriptional regulatory modules in mammalian cells.* Nat Methods, 2014. **11**(5): p. 559-65.

37.     Patwardhan, R.P., et al., *Massively parallel functional dissection of mammalian enhancers in vivo.* Nat Biotechnol, 2012. **30**(3): p. 265-70.

38.     Yanez-Cuna, J.O., et al., *Dissection of thousands of cell type-specific enhancers identifies dinucleotide repeat motifs as general enhancer features.* Genome Res, 2014. **24**(7): p. 1147-56.

39.     Shlyueva, D., G. Stampfel, and A. Stark, *Transcriptional enhancers: from properties to genome-wide predictions.* Nat Rev Genet, 2014. **15**(4): p. 272-86.

40.     Maston, G.A., et al., *Characterization of enhancer function from genome-wide analyses.* Annu Rev Genomics Hum Genet, 2012. **13**: p. 29-57.

41.     Thurman, R.E., et al., *The accessible chromatin landscape of the human genome.* Nature, 2012. **489**(7414): p. 75-82.

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

**Figure Captions**

**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 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 identified using FIREWACh are shown.

**Figure 5: Enhancer Validation Experiments.** A) A schematic of the enhancer validation scheme is show. 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: 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 S19. 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 S20.

Figure 1

Figure 2

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) |

- - - Single Core Promoter
——— Multiple Core Promoters

# Figure 3

## A)

| Feature | AUROC | AUPR |
|---|---|---|
| H3K27ac | 0.92 (0.96) | 0.55 (0.71) |
| H3K9ac | 0.87 (0.95) | 0.19 (0.69) |
| DHS | 0.83 (0.89) | 0.28 (0.59) |
| H3K4me2 | 0.85 (0.92) | 0.15 (0.49) |
| H3K4me3 | 0.63 (0.93) | 0.06 (0.64) |
| H3K4me1 | 0.90 (0.59) | 0.36 (0.16) |
| Integrated | 0.94 (0.97) | 0.66 (0.78) |

Figure 4

**Fly-based models on mouse**



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) |

Figure 5

# Figure 6

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


**Methods**

**Creation of Metaprofile:**

We utilized the smoothed histone signal tracks provided for the S2 cell-line by the modENCODE consortium [1] to aggregate the corresponding histone signals around the STARR-seq peaks [2]. This aggregation was performed to remove noise before using the metaprofile *s(n)* for identifying active regulatory regions in the genome. The genome-wide profile for open chromatin (DNase-seq or DHS) for the S2 cell-line was calculated based on the experiments by the Stark lab [2]. To create the smoothened metaprofile, we aggregated the H3K27ac signal of active STARR-seq peaks with a noticeable "double peak" pattern within the H3K27ac signal in the S2 cell-line. All the STARR-seq peaks that overlap with DHS or H3K27ac peaks are assumed to be active regulatory regions in the genome.

To identify double peak regions, we initially identified the minimum in the H3K27ac signal track closest to the middle of the STARR-seq peaks. A minimum is accepted if it has the lowest signal within a 100 base pair region in the H3K27ac signal track. Then we proceed to identify the flanking maxima (both sides of the minimum) within a total of 2-kilo base pair region of the STARR-seq peak (1kb on each direction from the center of the STARR-seq peak). These maxima are accepted only if they have the highest signal within a 100 base pair region in the H3K27ac signal track. Approximately 70% of the active STARR-seq peaks contained an identifiable double peak within the H3K27ac signal.

After identifying the double peaks surrounding STARR-seq peaks, we aggregated the signal after aligning the maxima flanking the regulatory region. The signal track is interpolated with a cubic spline fit so that the signal track contains equal number of points for each double peak region. All interpolation and smoothing steps were performed using the scipy module in python. The aggregated signal tracks are averaged to create the metaprofile for the double peak regions. While the signal tracks are aggregated based on identifying the double peak regions in the H3K27ac signal track, the same set of operations can be performed with any epigenetic mark expected to have the double peak pattern flanking regulatory regions.

In addition, while creating the metaprofile for H3K27ac signal close to active STARR-seq peaks, we also performed the same set of transformations on other dependent epigenomic datasets (other histone marks and/or DHS signal). In this study (Figures 1 and S2), the dependent profiles for all other epigenetic datasets are calculated by averaging the corresponding signal based on identifying double peak regions within H3K27ac signal. If the signal tracks of the other epigenetic marks also tend to contain a double peak pattern in the same regions, the metaprofiles for the corresponding epigenetic marks will also contain a double peak pattern as observed in Figure S2A. However, as DHS and repressive histone marks do not contain a double peak pattern (Figure S2), these regions do not have the same epigenetic template associated with enhancers.

