

METHOD

Open Access



# MTSplice predicts effects of genetic variants on tissue-specific splicing

Jun Cheng<sup>1\*</sup>, Muhammed Hasan Çelik<sup>1</sup>, Anshul Kundaje<sup>2,3</sup> and Julien Gagneur<sup>1,4,5\*</sup> 

\*Correspondence:  
s6juncheng@gmail.com;  
gagneur@in.tum.de

<sup>1</sup>Department of Informatics,  
Technical University of Munich,  
Boltzmannstraße, 85748 Garching,  
Germany

<sup>4</sup>Institute of Computational Biology,  
Helmholtz Zentrum München,  
Neuherberg, Germany

Full list of author information is  
available at the end of the article

## Abstract

We develop the free and open-source model Multi-tissue Splicing (MTSplice) to predict the effects of genetic variants on splicing of cassette exons in 56 human tissues. MTSplice combines MMSplice, which models constitutive regulatory sequences, with a new neural network that models tissue-specific regulatory sequences. MTSplice outperforms MMSplice on predicting tissue-specific variations associated with genetic variants in most tissues of the GTEx dataset, with largest improvements on brain tissues. Furthermore, MTSplice predicts that autism-associated de novo mutations are enriched for variants affecting splicing specifically in the brain. We foresee that MTSplice will aid interpreting variants associated with tissue-specific disorders.

## Introduction

Splicing is a fundamental biological process in which introns are cut out from precursor RNAs and exons are joined together. Alternative splicing refers to alternative usage of exons. It is estimated that approximately 95% of human multi-exon genes undergo alternative splicing [1]. Exon skipping (of so-called cassette exons) is the most common alternative splicing pattern [2]. Skipping level of an exon is commonly quantified with the percent spliced-in (PSI or  $\Psi$ ) [3]. Percent spliced-in can be estimated from RNA-sequencing (RNA-Seq) data as the number of split RNA-Seq reads supporting the inclusion of the exon divided by the total number of split reads supporting the skipping or the inclusion of the exon. Splicing is a complex process which involves regulation by sequence elements in the exons and flanking introns [4, 5]. Moreover, alternative splicing is often tissue-specific [2, 3, 6, 7]. This means that certain splicing isoforms are only present in certain tissues or that the relative abundances of splice isoforms differ across tissues. Alternative splicing plays an important role in tissue development and shaping tissue identity [8, 9]. Analyzing the protein-coding roles of tissue-specific exons revealed their critical role in rewiring protein interaction networks in different tissues [10]. Tissue-specific splicing patterns are associated with short RNA motifs [2, 11–14]. These short RNA motifs encode tissue-specific splicing regulatory elements, typically intronic or



exonic binding sites for splicing factors with a tissue-specific activity. Mammalian tissue-specific splicing factors include Nova1, Nova2, PTB/nPTB, and RBFOX1 for nervous tissues, and MBNL1 for muscles, among others. For a review, see Chen and Manley [15].

Splicing defects account for an important fraction of the genetic basis of human diseases [16–18]. Some of these splicing defects are specific to disease-relevant tissues. For instance, individuals affected by autism spectrum disorder (ASD) frequently present missplicing of brain-specific exons [19–21] as well as an enrichment of de novo mutations in brain-specific exons [22]. Hence, computational tools that can predict the tissue-specific effects of genetic variants on splicing would be relevant for understanding the genetic basis of tissue-specific diseases such as ASD.

Many computational tools have been developed to predict splice sites or splicing strength from sequence [23–33]. However, tools are lacking for predicting tissue-specific effects of human genetic variants on splicing. Barash et al. developed the first sequence-based model predicting tissue-specific splicing in mouse cells [34]. The model integrates regulatory sequence elements to qualitatively predict whether the inclusion of a cassette exon increases, decreases, or remains at a similar level from one tissue to another tissue. This model was further improved to predict directional changes between tissues along with discretized  $\Psi$  categories (low, medium, and high) within a tissue by using a Bayesian neural network with hidden variables [35]. In a subsequent study, a similar Bayesian neural network (SPANR) was trained on human data [29]. However, SPANR was evaluated only for predicting the largest effect across all investigated tissues. Hence, the performance of SPANR on any given tissue is unclear. Moreover, the publicly available SPANR does not allow performing tissue-specific predictions.

We previously developed MMSplice, a neural network with a modular design that predicts the effect of variants on splicing [30, 31]. Unlike SPANR, which has been trained on natural endogenous genomic sequence, MMSplice leverages perturbation data from a recently published massively parallel reporter assay [28]. MMSplice outperformed SPANR and many other splicing predictors in predicting  $\Psi$  variations associated with naturally occurring genetic variants as well as effects of variants on percent spliced-in measured on reporter assays [30, 36]. MMSplice models the odds ratio of a cassette exon to be spliced-in when comparing an alternative sequence to a reference sequence. The predicted odds ratios are the same for all tissues because MMSplice has been trained in a tissue-agnostic fashion and therefore does not capture effects of variants affecting tissue-specific regulatory elements.

Deep learning models of tissue-specific regulatory elements have been developed for other biological processes. These models include DeepSEA for chromatin-profiles [37], Basset for DNase I hypersensitivity [38], ExPecto for tissue-specific gene expression [39], FactorNet for transcription factor binding [40], and ChromDragoNN for chromatin accessibility [41]. A common denominator of these models is that they are trained by multi-task learning, i.e., the models make joint predictions for all tissues or cell types using a common set of underlying predictive features. This strategy allows models to efficiently pool information about regulatory elements that are shared across cell types or tissues.

Here, we developed MTSplice (Multi-tissue Splicing), a model that predicts tissue-specific splicing effects of human genetic variants. MTSplice adjusts the MMSplice predictions with the predictions of TSplice (Tissue-specific Splicing), a novel deep neural network predicting tissue-specific variations of  $\Psi$  from sequence which we trained on 56

human tissues using multi-task learning. Performance of MTSplice is demonstrated by predicting tissue-specific variations of  $\Psi$  associated with naturally occurring genetic variants of the GTEx dataset as well as investigating brain-specific splicing effect predictions for autism-associated variants. MTSplice is open-source and freely available at the model repository Kipoi [42].

## Results

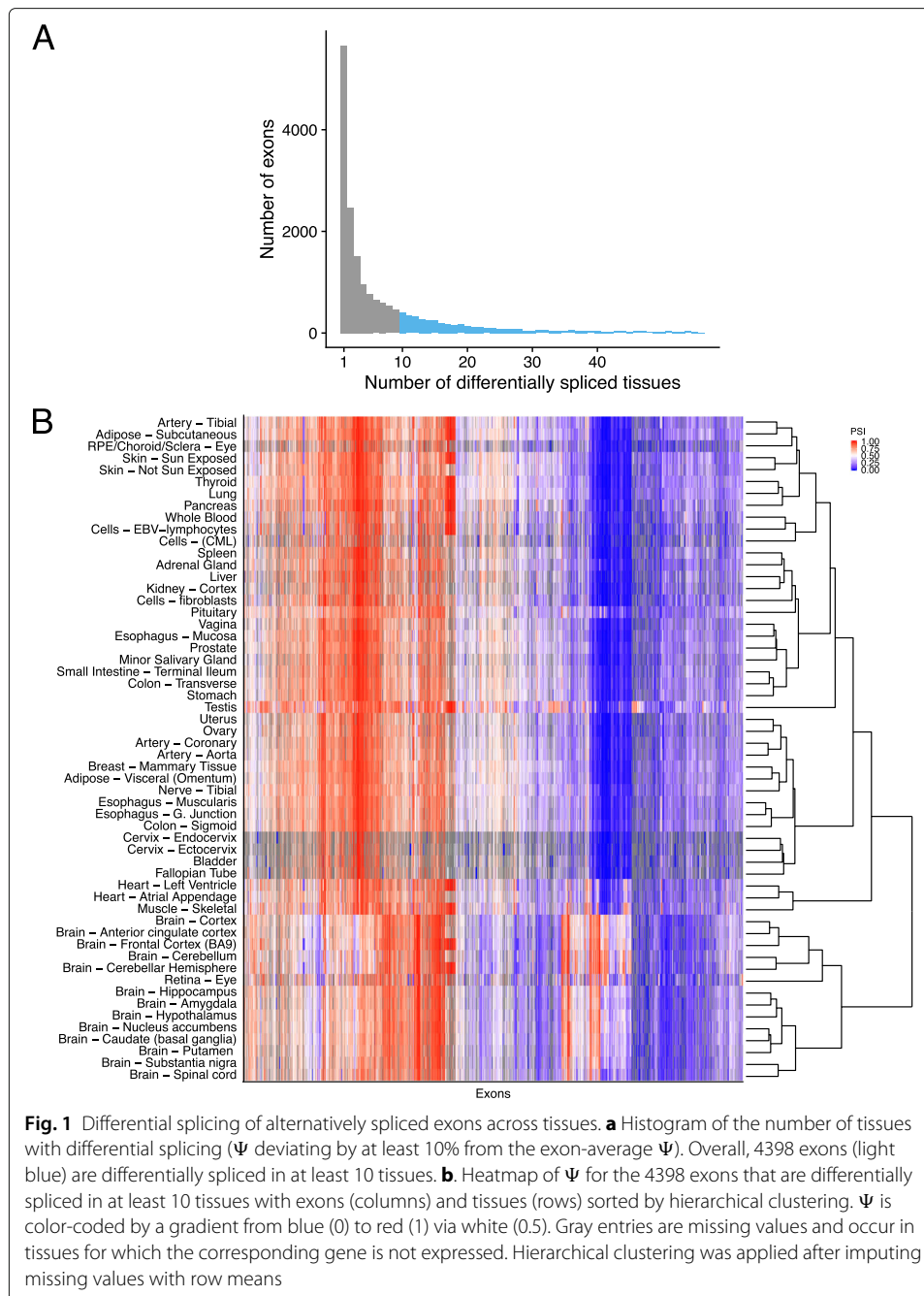
### Tissue-specific alternatively spliced exons

To train a tissue-specific model of splicing, we considered the alternative splicing catalog of the transcriptome ASCOT [43]. Because the ASCOT annotation and quantification pipeline is annotation-free, it also covers non-annotated exons. Altogether, ASCOT provides  $\Psi$  values for 61,823 cassette exons across 56 tissues including 53 tissues from the GTEx dataset [44] and additional RNA-Seq data from peripheral retina. Of note, these tissue-specific values are flagged as missing when the corresponding gene is not expressed [43].

