



Deciphering RNA splicing logic with interpretable machine learning

Susan E. Liao^{a,1}, Mukund Sudarshan^{a,1}, and Oded Regev^{a,2}

Edited by Bonnie Berger, Massachusetts Institute of Technology, Cambridge, MA; received January 6, 2023; accepted August 29, 2023

Machine learning methods, particularly neural networks trained on large datasets, are transforming how scientists approach scientific discovery and experimental design. However, current state-of-the-art neural networks are limited by their uninterpretability: Despite their excellent accuracy, they cannot describe how they arrived at their predictions. Here, using an “interpretable-by-design” approach, we present a neural network model that provides insights into RNA splicing, a fundamental process in the transfer of genomic information into functional biochemical products. Although we designed our model to emphasize interpretability, its predictive accuracy is on par with state-of-the-art models. To demonstrate the model’s interpretability, we introduce a visualization that, for any given exon, allows us to trace and quantify the entire decision process from input sequence to output splicing prediction. Importantly, the model revealed uncharacterized components of the splicing logic, which we experimentally validated. This study highlights how interpretable machine learning can advance scientific discovery.

RNA splicing | interpretable machine learning | artificial intelligence

Machine learning algorithms, in particular neural networks, capture complex quantitative relationships between input and output. However, as neural networks are typically black boxes, it is difficult to extract post hoc insights into how they achieve their predictive success. Furthermore, they easily capture artifacts or biases in the training data, often fail to generalize beyond the datasets used for training and testing, and do not lead to insights into the underlying processes (1).

In recent years, neural networks have been used to tackle challenging biological questions. One outstanding question in genomics is understanding the regulatory logic of RNA splicing, which plays a critical role in the fundamental transfer of information from DNA to functional RNA and protein products. Splicing removes introns and ligates exons together to form mature RNA transcripts. While some canonical sequence features are necessary for exon definition (splice sites delimiting exons and branch points used during intron removal), exon definition is also facilitated by exon sequence (2, 3). Despite recent success using neural networks to predict splicing outcomes (4, 5), understanding how exon sequence dictates inclusion or skipping remains an open challenge. The challenge is further underscored by the sensitivity of splicing logic, where almost all single-nucleotide changes along an exon can lead to dramatic changes in splicing outcomes (6, 7).

To enable scientific progress, machine learning models should not only accurately predict outcomes but also describe how they arrive at their predictions. Here, we demonstrate that an “interpretable-by-design” model achieves predictive accuracy without sacrificing interpretability, captures a unifying decision-making logic, and reveals previously uncharacterized splicing features.

Results

Generating a Synthetic Dataset for Interpretable Machine Learning. As neural network performance and interpretability are inextricable from the data it is trained on, we began by generating a large, high-quality synthetic splicing dataset. The use of synthetic datasets offers several advantages over genomic data used in previous work. First, genomic datasets are limited by the number of exons in the genome. In contrast, synthetic assays can dramatically increase the number of data points by orders of magnitude (8, 9). Second, genomic exons are flanked by varying sequences (splice sites, introns, promoters) that also participate in splicing decisions (10), greatly complicating attempts at interpretability. In contrast, synthetic datasets fix all but one variable region, allowing one to focus on the region of interest. Third, genomic exons contain overlapping RNA codes (e.g., protein-coding sequences). In contrast, sequences in synthetic datasets are devoid of

Significance

Machine learning approaches are increasingly applied to advancing discovery in the biological sciences. However, despite achieving predictive accuracy, many machine learning models cannot explain how they achieve their predictive success. Here, we demonstrate that bespoke data generation coupled with model design that infuses foundational biological knowledge enables an “interpretable-by-design” approach that advances our understanding of RNA splicing. Our model not only accurately predicts the quantitative splicing outcomes but also explains how specific combinations of RNA features dictate splicing outcomes. We validate the network predictions and interpretations through additional data generation and experimental validation. These results demonstrate that “interpretable-by-design” machine learning represents a powerful approach to harnessing the potential of machine learning toward advancing our understanding of biological processes.

The authors declare no competing interest.

This article is a PNAS Direct Submission.

Copyright © 2023 the Author(s). Published by PNAS. This open access article is distributed under [Creative Commons Attribution License 4.0 \(CC BY\)](https://creativecommons.org/licenses/by/4.0/).

¹S.E.L. and M.S. contributed equally to this work.

²To whom correspondence may be addressed. Email: regev@cims.nyu.edu.

This article contains supporting information online at <https://www.pnas.org/lookup/suppl/doi:10.1073/pnas.2221165120/-/DCSupplemental>.

Published October 5, 2023.

