1 **Predicting splicing patterns from the transcription factor binding sites in the**

2 **promoter with deep learning**

3

4 Tzu-Chieh Lin[1,†], Cheng-Hung Tsai[1,†], Cheng-Kai Shiau[1], Jia-Hsin Huang[1,2,*], Huai-Kuang

5 Tsai[1,2,*]

6

7 [1]Institute of Information Science, Academia Sinica, Taipei 11529, Taiwan

8 [2]Taiwan AI Labs & Foundation, Taipei 10351, Taiwan

9

10 [†]Both authors contributed equally to this work.

11 [*]Co-corresponding authors: Jia-Hsin Huang (jiahsin.huang@gmail.com) and Huai-Kuang Tsai

12 (hktsai@iis.sinica.edu.tw)

13

14

## Abstract

**Background**

17 Alternative splicing is a crucial mechanism of post-transcriptional modification responsible for

18 the transcriptome plasticity and proteome diversity of a metazoan cell. Although many splicing

19 regulations around the exon/intron regions have been discovered, the relationship between

20 promoter-bound transcription factors and the downstream alternative splicing remains largely

21 unexplored.

**Results**

23 In this study, we present computational approaches to decipher the regulation relationship

24 connecting the promoter-bound transcription factor binding sites (TFBSs) and the splicing

25 patterns. We curated a fine data set, including DNase I hypersensitive sites sequencing and

26 transcriptome in fifteen human tissues from ENCODE. Specifically, we proposed different

27 representations of TF binding context and splicing patterns to tackle the associations between the

28 promoter and downstream splicing events. Our results demonstrated that the convolutional neural

29 network (CNN) models learned from the TF binding changes in the promoter to predict the

30 splicing pattern changes. Furthermore, through an *in silico* perturbation-based analysis of the

31 CNN models, we identified several TFs that considerably reduced the model performance of
32 splicing prediction.

**Conclusion**

34 In conclusion, our finding highlights the potential role of promoter-bound TFBSs in influencing
35 the regulation of downstream splicing patterns and provides insights for discovering alternative
36 splicing regulations.

37

38

## Keywords

40 Alternative splicing; TFBS; deep learning; CNN; transcriptional regulation

41

42

## Background

44 Gene splicing endows the transcriptional diversity of the metazoan genome. Splicing is the
45 process by which introns are removed from the nascent pre-mRNA and exons are joined,
46 generating the functional mRNA. Alternative splicing (AS), the selective removal of exons and
47 reconnection of exons by multiple processes, is known to play a pivotal role in regulatory
48 pathways from invertebrates to mammals [1, 2]. By the regulatory mechanism of AS, a single
49 gene is capable of generating multiple RNA molecules encoding proteins with different functions
50 [3]. The importance of AS lies in the evidence that the human genome has been estimated more
51 than 95% of multi-exon genes undergo alternative splicing in an underlying tissue-specific
52 manner [4]. Moreover, the variations in splicing patterns are prevalent to associate with many
53 complex diseases in humans [5, 6], and one-third of all disease-associated alleles have been
54 estimated to alter splicing [7].

55

56 Studies on AS regulation have mainly focused on the sequence information of spliced exons and
57 flanked introns. Machine learning has unprecedented performance in predicting exon-
58 inclusion/skipping levels in bulk tissues or single cells. Several computational models to derive
59 "splicing codes" that predict splice site selection in a genomic sequence successfully capture
60 patterns around the skipped exon and elucidate complex regulatory mechanisms from genomic

61  and epigenomic features [8–12]. Despite many efforts to characterize the splicing regulatory

62  codes within the splice sites, the extent and effects of transcription machinery at the relatively

63  distant promoter regions in splicing regulation remain unsolved.

64

65  In the past decades, AS has been generally accepted to be tightly coupled with RNA polymerase

66  transcription of the nascent pre-mRNA [13, 14]. Two prevailing models have been proposed to

67  explain the coupling between alternative splicing and transcription: the recruitment model [15,

68  16] and the kinetics model [14]. Notably, the chromatins are mostly not in linear form; the

69  transcription complex on a promoter affects the recruitment of splicing factors and elongation of

70  RNA polymerase II to promote exon exclusion through chromatin looping [17]. In addition,

71  various DNA-binding proteins have been reported to influence the AS patterns by changing

72  epigenetic conditions in the promoter [18].

73

74  Each gene contains a set of unique combinations of TF binding sites (TFBSs) in the promoter

75  that determines its temporal and spatial expression. Transcriptional regulation is usually a

76  combinatorial effect of multiple TFs binding to *cis*-regulatory elements located in the proximate

77  and distal regions from transcription start sites [19]. Date to 20 years ago, the regulation of exon

78  splicing patterns was demonstrated directly through the specific TFBS occupancy in the

79  promoter [20, 21]. Moreover, the coupling of promoter and splicing is later proposed with

80  extensive regulator mechanisms [22, 23]. Given the three-dimensional folding of chromatin

81  loops, the proximal promoter- or distal enhancer-bound factors joined into transcription

82  compartments correlate with alternative splicing of exons [24]. Although the biological findings

83  connect the promoter with AS by focusing on a few gene models, the hypothesis that promoter

84  architecture in terms of TFBS composition regulates AS remains unexplored at the genome-wide

85  level.

86

87  In this study, we developed analytical strategies to approach this question using data of both

88  RNA-seq and DNase-seq in pairs across the different human tissues from the ENCODE project.

89  We first considered the associations between the occurrences of more than 300 TF binding

90  motifs in the promoter and the corresponding splicing patterns. Secondly, we examined whether

91  the changes in TF binding condition were able to predict the splicing change by studying the

92    relative changes of the splice-in percent (PSI) values between any paired tissues. Then, we

93    conducted machine learning methods and deep learning neural networks to predict the splicing

94    patterns. Notably, the convolutional neural network (CNN) models that took complete TF

95    occupancy information in promoter regions as input achieved the highest performance at 0.889

96    of the area under receiver operating characteristic curve (AUROC). Lastly, we applied the

97    importance analysis of the CNN models for each TF and identified some important TFs that

98    affecting the splicing prediction genome-widely.

99

100

## Results

102    In this study, we considered the cassette exon splicing, which is the most frequent alternative

103    splicing type in the human genome [36]. We proposed two scenarios to examine the relationship

104    between TFBSs in the promoter and the splicing patterns of the gene. First, we asked if

105    compositions of TFBS occupancies, which were defined as the expressed TFs (TPM > 1) in the

106    given tissues and their binding motifs in the open chromatin regions, are associated with the

107    splicing patterns of the gene. Second, we asked if the changes of TF binding condition in the

108    promoter modify the splicing efficiency of the cassette exon usage by comparing their PSI values.

109    The data preprocessing procedures for TFBS identification in the promoter and exon-skipped

110    events are illustrated in Fig. 1A. The TF binding profiles of each promoter were curated by

111    integration of DNase-seq for open-chromatin regions, human TF motif scan, and expression

112    profile across 15 tissues. The splicing patterns of each gene were analyzed based on the

113    transcriptome in different tissues.

114

**Characterizing the TFBS occupancies in the promoter and first cassette exons across tissues**

117    We investigated the associated relationship between the TFBSs in the promoter and the first

118    cassette exon, which is relatively closed to the promoter. The distribution of the PSI values as

119    exon usage levels was bimodal across 15 human tissues (Fig. 1B). Here, we defined the PSI

120    values smaller than 0.2 and larger than 0.8 as the exclusion form and inclusion form, respectively.

121    Based on the criteria, the usage of the first cassette exons of human genes across 15 tissues was

122    mostly skewed in either one of the categories, *i.e.,* exclusion or inclusion forms (Fig. 1C). There

123    were only 4.6% of genes having both splicing forms in different tissues.

124

125    Experimental studies have shown that the promoter architecture, by using different gene

126    promoters, affects the splicing patterns of the exon skipping in the gene bodies [37, 38].

