

Machine learning for discovery: deciphering RNA splicing logic

SUSAN E. LIAO^{1,2}, MUKUND SUDARSHAN^{1,2} & ODED REGEV^{1*}

Summary

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 novel components of the splicing logic, which we experimentally validated. This study highlights how interpretable machine learning can advance scientific discovery.

Introduction

Machine learning algorithms, in particular neural networks, capture complex quantitative relationships between input and output. However, as neural networks are typically black box, it is difficult to extract post-hoc insights on 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 new insights on the underlying processes¹.

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 novel splicing features.

Generating a synthetic dataset for interpretable machine learning

As neural network performance and interpretability is inextricable from the data it is trained on, we began by generating a large, high-quality synthetic splicing dataset. The use of synthetic datasets

¹Department of Computer Science, Courant Institute of Mathematical Sciences, New York University, NY, USA.

²These authors contributed equally: Susan E. Liao, Mukund Sudarshan.

*e-mail: regev@cims.nyu.edu

40 offers several advantages over genomic data used in previous work. First, genomic datasets are
 41 limited by the number of exons in the genome. In contrast, synthetic assays can dramatically
 42 increase the number of data points by orders of magnitude^{8,9}. Second, genomic exons are flanked
 43 by varying sequences (splice sites, introns, promoters) that also participate in splicing decisions¹⁰,
 44 greatly complicating attempts at interpretability. In contrast, synthetic datasets fix all but one variable
 45 region, allowing to focus on the region of interest. Third, genomic exons contain overlapping RNA
 46 codes (e.g., protein coding sequences). In contrast, sequences in synthetic datasets are devoid of
 47 overlapping codes by design. In summary, from both a quantity and quality perspective, synthetic
 48 datasets provide crucial advantages for machine learning over genomic datasets.

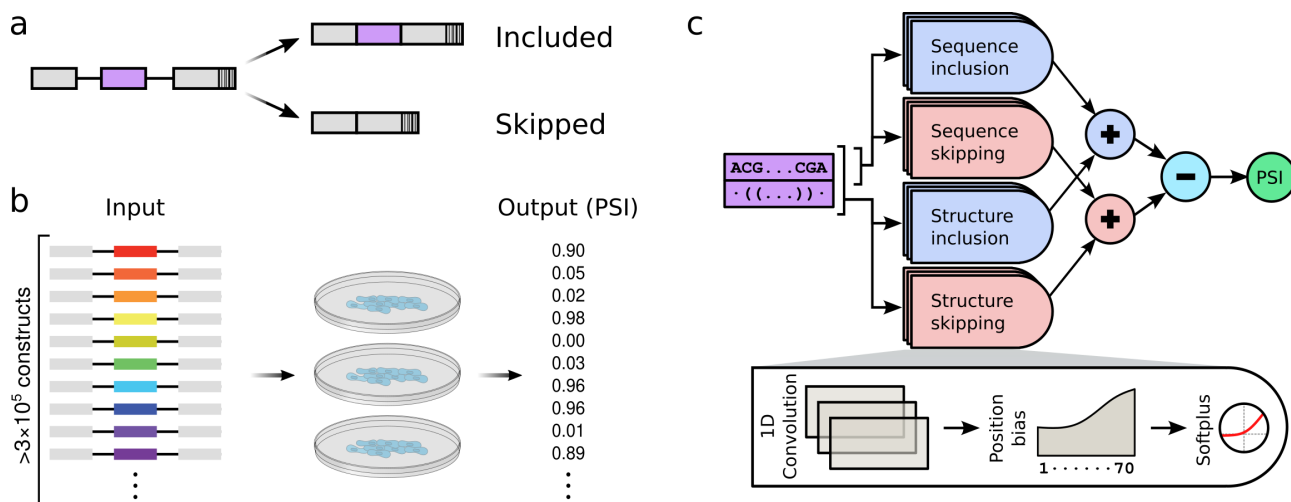


Figure 1 | Data generation and interpretable-by-design machine learning model. **a**, All reporters in the assay share the same three-exon design, and differ only in their middle exon, which contains a random 70 nucleotide-long sequence. Depending on its sequence, an exon might be included, skipped, 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. The difference between the total strength of the inclusion filters and the total strength of the skipping filters is used to compute the output predicted PSI.