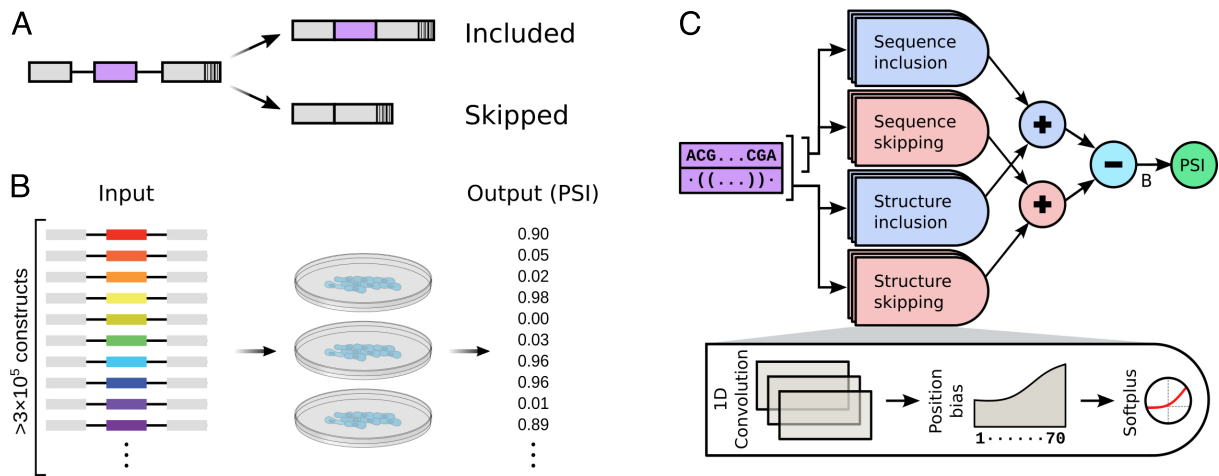


Fig. 1. Data generation and interpretable-by-design machine learning model. (A) All reporters in the assay share the same three-exon design and differ only in their middle exon, which contains a random 70-nucleotide-long sequence. Depending on its sequence, an exon might be included, skipped, or a probabilistic mix of the two. Each reporter includes a unique barcode at the end of the third exon so that exon identity can be inferred in exon-skipping products. (B) The assay includes over 3×10^5 different reporters. The reporters were transfected into HeLa cells in a pooled fashion in three biological replicates. High-throughput sequencing then provides a “percent spliced in” (PSI) value to each reporter. (C) The machine learning model consists of both short convolution filters (applied to exon sequence only) and long convolution filters (applied to both exon sequence and predicted structure). The output of these filters (strength) can depend on the position along the exon. Half of the filters are designated as inclusion filters, and the rest are skipping filters. Predicted PSI is computed from the difference between the total strength of inclusion filters and the total strength of skipping filters, after adding an initial basal strength (B).

overlapping codes by design. In summary, from both a quantity and quality perspective, synthetic datasets provide crucial advantages for machine learning over genomic datasets.

The synthetic dataset we generated includes hundreds of thousands of input–output data points. Each data point is a different random 70-nucleotide exon sequence, paired with a measured percent spliced in (PSI) output, which is a number between 0 (always skipped) and 1 (always included) (Fig. 1A). The dataset is generated by a massively parallel reporter assay that allows for PSI quantification for hundreds of thousands of unique sequences in a single biological experiment (Fig. 1B). Splicing outcomes for the parallel reporter assay were measured after transfection into human HeLa cells using high-throughput sequencing. We confirmed that reporters are evenly represented in the reporter assay (*SI Appendix, Fig. S1A*). The vast majority of splicing products corresponded to exon inclusion or exon-skipping products (*SI Appendix, Fig. S1B*), and we filtered our data to exclude spurious splicing products. PSI values are calculated as the number of inclusion reads divided by the total number of inclusion and skipping reads. Three biological replicates of the assay showed excellent agreement (*SI Appendix, Fig. S1C*), and their sequencing results were combined for all downstream analyses. High-throughput sequencing measurements were consistent with semiquantitative measurements of individual reporters (*SI Appendix, Fig. S1D*).

An Interpretable-by-Design Model Accurately Predicts Splicing Outcomes. We first compared the predictive accuracy of three off-the-shelf machine learning algorithms on our dataset: a k -mer splicing scoring algorithm (8) (*SI Appendix, Fig. S2A*), a gated recurrent unit neural network (11) (*SI Appendix, Fig. S2B*), and a transformer neural network (12) (*SI Appendix, Fig. S2C*). We found that the two neural networks outperformed the k -mer scoring algorithm. This gap in predictive accuracy suggests that the more complex neural networks capture additional features affecting splicing outcomes. However, as these models are not interpretable, we were unable to pinpoint which specific features contribute to the improved predictions.

We therefore designed a neural network model with the explicit goal of being interpretable (13). The predictive accuracy of our interpretable-by-design model is comparable to that of the two state-of-the-art neural networks (*SI Appendix, Fig. S2D*). This suggests that interpretability need not come at the expense of accuracy.

In addition to our own dataset, the model accurately predicts splicing outcomes from other splicing datasets (7, 8, 14–17) (*SI Appendix, Table S1 and Fig. S3*). Importantly, unlike our random exons, these datasets were modeled on specific genomic exons, with each dataset differing in splice sites, introns, and flanking exons. Furthermore, these datasets were generated in different immortalized cell lines. Encouragingly, despite these dramatic differences in RNA architecture and cell types, our model performed well on these datasets, suggesting that our model generalizes and captures critical aspects of splicing regulatory logic.

Model Architecture Reveals Unifying Decision-Making Process. Our interpretable-by-design model incorporates domain knowledge throughout its architecture (Fig. 1C). Specifically, we reasoned that short six-nucleotide sequence filters would capture motifs previously demonstrated to play an important role in splicing decisions (18, 19). We therefore introduced one-dimensional convolutional filters applied to the input RNA sequence. Next, since RNA secondary structure was previously implicated in splicing outcomes (16, 20), we also provided the network with predicted structure (21). We then introduced longer (30-nucleotide) one-dimensional convolutional filters to the structure-augmented sequence. Crucially, while we fixed filter lengths using minimal domain knowledge, we did not explicitly specify sequences and structures, allowing the network flexibility to learn filters in an unbiased manner. Furthermore, our model explicitly quantifies the strength (in network-defined arbitrary units) of each activated filter to the inclusion or skipping decision. Importantly, we allowed the strength of any filter to vary along the length of an exon, providing the network the flexibility to capture position-dependent effects of RNA features on splicing outcomes.

To arrive at its output, the network computes the difference in the sum total of exon inclusion strengths and exon skipping strengths (Δ strength), which is then converted to predicted PSI. The greater the magnitude of this difference, the closer the PSI is to 0 (difference \ll 0) or 1 (difference \gg 0). This additive combinatorial behavior is consistent with the previous literature (8, 22).

Model Extends Understanding of Splicing Regulatory Logic. Even though our model was trained on a synthetic dataset, it recapitulates and extends domain knowledge from previous genomic and biochemical studies.

Many filters in the model match binding motifs of RNA-binding proteins implicated in splicing regulation (splicing factors) (23, 24) (Fig. 2 and *SI Appendix, Fig. S4*). Consistent with previous studies, network inclusion filters match binding sites for SR proteins known to promote exon inclusion (26, 27), whereas network skipping filters match binding sites for hnRNP proteins known to promote exon skipping (28).

However, while the directionality of these RNA features toward splicing was established, their magnitude was not clear. Importantly, the model addresses this issue by assigning a quantitative strength to each filter. Moreover, some filters exhibit striking position-dependent strengths, suggesting that the position of an RNA feature along an exon affects its strength. This is consistent with previous experimental reports demonstrating position-dependent effects of RNA sequences within exons (8).

