

RESEARCH ARTICLE

ReFeaFi: Genome-wide prediction of regulatory elements driving transcription initiation

Ramzan Umarov^{1*}, Yu Li², Takahiro Arakawa³, Satoshi Takizawa³, Xin Gao^{4*}, Erik Arner^{1,3*}

1 Graduate School of Integrated Sciences for Life, Hiroshima University, Higashi-Hiroshima, Japan, **2** Department of Computer Science and Engineering (CSE), The Chinese University of Hong Kong (CUHK), Hong Kong, People's Republic of China, **3** Laboratory for Applied Regulatory Genomics Network Analysis, RIKEN Center for Integrative Medical Sciences, Yokohama, Kanagawa, Japan, **4** King Abdullah University of Science and Technology, Computational Bioscience Research Center, Computer, Electrical and Mathematical Sciences and Engineering Division, Thuwal, Saudi Arabia

* umarov@hiroshima-u.ac.jp (RU); xin.gao@kaust.edu.sa (XG); erik.arner@riken.jp (EA)



OPEN ACCESS

Citation: Umarov R, Li Y, Arakawa T, Takizawa S, Gao X, Arner E (2021) ReFeaFi: Genome-wide prediction of regulatory elements driving transcription initiation. *PLoS Comput Biol* 17(9): e1009376. <https://doi.org/10.1371/journal.pcbi.1009376>

Editor: Andrey Rzhetsky, University of Chicago, UNITED STATES

Received: June 22, 2021

Accepted: August 23, 2021

Published: September 7, 2021

Copyright: © 2021 Umarov et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: Trained models, code to generate them, and code for analysis and figures described in this study are available at the following GitHub repository: <https://github.com/umarov90/ReFeaFi>. All the data used in training and validation are publicly available through databases referenced in the manuscript.

Funding: The author(s) received no specific funding for this work.

Competing interests: The authors have declared that no competing interests exist.

Abstract

Regulatory elements control gene expression through transcription initiation (promoters) and by enhancing transcription at distant regions (enhancers). Accurate identification of regulatory elements is fundamental for annotating genomes and understanding gene expression patterns. While there are many attempts to develop computational promoter and enhancer identification methods, reliable tools to analyze long genomic sequences are still lacking. Prediction methods often perform poorly on the genome-wide scale because the number of negatives is much higher than that in the training sets. To address this issue, we propose a dynamic negative set updating scheme with a two-model approach, using one model for scanning the genome and the other one for testing candidate positions. The developed method achieves good genome-level performance and maintains robust performance when applied to other vertebrate species, without re-training. Moreover, the unannotated predicted regulatory regions made on the human genome are enriched for disease-associated variants, suggesting them to be potentially true regulatory elements rather than false positives. We validated high scoring “false positive” predictions using reporter assay and all tested candidates were successfully validated, demonstrating the ability of our method to discover novel human regulatory regions.

Author summary

Identification of regulatory elements (promoters and enhancers) is important for understanding gene expression patterns. The set of promoters and enhancers is not complete for non-model organisms and even for the human genome there are still unannotated regions, such as alternative promoters for the known genes or promoters that are only expressed in a small fraction of cells or under specific conditions. Despite the development of experimental techniques, the regulatory regions annotation remains expensive and

laborious and computational methods can speed up this process by providing candidates for the validation. We developed an easy-to-use tool capable of regulatory regions annotation in eukaryotic genomes. The developed method reduces the number of false positives made by including difficult samples in the training set. The method consists of two deep learning models, where one model scans the genome and identifies putative regulatory regions while the other model pinpoints the Transcription Start Site (TSS) location within the identified region. The predicted regions were validated using reporter assay, finding previously unknown regulatory regions in the human genome. The trained model achieved good genome-wide performance and was supported by meaningful extracted biological features.

Introduction

The study of gene regulation is primarily concerned with two classes of regulatory elements: promoters, which define the Transcription Start Site (TSS), and enhancers, that amplify the transcription [1]. The TSS is the first nucleotide that is copied at the 5' end of the corresponding mRNA. A core promoter is a minimal promoter region that typically spans several hundred bp up- and downstream of a TSS and is capable of initiating basal transcription [2]. Core promoters have complex and gene-specific architectures consisting of unique compositions of binding sites for Transcription Factors (TFs) involved in specific regulation of transcription. Transcription is further stimulated by enhancer elements, which can be located at a long distance from the target core promoter. These distal locations are able to affect the transcription due to a favorable folding of the genome in the three-dimensional space [3]. Enhancer sequences are also capable of bidirectional transcription of RNAs (eRNAs) at a large scale. This means that gene promoters and enhancers share a similar promoter architecture, each bound by RNA pol II when active [4]. Recent studies have shown that promoters and enhancers share several properties and functions related to their chromatin and sequence architectures [5]. The distinction between these regulatory elements is further reduced by the fact that promoters can enhance transcription at a distant site [6] while enhancers can drive local transcription initiation [7]. This provides a motivation to consider these elements together when trying to understand transcription regulation and build models for their identification.

Thanks to the development of advanced experimental techniques, significant progress has been made in identifying gene regulatory sequences [8–10]. However, a detailed experimental exploration of transcripts is still an expensive and difficult procedure. Therefore, in addition to experimental efforts, accurate computational identification of putative regulatory regions, for both individual genes and entire genomes, remains an important challenge of genomics studies. Computational prediction is important for guiding experimental biologists, providing putative regions which can be validated using reporter gene assays.

Accurate computational identification of regulatory elements remains a difficult task due to the high diversity of DNA sequence features and tissue specificity of the transcriptional regulation. Some promoters and enhancers are only active in a small fraction of cells or under specific conditions. This makes it difficult to get a complete list of regulatory regions in the human genome, such as alternative promoters for the known genes. There are numerous computational tools developed in an attempt to predict promoters and enhancers [11–14]. However, many of them focus on discrimination between fixed sets of promoter/enhancer sequences and random genomic sequences. The reported performance deduced from such

small and balanced test sets does not hold when evaluated at the level of the whole genome, which is a much more difficult task due to the huge number of tested locations [15].

It has been shown that accurate regulatory element prediction can be achieved on the genome wide scale by using epigenetic data, such as DNA methylation and histone modification profiling [16]. For example, H3K4me1, H3K4me3, and H3K27ac marks are associated with promoter and enhancer activities [17]. In recent years, a number of methods have been developed that utilize local epigenetic marks for regulatory element prediction, based on machine learning algorithms such as random forests [16,18], support vector machines [19,20], and deep learning [21,22]. However, there is still a need for methods that can provide accurate predictions based on DNA sequences, in particular for annotating the genomes of species where epigenetic data are not widely available.

In this study, we developed ReFeaFi (Regulatory Feature Finder), a general genome-wide promoter and enhancer predictor, using the DNA sequence alone. Using Cap Analysis Gene Expression (CAGE) data [23] as the ground truth for promoters and enhancers, we used a dynamic training set updating scheme to train the deep learning model, which allows us to have high recall while keeping the number of false positives low, improving the discrimination and generalization power of the model. ReFeaFi achieved comparable performance when the model trained on the human genome was applied to other vertebrate species, showing the generality of our model. We found that unannotated regulatory regions predicted by our method are enriched for genetic variants associated with disease [24,25], suggesting that they might be real regulatory elements. High scoring unannotated predictions were validated using reporter assays and all the candidates showed strong luciferase signal. By analyzing synthetic promoters, we found that the predicted score strongly correlates with measured expression strength. We used the trained deep learning model to study the architecture of regulatory elements and to find out how conserved elements affect transcription strength. Finally, we have developed a novel model analysis technique that reveals related positions around the TSS.

