1

2

3

# Epigenome-based Splicing Prediction using a Recurrent Neural Network

5

6

7

8  Donghoon Lee[1,2], Jing Zhang[1,2], Jason Liu[2], and Mark B Gerstein[1,2,3,4*]

9

10

11

12

13  1 Program in Computational Biology and Bioinformatics, Yale University, New Haven, CT

14  06520, USA

15  2 Department of Molecular Biophysics and Biochemistry, Yale University, New Haven, CT

16  06520, USA

17  3 Department of Computer Science, Yale University, New Haven, CT 06520, USA

18  4 Department of Statistics and Data Science, Yale University, New Haven, CT 06520, USA

19  * Corresponding author

20  E-mail: pi@gersteinlab.org

21

22

23    **Abstract**

24    Alternative RNA splicing provides an important means to expand metazoan transcriptome

25    diversity. Contrary to what was accepted previously, splicing is now thought to predominantly

26    take place during transcription. Motivated by emerging data showing the physical proximity of

27    the spliceosome to Pol II, we surveyed the effect of epigenetic context on co-transcriptional

28    splicing. In particular, we observed that splicing factors were not necessarily enriched at exon

29    junctions and that most epigenetic signatures had a distinctly asymmetric profile around known

30    splice sites. Given this, we tried to build an interpretable model that mimics the physical layout

31    of splicing regulation where the chromatin context progressively changes as the Pol II moves

32    along the guide DNA. We used a recurrent-neural-network architecture to predict the inclusion

33    of a spliced exon based on adjacent epigenetic signals, and we showed that distinct spatio-

34    temporal features of these signals were key determinants of model outcome, in addition to the

35    actual nucleotide sequence of the guide DNA strand. After the model had been trained and tested

36    (with >80% precision-recall curve metric), we explored the derived weights of the latent factors,

37    finding they highlight the importance of the asymmetric time-direction of chromatin context

38    during transcription.

39

40    **Author Summary**

41    In humans, only about 2% of the genome is comprised of so-called coding regions and can give

42    rise to protein products. However, the human transcriptome is much more diverse than the

43    number of genes found in these coding regions. Each gene can give rise to multiple transcripts

44    through a process during transcription called alternative splicing. There is a limited

45    understanding of the regulation of splicing and the underlying splicing code that determines cell-

46    type-specific splicing. Here, we studied epigenetic features that characterize splicing regulation

47    in humans using a recurrent neural network model. Unlike feedforward neural networks, this

48    method contains an internal memory state that learns from spatiotemporal patterns – like the

49    context in language – from a sequence of genomic and epigenetic information, making it better

50    suited for characterizing splicing. We demonstrated that our method improves the prediction of

51    spicing outcomes compared to previous methods. Furthermore, we applied our method to 49 cell

52    types in ENCODE to investigate splicing regulation and found that not only spatial but also

53    temporal epigenomic context can influence splicing regulation during transcription.

54

55    **Introduction**

56    Alternative splicing of pre-messenger RNA plays an integral role in diversifying the

57    transcriptome. This process is more pervasive in higher eukaryotes and is estimated to affect

58    approximately 95% of protein-coding genes in humans [1,2]. Accurate characterization of the

59    process by which multiple functional protein products are produced from a single gene is crucial

60    for understanding the function of the transcriptome [3].

61    Recent discoveries have revealed that splicing occurs predominantly during transcription in

62    humans [4–8]. Nascent RNA is almost immediately spliced upon transcription [9,10] and introns

63    are mostly spliced out during transcript elongation. This timing suggests that the recruitment of

64    splicing factors and spliceosome assembly, detection of exon-intron boundaries, and modulation

65    of alternative splicing must occur at the same time scale as transcription [9].

66    Co-transcriptional splicing indicates a key observation that splicing takes place progressively in

67    the direction of RNA transcription, rather than processed simultaneously after transcription. As a

68    result, the contexts of guide DNA, nascent RNA, and its immediate folded structure

69    progressively change as RNA polymerase II (Pol II) moves along the guide DNA strand [11] and

70    may influence splicing regulation. Furthermore, co-transcriptional splicing signifies the physical

71    proximity of the spliceosome assembly to Pol II and other transcriptional machinery [9]. Pol II

72    physically interacts with nucleosomes and its histone modifications around them, modulating the

73    transcription rate [12].

74    DNA sequence alone may not contain sufficient information to process alternative splicing

75    deterministically [13]. Djebali et al. [4] and many others have shown that there is an enrichment

76    of chromatin marks around spliced exons, suggesting the role of epigenetic modifications during

77    context-dependent modulation of alternative splicing [14,15]. For example, exonic boundaries

78    are characterized by increased levels of nucleosome density and positioning [16–18], DNA

79    methylation [19,20], and strong enrichment of specific histone modifications including

80    H3K36me3, H3K79me1, H2BK5me1, H3K27me1, H3K27me2, and H3K27me3 [16,17,21–23].

81    In addition, a recent genome-wide survey of alternative splicing showed that DNA methylation

82    can either enhance or silence exon recognition in a context-dependent manner [24]. Furthermore,

83    studies have shown that there is significant regulatory crosstalk between histone modifications

84    during transcriptional elongation [12].

85    Despite many efforts to characterize the splicing regulatory code both experimentally and

86    computationally, we have yet to understand how the cell type-specific epigenomic context is

87    utilized during co-transcriptional splicing. Previous computational methods on splicing have

88    largely focused on discovering novel splice junctions based on RNA sequencing (RNA-seq)

89    alignments [25,26], utilizing machine learning approaches [27,28] including deep neural

90    networks [29]. Only a limited set of tools can model splicing regulation based on genomic

91    sequences and select RNA features [30–32]. Moreover, studies on splicing regulation have

92    focused heavily on identifying mutations that land within splice sites (SSs), cis-acting splicing

93    regulatory elements, and trans-acting splicing factors [30,33]. The extent, nature, and effects of

94    the epigenetic context in splicing regulation remain unsolved.

95    In this study, we propose a new computational approach to characterize the role of epigenetic

96    modifications during co-transcriptional splicing. To build an interpretable model, we adopted a

97    recurrent neural network (RNN) architecture, which to some degree resembles the physical

98    characteristics of co-transcriptional splicing (Figure 1). The model can learn from a temporal

99    sequence of epigenetic contexts, similar to how epigenetic contexts progressively change as Pol

100   II moves forward along the guide DNA strand during co-transcriptional splicing. The RNN

101   model allows us to predict the inclusion of exons based on adjacent DNA sequences and

102   epigenetic modifications. Moreover, the physical resemblance of the model allows us to interpret

103   the trained model weight parameters and explore the spatio-temporal links between the guide

104   DNA elements and the surrounding epigenetic modifications. In summary, we leveraged the

105   mechanistic properties of co-transcriptional splicing to build an interpretable splicing model, and

106   we explored the trained model to understand the underlying characteristics of the epigenetic

107   context during co-transcriptional splicing.

108

109   **Results**

110   We first explore the epigenetic data context around known splice sites in depth. We then describe

111   the model and rationale for applying the specific architecture. Finally, we use the model to

112   further examine the effect of epigenetic context during co-transcriptional splicing.

113

114   **Distinct epigenomic signatures characterize splicing regulation**

115   We studied the epigenetic context of alternative splicing by examining the enrichment of

116   multiple histone modifications and DNA methylations around the exon-intron boundary. We

117   mapped the epigenomic signatures around SSs of cassette exons at a base-pair resolution. We

118   aggregated multiple histone modifications across 49 cell types in ENCODE and observed their

119   enrichment as a function of distance from SSs (Figure 2A, B, Supplementary Figure 1, 2A, B).

120   We found the most interesting trend within 100 bp of SSs for both the 3' acceptor and 5' donor.

121   A strong enrichment pattern of H3K36me3 and H3K27me3 appeared around the exon boundary.

122   Although studies have demonstrated a role for H3K36me3 in defining the exon-intron boundary

123   [22,34], the dynamic interplay between other histone modifications has been overlooked. From

124   the 3' acceptor, peak enrichment occurred around 100 bp into the exon; at the 5' donor, it was

125   closer, at around 50 bp into the exon. We also observed a slight depletion of H3K27ac and

126   H3K4me3 marks within 100 bp of the intron at the 3' acceptor SS but not within the 5' donor SS.

127   As this region contains a branch site, these histone marks may indicate a role in defining the

128   branch point.