**Matched Filter Algorithm:**

The epigenetic signal at enhancers and promoters can be approximated as the linear superposition of background noise and the metaprofile *s(n)* learned in Figure 1 (Figure S2) for the corresponding experimental dataset. The matched filter *h(n)* is used to scan the epigenetic signal to identify the occurrence of the metaprofile pattern within different regions of the genome.  Before calculating the matched filter score, interpolation of signal is used to ensure that the scanned region contains the same number of points as the metaprofile. The matched filter process is equivalent to the computation of the cross correlation between the signal *y(n)* and the reverse of the transformed metaprofile template *s\*(N-n)* (where *N* is the total number of points in the template). In other words:

$$r(n) = \sum_{i=1}^{N} y(i) * h(i)$$

where *h(i)* is the matched filter and can be written as:
$$h(i) = s^*(N - i)$$

As shown in Figure S1, there is a large amount of variability in the span (distance between the two peaks in the histone signal) of the regulatory region in the epigenetic signal. As a result, we scan the genome with the matched filter scanning different spans of the genome (distance between the two peaks allowed to vary between 300 and 1100 base pairs) and take the highest score as the matched filter score for that region. The matched filter is the filter that recognizes any given template in the presence of noise in a signal with the highest signal-to-noise ratio [3]. In the presence of white noise alone, the matched filter score is low and follows a Gaussian distribution (negatives). The presence of the metaprofile within the signal leads to higher matched filter scores for positives.

**Statistical Learning Models**
The matched filter scores for negatives for different histone marks are unimodal that can be fit using separate Gaussian distributions. The Z-scores of matched filter scores with respect to the negatives (random regions of genome) are used as input features for training different statistical learning models. The Z-score of the matched filter score for a region (*z(i)*) is:
$$z(i) = \frac{r(i) - \mu}{\sigma}$$

where *r(i)* is the matched filter score for region *i* while *μ and σ* are the mean and standard deviation of the Gaussian fit to the matched filter scores for random regions in genome. In the main text, we discuss our results of the Support Vector Machine (SVM) model, which is one of the most versatile and successful binary classifiers [4]. We utilized a linear kernel to distinguish between the positives and negatives. The linear SVM identifies a decision boundary that maximally discriminates the epigenetic features of regulatory regions from random regions of the genome in the SVM feature vector space.

In Figure S5, we also present results for Ridge Regression [5], Random Forest [6], and Gaussian Naïve Bayes [7] models and the accuracy of different models are comparable.

Ridge regression is a linear regression technique that prevents over fitting by penalizing large weights for each feature. Random Forest is an ensemble learning method that operates by constructing a large number of decision trees and outputting the mean prediction of different decision trees. We used thousand trees for creating our enhancer and promoter prediction models. The naïve Bayes classifier is a family of simple probabilistic classifiers that assumes that all the features are independent of one another. We used scikit-learn [8] with default parameters for training and assessing the performance of all the statistical models. In general, the SVM and random forest models performed the best over all the tests and were the most flexible models.


**Assessing the Models:**

In order to assess the accuracy of matched filter for predicting enhancers and promoters, we used 10-fold cross validation. During 10-fold cross validation, the positives and negatives are randomly divided in to 10 groups each. Nine of the 10 groups are randomly combined to train the model and the predictions are tested on the 10[th] group. To evaluate the performance of trained classifiers, we performed 10-fold cross-validation on the training data and quantified our results with area under receiver-operating characteristic (ROC), and area under precision-recall (PR) curves.

In the ROC curve [9], the true positive (TP) rate is plotted against the false positive (FP) rate at different thresholds in the statistical model. The TP rate is defined as the fraction of positives identified correctly by the model (i.e., ratio of number of true positives identified by the model to the total number of positives). The FP rate is defined as the fraction of negatives identified correctly by the model (i.e., ratio of number of negatives misclassified by the model to the total number of negatives). While comparing the performance of two different classifiers in the ROC curve, the classifier with higher TP rate at the same FP rate is considered to be a better classifier. The area under the ROC is a single measure for the accuracy of a model as models with higher area under ROC are generally considered to be better models.

In the PR curve, the precision is plotted against recall at different thresholds in the statistical model. The recall is the same as the TP rate of the model (i.e., ratio of number of true positives identified by the model to the total number of real positives). The precision is the fraction of positives in the model that are correct (i.e., ratio of number of true positives identified by the model to the total number of positives according to the model). In skewed datasets with large number of negatives in comparison to positives, the FP rate can be low even when the number of false positives misclassified by the model is comparable to the number of true positives. For such skewed datasets, te area under ROC for two different models may be very similar even though they actually differ in performance with respect to their precision. Hence, the area under the PR curve is a better reflection of the performance difference between two models with similar area under ROC in skewed datasets.

In Figure 2, the positives are defined as the active peaks (intersecting with DHS or H3K27ac peaks) from a single STARR-seq experiment (singe core promoter) or the union of active peaks from multiple STARR-seq experiments (multiple core promoters). The negatives are randomly chosen regions in the genome with H3K27ac signal that had the same width distribution as the distribution of distance between double peaks near STARR-seq peaks (shown in Figure S1). We typically chose between 5 to 10x

number of negatives as compared to number of positives in Figures 2, 3, and 4 as the number of enhancers and promoters in the genome (positives) are far lesser than the number of negatives and area under PR curve is dependent on the ratio of negatives to positives during 10-fold cross validation. The matched filter score for each region is chosen as the best matched filter score with a 1500 bp region centered on each positive and negative. The matched filters are scanned with distances between 300-1100 bp before choosing the best score. While comparing the performance of the matched filter to the peak-based models of the different epigenetic marks (Figure S4), we assumed that histone (DHS) peaks that overlapped with at least 50% (10%) of the STARR-seq peak is used to rank that prediction. We used a smaller threshold for DHS peaks as they are much smaller than histone peaks. We achieved similar results with thresholds of 25% for both histone and DHS peaks. The p-value of the intersecting peak is used to rank the peak-based predictions. The modENCODE histone peaks [1] and DHS peaks [2] were compared to the matched filter scores in Figure S4.