127    Following this idea, we sought to examine whether the promoter architecture in terms of TFBS

128    occupancies as the features determine the inclusion or exclusion of the first cassette exon. First,

129    we asked which TFBSs were predominant within the promoters of these genes with different

130    splicing patterns of their first cassette exon. In order to address this, the discrepancy between the

131    frequency of individual TFBS on the promoters of the exclusion sets and that of the inclusion

132    sets was evaluated independently by using a chi-squared ($\chi^2$) test for each tissue. Considering an

133    adjusted significance level of $p$-value $< 0.001$ after Bonferroni correction, more than half of TF

134    binding motifs are significantly enriched in the promoter of either exclusion or inclusion sets. In

135    addition, we calculated the gene expression specificity index tau [31, 32] for each TF and set 0.8

136    as the cut-off for tissue-specific TFs. However, there is no particular enrichment of TFs showing

137    more enriched across statistical significance ranks (Fig. 1D, right panel).

138

139    Next, we considered the complex relationship among TFBSs within promoters on the prediction

140    of splicing patterns by using a machine learning approach. We employed the XGBoost method

141    [39], a decision-tree-based ensemble model, and used the presence of TFBSs within the open

142    chromatins of promoter as input data to predict the inclusion or exclusion of the first cassette

143    exons. Due to the coarser resolution of DNase-seq and *in silico* motif scanning to profile the

144    TFBS occupancies, we noticed that some genes share identical features in different tissues. We

145    thus removed the samples that share identical features in the training data from the testing data of

146    the given tissues to avoid the fallacy of prediction accuracy in the cross-tissue evaluation.

147    Herein, we proposed three different cross-fold validation schemes in order to properly evaluate

148    prediction performance (Fig. 1E). For event-wise scheme, we randomly left 10% of promoter-

149    splicing pairs as the independent testing data and performed a 10-fold cross-validation (CV). For

150    tissue-wise scheme, we conducted leave-one-tissues-out cross-validation by treating the

151    promoter-splicing pairs from a single tissue as the independent testing data. For gene-wise

152    scheme, we used 90% of genes with all promoter-splicing pairs across tissues to train model and

153    remained 10% of genes were for an independent testing set. In Fig. 1F, three evaluation metrics,

5

154    including F1-score, AUROC, and accuracy, were shown to compare the prediction performance

155    in different CV schemes. Interestingly, the prediction performance using event-wise scheme

156    achieved an F1-score and AUROC closed to 0.80 (Fig. 1F, green bars). In the cross-tissue

157    validation results, we further observed that the overall performance of the models obtained an

158    average AUROC of 0.84 (Fig. 1F, purple bar). However, all three metrics underlying gene-wise

159    CV could yield slightly better than random guess at 0.50 (Fig. 1F, yellow bars).

160

161    It is worth noting that the gene-wise CV scenario indeed examined whether the generalization of

162    a trained model enables to classify the splicing events using the unseen promoter information

163    about TF binding profiles, which were not included in the training dataset. We later addressed a

164    following question if the same gene promoter in different tissues both present in the training and

165    testing sets was critical for prediction performance. Subsequently, we split the genes into three

166    groups, i.e., one-sided, both-sided, and singleton, according to their splicing forms across all

167    tissues and re-examined the results of prediction accuracy in the individual tissues. In contrast to

168    the genes with one-sided and both-sided splicing forms, the trained models using data from other

169    tissues did not predict the splicing forms of the singleton genes correctly in the given tissue (Fig.

170    1G, left panel). Furthermore, we counted the number of genes in the respective groups (Fig. 1G,

171    right panel), and found that a good overall performance of the models underlying tissue-wise CV

172    was dominant by the large number of genes with one-sided splicing form across all tissues. The

173    poor prediction on those small portions of singleton genes (less than 200) did not cause a drastic

174    drop in overall prediction accuracy. In summary, our current approach failed to construct the

175    models with generalization ability to infer the splicing forms using promoter information that

176    pertains to TF binding profiles.

177

178    **Changes of TF binding to the promoter reflect the distinct exon splicing phases**

179    In this section, we sought to examine whether changes of individual TF binding to promoter alter

180    the splicing efficiency that was estimated by PSI values. The PSI value summarizes the splicing

181    condition of the constitutive exons that are included in all or part of transcripts from expressed

182    isoforms [40]. As the fact that ranges of PSI values of different genes are varied across 15 tissues,

183    the genes differ from each other in terms of their efficiency of splicing first cassette exon into the

184    expressed isoforms. As a result, the efficiency of exon usage should be considered for each gene

185     itself instead of the absolute PSI ($\Psi$) value. To this end, we applied the Z-score transformation to

186     normalize the absolute PSI scores of all genes. Of note, some genes that had a smaller PSI range

187     ($< 0.2$) and/or expressed in less than three tissues were discarded in the following experiments.

188     We then defined the top 20% and last 20% of transformed $Z_\Psi$ scores in each gene as the two

189     distinct phases of exon usage, i.e., low and high splicing efficiency respectively (Fig. 2A). To

190     test the hypothesis that changes of TFBS in the open chromatin of the promoter are associated

191     with splicing phase change, the differences of two $Z_\Psi$ and their TF binding occupancies in a

192     given paired tissues for each gene were calculated (Fig. 2B). The distribution of delta $Z_\Psi$ scores

193     was shown in Fig. 2C, where the unchanged group (same splicing phase) was below 1 and the

194     changed group (different splicing phase) was larger than 1.8. Of note, no overlapped events were

195     observed between concordance and discordance groups.

196

197     To examine the association between TFBS-occupied difference and splicing phase for individual

198     TFs, we constructed a $2 \times 2$ contingency table for each TF. Specifically, for each tissue pair in

199     one gene, we assigned the pair into groups according to whether its TFBS occupancy is changed

200     ($\Delta TF\alpha = 0$ or $\Delta TF\alpha = 1$), and whether the splicing phase is changed (concordance or

201     discordance). We thus calculated the odds ratio from contingency table and applied chi-squared

202     test. About two-third of TFs, their binding occupancy changes were significantly associated with

203     splicing phase changes (N = 203, adj. $p$-value $< 10^{-3}$, Fig. 2D). Since every tissue usually

204     expresses different sets of TFs to control the cell fate [41, 42], we estimated the tissue specificity

205     of TF expression by tau score [32]. More than half (53%) of TFs among those non-significant

206     groups were ubiquitously expressed, while most of the TFs (75%) among those significant

207     associations with splicing phase change were tissue-specifically expressed (Fig. 2D). Of note, the

208     open chromatin regions in the promoter of the same gene in different tissues show less variations.

209     Thus, TFBSs without filtered by expression profiles of given TFs did not show any significant

210     association. Therefore, although the DNA sequences of the promoter are identical, the

211     divergence on the TF expression across different tissues is a likely regulating mechanism to

212     affect the splicing phase change.

213

214     **Machine learning confirm the association between TF binding changes and splicing phase**

215     **shift**

216   Next, we employed different machine learning algorithms, including logistic regression,

217   XGBoost (ensemble tree algorithm), and deep neural network methods, to test whether the

218   combinations of TF binding changes predict the splicing phase changes. To monitor sensitivity

219   and specificity simultaneously, we assessed the models using the AUROC in the plot of the true

220   positive rate (TPR) against the false positive rate (FPR) for five-fold cross-validation tests (Fig.

221   3A). Three classifiers achieved an average AUROC of 0.691, 0.766, 0.771 for logistic regression

222   (LReg), deep neural network (DNN), and XGBoost (XGB) models, respectively on all the events

223   of the dataset. Since there were imbalanced data sets in the changed and unchanged groups, the

224   area under the precision-recall curve (AUPRC) is also instructive to assess the model

225   performance (Fig. 3B). The XGB models also achieved a greater mean AUPRC of 0.630 than

226   0.531 and 0.624 respectively for LReg and DNN. Because there is often more than one binding

227   site in the promoter for each TF, we also constructed other ML models using frequencies of all

228   possible TF binding site changes between promoters as the features. The overall performance of

229   prediction of splicing phase change was decreased about 6% based on AUROC. This indicates

230   that the decision tree-based ML method could not deal with the frequencies of TFBSs change

231   properly.