Surprisingly, our network accurately predicted splicing outcomes using a concise list of filters (Fig. 2). This contrasts with previous studies suggesting that splicing outcomes result from the combinatorics of hundreds of unique RNA features (8, 29, 30).

Using the local interpretability of our model, we introduce a visualization (balance plot) that enables explicit examination and quantification of how multiple RNA features lead to splicing outcomes for any given exon from our dataset (Fig. 3, *SI Appendix, Fig. S6*) and other datasets (*SI Appendix, Fig. S5*). For a given exon, the total strengths of activated filters are represented as bars of the appropriate height. Total inclusion

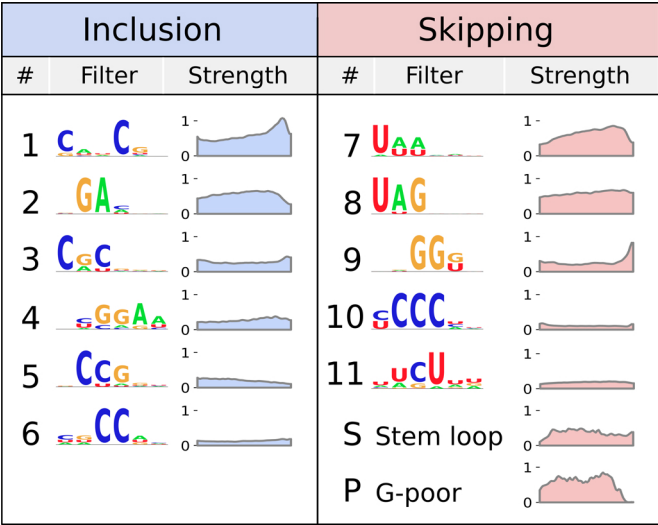


Fig. 2. Model expands on known splicing logic. Splicing features detected by the model's filters, represented by their sequence logo (25). Filters either contribute to inclusion (blue) or skipping (red). Plots show the average strength in our dataset of each filter as a function of position along the exon. The model also identified short stem loops and long G-poor stretches as contributing to exon skipping.

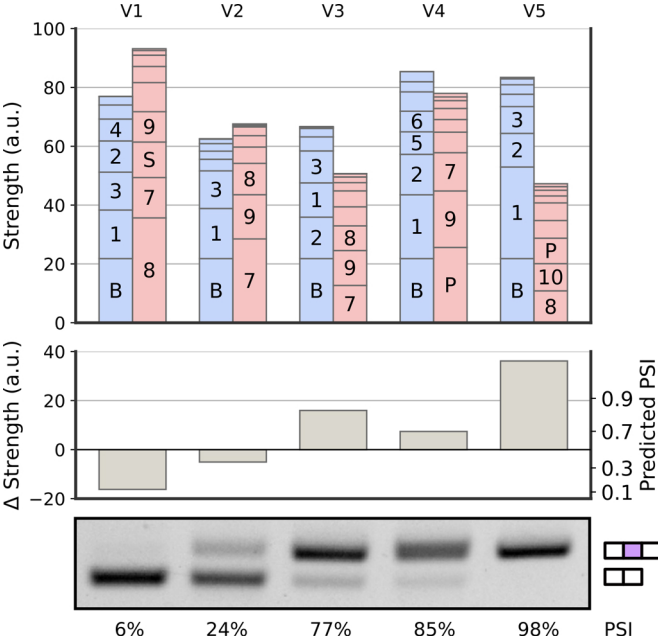


Fig. 3. Model predictions can be interpreted using balance plots. Balance plots used to visualize the logic leading to PSI prediction for five randomly picked exons (V1–V5). Bar plot showing the total strength contributed by each filter (Top). Bars are labeled by filter numbers from Fig. 2. The bar labeled B represents a fixed initial basal strength. Labels are not shown for smaller bars. The difference between total inclusion and total skipping strengths (Δ strength) leads to predicted PSI (Middle). PSI as measured by semiquantitative RT-PCR matches the machine learning predictions (Bottom).

strength (blue) and skipping strength (red) are then visible as the height of the stacked bars. The Δ strength is represented by the difference in heights between the stacked inclusion and skipping filters. These visualizations provide an intuitive tool to understand the contributions of individual sequence and structure features leading to each exon's predicted PSI. They emphasize that splicing logic results from contributions of many RNA features along the exon and that a single nucleotide can be part of multiple overlapping filters (6, 8).

Identification and Validation of Uncharacterized Splicing Features. Next, we asked whether our interpretable-by-design model could identify uncharacterized splicing features. While most network filters were consistent with previously described splicing features, two uncharacterized long skipping filters with strong influence on splicing predictions stood out (Fig. 2). We confirmed that these filters were robustly identified across multiple initialization seeds and training/testing splits, suggesting that they are not training artifacts. We then turned our attention to characterizing and experimentally validating these features.

Examining the first uncharacterized filter revealed that it identifies stem loop structures with short, GC-rich, 5–7 nucleotide double-stranded regions (Fig. 4). Next, we experimentally validated that these stem loops contribute to exon skipping and are not artifacts. We introduced mutations that disrupt double-stranded base pairing in an exon with such a stem loop. First, we introduced single-nucleotide mutations predicted to abolish the stem by disrupting base pairing. Notably, these mutations were designed to minimize disruptions of other filters, ensuring that prediction differences are mainly due to altered secondary structure, and not due to the introduction or disruption of other sequence features. In addition to two such mutations, we also introduced both compensatory mutations

Discussion

In this study, we demonstrate that an interpretable-by-design model advanced scientific discovery. Our model accurately predicts splicing outcomes on both our assay and on previously published assays, demonstrating that interpretability need not come at the expense of accuracy or generalizability. Model interpretability enabled a systematic understanding of RNA splicing logic, including the identification of two candidate exon-skipping features which were subsequently experimentally validated. The model's ability to quantify contributions of specific features to splicing outcomes for individual exons has considerable potential for a range of medical and biotechnology applications, including genome- or RNA-editing of target exons to correct splicing behavior or guiding rational design of RNA-based therapeutics like antisense oligonucleotides (33).