During STARR-seq, each peak is functioning as an enhancer within the plasmid environment in S2 cell-line. However, to delineate the native role of the region, we classify them as promoters and enhancers based on their distance to the transcription start sites in the genome. In Figure 3, the active promoters are defined as active STARR-seq peaks (multiple core promoter) within 1 kb of TSS (Ensembl release 78) while enhancers were active STARR-seq peaks more than 1kb from any TSS in *Drosophila melanogaster*. While calculating the matched filter for positives and negatives, we considered the best scoring matched filter score after padding each region to 1.5kb width.

In Figure 4, the promoters are defined as FIREWACh peaks within 2 kb of TSS (GENCODE release vM4) while enhancers were FIREWACh peaks more than 2kb from any TSS. The larger distance (2 kb) for defining promoters was used because of the larger size of the mouse genome. The FIREWACh assay is performed in a transduction assay and was based on ChIP-seq peaks of a few key TFs. Hence, we did not split the FIREWACh peaks in to active and poised enhancers and promoters. The ENCODE histone and DHS datasets for mESC were used to predict enhancers and promoters in Figure 4.

**H1-hESC whole genome prediction**

To predict enhancers and promoters on the whole genome, we utilized the 6 parameter machine learning model shown in Figure 2. The histone and DHS signals from ENCODE consortium [10] were used to predict enhancers and promoters in H1-hESC. The histone signals were converted to log fold enrichment (with respect to control signal) before we scanned it with the matched filter. There were 43463 active regulatory regions predicted in the human genome (< 2% of genome). All regions within 2kb of TSS were annotated as promoters while active regulatory regions that were more than 2kb from TSS were annotated as enhancers. The distribution of the expression of closest gene (GENCODE v19 TSS) from ENCODE RNA-seq dataset [10] for H1-hESC was compared to the expression of all genes from H1-hESC. The Wilcoxon test was used to measure the significance of changes in gene expression.

**Overlap with chromatin state predicted by chromHMM and SegWay**

4

We compared the promoter and enhancer predictions for the H1-hESC cell-line with the chromatin states for the H1-hESC cell-line predicted by chromHMM and SegWay. The chromatin states for H1-hESC were downloaded from the ENCODE portal. The prediction is considered to be overlapping with the corresponding chromatin state if more than 50% of the predicted enhancer or promoter is labeled as the same chromatin state.

**Enhancer Validation Experiment**

**Cell lines**

WA01 or H1 hESC was obtained from WiCell and maintained feeder-free on matrigel-coated plates in mTESR1medium (StemCell Technologies) supplemented with penicillin and streptomycin. Roughly once weekly cell colonies were dissociated using dispase and absence of differentiation was confirmed by visual inspection and periodically staining cells using anti-SSEA4 conjugated to FITC and performing flow cytometry. Other cell types (HOS and A549 obtained from ATCC and TZMbl from the AIDS Reagent Repository) were maintained in DMEM supplemented with 10% fetal calf serum and passaged twice weekly using trypsin.

**Preparation of HIV vector, cellular transduction and analysis**

Self-inactivating (SIN) HIV vector pFG12 was modified in that the UBC promoter driving eGFP along with the WPRE was removed and replaced with a 1.4 kb IRES-eGFP cassette. Upstream of the IRES a 142 bp basal Oct 4 promoter (5'-CCTCCCTCTCCTCCACCCATCCAGGGGGCGGGGCCAGAGGTCAAGGCTAGTGGG TGGGACTGGGGAGGGAGAGAGGGGTTGAGTAGTCCCTTCGCAAGCCCTCATTTCA CCAGGCCCCCGGCTTGGGGCGCCTTCCTTCCCC-3'; coordinates on chromosome 6, negative strand: 31138398-31138539) was inserted, which overlaps with the TSS of *Oct4* but not with the coding sequence. A unique Xba 1 site was present just upstream of the basal Oct4 promoter, for cloning of test insert DNA fragments. Each test DNA fragment was amplified from genomic DNA using nested PCR and Takara LA enzyme. Typical initial PCR amplification conditions were $98^{o}$C for 10 s, $55^{o}$C for 15 s, and $68^{o}$C for 3 min for 30 cycles using 100-200 ng of genomic DNA, with the annealing temperature being variable depending upon the $T_m$ of the primer pair. For the second (internal) round of PCR, only 1-2% of the original product was used under similar PCR conditions, but for 15 cycles.

PCR products were individually cloned into TOPO pCRII-blunt vector (Invitrogen) and insert identity confirmed by both restriction digests and dideoxy sequencing. All DNA inserts were cloned into the unique Xba 1 site of the HIV vector described above using compatible cohesive ends, and in each case both orientations of the insert within the vector were confirmed by appropriate restriction digests.

HIV vector supernatants were prepared by co-transfecting 35 mm tissue culture wells of 293T cells (~75-80% confluence), each with 5 $\mu$g of HIV transfer vector (HIV-TV) with DNA element of interest, HIV packaging vector, and pME VSV G (encoding Indiana strain VSV G). After 48-72 hours, vector supernatant was harvested, centrifuged at 3000 x *g* for 10 min, and stored at $-80^{o}$ C until use.