232

233   **Integration of TFBS locations in the promoter using deep learning models improve**

234   **prediction performance**

235   We next integrated the position information of TFBSs in the promoter as the features to train the

236   deep neural network (DNN) and convolution neural network (CNN) models respectively. The

237   two-dimension array consisting of 2,500 bp and 345 TF binding changes were used as the input

238   features as shown in Fig. 4A. The architecture of the CNN model includes the one-dimensional

239   convolutions kernels, which are designed as the filters for revealing the combinations of TF

240   binding changes. The convolution layers are followed by a max-pooling layer with sliding

241   window size and a stride step of 10 units. And a single flatten layer with 256 neurons was used to

242   summarize all features and followed by three hidden layers. To prevent overfitting, the dropout

243   technique was applied to remove 25% of the connected neurons in the flatten and hidden layer

244   during the training (26).

245

246   Training the network with input matrices including both TFBS and their interactions with other

8

247 TFs markedly impacts the performance of the splice predictions. In contrast to the performance

248 of previous DNN models using only TFBS changes input (Fig. 3A), current DNN classifiers

249 achieved greater AUROC, increasing from the average 0.766 to 0.853 (Fig. 4B). The CNN

250 classifiers achieved an even greater AUROC of 0.889 (Fig. 4B). Additionally, CNN models

251 achieved greater AUPRC for all five-fold experiments than DNN models, increasing the average

252 from 0.730 to 0.782 (Fig. 4C).

253

254 **Evaluation of TF changes on the splicing patterns**

255 We next to understand the importance of TF motifs on splicing patterns utilized by the network

256 to achieve its remarkable accuracy. In brief, we performed systemic *in silico* substitution of each

257 TF change as zero, then measured the effects on the CNN model's prediction. The importance of

258 each TF was estimated by the fraction of changed prediction under the *in silico* substitution. The

259 underlying idea is if assume a TF plays a key role in regulating splicing patterns, the prediction

260 output of the machine learning model should change dramatically after substitution rather than

261 other TF. We performed importance analysis on each TF and ranked them by their importance

262 measurement, and found that a small proportion of TFs resulted in dramatical changes in the

263 splicing prediction (Fig. 5A). As most of TFs had a little effect on the CNN model performance,

264 we highlighted top-ranked 19 TFs with outlier values based on the interquartile range rule (Q3 +

265 $1.5 \times$ IQR) as the candidate splicing regulators.

266

267 Previous studies have demonstrated that binding of the acetyltransferase p300 at promoter

268 regions modifies acetylation of splicing factors, and thereby modulate the alternative splicing

269 pattern of the gene [43, 44]. We thus submitted our 19 candidate TFs and p300 to the STRING

270 database [45] for identification of their interactions. We applied default settings to search both

271 functional and physical protein associations with medium confidence score of 0.400 in the

272 STRING database (ver. 11.5). Then, we configured the network between query proteins only to

273 reveal the associations among them. Interestingly, the network was relatively less complex and

274 p300 were thought of as a hub gene associated with nine out of 19 top-ranked TFs (Fig. 5B).

275 Moreover, the interaction between KLF14 and p300 is experimentally and functionally

276 confirmed that the binding of KLF14 to the promoter recruits p300 to increase the levels of

277 acetylation associated with transcriptional activation [46]. Although the interaction between

278   KLF14 and p300 on the gene activation was not investigated in the context of splicing,
279   compelling evidence showing a direct link between histone modification and splicing [17, 18]
280   raises the intriguing possibility of KLF14/p300 complex in modulating exon splicing. Similarly,
281   some top-ranked TFs might share a common mechanism in regulating RNA splicing via
282   recruitment of p300 to promote the deposition of histone acetylation at the promoter.

283

284   Lastly, to further confirm our *in silico* prediction for potential splicing regulators, we obtained
285   the K562 CTCFL shRNA knock-down RNA-Seq data [47] and its control from previous
286   research [48]. We re-analyzed the splicing status by calculating PSI through MISO and applied
287   Z-score transformation using the previous method in machine learning model training. We
288   observed the $\Delta Z_\Psi$ values of CTCFL-target genes were higher than that of non-target genes
289   significantly (Fig. 5B, with *p*-value < 0.0001, Wilcoxon rank-sum test). This revealed in the
290   CTCFL deplete condition, genes targeted by CTCFL change their first skipped exon usage thus
291   influence $\Delta Z_\Psi$. We further seek for case studies to investigate how splicing status changed in
292   CTCFL-target genes under CTCFL depletes (Fig. 5C). The first skipped exon in
293   ENSG00000101096 has a higher skipped exon usage and increases the average PSI value. In
294   contrast, in ENSG00000147364 the first skipped exon usage reduced in the CTCFL deplete
295   condition thus has a lower average PSI value. These results suggest that CTCFL can influence
296   the splicing pattern. Nevertheless, CTCFL shows a dual function in splicing regulation, not only
297   increase skipped exon usage but also reduce usage in some genes. This result also matches the
298   previous study on CTCFL-depletion mediate alternative splicing change in MCF7 cell line [49].
299   In the CTCFL-depletion they detect exclusion of 361 and the inclusion of 221 alternative exons
300   compared to the normal condition. The CTCFL can influence the recruitment of RNAPII and
301   thus impact the RNAPII elongation speed and finally alter the splicing result of pre-mRNA.
302   Overall, these results support the feasibility of our modeling and importance analysis approaches
303   for *in silico* prediction.

304

305

## Discussion

307   The applications of machine learning methods to characterize the regulatory potential of genomic
308   sequences on alternative splicing have been a subject of interest for over a decade [8, 50]. Instead

309 of using the genomic information around the splicing exons, in this study, we focused on the
310 upstream promoter region for predicting downstream exon-skipped events genome-widely. In
311 contrast to some previous study using the DNA sequences directly [8, 9, 11], one major
312 difference of our approach is that we applied TF binding motif scan with prior domain
313 knowledge to represent the sequence information in the promoter. We demonstrate how the
314 promoter signals in terms of TFBS profiles can be integrated using machine learning approaches
315 for the further implication of association between the promoter and alternative splicing. Our
316 results showed that the prediction accuracy differed among the different algorithms and input
317 information. Notably, one-dimensional CNN architecture is highly capable of learning the
318 regulatory code from the TF binding changes in the promoter to discriminate the splicing
319 patterns (Fig 4).

320

321 The main drawback of this study is the limited number of tissues because we aimed to use a
322 high-quality dataset to avoid the noise and artifacts in the DNase-seq and RNA-seq datasets
323 conducted by different labs. Thus, we excluded any experiments that did not meet every quality
324 standard defined by ENCODE. When conducting the data analyses, we noticed that the splicing
325 forms for most of the gene were not varied extensively in these 15 tissues (Fig. 1C). Inspired by
326 the previous study to avoid fallacy of model performance using alternative cross-fold validation
327 schemes properly [51], we implemented three different CV schemes, i.e., event-wise, tissue-wide,
328 and gene-wise, to evaluate generation performance carefully. In the course of examining the
329 difference across three CV schemes to find possible reasons for high performance in the tissue-
330 wise evaluation, we noticed that majority of genes were expressed in more than two tissues and
331 displayed same splicing form. Because every gene promoter in different tissues shares most
332 TFBS features, the event- and tissue-wise schemes are subject to the problem of test set
333 contamination and could lead to an artificially inflated accuracy in this study. On the bright side,
334 there is considerable room for improvement in model generalization by collecting varied splicing
335 forms of every gene from different tissues extensively to evaluate promoter-splicing interactions.

336

337 To address the problem of shared TFBSs in promoter across tissues, we turned to look at the TF
338 binding changes in promoter (Fig 2B). Notably, this approach diminished the high similarity of
339 TFBS features in tissues and making a comparison in any given paired tissues also augmented

11

340    the datasets incrementally for improvement of the model training. On the other hand, we

341    considered the changes in splicing efficiency ($\Delta Z_\Psi$) by introducing a transformation procedure of

342    absolute PSI values into the efficiency of exon usage. Our computational method is different

343    than a previous study using the absolute PSI values to estimate splicing efficiency directly [52].