Results

Genome-wide identification of regulatory elements

The overview of our method is shown in Fig 1. The method is two-tiered and consists of a scan model and a prediction model, which are trained iteratively to reduce false positives (FPs). Briefly, the scan model picks candidate regions for the prediction model, which then makes a decision if these regions contain one or more TSS. If an unannotated region receives a score above the threshold, it is added to the negative set. The whole process is repeated several times to generate a difficult negative set which forces the model to learn more complex features for the identification of regulatory regions. We trained our model on the human genome, with chromosome 1 used as the test set and the rest of the genome for training and validation. The positive set of regulatory region annotations used was constructed by merging the human robust promoter set and human permissive enhancer set downloaded from the FANTOM5 website [26]. We initially compared our method against several previously published promoter predictors (Basset [27], PromPredict [28], and EP3 [11]), which are capable of genome-wide TSS predictions. Using human chromosome 1 for evaluation (see Methods), our developed model achieved good performance, significantly outperforming other methods (Table 1 and Fig 2A). In particular, ReFeaFi generated substantially less FPs than a deep learning-based method for regulatory elements prediction (Basset), for all recall values, see Fig 2A. For the general recall rate of 0.50, promoters of the protein coding genes were predicted with recall 0.77 and 0.41 recall was obtained for promoters of the long non-coding RNAs. Since enhancers

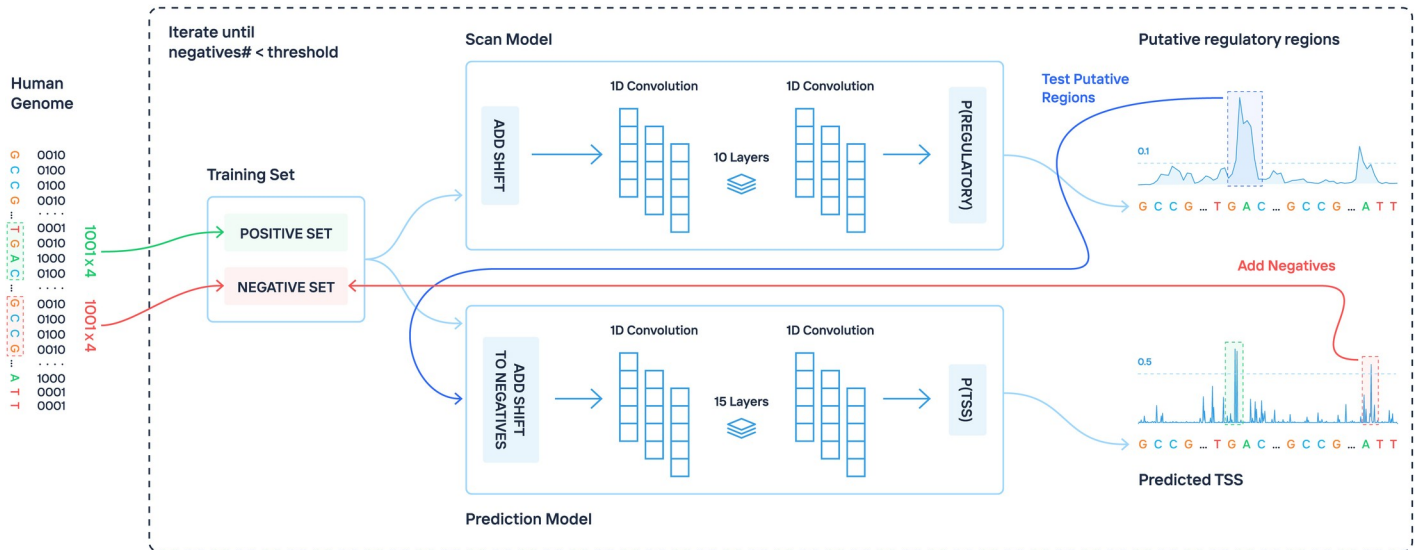


Fig 1. Workflow of the proposed method for genome-wide regulatory elements prediction. The scan model uses a sliding window approach to pick putative regulatory regions. The prediction model finds TSS positions inside these regions by testing each position. The false positive predictions made by the second model are added to the negative set for the next round of training. The whole process is repeated iteratively to generate a difficult negative set which forces the model to learn how to distinguish the difficult negatives from the real regulatory sequences.

<https://doi.org/10.1371/journal.pcbi.1009376.g001>

have weaker transcriptional output, only 9 percent of the permissive enhancers were detected using the strict decision threshold (S1 Table).

We subsequently applied the model trained on human CAGE data to genomes of other species, without re-training (Fig 2B). The obtained result for the mouse genome was similar to the human performance, achieving comparable recall and FP rate. For other species, the recall was high, but the FP rate was significantly higher than for human and mouse. This is expected and due to the fact that the transcriptomes of these species have been profiled at considerably lower depth, and there are thus much fewer CAGE peaks known for these species (S2 Table, most of them are related to housekeeping genes which are easier to detect due to higher GC content [29]). This gives rise to a high recall, where additional FPs are most likely related to tissue-specific promoters, which have not been detected in the currently available CAGE datasets for these organisms.

Validation of unannotated predicted regions by reporter assay

We next set out to investigate the extent to which predicted regulatory regions marked as false positives in the initial evaluation might represent regions with true regulatory potential, but not yet discovered by the current experimental annotations. To this end, we performed reporter assays of 17 high scoring regions and three regions with a low score but still predicted

Table 1. Comparison of the performance of different TSS identification methods. The decision thresholds for these methods were adjusted to achieve the same recall of 0.50.

Method	Recall	Precision	F1 score	FP per correct	FP per 1 Mb
ReFeaFi	0.50	0.52	0.51	0.93	46.44
EP3	0.50	0.20	0.29	4.01	198.26
Basset	0.50	0.07	0.12	13.48	664.92
PromPredict	0.50	0.07	0.12	13.44	666.26