49 The synthetic dataset we generated includes hundreds of thousands of input-output data points.
 50 Each data point is a different random 70-nucleotide exon sequence, paired with a measured percent
 51 spliced in (PSI) output, which is a number between 0 (always skipped) and 1 (always included)
 52 (Fig. 1a). The dataset is generated by a massively parallel reporter assay that allows for PSI quanti-
 53 fication for hundreds of thousands of unique sequences in a single biological experiment (Fig. 1b).
 54 Splicing outcomes for the parallel reporter assay were measured after transfection into human HeLa
 55 cells using high-throughput sequencing. We verified that reporters are evenly represented in the
 56 reporter assay (Extended Data Fig. 1a). The vast majority of splicing products corresponded to
 57 exon inclusion or exon skipping products (Extended Data Fig. 1b), and we filtered our data to
 58 exclude spurious splicing products. PSI values are calculated as the number of inclusion reads
 59 divided by the total number of inclusion and skipping reads. Three biological replicates of the assay
 60 showed excellent agreement (Extended Data Fig. 1c) and their sequencing results were combined
 61 for all downstream analysis. High-throughput sequencing measurements were consistent with
 62 semi-quantitative measurements of individual reporters (Extended Data Fig. 1d).

63 **An interpretable-by-design model accurately predicts splicing outcomes**

64 State-of-the-art neural networks (based on gated recurrent units¹¹ and transformers¹²) trained on this
65 dataset provided excellent prediction accuracy on a held-out test set (RMSE=0.165 and RMSE=0.183,
66 respectively). However, these models are not interpretable, and do not provide any biological insights.
67 We therefore designed a novel model with the explicit goal of being interpretable.

68 The predictive accuracy of our interpretable-by-design model is comparable to that of state-of-the-
69 art models trained on the same synthetic dataset (RMSE=0.180; Extended Data Fig. 2a). This suggests
70 that interpretability need not come at the expense of accuracy. In addition to our own dataset, the
71 model accurately predicts splicing outcomes from other splicing datasets^{7,8,13-16} (Extended Data
72 Fig. 2b). Importantly, unlike our random exons, these datasets were modeled on specific genomic
73 exons, with each dataset differing in splice sites, introns, and flanking exons. Furthermore, these
74 datasets were generated in different immortalized cell lines. Encouragingly, despite these dramatic
75 differences in RNA architecture and cell types, our model tested well on these datasets, suggesting
76 that our model generalizes and captures critical aspects of splicing regulatory logic.

77 **Model architecture reveals unifying decision-making process**

78 Our interpretable-by-design model incorporates domain knowledge throughout its architecture
79 (Fig. 1c). Specifically, we reasoned that short six nucleotide sequence filters would capture motifs
80 previously demonstrated to play an important role in splicing decisions^{17,18}. We therefore introduced
81 one-dimensional convolutional filters applied to the input RNA sequence. Next, since RNA secondary
82 structure was previously implicated in splicing outcomes^{15,19}, we also provided the network with
83 predicted structure²⁰. We then introduced longer (30 nucleotide) one-dimensional convolutional
84 filters to the structure-augmented sequence. Crucially, while we fixed filter lengths using minimal
85 domain knowledge, we did not explicitly specify sequences and structures, allowing the network
86 flexibility to learn filters in an unbiased manner. Furthermore, our model explicitly quantifies the
87 strength (in network-defined arbitrary units) of each activated filter to the inclusion or skipping
88 decision. Importantly, we allowed the strength of any filter to vary along the length of an exon,
89 providing the network the flexibility to capture position-dependent effects of RNA features on
90 splicing outcomes.

91 To arrive at its output, the network computes the difference in the sum total of exon inclusion
92 strengths and exon skipping strengths (Δ strength), which is then converted to predicted PSI.
93 The greater the magnitude of this difference, the closer the PSI is to 0 (difference $\ll 0$) or 1
94 (difference $\gg 0$). This additive combinatorial behavior is consistent with previous literature^{8,21}.

95 **Model extends understanding of splicing regulatory logic**

96 Even though our model was trained on a synthetic dataset, it recapitulates and extends domain
97 knowledge from previous genomic and biochemical studies.

98 Many filters in the model match binding motifs of RNA binding proteins implicated in splicing
99 regulation (splicing factors)^{24,25} (Fig. 2a). Consistent with previous studies, network inclusion filters
100 match binding sites for SR proteins known to promote exon inclusion^{23,26}, whereas network skipping
101 filters match binding sites for hnRNP proteins known to promote exon skipping²⁷.