344    The fact that the ranges of the PSI values in a particular gene across 15 tissues are mostly

345    ununiformed distribution is evident as the averaged PSI values of genes from closed to 0 or 1

346    (Fig 2A). The Z-transform method could remain commensurate in the scale to measure splicing

347    efficiency for each gene accordingly. In addition, instead of using fixed arbitrary cutoff values

348    (e.g., $\Psi < 0.2$ and $\Psi > 0.8$) to subsect the splicing status, we applied a percentile threshold to

349    divide genes into two tendencies, i.e., "splice-in" or "splice-out". This approach avoids that those

350    small-PSI-ranged genes are skew to be classified into a single group of splice-in or splice-out.

351    Based on our observation, it is perhaps noteworthy to rethink about the definition of the splicing

352    status using PSI as a metric to explore alternative solutions in discovery of splicing mechanisms.

353    By carefully considering the fundamental issues in our preprocessing procedures on data, this

354    study provides a different perspective to study how TFs in promoter affects the exon splicing

355    genome-widely.

356

357    To train the prediction model of splicing phase shift, we used two different input data, *i.e.*, an

358    array of TF binding changes and a matrix of full TF binding changes along with the promoter

359    regions. Our results demonstrated that training the DNN models with varying input of TF

360    binding context noticeably impacts the accuracy of the splicing phase shift prediction (Fig. 3 and

361    4). Despite amount of trainable network parameters drastically are increased when using an input

362    of TF binding context, DNN models is capable to automatically learn the task from the training

363    data. Remarkably, CNNs achieved even higher prediction performance than DNNs with matrices

364    of TF binding context (Fig. 4). In contrast to DNNs, CNNs indeed are designed to deal with

365    high-dimensional inputs by applying of a serious of convolutional and pooling steps [53, 54]. A

366    likely explanation for high accuracy boosting in CNNs is the convolutional operations, which

367    learned higher-level features from the combinations of different TF changes. With the good

368    prediction performance of CNN models, the importance analysis experiments allowed us to

369    identify a couple of TFs that potentially involve in splicing regulation. To our knowledge, our

370    study is the first genome-wide effort to investigate that the splicing pattern changes across tissues

371    were accurately predicted from the TF binding occupancies in the promoter.

372

373

## Materials and Methods

**Data processing and sample selection**

We downloaded both the DNase-seq peak BED files and the RNA-seq data for 15 human tissues from the ENCODE data portal [25]. To obtain high quality of data, the data without any flags, such as insufficient read depth, in the experimental metadata that were reported by the ENCODE Data Coordination Center are used in the following experiments. For DNase-seq datasets, the standard pipeline (accession: ENCPL201DNS for single-ended data, ENCPL202DNS for paired-ended data) from ENCODE called the peaks using hotspot2 algorithm with 1% false-discovery rate. For RNA-seq data, the ENCODE RNA-seq pipeline for long RNAs (accession: ENCPL002LSE for single-ended data, ENCPL002LPE for paired-ended data) used the STAR program for mapping the reads and the RSEM algorithm for quantification of genes. We used genomic and annotation files of the human reference genome version GRCh37 as provided by release V19 of GENCODE [26].

387

**Identification of putative *in vivo* TF binding sites**

The DNase-seq peaks were used to define the open chromatin regions in the promoter regions (−2 kb to +500 bp from the transcription start site). We downloaded TF motifs from the JASPAR database (ver. 2018) [27] and excluded the fusion TF (i.e., EWSR1/FLI1 fusion) and older versions of motifs from the same TF, as a result, we obtained 407 TF binding motifs from JASPAR. Later, we scanned the sequence from each open chromatin region for each TF binding motif in position-weight-matrix (PWM) format, using FIMO from the MEME (Motif-based sequence analysis tools) suite [28]. Of note, we applied the FIMO with a threshold false discovery rate of $< 10^{-3}$, which is less stringent than the general recommended parameter ($< 10^{-4}$) for putative *cis*-regulatory elements detection. Since we only considered TF binding sites located in the open chromatin regions, the general parameter is too stringent for our purpose.

399

**RNA-seq processing and calculation of cassette exon usage (PSI)**

To estimate the splicing level for each exon and tissue, we first used CATANA [29] to annotate

13

402 AS events in all human transcripts for the AS annotation index file. The BAM files of RNA-seq

403 data generated by the ENCODE were used to estimate the percent spliced-in (PSI) values for the

404 first cassette exon of the protein-coding genes using the MISO (Mixture of Isoforms) tool [30].

405 For the calculation of the $Z_\Psi$ score, we first selected the genes that PSI range is larger than 0.2

406 across different tissues and then standardized their PSI by z-score transformation for each gene.

407

408 **Enrichment analysis**

409 We analyzed the association of TF binding occupancies and splicing patterns from $2 \times 2$

410 contingency tables categorizing all human genes according to the occurrences of binding sites for

411 a given TF and splicing patterns (exclusion or inclusion in Fig. 1D). In parallel, we built the

412 contingency table to analyze the association between TFBS-occupied differences and splicing

413 phases (concordance or discordance in Fig. 2D) for each TF. The odds ratio (OR) based on the

414 contingency table was calculated for each TF and a chi-squared ($\chi 2$) test was applied to

415 determine the statistical significance of the association. The *p*-value is adjusted by Bonferroni

416 correction (and its −log10 transformation) for the association, and the odds ratio with log2

417 transformation is a measure of the effect size. The adjusted *p*-value < 0.001 is considered as

418 significant.

419

420 **Tau index of TF tissue specificity**

421 We calculated the tissue specificity index tau [31, 32] using the gene expression of each TFs

422 across different tissues, as follows:

$$\text{tau} = \frac{\sum_{i=1}^{n}(1-\hat{x}_i)}{n-1}; \; \hat{x}_i = \frac{x_i}{\max_{1 \le i \le n} x_i}$$

424 where $x_i$ represents the gene expression of TF $x$ in tissue $i$ ; and $n$ is the number of tissues

425 expressing the TF (TPM > 1). We then adopted the cut-off of tau based on a previous study [33]

426 and defined the TFs with tau ≥ 0.8 as tissue-specifically expressed.

427

428 **Machine learning and deep learning models**

429 In order to get a better prediction power, we compared the accuracy between four methods,

430 logistic regression, XGBoost, deep neural network (DNN), and convolutional neural network

431 (CNN). To avoid biases caused by imbalanced data, we applied a balanced sampler as the

432    concept            described          on          the          imbalanced-dataset-sampler            (from

433    https://github.com/ufoym/imbalanced-dataset-sampler) to our training dataset before model

434    training. We trained the basic logistic regression model with default parameter settings described

435    in the scikit-learn [34]. For the XGBoost model, we limited the max tree depth to 6, set the eta by

436    1, and used gbtree as a booster.

437

438    In this research, we implemented our DNN and CNN models using the PyTorch framework [35].

439    The architecture of DNN began with flattening the input data and followed by 3 dense layers,

440    with 512, 256, and 128 nodes, respectively. ReLU activation function was applied on the output

441    of each dense layer and then followed by a dropout layer to randomly set 25 percent of input

442    units to 0. The sigmoid function was applied to the final output of the tensor to generate binary

443    classification predictions.

444

445    The architecture of CNN is similar to DNN with some modifications. The input data was first

446    processed through a convolution layer which followed by the ReLU activation function, max

447    pooling layer and a dropout layer, and then connected to 2 dense-ReLU-dropout units as

448    described above, both with 128 nodes. The sigmoid function is also used to do the binary

449    classification task.