Overall,  $\Psi$  of 17,991 exons (29%) of the ASCOT dataset deviate by at least 10% in at least one tissue from its exon-specific average across tissues. These deviations from the exon-specific average  $\Psi$  by 10% often occurred in a single tissue (5658 exons, 31%) and in at least 10 tissues for 4398 exons (25%, Fig. 1a). We investigated co-variations between tissues using these 4398 exons (Fig. 1b). This revealed that samples from the central nervous system (brain, spinal cord, and retina) have very distinct splicing patterns compared to other tissues, in agreement with previous reports [24]. Moreover, skeletal muscle and the two heart tissues (left ventricle and aortic appendage) also clustered together with shared splicing patterns. Altogether, this analysis indicates that the ASCOT dataset provides thousands of tissue-specific splicing events that could be used to train a sequence-based predictive model. Also, the ASCOT dataset provides the possibility for a multi-task model to exploit shared splicing regulation of tissues of the central nervous system and, to a lower extent, between skeletal muscle and cardiac tissues.

### Differential splicing associated with genetic variants show little tissue-specific variations

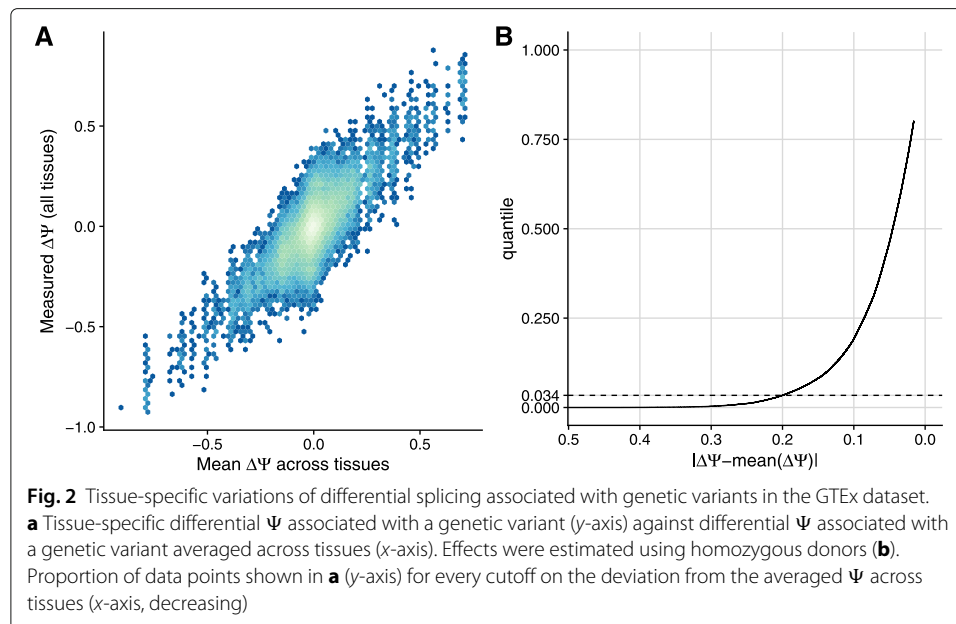
The ASCOT dataset consists of data aggregated per tissue. In principle, the genetic variations between donors of the original GTEx dataset provide further information that a sequence-based model could exploit. We therefore next asked how much genetic variation among individuals in GTEx associated with tissue-specific splicing variations. To this end, we computed  $\Delta\Psi$ , the difference between  $\Psi$  averaged across individuals homozygous for the alternative allele and  $\Psi$  averaged across individuals homozygous for the reference allele for exons with a single variant within the exon body and 300 nucleotides flanking the exon either side ("Materials and methods" section). We estimated  $\Psi$  using the software for estimating splice isoform abundances MISO [45], which only takes annotated and alternatively spliced exons into account. Over all these 1767 single-nucleotide variants, little tissue-specific deviation of  $\Delta\Psi$  compared to its average across tissues was observed (Fig. 2a). Specifically, less than 1476 instances (3.4% of exon-variant-tissue pairs) of tissue-specific  $\Delta\Psi$  deviated by 20% from the tissue-averaged  $\Delta\Psi$  (Fig. 2b). This observation is consistent with the fact that only a limited fraction (between 7 and 21%) of splicing QTLs are tissue-specific [46]. Since GTEx samples are derived from healthy donors, this observation, however, does not rule out the possibility that some disease-



causing variants do alter splicing in a tissue-specific way. Due to the small amount of tissue-specific splicing variation associated with genetic variants in GTEx, we decided to train a sequence-based model solely based on the variations between exons using the ASCOT aggregated data and to keep the genetic variations between donors of the GTEx dataset to independently assess the model afterward.

#### TSplice predicts tissue-specific $\Psi$

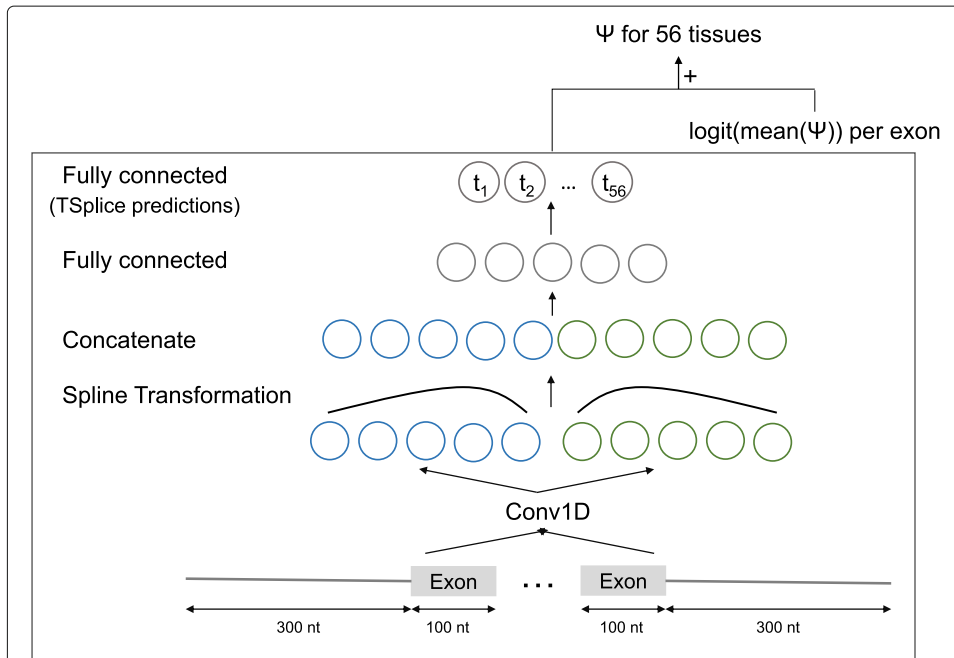
We next developed a neural network, TSplice, to predict tissue-specific  $\Psi$  values from sequence and tissue-averaged  $\Psi$  (“Materials and methods” section). TSplice considers the



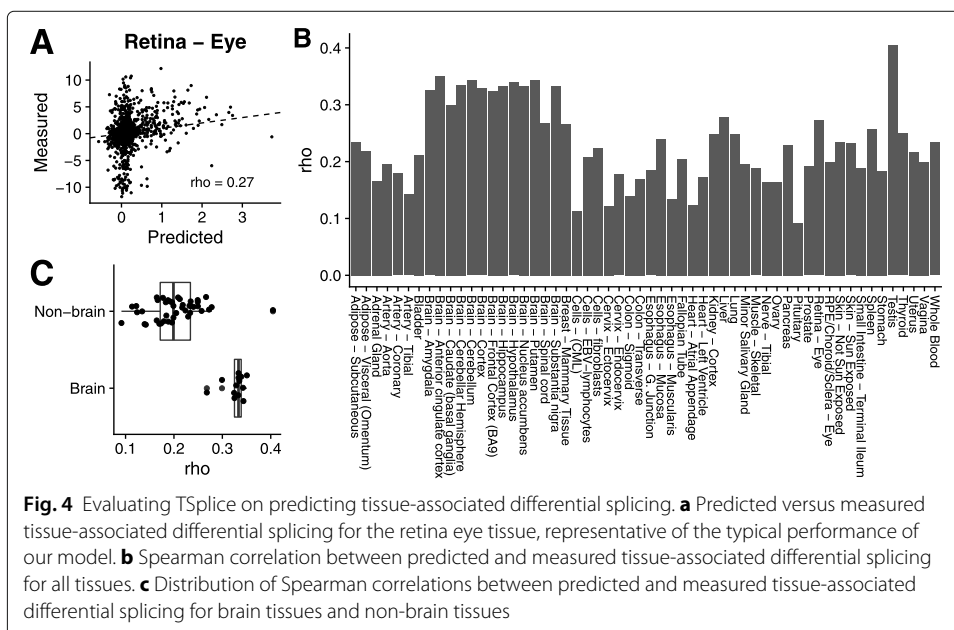
300 nt flanking either side of the exon and the first and last 100 nt of the exon body. TSplice is a convolutional neural network (Fig. 3) in which positional effects of sequence elements relative to splice sites are modeled using spline transformations [47]. TSplice was trained on the ASCOT dataset using all chromosomes except for chromosome 2, 3, and 5. We report our model prediction performances on these held-out chromosomes.

The performance of TSplice was first assessed on test data by comparing the observed against the predicted log odds ratios of tissue-specific  $\Psi$  for 1621 exons (“variable exons”) with  $\Psi$  deviating from the tissue-averaged  $\Psi$  by at least 0.2 in at least one tissue and for which the gene is expressed in at least 10 tissues (Fig. 4a for the retina eye as an example, Spearman  $\rho = 0.27$ ). The predictions positively correlated with the measurements in all tissues and showed a median Spearman correlation of 0.22 (Fig. 4b, Additional file 1: Fig. S1). The performance was higher for tissues of the central nervous system (Fig. 4c), possibly because central nervous system tissues harbor similar splicing patterns and because they are well represented in the ASCOT dataset.