102 However, while the directionality of these RNA features towards splicing was established, their
103 magnitude was not clear. Importantly, the model addresses this issue by assigning a quantitative
104 strength to each filter. Moreover, some filters exhibit striking position dependent strengths, suggesting
105 that the position of an RNA feature along an exon affects its strength.

106 Surprisingly, our network accurately predicted splicing outcomes using a concise list of filters
107 (Fig. 2a). This contrasts with previous studies suggesting that splicing outcomes result from the

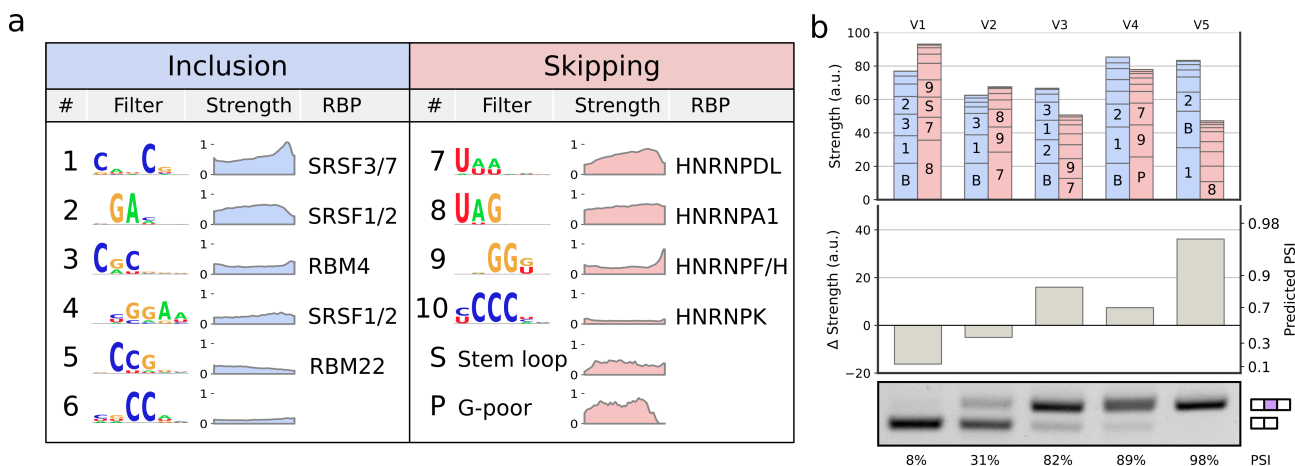


Figure 2 | Model expands on known splicing logic and its predictions can be interpreted using balance plots. a, Splicing features detected by the model’s filters, represented by their sequence logo²². Filters either contribute to inclusion (blue) or skipping (red). Plots show the average strength in our dataset of each filter as a function of position along the exon. RNA binding proteins (RBP) with a similar binding motif, as reported in previous work^{23–25}. The model also identified short stem loops and long G-poor stretches as contributing to exon skipping. **b,** 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 panel a. Bar labeled B represents a constant initial inclusion strength. Labels are not shown for smaller bars. The difference between total inclusion and total skipping strengths (Δ strength) leads to predicted PSI (center). PSI as measured by semi-quantitative RT-PCR matches the machine learning predictions (bottom).

108 combinatorics of hundreds of unique RNA features^{8,28,29}.

109 Using the local interpretability of our model, we introduce a visualization (balance plot) that
 110 enables explicit examination and quantification of how multiple RNA features lead to splicing
 111 outcomes for any given exon from our dataset (Fig. 2b, Extended Data Fig. 3). For a given exon, the
 112 total strengths of activated filters are represented as bars of the appropriate height. Total inclusion
 113 strength (blue) and skipping strength (red) are then visible as the height of the stacked bars. The
 114 Δ strength is represented by the difference in heights between the stacked inclusion and skipping
 115 filters. These visualizations provide an intuitive tool to understand the contributions of individual
 116 sequence and structure features leading to each exon’s predicted PSI. They emphasize that splicing
 117 logic results from contributions of many RNA features along the exon, and that a single nucleotide
 118 can be part of multiple overlapping filters^{6,8}.

119 Discovery and validation of novel splicing features

120 Next, we asked whether our interpretable-by-design model could advance scientific discovery by
 121 identifying novel splicing features. While most network filters were consistent with previously-
 122 described splicing features, two uncharacterized long skipping filters with strong influence on
 123 splicing predictions stood out (Fig. 2a). We confirmed that these filters were robustly identified
 124 across multiple initialization seeds and training/testing splits, suggesting that they are not training
 125 artifacts. We then turned our attention to characterizing and experimentally validating these features.