In order to transduce the WA01 hESC, cells were first lifted using dispase, washed extensively, and plated in the presence of ROCK Inhibitor Y-27632 (StemCell Technologies) on matrigel-coated plates. After a few hours, cells were transduced for 4-6 h with lentiviral vector supernatant, After 48-72 h single cell suspensions were again prepared using dispase and Y-27632 and cells were analyzed for eGFP expression as described above, collecting 10,000 events. For all other cell lines, cells were plated the day before in 12 well format, transduced using the indicated amounts of vector supernatant, refed the following day, and analyzed for eGFP expression 48-72 h later, as described above.

The fold change of inactive elements was used to calculate the background distribution of inactive elements. This was fit to a normal distribution and putative enhancers that displayed higher activity than expected by chance (p-value < 0.05) were considered to be active in the cell-line. This was done for the forward and reverse directions separately and elements that were positive in either orientation were considered to be active.

**H1-hESC TF binding**

To measure the differences in TF binding and co-binding patterns at promoters and enhancers, we overlapped the ChIP-seq peaks from ENCODE with our predicted enhancers and promoters using intersectBed. The two regions were considered to be overlapping if at least 25% of the ChIP-seq peak was overlapping with the predicted enhancer or promoter.

**Table S1 – Performance of matched filter models with single epigenetic feature for all STARR-seq peaks (multiple core promoters)**

| Feature | AUROC | AUPR |
|---|---|---|
| H3K27ac | 0.95 | 0.90 |
| H3K4me1 | 0.70 | 0.59 |
| H3K4me2 | 0.91 | 0.79 |
| H3K4me3 | 0.84 | 0.76 |
| H3K9ac | 0.92 | 0.85 |
| H4K12ac | 0.92 | 0.86 |
| H3 | 0.80 | 0.70 |
| H1 | 0.88 | 0.81 |
| H2BK5ac | 0.94 | 0.90 |
| H4K8ac | 0.88 | 0.79 |
| H4K5ac | 0.87 | 0.79 |
| H4K16ac | 0.89 | 0.72 |
| H3K18ac | 0.90 | 0.84 |
| H3K9me1 | 0.71 | 0.61 |
| H3K79me2 | 0.79 | 0.58 |
| H4K27me2 | 0.81 | 0.68 |
| H2Av | 0.66 | 0.57 |
| H3K27me3 | 0.83 | 0.64 |
| H3K23ac | 0.66 | 0.46 |
| H3K79me3 | 0.70 | 0.51 |
| H3K27me1 | 0.64 | 0.43 |
| H4 | 0.67 | 0.49 |
| H3K36me1 | 0.54 | 0.41 |
| H3K9me3 | 0.59 | 0.42 |
| H3K9me2 | 0.60 | 0.41 |
| H3K36me3 | 0.57 | 0.38 |
| H4K20me1 | 0.47 | 0.31 |
| H3K79me1 | 0.47 | 0.30 |

**Table S2 – Performance of matched filter models with single epigenetic feature for promoters and enhancers (multiple core promoters). Numbers within (outside) parenthesis are accuracy of models for predicting promoters (enhancers).**

| Feature | AUROC | AUPR |
|---|---|---|
| H3K27ac | 0.91 (0.96) | 0.60 (0.73) |
| H3K4me1 | 0.88 (0.60) | 0.42 (0.16) |
| H3K4me2 | 0.84 (0.92) | 0.21 (0.48) |
| H3K4me3 | 0.62 (0.92) | 0.09 (0.65) |
| H3K9ac | 0.85 (0.94) | 0.24 (0.70) |
| H4K12ac | 0.90 (0.93) | 0.33 (0.58) |
| H3 | 0.78 (0.83) | 0.26 (0.48) |
| H1 | 0.83 (0.92) | 0.36 (0.61) |
| H2BK5ac | 0.91 (0.96) | 0.59 (0.70) |
| H4K8ac | 0.90 (0.86) | 0.55 (0.37) |
| H4K5ac | 0.89 (0.86) | 0.52 (0.41) |
| H4K16ac | 0.90 (0.90) | 0.52 (0.40) |
| H3K18ac | 0.90 (0.88) | 0.60 (0.47) |
| H3K9me1 | 0.53 (0.81) | 0.09 (0.44) |
| H3K79me2 | 0.70 (0.83) | 0.10 (0.27) |
| H4K27me2 | 0.68 (0.85) | 0.19 (0.44) |
| H2Av | 0.63 (0.78) | 0.15 (0.36) |
| H3K27me3 | 0.81 (0.86) | 0.20 (0.36) |
| H3K23ac | 0.55 (0.71) | 0.07 (0.20) |
| H3K79me3 | 0.61 (0.74) | 0.08 (0.23) |
| H3K27me1 | 0.72 (0.57) | 0.12 (0.12) |
| H4 | 0.69 (0.68) | 0.13 (0.21) |
| H3K36me1 | 0.75 (0.58) | 0.19 (0.18) |
| H3K9me3 | 0.59 (0.64) | 0.11 (0.15) |
| H3K9me2 | 0.62 (0.63) | 0.09 (0.15) |
| H3K36me3 | 0.60 (0.62) | 0.09 (0.14) |
| H4K20me1 | 0.55 (0.50) | 0.07 (0.10) |
| H3K79me1 | 0.54 (0.58) | 0.06 (0.12) |

