1  **DeeReCT-APA: Prediction of Alternative Polyadenylation Site**

2  **Usage Through Deep Learning**

3  Zhongxiao Li[1,a], Yisheng Li[2,b], Bin Zhang[3,c], Yu Li[1,d], Yongkang Long[1,2,e], Juexiao

4  Zhou[2,f], Xudong Zou[2,g], Min Zhang[2,h], Yuhui Hu[2,*,i], Wei Chen[2,*,j], Xin Gao[1,*,k]

5  *[1]King Abdullah University of Science and Technology (KAUST), Computational*

6  *Bioscience Research Center (CBRC), Computer, Electrical and Mathematical Sciences*

7  *and Engineering (CEMSE) Division, Thuwal, 23955-6900, Saudi Arabia.*

8  *[2]Department of Biology, Southern University of Science and Technology (SUSTech),*

9  *Shenzhen, 518055, China.*

10  *[3]Cancer Science Institute of Singapore, Singapore 117599, Singapore.*

11  [*] Corresponding author.

12  E-mail:  huyh@sustech.edu.cn(Hu  Y),  chenw@sustech.edu.cn(Chen  W),

13  xin.gao@kaust.edu.sa (Gao X)

14  **Running title:**

15  *Li Z et al / DeeReCT-APA: Deep Learning Prediction of APA*

16  [a] ORCID: 0000-0003-2480-0750

17  [b] ORCID: 0000-0001-8015-7128

18  [c] ORCID: 0000-0001-8835-8370

19  [d] ORCID: 0000-0002-3664-6722

20  [e] ORCID: 0000-0003-0953-7325

21  [f] ORCID: 0000-0002-6739-6236

22  [g] ORCID: 0000-0002-2958-0438

23  [h] ORCID: 0000-0002-3462-3711

24  [i] ORCID: 0000-0002-5210-5301

25  [j] ORCID: 0000-0003-3263-1627

26  [k] ORCID: 0000-0002-7108-3574

27  Word count: 8762

28  Figure count: 7

29  Table count: 2

30  Supplementary Figures: 4

31  Supplementary Tables: 7

32  Reference Count: 42

33    Reference Count after 2014: 21

34    Count of letters in the article title: 78

35    Count of letters in the running title: 47

36    Count of keywords: 4

37    Count of words in abstract: 246

38

39

40 **Abstract**

41 Alternative polyadenylation (APA) is a crucial step in post-transcriptional regulation.

42 Previous bioinformatic works have mainly focused on the recognition of

43 polyadenylation sites (PAS) in a given genomic sequence, which is a binary

44 classification problem. Recently, computational methods for predicting the usage level

45 of alternative PAS in a same gene have been proposed. However, all of them cast the

46 problem as a non-quantitative pairwise comparison task and do not take the competition

47 among multiple PAS into account. To address this, here we propose a deep learning

48 architecture, DeeReCT-APA, to quantitatively predict the usage of all alternative PAS

49 of a given gene. To accommodate different genes with potentially different numbers of

50 PAS, DeeReCT-APA treats the problem as a regression task with a variable-length

51 target. Based on a CNN-LSTM architecture, DeeReCT-APA extracts sequence features

52 with CNN layers, uses bidirectional LSTM to explicitly model the interactions among

53 competing PAS, and outputs percentage scores representing the usage levels of all PAS

54 of a gene. In addition to the fact that only our method can predict quantitatively the

55 usage of all the PAS within a gene, we show that our method consistently outperforms

56 other existing methods on three different tasks for which they are trained: pairwise

57 comparison task, highest usage prediction task and ranking task. Finally, we

58 demonstrate that our method can be used to predict the effect of genetic variations on

59 APA patterns and shed light on future mechanistic understanding in APA regulation.

60 Our code and data are available at https://github.com/lzx325/DeeReCT-APA-repo.

61

62 **KEYWORDS:** Polyadenylation; Gene regulation; Deep learning; Bioinformatics

## Introduction

In eukaryotic cells, the termination of Pol II transcription involves 3'-end cleavage followed by addition of a poly(A) tail, a process termed as "polyadenylation". Often, one gene could have multiple polyadenylation sites (PAS). The so-called alternative polyadenylation (APA) could generate from the same gene locus different transcript isoforms with different 3'-UTRs and sometimes even different protein coding sequences. The diverse 3'-UTRs generated by APA may contain different sets of *cis*-regulatory elements, thereby modulating the mRNA stability [1–3], translation [4], subcellular localization of mRNAs [5–7], or even the subcellular localization and function of the encoded proteins [8]. Importantly, it has been shown that dysregulation of APA could result in various human diseases [9–12].

APA is regulated by the interaction between *cis*-elements located in the vicinity of PAS and the associated *trans*-factors [13]. The most well-known *cis*-element that defines a PAS is the hexamer AAUAAA and its variants located 15-30nt upstream of the cleavage site, which is directly recognized by the cleavage and polyadenylation specificity factor (CPSF) components: CPSF30 and WDR33 [14]. Other auxiliary *cis*-elements located upstream or downstream of the cleavage site include upstream UGUA motifs bound by the cleavage factor Im (CFIm) and downstream U-rich or GU-rich elements targeted by the cleavage stimulation factor (CstF) [14]. The usage of individual PAS for a multi-PAS gene depends on how efficiently each alternative PAS is recognized by these 3' end processing machineries, which is further regulated by additional RNA binding proteins (RBPs) that could enhance or repress the usage of distinct PAS signals through binding in their proximity. In addition, the usage of alternative PAS is mutually exclusive. In particular, once an upstream PAS is utilized, all the downstream ones would have no chance to be used no matter how strong their PAS signals are. Therefore, proximal PAS, which are transcribed first, have positional advantage over the distal competing PAS [15]. Indeed, it has been observed that the terminal PAS more often contain the canonical AAUAAA hexamer, which is considered to have higher affinity than the other variants, which possibly compensates for their positional disadvantage [16].

There has been a long-standing interest in predicting PAS based on genomic sequences using purely computational approaches. The so-called "PAS recognition

96    problem" aims to discriminate between nucleotide sequences that contain a PAS and

97    those do not. A variety of hand-crafted features have been proposed and statistical

98    learning algorithms, *e.g.*, random forest (RF), support vector machines (SVM) and

99    hidden Markov models (HMM), are then applied on these features to solve the binary

100   classification problem [17–19]. Very recently researchers started investigating the

101   "PAS quantification problem", which aims to predict a score that represents the strength

102   of a PAS [20, 21]. This is much more difficult than the recognition one.

103

104   Recent developments in deep learning have made great improvements on many tasks

105   [22]. With remarkable success, it has also been applied to bioinformatics tasks such as

106   protein-DNA binding [23], RNA splicing pattern prediction [24], enzyme function

107   prediction [25, 26], Nanopore sequencing [27, 28], and promoter prediction [29]. Deep

108   learning is favored due to its automatic feature extraction ability and good scalability

109   with large amount of data. As for polyadenylation prediction, deep learning models

110   have been applied on the PAS recognition problem and they outperformed existing

111   feature-based methods by a large margin [30]. Recently, deep learning models have

112   also been applied on the PAS quantification problem, where Polyadenylation Code [20]

113   was developed to predict the stronger one from a given pair of two competing PAS.

114   Very recently, another model, DeepPASTA [21] has been proposed. DeepPASTA

115   contains four different modules that deal with both the PAS recognition problem and

116   PAS quantification problem. Similar as Polyadenylation Code, DeepPASTA also casts

117   the PAS quantification problem into a pairwise comparison task.

118

119   In this paper, we propose a novel deep learning method, DeeReCT-APA (Deep

120   Regulatory Code and Tools for Alternative Polyadenylation), for the PAS

121   quantification problem. DeeReCT-APA can quantitatively predict the usage of all the

122   competing PAS from a same gene simultaneously, regardless of the number of PAS.

123   The model is trained and evaluated based on the dataset from a previous study [31],

124   which consists of a genome-wide PAS measurement of two different mouse strains

125   (C57BL/6J (BL) and SPRET/EiJ (SP)), and their F1 hybrid. After training our model

126   on the dataset, we comprehensively evaluate our model based on a number of criteria.

127   We demonstrate the necessity of modeling the competition among multiple PAS

128   simultaneously. Finally, we show that our model can predict the effect of genetic

129    variations on APA patterns, visualize APA regulatory motifs and potentially facilitate

130    the mechanistic understanding of APA regulation.