129

130   **Enrichment of RNA-binding factors around splice sites**

131   Alternative splicing regulation is an elaborate process that requires precise coordination of

132   multiple splicing factors and enzymes. Studies have shown that RNA-binding proteins (RBPs)

133   facilitate splicing regulation during transcription [35]. For example, the serine/arginine-rich

134   splicing factor family member SRSF7 binds to poised exons and promotes the inclusion rate

135   [36][37]. Another member of the serine/arginine-rich splicing factor family, U2AF1, is

136   responsible for mediating the binding of U2 small nuclear ribonucleoprotein to the pre-mRNA

137   branch site [38]. The recent release of the ENCODE project included enhanced CLIP

138    experiments (eCLIP) datasets that span 112 RBPs from K562 and HepG2 cell types. As

139    sequence-specific RBPs have been shown to facilitate splicing regulation in a context-specific

140    manner [15], we investigated their spatial relationship to both the 5' donor and 3' acceptor

141    splicing sites. Specifically, we investigated the enrichment of splicing factors (n=29) and their

142    relative distance to these sites. We observed that, on average, splicing factors show preferential

143    binding to the intronic side of the splicing site in both 3' acceptor and 5' donor SSs

144    (Supplementary Figure 2C). Furthermore, we found that splicing factors may show slightly

145    different patterns in their spatial binding preferences. In particular, hnRNP A1 and SRSF1 were

146    enriched in the intronic region outside 3' SSs whereas SF3B4 and hnRNP C were enriched in the

147    exonic region (Figure 2C). At 5' SSs, RBM22 and PRPF8 were bound at the exonic end, which

148    has been shown to be critical for splicesome assembly [39,40].

149

150    **Correlating epigenomic signatures to exonic expression**

151    We tested whether histone modifications have any effect on inclusion and expression of

152    alternative exons. We observed a trend where enrichment of H3K36me3 at the exon-intron

153    boundary was positively correlated with exonic expression, whereas H3K27me3 marks showed

154    the opposite trend (Figure 3A, B, Supplementary Figure 3). Compared to excluded or nominally

155    expressed alternative exons, highly expressed spliced exons had statistically significant

156    enrichment of H3K36me3 and depletion of H3K27me3 at their exon-intron boundary (Figure

157    3C). The contrasting trend and the correlation of these histone methylations to exonic expression

158    suggest that the splicing code may be directly or indirectly encoded within the epigenomic

159    context.

160

**Clustering biosamples based on splicing patterns**

Previous studies have shown that various epigenomic marks are correlated across similar tissues and cell types [41]. It is now widely accepted that the transcriptional regulatory circuitry of a particular cell type is reflected in its epigenetic landscape. To explore the potential linkage between epigenetic regulation and tissue-specific splicing, we examined splicing patterns across 49 ENCODE biosamples. Based on a similarity of percent-splice-in (PSI) values for all coding exons (n=185,405), we clustered biosamples into five categories using hierarchical clustering (Figure 3D). Splicing patterns were highly correlated among tissue types from the same cell-of-origin, reproducing similar clustering results based on epigenetic marks. For example, blood-lineage cell types formed cluster C2 whereas brain and neural cells were clustered in cluster C4.

In addition to using the PSI similarity matrix to cluster cell types into categories, we can project the cells onto a low-dimensional cell space using principal component analysis (PCA). We measured alternative splicing patterns in terms of exonic expression level (fragment per kilobase per million reads mapped, FPKM) across diverse ENCODE cell types and examined how cells are placed in the context of others. Interestingly, we observed that cancer-related cell lines were located proximal to each other in the PCA cell space (Supplementary Figure 4).

**Modeling splicing regulation: key characteristics of an RNN architecture**

To investigate the latent representation of splicing instruction encoded within the epigenomic context, we aimed to construct a predictive model of splicing. We opted for an RNN architecture, which has proven successful in various sequential information processing and prediction tasks

183    such as natural language processing and translation [42–44], to explore the contribution of the

184    epigenomic context to the regulation of alternative splicing.

185    We start by describing a simple RNN, which shares many of the features we intend to model. A

186    simple RNN is made of many recurrent neurons that are sequentially linked to each other. A

187    neuron at specific time point $t$ is influenced by previous time point $t-1$, combining some

188    relationship of the current input $x_t$ with the previous hidden state $h_{t-1}$.

189

$$h_t = f(h_{t-1}, x_t)$$

190

191    where $h_t$ is hidden state at time $t$ and $x_t$ is input variable at time $t$. If we suppose the activation

192    function as a hyperbolic tangent for a simple RNN, the state at time $t$ can be represented as

193

$$h_t = \tanh(W_h{}^T h_{t-1} + W_x{}^T x_t + \mathrm{b})$$

194

195    where $W_h$ and $W_x$ are the weight of the hidden state and input variable, respectively, and b is the

196    bias vector. The output can be expressed in terms of an output weight matrix, $W_y$, and a hidden

197    state at time $t$, $h_t$:

198

$$\hat{y}_t = \mathrm{S}(W_y{}^T h_t)$$

199

200    where $S$ is sigmoid function:

201

$$S(x) = \frac{e^x}{e^x + 1}$$

202

203  This time-dependency allows us to explore the complex contextual relationship between features.

204  In particular, we adopted the long short-term memory (LSTM) [45] model to describe an RNN

205  architecture. In principal, a simple RNN allows us to model a time-dependent task from

206  sequential data. However, in practice, the simple model suffers from the problem of vanishing

207  gradients, where the gradients responsible for updating weights with respect to the partial

208  derivative of error function becomes negligible in a long sequence and hampers the model from

209  learning long-term time dependencies. Therefore, we used both LSTM and gated recurrent unit

210  (GRU), which have many of the same simple intuitive properties of the simple RNN but allow

211  learning from longer sequences. The LSTM is an extension of the same idea that includes more

212  sophisticated gates, which allows the cell to retain long-term memory between cells while

213  avoiding the problem of vanishing gradients when training the network. The specific equations

214  for the LSTM model we adopted is shown in the Methods.

215

216  **Modeling splicing regulation: How the RNN architecture fits the problem**

217  The rationale for applying an RNN to our model is that (1) an RNN is optimized for processing

218  sequential information like genomic sequences and epigenomic profiles along genomic

219  coordinates, (2) an RNN has a time-direction resembling how RNA is transcribed by RNA

220  polymerase in the 5' to 3' direction, (3) temporal memory cells of an RNN allow the model to

221  learn about complex context-dependent relationships among epigenomic features, such as the

222  influence of features and input seen at $t$-1 on the neural cell at time $t$, and (4) an RNN is very

223  flexible with the type of input and output data and therefore can easily integrate heterogeneous

224  sequential information. Not surprisingly, researchers recently have applied RNN models to the

225    area of genomics to predict non-coding DNA function [46] and to detect exon junctions [47].

226    Moreover, since the mechanics of the RNN calculation is somewhat parallel to the actual spatial

227    and temporal dependency found in co-transcriptional splicing, the overall results from the trained

228    model are more readily interpretable. The data processing and implementation of the predictive

229    models are collected in a package named Epigenome-based Splicing Prediction using Recurrent

230    Neural Network (ESPRNN; available at https://github.com/gersteinlab/esprnn). Using our

231    method, we attempted to decipher context-dependent effects of various epigenomic features on

232    splicing for both canonical (e.g., dinucleotide GT for 5' donors and AG for 3' acceptors) and

233    non-canonical SSs. Our model is especially useful since splicing signals are not only enriched at

234    the splice site but often found up and downstream of splice sites.

235

236    **Modeling splicing regulation: Initial evaluation**

237    We used ESPRNN to predict alternate usages of cassette exons (inclusion or exclusion of exons),

238    the most common form of alternative splicing events [48], using DNA sequences and

239    epigenomic signals adjacent to SSs (Figure 4A). We used the exon definition of splicing, which

240    is considered to be the dominant mechanism in higher eukaryotes [49]. Our model had an

241    average F1 score (harmonic mean of the precision and recall) of 0.8472 for the LSTM-based

242    model across cell types [0.8757 for the GRU-based model] using five core histone modification

243    tracks (Figure 4B). The average F1 score marginally increased to 0.8573 when using 17 histone,

244    chromatin accessibility, DNA methylation, and nucleosome density profiles.