**Table S3 The results of the validation experiment for 25 putative enhancers in four different cell lines**

| Region | H1-hESC | HOS | A549 | TZMBL |
|---|---|---|---|---|
| chr1:1953310-192546069 | Positive | Positive | Positive | Positive |
| chr2:231809337-231809988 | Negative | Positive | Positive | Positive |
| chr9:134224987-134225644 | - | - | - | - |
| chr11:65679112-61679919 | Positive | Positive | Positive | Positive |
| chr12:125039037-125040700 | Positive | Positive | Positive | Positive |
| chr13:113921562-113922944 | Positive | Positive | Positive | Positive |
| chr14:77422602-77423265 | Positive | Positive | Positive | Positive |
| chr17:2929462-2930394 | Positive | Positive | Positive | Positive |
| chr17:72390462-72391344 | - | - | - | - |
| chr22:31662162-31663116 | Negative | Positive | Positive | Positive |
| chr1:54839458-54841157 | Negative | Positive | Negative | Positive |
| chr3:128150669-128152511 | Positive | Negative | Negative | Negative |
| chr4:6246837-6247511 | Positive | Positive | Positive | Positive |
| chr7:1956626-1958036 | Positive | Negative | Positive | Positive |
| chr7:73448387-73448811 | Negative | Negative | Positive | Negative |
| chr9:132976212-132977003 | Negative | Positive | Positive | Positive |
| chr9:138892812-1338893419 | Positive | Negative | Negative | Negative |
| chr11:44307337-44308437 | Negative | Negative | Positive | Negative |
| chr12:52536500-52539000 | Negative | Negative | Negative | Negative |
| chr13:24121112-24121886 | Positive | Positive | Positive | Positive |
| chr14:75905362-75907344 | Positive | Negative | Positive | Negative |
| chr18:12271615-12272169 | Negative | Positive | Positive | Positive |
| chr19:6235287-6237180 | Positive | Negative | Positive | Negative |
| chr22:44243837-44244786 | Negative | Negative | Negative | Negative |
| chr22:45986287-45987069 | Negative | Negative | Negative | Negative |
| Overall | 13/23 | 13/23 | 16/23 | 13/23 |

**Table S4 The fold change of gene expression as compared to control sequences in the forward as well as reverse directions for the 25 putative enhancers.**

| Element | H1-hESC | HOS | A549 | TZMBL |
|---|---|---|---|---|
| chr1:1953310-192546069 | 3.06, 7.55 | 18.67, 60.75 | 3, 19.9 | 5.67, 9.67 |
| chr2:231809337-231809988 | 0. 1.06 | 6.33, 3.83 | 3.21, 0.48 | 3.58, 2.08 |
| chr9:134224987-134225644 | - | - | - | - |
| chr11:65679112-61679919 | 2.86, 2.45 | 8.17,25.83 | 14.2, 2.42 | 5.17, 9.75 |
| chr12:125039037-125040700 | 0, 2.24 | 11.17, 11.67 | 1.31, 4.9 | 6.58, 8.25 |
| chr13:113921562-113922944 | 1.20, 4.49 | 18.67, 9.83 | 6.1, 1.1 | 8.25, 5.75 |
| chr14:77422602-77423265 | 11.84, 2.04 | 34.58, 3.5 | 0.24, 0.24 | 10, 0.55 |
| chr17:2929462-2930394 | 0, 11.63 | 0.92, 37.5 | 0.71, 54.5 | 0.33, 6.92 |
| chr17:72390462-72391344 | - | - | - | - |
| chr22:31662162-31663116 | 0, 1.02 | 1.83, 7.0 | 2.4, 2.1 | 0.92, 1.25 |
| chr1:54839458-54841157 | 0, 1.80 | 10.58, 1.33 | 1.8, 0.12 | 2.58, 0.12 |
| chr3:128150669-128152511 | 2.24, 1.78 | 2.17, 1.42 | 0.24, 0.25 | 0.48, 1.17 |
| chr4:6246837-6247511 | 11.63, 0.88 | 40.75, 1 | 43.75, 0.79 | 5.5, 0.16 |
| chr7:1956626-1958036 | 6.53, 0 | 0.83, 1.19 | 29.73, 1.11 | 14.3, 0 |
| chr7:73448387-73448811 | 0, 1.73 | 0.97, 1.36 | 1.64, 2.19 | 0.57, 1.21 |
| chr9:132976212-132977003 | 0.90, 0.88 | 0.51, 6.71 | 0.36, 14.93 | 0.93, 6.3 |
| chr9:138892812-1338893419 | 1.82, 0 | 0.66, 0.51 | 0.88, 0.72 | 0.46, 0.34 |
| chr11:44307337-44308437 | 0, 0 | 0.89, 0.85 | 0, 5.48 | 0, 1.2 |
| chr12:52536500-52539000 | 0. 0.42 | 0.16, 1.34 | 0.53, 0.52 | 1, 0.93 |
| chr13:24121112-24121886 | 3.24, 0.39 | 4.79, 7.34 | 11.09, 38.36 | 4.8, 4.6 |
| chr14:75905362-75907344 | 4.06, 0 | 2.05, 1.78 | 7.34, 2.19 | 1, 1.1 |
| chr18:12271615-12272169 | 0.42, 0.44 | 2.74, 3.15 | 6.44, 4.38 | 2.5, 4.1 |
| chr19:6235287-6237180 | 6.72, 0.97 | 1.15, 0.16 | 23.97, 0.68 | 0.81, 0 |
| chr22:44243837-44244786 | 0.82, 0.89 | 0.12, 0 | 0.20, 0.01 | 0.99, 1.02 |
| chr22:45986287-45987069 | 1.88, 0.46 | 0.19, 0 | 0.16, 0.07 | 1.08, 0.87 |

**Figure Captions:**

**Figure S1: Variability in double peak pattern.** A) The frequency of distance between the two maxima in a double peak flanking active STARR-seq peaks is plotted. B) The symmetricity of the double peak pattern is plotted. The ratio of the distance between the two peaks to the ratio between one of the maxima and the minima is plotted. While there is large amount of variability in the distance between the two peaks (mostly between 300-1100 bp), the trough in the double peak tends to occur in the center of the two peaks.