131

## Methods

132

### Description of DeeReCT-APA architecture

133

134    The DeeReCT-APA method is based on a deep learning architecture that contains a set

135    of neural network models composed of base networks (Base-Net, one for each

136    competing PAS) and upper-level interaction layers. Each base network takes a 455nt

137    long genomic DNA sequence centered around one competing PAS cleavage site as

138    input and gives as output a vector which can be interpreted as the distilled features of

139    that sequence. There are two types of base networks in our design, based on: (1) hand-

140    engineered feature extractor and (2) convolutional neural networks (CNN). The output

141    of the lower-level base network is then passed to the upper-level interaction layers,

142    which computationally model the process of choosing competing PAS. The interaction

143    layers of DeeReCT-APA are based on Long Short Term Memory Networks (LSTM)

144    [32], which have achieved remarkable success in natural language processing and can

145    naturally handle sentences with an arbitrary length, therefore suitable for handling any

146    number of alternative PAS from a same gene locus. The interaction layers then output

147    the percentage values of all the competing PAS of the gene. The architecture is

148    illustrated in **Figure 1**. The design of each part of the network is further explained in

149    the following subsections.

150

151      We use three different base network designs: deep neural network architectures

152    based on a single 1D convolution layer (Single-Conv-Net), multiple 1D convolution

153    layers (Multi-Conv-Net) and a handcrafted feature extractor with fully-connected

154    layers (Feature-Net). Single-Conv-Net and Multi-Conv-Net are two convolutional

155    neural network (CNN) structures for Base-Net. The Single-Conv-Net consists of only

156    one layer of the 1D convolutional layer and takes directly the one-hot encoded

157    sequences as input. The convolutional layer has a number of convolution filters which

158    become automatically-learned feature extractors after training. A rectified linear unit

159    (ReLU) is used as the activation function. The max-pooling operation after that allows

160    only values from highly-activated neurons to pass to the upper fully-connected layers.

161    The three operations: convolution, ReLU and max-pooling form a convolution block.

162    While the Single-Conv-Net uses one convolution block, the Multi-Conv-Net uses two

163    convolution blocks before fully-connected layers. The increased depth of the network

164    makes it possible for the network to learn more complex representations. The structures

165    of Single-Conv-Net and Multi-Conv-Net are shown in **Figure 2A** and **Figure 2B**,

166    respectively.

167

168    As a comparison, we also design a base network that works with hand engineered

169    features which we call Feature-Net. The Feature-Net only consists of multiple fully-

170    connected layers and takes as input multiple types of features extracted from the

171    sequences of interest. The features, described in [20], include polyadenylation signals,

172    auxiliary upstream elements, core upstream elements, core downstream elements,

173    auxiliary downstream elements [33], RNA-binding protein motifs, as well as 1-mer, 2-

174    mer, 3-mer, and 4-mer features (detailed in Supplementary Materials Section S1 and

175    Supplementary Table S1). Each feature value corresponds to the occurrence of each

176    motif. The extracted features are then z-score normalized. The architecture is illustrated

177    in **Figure 2C**.

178

179    **Design of the interaction layers**

180    The utilization of alternative PAS is intrinsically competitive. On the one hand, as a

181    multi-PAS gene is transcribed, any one of its PAS along the already transcribed region

182    is possible to be used. But if one of them has already been used, it will make other PAS

183    impossible to be chosen. On the other hand, given that the same polyadenylation

184    machinery is used by all the alternative PAS, such competition of resources also

185    contributes to the competitiveness of this process. However, previous work in

186    polyadenylation usage prediction did not take this important point into account [20, 21].

187    Both existing models, Polyadenylation Code and DeepPASTA (tissue-specific

188    relatively dominant poly(A) sites prediction model, Section 2.3 in [21]) can only take

189    in two PAS at a time, ignoring the competition with others. Here, to overcome this

190    limitation, we consider all the competing PAS at the same time and take as input all the

191    PAS in a gene simultaneously into our model, then jointly predict the usage levels of

192    all of them.

193

194    To fulfil this, we design the interaction layers above the base networks to model the
195    interaction between different PAS. In neural networks, the most common way to model
196    interactions among inputs is to introduce a recurrent neural network (RNN) layer, which
197    can capture the interdependencies among inputs corresponding to each time step. We
198    decide to choose the LSTM [32] as the foundation of interaction layers. LSTM is a type
199    of RNN that has hidden memory cells which are able to remember a state for an
200    arbitrary length of time steps, making it one of the most popular RNNs. To fit into the
201    PAS usage level prediction task, each time step of LSTM corresponds to one PAS, at
202    which the LSTM takes the extracted features of that PAS from the lower-level base
203    network. As there is both influence from upstream PAS to downstream PAS and vice
204    versa, we decide to use a bidirectional LSTM (BiLSTM), in which one LSTM's time
205    step goes from upstream PAS to downstream one and the other from downstream to
206    upstream. The outputs of the two LSTMs at the same PAS are then concatenated and
207    sent to the upper fully-connected layer. The fully-connected layer transforms the LSTM
208    output to a scalar value representing the log-probability of that PAS to be used. After
209    the log-probabilities of all competing PAS pass through a final SoftMax layer, they are
210    transformed to properly normalized percentage scores, which sum up to one,
211    representing their probability of being chosen. The detailed architecture is shown in
212    Figure 1. We point out that, although DeepPASTA also contains a BiLSTM component,
213    their BiLSTM layer is to process the sequence of one of the two competing PAS that
214    are given as input. The time steps of the BiLSTM correspond to different positions in
215    one particular sequence rather than to different PAS, and therefore the BiLSTM is not
216    to model the interactions between different PAS, which is clearly different from the
217    design in DeeReCT-APA.

218

219    As mentioned above, the aim of our model is to take all PAS of a gene as a whole
220    and try to predict the usage level of each PAS as accurate as possible. Therefore, at one
221    time, we must take all PAS in a gene as input. Considering that the number of PAS
222    within a gene is not a constant, we design our model to take inputs of a variable length.
223    Since most genes have a small number of PAS, we choose not to pad all the genes with
224    dummy PAS to make them of the same length, otherwise it will be highly inefficient.
225    Instead, we design the interaction layers in a way that it can take an arbitrary number
226    of Base-Nets.

227    We further design two experiments for ablation study of DeeReCT-APA's BiLSTM

228    interaction layer. The first is to remove the BiLSTM layer and only keep the fully-

229    connected layer and the SoftMax layer. In this scenario, the network still considers all

230    PAS of a gene simultaneously, but with a non-RNN interaction layer. The second is to

231    remove the interaction layer altogether and use comparison-based training (like in

232    Polyadenylation Code) to train a Base-Net. We show their performance separately in

233    the "Overall Performance" section.

234    **A genome-wide PAS quantification dataset derived from fibroblast cells of**

235    **C57BL/6J (BL) and SPRET/EiJ (SP) mouse and their F1 hybrid**

236    A genome-wide PAS quantification dataset derived from fibroblast cells of C57BL/6J

237    (BL) and SPRET/EiJ (SP), as well as their F1 hybrid is obtained from the previous

238    study [31]. In the F1 cells, the two alleles have the same *trans* environment and the PAS

239    usage difference between two alleles can only be due to the sequence variants between

240    their genome sequences, making it a valuable system for APA *cis*-regulation study.

241    Apart from APA, this kind of systems have also been used in the study of alternative

242    splicing and translational regulation [34, 35].

243

244    The detailed description of the sequencing protocol and data analysis procedure can

245    be found in [31]. As a brief summary, the study uses fibroblast cell lines from BL, SP

246    and their F1 hybrids. The total RNA is extracted from fibroblast cells of BL and SP

247    undergoes 3'-Region Extraction and Deep Sequencing (3'READS) [16] to build a good

248    PAS reference of the two strains. The 3'-mRNA sequencing is then performed in all

249    three cell lines to quantify those PAS in the reference. In the F1 hybrid cell, reads are

250    assigned to BL and SP alleles according to their strain specific SNPs. The PAS usage

251    values are then computed by counting the sequencing reads assigned to each PAS. The

252    sequence centering around each PAS cleavage site (448nt in total) is extracted and

253    undergoes feature extraction or one-hot encoding before training the model. The

254    extracted features are then inputted to Feature-Net, while the one-hot encoded

255    sequences are inputted to Single-Conv-Net and Multi-Conv-Net.