450

451    **Importance analysis**

452    To extract informative TF binding features from the CNN model, we performed an *in silico*

453    perturbation-based analysis to observe the impact on the perturbed input data. Similar to the

454    previous method, we perturb the input by assigning a zero value for a given TF of the input

455    feature (zero-out operation) and perform inference on the trained model. The feature importance

456    through zero-out operation was measured by the output changing ratio. Output changing ratio

457    was defined as $N_{changed}$ / $N_{total}$, where $N_{changed}$ represents the count of changed output label after

458    zero-out and $N_{total}$ represents the total input delta instances number with corresponding TF

459    binding site.

460

461

462    **Availability of data and materials**

463 The source codes supporting the conclusions of this study are available at GitHub repository

464 (https://github.com/bio-it-station/DoTA).

465

466

## References

468 1. Cáceres JF, Kornblihtt AR. Alternative splicing: multiple control mechanisms and

469 involvement in human disease. Trends in Genetics. 2002;18:186–93.

470 2. Ule J, Blencowe BJ. Alternative Splicing Regulatory Networks: Functions, Mechanisms, and

471 Evolution. Molecular Cell. 2019;76:329–45.

472 3. Kelemen O, Convertini P, Zhang Z, Wen Y, Shen M, Falaleeva M, et al. Function of

473 alternative splicing. 2013.

474 4. Pan Q, Shai O, Lee LJ, Frey BJ, Blencowe BJ. Deep surveying of alternative splicing

475 complexity in the human transcriptome by high-throughput sequencing. Nature Genetics.

476 2008;40:1413–5.

477 5. Tazi J, Bakkour N, Stamm S. Alternative splicing and disease. Biochimica et biophysica acta.

478 2009;1792:14–26.

479 6. Daguenet E, Dujardin G, Valcárcel J. The pathogenicity of splicing defects: mechanistic

480 insights into pre☐ mRNA processing inform novel therapeutic approaches. EMBO reports.

481 2015;16:1640–55.

482 7. Havens MA, Duelli DM, Hastings ML. Targeting RNA splicing for disease therapy. 2013.

483 8. Barash Y, Calarco JA, Gao W, Pan Q, Wang X, Shai O, et al. Deciphering the splicing code.

484 Nature. 2010;465:53–9.

485 9. Bretschneider H, Gandhi S, Deshwar AG, Zuberi K, Frey BJ. COSSMO: Predicting

486 competitive alternative splice site selection using deep learning. In: Bioinformatics. Oxford

487 University Press; 2018. p. i429–37.

488 10. Bao S, Moakley DF, Zhang C. The Splicing Code Goes Deep. Cell. 2019;176:414–6.

489    11. Louadi Z, Oubounyt M, Tayara H, Chong KT. Deep Splicing Code: Classifying Alternative

490    Splicing Events Using Deep Learning. Genes. 2019;10:587.

491    12. Baeza-Centurion P, Miñana B, Schmiedel JM, Valcárcel J, Lehner B. Combinatorial

492    Genetics Reveals a Scaling Law for the Effects of Mutations on Splicing. Cell. 2019;176:549-

493    563.e23.

494    13. Carrocci TJ, Neugebauer KM. Pre-mRNA Splicing in the Nuclear Landscape. Cold Spring

495    Harbor Symposia on Quantitative Biology. 2020;:040402.

496    14. Kornblihtt AR, Schor IE, Alló M, Dujardin G, Petrillo E, Muñoz MJ. Alternative splicing: a

497    pivotal step between eukaryotic transcription and translation. Nature reviews Molecular cell

498    biology. 2013;14:153–65.

499    15. Muñoz MJ, la Mata M, Kornblihtt AR. The carboxy terminal domain of RNA polymerase II

500    and alternative splicing. 2010.

501    16. Huang Y, Li W, Yao X, Lin Q-J, Yin J-W, Liang Y, et al. Mediator complex regulates

502    alternative mRNA processing via the MED23 subunit. Molecular cell. 2012;45:459–69.

503    17. Rambout X, Dequiedt F, Maquat LE. Beyond Transcription: Roles of Transcription Factors

504    in Pre-mRNA Splicing. Chemical Reviews. 2018;118:4339–64.

505    18. Kolathur KK. Role of promoters in regulating alternative splicing. Gene. 2021;782:145523.

506    19. Komili S, Silver PA. Coupling and coordination in gene expression processes: A systems

507    biology view. 2008.

508    20. Monsalve M, Wu Z, Adelmant G, Puigserver P, Fan M, Spiegelman BM. Direct coupling of

509    transcription and mRNA processing through the thermogenic coactivator PGC-1. Mol Cell.

510    2000;6:307–16.

511    21. Auboeuf D, Hönig A, Berget SM, O'Malley BW. Coordinate regulation of transcription and

512    splicing by steroid receptor coregulators. Science. 2002;298:416–9.

513    22. Kornblihtt AR. Promoter usage and alternative splicing. Current Opinion in Cell Biology.

514    2005;17:262–8.

515    23. Maniatis T, Reed R. An extensive network of coupling among gene expression machines.

516    Nature. 2002;416:499–506.

517    24. Mercer TR, Edwards SL, Clark MB, Neph SJ, Wang H, Stergachis AB, et al. DNase I-

518    hypersensitive exons colocalize with promoters and distal regulatory elements. Nature Genetics.

519    2013;45:852–9.

520    25. Luo Y, Hitz BC, Gabdank I, Hilton JA, Kagda MS, Lam B, et al. New developments on the

521    Encyclopedia of DNA Elements (ENCODE) data portal. Nucleic Acids Res. 2020;48:D882–9.

522    26. Harrow J, Frankish A, Gonzalez JM, Tapanari E, Diekhans M, Kokocinski F, et al.

523    GENCODE: the reference human genome annotation for The ENCODE Project. Genome Res.

524    2012;22:1760–74.

525    27. Khan A, Fornes O, Stigliani A, Gheorghe M, Castro-Mondragon JA, van der Lee R, et al.

526    JASPAR 2018: update of the open-access database of transcription factor binding profiles and its

527    web framework. Nucleic Acids Res. 2018;46:D260–6.

528    28. Bailey TL, Johnson J, Grant CE, Noble WS. The MEME Suite. Nucleic Acids Res.

529    2015;43:W39-49.

530    29. Shiau C-K, Huang J-H, Tsai H-K. CATANA: a tool for generating comprehensive

531    annotations of alternative transcript events. Bioinformatics. 2019;35:1414–5.

532    30. Katz Y, Wang ET, Airoldi EM, Burge CB. Analysis and design of RNA sequencing

533    experiments for identifying isoform regulation. Nat Methods. 2010;7:1009–15.

534    31. Kryuchkova-Mostacci N, Robinson-Rechavi M. A benchmark of gene expression tissue-

535    specificity metrics. Briefings in Bioinformatics. 2017;18:205–14.

536   32. Yanai I, Benjamin H, Shmoish M, Chalifa-Caspi V, Shklar M, Ophir R, et al. Genome-wide

537   midrange transcription profiles reveal expression level relationships in human tissue

538   specification. Bioinformatics. 2005;21:650–9.

539   33. Guschanski K, Warnefors M, Kaessmann H. The evolution of duplicate gene expression in

540   mammalian organs. Genome Res. 2017;27:1461–74.

541   34. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn:

542   Machine Learning in Python. Journal of Machine Learning Research. 2011;12:2825–30.

543   35. Paszke A, Gross S, Massa F, Lerer A, Bradbury J, Chanan G, et al. PyTorch: An Imperative

544   Style, High-Performance Deep Learning Library. 2019.

545   36. Yeo G, Holste D, Kreiman G, Burge CB. Variation in alternative splicing across human

546   tissues. Genome Biology. 2004;5:R74.

547   37. Cramer P, Pesce CG, Baralle FE, Kornblihtt AR. Functional association between promoter

548   structure and transcript alternative splicing. Proceedings of the National Academy of Sciences of

549   the United States of America. 1997;94:11456–60.

550   38. Pagani F, Buratti E, Stuani C, Bendix R, Dörk T, Baralle FE. A new type of mutation causes

551   a splicing defect in ATM. Nature genetics. 2002;30:426–9.