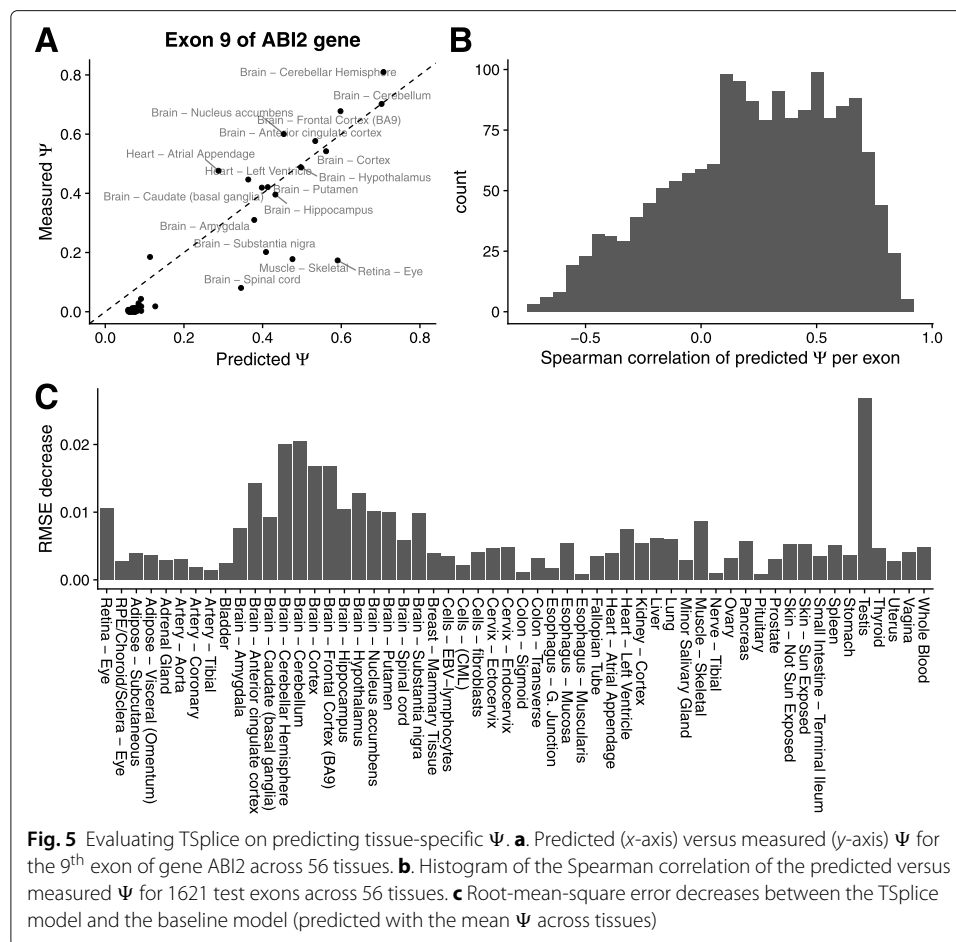
We had first assessed log odds ratio predictions because these are the actual quantities the model was trained for. However, percent spliced-ins on the natural scale often matter more for biological and medical applications. We hence next evaluated how well TSplice performs on predicting tissue-specific  $\Psi$  on test exons. A successful example is the 9<sup>th</sup> exon of the gene *ABI2*, which is included in brain, heart, muscle, and retina tissues and for which TSplice predicts well the order of the tissues (Fig. 5a, Spearman  $\rho = 0.8$ ) and the absolute values of tissue-specific  $\Psi$  per-tissue (root-mean-square error, short RMSE, 0.11). For the majority of the variable exons (73.9%, 1198 out of 1621), TSplice ranked tissue-specific  $\Psi$  in the right direction (median  $\rho = 0.25$ , Fig. 5b). We benchmarked TSplice against a  $L_2$  regularized linear model based on known splice-regulating motifs and splice site sequences (“Materials and methods” section). Although the performance (Spearman correlation of predicted versus measured  $\Psi_{e,t}$ ) of the alternative model correlates ( $R = 0.373$ ) with TSplice across tissues, TSplice outperforms the alternative model for all 56 tissues (Additional file 1: Fig. S2A). Furthermore, when evaluated per



**Fig. 3** Model architecture to predict tissue-specific percent spliced-in. The model TSplice consists of one convolution layer with 64 length-9 filters capturing sequence elements from one-hot encoded input sequences. This is followed by two spline transformation layers modulating the effect of sequence elements depending on their position relative to the acceptor splice sites (leftmost layer) and the donor (rightmost layer). The outputs of the two spline transformation layers are concatenated, and global average pooling is applied along the sequence dimension. This is then followed by feeding two consecutive fully connected layers. The last fully connected layer outputs a 56-dimension vector which are the predicted log odds ratios of tissue-specific  $\Psi$  versus tissue-averaged  $\Psi$  for the 56 tissues of the ASCOT dataset. Natural scale tissue-specific  $\Psi$  are obtained by adding predicted odds ratios with measured tissue-averaged  $\Psi$  on the logit scale. Batch normalization was used after all layers with trainable parameters except the last fully connected layer. In total, the model has 8024 trainable parameters



**Fig. 4** Evaluating TSplice on predicting tissue-associated differential splicing. **a** Predicted versus measured tissue-associated differential splicing for the retina eye tissue, representative of the typical performance of our model. **b** Spearman correlation between predicted and measured tissue-associated differential splicing for all tissues. **c** Distribution of Spearman correlations between predicted and measured tissue-associated differential splicing for brain tissues and non-brain tissues



exon, TSplice had higher Spearman correlation than the alternative model for 63.5% of the exons ( $P < 2.2 \times 10^{-16}$ , paired Wilcoxon test, Additional file 1: Fig. S2B).

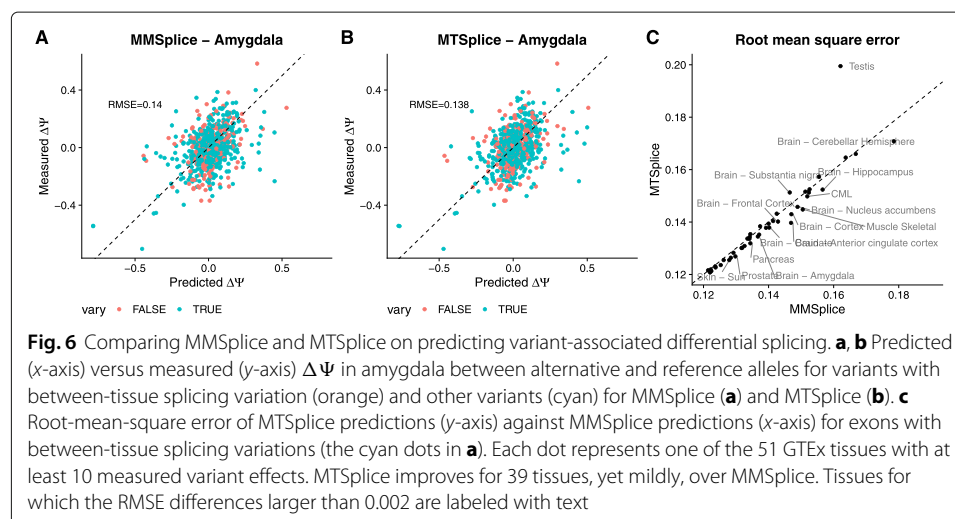
Visualization of the positional weights learned by the splines of TSplice showed that some filters were important for the 5' half of the model, others for the 3' half, while about a third of them were important for both halves. Moreover, positional effects were particularly marked near the splice sites (Additional file 1: Fig. S3). Visualizing the model gradient with respect to the input sequence indicated that the model activates at sequences matching binding site motifs of the splicing factors PTBP1/2, NOVA1/2, and MBNL1 (Additional file 1: Fig. S4 for examples, "Materials and methods" section). To study the role of these motifs systematically, we next generated *in silico* mutated sequences by injecting the consensus sequences of these splicing factor binding sites at various positions along 1000 randomly selected sequences of the test set. We then computed the TSplice score difference between each mutated sequence and its original counterpart. TSplice predicted that cassette exons with a NOVA1/2 binding site consensus sequence in the upstream intron are less spliced-in (i.e., more skipped) in the brain compared to other tissues on average (Additional file 1: Fig. S5). Since the RNA-binding protein NOVA1 is neuron-specifically expressed, these TSplice predictions are consistent with a NOVA1 repressive role when binding in the upstream intron [11]. Moreover, exons with an MBNL1 binding site consensus sequence in the upstream intron are predicted to be more spliced-in

in the brain and muscle than in other tissues on average (Additional file 1: Fig. S5). This is consistent with the repressive role of MBNL1 when binding to upstream introns and with MBNL1 being less expressed in the brain and muscle than in other tissues [48]. The interpretation of the effects of PTBP1/2 binding site consensus sequence is more complex since it is recognized by two competing factors with anti-correlated expression during neuronal differentiation [49].

Altogether, these results show that TSplice captured sequence features predictive of  $\Psi$  changes across tissues.

### Tissue-specific variant effect prediction

We next considered combining MMSplice, which models tissue-independent effects together with TSplice, which models differential effects between tissues, to predict the effects associated with genetic variants for any GTEx tissue (“Materials and methods” section). We name this combined model MTSplice. For amygdala, taken as a representative tissue, the MTSplice predictions correlate well ( $\rho = 0.42$ , Fig. 6a) with differences of  $\Psi$  observed between homozygous donors (“Materials and methods” section). This is consistent with the observation that most variants have similar effects across tissues. Nevertheless, MTSplice further improved the prediction accuracy when evaluated on 1030 variants with  $\Delta\Psi$  varying by at least 0.2 in at least one tissue (RMSE = 0.140 for MMSplice alone, RMSE = 0.138 for MTSplice, RMSE = 0.141 versus 0.139 when evaluated on all variant, Fig. 6). When evaluated on the 51 tissues with at least 10 measured variant effects, MTSplice outperformed MMSplice for 39 out of 51 tissues in terms of root-mean-square error ( $P = 1.76 \times 10^{-5}$ , paired Wilcoxon test, Fig. 6c). Notably, MTSplice outperformed MMSplice in 10 out of 12 brain tissues (Additional file 1: Fig. S6A). Although the improvement of MTSplice over MMSplice are significant, the relative decrease of RMSE remains modest. The relative increases were more pronounced when restricting the analysis to those measurements harboring large tissue-specific effects (Additional file 1: Fig. S6B).