**Figure S2: Metaprofile for different epigenetic marks.** The metaprofile around active STARR-seq peaks is plotted for different epigenetic marks. Histone marks that are enriched near STARR-seq peaks display the characteristic double peak pattern shown in A) due to the depletion of histone proteins at active regulatory regions. In addition, DHS displays a single peak at the center of these regulatory regions as shown in A). B) On the other hand, no such double peak pattern is observed on depleted histone marks at STARR-seq peaks.

**Figure S3: Histogram of matched filter scores.** The probability density of matched filter scores for different epigenetic marks for STARR-seq peaks (positives) and random regions of the genome (negatives) with H3K27ac signal. In most cases, the matched filter scores for positives and negatives are Gaussian curves. The amount of overlap between these two curves determines the accuracy of the matched filter for predicting STARR-seq peaks using thematched filters for the corresponding epigenetic feature.

**Figure S4: Accuracy of matched filter and peak-based models.** The performance of the matched filters of different epigenetic marks and the peak-based models for predicting all STARR-seq peaks is compared here using 10-fold cross validation. A) The numbers within the parentheses refer to the AUROC and AUPR for predicting the STARR-seq peaks (multiple core promoters) with histone peaks while the numbers outside the parentheses refer to the AUROC and AUPR for the matched filter model. B) The individual ROC and PR curves for each matched filter and the peak-based model are shown.

**Figure S5: Comparison of different statistical models.** The performance of the different statistical models to integrate the information from six epigenetic features is shown. A) The numbers within the parentheses refer to the AUROC and AUPR for predicting the STARR-seq peaks (single core promoter) with histone peaks while the numbers outside the parentheses refer to the AUROC and AUPR for predicting STARR-seq peaks identified after combining multiple core promoters. B) The individual ROC and PR curves for each statistical model. C) The contribution of the matched filter score for each epigenetic feature to the different integrated models.

**Figure S6: Comparison of different statistical models for 30-feature model.** The performance of the different statistical models to integrate the information from 30 epigenetic features is shown. A) The numbers within the parentheses refer to the AUROC and AUPR for predicting the STARR-seq peaks (single core promoter) with histone peaks while the numbers outside the parentheses refer to the AUROC and AUPR for predicting STARR-seq peaks identified after combining multiple core promoters. B) The individual ROC and PR curves for each statistical model. C) The contribution of the matched filter score for each epigenetic feature to the different integrated models.

**Figure S7: Histogram of matched filter scores for chosen features in promoters and enhancers.** A) The histogram of matched filter scores for small set of epigenetic features on

promoters is compared to random regions of the genome. B) The histogram of matched filter scores for small set of epigenetic features on enhancers is compared to random regions of the genome.

**Figure S8: Comparison of different statistical models for predicting enhancers and promoters.** The performance of the different statistical models to integrate the information from six epigenetic features for promoter and enhancer prediction is shown. A) The numbers within the parentheses refer to the AUROC and AUPR for predicting the promoters with histone peaks while the numbers outside the parentheses refer to the AUROC and AUPR for predicting enhancers. The promoters and enhancers from multiple STARR-seq experiments with different core promoters are merged in this analysis. B) The individual ROC and PR curves for each statistical model is shown. The contribution of the matched filter score for each epigenetic feature to the different integrated models for promoter prediction (C) and enhancer prediction (D) are shown.

**Figure S9: Comparison of different statistical models for predicting enhancers and promoters.** The performance of the different statistical models to integrate the information from thirty epigenetic features for promoter and enhancer prediction is shown. A) The numbers within the parentheses refer to the AUROC and AUPR for predicting the promoters with histone peaks while the numbers outside the parentheses refer to the AUROC and AUPR for predicting enhancers. The promoters and enhancers from multiple STARR-seq experiments with different core promoters are merged in this analysis. B) The individual ROC and PR curves for each statistical model is shown. The contribution of the matched filter score for each epigenetic feature to the different integrated models for promoter prediction (C) and enhancer prediction (D) are shown.

**Figure S10: Accuracy of enhancer-trained matched filter and statistical models for promoter prediction.** The performance of the enhancer-trained matched filters of different epigenetic marks and statistical models for predicting active promoters is compared. A) The AUROC and AUPR for each matched filter and statistical model are tabulated. The individual ROC and PR curves for each matched filter (B) and each statistical model (C) are shown.

**Figure S11: Accuracy of promoter-trained matched filter and statistical models for enhancer prediction.** The performance of the promoter-trained matched filters of different epigenetic marks and statistical models for predicting active enhancers is compared. A) The AUROC and AUPR for each matched filter and statistical model are tabulated. The individual ROC and PR curves for each matched filter (B) and each statistical model (C) are shown.