256    **Training and evaluation of the model**

257    We train the DeeReCT-APA models based on the parental BL/SP PAS usage level

258    dataset. For F1 hybrid data, however, we choose to start from the pre-trained parental

259  model (which we use either the BL parental model or the SP parental model and the
260  results are shown separately) and fine-tune the model on the F1 dataset. This is because,
261  due to the read assignment problem, the usage of many PAS in F1 cannot be
262  unambiguously characterized by 3'-mRNA sequencing [31]. As a result, the F1 dataset
263  does not contain enough number of PAS to train our model from scratch. At the training
264  stage, genes are randomly selected from the training set and the sequences of their PAS
265  flanking regions are fed into the network. Each sequence of PAS in a gene passes
266  through one Base-Net. The parameters of the Base-Net that are responsible for each
267  PAS are all shared. The Base-Net then each outputs a vector representing distilled
268  features for each PAS, which is then sent to the interaction layers. The interaction layers
269  generate a percentage score of each PAS of this gene. Cross-entropy loss between the
270  predicted usage and the actual usage is used as the training target. During back-
271  propagation, the gradients are back-propagated through the passage originated from
272  each PAS. As the model parameters are shared between base networks, the gradients
273  are then summed up to update the model parameters. We use several techniques to
274  reduce overfitting: (1) Weight decay is applied on weight parameters of CNN and all
275  fully-connected layers. (2) Dropout is applied on BiLSTM. (3) We stop training as soon
276  as the mean absolute error of the predicted usage value does not improve on the
277  validation set. (4) While fine-tuning the model on F1 dataset, we use a learning rate that
278  is ~100 times smaller than the one used when training from scratch.

279

280  The network is trained with the adaptive moment estimation (Adam) optimizer [36].
281  A detailed list of hyperparameters we used is specified in Supplementary Materials
282  Section S2 and Supplementary Table S2. We construct the network using the PyTorch
283  deep learning framework [37] and utilize one NVIDIA GeForce GTX 980 Ti as the
284  GPU hardware platform.

285

286  To evaluate the performance of the model, we conduct a 5-fold cross validation at
287  the gene level using all the genes in our dataset for each strain. That is, if a gene is
288  selected as a training (testing) sample, all of its PAS are in the train (test) set.  At each
289  time, four folds are used for training and the remaining one is used for testing. To make
290  a fair comparison with Polyadenylation Code and DeepPASTA in Section 3.1, we also
291  train (fine-tune) the two models and optimize their model parameters on the parental
292  and F1 datasets.

293

**Performance measures**

To comprehensively evaluate DeeReCT-APA and compare it against baseline and state-of-the-art methods, we use the following performance measures.

*Mean Absolute Error (MAE).* This metric is defined as the mean absolute error (MAE) of the usage prediction of each PAS, which is

$$MAE = \frac{1}{M} \sum_{i=1}^{M} |p_i - t_i| \qquad (1)$$

where $p_i$ stands for the predicted usage, $t_i$ stands for the experimentally determined ground truth usage for PAS i and M is the total number of PAS across all genes in the test set. This is the most intuitive way of measuring the performance of DeeReCT-APA. However, it is not applicable to Polyadenylation Code [20] or DeepPASTA [21] as they do not have quantitative outputs that can be interpreted as the PAS usage values. For the same reason, it is not applicable to DeeReCT-APA either, when its interaction layers are removed and use comparison-based training (Section "Design of the interaction layers").

*Comparison Accuracy.* We here define the Pairwise Comparison Task. We enumerate all the pairs of PAS in a given gene and keep those pairs with PAS usage level difference greater than 5%. We then ask the model to predict which PAS in the pair is of the higher usage level. The accuracy is defined as,

$$Comparison\ Accuracy = \frac{\#\ Pairs\ Correctly\ Predicted}{\#\ All\ Pairs}. \qquad (2)$$

Note that the primary reason that we use this metric is to compare with Polyadenylation Code and DeepPASTA, as they were designed for predicting which one is stronger between the two competing PAS.

*Highest Usage Prediction Accuracy.* We here define the Highest Usage Prediction Task. This task aims to test the model's ability of predicting which PAS is of the highest usage level in a single gene. We select all the genes which has its highest PAS usage level greater than its second highest one by at least 15% in the test set for evaluation. For DeeReCT-APA, the predicted usage in percentage is used for ranking the PAS. For Polyadenylation Code and DeepPASTA, as they do not provide a predicted value in percentage, the logit value before the SoftMax layer is used instead. The logit values, though not in the scale of real usage percentage values, can at least give a ranking of

322     different PAS sites. The highest usage prediction accuracy is the percentage of genes

323     whose highest-usage PAS are correctly predicted.

324     *Averaged Spearman's Correlation.* We here define the Ranking Task. We convert

325     the predicted usage levels by each model into a ranking of PAS sites in that gene. We

326     then compute the Spearman's correlation between the predicted ranking and ground

327     truth ranking. The correlation values for all genes are then averaged together to give an

328     aggregated score. In other words,

$$Averaged\ Spearman's\ Correlation$$

$$= \frac{1}{N}\sum_{i=1}^{N}\frac{\sum_{p=1}^{P_i}(pr_{ip} - \overline{pr}_i)(gr_{ip} - \overline{gr}_i)}{\sqrt{\sum_{p=1}^{P_i}(pr_{ip} - \overline{pr}_i)^2}\sqrt{\sum_{p=1}^{P_i}(gr_{ip} - \overline{gr}_i)^2}} \quad (3)$$

329     where $N$ is the total number of genes, $P_i$ is the number of PAS in gene $i$, $pr_{ip}$ is the

330     predicted rank of PAS $p$ in gene $i$, $gr_{ip}$ is the ground truth rank of PAS p in gene i, $pr_i$

331     and $gr_i$ are averaged predicted and ground truth ranks in gene i, respectively.

332

333     **Results**

334     **Overall performance**

335     First, to compare the performance of different Base-Net designs, we evaluated

336     DeeReCT-APA with different Base-Nets: Feature-Net, Single-Conv-Net, and Multi-

337     Conv-Net. As shown in Supplementary Table S3, both on the parental BL dataset and

338     on the F1 dataset, DeeReCT-APA with Multi-Conv-Net performs the best, followed by

339     that with Single-Conv-Net. This is expected, as deeper neural networks have higher

340     representation learning capacity.

341

342     We then compared the performance of DeeReCT-APA with Multi-Conv-Net to

343     Polyadenylation Code and DeepPASTA. As shown in **Table 1**, both on the parental BL

344     dataset and on the F1 dataset, DeeReCT-APA with Multi-Conv-Net consistently

345     performs the best across all four metrics. The standard deviation across 5-fold cross

346     validation is higher in the F1 dataset than in the parental dataset, indicating the

347     increased instability in F1 prediction which is probably due to the limited amount of F1

348     data. As we have a rather small dataset, a very complex model like DeepPASTA is

349     prone to overfitting, which is probably the reason why it performs the worst here.

350     Indeed, for the smaller F1 dataset, DeepPASTA lags even more behind other methods.

351    Similar results on the parental SP dataset and the performance of F1 model that is fine-

352    tuned from the SP parental model are shown in Supplementary Materials Section S3

353    and Supplementary Table S4. Unless otherwise stated, the F1 model that we use in the

354    remaining part of the paper is the one fine-tuned from the parental BL model and using

355    the training set folds that do not include the gene or PAS to be tested.

356

357     Next, we show that, in terms of comparison accuracy, the improvement made by

358    DeeReCT-APA is statistically significant, even though the performance improvement

359    is not numerically substantial. For this purpose, we repeat the experiment for five times,

360    with each of them having the dataset randomly split in a different way, and report the

361    accuracy   of   DeeReCT-APA   (Multi-Conv-Net),   Polyadenylation   Code,   and

362    DeepPASTA after 5-fold cross validation (Supplementary Materials Section S4 and

363    Supplementary Table S5). The performance of three tools is then compared with p-

364    value   computed   by   t-test.   As   shown   in   Supplementary   Table   S5,   indeed   the

365    improvement of DeeReCT-APA over the other two methods is statistically significant.