245    We performed the splicing prediction with or without the RBP profile and measured how much

246    predictive performance is gained from additional information. We observed a marginal

247    improvement in predictive performance when RBP binding profiles were added to the baseline

248   model (measured in improvement of F1 score from 0.84 to 0.86) (Supplementary Figure 9A, B).

249   This suggests RBP binding information may be redundant and already represented in the

250   epigenetic features. We also compared prediction results from normal cell types to those from

251   cancerous cell lines. Since previous studies on cancer-specific alternative splicing [50,51] have

252   suggested potential linkage of aberrant splicing events to the disease risk [52–55], we expected

253   to see differences in splicing regulation between normal and cancerous cell types. However, we

254   did not observe a significant difference in prediction performance between normal and cancerous

255   cell types (average F1 score for normal biosamples: 0.8465, cancerous biosamples: 0.8765). We

256   also cross-tested a model trained from one cell type to another. After we fit our model to one cell

257   type, we transferred the fitted weights and model parameters to predict splicing on other cell

258   types. When we tested between cell types from the same cell-of-origin (e.g., train on adult liver

259   model and test on HepG2 data, train on lung model and test on A549 data), we did not observe a

260   significant difference in predictive performance. However, we observed a moderate reduction in

261   splicing prediction performance when we cross-tested cells from different cell-of-origin

262   (Supplementary Figure 5B, F1 score is better metric for comparing cross-cell testing due to class

263   imbalance across cell types). Thus, the epigenomic regulatory landscape around SSs appears to

264   be generally conserved across cell types. Moreover, we compared the classification performance

265   to other models based on random forest and k-nearest neighbors and found that our model was

266   superior in terms of classification accuracy (Figure 4D, Supplementary Figure 7).

267   We tried to measure the contribution of each individual epigenetic feature to splicing in a

268   number of ways. (1) We performed an empirical analysis via a leave-one-out strategy. Using

269   GM12878 as an example, we first built a reference model based on all available epigenetic

270   features. By removing one variable at a time, we then measured the mean decrease in F1 score

271   and area under the receiver operating characteristic curve (ROC AUC), as an indicator of

272   variable importance (Figure 4C). (2) Alternatively, we trained a DNA-only model using DNA

273   sequence features only and compared to a "baseline model." The baseline model was trained

274   using DNA sequence features plus additional chromatin accessibility (DHS) and 6 histone marks.

275   Here, we observed a significant loss of predictive performance in the DNA-only model (13%

276   reduction in F1 score) (Supplementary Figure 6A). (3) Next, starting from the DNA-only model,

277   we added one epigenetic feature at a time to measure the information gain from each feature

278   (Supplementary Figure 6B). While the addition of some epigenetic features like H3K27ac

279   increased the variability in prediction performance, an active mark H3K36me3 or a repressive

280   mark H3K27me3 was the most informative at predicting splicing. Moreover, the combination of

281   both H3K36me3 and H3K27ac further improved the prediction performance compared to other

282   pairs (Supplementary Figure 6C). We observed that the combination of H3K36me3 and

283   H3K27ac features together contributed more than when they were used individually

284   (Supplementary Figure 6D).

285   Overall, we found H3K36me3 to be the most important variable in predicting splicing. This

286   observation coincides with previous studies reporting that H3K36me3 recruits the splicing

287   factors PTB [34] and SRSF1 [56] to facilitate splicing. Interestingly, one of the top predictors of

288   splicing was H3K79me2, which was previously shown to associate with H3K36me3 at gene

289   bodies [57]. H3K9me3, a histone modification that can recruit adaptor proteins like HP1 to

290   facilitate splicing factors [24], was also ranked among the top predictors.

291

292   **Interpretation of weights of the splicing model**

293    Since the model follows the physical layout of splicing regulation, one can examine the trained

294    model and learn from the trained weights how each epigenetic feature contributes to splicing

295    regulation. To interpret the splicing model, we designed an LSTM-based model composed of

296    only one hidden state and trained for a longer period (400 epochs). We made sure that this

297    simplified model performs nearly as well at predicting splicing as our main model (usually after

298    >20 epochs of training, Supplementary Figure 8A). We also made sure that the overall predictive

299    performance of the simplified model is stable after approximately 100 epochs (Supplementary

300    Figure 8B, C). When we analyzed the simplified model, we found that the trained weights of

301    various gates at the recurrent unit showed that open chromatin (DHS), H3K27ac, K3K36me3,

302    and H3K4me1 are weighted more highly than other epigenetic features -- as expected

303    (Supplementary Figure 8D). We also noticed that H3K27me3 and K3K9me3 were negatively

304    weighted at the input gate, suggesting that these features have a negative impact on exon

305    inclusion, consistent with our previous findings.

306

307    **Influence of temporal epigenetic context on splicing regulation**

308    We specifically designed our splicing model to represent the physical layout of splicing

309    regulation, where a sequence of chromatin contexts is fed progressively to the model. Therefore,

310    the model takes into account the temporal direction (progression from 5' to 3' in direction). To

311    show that model has learned this asymmetric temporal relationship of epigenetic features, we

312    first trained a baseline model (in the normal 5' to 3' direction) and then fed a series of epigenetic

313    signals in a "reverse" order (3' to 5' in direction) as input to it. We observed how the model

314    prediction behaved in this context. If the model was agnostic to the temporal direction of features,

315    both forward and reverse input features should give the same predictive power. In fact, we saw a

316    moderate decrease in prediction performance (Supplementary Figure 9A, B), with an F1 score

317    decreasing from 0.871 to 0.865 and ROC AUC decreasing from 0.886 to 0.865.

318

319    **Discussion**

320    Our prediction model revealed that the epigenomic signature of an SS plays a large role in

321    determining the splicing outcome. In addition, the positive results suggest that our model can be

322    extended to predict the full transcriptomic composition from a genomic and epigenomic context.

323    We expect that we could further improve the proposed model by adding more deep hidden layers

324    and increasing the number of training samples by utilizing the full set of available epigenomic

325    data in the ENCODE project. Our approach does contain some limitations, as it is still

326    challenging to visualize and evaluate the multi-dimensional context of the weight matrix in the

327    trained model. We could apply dimensionality reduction techniques to probe the latent

328    representation of relationships between various epigenomic signals.

329    In this study, we used ENCODE polyA RNA-seq assays to measure splicing and exon-level

330    expression; we note that this is an indirect measure of what is actually happening during

331    transcription. RNAs are often unstable and may be subjected to many post-transcriptional

332    modifications. RNA-seq measures the steady-state level of the transcript, accounting for both

333    mRNA synthesis and decay. Future studies with a more direct measure of transcriptional rates,

334    such as nuclear run-on assays like global run-on (GRO-seq) or bromouridine sequencing (Bru-

335    seq), will allow us to accurately measure the effect of epigenomic context on splicing and,

336    ultimately, on the transcriptional rate.

337    Future studies should focus on comparing splicing models from normal and cancer samples in

338    the hope of illuminating the differences in the epigenomic landscapes of splicing regulation.

339     Although splicing is an elaborate process, it could become pathogenic when misregulated [58,59].

340     Unsurprisingly, aberrant splicing events, which collectively referred to splicing events that could

341     confer the risk of a disease, are often implicated in systemic diseases like cancer [51,60].

342     Aberrant splicing events based on mutations are relatively well characterized [54,60–62];

343     however, a large fraction of aberrant splicing events that have no direct mutational cause still

344     remain unknown. Although our understanding of epigenomic context on splicing regulation is

345     incomplete, our prediction model highlights that splicing is elaborately regulated via various

346     epigenomic signatures. This suggests that epigenomic dysregulation may be closely linked to the

347     onset of aberrant splicing. Thus, even though aberrantly spliced RNAs in healthy cells may be

348     degraded by the mRNA surveillance system, epigenomic dysregulation may render this

349     checkpoint system useless. Further studies on cell-type-specific and context-dependent splicing

350     regulation will reveal whether epigenetic modulation can serve as a therapeutic method of

351     complex disease in the future.