**Figure S12: Transferability of models across cell-lines.** The performance of the BG3-trained matched filters of different epigenetic marks and statistical models for predicting active promoters and enhancers are compared. A) The AUROC and AUPR for each matched filter and statistical model are tabulated. The individual ROC and PR curves for each matched filter (B) and each statistical model (C) are shown.

**Figure S13: Location of H1-hESC predictions.** A) The probability density of the distance of the predicted promoter and enhancer from the closest TSS is shown. B) The location of the enhancers and promoters on genomic elements are shown. Promoters are defined as TSS +/- 2kb. All TSS, UTR, exons, introns, and intergenic elements are calculated based on GENCODE 19 definitions [11]. A regulatory region is considered to overlap with the elements if more than 50% of the matched filter region overlaps with the corresponding element in B.

**Figure S14: Gene expression of closest gene.** The distribution of gene expression of gene closest to the enhancer/promoters are plotted and compared to the gene expression of all genes in H1-hESC. A Wilcoxon test shows that P-value for differences in gene expression of genes close to enhancers and promoters are significantly higher than expression of all genes in H1-hESC (< $10^{-100}$ each).

**Figure S15: Overlap of predicted promoters with chromatin states predicted by ChromHMM.** The promoters predicted to be active by matched filter in H1-hESC cell line are compared with the chromatin states predicted using chromHMM. Most of the matched filter promoters are also predicted to be either strong or weak promoters by chromHMM while some of the other matched filter promoters are labeled as weak enhancers or transcription related elements in chromHMM. However, very few inactive regions and insulators are predicted to be promoters by matched filter. However, the boundaries of the elements can be very different as chromHMM promoters can also be tens of kilobases in length.

**Figure S16: Overlap of predicted enhancers with chromatin states predicted by ChromHMM.** The enhancers predicted to be active by matched filter in H1-hESC cell line are compared with the chromatin states predicted using chromHMM. Most of the matched filter enhancers are also predicted to be either strong or weak enhancers by chromHMM while some of the other matched filter promoters are labeled as transcription related elements in chromHMM. However, very few inactive regions and insulators are predicted to be promoters by matched filter.

**Figure S17: Overlap of predicted promoters with chromatin states predicted by SegWay.** The promoters predicted to be active by matched filter in H1-hESC cell line are compared with the chromatin states predicted using SegWay. Most of the matched filter promoters are also predicted to be either active promoters by SegWay while some of the other matched filter promoters are labeled as promoter flanking or transcription related elements in SegWay. However, very few inactive regions and insulators are predicted to be promoters by matched filter. However, the boundaries of the elements can be very different.

**Figure S18: Overlap of predicted enhancers with chromatin states predicted by SegWay.** The enhancers predicted to be active by matched filter in H1-hESC cell line are compared with the chromatin states predicted using SegWay. Most of the matched filter enhancers are also predicted to be promoters or enhancers by SegWay while some of the other matched filter enhancers are labeled as either promoter flanking or transcription related elements in SegWay. However, very few inactive regions and insulators are predicted to be promoters by matched filter.

**Figure S19. Activity of putative enhancers in three different cell-lines.** While the enhancers were predicted in H1-hESC, the activity of these enhancers is compared in three other cell-lines and the enhancers are active in these cell-lines too.

**Figure S20: Overlap of TF binding site with predicted promoters/enhancers.** The fraction of promoters and enhancers that overlap with different TF ChIP-seq peaks in H1-hESC are plotted. The color of the bar is plotted based on the fraction of ChIP-seq peaks for corresponding TF that overlap with the promoter/enhancer. The difference in patterns of TF binding was used to create models that distinguish enhancers from promoters (Figure 5B).

**Figure S21: Patterns of co-TF binding on enhancers and promoters.** The patterns of TF co-occurrence on a single matched filter prediction around promoters and enhancers are plotted. The differences between co-TF binding at enhancers and promoters can be used to gain some mechanistic insight into TF cooperativity.

**References:**

1. mod, E.C., et al., *Identification of functional elements and regulatory circuits by Drosophila modENCODE.* Science, 2010. **330**(6012): p. 1787-97.
2. Arnold, C.D., et al., *Genome-wide quantitative enhancer activity maps identified by STARR-seq.* Science, 2013. **339**(6123): p. 1074-7.
3. Kumar, V.B.V.K., A. Mahalanobis, and R.D. Juday, *Correlation Pattern Recognition.* 2005.
4. Blanchard, G., O. Bousquet, and P. Massaer, *Statistical performance of support vector machines.* Ann. Statist., 2008. **36**: p. 489-531.
5. Hoerl, A.E. and R.W. Kennard, *Ridge Regression: Biased Estimation for Nonorthogonal Problems.* Technometrics, 1970. **12**(1): p. 55--67.
6. Breiman, L., *Random Forests.* Machine Learning, 2001. **45**(1): p. 5--32.
7. Stuart, R. and P. Norvig, *Artificial Intelligence: A Modern Approach.* 2nd ed. 2003.
8. Pedregosa, F., et al., *Scikit-learn: Machine Learning in Python.* Journal of Machine Learning Research, 2011. **12**: p. 2825--2830.
9. Davis, J. and M. Goadrich, *The Relationship Between Precision-Recall and ROC Curves.* Proceedings of the 23rd international conference on Machine Learning, 2006: p. 233-240.
10. Consortium, E.P., *An integrated encyclopedia of DNA elements in the human genome.* Nature, 2012. **489**(7414): p. 57-74.
11. Harrow, J., et al., *GENCODE: the reference human genome annotation for The ENCODE Project.* Genome Res, 2012. **22**(9): p. 1760-74.