In addition, model-identified features hint at novel biochemical mechanisms that warrant further study. For example, the fact that splicing decisions are modeled well by an additive quantity (Δ strength) supports a biochemical mechanism involving the nuclear spatial organization of SR and hnRNP proteins (34). Furthermore, our model identified two unusual exon-skipping features. These features may be recognized by an uncharacterized RNA-binding protein or complex. Alternatively, the introduction of highly structured or unstructured regions may change the physical distance between splice sites, enhancing exon skipping. These open questions further underscore how interpretable-by-design models can advance scientific discovery by aiding hypothesis generation.

Our model performs well on synthetic datasets from immortalized cell lines, yet further work is needed to capture the dynamics of developmentally regulated splicing logic (35–37). Importantly, splicing outcomes change depending on the expression level of cell type-specific RNA-binding proteins (38). These questions can be addressed by generation of additional synthetic splicing datasets in developmentally relevant cell types paired with interpretable-by-design models that capture cell type-specific regulatory features.

Beyond the context of splicing, the interpretable-by-design framework can be used to decipher the multiple, complex, and overlapping codes that dictate biomolecular processing. Importantly, many rich synthetic datasets that address RNA untranslated 5' (39) and 3' (40) region regulation, methylation (41), and small RNA biogenesis (42) have already been generated. We expect that additional data generation efforts paired with the interpretable-by-design framework will stimulate advances in understanding biological codes more broadly.

Materials and Methods

Reporter Assay Design and Cloning. The splicing reporter is based on a three-exon beta globin minigene (43) under the control of a truncated mammalian CAG promoter. The massively parallel splicing assay allows for high-throughput characterization of exon variants on splicing outcomes (44) using Gibson assembly and ligation cloning. The assay replaces the middle beta globin exon with 70-nt random sequences flanked by weak splice sites [MaxEnt scores (45): 3'ss 9.41, 5'ss 5.06]. Each 70-nt exon is coupled with a 20-nt barcode downstream of the third exon, allowing for identification of middle exon identity in exon-skipping products. Random exons and barcodes were synthesized separately as degenerate single-stranded oligonucleotides (IDT) and were joined using an anneal-extend procedure as follows. In a 100 μ L reaction (Phusion® Hot Start Flex 2X Master Mix, NEB), 200 nM exon and barcode oligonucleotides were denatured at 98 °C for 10 min, cooled slowly to 60 °C (0.1 °C/s), annealed at 60 °C for 5 min, and extended at 72 °C for 60 min. Single-stranded products were removed from pooled double-stranded exon-

barcode using a silica column purification according to the manufacturer's specifications (ZymoPURE Plasmid Miniprep Kit). Pooled exon-barcode products were cloned into a backbone digested with BsmBI and XbaI and expanded using electrocompetent bacterial cells (ElectroMAX™ DH10B Cells, ThermoFisher) on large solid agar Bioassay plates (Nunc™ Square BioAssay Dishes, ThermoFisher). After resuspending pooled bacteria in 1X PBS, DNA was recovered using silica column purification (ZymoPURE II Plasmid Maxiprep Kit, Zymo Research) following the manufacturer's specifications. The resulting pooled library (Lib1) includes the truncated CAG promoter, followed by the first minigene exon and intron, and the exon-barcode insertion. High-throughput amplicon sequencing of Lib1 was used to match exon-barcode pairs. To generate the final splicing reporter assay (Lib2), a fixed sequence, containing the second intron and third exon, was introduced to separate exons from their barcodes. Lib1 was digested with Esp3I (NEB) to introduce overhangs between the exons and barcodes; the digested product was gel-purified to facilitate downstream cloning (Zymoclean Gel DNA Recovery Kits). A segment containing the second intron and third exon was ligated into the digested Lib1 product (NEB Quick Ligation). Lib2 library was expanded using electrocompetent bacteria cells resulting in about 10 times as many colonies as Lib1 to ensure even representation across reporters and recovered using silica column purification as described for Lib1. DNA was quantified using a spectrophotometer (NanoDrop™ One^C, Fisher Scientific).

Individual Reporter Cloning. To validate consequences of point mutations on splicing outcomes, individual exons were synthesized as two single-stranded oligonucleotides (IDT DNA Technologies) and joined using an anneal-extend procedure. Briefly, 200 nM of each oligonucleotide were joined in a 100 μ L reaction with 5U DNA polymerase (NEB Klenow). Oligonucleotides were denatured at 98 °C for 10 min, annealed after cooling slowly to 25 °C (1 °C/s), and extended at 25 °C for at least 2 h. Reactions were heat inactivated at 75 °C for 20 min and used directly for Gibson assembly into a digested receiving plasmid with a fixed barcode.

Cell Culture. HeLa cells (ATCC) were grown in high-glucose DMEM supplemented with 10% fetal bovine serum and penicillin and streptomycin (ThermoFisher). All cells were grown at 37 °C, 5% CO₂, and 95% relative humidity.

Transfection, RNA Extraction, and Reverse Transcription. Cells were transfected at 60 to 80% confluence with FuGENE HD® according to the manufacturer's protocol at a 3:1 FuGENE HD® to DNA ratio. For high-throughput measurements of splicing outcomes, 10 μ g pooled reporter assay DNA was transfected in three 100-mm plates. For biochemical analysis of individual reporters, 1 μ g or 2.5 μ g individual reporter DNA was transfected into each well of a 12- or 6-well plate (respectively). Then, 24 h after transfection, total RNA was isolated from detached cells (Accutase®, ThermoFisher). For amplicon sequencing, total RNA was isolated using phenol-chloroform (Ambion) extraction (5PRIME Phase Lock Gel, Quantabio) followed by DNase treatment (TURBO DNase). For biochemical analysis, RNA was isolated using a silica column (illustra™ RNAspin Mini RNA Isolation Kit, GE Healthcare) with on-column DNase digestion following the manufacturer's automated protocol. DNase-treated RNA was reverse transcribed using a reporter-specific primer following the manufacturer's specifications (SuperScript IV Reverse Transcriptase, Thermo Fisher) with RNase H treatment. Reverse transcription primers included degenerate nucleotides to serve as unique molecular identifiers (UMIs) during amplicon sequencing (46, 47). cDNA products were used for amplicon sequencing or biochemical analysis.