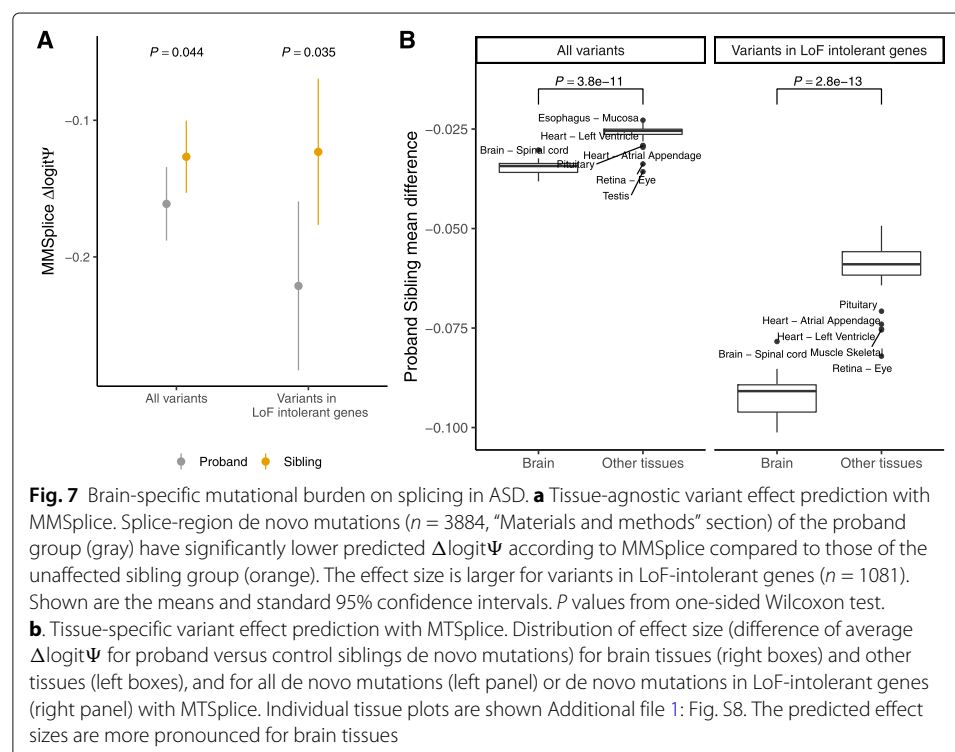




### MTSplice predicts brain-specific signals for autism patients

To assess the potential of MTSplice on scoring tissue-specific disease variants, we considered de novo mutations that were reported for 1790 autism spectrum disorder (ASD) simplex families from the Simons Simplex Collection [50–54] and as provided by Zhou et al. [55]. The data consists of 127,140 de novo mutations, with 65,147 from the proband group and 61,993 from the unaffected siblings. Of those, we further considered the 3884 mutations lying in exons or in their 300-nt flanking intronic regions and predicted with MMSplice with a  $\Delta\text{logit}(\Psi)$  magnitude greater than 0.05. Overall, MMSplice predicted that variants of the proband group would disrupt splicing more strongly than variants of the control siblings (negative MMSplice scores, Fig. 7a,  $P = 0.042$ , Wilcoxon rank-sum test). The effect was even stronger for the 1081 loss-of-function (LoF) intolerant genes (Fig. 7a,  $P = 0.0035$ , Wilcoxon rank-sum test, “Materials and methods” section). This result is consistent with the report that LoF-intolerant genes are vulnerable to noncoding disruptive mutations in ASD [55] and points to an important contribution of splicing.

We then asked whether MTSplice was able to identify tissue-specific effects of ASD-associated de novo mutations. Consistent with the MMSplice results, the de novo mutations of the proband group were predicted by MTSplice to more severely disrupt splicing than the de novo mutations of the control group for all tissues (Fig. 7b). The effect size was larger for the brain tissues (Fig. 7b). Since autism is a neurological disorder, these results indicate that MTSplice may be used to prioritize variants that could play a tissue-specific pathogenic role. Besides the brain tissues, the tissues with most pronounced differences were the retina, which is also part of the central nervous systems and muscle, which has been associated with autism as well [56]. These differences were



further amplified when restricting the analysis to the de novo mutations in LoF-intolerant genes (Fig. 7b).

We next asked whether MTSplice could capture tissue-specific disease signal that would otherwise be missed by MMSplice. Among the mutations predicted to have small effects by MMSplice ( $\Delta\text{logit}(\Psi)$  magnitude smaller than 0.05), the predicted effect for the proband group is indeed no severe than the control sibling group (Additional file 1: Fig. S7A). Moreover, when considering the tissue-specific effect predicted by MTSplice, the de novo mutations from the proband group were predicted to disrupt more severely splicing in the frontal cortex (Additional file 1: Fig. S7A,  $P = 0.036$ , one-sided Wilcoxon test). Furthermore, these ASD-associated signals were generally found in the brain, heart, muscle, and retina but not in other tissues (Additional file 1: Fig. S7B). This shows that MTSplice is able to capture tissue-specific disease signals that would have been missed by MMSplice.

Altogether, these analyses demonstrate the value of MTSplice on predicting tissue-specific effects of potentially disease-causing mutations.

## Discussion

We introduced the model MTSplice which quantitatively predicts effects of human genetic variants on RNA splicing in 56 tissues. MTSplice has two components. One component, MMSplice, models constitutive splicing regulatory sequences. The other component, TSplice, models tissue-specific splicing regulatory sequences. The combined model MTSplice outperforms MMSplice on predicting tissue-specific variations in percent spliced-in associated with naturally occurring genetic variants in most tissues of the GTEx dataset. Applying MTSplice to de novo mutations from autism spectrum disorder simplex families [55], we found a significantly higher burden for the proband group compared to the control siblings, particularly in brain tissues. These results suggest that MTSplice could be applied for scoring variants with a tissue-specific pathogenic role.

Various lines of evidence indicated that our model performed best for brain tissues. This may reflect the fact that brain tissues are well represented in GTEx but also that tissue-specific alternative splicing is particularly strong in brain tissues, giving more useful sequences to train on. Supportive of this, model interpretation revealed that sequence elements known to be bound by brain-specific splicing factors contributed to TSplice predictions. For other tissues, the improvements were more moderate yet consistent. One exceptional tissue is the testis, for which MMSplice still has a substantially better RMSE than MTSplice for variant effect predictions. We could not rationalize that observation. Perhaps, this could be due to the unique transcriptional state of the testis which may affect splicing in a way that the model failed to learn [57].

The TSplice component was trained from tissue-specific alternative splicing observed in the ASCOT dataset. This approach has two main limitations. First, only less than ten thousand exons show tissue-specific alternative splicing in the ASCOT dataset. This amount of data prohibits training of more complex models. In comparison, MMSplice was trained using over 2 million sequences of a massively parallel reporter assay and over half a million naturally occurring splice sites. To overcome this limitation, one could leverage complementary data notably tissue-specific expression of splicing-related RNA-binding proteins (RBPs) combined with transcriptome-wide RBP binding profiles [58]. One example of a transfer learning approach in this context is given by Jha et al. [59],

who showed the benefits of integrating CLIP-Seq data to predict splicing. The second limitation is that the ASCOT dataset is an observational dataset. Models trained from observational data with genomic sequences may learn sequence features that are correlative but not causal, preventing the models from correctly predicting the effect of genetic variants. This could lead to limited predictive performances of our current model.

One approach to overcome the issue of observational data is to perform massively parallel reporter assays (MPRA) for different cell types. MPRA for human splicing have been performed in HEK293 cells [26, 28, 60–62], K562 cells [63, 64], HepG2 cells [63], and HELA and MCF7 cells [65]. These data provide powerful resources to train complex models on splicing, but tissue and cell-type diversity is still lacking. Tissue-specific MPRA data would also be of prime importance for benchmarking models. Here we had to rely on naturally occurring variants in GTEx for benchmarking. Tissue-specific alteration of splicing can be the outcome of genetic variation affecting either (i) constitutive splicing regulatory elements of tissue-specific exons or (ii) tissue-specific splicing regulatory elements. Very few GTEx variants were from the latter class. Hence, the mean square error differences in GTEx between MTSplice and MMSplice could only be very mild. Previous two-cell-line splicing MPRA experiment did not find tissue-specific variant effects between K562 and HepG2 cells [63], maybe also because the variants tested were selected randomly. A designed MPRA, however, could specifically engineer variations of tissue-specific splicing regulatory elements by using prior knowledge in order to more deeply probe the effect of variants on tissue-specific splicing regulation. The generation of large-scale tissue or cell-type-specific perturbation data could therefore be instrumental for probing tissue-specific regulatory elements and could yield more sensitive benchmarks of predictive models. Other than improvement on the training data side, future models might be able design better architecture and data augmentation techniques to further improve performance. Finally, because the approach to predict tissue-specific variant effects by combining MMSplice with a tissue splicing-level prediction model is general, any model that outputs tissue-specific  $\text{logit}(\Psi_{e,t})$  could substitute to TSplice and be combined with MMSplice to predict tissue-specific variant effect on splicing.

## Materials and methods

### Dataset

We split the 61,823 cassette exons from ASCOT into a training, a validation, and a test set. The training set consisted of 38,028 exons from chromosome 4, 6, 8, 10–23, and the sex chromosomes. The 11,955 exons from chromosome 1, 7, and 9 were used as the validation set, and the remaining 11,840 exons were used as the test set (chromosomes 2, 3, and 5). Models are evaluated based on their performance on the test set.