<https://doi.org/10.1371/journal.pcbi.1009376.t001>

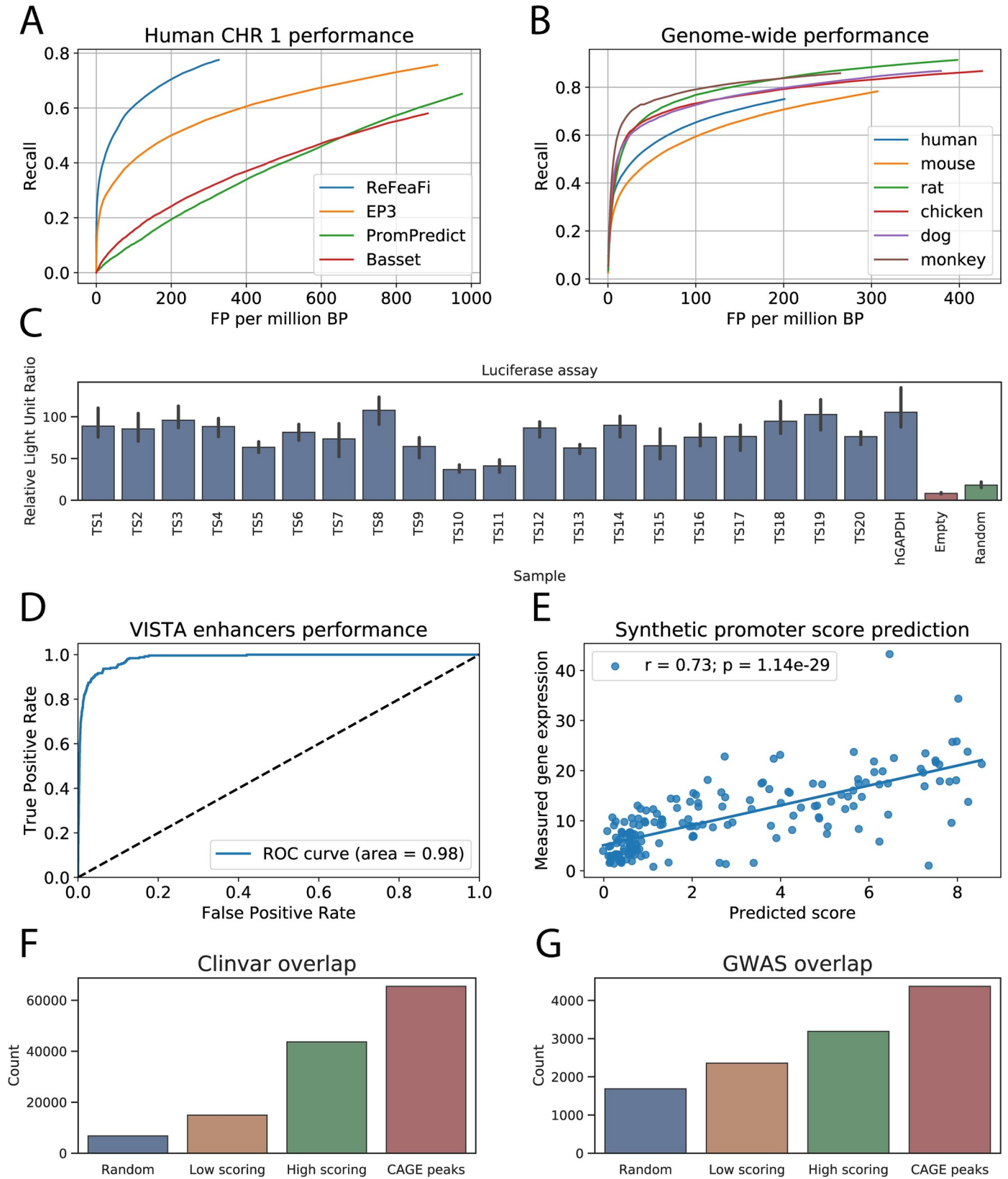


Fig 2. The performance evaluation of the proposed method. (A): Performance of TSS predictors on human chromosome 1. Our method, ReFeaFi, significantly outperforms other predictors. (B): Performance of the proposed method applied on different vertebrate species, without re-training the model. Due to the small number of known CAGE peaks, there is a significant difference between human/mouse and other organisms. (C): Results of the reporter assay on the predicted unannotated regulatory regions. Each chosen sample is at least 5kb away from any known CAGE peaks or GENCODE annotations. Empty vector and random sequences were used as negative controls, while GAPDH promoter is used as positive control. All the candidates passed the validation, with several candidates showing intensity similar to the positive control. (D): ReFeaFi achieves outstanding performance on discriminating VISTA enhancers and 100 times as many random genomic regions. (E): Correlation between score predicted by our method and measured mean expression of the synthetic promoters. (F): The number of Clinvar genetic variants located in the vicinity of low and high scoring predicted regions. (G): GWAS genetic variants overlap with the predicted regions.

<https://doi.org/10.1371/journal.pcbi.1009376.g002>

as regulatory by our model. The tested sequences included regions with and without classic promoter elements (INR and TATA box, [S3 Table](#)) and were chosen so that they were at least 5kb away from any known CAGE peaks or GENCODE annotations. All the chosen candidates showed significantly higher luciferase signals than negative controls (empty vector and random sequence, [Fig 2C](#)), with several sequences showing intensity similar to the GAPDH promoter used as the positive control.

Highly accurate prediction of VISTA enhancers

The VISTA Enhancer Browser contains experimentally validated human and mouse enhancers with their activity measured in transgenic mice [30]. VISTA enhancers were previously used by Yang *et al.* to validate their enhancer predictor BiRen [13], where like our method, the model is trained using DNA sequences alone, unlike methods that use epigenetic information to predict enhancers [19,20]. In this validation approach, a test set is constructed by using VISTA enhancers as a positive set, with a negative set that contains ten times as many non-enhancer sequences. When evaluated, BiRen achieved AUC of 0.945, improving over DEEP [14] and SVM [31] enhancer predictors which achieved AUC of 0.883 and 0.621, respectively. We adopted the same experimental setup, constructing a positive set with all the VISTA human and mouse enhancers from chromosome 1 and 10 times more random genomic regions for the negative set. Our model achieved almost perfect accuracy, with $AUC > 0.99$. The same result was obtained when repeating the experiment with 100 times more negative sequences (AUC 0.98, [Fig 2D](#)).

Predicted score correlates with measured expression

Weingarten-Gabbay *et al.* devised a high-throughput assay to quantify the activity of fully designed sequences that were integrated and expressed from a fixed location within the human genome [32], using the method to investigate binding regions of core promoters. We applied our model to the designed sequences and computed the correlation between the measured expression and our predicted scores ([Fig 2E](#)). Correlation between the score and measured expression was 0.73, further showing the generality of our model, since it had not seen any of these synthetic sequences during training and suggesting that our model may be useful for designing new promoters with desired strength by screening potential candidate sequences.

Non-annotated predicted regions are enriched for disease-associated variants

Nucleotide variation in enhancers and promoter regions has been shown to be associated with human diseases. Indeed, most of the disease-associated genome-wide association studies (GWAS) hits fall in the non-coding regions of the human genome [33]. GWAS helps to understand disease mechanisms and provides the starting point for the development of medical diagnosis, prognosis, and treatments. We tested if our genome-wide predictions are enriched

for genetic variants from GWAS [33] and Clinvar [25] databases. First, we computed the overlap between the known CAGE peaks and GWAS and Clinvar sets with a margin of 200 bp relative to the TSSs. Next, we compared the overlap with the variant sets for the same number of low scoring and high scoring predicted regions, excluding predictions overlapping known TSS or enhancers. The Clinvar set overlapped the high scoring regions three times more compared to the low scoring regions (Fig 2F), suggesting that they represent novel regulatory elements. The GWAS disease-associated SNPs were overrepresented in the high scoring regions as well, overlapping them two times more than expected by chance, Fig 2G. Together with the experimental results, this shows that the unannotated predicted regulatory regions we detected might have regulatory potential to a significant degree, and that they may be useful for prioritization of further studies of disease-associated variants.

Non-annotated predicted regions are enriched for meaningful epigenomic marks

To further validate the regulatory potential of the predicted unannotated regions, we decided to check if they are enriched for epigenomic marks associated with regulatory activity. We have downloaded 129 ROADMAP [34] epigenome datasets, keeping marks associated with promoter and enhancer activity. They were averaged per epigenomic mark and then normalized. Next, we computed the overlap with our low scoring and high scoring predicted regions. As shown in Table 2, the high scoring false positive regions are more enriched for the marks associated with regulatory activity compared to the low scoring regions.

Model analysis reveals known and novel regulatory features

Given the limitations of the biochemical assays used to identify and characterize core promoter elements, it has been difficult to assign a clear function to each of these elements [2]. By feeding in modified regulatory sequences to our deep learning model, we can study how each core promoter element fine-tunes the gene expression. To measure the effects of each known motif on the promoter score, we searched for them using previously constructed PWMs for core promoter elements [35] and replaced them with random nucleotides. New promoter scores were computed for the modified sequences and the change in the predicted score was recorded for each motif. The most important motifs were TATA-box and INR with 35% and 9% effect on the score, respectively (Fig 3). Mutation of DPE and CCAT motifs had no strong effect on the predicted score, suggesting that they do not play a significant role in regulating human promoter activity. The same observations were made in a recent human promoter study [32]. In fact, the order in terms of gene expression impact is the same for the motifs we analyzed.

We next used the model to identify important locations in the input sequences that contribute most to the predicted score. Using a sliding window moving from the beginning of a regulatory sequence, we built a performance profile that reflects an effect of a random sequence, inserted in each sequence position in place of an original sequence, on the predicted score. The results for promoter and enhancer sequences (shown in Fig 4A and 4B respectively) revealed that for promoters, the TSS and TATA regions are extremely important, while for

Table 2. Enrichment for epigenomic marks of low scoring and high scoring predictions.