Amplicon Sequencing. Amplicon sequencing was used to identify exon-barcode pairings in Lib1 and to quantify splicing products from reverse-transcribed cDNA. Second-strand synthesis added additional UMIs in a single anneal-extend cycle of 98 °C for 10 min, cooled slowly to 60 °C (0.1 °C/s), annealed at 60 °C for 5 min, and extended at 72 °C for 5 min (Phusion® Hot Start Flex 2X Master Mix, NEB). Resulting double-stranded amplicons were amplified using a two-stage procedure. In the first stage, targets were amplified by PCR primers. PCR was performed using the following protocol: 98 °C for 30-s initial denaturation, then 16 cycles of 98 °C denaturation for 10 s, 60 °C

annealing for 15 s, 72 °C extension for 1 min 45 s, and a final extension step at 72 °C for 5 min (Phusion®Hot Start Flex 2X Master Mix, NEB). Longer extension times and a minimal number of PCR cycles were used to avoid recombination across exons and barcodes. The number of cycles was determined for each sample by first running 10 µL qPCR reactions (LightCycler®480 SYBR Green I Master, Roche). In the second stage, index primers were added using 5 PCR cycles. PCR was performed using the following protocol: 98 °C for 30 s initial denaturation, then 5 cycles of 98 °C denaturation for 10 s, 71 °C annealing for 15 s, 72 °C extension for 1 min 45 s, and a final extension step at 72 °C for 5 min (Phusion®Hot Start Flex 2X Master Mix, NEB). Final DNA concentrations were measured using fluorometric measurements (Qubit 1X dsDNA HS Assay, Thermo Fisher) on a Qubit™ 3 Fluorometer. Paired-end sequencing was carried out on an Illumina NextSeq 550 with 10% PhiX spiked in, with 54 cycles in read 1 (reverse) and 106 in read 2 (forward). About 4M paired-end reads (> 10X coverage) were acquired for Lib1 exon-barcode sequencing and an average of 22M paired-end reads (> 50X coverage) for each PSI quantification replicate.

Biochemical Analysis. PCR amplification reactions to determine splicing products were carried out in 20 µL reactions containing 10 µL OneTaq®2X Master Mix with Standard Buffer (NEB), 200 nM each forward and reverse primers (IDT), and 1 µL cDNA. PCR was performed using the following protocol: 94 °C for 30-s initial denaturation, then 25 cycles of 94 °C denaturation for 10 s, 62 °C annealing for 15 s, 68 °C extension for 20 s, and a final extension step at 68 °C for 1 min. Then, 5 µL final PCR product was run out on 2.0% agarose (Denville Scientific) Tris-acetate-EDTA (TAE) gel with ethidium bromide and visualized on a Bio-Rad imager. Densitometry measurements to calculate PSI were measured using Bio-Rad Image Lab (Windows v6.1).

Reporter Assay Preprocessing. The list of all exons in the reporter assay with their corresponding barcodes was extracted from DNA sequencing of Lib1. To ensure unique coupling of barcodes to exons, barcodes appearing with more than one exon sequence were filtered out. This step ignored exon sequences appearing only once, as those are likely due to sequencing errors. Barcodes with fewer than two DNA reads in total were also filtered out.

Next, splicing outcomes were extracted from RNA sequencing of each of the three replicate transfections of Lib2. For each replicate, each read was identified by barcode and was assigned a splicing outcome (exon skipping, exon inclusion, intron retention, splicing inside exon, or unknown splicing). Carryover from Lib1 was filtered out, as were reads for which exon 1 could not be identified. Using unique molecular identifiers (UMIs), the fraction of duplicate reads in each replicate was estimated to be below 23%. The counts from all three replicates were merged for downstream analysis. Barcodes with fewer than 60 total reads, barcodes that contained an Esp3I restriction site in either strand of the exon or its barcode, and barcodes where inclusion and skipping made up less than 80% of all reads were filtered out.

Finally, the dataset was generated by computing PSI for each barcode as

$$\text{PSI} = \frac{n_{\text{inclusion}}}{n_{\text{skipping}} + n_{\text{inclusion}}},$$

where $n_{\text{inclusion}}$ is the total number of exon inclusion reads, and n_{skipping} is the total number of exon skipping reads. In addition to the measured PSI, the dataset includes for each barcode: 1) a 90-nucleotide sequence, containing the 70-nucleotide variable exon sequence plus the 10 fixed flanking nucleotides on each side; 2) structure in dot-bracket notation predicted by RNAFold [Vienna RNA (21), version 2.4.17], using default parameters; 3) an indicator vector indicating which nucleotide participates in a predicted G-U wobble base pair. The dataset was split randomly into a training set and a test set in an 80/20 split, using a fixed seed for reproducibility.

Model Design. The model's input is a triple of vectors (x_{seq} , x_{struct} , x_{wobble}),

$$\begin{aligned} x_{\text{seq}} &\in \{A, C, G, U\}^d && [4\text{-category sequence input}] \\ x_{\text{struct}} &\in \{(\cdot, \cdot)\}^d && [3\text{-category structure input}] \\ x_{\text{wobble}} &\in \{0, 1\}^d, && [\text{wobble pair indicator input}], \end{aligned}$$

where $d = 90$. The neural network contains four "strength-computation modules" (SCM) defined as

$$f_a^b : x \mapsto \text{Sum}(\text{Softplus}(\text{Position-Bias}(\text{Convolution}(x; \alpha_a^b; \beta_a^b))) \quad [\text{SCM}]$$

$$\alpha_a^b \in \mathbb{R}^{w_a^b \times c_a^b \times k_a^b}, \quad \beta_a^b \in \mathbb{R}^{(d-w_a^b+1) \times k_a^b},$$

where $a \in \{\text{incl}, \text{skip}\}$, and $b \in \{\text{seq}, \text{struct}\}$. The input is either $x = [x_{\text{seq}}]$ (sequence SCM) or $x = [x_{\text{seq}}, x_{\text{struct}}, x_{\text{wobble}}]$ (structure SCM), and the output is a scalar. The 1D convolutional layer contains $k_a^b = 20$ convolutional filters of width $w_a^b = 6$ for each sequence SCM ($b = \text{seq}$), and $k_a^b = 8$ convolutional filters of width $w_a^b = 30$ for each structure SCM ($b = \text{struct}$). The number of input channels is $c_a^b = 4$ for sequence SCM (corresponding to the one-hot encoded four nucleotides) and $c_a^b = 8$ for structure SCM (corresponding to sequence, structure, and wobble indicator). The output of the convolution layer is a $(d - w_a^b + 1) \times k_a^b$ matrix z of "raw" strengths. The position bias layer maps inputs z to $z + \beta_a^b$, adjusting the raw strengths based on position along the exon. Finally, each position-adjusted raw strength is passed through a softplus activation, and the resulting strengths are all summed up to form the output of the SCM f_a^b .