352

353     **Methods**

354     **Dataset**

355     The current release of the ENCODE dataset provides an unprecedented number of functional

356     assays across broad biosample types, including primary cells and tissues. In this study, we

357     leveraged both the breadth and depth of ENCODE, including assays for histone modification

358     (chromatin immunoprecipitation sequencing, ChIP-seq), chromatin accessibility (DNase I

359     hypersensitive sites sequencing, DNase-seq), RBPs (eCLIP), methylations (WGBS and RRBS)

360     and gene expression (RNA-seq), to systematically probe the data-rich context of alternative

361    splicing and its regulation. The list of accessions for experiments used in this study is found in

362    Supplementary Table 1.

363

364    **Processing of RNA-seq data**

365    To quantify levels of exon expression from RNA-seq data, we collected all raw sequencing reads

366    from experiments tagged as reference epigenome series from the ENCODE portal. These reads

367    were polyA plus long RNA-seq (200 bp or larger) from whole-cell fractions rather than nuclear

368    or cytosolic fractions. To minimize potential batch effects and sample bias, we carefully selected

369    untreated experiments from the reference epigenome series. As of November 2019, there are 81

370    cell and tissue types (covering 49 unique biosamples) in the reference epigenome series,

371    including both RNA-seq and ChIP-seq of H3K4me1, H3K4me3, H3K36me3, H3K27ac,

372    H3K27me3, and H3K9me3. We first aligned all RNA-seq data to the GRCh38 genome using

373    RNA STAR (v 2.7.0). Since the model requires splice site annotation, we constructed exon

374    annotation from GENCODE version 24 (to synchronize with ENCODE annotation) by extracting

375    all unique exons with known protein-coding transcripts. We excluded exons that could

376    ambiguously map to both chromosome X and Y. This analysis included 597,937 exons (185,405

377    unique exons after removing duplicates from isoforms) that averaged 28.01 exons per gene and

378    296.49 bp in length (150.92 bp in length for unique exons). We obtained read counts at each

379    exon using HTSeq (v0.11.2) [63]. Based on read counts, we used a custom script

380    (esprnn/preproc_calcExonFPKM.py) to calculate normalized exonic expression levels in FPKM.

381    Our rationale for using the exonic expression was to intentionally make the model agnostic to the

382    overall transcript level. Each exon was evaluated independently from other exons, and we

383    counted the number of sequencing reads supporting the inclusion of a particular exon. The

384 counts were normalized similar to how a gene's expression is normalized by size of annotation

385 and total number of mapped reads (FPKM). We binarized the exonic expression level (FPKM)

386 using a threshold of one. Therefore, we only considered whether an exon has enough evidence

387 supporting exon inclusion.

388

389 In addition to the exonic expression level, alternatively, we calculated a metric, PSI, to measure

390 the level of splicing. PSI represents the fraction of the reads supporting exon inclusion from the

391 split reads at the splice junction. We used a custom script (esprnn/scripts/calcPSI.sh) based on

392 equations from Schafer et al. [64] to calculate PSI normalized by the size of read and exon

393 annotation.

394

$$\widetilde{F}_i^{incl} = \frac{F_i^{incl}}{L_i + L_f}$$

$$\widetilde{F}_i^{excl} = \frac{F_i^{excl}}{L_f}$$

$$PSI\ (\Psi) = \frac{\widetilde{F}_i^{incl}}{\widetilde{F}_i^{incl} + \widetilde{F}_i^{excl}}\%$$

395

396 $F_i^{incl}$ number of reads or fragments supporting the inclusion of $i$-th exon; $F_i^{excl}$ number of reads

397 or fragments supporting the exclusion of $i$-th exon; $L_f$ fragment length; $L_i$ size of $i$-th exon.

398

399 **RNA-binding proteins**

400 RBP enrichment was calculated based on the peaks identified from the eCLIP experiments. We

401 downloaded the ENCODE eCLIP uniformly processed peaks from K562 and HepG2 cell types

402   (see Supplementary Table 1 for eCLIP data accession). The peak was called using CLIPPER

403   software [65] and filtered for peaks having a score of 1,000. We then counted numbers of RBP

404   binding events at a base-pair resolution, agnostic to cell type.

405   To examine preferential binding patterns of splicing factors around SSs, RBP peaks were

406   annotated as splicing-related factors if they belong to hnRNP- and SR-families (n=29). We

407   extended both 3' acceptor and 5' donor SS by 1,000 bp in both up and downstream direction and

408   binned the region into 100 bp intervals. We defined the position relative to the distance to the SS,

409   in the 5' to 3' direction. For each interval, we calculated the frequency of splicing factor binding

410   normalized to the size of the interval. The value of RBP enrichment means the normalized

411   binding frequency of splicing-related factors.

412

413   **LSTM model**

414   We adopted the following equations for the modeling of splicing using LSTM. $\sigma$ function

415   denotes sigmoid function. $\otimes$ denotes Hadamard product where two matrices are multiplied in a

416   pair-wise fashion. $x_t$ denotes input vector and $h_t$ denotes output vector, $f_t$ denotes forget gate

417   vector, $i_t$ denotes input or update gate vector, $o_t$ denotes output gate vector, $c_t$ denotes cell state

418   vector.

419

$$f_t = \sigma(W_{hf}^{T} h_{t-1} + W_{xf}^{T} x_t + b_f)$$

$$i_t = \sigma(W_{hi}^{T} h_{t-1} + W_{xi}^{T} x_t + b_i)$$

$$o_t = \sigma(W_{ho}^{T} h_{t-1} + W_{xo}^{T} x_t + b_o)$$

$$g_t = \tanh(W_{hg}^{T} h_{t-1} + W_{xg}^{T} x_t + b_g)$$

$$c_t = f_t \otimes c_{t-1} + i_t \otimes g_t$$

$$h_t = o_t \otimes \tanh(c_t)$$

420

421 **GRU model**

422 We adopted the following equations for the modeling of splicing using GRU. $x_t$ denotes input

423 vector and $h_t$ denotes output vector, $z_t$ denotes update gate vector and $r_t$ denotes reset gate

424 vector.

425

$$z_t = \sigma\left(W_{hz}{}^T h_{t-1} + W_{xz}{}^T x_t + b_z\right)$$

$$r_t = \sigma(W_{hr}{}^T h_{t-1} + W_{xr}{}^T x_t + b_r)$$

$$h_t = (1 - z_t)h_{t-1}\otimes + z_t\otimes\tanh(W_{hh}{}^T(r_t\otimes h_{t-1}) + W_{xh}{}^T x_t + b_h)$$

426

427 **Pre-processing of data for the training model**

428 We selected six normal and three cancer samples from the reference epigenome series. The

429 dataset contains consolidated epigenomes from the Roadmap Epigenomics Consortium [41] and

430 the ENCODE Consortium. All datasets were uniformly processed and mapped to the GRCh38

431 human reference genome. All samples contained a core set of histone modification tracks

432 (H3K4me1, H3K4me3, H3K36me3, H3K27ac, H3K27me3, and H3K9me3) as well as RNA-seq

433 data. We used additional histone modification tracks, as well as DNase I hypersensitivity, DNA

434 methylation, and nucleosome positioning tracks, to predict alternative splicing upon availability.

435 Detailed information on datasets used can be found in Supplementary Table 1. For each exon, we

436 obtained DNA sequences at intron-exon boundaries (3' acceptors) and exon-intron boundaries (5'

437 donors), as well as 100 bp upstream and downstream of SSs. Splice junctions included both

438 canonical and non-canonical SSs. We processed all sequences to read in the 5' to 3' direction

439    using strand information from each gene. Each 400 bp DNA sequence was encoded into a 1,000

440    by 4 binary array using one-hot encoding. We used RNA-seq expression profiles to indicate

441    tissue-specific alternative splicing patterns. Genes having fewer than two exons were discarded

442    and the first and last exons were excluded from the analysis. We classified an exon as being

443    expressed if its FPKM was greater than or equal to 1. We normalized all ChIP-seq histone

444    modification tracks and DNase-seq tracks over corresponding input signal tracks using MACS

445    v2.0.10 (https://github.com/taoliu/MACS) [66]. We used negative log10 of the Poisson p-value

446    to measure the enrichment level over the background. Due to the wide dynamic range observed,

447    we used a p-value threshold of 1e-2 for the upper limit. We processed all feature tracks including

448    DNA methylation and nucleosome signal tracks to read in the 5' to 3' direction and scaled them

449    to a range of 0 to 1.