552   39. Chen T, Guestrin C. XGBoost: A Scalable Tree Boosting System. In: Proceedings of the

553   22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. New

554   York, NY, USA: Association for Computing Machinery; 2016. p. 785–94.

555   40. Schafer S, Miao K, Benson CC, Heinig M, Cook SA, Hubner N. Alternative Splicing

556   Signatures in RNA-seq Data: Percent Spliced in (PSI). Current Protocols in Human Genetics.

557   2015;87:11.16.1-11.16.14.

558   41. Sonawane AR, Platig J, Fagny M, Chen C-Y, Paulson JN, Lopes-Ramos CM, et al.

559   Understanding Tissue-Specific Gene Regulation. Cell Reports. 2017;21:1077–88.

560  42. Vaquerizas JM, Kummerfeld SK, Teichmann SA, Luscombe NM. A census of human

561  transcription factors: function, expression and evolution. Nat Rev Genet. 2009;10:252–63.

562  43. Siam A, Baker M, Amit L, Regev G, Rabner A, Najar RA, et al. Regulation of alternative

563  splicing by p300-mediated acetylation of splicing factors. RNA. 2019;25:813–24.

564  44. Dušková E, Hnilicová J, Staněk D. CRE promoter sites modulate alternative splicing via

565  p300-mediated histone acetylation. RNA Biology. 2014;11:865–74.

566  45. Szklarczyk D, Gable AL, Nastou KC, Lyon D, Kirsch R, Pyysalo S, et al. The STRING

567  database in 2021: customizable protein–protein networks, and functional characterization of

568  user-uploaded gene/measurement sets. Nucleic Acids Research. 2021;49:D605–12.

569  46. de Assuncao TM, Lomberk G, Cao S, Yaqoob U, Mathison A, Simonetto DA, et al. New role

570  for Kruppel-like factor 14 as a transcriptional activator involved in the generation of signaling

571  lipids. J Biol Chem. 2014;289:15798–809.

572  47. Teplyakov E, Wu Q, Liu J, Pugacheva EM, Loukinov D, Boukaba A, et al. The

573  downregulation of putative anticancer target BORIS/CTCFL in an addicted myeloid cancer cell

574  line modulates the expression of multiple protein coding and ncRNA genes. Oncotarget.

575  2017;8:73448–68.

576  48. Pugacheva EM, Teplyakov E, Wu Q, Li J, Chen C, Meng C, et al. The cancer-associated

577  CTCFL/BORIS protein targets multiple classes of genomic repeats, with a distinct binding and

578  functional preference for humanoid-specific SVA transposable elements. Epigenetics Chromatin.

579  2016;9:35.

580  49. Singh S, Narayanan SP, Biswas K, Gupta A, Ahuja N, Yadav S, et al. Intragenic DNA

581  methylation and BORIS-mediated cancer-specific splicing contribute to the Warburg effect. Proc

582  Natl Acad Sci U S A. 2017;114:11440–5.

583  50. Jaganathan K, Panagiotopoulou SK, McRae JF, Darbandi SF, Knowles D, Li YI, et al.

584  Predicting Splicing from Primary Sequence with Deep Learning. Cell. 2019;176:535-548.e24.

585    51. Xi W, Beer MA. Local epigenomic state cannot discriminate interacting and non-interacting

586    enhancer–promoter pairs with high accuracy. PLOS Computational Biology. 2018;14:e1006625.

587    52. Adamson SI, Zhan L, Graveley BR. Vex-seq: high-throughput identification of the impact of

588    genetic variation on pre-mRNA splicing efficiency. Genome Biol. 2018;19:71.

589    53. LeCun Y, Bengio Y, Hinton G. Deep learning. Nature. 2015;521:436–44.

590    54. Zou J, Huss M, Abid A, Mohammadi P, Torkamani A, Telenti A. A primer on deep learning

591    in genomics. Nat Genet. 2019;51:12–8.

592

593

# Funding

598

599

# Authors' information

601    Affiliations

602    **Institute of Information Science, Academia Sinica, Taipei, Taiwan**

603    Tzu-Chieh Lin, Cheng-Hung Tsai, Cheng-Kai Shiau, Jia-Hsin Huang, Huai-Kuang Tsai

604

605    **Taiwan AI Labs & Foundation, Taipei, Taiwan**

606    Jia-Hsin Huang, Huai-Kuang Tsai

607

608    **Corresponding authors**

609    Correspondence to Jia-Hsin Huang (jiahsin.huang@gmail.com) or Huai-Kuang Tsai

610    (hktsai@iis.sinica.edu.tw).

611

612    **Authors' contributions**

613     TCL and CHT conceived the study and carried out the bioinformatics pipelines. TCL prepared

614     the initial draft of the manuscript. CKS assisted the data curation. JHH and HKT conceived and

615     designed the research, interpreted the results, and drafted the manuscript. All authors read and

616     approved the final manuscript.

617

618

## Ethics declarations

**Ethics approval and consent to participate**

Not applicable.

**Consent for publication**

Not applicable.

**Competing interests**

The authors declare that they have no competing interests.

## Figure and Figure legends

631

**Figure 1.** (A) The workflow schema and the experiment design. We obtain 15 tissues that have matched DNase-seq and RNA-seq from ENCODE. DNase-seq data was used to identify open chromatin regions and followed by TF motif scanning to identify TF binding profile in promoter. RNA-seq data was processed by the MISO program to obtain percent splice in (PSI) metrics which represent the splicing pattern of the first skipped exon. (B) PSI distribution histogram. The horizontal axis represents the PSI value and the vertical axis represents the number of skipping exon events. (C) Venn diagram of the exclusion group gene and inclusion group gene. The exclusion group gene defined as $PSI < 0.2$ and the inclusion group gene defined as $PSI > 0.8$. (D) Volcano plot of the chi-square test results and the TF expression tissue specificity distribution along with ranking $p$-values of the chi-squared test. The horizontal axis of the volcano plot represents the $-\log_{10}$ (adjust $p$-value) and the $\log_2$ (OR). The chi-square test $p$-value is corrected by Bonferroni multiple test correction. The blue dot denoted the ubiquitously expressed TFs (tau $< 0.8$) and the red dot denoted the tissue-specific expressed TFs (tau $\geq 0.8$). (E) The schema of validation strategies. From left to right represents event-wise, tissue-wise, and gene-wise validation schema, respectively. (F) The model performance of event-wise, tissue-wise, and gene-wise validation schema. For left panel to right panel represents F1-score, AUROC, and accuracy, respectively. (G) The gene were assigned into three groups according to the splicing forms across all tissues. One-sided denotes the genes belonging to same splicing form in more than two tissues; both-sided denotes the genes having both inclusion and exclusion forms in 15 tissues; singleton denotes the genes expressed in a particular tissue only. The accuracies of prediction and number of genes in three groups were calculated respectively for each tissue from the tissue-wise validation experiments.
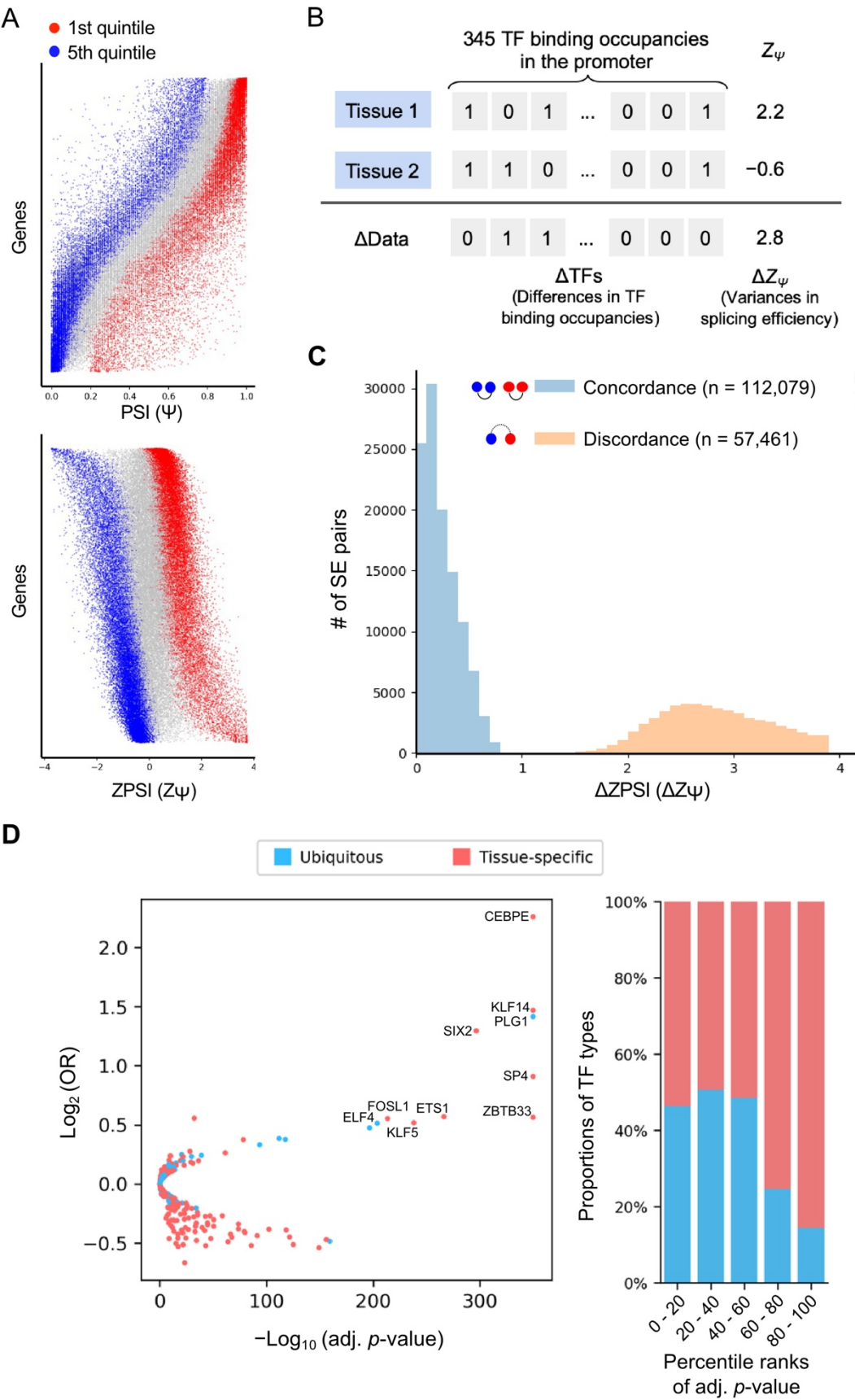
24