The splicing prediction model $m(x_{\text{seq}}, x_{\text{struct}}, x_{\text{wobble}}; \theta)$ is then defined as

$$\begin{aligned} m(x_{\text{seq}}, x_{\text{struct}}, x_{\text{wobble}}; \theta) &= \text{Tuner}(f_{\text{incl}}^{\text{seq}}([x_{\text{seq}}]) + f_{\text{incl}}^{\text{struct}}([x_{\text{seq}}, x_{\text{struct}}, x_{\text{wobble}}])) \\ &\quad - f_{\text{skip}}^{\text{seq}}([x_{\text{seq}}]) - f_{\text{skip}}^{\text{struct}}([x_{\text{seq}}, x_{\text{struct}}, x_{\text{wobble}}]); \gamma). \quad [1] \end{aligned}$$

This model computes the total strength for inclusion and for skipping and uses their difference to predict splicing outcomes. The function $\text{Tuner}(\cdot; \gamma) : \mathbb{R} \rightarrow [0, 1]$ is a learned nonlinear activation function that maps this difference to a PSI prediction. It consists of a 3-layer fully connected network with a residual connection from the input to the output layer, followed by a sigmoid activation. The parameter set θ contains the parameters of all SCMs and the parameter γ .

Model Training. The model was implemented in Python 3.8 (48) using Tensorflow 2.6 (49) and Numpy 1.20 (50). Batched gradient descent was used to optimize the model's parameters using the Adam optimizer, with KL divergence as the loss function. Hyperparameters such as regularization parameters were tuned with grid search. Training the model took about 2 h on a mid-range 4-core with 16 GB of RAM.

To improve interpretability, the model was trained in steps (custom training schedule), progressively adding learnable parameters in each step. First, a simplified model given by

$$\text{Tuner}'(f_{\text{incl}}^{\text{seq}}([x_{\text{seq}}]) - f_{\text{skip}}^{\text{seq}}([x_{\text{seq}}]); \nu, \eta)$$

was trained. Here, $\text{Tuner}'(\cdot; \nu, \eta) : \mathbb{R} \rightarrow [0, 1]$ is a learned nonlinear activation function defined by $x \mapsto \sigma(\nu x + \eta)$, where σ is the sigmoid function, and ν and η are two real parameters. This step ensures that short sequence motifs are captured by the sequence SCMs and not the more complex structure SCMs. In the second step, the structure SCMs were added, leading to a model identical to the final one Eq. 1, except for the use of Tuner' instead of Tuner . The sequence SCM weights were initialized to those of the previous model. In the third and last step, the Tuner function was introduced, leading to the final model Eq. 1. SCM weights were initialized to those of the previous model.

To further improve the model's interpretability, regularization terms were added. First, to obtain a concise list of filters, an activity regularization loss term was used. The term consists of the ℓ_1 norm of all the strengths. Second, a smoothness regularization loss term was applied to position bias layer weights. This term consists of the ℓ_2 norm of the discrete derivative of the weight vectors (defined as the difference between the vector and itself shifted by one along the sequence dimension). Each of the two loss terms was multiplied by a hyperparameter.

Hyperparameters were optimized based on two criteria: held-out KL divergence and sparsity of activations. Sparsity was measured as the minimum number of activations needed per exon to achieve KL below a threshold. Among all hyperparameters leading to sufficiently high accuracy and sparsity, the one with the highest smoothness regularization was chosen.

Prediction Accuracy on Other Assays. Exon sequences and PSI measurements were obtained from previous publications. Exons including indel mutations or differing from WT sequence in the first or last three nucleotides were filtered out. In order to generalize the model to exons of varying lengths, we applied Lanczos resampling (with parameter $a = 3$) to the position bias weights. To account for differences in splice sites, flanking sequences, and cell types, one scalar correction term was introduced per assay, effectively adjusting the basal strength (B) (17).

Filter Visualization. To avoid reporting redundant sequence filters, hierarchical clustering using SciPy (51) was applied. Each sequence filter was represented by a vector containing its total strength for each of the exons in the dataset. The strongest filter in each cluster was then used to generate a sequence logo (25). The logo represents the set of 6-mers that lead to positive filter activation.

The structure filters included one G-poor filter and three stem loop filters. Since enumerating all 30-mers is not tractable, the G-poor sequence logo was computed by evaluating the filter on a subset of sequences from our dataset. As the three stem loop filters differed in the length of the loop (short, medium, long) but were otherwise very similar, they were considered as one cluster. Layer-wise relevance propagation was used to visualize individual nucleotide contributions to filter strength (52).

Ruling Out Sequencing Artifacts. Our model identified two uncharacterized exon-skipping features (stem loop structures and G-poor sequences). As our dataset was generated using multiple enzymatic reactions (from reverse transcription to polymerase amplification), we were concerned that rather than learning novel splicing features, our model instead learned artifacts introduced during data generation. Indeed, previous work noted that both structured and unstructured nucleic acid regions can impair enzymatic reactivity (53). If these features impaired enzymatic reactivity, we would expect to observe distinct biases in the absolute number of sequencing reads for such exons. Specifically, for exons containing stem loop structures or G-poor sequences, we would expect an undercounting of inclusion reads $n_{\text{inclusion}}$ (due to the inability to sequence reads containing the feature in the included exon) and, crucially, no change in the number of skipping reads n_{skipping} (since such reads do not contain the exon).

In contrast, a bona fide exon skipping feature should appear as a reduction in $n_{\text{inclusion}}$ accompanied by an increase in n_{skipping} . An analysis of our sequencing data supports that the stem loops and G-poor sequences are bona fide exon-skipping features rather than sequencing artifacts (SI Appendix, Fig. S9).