450

451    **Performance evaluation of the model**

452    There is no single metric that can give you a measure of performance in a binary classification

453    problem. Relying on one metric can be misleading especially when there is high class imbalance.

454    Therefore, we employed various metrics to measure the performance of the predictive model.

455    ROC curve explains the tradeoff between true-positive rate (TPR) and false-positive rate (FPR).

456    PR curve visualizes the tradeoff between positive predictive value (PPV) and true-positive rate

457    (TPR).

458

$$Precison = \frac{tp}{tp + fp}$$

$$Recall = \frac{tp}{tp + fn}$$

$$Accuracy = \frac{tp + tn}{tp + tn + fp + fn}$$

459

460    In addition, we used F1-score, which is the harmonic mean of precision and recall, to measure

461    the performance of the splicing model.

$$F_1 = 2 \cdot \frac{precison \cdot recall}{precison + recall}$$

462

463    **Hyperparameter tuning of splicing model and training**

464    We tested a range of dimensions and depths of RNN models and network design

465    hyperparameters to optimize the alternative splicing model. We chose optimal hyperparameters

466    by tuning one parameter at a time while fixing the rest. Hyperparameters included but were not

467    limited to the number of recurrent layers, size of neurons in each layer, pooling strategy, dropout

468    rate, choice of activation function and loss function, optimizer, and number of the epoch. We

469    shuffled the order of the data and split the dataset into training and test sets using an 80 to 20%

470    ratio. 20% of test data was set aside for the performance evaluation. 80% of training data was

471    split again between 80 to 20% (64 and 16% of the original data) for fitting the model and

472    validating the model fit during the training phase. We fed a range of sequences from 50 to 1,000

473    bp within each SS and found the 400 bp span to be the ideal size for the model. For the neural

474    network architecture, we achieved the best result when two RNN units were stacked together,

475    which allowed the model to learn higher-level temporal representations. We used a hidden state

476    size of two by default and we recommend not using a hidden state size greater than 128 to avoid

477    overfitting problems (Supplementary Figure 8A). We applied three variants of the RNN model,

478    LSTM [45], GRU [67], and simple RNN. To compare the performance of memory-based units

479    (LSTM and GRU), we implemented a simple RNN model using the same network architecture.

480    We found that both LSTM and GRU were capable of learning long-term dependencies and were

481    effective in learning high-dimensional contextual relationships between epigenomic features

482    around the SSs. We split the input sequences into two parts where the first half represented a 3'

483    acceptor SS and the latter half represented a 5' donor SS. We fed these sequences into two

484    separate RNN units of size 200 and merged them into another RNN unit of size 400. The last

485    RNN layer was followed by a dropout layer to prevent overfitting of the training dataset. The last

486    fully-connected layer contained the softmax activation function for classifying exons as either

487    spliced or unspliced. To train the model, we used a binary cross-entropy objective function with

488    the Adam optimizer [68]. For each dataset, we trained the model for 20 epochs. We tested the

489    implementation of ESPRNN using TensorFlow v2.0 (https://www.tensorflow.org). Our

490    implementation also works with Keras v1.0.3 or v2.2.4 (https://github.com/fchollet/keras) with

491    either TensorFlow v1.15 and Theano v0.8.2 [69] backend with a minor tweak. We used various

492    Nvidia GPUs (Titan K20m, K80, GTX 1080ti, RTX2080, P100, and Titan V) to train the model.

493

494    **Acknowledgements**

495    We would like to acknowledge Steve Weston from the Yale Center for Research Computing for

496    technical support in setting up our GPU computing infrastructure.

497

498    **References**

499    1.    Wang ET, Sandberg R, Luo S, Khrebtukova I, Zhang L, Mayr C, et al. Alternative isoform
500          regulation in human tissue transcriptomes. Nature. 2008;456: 470–6.
501          doi:10.1038/nature07509
502    2.    Pan Q, Shai O, Lee LJ, Frey BJ, Blencowe BJ. Deep surveying of alternative splicing
503          complexity in the human transcriptome by high-throughput sequencing. Nat Genet.
504          2008;40: 1413–5. doi:10.1038/ng.259

505    3.    Graveley BR. Alternative splicing: Increasing diversity in the proteomic world. Trends in
506          Genetics. 2001. pp. 100–107. doi:10.1016/S0168-9525(00)02176-4
507    4.    Djebali S, Davis CA, Merkel A, Dobin A, Lassmann T, Mortazavi A, et al. Landscape of
508          transcription in human cells. Nature. 2012;489: 101–108. doi:10.1038/nature11233
509    5.    Listerman I, Sapra AK, Neugebauer KM. Cotranscriptional coupling of splicing factor
510          recruitment and precursor messenger RNA splicing in mammalian cells. Nat Struct Mol
511          Biol. 2006;13: 815–822. doi:10.1038/nsmb1135
512    6.    Wada Y, Ohta Y, Xu M, Tsutsumi S, Minami T, Inoue K, et al. A wave of nascent
513          transcription on activated human genes. Proc Natl Acad Sci. 2009;106: 18357–18361.
514          doi:10.1073/pnas.0902573106
515    7.    Ameur A, Zaghlool A, Halvardson J, Wetterbom A, Gyllensten U, Cavelier L, et al. Total
516          RNA sequencing reveals nascent transcription and widespread co-transcriptional splicing
517          in the human brain. Nat Struct Mol Biol. 2011;18: 1435–1440. doi:10.1038/nsmb.2143
518    8.    Girard C, Will CL, Peng J, Makarov EM, Kastner B, Lemm I, et al. Post-transcriptional
519          spliceosomes are retained in nuclear speckles until splicing completion. Nat Commun.
520          2012;3: 994. doi:10.1038/ncomms1998
521    9.    Carrillo Oesterreich F, Herzel L, Straube K, Hujer K, Howard J, Neugebauer KM.
522          Splicing of Nascent RNA Coincides with Intron Exit from RNA Polymerase II. Cell.
523          2016;165: 372–381. doi:10.1016/j.cell.2016.02.045
524    10.   Alpert T, Herzel L, Neugebauer KM. Perfect timing: splicing and transcription rates in
525          living cells. Wiley Interdisciplinary Reviews: RNA. Blackwell Publishing Ltd; 2017.
526          doi:10.1002/wrna.1401
527    11.   Herzel L, Ottoz DSM, Alpert T, Neugebauer KM. Splicing and transcription touch base:
528          Co-transcriptional spliceosome assembly and function. Nature Reviews Molecular Cell
529          Biology. Nature Publishing Group; 2017. pp. 637–650. doi:10.1038/nrm.2017.63
530    12.   Tanny JC. Chromatin modification by the RNA polymerase II elongation complex.
531          Transcription. 2014;5. doi:10.4161/21541264.2014.988093
532    13.   Tilgner H, Knowles DG, Johnson R, Davis CA, Chakrabortty S, Djebali S, et al. Deep
533          sequencing of subcellular RNA fractions shows splicing to be predominantly co-
534          transcriptional in the human genome but inefficient for lncRNAs. Genome Res. 2012;22:
535          1616–1625. doi:10.1101/gr.134445.111
536    14.   Motta-Mena LB, Heyd F, Lynch KW. Context-Dependent Regulatory Mechanism of the
537          Splicing Factor hnRNP L. Mol Cell. 2010;37: 223–234. doi:10.1016/j.molcel.2009.12.027
538    15.   Fu X-DD, Ares M. Context-dependent control of alternative splicing by RNA-binding
539          proteins. Nat Rev Genet. 2014;15: 689–701. doi:10.1038/nrg3778
540    16.   Andersson R, Enroth S, Rada-Iglesias A, Wadelius C, Komorowski J. Nucleosomes are
541          well positioned in exons and carry characteristic histone modifications. Genome Res.
542          2009;19: 1732–1741. doi:10.1101/gr.092353.109
543    17.   Schwartz S, Meshorer E, Ast G. Chromatin organization marks exon-intron structure. Nat
544          Struct Mol Biol. 2009;16: 990–995. doi:10.1038/nsmb.1659
545    18.   Tilgner H, Nikolaou C, Althammer S, Sammeth M, Beato M, Valcárcel J, et al.
546          Nucleosome positioning as a determinant of exon recognition. Nat Struct Mol Biol.
547          2009;16: 996–1001. doi:10.1038/nsmb.1658
548    19.   Shukla S, Kavak E, Gregory M, Imashimizu M, Shutinoski B, Kashlev M, et al. CTCF-
549          promoted RNA polymerase II pausing links DNA methylation to splicing. Nature.
550          2011;479: 74–9. doi:10.1038/nature10442