654

655  **Figure 2.** (A) The distribution of PSI and ZPSI. Each horizontal line represents PSI values of a

656  gene and the vertical axis was sorted by gene median PSI. The blue dots denote the first quintile

657  (top 20%) of PSI and the red dots denote the fifth quintile (latest 20%) of PSI. (B) The "delta"

658  schema of splicing events. For each gene, we enumerate all tissue pairs and perform exclusive-or

659  (XOR) operation on the TF binding occupancies and yield $\Delta$Data representation which means the

660  differences in TF binding occupancies. For the splicing pattern, we calculate the absolute

661  difference of the ZPSI and yield $\Delta Z_\Psi$, which represents the variances in splicing efficiency. (C)

662  The distribution of $\Delta Z_\Psi$ among splicing status unchanged group (concordance) and changed

663  group (discordance). The distribution showed a clear bimodal pattern, that the discordance $\Delta Z_\Psi$

664  is distinctly higher than the concordance $\Delta Z_\Psi$. (D) The chi-squared test of association between

665  TFBS-occupied differences and splicing phases. The left panel is the volcano plot of the chi-

666  square test; the horizontal axis represents the $-\log_{10}$ (adjusted $p$-value) and the vertical axis

667  represents the $\log_2$ (OR). Top 10 significant TFs are shown in their names. The right panel is the

668  ratio of tissue-specific and ubiquitous TFs among adjusted $p$-value rankings.
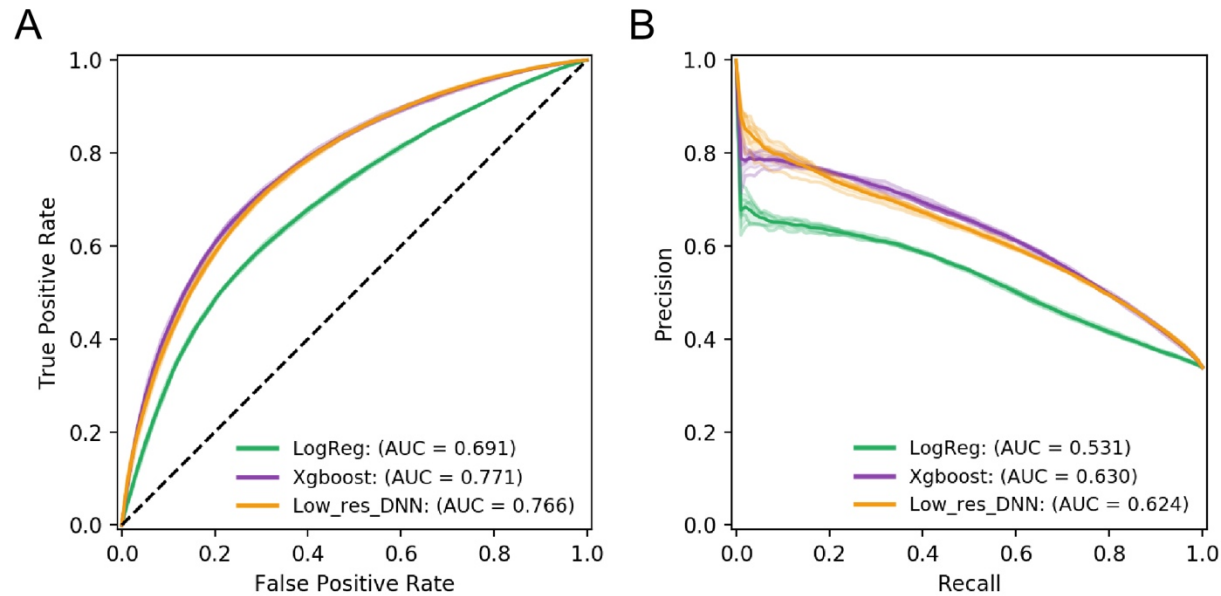
669