Epigenomic mark	Low scoring	High scoring	Increase
H3K4me1	0.074	0.123	1.662
H3K4me3	0.010	0.046	4.600
H3K27ac	0.027	0.052	1.926

<https://doi.org/10.1371/journal.pcbi.1009376.t002>

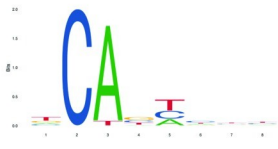
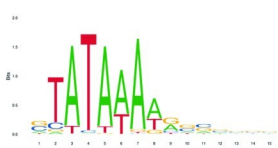


Name	Motif	Start	Frequency	Score Effect
Inr		-2	28.0%	-9%
TATA		[-30 : -31]	8.3%	-35%
CATT		[-212 : -57]	9.2%	-2%
DPE		[+23 : + 32]	28.0%	-3%

Fig 3. Effect of substituting promoter conserved motifs with random nucleotides on the score predicted by the deep learning model.

<https://doi.org/10.1371/journal.pcbi.1009376.g003>

enhancers, TSS regions contribute less to the total score compared to other regions. For example, mutation of the eRNA TSS nucleotides only reduces the score by two percent on average, while for promoters this value is 35 percent.

Despite the main effect on score coming from the core promoter region, in some cases replacing part of the sequences outside core promoter had a significant effect on the score, reducing it up to 50%. Motifs within core promoters have been extensively studied, but not much is known about distant motifs important for promoter activity. We tested if transcription factor binding motifs from the JASPAR database [36] have a significant effect on the predicted score, and observed that even though there was often a perfect match of a known binding motif inside the promoter region, replacing it with random nucleotides did not change the output of the model in many cases. However, some of the motifs affected the score significantly, even when they were located far away from the TSS. Interestingly, the set of most influential motifs for promoters and enhancers were close to each other, suggesting that they are regulated by a similar set of TFs (S4 Table).

To assess the contributions of different nucleotides in different positions of the regulatory sequences, we employed a modification of a feature mutation map for our test set. Fig 4C shows the mutation maps for the core promoter in the promoter sequences, while the mutation maps for the enhancer sequences is shown in Fig 4D. To build these maps, we studied

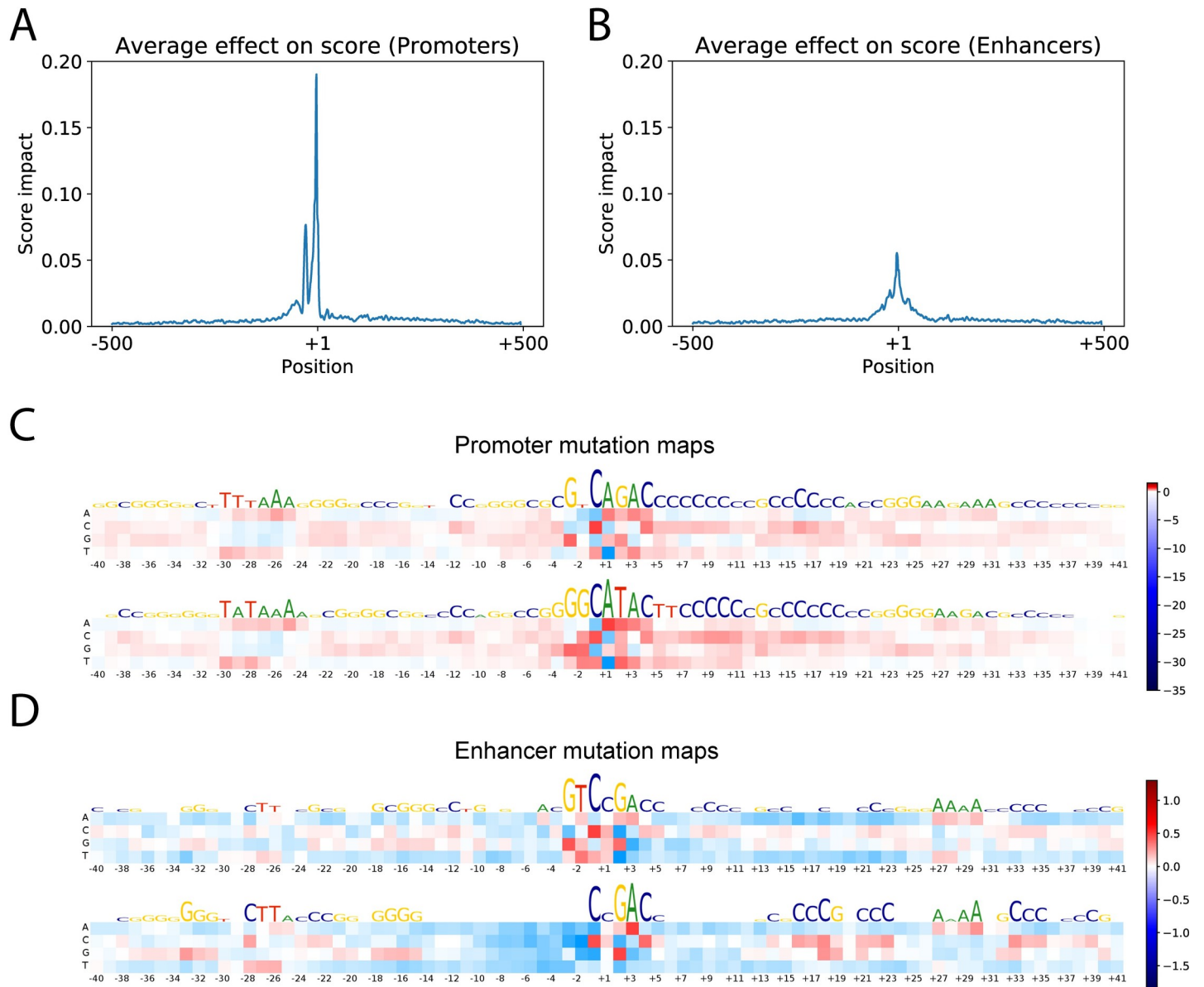


Fig 4. Analysis of the trained model. (A): Effect of 7 bp sequence window substitution by a random sequence on the predicted score of the promoters. (B): The same substitution procedure applied to the enhancer sequences. (C): Mutation maps for the promoter sequences. K-means with $k = 10$ was applied to the mutation matrices before averaging and the two biggest clusters are shown. (D): Mutation maps for the two biggest clusters of enhancer mutation matrices.

<https://doi.org/10.1371/journal.pcbi.1009376.g004>

how nucleotide substitutions will change the output score computed by our model. At each position of the tested sequences, we replaced a nucleotide with a different one in all the sequences and computed their average score change. The rows represent nucleotides that are used for replacement and the columns show different positions inside the regulatory regions. If the new score on average increases, it is represented by a red-colored square. Decreasing the score is shown by using a blue-colored square. The intensity of color is proportional to the effect of substitution on the score. Since core promoters are very diverse, if we draw a mutation map for all the sequences, we cannot capture all the information from our model. To alleviate this, we clustered the mutation matrices based on their similarity (Frobenius norm) into several clusters using the k-means algorithm with $k = 10$.

Analysis of the mutation maps showed that promoters have higher GC content than enhancers and conserved C and A nucleotides at positions -1 and +1 respectively, which corresponds to the Initiator motif (Fig 3). Enhancers do not have a conserved TSS nucleotide. However, clustering reveals that many of the enhancers have a CCGACC motif around the TSS. For promoters, the revealed motifs in the TSS region are GGCATAC and GTCAGAC.