20. Lev Maor G, Yearim A, Ast G. The alternative role of DNA methylation in splicing regulation. Trends Genet. 2015;31: 274–280. doi:10.1016/j.tig.2015.03.002

21. Hon G, Wang W, Ren B. Discovery and annotation of functional chromatin signatures in the human genome. Segal E, editor. PLoS Comput Biol. 2009;5: e1000566. doi:10.1371/journal.pcbi.1000566

22. Kolasinska-Zwierz P, Down T, Latorre I, Liu T, Liu XS, Ahringer J. Differential chromatin marking of introns and expressed exons by H3K36me3. Nat Genet. 2009;41: 376–381. doi:10.1038/ng.322

23. Spies N, Nielsen CB, Padgett RA, Burge CB. Biased Chromatin Signatures around Polyadenylation Sites and Exons. Mol Cell. 2009;36: 245–254. doi:10.1016/j.molcel.2009.10.008

24. Yearim A, Gelfman S, Shayevitch R, Melcer S, Glaich O, Mallm JP, et al. HP1 Is Involved in Regulating the Global Impact of DNA Methylation on Alternative Splicing. Cell Rep. 2015;10: 1122–1134. doi:10.1016/j.celrep.2015.01.038

25. Trapnell C, Pachter L, Salzberg SL. TopHat: Discovering splice junctions with RNA-Seq. Bioinformatics. 2009;25: 1105–1111. doi:10.1093/bioinformatics/btp120

26. Au KF, Jiang H, Lin L, Xing Y, Wong WH. Detection of splice junctions from paired-end RNA-seq data by SpliceMap. Nucleic Acids Res. 2010;38: 4570–4578. doi:10.1093/nar/gkq211

27. Pertea M, Lin X, Salzberg SL. GeneSplicer: a new computational method for splice site prediction. Nucleic Acids Res. 2001;29: 1185–90. doi:10.1093/nar/29.5.1185

28. Sonnenburg S, Schweikert G, Philips P, Behr J, Rätsch G. Accurate splice site prediction using support vector machines. BMC Bioinformatics. 2007;8: S7. doi:10.1186/1471-2105-8-S10-S7

29. Jaganathan K, Kyriazopoulou Panagiotopoulou S, McRae JF, Darbandi SF, Knowles D, Li YI, et al. Predicting Splicing from Primary Sequence with Deep Learning. Cell. 2019;0. doi:10.1016/j.cell.2018.12.015

30. Barash Y, Calarco JA, Gao W, Pan Q, Wang X, Shai O, et al. Deciphering the splicing code. Nature. 2010;465: 53–59. doi:10.1038/nature09000

31. Xiong HY, Barash Y, Frey BJ. Bayesian prediction of tissue-regulated splicing using RNA sequence and cellular context. Bioinformatics. 2011;27: 2554–2562. doi:10.1093/bioinformatics/btr444

32. Barash Y, Vaquero-Garcia J, González-Vallinas J, Xiong HY, Gao W, Lee LJ, et al. AVISPA: a web tool for the prediction and analysis of alternative splicing. Genome Biol. 2013;14: R114. doi:10.1186/gb-2013-14-10-r114

33. Cartegni L, Chew SL, Krainer AR. Listening to silence and understanding nonsense: exonic mutations that affect splicing. Nat Rev Genet. 2002;3: 285–298. doi:10.1038/nrg775

34. Luco RF, Pan Q, Tominaga K, Blencowe BJ, Pereira-Smith OM, Misteli T. Regulation of alternative splicing by histone modifications. Science (80- ). 2010;327: 996–1000. doi:10.1126/science.1184208

35. Witten JT, Ule J. Understanding splicing regulation through RNA splicing maps. Trends Genet. 2011;27: 89–97. doi:10.1016/j.tig.2010.12.001

36. Lareau LF, Inada M, Green RE, Wengrod JC, Brenner SE. Unproductive splicing of SR genes associated with highly conserved and ultraconserved DNA elements. Nature. 2007;446: 926–929. doi:10.1038/nature05676

597  37.   Pervouchine D, Popov Y, Berry A, Borsari B, Frankish A, Guigó R. Integrative
598        transcriptomic analysis suggests new autoregulatory splicing events coupled with
599        nonsense-mediated mRNA decay. Nucleic Acids Res. 2019;47: 5293–5306.
600        doi:10.1093/nar/gkz193
601  38.   Ruskin B, Zamore PD, Green MR. A factor, U2AF, is required for U2 snRNP binding and
602        splicing complex assembly. Cell. 1988;52: 207–219. doi:10.1016/0092-8674(88)90509-0
603  39.   Rasche N, Dybkov O, Schmitzová J, Akyildiz B, Fabrizio P, Lührmann R. Cwc2 and its
604        human homologue RBM22 promote an active conformation of the spliceosome catalytic
605        centre. EMBO J. 2012;31: 1591. doi:10.1038/EMBOJ.2011.502
606  40.   Wickramasinghe VO, Gonzàlez-Porta M, Perera D, Bartolozzi AR, Sibley CR, Hallegger
607        M, et al. Regulation of constitutive and alternative mRNA splicing across the human
608        transcriptome by PRPF8 is determined by 5′ splice site strength. Genome Biol. 2015;16:
609        201. doi:10.1186/s13059-015-0749-3
610  41.   Consortium RE, Kundaje A, Meuleman W, Ernst J, Bilenky M, Yen A, et al. Integrative
611        analysis of 111 reference human epigenomes. Nature. 2015;518: 317–330.
612        doi:10.1038/nature14248
613  42.   Graves A, Mohamed A, Hinton G. Speech Recognition with Deep Recurrent Neural
614        Networks. IEEE Int Conf Acoust Speech Signal Process. 2013; 6645–6649.
615        doi:10.1109/ICASSP.2013.6638947
616  43.   Cho K, van Merrienboer B, Gulcehre C, Bahdanau D, Bougares F, Schwenk H, et al.
617        Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine
618        Translation. Proc 2014 Conf Empir Methods Nat Lang Process. 2014; 1724–1734.
619        doi:10.3115/v1/D14-1179
620  44.   Bahdanau D, Cho K, Bengio Y. Neural Machine Translation By Jointly Learning To
621        Align and Translate. Iclr 2015. 2014; 1–15. doi:10.1146/annurev.neuro.26.041002.131047
622  45.   Hochreiter S, Schmidhuber J, Hochreiter S, Schmidhuber J, Schmidhuber J. Long short-
623        term memory. Neural Comput. 1997;9: 1735–80. doi:10.1162/neco.1997.9.8.1735
624  46.   Quang D, Xie X. DanQ: A hybrid convolutional and recurrent deep neural network for
625        quantifying the function of DNA sequences. Nucleic Acids Res. 2016;44: gkw226.
626        doi:10.1093/nar/gkw226
627  47.   Lee B, Lee T, Na B, Yoon S. DNA-Level Splice Junction Prediction using Deep
628        Recurrent Neural Networks. arXiv e-prints. 2015; 1–6. Available:
629        http://arxiv.org/abs/1512.05135
630  48.   Koscielny G, Texier V Le, Gopalakrishnan C, Kumanduri V, Riethoven JJ, Nardone F, et
631        al. ASTD: The Alternative Splicing and Transcript Diversity database. Genomics. 2009;93:
632        213–220. doi:10.1016/j.ygeno.2008.11.003
633  49.   Keren H, Lev-Maor G, Ast G. Alternative splicing and evolution: diversification, exon
634        definition and function. Nat Rev Genet. 2010;11: 345–55. doi:10.1038/nrg2776
635  50.   Liu S, Cheng C. Alternative RNA splicing and cancer. Wiley Interdiscip Rev RNA.
636        2013;4: 547–566. doi:10.1002/wrna.1178
637  51.   Oltean S, Bates DO. Hallmarks of alternative splicing in cancer. Oncogene. 2014;33:
638        5311–5318. doi:10.1038/onc.2013.533
639  52.   Jiang P, Freedman ML, Liu JS, Liu XS. Inference of transcriptional regulation in cancers.
640        Proc Natl Acad Sci U S A. 2015. doi:10.1073/pnas.1424272112
641  53.   Ntziachristos P, Abdel-Wahab O, Aifantis I. Emerging concepts of epigenetic
642        dysregulation in hematological malignancies. Nature Immunology. 2016.