Design of Mutant Constructs. To validate the stem loop feature, candidate exons with high medium-length stem loop filter strength (top percentile) but with no other stem loop activations elsewhere in the exon were selected. Three mutants of each such exon were then generated. To ensure these mutants do not introduce or disrupt other features, exons where this mutation significantly changed strengths of other filters were filtered out.

To validate the G-poor stretch feature, candidate exons that strongly activate the G-poor filter exactly once along the exon were selected. For each candidate exon, a C-to-G mutation in the middle of the activated filter's window was introduced. As before, to ensure this does not introduce or disrupt other features, exons where this mutation significantly changed strengths of other filters were filtered out.

Data, Materials, and Software Availability. Sequence data that support the findings of this study have been deposited in the Gene Expression Omnibus under accession number GSE200096 (54). Custom code, preprocessed datasets, and trained model are available on GitHub (<https://github.com/regev-lab/interpretable-splicing-model>) (55). Plasmids used in this study will be available through AddGene.

ACKNOWLEDGMENTS. We thank members of the Regev laboratory and Lawrence Chasin for feedback; and Megan S. Hogan and Matthew T. Maurano [New York University (NYU) Institute for Systems Genetics] for technical assistance with amplicon sequencing. We thank Georg Seelig, Jef Boeke, and Brenton Graveley for plasmids used to construct the reporter assay. We acknowledge support by NYU Information Technology High Performance Computing resources, services, and staff expertise. We acknowledge funding from PhRMA Fellowship (M.S.); Lalor Foundation Fellowship (S.E.L.); Life Sciences Research Foundation Fellowship from Additional Ventures (S.E.L.); Simons Investigator Award (O.R.); NSF MCB-2226731 (O.R.).

Author affiliations: ^aDepartment of Computer Science, Courant Institute of Mathematical Sciences, New York University, New York, NY 10012

Author contributions: S.E.L. and O.R. designed research; S.E.L., M.S., and O.R. performed research; S.E.L., M.S., and O.R. analyzed data; and S.E.L., M.S., and O.R. wrote the paper.

1. R. Geirhos *et al.*, Shortcut learning in deep neural networks. *Nat. Mach. Intell.* **2**, 665–673 (2020).
2. T. Kashima, J. L. Manley, A negative element in SMN2 exon 7 inhibits splicing in spinal muscular atrophy. *Nat. Genet.* **34**, 460–463 (2003).
3. R. Cheung *et al.*, A multiplexed assay for exon recognition reveals that an unappreciated fraction of rare genetic variants cause large-effect splicing disruptions. *Mol. Cell* **73**, 183–194.e8 (2019).
4. N. Scalzitti *et al.*, Multi-species splice site prediction using convolutional neural networks. *BMC Bioinf.* **22**, 561 (2021).
5. K. Jaganathan *et al.*, Predicting splicing from primary sequence with deep learning. *Cell* **176**, 535–548.e24 (2019).
6. P. Julien, B. Miñana, P. Baeza-Centurion, J. Valcárcel, B. Lehner, The complete local genotype-phenotype landscape for the alternative splicing of a human exon. *Nat. Commun.* **7**, 11558 (2016).
7. S. Ke *et al.*, Saturation mutagenesis reveals manifold determinants of exon definition. *Genome Res.* **28**, 11–24 (2018).
8. A. B. Rosenberg, R. P. Patwardhan, J. Shendure, G. Seelig, Learning the sequence determinants of alternative splicing from millions of random sequences. *Cell* **163**, 698–711 (2015).
9. C. G. de Boer *et al.*, Deciphering eukaryotic gene-regulatory logic with 100 million random promoters. *Nat. Biotechnol.* **38**, 56–65 (2020).
10. K. J. Hertel, Combinatorial control of exon recognition. *J. Biol. Chem.* **283**, 1211–1215 (2008).
11. K. Cho, B. van Merriënboer, D. Bahdanau, Y. Bengio, "On the properties of neural machine translation: Encoder-decoder approaches" in *Syntax, Semantics and Structure in Statistical Translation* (2014), p. 103.
12. A. Vaswani *et al.*, "Attention is all you need" in *Advances in Neural Information Processing Systems*, U. von Luxburg, I. Guyon, S. Bengio, H. Wallach, R. Fergus, Eds. (MIT Press, 2017), vol. 30.
13. G. Novakovsky, N. Dexter, M. W. Libbrecht, W. W. Wasserman, S. Mostafavi, Obtaining genetics insights from deep learning via explainable artificial intelligence. *Nat. Rev. Genet.* **24**, 125–137 (2023).
14. N. N. Singh, E. J. Androphy, R. N. Singh, An extended inhibitory context causes skipping of exon 7 of SMN2 in spinal muscular atrophy. *Biochem. Biophys. Res. Commun.* **315**, 381–388 (2004).
15. N. N. Singh, E. J. Androphy, R. N. Singh, In vivo selection reveals combinatorial controls that define a critical exon in the spinal muscular atrophy genes. *RNA (New York, N.Y.)* **10**, 1291–1305 (2004).
16. N. N. Singh, R. N. Singh, E. J. Androphy, Modulating role of RNA structure in alternative splicing of a critical exon in the spinal muscular atrophy genes. *Nucleic Acids Res.* **35**, 371–389 (2007).
17. P. Baeza-Centurion, B. Miñana, J. M. Schmiedel, J. Valcárcel, B. Lehner, Combinatorial genetics reveals a scaling law for the effects of mutations on splicing. *Cell* **176**, 549–563.e23 (2019).
18. W. G. Fairbrother, R.-F. Yeh, P. A. Sharp, C. B. Burge, Predictive identification of exonic splicing enhancers in human genes. *Science (New York, N.Y.)* **297**, 1007–1013 (2002).
19. S. Ke *et al.*, Quantitative evaluation of all hexamers as exonic splicing elements. *Genome Res.* **21**, 1360–1374 (2011).
20. B. R. Graveley, Mutually exclusive splicing of the insect Dscam pre-mRNA directed by competing intronic RNA secondary structures. *Cell* **123**, 65–73 (2005).
21. R. Lorenz *et al.*, ViennaRNA package 2.0. *Algorithms Mol. Biol.* **6** (2011).
22. J. Zhu, A. Mayeda, A. R. Krainer, Exon identity established through differential antagonism between exonic splicing silencer-bound hnRNP A1 and enhancer-bound SR proteins. *Mol. Cell* **8**, 1351–1361 (2001).
23. D. Ray *et al.*, A compendium of RNA-binding motifs for decoding gene regulation. *Nature* **499**, 172–177 (2013).
24. D. Domínguez *et al.*, Sequence, structure, and context preferences of human RNA binding proteins. *Mol. Cell* **70**, 854–867.e9 (2018).
25. T. D. Schneider, R. M. Stephens, Sequence logos: A new way to display consensus sequences. *Nucleic Acids Res.* **18**, 6097–6100 (1990).
26. Y. Cavaloc, C. F. Bourgeois, L. Kister, J. Stévenin, The splicing factors 9G8 and SRp20 transactivate splicing through different and specific enhancers. *RNA (New York, N.Y.)* **5**, 468–483 (1999).
27. T. D. Schaal, T. Maniatis, Selection and characterization of pre-mRNA splicing enhancers: Identification of novel SR protein-specific enhancer sequences. *Mol. Cell. Biol.* **19**, 1705–1719 (1999).