Mutation/saliency maps and convolutional filter analysis are often used to understand deep learning models, but they cannot detect long-range interactions in the input sequence. We attempted to overcome this limitation by developing a novel technique to measure dependency between each pair of nucleotides—pair dependency map. Fig 5A shows the dependency between positions inside the core promoter region for the promoter sequences. Here every element of the matrix $V(i,j)$ is the difference between two values: the change in the predicted score when removing the pair (i,j) compared to the sum of changes caused by removing i and j separately. Every input sequence generates a symmetric matrix, and we draw the averaged matrix. The red color is used for positive elements of V , and the blue color is for negative elements. Intensity shows the absolute value of the element. Negative matrix elements occur typically when the pair i and j is involved in a motif where these two positions are correlated.

The pair maps for the enhancer sequences are shown in Fig 5B. As shown in previous experiments, interactions between TSS and the TATA-box are very important for the promoter sequences, while interactions in enhancer sequences are more uniformly distributed. This is especially easy to see after performing the k-means clustering of the matrices before averaging.

The proposed pair map analysis reveals related positions. Unlike the single position mutation map analysis, if a position changes a score significantly but independent of other positions, pairs involving it will be shown with a light color. The pair dependency map strongly suggests that our model can capture long range dependencies between different elements in the regulatory sequences, and that the trained model does not simply detect conserved core promoter elements but also complex interactions between them. Based on the pair map idea, we created a tool that can test whether two given positions (or two regions) are independent according to our model, which may be useful for creating hypotheses such as if two transcription factors regulate a subset of promoters together. Using the described approach, we analyzed interaction between JUND and BATF TFs, which was reported in the previous studies [37,38]. Searching for their motifs inside the regulatory regions using FIMO [39] with default parameters, the interaction between JUND and BATF was compared to the interaction of JUND with random locations of the same motif length. As expected, a strong relationship between the two binding motifs was detected ($p = 3.04e-11$, t-test).

Discussion

We have demonstrated that by using a dynamic training set, it is possible to tackle the problem of genome-wide regulatory elements prediction. This is a very difficult problem because many non-regulatory regions are very similar to the true ones in terms of their nucleotide sequences. Unlike common machine learning problems, prediction at the genome-wide scale is extremely unbalanced, which makes it difficult to achieve high sensitivity. Despite all these challenges, the proposed model achieves a good performance on the human genome and can be directly applied to genomes of other animals as well without significant loss of performance. It can also be used to pinpoint promoter and enhancer candidates in regions of interest in the human genome. Many regulatory regions might still be unannotated which is supported by the results of the reporter assay experiment and epigenomic mark enrichment analysis we performed. We showed that the false positive regions are enriched with GWAS signals. For the SNPs which

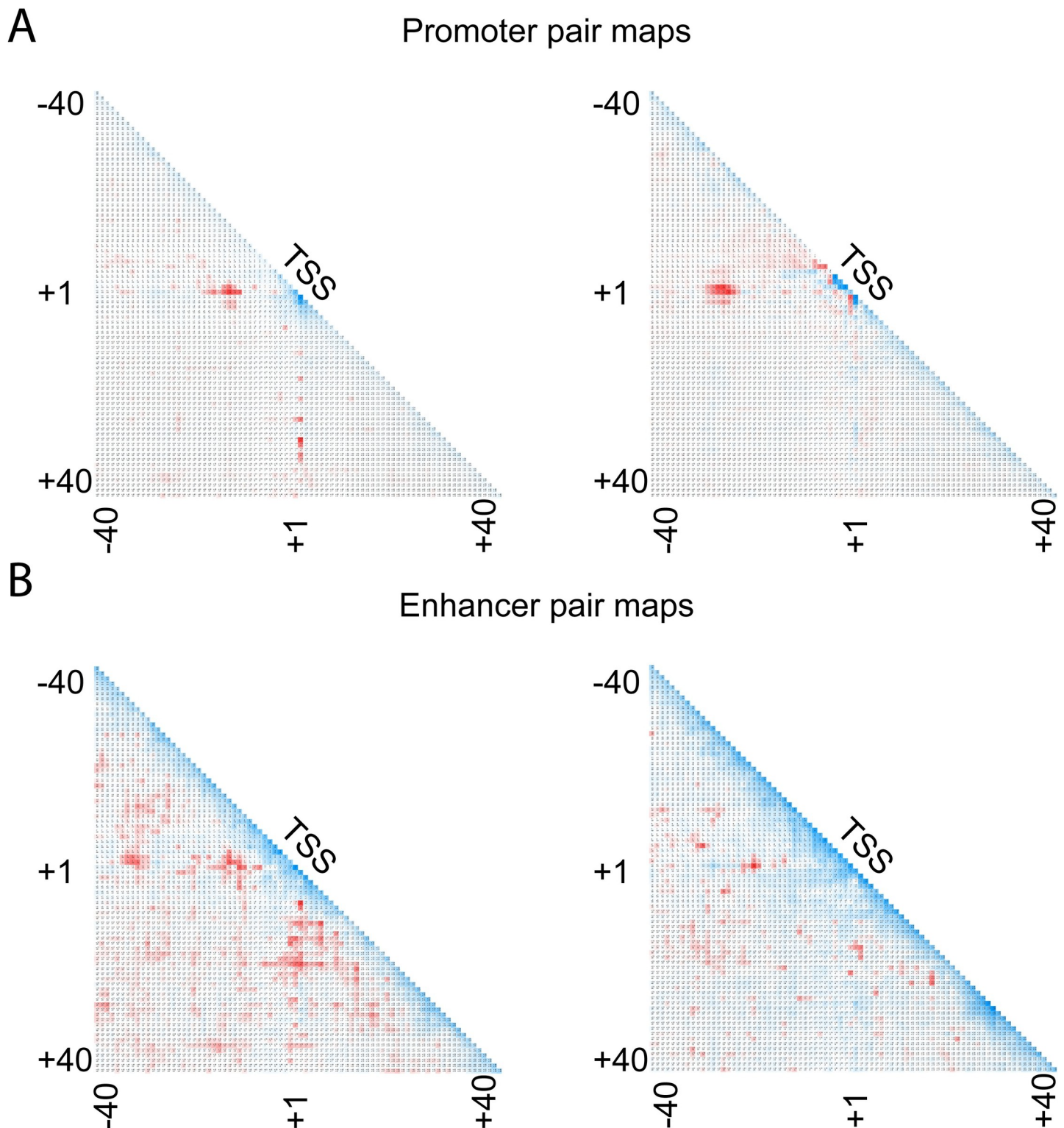


Fig 5. Pair dependency map reveals long-range interactions captured by our model. The red color means interaction, the blue color shows correlation, and the white color represents that the positions are independent. The results are clustered into ten groups using k-means algorithm and the two biggest groups are shown. (A): Promoter pair maps show relationships between conserved promoter elements, which have much stronger interactions than other pairs. They include TSS interactions with TATA and BRE elements in the upstream region, DCE and DPE elements in the downstream region. (B): Dependencies captured by our model for predicting enhancers are spread out in the [-40: +40] region. Unlike promoter sequences, there are no strong short-range interactions in the TSS region.

<https://doi.org/10.1371/journal.pcbi.1009376.g005>

fall into previously unknown regulatory regions, our predictions provide a hypothesis which can be tested in a further analysis.

In this study, we have shown that the trained model can identify motifs and positions important for the activity of regulatory elements, which allows for further exploration of the promoter architecture. Furthermore, we identified long-range interactions in enhancer and core promoter regions using a pair dependency map. The next step would be to validate these findings using reporter assays by constructing synthetic promoters and enhancers. This would help to better understand transcriptional regulation and also design new regulatory sequences with desired strength.