### Variant effect estimation

To compute variant effect, we first computed  $\Psi$  with MISO for all annotated alternatively spliced exons (MISO annotation v2.0, [http://genes.mit.edu/burgelab/miso/annotations/ver2/miso\\_annotations\\_hg19\\_v2.zip](http://genes.mit.edu/burgelab/miso/annotations/ver2/miso_annotations_hg19_v2.zip)) in all GTEx RNA-Seq samples. This led to  $\Psi$  estimates for 4686 samples from 53 tissues. Second, for each exon, we estimated variant effects using only those samples with a single variant within the exon body and 300 nt flanking of the exon. Third, we estimated the effect associated with the variants as the difference between  $\Psi$  averaged across samples homozygous for the alternative allele and

$\Psi$  averaged across samples homozygous for the reference allele. We required at least 2 samples in each of these two groups. For simplicity, we did not consider heterozygous samples for estimating the effects because  $\Psi$  of heterozygous samples is confounded by allele-specific RNA expression. Also, we did not consider indels.

### The TSplice model

We denote  $\Psi_{e,t}$  the percent spliced-in value of the cassette exon  $e$  in tissue  $t$ . The goal of the multi-tissue splicing model is to predict tissue-specific  $\Psi_{e,t}$  from the nucleotide sequence of the given exon  $S_e$ . We train the tissue-specific splicing model with multi-task learning, where each task corresponds to a tissue. The model has two input branches. The first input branch consists of the sequence 300 nt upstream of the acceptor and 100 nt downstream of the acceptor (Fig. 3). In a symmetric fashion, the second input branch consists of the sequence from the donor side, with 100 nt upstream of the donor and 300 nt downstream of the donor. All input sequences are one-hot encoded. The input layer is followed by a 1D convolution layer with 64 filters of length 9. Parameters of the convolution layer are shared by the two input branches, based on the assumption that many sequence motifs are presented both upstream and downstream of the exons. To model the positional effects of splicing motifs, spline transformations [47] are fitted for each of the convolution filters to weight the convolution activations based on the relative input position to donor and acceptor sites. The spline transformations are fitted differently for the two input branches to account for potential different positional effects of the upstream and downstream introns. The weighted activations are then concatenated along the sequence dimension. Two fully connected layers are followed after the concatenated outputs. The last fully connected layer output number of predictions equals the number of tissues ( $T$ ), corresponding to predictions for each tissue. These are the predictions of the TSplice model mentioned in the manuscript. During training, logit of the mean  $\Psi$  per exon ( $\text{logit}(\bar{\Psi}_e)$ ) was added to these prediction outputs followed by a sigmoid function. This encourages the model to learn sequence features associated with differential splicing across tissues.

Formally, for each exon, TSplice predicts for each tissue its  $\Psi_{e,t}$  deviation from the mean  $\bar{\Psi}_e$  across tissues on logit level. Specifically, we define the tissue-associated differential splicing as  $\Delta_{\text{tissue}}\text{logit}(\Psi_{e,t})$

$$\Delta_{\text{tissue}}\text{logit}(\Psi_{e,t}) := \text{logit}(\Psi_{e,t}) - \text{logit}(\bar{\Psi}_e) \quad (1)$$

as the logit  $\Psi$  deviation for tissue  $t$  and exon  $e$  from the logit of  $\bar{\Psi}_e := \frac{1}{T} \sum_{t=1}^T \Psi_{e,t}$ , the mean  $\Psi$  across tissues.

For exon  $e$  with input sequence  $S_e$ , TSplice predicts the target in  $\mathbb{R}^T$ :  $\text{TSplice}(S_e) := (\Delta_{\text{tissue}}\text{logit}(\Psi_{e,1}), \dots, \Delta_{\text{tissue}}\text{logit}(\Psi_{e,T}))$  corresponding to  $T$  tissues.

The tissue-specific  $\Psi_{e,t}$  can be predicted with TSplice and the given  $\text{logit}(\bar{\Psi}_e)$  computed from the data as:

$$\hat{\Psi}_{e,t} = \sigma(\text{TSplice}(S_e)_t + \text{logit}(\bar{\Psi}_e)) \quad (2)$$

where  $\text{TSplice}(S_e)_t$  is the TSplice predicted  $\Delta_{\text{tissue}}\text{logit}(\Psi_{e,t})$ , and  $\sigma$  is the sigmoid function:  $\sigma(x) = \frac{1}{1+e^{-x}}$ . Note that in Eq. 1 and elsewhere the average was computed before and not after logit-transformation because it gave more robust results.

### Model training and selection

The model was implemented with keras (version 2.2.4). The Kullback–Leibler (KL) divergence between the predicted and measured  $\Psi$  distribution was used as the loss function (Eq. 3), by considering the percent spliced-in as the probability of the cassette exon to be included in any given transcript.

$$\text{Loss} = \frac{1}{T \cdot E} \sum_{t=1}^T \sum_{e=1}^E \gamma_{e,t} \left( \Psi_{e,t} \log\left(\frac{\Psi_{e,t}}{\hat{\Psi}_{e,t}}\right) + (1 - \Psi_{e,t}) \log\left(\frac{1 - \Psi_{e,t}}{1 - \hat{\Psi}_{e,t}}\right) \right), \quad (3)$$

where

$$\gamma_{e,t} = \begin{cases} 1, & \text{if } \Psi_{e,t} \text{ observed} \\ 0, & \text{otherwise} \end{cases} \quad (4)$$

Missing values, which typically correspond to tissues in which the gene is not expressed, were masked out in the loss function.  $\Psi$  values were clipped to be between  $[10^{-5}, 1 - 10^{-5}]$ . Adam optimizer [66] with default parameters was used to optimize the model. Network weights were initialized with the He Normal initialization [67]. Hyperparameter search was performed with hyperopt [68] with the Tree Parzen Estimators method along with the package kopt (<https://github.com/Avsecz/kopt>). Hyperparameters were selected based on the loss on the validation set.

After finding the best hyperparameter combination, 20 models were trained with the best hyperparameters but different random initialization. A forward model selection strategy was used to select a set of models whose average predictions gives the smallest loss on the validation set. To this end, models were first sorted based on their performance on the validation set. Next, models were successively added to an ensemble model, defined as the average over the selected models, until the validation set performance no longer improved. This procedure yielded an ensemble model composed of 8 individual models. TSplice predictions are made by this ensemble model.

### Alternative tissue-specific $\Psi_{e,t}$ prediction model

The following model was considered as alternative model to predict tissue-specific  $\hat{\Psi}_{e,t}$ : For each tissue  $t$ , we train a  $L_2$  regularized linear regression model (ridge regression):

$$\text{logit}(\Psi_{e,t}) = \beta_{0,t} + \sum_i \beta_{i,t} X_i + \text{logit}(\Psi_{e,\text{average}}) + \epsilon_t, \quad (5)$$

where  $X_i$  are the sequence features shared for all tissues. Sequence features from five regions were considered separately: upstream intron, exon, downstream intron, donor and acceptor. For upstream and downstream introns, we considered TCAT (NOVA1/2), TGCAIG (RBF0X1), GCTTGC (MBNL1), and all 6-mers with A/T (PTBP1/2) [69–72]. For exons, we considered NOVA1/2, RBF0X1 and MBNL1 motifs and 2272 exonic splicing regulators identify by [26]. For all motifs,  $X_i$  are the vectors of counts of motif instances in the considered region. Acceptor and donor splice site contexts were also considered as for MaxEntScan [24]. On the acceptor side, 20 nt in the intron and 3 nt in the exon were considered. On the donor side, 6 nt in the intron and 3 nt in the exon were considered.

The alternative model was trained on the training and evaluation set, with hyperparameters chosen by cross-validation on the training and evaluation set. The performances were assessed for the same test than TSplice.

### Splicing motif analysis

Splicing motif logos were visualized with contribution scores computed with gradient times input. Examples of motifs were manually selected. Systematic motif discovery with contribution scores was unsuccessful. Motif instances were picked by manually inspecting the activation scores on the input sequences. Motif binding protein was determined by searching motif instances in the ATtRACT database [73]. Position weight matrix (PWM) for PTBP1/2 and MBNL1 were downloaded from RNAComete [74]. PWM for NOVA1/2 which is missing from RNAComete was downloaded from the RBPDB [75].  $P$  values for motif matching were computed with TOMTOM [76] and motif database from RNAComete [74].

To visualize the predicted effects of motifs across tissues, sequence motifs were inserted into the native sequences of 1000 randomly selected exons. Scores were computed by subtracting the TSplice predictions with the inserted motifs and without the motifs. Scores were computed for all tissues and all positions. The heatmaps visualize the scores averaged across 1000 exons by tissues and positions.

### Tissue-specific variant effect prediction

Tissue-specific variant effect  $\Delta\Psi_{e,t}$  is predicted as follows (we considered in this study only homozygous cases as described in the “Variant effect estimation” section in the “Materials and methods” section):

$$\Delta\Psi_{e,t} = \Psi_{e,t}^{\text{alt}} - \Psi_{e,t}^{\text{ref}} \quad (6)$$

where  $\Psi_{e,t}^{\text{ref}}$  is the measured  $\Psi$  for exon  $e$  and tissue  $t$  with the reference sequence, and  $\Psi_{e,t}^{\text{alt}}$  is the tissue-specific  $\Psi$  with the alternative sequence. We model the logit level of  $\Psi_{e,t}^{\text{alt}}$  with the following linear model:

$$\text{logit}(\Psi_{e,t}^{\text{alt}}) = \beta_0 + \beta_{\text{tissue}} + \beta_{\text{alt}} + \beta_{\text{alt} \times \text{tissue}} + \epsilon, \quad (7)$$

where  $\beta_0$  is intercept,  $\beta_{\text{tissue}}$  is the tissue effect,  $\beta_{\text{alt}}$  is the effect of the variant on an average tissue,  $\beta_{\text{alt} \times \text{tissue}}$  is the interaction term which we model the interaction of the variant effect and the given tissue. We model each of the terms as follows:

$$\begin{aligned} \beta_0 &= \text{logit}(\Psi_{e,\text{average}}^{\text{ref}}) \\ \beta_{\text{tissue}} &= \text{TSplice}(S_{\text{ref}}, \text{tissue}) \\ \beta_{\text{alt}} &= \text{MMSplice}(S_{\text{ref}}, S_{\text{alt}}) \\ \beta_{\text{alt} \times \text{tissue}} &= \text{TSplice}(S_{\text{alt}}, \text{tissue}) - \text{TSplice}(S_{\text{ref}}, \text{tissue}) \end{aligned} \quad (8)$$