366

367

368     To demonstrate that the results of our comparison is independent of the datasets, we

369    train and test DeeReCT-APA also on another dataset used in [20]. Since it consists of

370    polyadenylation quantification data from multiple human tissues, we report the

371    performance (comparison accuracy) of DeeReCT-APA for each tissue separately

372    (Supplementary Materials Section S4 and Supplementary Table S6). The performance

373    metrics of Polyadenylation Code and DeepPASTA is adapted from [20] and [21]

374    accordingly. For 6 out of 8 tissues, DeeReCT-APA achieves higher accuracy than the

375    other two methods.

376

377     We finally show through ablation study that the usage of BiLSTM interaction layer

378    contributes to the performance of DeeReCT-APA. As shown in **Table 2**, we compare

379    the performance of DeeReCT-APA with Multi-Conv-Net (1) without interaction layer,

380    to (2) with interaction layer but without BiLSTM, and (3) with interaction layer and

381    with BiLSTM (The detailed architectures are shown in Supplementary Figure S1). In

382    terms of all metrics, both the usage of interaction layer and BiLSTM improve the

383    performance. Although not numerically substantial, the improvements are in general

384    statistically significant after performing a similar experiment as we have done earlier

385 (Supplementary Table S7). The improvement of (2) over (1) (p=2.5e-6 for parental and

386 p=1.1e-3 for F1) is more statistically significant than (3) over (2) (p=3.7e-3 for parental

387 and p=9.9e-2 for F1) indicating that the majority of the performance gain of DeeReCT-

388 APA comes from using the interaction layers and the simultaneous consideration of all

389 PAS. This concludes that DeeReCT-APA, with an RNN interaction layer that considers

390 all PAS of a gene at the same time, can achieve better performance on the PAS

391 quantification task.

392

393 **Benefits of modelling all PAS jointly—one example**

394     To illustrate DeeReCT-APA's ability of modeling all PAS of a gene jointly, we use

395 the gene *Srr* (Ensembl Gene ID: ENSMUSG00000001323) as an example. As shown

396 in **Figure 3A**, the gene *Srr* use four different PAS, whereas **Figure 3B, 3C, 3D** shows

397 the ground truth usage level, the prediction of DeeReCT-APA with Multi-Conv-Net

398 and Polyadenylation Code, in the F1 hybrid cell for those four PAS, for both its BL

399 allele (blue bars) and SP allele (green bars), respectively. As before, the logits values

400 before the SoftMax layer of Polyadenylation Code are used as surrogates for predicted

401 usage values (and therefore not in the range from 0 to 1). As shown in **Figure 3**, the

402 prediction of DeeReCT-APA is much more consistent with the ground truth than that

403 of Polyadenylation Code and the relative magnitude between the BL allele and SP allele

404 for the prediction of DeeReCT-APA is correct for all four PAS. In comparison,

405 Polyadenylation Code model predicted PAS 4 in the BL allele to be of slightly *higher*

406 usage than the one in the SP allele whereas both in ground truth and the prediction made

407 by DeeReCT-APA, the reverse is true. We hypothesize in this case that the genetic

408 variants between the BL allele and SP allele in the sequences flanking PAS 4 alone

409 might make the BL allele a *stronger* PAS than the SP allele because Polyadenylation

410 Code only considers which one between the two is stronger and predicts the strength of

411 a PAS solely by its own sequence, without considering those of the others. However,

412 when simultaneously considering genetic variations in PAS 1, PAS 2, and PAS 3, which

413 probably have *stronger* effects, the usage of PAS 4 becomes *lower* in BL than in SP.

414     To test our hypothesis, we design an *in-silico* experiment by constructing a

415 hypothetical allele of gene *Srr* (hereafter referred to as "mixed allele") that has the BL

416 sequence of PAS 1, PAS 2, and PAS 3, and SP sequence of PAS 4. We then ask the

417 DeeReCT-APA model to predict the usage level of each PAS in the "mixed allele",

418    where the usage differences between the BL allele and the "mixed allele" should then

419    be purely due to the sequence variants in PAS 4 because the two alleles are exactly the

420    same on the other PAS. As shown in **Figure 3E**, consistent with our hypothesis, the

421    usage level of PAS 4 in the BL allele is indeed *higher* than that in the "mixed allele".

422    This example nicely demonstrates the benefit of jointly modeling all the PAS in a gene

423    simultaneously.

424    **Allelic difference in PAS usage between BL and SP**

425    One primary goal of developing DeeReCT-APA is to determine the effect of sequence

426    variants on APA patterns. The F1 hybrid system we choose here is ideal to test how

427    well such a goal is achieved, since in the F1 cells, the allelic difference in PAS usage

428    can only be due to the sequence variants between their genome sequences.

429

430    **Figure 4** shows two examples: gene *Zfp709* (Ensembl Gene

431    ID:ENSMUSG00000056019) and *Lpar2* (Ensembl Gene ID:

432    ENSMUSG00000031861), where previous analysis demonstrated that in the distal PAS

433    of gene *Zfp709*, a substitution (from A to T) in the SP allele relative to the BL allele

434    disrupted the PAS signal (ATTAAA to TTTAAA) (**Figure 4A**); in the distal PAS of

435    gene *Lpar2*, a substitution (from A to G) in the SP allele relative to the BL allele

436    disrupted another PAS signal (AATAAA to AATAAG) (**Figure 4B**), causing both of

437    them to be of lower usage in the SP allele than in the BL allele.

438

439    To check whether our model could be used to identify the effects of these variants,

440    we plot a "mutation map" for the two genes. In brief, for each gene, given the sequence

441    around the most distal PAS (suppose it is of length L), we generate 3L "mutated

442    sequences". Each one of the 3L sequences has exactly one nucleotide mutated from the

443    original sequence. These 3L sequences are then fed into the model along with other

444    PAS sequences from that gene and the model then predicts usage for all sites and for

445    each of the 3L sequences, separately. The predicted usage values of the original

446    sequence are then subtracted from each of the 3L predictions and plotted in a heatmap,

447    the "mutation map".

448

449    As shown in **Figure 4C** and **Figure 4D**, the heatmap entries that correspond to the

450    sequence variants between BL and SP is consistent with experimental findings from

451    [31] (**Figure 4A and Figure 4B**). In addition, the mutation maps can also show the

452    predicted effect of sequence variants other than those between BL and SP, giving an

453    overview of the effects from all potential mutations.

454

455    Obviously, the two examples described above involved sequence variants disrupting

456    PAS signals, which makes the prediction relatively trivial. To check whether our model

457    could be used for the variants with more subtle effect, we choose a third example, gene

458    *Alg10b*. Previous experiments showed that the usage of the most distal PAS of its BL

459    allele is higher than its SP allele (**Figure 5A**). Using reporter assays (**Figure 5B**), it has

460    been demonstrated that [31] an insertion of UUUU in the SP allele relative to the BL

461    allele contributes to this reduction in usage (**Figure 5C**). To check whether DeeReCT-

462    APA could reveal such effects, we also construct the same four *in silico* sequences as

463    in [31] : BL, SP, BL2SP, and SP2BL. Together with other PAS of gene *Alg10b*, the

464    four sequences are feed to the DeeReCT-APA model, separately. As shown in **Figure**

465    **5D**, comparing BL with BL2SP and SP with SP2BL, our model is able to reveal the

466    negative effect of poly(U) tract.

467

468    To globally evaluate the performance of DeeReCT-APA on predicting the allelic

469    difference in PAS usage, we compare the predicted allelic difference versus

470    experimentally measured allelic difference in a genome-wide manner (**Figure 6A**). As

471    a baseline control, we do the same for the prediction made by the Polyadenylation Code

472    where logit values before SoftMax are again used as surrogates for the predicted allelic

473    difference in PAS usage (**Figure 6B**). Here, the F1 model fine-tuned from the BL

474    parental model is used. Similar results of the F1 model fine-tuned from the SP parental

475    model are shown in Supplementary Materials Section S3 and Supplementary Figure S2.

476    It is worth noting that this is a very challenging task because the training data do not

477    well represent the complete landscape of genetic mutations. That is, the BL dataset only

478    contains invariant sequences from different PAS, and the F1 dataset contains a limited

479    number of genetic variants.