To further improve the results, it is necessary to consider the cell type when predicting promoters and enhancers. It has been shown that promoter prediction can also be improved using extra information as input, for example chromatin data [40] or possibly applying non-parametric methods as described and tested on promoter regions of a model dicot plant *Arabidopsis thaliana* [41]. However, the approach described in this paper has the advantage that it is very general. It will be straightforward to apply it to many different organisms where additional information like chromatin profiles might not be available.

Methods

Data sets used

The data we used for training our models consist of 210250 promoters and 65423 enhancers downloaded from FANTOM5 website. Sequences of length 1001 bp were extracted from the human genome (GRCh37) centered around the TSSs. All the regulatory regions from chromosome 1 were used for testing, 90% of the remaining peaks were used for training, and the rest for validation to choose deep learning parameters and perform early stopping during training. GENCODE version 34 was used when picking candidates for reporter assay validation to avoid picking regions near any known gene. GWAS and Clinvar sets were downloaded from their respective websites on 2020-07-06. The motif analysis used PWMs from the JASPAR 2020 version.

Negative set generation

One of the reasons the reported performance in previous studies does not extrapolate to the whole human genome is because of an inadequately chosen set of negative sequences [15]. Often the negative set consists of random genomic sequences which are very different from the regulatory regions or sequences from specific positions (e.g., fixed distance away from TSS) which introduces bias. To tackle this problem, we used an iterative approach that updates the set of non-promoter sequences used in the training set based on the false positive errors made in the previous iteration [42]. By including difficult non-promoter sequences in the training set, the predictor is forced to learn promoter patterns to rule out such sequences. This scheme allowed us to achieve high sensitivity while keeping the number of false positive predictions low. [S1 Fig](#) illustrates that without using these difficult negatives, the model mostly uses the GC content to make a prediction.

When a difficult negative set is obtained, the neural network struggles to discriminate provided positive and negative sequences. This can result in over-fitting, when the model simply memorizes some difficult negatives despite the regularization methods that are used. To avoid this issue, we have developed a new regularization technique. During each epoch, we added a small shift to the negative sequences, moving them upstream or downstream from the initial positions by a random distance, see [S2 Fig](#). This virtually makes the negative set very large and impossible to memorize. The neural network can only learn very general patterns to rule out

the negatives. We have found that this specific technique is superior to any other changes that one can do to the negative set. If parts of DNA sequence are manually changed, it is very easy for the neural network to detect such alterations, e.g., replacement, insertions, mirroring, generation of a new DNA sequence with specific nucleotide frequency. When such alterations are used, the network is trained to distinguish the original DNA sequence from an altered one, instead of discriminating promoters from non-promoters.

Genome-wide regulatory elements prediction

Since it is difficult to test every possible TSS location in the genome, we used a two-fold prediction procedure. One model is used to scan each chromosome using a sliding window approach. Because of the random shift added to both positive and negative sequences during the training, the input for the scan model does not need to be centered perfectly around the TSS. This allowed us to use a relatively large step for the sliding window, 50 bp. When the model outputs values above the threshold (0.5), the 100 bp region centered around the current position is scanned using position specific second model with the step equal to one. Predictions of the second model with scores higher than 0.5 are sorted and filtered based on the distance between them. The remaining results are output by the method. The models are trained to detect regulatory elements on both strands without explicitly performing reverse complement. The strand is decided by a special model, which decides the direction of transcription: either positive (+) or negative (-) for promoters and both (.) for enhancers. This design decision was made to make genome-wide TSS identification faster.

Evaluation criteria

During the evaluation, all the predictions more than 500 bp away from a known CAGE peak were considered as FPs. The margin for error is rather big to deal with the alternative promoters for the same gene, since the TSSs can be located a large distance from each other [43]. One or more TSS predictions inside the regulatory region count as a TP, however if there is not even one prediction in this region, we count it as a FN. We measured our performance using F1 score:

$$Precision = \frac{TP}{TP + FP},$$

$$Recall = \frac{TP}{TP + FN},$$

$$F1 = 2 * \frac{Recall * Precision}{Recall + Precision}.$$

We also calculated two additional metrics suitable for genome-wide setting, FP per correct prediction and FP rate per 1 Mb:

$$FP \text{ per correct} = \frac{FP}{TP},$$

$$FP \text{ per } 1 \text{ Mb} = 1000000 * \frac{FP}{chrLsize}.$$

Deep learning architecture

The deep learning architecture is based on deep residual networks [44]. The network consists of 5 residual blocks followed by a Softmax layer, S3 Fig. To avoid the problem of vanishing gradient, we employed batch normalization [45]. The activation function used throughout the model is Leaky ReLU [46]. Weight decay and dropout are used to improve the generalization capability of the model. Weight decay effectively limits the number of free parameters in the model to avoid over-fitting. Introducing weight decay makes it possible to regularize the cost function by penalizing large weights. The main idea of dropout is to randomly set some nodes of the neural network to zero during training to prevent co-dependency among them. During the training, the dropout for the feature vector with keep probability of 0.5 is used. The Adam optimization algorithm is used to train the weights [47], which is an improved version of stochastic gradient descent. TensorFlow [48] is used as the framework to construct the deep neural network. The training was performed on a workstation with four NVIDIA Quadro RTX 6000 GPUs and took about 5 days to complete. The method requires eight hours to scan and predict regulatory regions on chromosome 1 with the standard parameters.

Reporter assays

We used reporter assays to validate the candidate sequences. Double strand DNA cassette was constructed by annealing the DNA oligo that has SpeI site/I-CeuI site and the DNA oligo that has BamHII site/I-SceI site. Specific DNA primers were designed for amplification of the test candidates. Using these primers and Human Genomic DNA as template, PCR was performed with KOD-Plus-Neo. After QC with electrophoresis, PCR products (insert DNAs) were digested by I-CeuI/I-SceI. We selected and prepared human GAPDH promoter region as positive control and random sequences (backbone of pMCS-Cypridina Luc vector (201 bases long)) as negative control. The prepared vector and insert DNA were ligated using Rapid DNA ligation Kit, and One Shot TOP10 Chemically Competent *E. coli* was transformed by the ligated vector. These transformed *E. coli* cells were cultured in large scale and assay vectors were extracted using QIAGEN Plasmid Mini Kit. For quality checking, PCR was performed with assay vector and primers for checking. After this, electrophoresis was performed to confirm the presence of the desired insert DNA fragment in the assay vector.

Until the day before doing assay, HEK293T cells were pre-cultured. At the day before doing assay, pre-cultured HEK293T cells were spread on 96-well plate and incubated at 37°C CO₂ 5% for 16-24h. After incubation, assay vector was transfected into HEK293T cells using Turbofect Transfection Reagent according to the kit protocol. After incubation, Cells were lysed using Cell lysis buffer included in Pierce Cypridina-Firefly Luciferase Dual Assay Kit. Prepared reagent mixture containing D-Luciferin was added to the cell lysate and measurement of luminescence from luciferin-luciferase reaction was performed by luminometer. Detailed description of the experimental setup is provided in the S1 Appendix.

Alternative methods

We have obtained Basset from the GitHub repository and followed the provided guide for peak prediction (<https://github.com/calico/basenji/tree/master/manuscripts/basset>). The data were replaced with two of our bed files for promoters and enhancers. After generating the data, we modified the target neurons number, setting it to 2. The training took 18 epochs and was stopped automatically after the validation AUROC did not improve (0.75301). The trained model was then applied to chromosome 1 with the default test stride of 192 using provided script: `basenji_predict_bed.py`. EP3 and PromPredict do not require extra training, which is why we applied them directly on both chromosome 1 and its reverse complement.

Synthetic promoters analysis

The synthetic sequences were inserted into the same genomic background as in the original publication before using them as an input for our model. To undo the effect of Softmax which makes all the output values very close to 0 or 1, our predicted score was adjusted as follows: $\log(1 + \text{Score}/(1 - \text{Score}))$.