When we plug Eq. 8 into Eq. 7, we obtain the MTSplice model which combines MMSplice and TSplice to model tissue-specific variant effect:

$$\text{logit}(\Psi_{e,t}^{\text{alt}}) = \text{logit}(\Psi_{e,\text{average}}^{\text{ref}}) + \text{MMSplice}(S_{\text{ref}}, S_{\text{alt}}) + \text{TSplice}(S_{\text{alt}}, \text{tissue}) + \epsilon \quad (9)$$

Finally, the tissue-specific  $\Delta\Psi_{e,t}$  is predicted as follows:

$$\begin{aligned} \Delta\Psi_{e,t} &= \sigma \left( \text{logit}(\Psi_{e,\text{average}}^{\text{ref}}) + \text{MMSplice}(S_{\text{ref}}, S_{\text{alt}}) + \text{TSplice}(S_{\text{alt}}, \text{tissue}) \right) \\ &\quad - \Psi_{e,t}^{\text{ref}} \end{aligned} \quad (10)$$

While the prediction with MMSplice is as follows:

$$\Delta\Psi_{e,t} = \sigma \left( \text{logit}(\Psi_{e,\text{average}}^{\text{ref}}) + \text{MMSplice}(S_{\text{ref}}, S_{\text{alt}}) \right) - \Psi_{e,t}^{\text{ref}} \quad (11)$$

### Benchmark variant effect prediction on GTEx

On the benchmark of tissue-specific variant effect prediction, we further applied four filters. First, we selected variants that have  $|\Delta\Psi_{e,t} - \Delta\Psi_{e,average}| > 0.2$  in at least one tissue. Second,  $\Delta\Psi$  can be computed for at least 3 tissues. Third, we only considered tissues with more than 10 variants satisfying the above criteria. Altogether, these filters led to 1030 variant-exon pairs and 51 tissues used for benchmarking tissue-specific variant effect predictions.

To benchmark variants only in tissues where their effects are tissue-specific, we selected for each variant only the tissues where  $|\Delta\Psi_{e,t} - \Delta\Psi_{e,average}| > 0.2$ . Only tissues with at least 10 valid variants are considered. In total, 48 tissues were considered for this analysis.

### Autism variants

The processed de novo mutations were downloaded from the link provided by Zhou et al. [55] (<https://hb.flatironinstitute.org/asdbrowser/>). The original whole genome sequencing data were accessed through the Simons Foundation Autism Research Initiative (SFARI) [50–54]. The data provides 127,140 single-nucleotide variants (SNVs) from non-repeat-region. The variants were derived from 7097 whole genomes from the Simons Simplex Collection (SSC) cohort, which consists of whole-genome sequencing data from 1790 families (with probands and matched unaffected siblings).

To predict variant effect on splicing, variants were mapped to exons if they are within the annotated (ensembl gene annotation v75) exon body or within 300 nt flanking. If a variant was mapped to multiple exons, the largest effect size was reported as the effect of the variant. A total of 13,415 variants were mapped to known exons and therefore were predicted by our models. Among those variants, 3884 have predicted  $|\Delta\logit(\Psi)| > 0.05$ . We classified the variants into loss-of-function (LoF) group and loss tolerant group based on the loss-of-function observed/expected (oe) upper bound fraction (LOEUF) scores [77]. We used the suggested cutoff of 0.35 on the upper bound of the oe confidence interval to group the variants.

### Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s13059-021-02273-7>.

**Additional file 1:** Supplementary Figures S1-S8.

**Additional file 2:** Review history.

### Peer review information

Barbara Cheifet and Tim Sands were the primary editors of this article and managed its editorial process and peer review in collaboration with the rest of the editorial team.

### Acknowledgements

We thank Nils Wagner, Florian Hölzlwimmer, and Leonhard Wachutka for their feedback on the manuscript.

### Review history

The review history is available as Additional file 2.

### Author's contributions

JC and JG designed the model with the help of AK. JC and MHÇ implemented the software. JC performed the analysis with the help of MHÇ. JG and AK supervised the project. JC and JG wrote the manuscript. The authors read and approved the final manuscript.

### Funding

This work was supported by the Competence Network for Technical, Scientific High Performance Computing in Bavaria KONWIHR (to J.C. and M.H.Ç.) and by the German Bundesministerium für Bildung und Forschung (BMBF) through the

project MERGE (031L0174A to J.G.). A.K. was supported by NIH grant DP2GM12348501. We thank the TUM Informatics graduate school CEDOSIA for supporting J.C. to have a 3-month visit to the laboratory of Anshul Kundaje. The Genotype-Tissue Expression (GTEx) project was supported by the Common Fund of the Office of the Director of the National Institutes of Health, and by NCI, NHGRI, NHLBI, NIDA, NIMH, and NINDS. The data used for the analyses described in this manuscript were obtained from the GTEx Portal on April 4, 2019, under accession number dbGaP: phs000424.v7.p2. Open Access funding enabled and organized by Projekt DEAL.

#### Availability of data and materials

MTSplice is integrated with MMSplice and is available from [https://github.com/gagneurlab/MMSplice\\_MTSsplice](https://github.com/gagneurlab/MMSplice_MTSsplice) [78] under the MIT license, and also deployed at kipoi <http://kipoi.org/models/MMSplice/mtsplice/> [79]. The source code version used in the manuscript is available from <https://doi.org/10.5281/zenodo.4255942> [80]. The software comes along with reference  $\Psi$  tables. For novel cassette exons, users can provide their own  $\Psi$  table. Analysis code is available from <https://gitlab.cmm.in.tum.de/gagneurlab/mtsplice> [81] under the MIT license. ASCOT data is available from <http://ascot.cs.jhu.edu/>.

#### Ethics approval and consent to participate

Not applicable.

#### Competing interests

The authors declare that they have no competing interests.

#### Author details

<sup>1</sup>Department of Informatics, Technical University of Munich, Boltzmannstraße, 85748 Garching, Germany. <sup>2</sup>Department of Computer Science, Stanford University, Stanford, CA, USA. <sup>3</sup>Department of Genetics, Stanford University, Stanford, CA, USA. <sup>4</sup>Institute of Computational Biology, Helmholtz Zentrum München, Neuherberg, Germany. <sup>5</sup>Institute of Human Genetics, Klinikum rechts der Isar, Technical University of Munich, Munich, Germany.

Received: 6 June 2020 Accepted: 14 January 2021

Published online: 31 March 2021

#### References

- Pan Q, Shai O, Lee LJ, Frey BJ, Blencowe BJ. Deep surveying of alternative splicing complexity in the human transcriptome by high-throughput sequencing. *Nat Genet.* 2008;40(12):1413–5.
- Castle JC, Zhang C, Shah JK, Kulkarni AV, Kalsotra A, Cooper TA, Johnson JM. Expression of 24,426 human alternative splicing events and predicted cis regulation in 48 tissues and cell lines. *Nat Genet.* 2008;40(12):1416–25.
- Wang ET, Sandberg R, Luo S, Khrebtkova I, Zhang L, Mayr C, Kingsmore SF, Schroth GP, Burge CB. Alternative isoform regulation in human tissue transcriptomes. 2008. <https://doi.org/10.1038/nature07509>.
- Matera AG, Wang Z. A day in the life of the spliceosome. *Nat Rev Mol Cell Biol.* 2014;15(2):108–21.
- Ule J, Blencowe BJ. Alternative splicing regulatory networks: functions, mechanisms, and evolution. *Mol Cell.* 2019;76(2):329–45.
- Yeo G, Holste D, Kreiman G, Burge CB. Variation in alternative splicing across human tissues. *Genome Biol.* 2004;5(10):74.
- Tapial J, Ha KCH, Sterne-Weiler T, Gohr A, Braunschweig U, Hermoso-Pulido A, Quesnel-Vallières M, Permanyer J, Sodaei R, Marquez Y, Cozzuto L, Wang X, Gómez-Velázquez M, Rayon T, Manzanares M, Ponomarenko J, Blencowe BJ, Irimia M. An atlas of alternative splicing profiles and functional associations reveals new regulatory programs and genes that simultaneously express multiple major isoforms. *Genome Res.* 2017;27(10):1759–68.
- Kalsotra A, Cooper TA. Functional consequences of developmentally regulated alternative splicing. *Nat Rev Genet.* 2011;12(10):715–729.
- Gabut M, Samavarchi-Tehrani P, Wang X, Slobodenic V, O'Hanlon D, Sung H-K, Alvarez M, Talukder S, Pan Q, Mazzoni EO, Nedelec S, Wichterle H, Woltjen K, Hughes TR, Zandstra PW, Nagy A, Wrana JL, Blencowe BJ. An alternative splicing switch regulates embryonic stem cell pluripotency and reprogramming. *Cell.* 2011;147(1):132–46.
- Buljan M, Chalancon G, Eustermann S, Wagner GP, Fuxreiter M, Bateman A, Babu MM. Tissue-specific splicing of disordered segments that embed binding motifs rewires protein interaction networks. *Mol Cell.* 2012;46(6):871–83.
- Ule J, Stefani G, Mele A, Ruggiu M, Wang X, Taneri B, Gaasterland T, Blencowe BJ, Darnell RB. An RNA map predicting Nova-dependent splicing regulation. 2006. <https://doi.org/10.1038/nature05304>.
- Barbosa-Morais NL, Irimia M, Pan Q, Xiong HY, Gueroussov S, Lee LJ, Slobodenic V, Kutter C, Watt S, Colak R, Kim T, Misquitta-Ali CM, Wilson MD, Kim PM, Odom DT, Frey BJ, Blencowe BJ. The evolutionary landscape of alternative splicing in vertebrate species. *Science.* 2012;338(6114):1587–93.
- Llorian M, Schwartz S, Clark TA, Hollander D, Tan L-Y, Spellman R, Gordon A, Schweitzer AC, de la Grange P, Ast G, Smith CWJ. Position-dependent alternative splicing activity revealed by global profiling of alternative splicing events regulated by PTB. *Nat Struct Mol Biol.* 2010;17(9):1114–23.
- Badr E, ElHefnawi M, Heath LS. Computational identification of tissue-specific splicing regulatory elements in human genes from RNA-Seq data. *PLoS ONE.* 2016;11(11):0166978.
- Chen M, Manley JL. Mechanisms of alternative splicing regulation: insights from molecular and genomics approaches. *Nat Rev Mol Cell Biol.* 2009;10(11):741–54.
- López-Bigas N, Audit B, Ouzounis C, Parra G, Guigó R. Are splicing mutations the most frequent cause of hereditary disease? 2005. <https://doi.org/10.1016/j.febslet.2005.02.047>.
- Li YI, van de Geijn B, Raj A, Knowles DA, Petti AA, Golan D, Gilad Y, Pritchard JK. RNA splicing is a primary link between genetic variation and disease. *Science.* 2016;352(6285):600–4.
- Scotti MM, Swanson MS. RNA mis-splicing in disease. 2016. <https://doi.org/10.1038/nrg.2015.3>.