480

481    We then compute the Pearson correlation between the experimentally measured

482    allelic usage difference and the ones predicted by the two models. Clearly, DeeReCT-

483    APA outperforms Polyadenylation Code. We further evaluate the Pearson correlation

484    values using six subsets of the test set, each filtering out PAS with allelic usage

485    difference less than 0.1, 0.2, 0.3, 0.4, 0.6, 0.8, respectively (**Figure 6, Panel C**). When

486    the allelic usage difference is small, their relative magnitudes are more ambiguous and

487    the experimental measurement of their allelic usage difference (used here as ground

488    truth) are less confident. Indeed, with the increasing allelic difference, the prediction

489    accuracy increased for both DeeReCT-APA and Polyadenylation Code. Importantly, in

490    all these groups, DeeReCT-APA shows consistently better performance.

491    **Visualization of convolutional filters**

492    To show the knowledge learned by the convolutional filters of DeeReCT-APA, we

493    follow a similar procedure as in [36] to visualize the convolutional filters of the model.

494    The aim of visualization is to reveal the important subsequences around

495    polyadenylation sites that activate a specific convolutional filter. In contrast to [38], in

496    which the researchers only used sequences in the test set for visualization, we use all

497    sequences in the train and test dataset of F1 for visualization due to the smaller size of

498    our dataset. In visualization, neither the model parameters nor the hyperparameters are

499    tuned on the test set, our usage of test set for visualization is therefore legitimate. For

500    all the learned filters in layer 1, we convolve them with all the sequences in the above

501    dataset, and for each sequence, its subsequence (having the same size as the filters) with

502    the highest activation on that filter is extracted and accumulated in a position frequency

503    matrix (PFM). The PFM is then ready for visualization as the knowledge learned by

504    that specific filter. For layer 2 convolutional filters, as they do not convolve with raw

505    sequences during training and testing, directly convolving it with the sequences in the

506    dataset as we did for layer 1 would be undesirable. Instead, the layer 2 activations are

507    calculated by a partial forward pass in the network and the subsequences of the input

508    sequences in the receptive field of the maximally-activated neuron is extracted and

509    accumulated in a PFM.

510    As shown in **Figure 7A** and **7B**, DeeReCT-APA is able to identify the two strongest

511    PAS hexmer, AUUAAA and AAUAAA [31]. In addition, one of the layer 2

512    convolutional filters is able to recognize the pattern of a mouse specific PAS hexamer

513    UUUAAA [30] (**Figure 7C**). Furthermore, a Poly-U island motif previously reported

514    in [38] could also be revealed by DeeReCT-APA (**Figure 7D**). A complete

515    visualization of all the 40 filters in layer 1 and 40 filters in layer 2 is shown in

516    Supplementary Figure S3 and Supplementary Figure S4.

517

## Discussion and conclusion

519 In this study, we made the first attempt to simultaneously predict the usage of all
520 competing PAS within a gene. Our method incorporates both sequence-specific
521 information through automatic feature extraction by CNN and multiple PAS
522 competition through interaction modeling by RNN. We trained and evaluated our
523 model on the genome-wide PAS usage measurement obtained from 3'-mRNA
524 sequencing of fibroblast cells from two mouse strains as well as their F1 hybrid. Our
525 model, DeeReCT-APA, outperforms the state-of-the-art PAS quantification methods
526 on the tasks that they are trained for, including pairwise comparison, highest usage
527 prediction and ranking task. In addition, we demonstrated that modeling all the PAS of
528 a gene simultaneously captures the mechanistic competition among the PAS and
529 reveals the genetic variants with regulatory effects on PAS usage.

530

531 A similar idea of using BiLSTM to model competitive biological processes was
532 proposed recently in [39]. The researchers used BiLSTM to model the usage level of
533 competitive alternative 5'/3' splice sites. Given the similarity of modeling competing
534 polyadenylation sites and splice sites, it is therefore not surprising that DeeReCT-APA,
535 which also incorporates BiLSTM to model the interactions among competing
536 polyadenylation sites, achieves the best performance on the PAS quantification task.

537

538 Although DeeReCT-APA provides the first-of-its-kind method to model all the PAS
539 of a gene, it still has room for improvement. As shown in Figure 3, the model has limited
540 accuracy when the usage is very high or very low (comparing Figure 3B and Figure
541 3C). In addition, for allelic comparison as shown in Figure 5, some PAS with high
542 allelic usage difference are predicted to be of low difference (false negatives, along X
543 axis) and vice versa (false positives, along Y axis). One of the main reasons for our
544 model's limitation, as well as for all the other PAS quantification methods, is that all
545 the existing genome-wide PAS quantification datasets used as training data could only
546 sample the limited number of naturally occurring sequence variants. Although in our
547 study the two parental strains from which the F1 hybrid mouse was derived are already
548 the evolutionarily most distant ones among all the 17 mouse strains with complete
549 genomic sequences, the number of genetic variants is still rather limited. To address
550 this limitation and provide a complementary dataset, we are working on establishing a

551    large-scale synthetic APA mini-gene reporter-based system which samples the

552    regulatory effect of millions of random sequences (manuscript in preparation). Another

553    limitation of our current model lies in the fact that it does not take all the factors with

554    potential PAS regulatory effects into consideration. For example, transcription kinetics,

555    i.e., the elongation rate of Pol II, which is not considered by the model in this study,

556    can also affect APA choice [40]. Similarly, DeeReCT-APA does not take the distance

557    between consecutive PAS into account, which, together with the transcription

558    elongation rate, can also affect APA [41]. All of them are potential directions for further

559    improvement.

560

561    Finally, very recently, Zhang et al. showed that effectively combining the power of

562    deep learning and the information in RNA-seq data can significantly boost the

563    performance for investigating the pattern of alternative splicing [42]. Indeed, our

564    preliminary results showed that also for the recognition of APA patterns, there are

565    substantial cases in which deep learning cannot make accurate prediction but utilizing

566    the pattern of RNA-seq coverage around the cleavage site could provide clear evidence,

567    and vice versa. Future work integrating the strength of deep learning on genomic

568    sequences and experimental RNA-seq data will for certain not only improve the model

569    performance, but also shed more light on the APA regulatory mechanisms.

570

## Data Availability

Our implementation of DeeReCT-APA using the PyTorch [37] library is available at the repository (https://github.com/lzx325/DeeReCT-APA-repo). The genome-wide PAS quantification dataset of parental and F1 mouse fibroblast cell is available in the subfolder `APA_ML`. As provided in [31], the raw sequencing data from which this dataset is derived is accessible at European Nucleotide Archive (http://www.ebi.ac.uk/ena) under the accession number PRJEB15336 (URL: https://www.ebi.ac.uk/ena/browser/view/PRJEB15336).

## Authors' contributions

ZL, YH, WC, and XG conceived the project. ZL developed the deep learning model and did the computational experiments. Yisheng Li and BZ provided and pre-processed the dataset. JZ, XZ, and MZ provided additional biological insights on the experimental results. ZL, YH, WC, and XG drafted the paper. ZL, Yisheng Li, BZ, Yu Li, Yongkang Long, JZ, XZ, MZ, YH, WC, and XG read and approved the final manuscript.

## Competing interests

The authors have declared no competing interests.

## Acknowledgements

600

601

# References

603

[1] Barreau C, Paillard L, Osborne HB. AU-rich elements and associated factors: are there unifying principles? Nucleic Acids Res. 2005;33:7138-50.

[2] Chen CY, Shyu AB. AU-rich elements: characterization and importance in mRNA degradation. Trends Biochem Sci 1995;20:465–70.

[3] Jonas S, Izaurralde E. Towards a molecular understanding of microRNA-mediated gene silencing. Nat Rev Genet 2015;16:421–33.

[4] Lau AG, Irier HA, Gu J, Tian D, Ku L, Liu G, et al. Distinct 3'UTRs differentially regulate activity-dependent translation of brain-derived neurotrophic factor (BDNF). Proc Natl Acad Sci U S A 2010;107:15945–50.

[5] Bertrand E, Chartrand P, Schaefer M, Shenoy SM, Singer RH, Long RM. Localization of ASH1 mRNA particles in living yeast. Mol. Cell 1998;2:437–45.

[6] Ephrussi A, Dickinson LK, Lehmann R. Oskar organizes the germ plasm and directs localization of the posterior determinant nanos. Cell 1991;66:37–50.

[7] Niedner A, Edelmann FT, Niessing D. Of social molecules: The interactive assembly of ASH1 mRNA-transport complexes in yeast. RNA Biol 2014;11:998–1009.

[8] Berkovits BD, Mayr C. Alternative 3' UTRs act as scaffolds to regulate membrane protein localization. Nature 2015;522:363–7.