28. C. D. Chen, R. Kobayashi, D. M. Helfman, Binding of hnRNP H to an exonic splicing silencer is involved in the regulation of alternative splicing of the rat beta-tropomyosin gene. *Genes Dev.* **13**, 593–606 (1999).
29. X.H.-F. Zhang, M. A. Arias, S. Ke, L. A. Chasin, Splicing of designer exons reveals unexpected complexity in pre-mRNA splicing. *RNA (New York, N.Y.)* **15**, 367–376 (2009).
30. Z. Wang *et al.*, Systematic identification and analysis of exonic splicing silencers. *Cell* **119**, 831–845 (2004).
31. C. L. Williamson, N. M. Desai, J. M. Burke, Compensatory mutations demonstrate that P8 and P6 are RNA secondary structure elements important for processing of a group I intron. *Nucleic Acids Res.* **17**, 675–689 (1989).
32. K. K. Dey, D. Xie, M. Stephens, A new sequence logo plot to highlight enrichment and depletion. *BMC Bioinf.* **19**, 473 (2018).
33. I. Pitout, L. L. Flynn, S. D. Wilton, S. Fletcher, Antisense-mediated splice intervention to treat human disease: The odyssey continues. *F1000Research* **8**, F1000 Faculty Rev-710 (2019).
34. S. E. Liao, O. Regev, Splicing at the phase-separated nuclear speckle interface: A model. *Nucleic Acids Res.* **49**, 636–645 (2021).
35. G. Yeo, D. Holste, G. Kreiman, C. B. Burge, Variation in alternative splicing across human tissues. *Genome Biol.* **5**, R74 (2004).
36. E. T. Wang *et al.*, Alternative isoform regulation in human tissue transcriptomes. *Nature* **456**, 470–476 (2008).
37. P. V. Mazin, P. Khaitovich, M. Cardoso-Moreira, H. Kaessmann, Alternative splicing during mammalian organ development. *Nat. Genet.* **53**, 925–934 (2021).
38. J. A. Calarco *et al.*, Regulation of vertebrate nervous system alternative splicing and development by an SR-related protein. *Cell* **138**, 898–910 (2009).
39. L. Jia *et al.*, Decoding mRNA translatability and stability from the 5' UTR. *Nat. Struct. Mol. Biol.* **27**, 814–821 (2020).
40. M. Rabani, L. Pieper, G.-L. Chew, A. F. Schier, A massively parallel reporter assay of 3' UTR sequences identifies in vivo rules for mRNA degradation. *Mol. Cell* **68**, 1083–1094.e5 (2017).
41. Z. Luo, J. Zhang, J. Fei, S. Ke, Deep learning modeling m⁶A deposition reveals the importance of downstream cis-element sequences. *Nat. Commun.* **13**, 2720 (2022).
42. Y.-Y. Lee, H. Kim, V. Narry, Kim., Sequence determinant of small RNA production by DICER. *Nature* **615**, 323–330 (2023).
43. S. A. Smith, K. W. Lynch, Cell-based splicing of minigenes. *Methods Mol. Biol. (Clifton, N.J.)* **1126**, 243–255 (2014).
44. S. I. Adamson, L. Zhan, B. R. Graveley, Vex-seq: High-throughput identification of the impact of genetic variation on pre-mRNA splicing efficiency. *Genome Biol.* **19**, 71 (2018).
45. G. Yeo, C. B. Burge, Maximum entropy modeling of short sequence motifs with applications to RNA splicing signals. *J. Comput. Biol.* **11**, 377–394 (2004).
46. J. König *et al.*, iCLIP reveals the function of hnmp particles in splicing at individual nucleotide resolution. *Nat. Struct. Mol. Biol.* **17**, 909–915 (2010).
47. T. Kivioja *et al.*, Counting absolute numbers of molecules using unique molecular identifiers. *Nat. Methods* **9**, 72–74 (2012).
48. G. Van Rossum, F. L. Drake, *Python 3 Reference Manual* (CreateSpace, Scotts Valley, CA, 2009).
49. M. Abadi *et al.*, TensorFlow: Large-scale machine learning on heterogeneous systems (2015). tensorflow.org. 9 December 2022.
50. C. R. Harris *et al.*, Array programming with NumPy. *Nature* **585**, 357–362 (2020).
51. Fundamental Algorithms for Scientific Computing in Python, Pauli Virtanen, et al. *SciPy 1.0. Nat. Methods* **17**, 261–272 (2020).
52. S. Bach *et al.*, On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS ONE* **10**, e0130140 (2015).
53. T. M. Snyder, B. N. Tse, D. R. Liu, Effects of template sequence and secondary structure on DNA-templated reactivity. *J. Am. Chem. Soc.* **130**, 1392–1401 (2008).
54. S. E. Liao, M. Sudarshan, O. Regev, Machine learning for discovery: Deciphering RNA splicing logic. Gene Expression Omnibus. <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE200096>. Deposited 4 April 2022.
55. M. Sudarshan, Code for Liao, Sudarshan, and Regev: "Deciphering RNA splicing logic with interpretable machine learning." GitHub. <https://github.com/regev-lab/interpretable-splicing-model>. Accessed 19 September 2023.