19. Parras A, Anta H, Santos-Galindo M, Swarup V, Elorza A, Nieto-González JL, Picó S, Hernández IH, Díaz-Hernández JL, Belloc E, Rodolosse A, Parikshak NN, Peñagarikano O, Fernández-Chacón R, Irimia M, Navarro P, Geschwind DH, Méndez R, Lucas JJ. Autism-like phenotype and risk gene mRNA deadenylation by CPEB4 mis-splicing. *Nature*. 2018;560(7719):441–6.
20. Irimia M, Weatheritt RJ, Ellis JD, Parikshak NN, Gonatopoulos-Pournatzis T, Babor M, Quesnel-Vallières M, Tapial J, Raj B, O'Hanlon D, Barrios-Rodiles M, Sternberg MJE, Cordes SP, Roth FP, Wrana JL, Geschwind DH, Blencowe BJ. A highly conserved program of neuronal microexons is misregulated in autistic brains. *Cell*. 2014;159(7):1511–23.
21. Voineagu I, Wang X, Johnston P, Lowe JK, Tian Y, Horvath S, Mill J, Cantor RM, Blencowe BJ, Geschwind DH. Transcriptomic analysis of autistic brain reveals convergent molecular pathology. *Nature*. 2011;474(7351):380–4.
22. Uddin M, Tammimies K, Pellicchia G, Alipanahi B, Hu P, Wang Z, Pinto D, Lau L, Nalpathamkalam T, Marshall CR, Blencowe BJ, Frey BJ, Merico D, Yuen RKC, Scherer SW. Brain-expressed exons under purifying selection are enriched for de novo mutations in autism spectrum disorder. 2014. <https://doi.org/10.1038/ng.2980>.
23. Pertea M, Lin X, Salzberg SL. GeneSplicer: a new computational method for splice site prediction. *Nucleic Acids Res*. 2001;29(5):1185–90.
24. Yeo G, Burge CB. Maximum entropy modeling of short sequence motifs with applications to RNA splicing signals. *J Comput Biol*. 2004;11(2-3):377–94.
25. Desmet F-O, Hamroun D, Lalande M, Collod-Béroud G, Claustres M, Béroud C. Human splicing finder: an online bioinformatics tool to predict splicing signals. *Nucleic Acids Res*. 2009;37(9):67.
26. Ke S, Shang S, Kalachikov SM, Morozova I, Yu L, Russo JJ, Ju J, Chasin LA. Quantitative evaluation of all hexamers as exonic splicing elements. *Genome Res*. 2011;21(8):1360–74.
27. Jian X, Boerwinkle E, Liu X. In silico prediction of splice-altering single nucleotide variants in the human genome. *Nucleic Acids Res*. 2014;42(22):13534–44.
28. Rosenberg AB, Patwardhan RP, Shendure J, Seelig G. Learning the sequence determinants of alternative splicing from millions of random sequences. *Cell*. 2015;163(3):698–711.
29. Xiong HY, Alipanahi B, Lee LJ, Bretschneider H, Merico D, Yuen RKC, Hua Y, Guerussov S, Najafabadi HS, Hughes TR, Morris Q, Barash Y, Krainer AR, Jovic N, Scherer SW, Blencowe BJ, Frey BJ. RNA splicing the human splicing code reveals new insights into the genetic determinants of disease. *Science*. 2015;347(6218):1254806.
30. Cheng J, Nguyen TYD, Cygan KJ, Çelik MH, Fairbrother WG, Avsec Ž., Gagneur J. MMSplice: modular modeling improves the predictions of genetic variant effects on splicing. *Genome Biol*. 2019;20(1):48.
31. Cheng J, Çelik MH, Nguyen TYD, Avsec Ž., Gagneur J. CAGI 5 splicing challenge: improved exon skipping and intron retention predictions with MMSplice. *Human Mutation*. 2019;40(9):1243–51.
32. Sonnenburg S, Schweikert G, Phillips P, Behr J, Rättsch G. Accurate splice site prediction using support vector machines. *BMC Bioinformatics*. 2007;8 Suppl 10:7.
33. Jaganathan K, Kyriazopoulou Panagiotopoulou S, McRae JF, Darbandi SF, Knowles D, Li Yi, Kosmicki JA, Arbelaez J, Cui W, Schwartz GB, Chow ED, Kanterakis E, Gao H, Kia A, Batzoglu S, Sanders SJ, Farh KK-H. Predicting splicing from primary sequence with deep learning. *Cell*. 2019;176(3):535–54824.
34. Barash Y, Calarco JA, Gao W, Pan Q, Wang X, Shai O, Blencowe BJ, Frey BJ. Deciphering the splicing code. *Nature*. 2010;465(7294):53–59.
35. Xiong HY, Barash Y, Frey BJ. Bayesian prediction of tissue-regulated splicing using RNA sequence and cellular context. *Bioinformatics*. 2011;27(18):2554–62.
36. Mount SM, Avsec Ž., Carmel L, Casadio R, Çelik MH, Chen K, Cheng J, Cohen NE, Fairbrother WG, Fenesh T, Gagneur J, Gotea V, Holzer T, Lin C-F, Martelli PL, Naito T, Nguyen TYD, Savojardo C, Unger R, Wang R, Yang Y, Zhao H. Assessing predictions of the impact of variants on splicing in CAGI5. *Hum Mutat*. 2019;40(9):1215–24.
37. Zhou J, Troyanskaya OG. Predicting effects of noncoding variants with deep learning-based sequence model. 2015. <https://doi.org/10.1038/nmeth.3547>.
38. Kelley DR, Snoek J, Rinn JL. Basset: learning the regulatory code of the accessible genome with deep convolutional neural networks. *Genome Res*. 2016;26(7):990–9.
39. Zhou J, Theesfeld CL, Yao K, Chen KM, Wong AK, Troyanskaya OG. Deep learning sequence-based ab initio prediction of variant effects on expression and disease risk. *Nat Genet*. 2018;50(8):1171–9.
40. Quang D, Xie X. FactorNet: a deep learning framework for predicting cell type specific transcription factor binding from nucleotide-resolution sequential data. 2019. <https://doi.org/10.1016/j.jymeth.2019.03.020>.
41. Nair S, Kim DS, Perricone J, Kundaje A. Integrating regulatory dna sequence and gene expression to predict genome-wide chromatin accessibility across cellular contexts. *Bioinformatics*. 2019;35(14):108–16.
42. Avsec Ž., Kreuzhuber R, Israeli J, Xu N, Cheng J, Shrikumar A, Banerjee A, Kim DS, Beier T, Urban L, Kundaje A, Stegle O, Gagneur J. The kipo repository accelerates community exchange and reuse of predictive models for genomics. *Nat Biotechnol*. 2019;37(6):592–600.
43. Ling JP, Wilks C, Charles R, Leavey PJ, Ghosh D, Jiang L, Santiago CP, Pang B, Venkataraman A, Clark BS, et al. Ascot identifies key regulators of neuronal subtype-specific splicing. *Nat Commun*. 2020;11(1):1–12.
44. GTEx Consortium. The Genotype-Tissue expression (GTEx) project. *Nat Genet*. 2013;45(6):580–5.
45. Katz Y, Wang ET, Airolidi EM, Burge CB. Analysis and design of RNA sequencing experiments for identifying isoform regulation. *Nat Methods*. 2010;7(12):1009–15.
46. Mele M, Ferreira PG, Reverter F, DeLuca DS, Monlong J, Sammeth M, Young TR, Goldmann JM, Pervouchine DD, Sullivan TJ, Johnson R, Segre AV, Djebali S, Niarchou A, T G Consortium, Wright FA, Lappalainen T, Calvo M, Getz G, Dermizakis ET, Ardlie KG, Guigo R. The human transcriptome across tissues and individuals. *Science*. 2015;348(6235):660–5.
47. Avsec Ž., Barekatin M, Cheng J, Gagneur J. Modeling positional effects of regulatory sequences with spline transformations increases prediction accuracy of deep neural networks. 2018. <https://doi.org/10.1093/bioinformatics/btx727>.