[9] Yasuda M, Shabbeer J, Osawa M, Desnick RJ. Fabry disease: novel alpha-galactosidase A 3'-terminal mutations result in multiple transcripts due to aberrant 3'-end formation. Am J Hum Genet 2003;73:162–73.

[10] Bennett CL, Brunkow ME, Ramsdell F, O'Briant KC, Zhu Q, Fuleihan RL, et al. A rare polyadenylation signal mutation of the FOXP3 gene (AAUAAA-->AAUGAA) leads to the IPEX syndrome. Immunogenetics 2001;53:435–9.

[11] Higgs DR, Goodbourn SE, Lamb J, Clegg JB, Weatherall DJ, Proudfoot NJ. Alpha-thalassaemia caused by a polyadenylation signal mutation. Nature 1983;306:398–400.

[12] Orkin SH, Cheng TC, Antonarakis SE, Kazazian HH, Jr. Thalassemia due to a mutation in the cleavage-polyadenylation signal of the human beta-globin gene. Embo j 1985;4:453–6.

[13] Elkon R, Ugalde AP, Agami R. Alternative cleavage and polyadenylation: extent, regulation and function. Nat Rev Genet 2013;14:496-506.

[14] Mandel CR, Bai Y, Tong L. Protein factors in pre-mRNA 3′-end processing. Cellular and Molecular Life Sciences 2008;65:1099–122.

[15] Shi Y. Alternative polyadenylation: new insights from global analyses. RNA 2012;18:2105-17.

[16] Hoque M, Ji Z, Zheng D, Luo W, Li W, You B, et al. Analysis of alternative cleavage and polyadenylation by 3' region extraction and deep sequencing. Nat Methods 2013;10:133–9.

[17] Kalkatawi M, Rangkuti F, Schramm M, Jankovic BR, Kamau A, Chowdhary R, et al. Dragon PolyA Spotter: predictor of poly(A) motifs within human genomic DNA sequences. Bioinformatics (Oxford, England) 2012;28:127–9.

[18] Magana-Mora A, Kalkatawi M, Bajic VB. Omni-PolyA: a method and tool for accurate recognition of Poly(A) signals in human genomic DNA. BMC Genomics 2017;18:620.

648  [19] Xie B, Jankovic BR, Bajic VB, Song L, Gao X. Poly(A) motif prediction using
649  spectral latent features from human DNA sequences. Bioinformatics 2013;29:i316–25.
650  [20] Leung MKK, Delong A, Frey BJ. Inference of the human polyadenylation code.
651  Bioinformatics 2018;34:2889–98.
652  [21] Arefeen A, Xiao X, Jiang T. DeepPASTA: deep neural network based
653  polyadenylation site analysis. Bioinformatics 2019;35:4577–85.
654  [22] LeCun Y, Bengio Y, Hinton G. Deep learning. Nature 2015;521:436–44.
655  [23] Alipanahi B, Delong A, Weirauch MT, Frey BJ. Predicting the sequence
656  specificities of DNA- and RNA-binding proteins by deep learning. Nat Biotech
657  2015;33:831–8.
658  [24] Leung MK, Xiong HY, Lee LJ, Frey BJ. Deep learning of the tissue-regulated
659  splicing code. Bioinformatics 2014;30:i121–9.
660  [25] Li Y, Wang S, Umarov R, Xie B, Fan M, Li L, et al. DEEPre: sequence-based
661  enzyme EC number prediction by deep learning. Bioinformatics 2018;34:760–9.
662  [26] Zou Z, Tian S, Gao X, Li Y. mlDEEPre: Multi-Functional Enzyme Function
663  Prediction With Hierarchical Multi-Label Deep Learning. Frontiers in Genetics 2019;9.
664  [27] Han R, Li Y, Wang S, Gao X, Bi C, Li M. DeepSimulator: a deep simulator for
665  Nanopore sequencing. Bioinformatics 2018;34:2899–908.
666  [28] Wang S, Li Z, Yu Y, Gao X. WaveNano: a signal-level nanopore base-caller via
667  simultaneous prediction of nucleotide labels and move labels through bi-directional
668  WaveNets. Quantitative Biology 2018;6:359–68.
669  [29] Umarov R, Kuwahara H, Li Y, Gao X, Solovyev V. Promoter analysis and
670  prediction in the human genome using sequence-based deep learning models.
671  Bioinformatics 2019;35:2730–7.
672  [30] Xia Z, Li Y, Zhang B, Li Z, Hu Y, Chen W, et al. DeeReCT-PolyA: a robust and
673  generic deep learning method for PAS identification. Bioinformatics 2019;35:2371–9.
674  [31] Xiao MS, Zhang B, Li YS, Gao Q, Sun W, Chen W. Global analysis of regulatory
675  divergence in the evolution of mouse alternative polyadenylation. Molecular Systems
676  Biology 2016;12:890.
677  [32] Hochreiter S, Schmidhuber J. Long Short-Term Memory. Neural Computation
678  1997;9:1735–80.
679  [33] Hu J, Lutz CS, Wilusz J, Tian B. Bioinformatic identification of candidate cis-
680  regulatory elements involved in human mRNA polyadenylation. RNA 2005;11:1485–
681  93.
682  [34] Gao Q, Sun W, Ballegeer M, Libert C, Chen W. Predominant contribution of cis-
683  regulatory divergence in the evolution of mouse alternative splicing. Mol Syst Biol
684  2015;11:816.
685  [35] Hou J, Wang X, McShane E, Zauber H, Sun W, Selbach M, et al. Extensive allele-
686  specific translational regulation in hybrid mice. Mol Syst Biol 2015;11:825.
687  [36] Kingma DP, Ba J. Adam: A Method for Stochastic Optimization. arXiv:1412.6980
688  2014.
689  [37] Paszke A, Gross S, Massa F, Lerer A, Bradbury J, Chanan G, et al. PyTorch: An
690  imperative style, high-performance deep learning library. Advances in Neural
691  Information Processing Systems 2019:8024–35.
692  [38] Bogard N, Linder J, Rosenberg AB, Seelig G. A Deep Neural Network for
693  Predicting and Engineering Alternative Polyadenylation. Cell 2019;178:91–106.e23.
694  [39] Zuberi K, Gandhi S, Bretschneider H, Frey BJ, Deshwar AG. COSSMO:
695  predicting competitive alternative splice site selection using deep learning.
696  Bioinformatics 2018;34:i429–i37.

697  [40] Pinto PAB, Henriques T, Freitas MO, Martins T, Domingues RG, Wyrzykowska
698  PS, et al. RNA polymerase II kinetics in polo polyadenylation signal selection. The
699  EMBO journal 2011;30:2431–44.
700  [41] Gruber AJ, Zavolan M. Alternative cleavage and polyadenylation in health and
701  disease. Nature Reviews Genetics 2019;20:599–614.
702  [42] Zhang Z, Pan Z, Ying Y, Xie Z, Adhikari S, Phillips J, et al. Deep-learning
703  augmented RNA-seq analysis of transcript splicing. Nature Methods 2019;16:307–10.

704

705

706

707  **Figure legends**

708  **Figure 1 Illustration of the DeeReCT-APA architecture (Using BiLSTM as**

709  **interaction layer)**

710  **Figure 2 Three designs of Base-Net.**

711  All three of them output a feature vector that represents distilled features of the input

712  sequence. **A.** Single-Conv-Net uses a single convolution layer for feature extraction. **B.**

713  Multi-Conv-Net uses multiple convolution layers for feature extraction. **C.** Feature-Net

714  contains a hand-crafted feature extractor before being processed by fully-connected

715  layers.