Supporting information

S1 Fig. Difference between random negatives model and hard negatives model in terms of their mutation map. a. Random negatives model. b. Difficult negatives model.

(TIF)

S2 Fig. Procedure of shifting negative sequences, which virtually increases the negative set size and prevents overfitting.

(TIF)

S3 Fig. Deep learning architecture that was used in building our model. The model employs CNN with residual connections which is followed by a softmax layer.

(TIF)

S1 Table. Performance of ReFeaFi for different decision thresholds.

(XLSX)

S2 Table. Number of CAGE peaks for different species.

(XLSX)

S3 Table. Predicted regulatory sequences that were chosen for reporter assay validation.

(XLSX)

S4 Table. Importance scores assigned to each TF from the JASPAR database.

(XLSX)

S1 Appendix. Reporter assay details.

(DOCX)

Acknowledgments

We would like to thank Andrew Tae-Jun Kwon and Bogumil Kaczkowski for insightful comments on the manuscript.

Author Contributions

Conceptualization: Erik Arner.

Formal analysis: Ramzan Umarov.

Investigation: Ramzan Umarov.

Methodology: Ramzan Umarov.

Project administration: Erik Arner.

Software: Ramzan Umarov.

Supervision: Xin Gao, Erik Arner.

Validation: Ramzan Umarov, Yu Li, Takahiro Arakawa, Satoshi Takizawa.

Visualization: Ramzan Umarov.

Writing – original draft: Ramzan Umarov.

Writing – review & editing: Xin Gao, Erik Arner.

References

1. Lenhard B, Sandelin A, Carninci P. Metazoan promoters: emerging characteristics and insights into transcriptional regulation. *Nat Rev Genet.* 2012 Mar 6; 13(4):233–45. <https://doi.org/10.1038/nrg3163> PMID: 22392219
2. Roy AL, Singer DS. Core promoters in transcription: old problem, new insights. *Trends Biochem Sci.* 2015 Mar; 40(3):165–71. <https://doi.org/10.1016/j.tibs.2015.01.007> PMID: 25680757
3. Schoenfelder S, Fraser P. Long-range enhancer-promoter contacts in gene expression control. *Nat Rev Genet.* 2019 Aug; 20(8):437–55. <https://doi.org/10.1038/s41576-019-0128-0> PMID: 31086298
4. Andersson R, Refsing Andersen P, Valen E, Core LJ, Bornholdt J, Boyd M, et al. Nuclear stability and transcriptional directionality separate functionally distinct RNA species. *Nat Commun.* 2014 Nov 12; 5:5336. <https://doi.org/10.1038/ncomms6336> PMID: 25387874
5. Core LJ, Martins AL, Danko CG, Waters CT, Siepel A, Lis JT. Analysis of nascent RNA identifies a unified architecture of initiation regions at mammalian promoters and enhancers. *Nat Genet.* 2014 Dec; 46(12):1311–20. <https://doi.org/10.1038/ng.3142> PMID: 25383968
6. Engreitz JM, Haines JE, Perez EM, Munson G, Chen J, Kane M, et al. Local regulation of gene expression by lncRNA promoters, transcription and splicing. *Nature.* 2016 Nov 17; 539(7629):452–5. <https://doi.org/10.1038/nature20149> PMID: 27783602
7. Kim T-K, Hemberg M, Gray JM, Costa AM, Bear DM, Wu J, et al. Widespread transcription at neuronal activity-regulated enhancers. *Nature.* 2010 May 13; 465(7295):182–7. <https://doi.org/10.1038/nature09033> PMID: 20393465
8. Mundade R, Ozer HG, Wei H, Prabhu L, Lu T. Role of ChIP-seq in the discovery of transcription factor binding sites, differential gene regulation mechanism, epigenetic marks and beyond. *Cell Cycle Georget Tex.* 2014; 13(18):2847–52. <https://doi.org/10.4161/15384101.2014.949201> PMID: 25486472
9. Suryamohan K, Halfon MS. Identifying transcriptional cis-regulatory modules in animal genomes. *Wiley Interdiscip Rev Dev Biol.* 2015 Apr; 4(2):59–84. <https://doi.org/10.1002/wdev.168> PMID: 25704908
10. Levati E, Sartini S, Ottonello S, Montanini B. Dry and wet approaches for genome-wide functional annotation of conventional and unconventional transcriptional activators. *Comput Struct Biotechnol J.* 2016; 14:262–70. <https://doi.org/10.1016/j.csbj.2016.06.004> PMID: 27453771
11. Abeel T, Saeys Y, Bonnet E, Rouzé P, Van de Peer Y. Generic eukaryotic core promoter prediction using structural features of DNA. *Genome Res.* 2008 Feb; 18(2):310–23. <https://doi.org/10.1101/gr.6991408> PMID: 18096745
12. Kalkatawi M, Magana-Mora A, Jankovic B, Bajic VB. DeepGSR: an optimized deep-learning structure for the recognition of genomic signals and regions. *Bioinforma Oxf Engl.* 2019 Apr 1; 35(7):1125–32. <https://doi.org/10.1093/bioinformatics/bty752> PMID: 30184052
13. Yang B, Liu F, Ren C, Ouyang Z, Xie Z, Bo X, et al. BiRen: predicting enhancers with a deep-learning-based model using the DNA sequence alone. *Bioinforma Oxf Engl.* 2017 Jul 1; 33(13):1930–6. <https://doi.org/10.1093/bioinformatics/btx105> PMID: 28334114
14. Klefogiannis D, Kalnis P, Bajic VB. DEEP: a general computational framework for predicting enhancers. *Nucleic Acids Res.* 2015 Jan; 43(1):e6. <https://doi.org/10.1093/nar/gku1058> PMID: 25378307
15. Khodabandelou G, Routhier E, Mozziconacci J. Genome annotation across species using deep convolutional neural networks. *PeerJ Comput Sci.* 2020 Jun 15; 6:e278. <https://doi.org/10.7717/peerj-cs.278> PMID: 33816929
16. Ramisch A, Heinrich V, Glaser LV, Fuchs A, Yang X, Benner P, et al. CRUP: a comprehensive framework to predict condition-specific regulatory units. *Genome Biol.* 2019 Nov 8; 20(1):227. <https://doi.org/10.1186/s13059-019-1860-7> PMID: 31699133
17. Karlić R, Chung H-R, Lasserre J, Vlahovicek K, Vingron M. Histone modification levels are predictive for gene expression. *Proc Natl Acad Sci U S A.* 2010 Feb 16; 107(7):2926–31. <https://doi.org/10.1073/pnas.0909344107> PMID: 20133639
18. He Y, Gorkin DU, Dickel DE, Nery JR, Castanon RG, Lee AY, et al. Improved regulatory element prediction based on tissue-specific local epigenomic signatures. *Proc Natl Acad Sci U S A.* 2017 Feb 28; 114(9):E1633–40. <https://doi.org/10.1073/pnas.1618353114> PMID: 28193886