26

**Figure 3.** (A) The area under receiver operating characteristic curve (AUROC) of Logistic regression, XGBoost, and low-resolution deep learning model. The input of the low-resolution deep learning model only contains a single array of TF occupancy information denote as low-resolution. Of note, the XGBoost model has the highest AUROC. (B) The area under precision-recall curve (AUPRC) of Logistic regression, XGBoost, and low-resolution deep learning model. With the same trend of AUROC, the XGBoost model has the highest AUPRC.
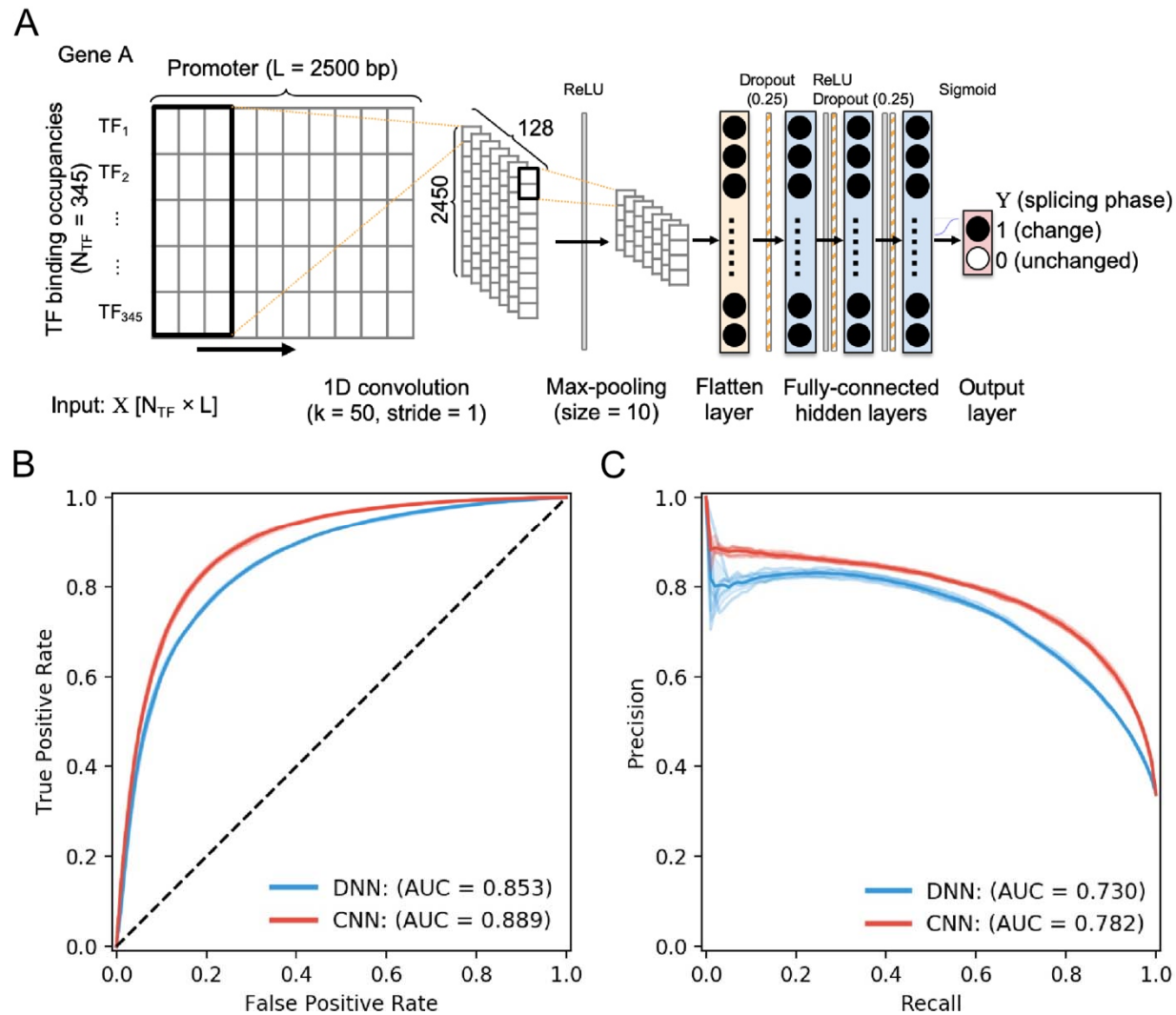
**Figure 4.** (A) The convolutional neural network schema. The first layer is a convolution layer with ReLU activation function and followed by a max-pooling layer. After pooling a flatten layer was applied to reshape the input. Then three dense layer is added followed by a sigmoid function to classified the output. (B) The area under receiver operating characteristic curve (AUROC) of convolutional neural network (CNN) and deep neural network (DNN). (C) The area under precision-recall curve (AUPR) of CNN and DNN. Both AUROC and AUPR suggest the CNN has the better performance.
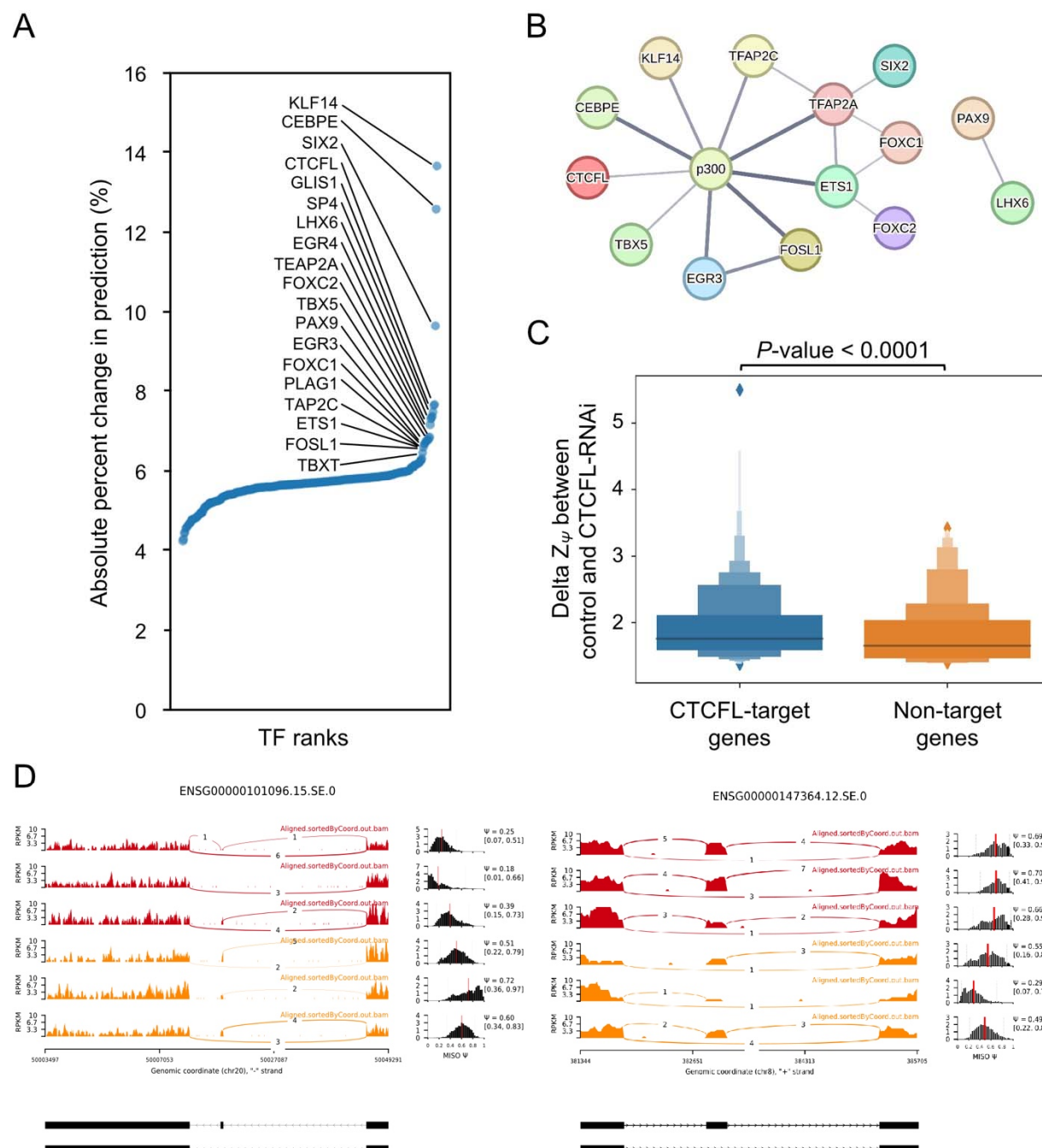
**Figure 5.** (A) The rank order plot of importance analysis. The horizontal axis represents the TF importance ranks. The vertical axis represents the importance measures (see importance analysis in method section). (B) The gene association network was constructed from the STRING database for top important TFs with p300. The thickness of edges denotes the strength of data support according to textmining, experiments, and databases. (C) The distribution of $\Delta Z_{\Psi}$ between control and CTCFL-RNAi experiment. The $\Delta Z_{\Psi}$ values of the CTCFL-target genes

694    show significant differences than that of non-CTCFL target genes with Wilcoxon rank-sum test

695    ($p$-value < 0.0001). (D) The sashimi plot and PSI distribution across control and CTCFL-RNAi

696    experiment. The left panel shows the first skipped exon event of ENSG00000101096. The right

697    panel shows the first skipped exon event of ENSG00000147364. Red samples were from the

698    control of CTCFL experiments and orange samples were from the CTCFL-RNAi samples.