716  **Figure 3 Prediction of gene *Srr***

717  This shows one example of the benefit of modelling all PAS jointly. Each panel shows

718  the predicted or ground truth usage of each of its four PAS: **A.**  PAS of gene *Srr*. **B.**

719  Ground Truth. **C.** DeeReCT-APA's (Multi-Conv-Net) prediction. **D.** Polyadenylation

720  Code's prediction. **E.** DeeReCT-APA's (Multi-Conv-Net) prediction of "mixed allele".

721  The prediction of DeeReCT-APA is much more consistent with ground truth compared

722  to Polyadenylation Code. Especially for PAS 4, DeeReCT-APA predicts the one of BL

723  allele to be of lower usage than the one of SP allele which is consistent with ground

724  truth. Polyadenylation Code, on the contrary, predicts the opposite. In Panel **E**, by

725  making prediction of the "Mixed Allele", we demonstrated that the increased usage of

726  PAS in SP allele is probably due to the concerted effects of the other three PAS.

727   **Figure 4 Previous experimental findings and mutation map of gene *Zfp709* and**

728   ***Lpar2***

729   Mutation map is consistent with previous experimental findings on two genes, *Zfp709*

730   (**A** & **C**) and *Lpar2* (**B** & **D**). Sequencing read coverage graphs (**A** & **B**) are adapted

731   from Figure 4H of [31]. The identified PAS are marked by red triangles on top of the

732   sequencing read coverage (black coverage graph). The sequence variants of the PAS

733   shaded in pink between BL and SP strains are shown on the top. The BL mutation map

734   (**C** & **D**) of the BL distal PAS sequence shows the effect of BL distal sequence mutation

735   on the usage of distal sites. The SP gene *Zfp709* and *Lpar2* can be viewed as undergoing

736   a substitution relative to BL. The four heatmap entries above each letter of the sequence

737   (**C** & **D**, bottom) show the relative change of usage level when the nucleotide at that

738   position is substituted with the nucleotide of the corresponding row. Darker red

739   indicates greater increase in usage and darker blue indicates more decrease in usage.

740   The entries that correspond to the genetic variants between BL and SP in **A** & **B** are

741   marked by red squares.

742   **Figure 5 Previous experimental findings and DeeReCT-APA's prediction of gene**

743   ***Alg10b***

744   *In silico* prediction for the *Alg10b* PAS reporter is consistent with previous

745   experimental findings. Similar to Figure 4A, the sequencing read coverage graph and

746   the sequence variants are shown in **A.** The red triangles mark the identified PAS sites.

747   The structures of PAS reporter constructs are shown in **B**, where "BL" is the original

748   BL version of the most distal PAS, "SP" is the original SP version, "BL2SP" is the BL

749   sequence only inserted with TTTT at the corresponding location and "SP2BL" is the

750   SP sequence only deleted with TTTT at the corresponding location. The experimental

751   results from PAS reporter assay for the four reporters are shown in **C.** and their *in silico*

752   predictions are shown in **D.** Considering the *in silico* prediction pairs, BL & BL2SP

753   and SP & SP2BL, it is clear that DeeReCT-APA is able to identify the negative

754   modulation of PAS usage by the poly(U) tract. Figure (**A**, **B** &**C**) are adapted from

755   Figure 4H of [31]. See text for more details.

756

757 **Figure 6 Comparison of the allelic usage difference predicted by DeeReCT-APA**

758 **and Polyadenylation Code**

759 F1 model fine-tuned from BL parental model is used. **A. B.** The horizontal axis is the

760 ground truth allelic usage difference (BL usage minus the SP usage). The vertical axis

761 shows the predicted allelic usage difference. The red line shows the perfect prediction.

762 In terms of Person correlation, DeeReCT-APA shows better correlation than

763 Polyadenylation Code. **C.** Pearson correlations (and their p-values) between two

764 quantities at different minimum allelic usage difference are shown in the table below.

765 The prediction of DeeReCT-APA still has better correlation than Polyadenylation Code

766 when the dataset is filtered at different thresholds.

767

768 **Figure 7 Visualization of learned convolutional filters in DeeReCT-APA**

769 Some visualization examples of the learned convolutional filters of DeeReCT-APA. **A.**

770 **B.** The most common polyadenylation motifs AUUAAA and AAUAAA are learned in

771 layer 1 convolutional filter #2 and #37, respectively. **C.** Visualization of a layer 2 filter,

772 #38 shows a mouse specific polyadenylation motif UUUAAA. **D.** Layer 2 filter #19

773 shows the Poly-U islands on polyadenylation. Note that the layer 2 filter visualization

774 PFMs are wider than the layer 2 filter (12nt) because the receptive field of neurons in a

775 deeper layer is in general greater than their corresponding filter width.

776

777 **Tables**

778 **Table 1 Performance summary for the BL parental model and the F1 model.**

779

780 **Table 2 The performance of DeeReCT-APA using different interaction layers**

781

782 **Supplementary material**

783 **Supplementary File**

784     DeeReCT-APA-Supplementary-File.pdf

785 **Figure S1 The structures of DeeReCT-APA models used in the ablation study.**

786     **A.** The structure of DeeReCT-APA with interaction layers but without BiLSTM. **B.**

787 The structure of DeeReCT-APA with interaction layers removed. Comparing **A** with

788 Figure 1 in the main text, it has BiLSTM removed and only has the affine layer in the

789 interaction layers. In **B**, the interaction layers are removed altogether and DeeReCT-

790 APA resorted to comparison-based training (to predict which one of the two PAS is of

791 higher usage). Note that an additional affine layer is added on top of the Base Networks

792 to cast the output of the base network (which is a vector) into a scalar.

793 **Figure S2 Comparison of the allelic usage difference prediction of DeeReCT-**

794 **APA and Polyadenylation Code.**

795     F1 model fine-tuned from SP parental model is used. **A. B.** The horizontal axis is

796 the ground truth allelic usage value difference between two homologous PAS (which

797 is the BL usage value minus the SP usage value). The vertical axis shows the predicted

798 allelic usage value difference. The scatter plot of DeeReCT-APA is shown in Panel **A**

799 and Polyadenylation Code is shown in Panel **B**. As DeeReCT-APA predicts the usage

800 value in percentage, we draw a red line that shows the perfect prediction. **C.** Pearson

801 correlations between two quantities at different minimum allelic usage difference are

802 shown in the table below.

803 **Figure S3 Visualization of convolutional filters in layer 1 of DeeReCT-APA.**

804     There are 40 convolutional filters in layer 1 of DeeReCT-APA. The model is trained

805 on parental BL dataset and fine-tuned on F1.

806 **Figure S4 Visualization of convolutional filters in layer 2 of DeeReCT-APA.**

807     There are 40 convolutional filters in layer 2 of DeeReCT-APA. The model is trained

808 on parental BL dataset and fine-tuned on F1.

809 **Table S1 List of features used in Feature-Net and their corresponding**
810 **dimensions.**

811 **Table S2 List of hyperparameters for the three DeeReCT-APA models.**

812 **Table S3 Performance summary for the BL parental model and the F1 model**
813 **fine-tuned from the BL parental model.**

814 **Table S4 Performance summary for the SP parental model and the F1 model**
815 **fine-tuned from the SP parental model.**

816 **Table S5 Replicated Experiments of 5-fold cross validation on 5 random splits.**