Figure S1

Figure S2

Figure S3

# Figure S4



A)

| Feature | AUROC | AUPR |
|---|---|---|
| H3K27ac | 0.92 (0.83) | 0.72 (0.63) |
| H3K9ac | 0.89 (0.77) | 0.52 (0.39) |
| DHS | 0.86 (0.77) | 0.58 (0.67) |
| H3K4me2 | 0.87 (0.75) | 0.41 (0.34) |
| H3K4me3 | 0.73 (0.64) | 0.32 (0.28) |
| H3K4me1 | 0.80 (0.72) | 0.46 (0.39) |

B)

Figure S5

A)

| Model | AUROC | AUPR |
|---|---|---|
| Random Forest | 0.96 (0.95) | 0.91 (0.79) |
| Ridge Regression | 0.95 (0.94) | 0.90 (0.77) |
| Linear SVM | 0.96 (0.95) | 0.91 (0.78) |
| Naive Bayes | 0.95 (0.93) | 0.89 (0.72) |

- - - Single Core Promoter
—— Multiple Core Promoters

B)



C)

# Figure S6

**A)**

| Model | AUROC | AUPR |
|---|---|---|
| Random Forest | 0.97 (0.98) | 0.92 (0.84) |
| Ridge Regression | 0.96 (0.97) | 0.92 (0.81) |
| Linear SVM | 0.97 (0.97) | 0.93 (0.83) |
| Naive Bayes | 0.94 (0.95) | 0.90 (0.72) |

**B)**



**C)**

Figure S7

A) Promoters

H3K27ac  H3K4me1  H3K4me3  H2Av  H2BK5ac

B) Enhancers

Probability density (x10⁻³)

Matched Filter score

# Figure S8

## A)

| Model | AUROC | AUPR |
|---|---|---|
| Random Forest | 0.95 (0.97) | 0.66 (0.78) |
| Ridge Regression | 0.94 (0.97) | 0.65 (0.78) |
| Linear SVM | 0.94 (0.97) | 0.66 (0.78) |
| Naive Bayes | 0.92 (0.97) | 0.54 (0.79) |

## B)



## C)

### Promoters



## D)

### Enhancers

Figure S9

A)

| Model | AUROC | AUPR |
|---|---|---|
| Random Forest | 0.95 (0.98) | 0.73 (0.75) |
| Ridge Regression | 0.94 (0.97) | 0.69 (0.81) |
| Linear SVM | 0.95 (0.98) | 0.74 (0.81) |
| Naive Bayes | 0.91 (0.95) | 0.62 (0.77) |

B)



C)

Promoters



D)

Enhancers

# Figure S10

A)

| Feature | AUROC | AUPR |
|---|---|---|
| H3K27ac | 0.94 | 0.92 |
| H3K9ac | 0.93 | 0.92 |
| DHS | 0.89 | 0.89 |
| H3K4me2 | 0.91 | 0.87 |
| H3K4me3 | 0.91 | 0.90 |
| H3K4me1 | 0.57 | 0.59 |
| Random Forest | 0.85 | 0.84 |
| Ridge Regression | 0.82 | 0.80 |
| Linear SVM | 0.79 | 0.80 |
| Naive Bayes | 0.95 | 0.93 |

Figure S11

A)

| Feature | AUROC | AUPR |
|---|---|---|
| H3K27ac | 0.88 | 0.86 |
| H3K9ac | 0.86 | 0.73 |
| DHS | 0.82 | 0.77 |
| H3K4me2 | 0.83 | 0.70 |
| H3K4me3 | 0.58 | 0.46 |
| H3K4me1 | 0.89 | 0.83 |
| Random Forest | 0.91 | 0.82 |
| Ridge Regression | 0.89 | 0.80 |
| Linear SVM | 0.90 | 0.86 |
| Naive Bayes | 0.88 | 0.83 |

B)



C)

Figure S12

A)

| Feature | AUROC | AUPR |
|---|---|---|
| H3K27ac | 0.88 (0.94) | 0.78 (0.87) |
| H3K9ac | 0.86 (0.94) | 0.56 (0.86) |
| H3K4me2 | 0.84 (0.92) | 0.53 (0.79) |
| H3K4me3 | 0.58 (0.91) | 0.28 (0.84) |
| H3K4me1 | 0.89 (0.58) | 0.74 (0.44) |
| Random Forest | 0.91 (0.94) | 0.81 (0.90) |
| Ridge Regression | 0.93 (0.96) | 0.84 (0.90) |
| Linear SVM | 0.92 (0.95) | 0.84 (0.90) |
| Naive Bayes | 0.92 (0.96) | 0.82 (0.91) |

B)



C)

Figure S14

# Figure S15

# Figure S16

Figure S17

# Figure S18



Active Promoter · Weak Promoter · Strong Enhancer · Weak Enhancer · Weak Enhancer (DHS) · Insulator · Low Activity · Transcription Associated · Promoter Flanking · Repressed · Heterochromatin · Repetitive

Overlap with SegWay categories

Enhancer Predictions (Ranked by Matched Filter)

Figure S19

A)



HOS

B)



A549

C)



TZMBL

Figure S20