48. Wang ET, Cody NA, Jog S, Biancolella M, Wang TT, Treacy DJ, Luo S, Schroth GP, Housman DE, Reddy S, et al. Transcriptome-wide regulation of pre-mRNA splicing and mRNA localization by muscleblind proteins. *Cell*. 2012;150(4):710–724.
49. Linares AJ, Lin C-H, Damianov A, Adams KL, Novitsch BG, Black DL. The splicing regulator PTBP1 controls the activity of the transcription factor pbx1 during neuronal differentiation. *Elife*. 2015;4:09268.
50. Yuen RKC, Thiruvahindrapuram B, Merico D, Walker S, Tammimies K, Hoang N, Chrysler C, Nalpathamkalam T, Pellicchia G, Liu Y, Gazzellone MJ, D'Abate L, Deneault E, Howe JL, Liu RSC, Thompson A, Zarrei M, Uddin M, Marshall CR, Ring RH, Zwaigenbaum L, Ray PN, Weksberg R, Carter MT, Fernandez BA, Roberts W, Szatmari P, Scherer SW. Whole-genome sequencing of quartet families with autism spectrum disorder. *Nat Med*. 2015;21(2):185–91.
51. C Yuen RKC, Merico D, Bookman M, L Howe J, Thiruvahindrapuram B, Patel RV, Whitney J, Deflaux N, Bingham J, Wang Z, Pellicchia G, Buchanan JA, Walker S, Marshall CR, Uddin M, Zarrei M, Deneault E, D'Abate L, Chan AJ, Koyanagi S, Paton T, Pereira SL, Hoang N, Engchuan W, Higginbotham EJ, Ho K, Lamoureux S, Li W, MacDonald JR, Nalpathamkalam T, Sung WWL, Tsoi FJ, Wei J, Xu L, Tasse A-M, Kirby E, Van Etten W, Twigger S, Roberts W, Drmic I, Jilderda S, Modi BM, Kellam B, Szego M, Cytrynbaum C, Weksberg R, Zwaigenbaum L, Woodbury-Smith M, Brian J, Senman L, Iaboni A, Doyle-Thomas K, Thompson A, Chrysler C, Leef J, Savion-Lemieux T, Smith IM, Liu X, Nicolson R, Seifer V, Fedele A, Cook EH, Dager S, Estes A, Gallagher L, Malow BA, Parr JR, Spence SJ, Vorstman J, Frey BJ, Robinson JT, Strug LJ, Fernandez BA, Elsabbagh M, Carter MT, Hallmayer J, Knoppers BM, Anagnostou E, Szatmari P, Ring RH, Glazer D, Pletcher MT, Scherer SW. Whole genome sequencing resource identifies 18 new candidate genes for autism spectrum disorder. *Nat. Neurosci*. 2017;20(4):602–611.
52. Michaelson JJ, Shi Y, Gujral M, Zheng H, Malhotra D, Jin X, Jian M, Liu G, Greer D, Bhandari A, Wu W, Corominas R, Peoples A, Koren A, Kang S, Lin GN, Estabillio J, Gadomski T, Singh B, Zhang K, Akshoomoff N, Corsello C, McCarroll S, Iakoucheva LM, Li Y, Wang J, Sebat J. Whole-genome sequencing in autism identifies hot spots for de novo germline mutation. *Cell*. 2012;151(7):1431–42.
53. Jiang Y-H, Yuen RKC, Jin X, Wang M, Chen N, Wu X, Ju J, Mei J, Shi Y, He M, Wang G, Liang J, Wang Z, Cao D, Carter MT, Chrysler C, Drmic IE, Howe JL, Lau L, Marshall CR, Merico D, Nalpathamkalam T, Thiruvahindrapuram B, Thompson A, Uddin M, Walker S, Luo J, Anagnostou E, Zwaigenbaum L, Ring RH, Wang J, Lajonchere C, Wang J, Shih A, Szatmari P, Yang H, Dawson G, Li Y, Scherer SW. Detection of clinically relevant genetic variants in autism spectrum disorder by whole-genome sequencing. *Am J Hum Genet*. 2013;93(2):249–63.
54. Kong A, Frigge ML, Masson G, Besenbacher S, Sulem P, Magnusson G, Gudjonsson SA, Sigurdsson A, Jonasdottir A, Jonasdottir A, Wong WSW, Sigurdsson G, Bragi Walters G, Steinberg S, Helgason H, Thorleifsson G, Gudbjartsson DF, Helgason A, Magnusson OT, Thorsteinsdottir U, Stefansson K. Rate of de novo mutations and the importance of father's age to disease risk. *Nature*. 2012;448(7412):471–5.
55. Zhou J, Park CY, Theesfeld CL, Wong AK, Yuan Y, Scheckel C, Fak JJ, Funk J, Yao K, Tajima Y, Packer A, Darnell RB, Troyanskaya OG. Whole-genome deep-learning analysis identifies contribution of noncoding mutations to autism risk. *Nat. Genet*. 2019;51(6):973–980.
56. Paquet A, Olliac B, Bouvard M-P, Golse B, Vaivre-Douret L. The semiology of motor disorders in autism spectrum disorders as highlighted from a standardized neuro-psychomotor assessment. *Front Psychol*. 2016;7:1292.
57. Xia B, Yan Y, Baron M, Wagner F, Barkley D, Chiodin M, Kim SY, Keefe DL, Alukal JP, Boeke JD, et al. Widespread transcriptional scanning in the testis modulates gene evolution rates. *Cell*. 2020;180(2):248–62.
58. Van Nostrand EL, Freese P, Pratt GA, Wang X, Wei X, Xiao R, Blue SM, Chen J-Y, Cody NA, Dominguez D, et al. A large-scale binding and functional map of human rna-binding proteins. *Nature*. 2020;583(7818):711–9.
59. Jha A, Gazzara MR, Barash Y. Integrative deep models for alternative splicing. *Bioinformatics*. 2017;33(14):274–82.
60. Julien P, Miñana B, Baeza-Centurion P, Valcárcel J, Lehner B. The complete local genotype-phenotype landscape for the alternative splicing of a human exon. *Nat Commun*. 2016;7:11558.
61. Soemedi R, Cygan KJ, Rhine CL, Wang J, Bulacan C, Yang J, Bayrak-Toydemir P, McDonald J, Fairbrother WG. Pathogenic variants that alter protein code often disrupt splicing. *Nat Genet*. 2017;49(6):848–55.
62. Ke S, Anquetil V, Zamalloa JR, Maity A, Yang A, Arias MA, Kalachikov S, Russo JJ, Ju J, Chasin LA. Saturation mutagenesis reveals manifold determinants of exon definition. *Genome Res*. 2018;28(1):11–24.
63. Adamson SI, Zhan L, Graveley BR. Vex-seq: high-throughput identification of the impact of genetic variation on pre-mRNA splicing efficiency. *Genome Biol*. 2018;19(1):71.
64. Mikl M, Hamburg A, Pilpel Y, Segal E. Dissecting splicing decisions and cell-to-cell variability with designed sequence libraries. *Nat Commun*. 2019;10(1):4572.
65. Linder J, Bogard N, Rosenberg AB, Seelig G. Deep exploration networks for rapid engineering of functional DNA sequences. *BioRxiv*. 2019864363. <https://doi.org/10.1101/864363>.
66. Kingma DP, Ba J. Adam: a method for stochastic optimization. 2014. arXiv preprint arXiv:1412.6980.
67. He K, Zhang X, Ren S, Sun J. Delving deep into rectifiers: surpassing human-level performance on ImageNet classification. 2015. <https://doi.org/10.1109/iccv.2015.123>.
68. Bergstra J, Komer B, Eliasmith C, Yamins D, Cox DD. Hyperopt: a Python library for model selection and hyperparameter optimization. 2015. <https://doi.org/10.1088/1749-4699/8/1/014008>.
69. Patton JG, Mayer SA, Tempst P, Nadal-Ginard B. Characterization and molecular cloning of polypyrimidine tract-binding protein: a component of a complex necessary for pre-mRNA splicing. *Genes Dev*. 1991;5(7):1237–51.
70. Ule J, Stefani G, Mele A, Ruggiu M, Wang X, Taneri B, Gaasterland T, Blencowe BJ, Darnell RB. An RNA map predicting Nova-dependent splicing regulation. *Nature*. 2006;444(7119):580–586.
71. Konieczny P, Stepniak-Konieczna E, Sobczak K. MBNL proteins and their target RNAs, interaction and splicing regulation. *Nucleic Acids Res*. 2014;42(17):10873–87.
72. Xue Y, Zhou Y, Wu T, Zhu T, Ji X, Kwon Y-S, Zhang C, Yeo G, Black DL, Sun H, et al. Genome-wide analysis of ptb-rna interactions reveals a strategy used by the general splicing repressor to modulate exon inclusion or skipping. *Mol Cell*. 2009;36(6):996–1006.

73. Giudice G, Sánchez-Cabo F, Torroja C, Lara-Pezzi E. ATTRACT-a database of RNA-binding proteins and associated motifs. Database. 2016;2016:. <https://doi.org/10.1093/database/baw035>.
74. Ray D, Kazan H, Cook KB, Weirauch MT, Najafabadi HS, Li X, Gueroussov S, Albu M, Zheng H, Yang A, et al. A compendium of RNA-binding motifs for decoding gene regulation. *Nature*. 2013;499(7457):172–7.
75. Cook KB, Kazan H, Zuberi K, Morris Q, Hughes TR. RbpdB: a database of RNA-binding specificities. *Nucleic Acids Res*. 2010;39(suppl\_1):301–8.
76. Gupta S, Stamatoyannopoulos JA, Bailey TL, Noble WS. Quantifying similarity between motifs. *Genome Biol*. 2007;8(2):24.
77. Karczewski KJ, Francioli LC, Tiao G, Cummings BB, Alföldi J, Wang Q, Collins RL, Laricchia KM, Ganna A, Birnbaum DP, et al. The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature*. 2020;581(7809):434–43.
78. Cheng J, Çelik MH, Kundaje A, Gagneur J. MTSplice predicts effects of genetic variants on tissue-specific splicing. GitHub. [https://github.com/gagneurlab/MMSplice\\_MTSplice](https://github.com/gagneurlab/MMSplice_MTSplice). Accessed 3 Oct 2020.
79. Cheng J, Çelik MH, Kundaje A, Gagneur J. MTSplice predicts effects of genetic variants on tissue-specific splicing. GitHub. <http://kipoi.org/models/MMSplice/mtsplice/>. Accessed 8 June 2020.
80. Cheng J, Çelik MH, Kundaje A, Gagneur J. MTSplice predicts effects of genetic variants on tissue-specific splicing. Zenodo. <https://doi.org/10.5281/zenodo.4255942>. Accessed 7 Nov 2020.
81. Cheng J, Çelik MH, Kundaje A, Gagneur J. MTSplice predicts effects of genetic variants on tissue-specific splicing. GitLab. <https://gitlab.cmm.in.tum.de/gagneurlab/mtsplice>. Accessed 7 June 2020.

### Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

**Ready to submit your research? Choose BMC and benefit from:**

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

**At BMC, research is always in progress.**

Learn more [biomedcentral.com/submissions](https://biomedcentral.com/submissions)