643          doi:10.1038/ni.3517

644   54.   Jung H, Lee D, Lee J, Park D, Kim YJ, Park W-Y, et al. Intron retention is a widespread
645          mechanism of tumor-suppressor inactivation. Nat Genet. 2015;47: 1242–1248.
646          doi:10.1038/ng.3414

647   55.   Obeng EA, Ebert BL. Charting the "Splice" Routes to MDS. Cancer Cell. 2015.
648          doi:10.1016/j.ccell.2015.04.016

649   56.   Pradeepa MM, Sutherland HG, Ule J, Grimes GR, Bickmore WA. Psip1/Ledgf p52 Binds
650          Methylated Histone H3K36 and Splicing Factors and Contributes to the Regulation of
651          Alternative Splicing. Reik W, editor. PLoS Genet. 2012;8: e1002717.
652          doi:10.1371/journal.pgen.1002717

653   57.   Huff JT, Plocik AM, Guthrie C, Yamamoto KR. Reciprocal intronic and exonic histone
654          modification regions in humans. Nat Struct Mol Biol. 2010;17: 1495–1499.
655          doi:10.1038/nsmb.1924

656   58.   Venables JP. Aberrant and alternative splicing in cancer. Cancer Research. 2004. pp.
657          7647–7654. doi:10.1158/0008-5472.CAN-04-1910

658   59.   Tazi J, Bakkour N, Stamm S. Alternative splicing and disease. Biochimica et Biophysica
659          Acta - Molecular Basis of Disease. 2009. pp. 14–26. doi:10.1016/j.bbadis.2008.09.017

660   60.   Sveen A, Kilpinen S, Ruusulehto A, Lothe RA, Skotheim RI. Aberrant RNA splicing in
661          cancer; expression changes and driver mutations of splicing factor genes. Oncogene.
662          2016;35: 2413–2427. doi:10.1038/onc.2015.318

663   61.   Faustino NA, Cooper TA. Pre-mRNA splicing and human disease. Genes and
664          Development. Cold Spring Harbor Lab; 2003. pp. 419–437. doi:10.1101/gad.1048803

665   62.   Xiong HY, Alipanahi B, Lee LJ, Bretschneider H, Merico D, Yuen RKC, et al. The
666          human splicing code reveals new insights into the genetic determinants of disease. Science
667          (80- ). 2014;347: 1254806-. doi:10.1126/science.1254806

668   63.   Anders S, Pyl PT, Huber W. HTSeq-A Python framework to work with high-throughput
669          sequencing data. Bioinformatics. 2015. doi:10.1093/bioinformatics/btu638

670   64.   Schafer S, Miao K, Benson CC, Heinig M, Cook SA, Hubner N. Alternative Splicing
671          Signatures in RNA☐seq Data: Percent Spliced in (PSI). Curr Protoc Hum Genet. 2015;87.
672          doi:10.1002/0471142905.hg1116s87

673   65.   Lovci MT, Ghanem D, Marr H, Arnold J, Gee S, Parra M, et al. Rbfox proteins regulate
674          alternative mRNA splicing through evolutionarily conserved RNA bridges. Nat Struct
675          Mol Biol. 2013;20: 1434–1442. doi:10.1038/nsmb.2699

676   66.   Zhang Y, Liu T, Meyer CA, Eeckhoute J, Johnson DS, Bernstein BE, et al. Model-based
677          Analysis of ChIP-Seq (MACS). Genome Biol. 2008;9: R137. doi:10.1186/gb-2008-9-9-
678          r137

679   67.   Cho K, van Merrienboer B, Bahdanau D, Bengio Y. On the Properties of Neural Machine
680          Translation: Encoder–Decoder Approaches. Proc SSST-8, Eighth Work Syntax Semant
681          Struct Stat Transl. 2014; 103–111. Available: http://arxiv.org/abs/1409.1259

682   68.   Kingma D, Ba J. Adam: A Method for Stochastic Optimization. Int Conf Learn Represent.
683          2014; 1–13. Available: http://arxiv.org/abs/1412.6980

684   69.   Theano Development Team. Theano: A Python framework for fast computation of
685          mathematical expressions. arXiv e-prints. 2016; 19. Available:
686          http://arxiv.org/abs/1605.02688

687

688

689     **<u>Supporting Information Legend</u>**

690     **Supplementary Table 1**

691     **List of datasets and the accession numbers used for the study.**

692

693     **Supplementary Table 2**

694     **Overview of dataset used for training the ESPRNN model.** The model was trained using the

695     CORE (highlighted in red) and FULL set based on the availability of data. The CORE set was

696     used to compare the predictive performance across cell types.

697

698     **Supplementary Table 3**

699     **ESPRNN model prediction performance measured by F1 score.** Predictive performance was

700     compared between the CORE and FULL set of genomic features. For each set, performance was

701     compared using LSTM, GRU, and simple RNN models. Predictive performance was measured

702     by F1 score.

703

704     **Supplementary Table 4**

705     **Comparison of models trained with 50 bp span and 100 bp span data.** Each model was

706     trained using genomic features derived from 50 bp span or 100 bp span data from splice sites

707     using the LSTM model. Performance was measured using F1 score and ROC AUC.

708

709     **Supplementary Figure 1**

710   (Shadow figure of the main Figure 2A) Enrichment of various epigenomic marks of HepG2 at

711   the exon-intron boundary. High PSI indicates exon inclusion, mid PSI indicates exons with 40-

712   60% PSI, and low PSI indicates exon skipping.

713

714   **Supplementary Figure 2**

715   (Shadow figure of the main Figure 2B) Comparison of epigenetic enrichment around different

716   segments of the 3' acceptor site for **(A)** K562 and **(B)** HepG2. High PSI indicates exon inclusion,

717   mid PSI indicates exons with 40-60% PSI, and low PSI indicates exon skipping. Mann-Whitney-

718   Wilcoxon two-sided test, ns: $0.05 < p <= 1$; *: $0.01 < p <= 0.05$; **: $0.001 < p <= 0.01$; ***:

719   $0.0001 < p <= 0.001$; ****: $p <= 0.0001$. **(C)** Fold enrichment of splicing-related RBPs to non-

720   splicing-related RBPs around the 3' acceptor splice site and 5' donor splice site.

721

722   **Supplementary Figure 3**

723   Correlation of exonic expression (FPKM) and histone enrichment of **(A)** HepG2 H3K36me3, **(B)**

724   HepG2 H3K27me3, **(C)** liver H3K36me3, and **(D)** liver H3K27me3. PCC: Pearson Correlation

725   Coefficient.

726

727   **Supplementary Figure 4**

728   Splicing patterns based on exonic expression level (FPKM) for diverse ENCODE cell types are

729   projected on a PCA cell space.

730

731   **Supplementary Figure 5**

732    **(A)** Difference in splicing prediction performance when RBP binding profiles were added as an

733    additional feature of the base model containing chromatin accessibility and histone marks. **(B)**

734    Cross-cell testing of model. Model was trained on HepG2 data and tested on K562 data, and vice

735    versa.

736

737    **Supplementary Figure 6**

738    **(A)** Comparison of the baseline model trained using chromatin accessibility and 6 histone marks

739    to a model using DNA sequence feature only **(B)** Measure of information gain from additional

740    epigenetic feature based on DNA sequence only model **(C)** Comparison of splicing prediction

741    performance using a pair of epigenetic features.