817 **Table S6 Comparison accuracy on dataset from [20]**

818 **Table S7 Replicated Experiments of ablation study.**

819

Predicted Usage Level

Interaction Layers

LSTM          LSTM          LSTM

Base Networks

TGGGATACAACCTCT          ATACTACACATAACA          GGGCATTACCAGAA

Polyadenylation Sites

A

Fully Connected
Layer

Flatten

Convolution

One-hot
encoding

TGGGATACAACCTCT

B

Fully Connected
Layer

Flatten

Convolution

Convolution

One-hot
encoding

TGGGATACAACCTCT

C

Feature Extractor

TGGGATACAACCTCT

Fully Connected
Layer

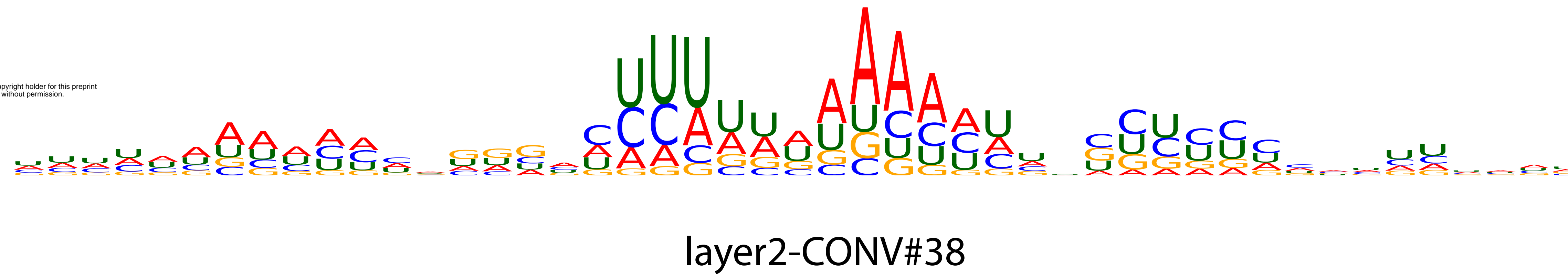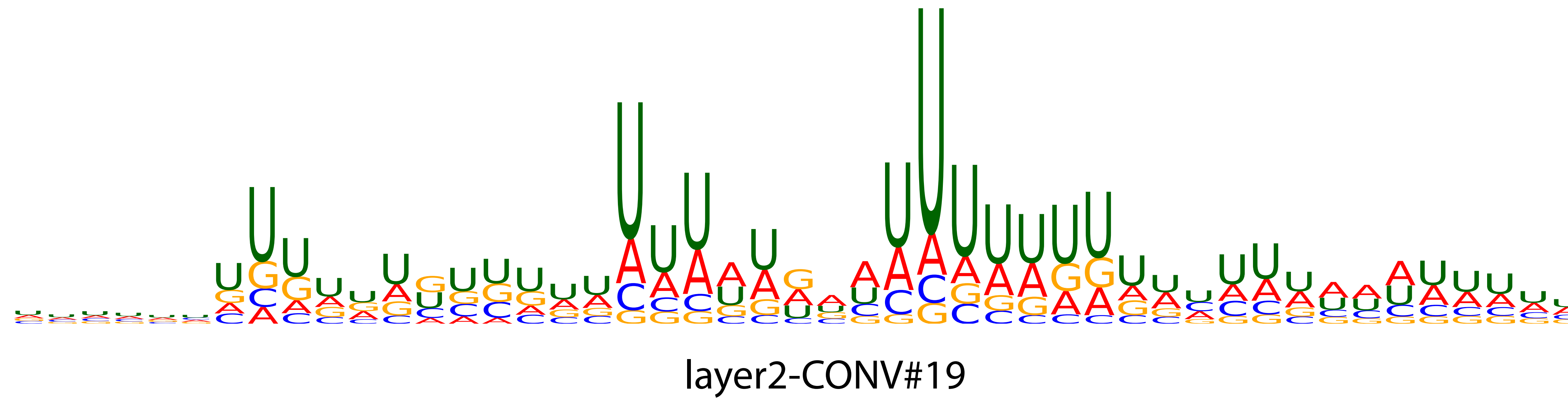| Min. Allelic Usage Difference | | 0.1 | 0.2 | 0.3 | 0.4 | 0.6 | 0.8 |
|---|---|---|---|---|---|---|---|
| DeeReCT-APA | Pearson R | 0.419 | 0.486 | 0.522 | 0.543 | 0.657 | 0.744 |
| | p value | $1.9 \times 10^{-93}$ | $1.2 \times 10^{-67}$ | $4.2 \times 10^{-45}$ | $5.0 \times 10^{-30}$ | $7.9 \times 10^{-16}$ | $1.5 \times 10^{-4}$ |
| Polyadenylation Code | Pearson R | 0.316 | 0.360 | 0.401 | 0.433 | 0.502 | 0.543 |
| | p value | $1.2 \times 10^{-53}$ | $5.5 \times 10^{-37}$ | $1.2 \times 10^{-22}$ | $1.9 \times 10^{-16}$ | $6.0 \times 10^{-8}$ | $4.7 \times 10^{-3}$ |

A

B

C    layer1-CONV#2    layer1-CONV#37

layer2-CONV#38

D

layer2-CONV#19

**Table 1 Performance summary for the BL parental model and the F1 model.**

A

| Model | Performance on Parental Dataset | | | |
|---|---|---|---|---|
| | MAE[*] | Comparison Accuracy[*] | Highest Usage Prediction Accuracy[*] | Averaged Spearman's Correlation |
| DeeReCT-APA (Multi-Conv-Net) | $\mathbf{17.22\% \pm 0.3\%}$ | $\mathbf{77.64\% \pm 0.4\%}$ | $\mathbf{63.48\% \pm 0.9\%}$ | $\mathbf{0.5140 \pm 0.021}$ |
| Polyadenylation Code | N/A | $75.88\% \pm 0.8\%$ | $59.82\% \pm 1.5\%$ | $0.4673 \pm 0.022$ |
| DeepPASTA | N/A | $74.08\% \pm 1.1\%$ | $58.78\% \pm 1.4\%$ | $0.4394 \pm 0.017$ |

*The values for a random predictor are 43.12%, 50.00% and 25.49% respectively. MAE, mean absolute error. Note that for MAE, it is the *lower* the better.

B

| Model | Performance on F1 Dataset | | | |
|---|---|---|---|---|
| | MAE[*] | Comparison Accuracy[*] | Highest Usage Prediction Accuracy[*] | Averaged Spearman's Correlation |
| DeeReCT-APA (Multi-Conv-Net) | $\mathbf{17.80\% \pm 0.3\%}$ | $\mathbf{77.14\% \pm 1.2\%}$ | $\mathbf{64.52\% \pm 0.7\%}$ | $\mathbf{0.4567 \pm 0.009}$ |
| Polyadenylation Code | N/A | $74.20\% \pm 0.1\%$ | $59.04\% \pm 0.9\%$ | $0.4224 \pm 0.014$ |
| DeepPASTA | N/A | $70.14\% \pm 1.5\%$ | $53.82\% \pm 1.7\%$ | $0.3693 \pm 0.018$ |

*The values for a random predictor are 40.96%, 50.00% and 28.56% respectively. MAE, mean absolute error. Note that for MAE, it is the *lower* the better.

*Note:* The parental model is trained from scratch and the F1 model is fine-tuned from the BL parental model. The table shows the performance of three models across four evaluation metrics. Results are shown in the mean±std format. The best performance is in bold. See Section "Overall performance" for details. A. Performance on the Parental Dataset (BL). B. Performance on the F1 Dataset (fine-tuned from parental BL model).

**Table 2 The performance of DeeReCT-APA using different interaction layers**

A

| Model | Performance on Parental Dataset | | | |
|---|---|---|---|---|
| | MAE[*] | Comparison Accuracy[*] | Highest Usage Prediction Accuracy[*] | Averaged Spearman's Correlation |
| DeeReCT-APA (Multi-Conv-Net) (No Interaction Layer) | - | 76.12% ± 0.5% | 60.02% ± 0.7% | 0.4988 ± 0.027 |
| DeeReCT-APA (Multi-Conv-Net) (w/o BiLSTM) | 17.54% ± 0.3% | 77.12% ±0.5% | 61.73% ± 0.6% | 0.5007 ±0.034 |
| DeeReCT-APA (Multi-Conv-Net) (BiLSTM) | **17.22% ± 0.3%** | **77.64% ± 0.4%** | **63.48% ± 0.9%** | **0.5140 ± 0.021** |

[*] The values for a random predict or are 43.12%, 50.00% and 25.49% respectively. MAE, mean absolute error. Note that for MAE, it is the lower the better.

B

| Model | Performance on F1 Dataset | | | |
|---|---|---|---|---|
| | MAE[*] | Comparison Accuracy[*] | Highest Usage Prediction Accuracy[*] | Averaged Spearman's Correlation |
| DeeReCT-APA (Multi-Conv-Net) (No Interaction Layer) | - | 76.28% ± 1.1% | 61.72% ± 0.8% | 0.4337 ± 0.019 |
| DeeReCT-APA (Multi-Conv-Net) (w/o BiLSTM) | 18.03% ± 0.2% | 76.77% ± 1.0% | 63.44% ± 0.3% | 0.4751 ± 0.011 |
| DeeReCT-APA (Multi-Conv-Net) (BiLSTM) | **17.80% ± 0.4%** | **77.14% ± 1.2%** | **64.52% ± 0.7%** | **0.4957 ± 0.009** |

[*]The values for a random predictor are 40.96%, 50.00% and 28.56% respectively. MAE, mean absolute error. Note that for MAE, it is the lower the better.

*Note:* The table shows the performance of DeeReCT-APA with different interaction layers. Note that for DeeReCT-APA without interaction layer, the model is trained based on comparison and its output cannot

be interpreted as a percentage score. Therefore, like for Polyadenylation Code and DeepPASTA earlier, we do not report its MAE value. A. Performance on the Parental Dataset (BL). B. Performance on the F1 Dataset (fine-tuned from parental BL model).