19. Fernández M, Miranda-Saavedra D. Genome-wide enhancer prediction from epigenetic signatures using genetic algorithm-optimized support vector machines. *Nucleic Acids Res.* 2012 May; 40(10):e77. <https://doi.org/10.1093/nar/gks149> PMID: 22328731
20. Sethi A, Gu M, Gumusgoz E, Chan L, Yan K-K, Rozowsky J, et al. Supervised enhancer prediction with epigenetic pattern recognition and targeted validation. *Nat Methods.* 2020 Aug; 17(8):807–14. <https://doi.org/10.1038/s41592-020-0907-8> PMID: 32737473
21. Williams J, Xu B, Putnam D, Thrasher A, Li C, Yang J, et al. MethylationToActivity: a deep-learning framework that reveals promoter activity landscapes from DNA methylomes in individual tumors. *Genome Biol.* 2021 Jan 19; 22(1):24. <https://doi.org/10.1186/s13059-020-02220-y> PMID: 33461601
22. Kim SG, Harwani M, Grama A, Chaterji S. EP-DNN: A Deep Neural Network-Based Global Enhancer Prediction Algorithm. *Sci Rep.* 2016 Dec 8; 6:38433. <https://doi.org/10.1038/srep38433> PMID: 27929098
23. Takahashi H, Lassmann T, Murata M, Carninci P. 5' end-centered expression profiling using cap-analysis gene expression and next-generation sequencing. *Nat Protoc.* 2012 Feb 23; 7(3):542–61. <https://doi.org/10.1038/nprot.2012.005> PMID: 22362160
24. Buniello A, MacArthur JAL, Cerezo M, Harris LW, Hayhurst J, Malangone C, et al. The NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted arrays and summary statistics 2019. *Nucleic Acids Res.* 2019 Jan 8; 47(D1):D1005–12. <https://doi.org/10.1093/nar/gky1120> PMID: 30445434
25. Landrum MJ, Lee JM, Benson M, Brown GR, Chao C, Chitipiralla S, et al. ClinVar: improving access to variant interpretations and supporting evidence. *Nucleic Acids Res.* 2018 Jan 4; 46(D1):D1062–7. <https://doi.org/10.1093/nar/gkx1153> PMID: 29165669
26. Lizio M, Abugessaisa I, Noguchi S, Kondo A, Hasegawa A, Hon CC, et al. Update of the FANTOM web resource: expansion to provide additional transcriptome atlases. *Nucleic Acids Res.* 2019 Jan 8; 47(D1):D752–8. <https://doi.org/10.1093/nar/gky1099> PMID: 30407557
27. Kelley DR, Snoek J, Rinn JL. Basset: learning the regulatory code of the accessible genome with deep convolutional neural networks. *Genome Res.* 2016 Jul; 26(7):990–9. <https://doi.org/10.1101/gr.200535.115> PMID: 27197224
28. Yella VR, Kumar A, Bansal M. Identification of putative promoters in 48 eukaryotic genomes on the basis of DNA free energy. *Sci Rep.* 2018 Mar 14; 8(1):4520. <https://doi.org/10.1038/s41598-018-22129-8> PMID: 29540741
29. Schug J, Schuller W-P, Kappen C, Salbaum JM, Bucan M, Stoeckert CJ. Promoter features related to tissue specificity as measured by Shannon entropy. *Genome Biol.* 2005; 6(4):R33. <https://doi.org/10.1186/gb-2005-6-4-r33> PMID: 15833120
30. Visel A, Minovitsky S, Dubchak I, Pennacchio LA. VISTA Enhancer Browser—a database of tissue-specific human enhancers. *Nucleic Acids Res.* 2007 Jan; 35(Database issue):D88–92. <https://doi.org/10.1093/nar/gkl822> PMID: 17130149
31. Lee D, Karchin R, Beer MA. Discriminative prediction of mammalian enhancers from DNA sequence. *Genome Res.* 2011 Dec; 21(12):2167–80. <https://doi.org/10.1101/gr.121905.111> PMID: 21875935
32. Weingarten-Gabbay S, Nir R, Lubliner S, Sharon E, Kalma Y, Weinberger A, et al. Systematic interrogation of human promoters. *Genome Res.* 2019 Feb; 29(2):171–83. <https://doi.org/10.1101/gr.236075.118> PMID: 30622120
33. Welter D, MacArthur J, Morales J, Burdett T, Hall P, Junkins H, et al. The NHGRI GWAS Catalog, a curated resource of SNP-trait associations. *Nucleic Acids Res.* 2014 Jan; 42(Database issue):D1001–1006. <https://doi.org/10.1093/nar/gkt1229> PMID: 24316577
34. Roadmap Epigenomics Consortium, Kundaje A, Meuleman W, Ernst J, Bilienky M, Yen A, et al. Integrative analysis of 111 reference human epigenomes. *Nature.* 2015 Feb 19; 518(7539):317–30. <https://doi.org/10.1038/nature14248> PMID: 25693563
35. Sloutskin A, Danino YM, Orenstein Y, Zehavi Y, Doniger T, Shamir R, et al. ElemeNT: a computational tool for detecting core promoter elements. *Transcription.* 2015; 6(3):41–50. <https://doi.org/10.1080/21541264.2015.1067286> PMID: 26226151
36. Fornes O, Castro-Mondragon JA, Khan A, van der Lee R, Zhang X, Richmond PA, et al. JASPAR 2020: update of the open-access database of transcription factor binding profiles. *Nucleic Acids Res.* 2019 Nov 8; gkz1001.
37. Newman JRS, Keating AE. Comprehensive identification of human bZIP interactions with coiled-coil arrays. *Science.* 2003 Jun 27; 300(5628):2097–101. <https://doi.org/10.1126/science.1084648> PMID: 12805554

38. Ravasi T, Suzuki H, Cannistraci CV, Katayama S, Bajic VB, Tan K, et al. An Atlas of Combinatorial Transcriptional Regulation in Mouse and Man. *Cell*. 2010 Mar; 140(5):744–52. <https://doi.org/10.1016/j.cell.2010.01.044> PMID: 20211142
39. Grant CE, Bailey TL, Noble WS. FIMO: scanning for occurrences of a given motif. *Bioinforma Oxf Engl*. 2011 Apr 1; 27(7):1017–8. <https://doi.org/10.1093/bioinformatics/btr064> PMID: 21330290
40. Kopp W, Monti R, Tamburrini A, Ohler U, Akalin A. Deep learning for genomics using Janggu. *Nat Commun*. 2020 Jul 13; 11(1):3488. <https://doi.org/10.1038/s41467-020-17155-y> PMID: 32661261
41. Tatarinova T, Kryshchenko A, Triska M, Hassan M, Murphy D, Neely M, et al. NPEST: a nonparametric method and a database for transcription start site prediction. *Quant Biol Beijing China*. 2013 Dec; 1(4):261–71. <https://doi.org/10.1007/s40484-013-0022-2> PMID: 25197613
42. Umarov R, Kuwahara H, Li Y, Gao X, Solovyev V. Promoter analysis and prediction in the human genome using sequence-based deep learning models. *Bioinforma Oxf Engl*. 2019 Aug 15; 35(16):2730–7. <https://doi.org/10.1093/bioinformatics/bty1068> PMID: 30601980
43. Andersson R, Sandelin A. Determinants of enhancer and promoter activities of regulatory elements. *Nat Rev Genet*. 2020 Feb; 21(2):71–87. <https://doi.org/10.1038/s41576-019-0173-8> PMID: 31605096
44. He K, Zhang X, Ren S, Sun J. Deep Residual Learning for Image Recognition. *ArXiv151203385 Cs [Internet]*. 2015 Dec 10 [cited 2021 Mar 15]; Available from: <http://arxiv.org/abs/1512.03385> PMID: 27218121
45. Ioffe S, Szegedy C. Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. *ArXiv150203167 Cs [Internet]*. 2015 Mar 2 [cited 2021 Mar 15]; Available from: <http://arxiv.org/abs/1502.03167> PMID: 27218121
46. Maas AL, Hannun AY, Ng AY. Rectifier nonlinearities improve neural network acoustic models. In: *Proc icml*. Citeseer; 2013. p. 3.
47. Kingma DP, Ba J. Adam: A method for stochastic optimization. *ArXiv Prepr ArXiv14126980*. 2014;
48. Abadi M, Barham P, Chen J, Chen Z, Davis A, Dean J, et al. Tensorflow: A system for large-scale machine learning. In: 12th {USENIX} symposium on operating systems design and implementation ({OSDI} 16). 2016. p. 265–83.