742

743    **Supplementary Figure 7**

744    Comparison of LSTM-based model with other machine learning algorithms. Four different

745    algorithms, k-Nearest neighbor (kNN), decision tree, random forest, and support vector machine

746    (SVM), were compared to the LSTM-based model across four different tissue types (A549,

747    HepG2, GM12878, K562).

748

749    **Supplementary Figure 8**

750    **(A)** Comparison of splicing prediction performance across different sizes of hidden state. **(B)**

751    Loss of training an LSTM model with 1 hidden layer for 400 epochs. **(C)** Accuracy of training an

752    LSTM model with one hidden layer for 400 epochs. **(D)** Trained weights of LSTM recurrent

753    cells.

754

755    **Supplementary Figure 9**

756    Comparison of splicing prediction performance when input features are reversed in time-

757    direction via a **(A)** precision-recall curve and a **(B)** ROC curve.

758

759    **<u>Figure Legend</u>**

760    **Figure 1**

761    Overview of the co-transcriptional splicing model. Depiction of co-transcriptional splicing in

762    terms of **(A)** biological context, **(B)** genomic and epigenomic data context, and how it relates to

763    the **(C)** RNN model.

764

765    **Figure 2**

766    **(A)** Enrichment of various epigenomic marks of K562 at the exon-intron boundary. We

767    aggregated histone modifications up to 500 bp upstream and downstream of intronic and exonic

768    regions flanking 3' and 5' SSs for cassette exons across ENCODE cell types. High PSI indicates

769    exon inclusion, mid PSI indicates exons with 40-60% PSI, and low PSI indicates exon skipping.

770    **(B)** Statistical significance testing of epigenetic mark enrichment. Average histone modification

771    enrichment at four exonic segments were compared based on PSI values. Mann-Whitney-

772    Wilcoxon two-sided test, ns: $0.05 < p <= 1$; *: $0.01 < p <= 0.05$; **: $0.001 < p <= 0.01$; ***:

773    $0.0001 < p <= 0.001$; ****: $p <= 0.0001$. **(C)** RBP enrichment across the exon-intron boundary.

774

775    **Figure 3**

776    Correlation of exonic expression to **(A)** H3K36me3 and **(B)** H3K27me3. The line represents a

777    linear regression model fit, and the shaded band represents 95% confidence interval. **(C)**

778    Alternative exons were grouped by expression level and their relative histone enrichment was

779    compared near the SSs. Asterisks represents statistical significance using the Wilcoxon rank sum

780    test; (*) P <= 0.05, (**) P <= 0.01, (***) P <= 0.001, (****) P <= 0.0001. **(D)** Hierarchical

781    clustering of similarity based on PSI across 49 ENCODE biosamples. The results are clustered

782    into five categories of cell types.

783

784    **Figure 4**

785    **(A)** Overview of the ESPRNN model. The model is composed of two recurrent layers. Inputs

786    from 3' and 5' SSs are separately processed in the first recurrent layer and then merged in the

787    next recurrent layer. A softmax classifier is used to determine the inclusion of the exon. Using

788    genomic sequences and epigenomic contexts as input, the alternative usage of the exon is

789    predicted. **(B)** Precision-recall curves from six different ENCODE cell types. **(C)** Epigenetic

790    features that contribute to splicing regulation. The order and magnitude of importance was

791    determined using leave-one-out analysis and loss of the ROC AUC was calculated when training

792    the model lacking a particular feature. **(D)** Comparison of LSTM model with other models based

793    on k-nearest neighbor, support vector machine, decision tree, and random forest algorithms.

A

K562 H3K36me3

POC = 0.092, P-val = 2.9e-64

Enrichment

Expression

B

K562 H3K27me3

POC = -0.052, P-val = 1.8e-27

Enrichment

Expression

C

Enrichment

High Expr.
Low Expr.

K562
H3K36me3

HepG2
H3K36me3

Liver
H3K36me3

K562
H3K27me3

HepG2
H3K27me3

Liver
H3K27me3

D

biosample

C1
C2
C3
C4
C5

A  HepG2 H3K36me3 — PCC = 0.021, P-val = 0.00022

B  HepG2 H3K27me3 — PCC = -0.02, P-val = 0.013

C  liver H3K36me3 — PCC = 0.0046, P-val = 0.43
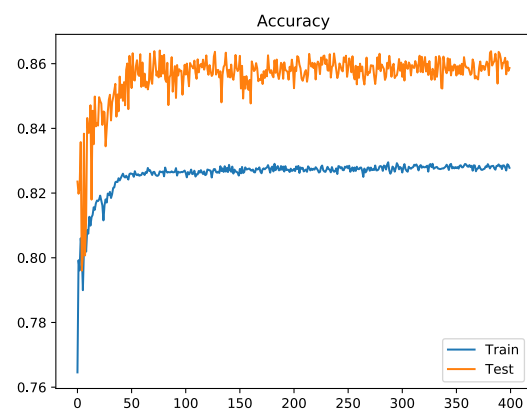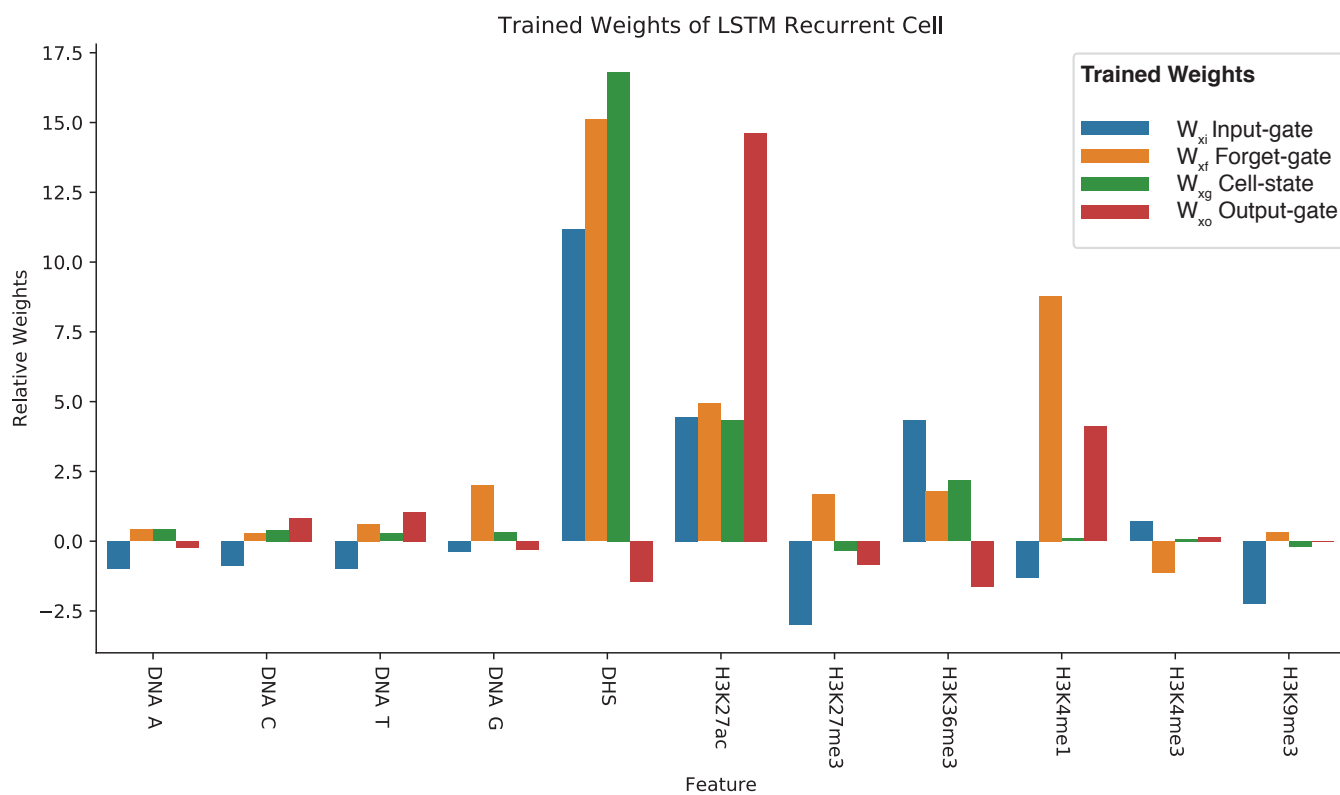
D  liver H3K27me3 — PCC = -0.0068, P-val = 0.27

**A**



**B**

**A**



**B**



**C**



**D**

**A**



**B**



**C**



**D**

**A**



**B**