126 Examining the first uncharacterized filter revealed that it identifies stem loop structures with short,
 127 GC-rich, 5-7 nucleotide double-stranded regions (Fig. 3). Next, we experimentally validated that these
 128 stem loops contribute to exon skipping and are not artifacts. We introduced mutations that disrupt
 129 double-stranded base pairing in an exon with such a stem loop. First, we introduced single nucleotide
 130 mutations predicted to abolish the stem by disrupting base pairing. Notably, these mutations were
 131 designed to minimize disruptions of other filters, ensuring that prediction differences are mainly

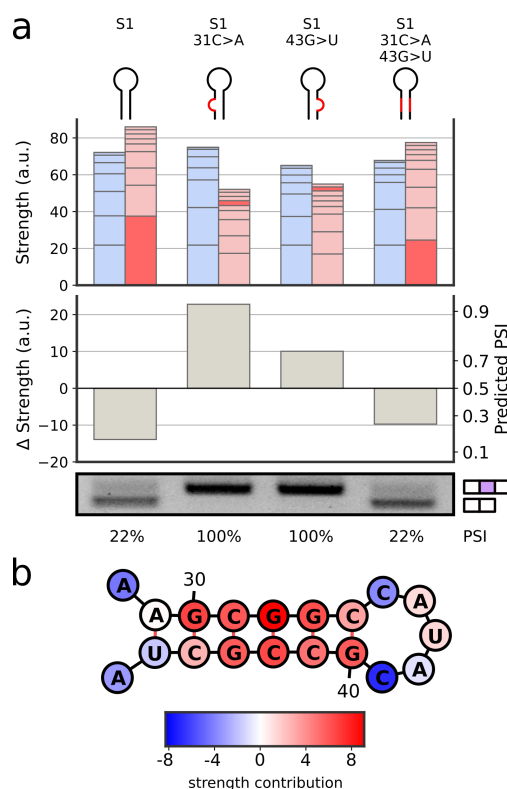


Figure 3 | Validation of novel stem loop feature. **a**, The machine learning model identifies a stem loop in an exon (S1) as having a strong skipping strength (dark red bar; top), leading to near complete skipping prediction (middle). Single nucleotide mutations disrupting a downstream or upstream stem base pair are predicted to significantly reduce the skipping strength and restore exon inclusion. Finally, including both single nucleotide mutations is predicted to restore the stem loop skipping strength and lead to skipping. RT-PCR validation (bottom) confirms the machine learning predictions. **b**, The stem loop identified in S1, with the individual contributions to its strength by each nucleotide.

132 due to altered secondary structure, and not due to the introduction or disruption of other sequence
 133 features. In addition to two such mutations, we also introduced both compensatory mutations
 134 together, restoring the original stem loop structure³⁰. We measured splicing outcomes for all four
 135 individual reporters (original, upstream mutation, downstream mutation, and double mutations) and
 136 observed that splicing outcomes matched our predictions (Fig. 3). Namely, PSI increased dramatically
 137 in both single nucleotide mutants, in agreement with the predicted decrease in filter strength.
 138 Furthermore, when both compensatory mutations are present and structure is restored, measured
 139 PSI was comparable to that of the original exon. We applied the same experimental validation
 140 scheme to two other stem loop-containing exons. In both cases, stem-disrupting single mutations
 141 increased exon inclusion and structure-restoring double compensatory mutations had minimal effects
 142 (Extended Data Fig. 4). Together, these experiments demonstrate that model-identified stem loops,
 143 rather than sequence, contribute to exon skipping.

144 In contrast, examining the second uncharacterized filter did not reveal any secondary structure
 145 preference. Instead, the filter exhibited a preference for long guanine depleted (G-poor) sequences
 146 (Fig. 4a). To validate that guanine depletion underlies filter behavior, we selected an exon with a
 147 G-poor sequence and introduced a single C>G mutation. As before, we ensured that the predicted
 148 strengths of other filters are only minimally disrupted (Fig. 4b). Strikingly, this single mutation led
 149 to marked increase in PSI. We applied the same validation scheme to three other exons with G-poor
 150 sequences; in every instance, a single C>G mutation increased exon inclusion (Extended Data Fig. 5).

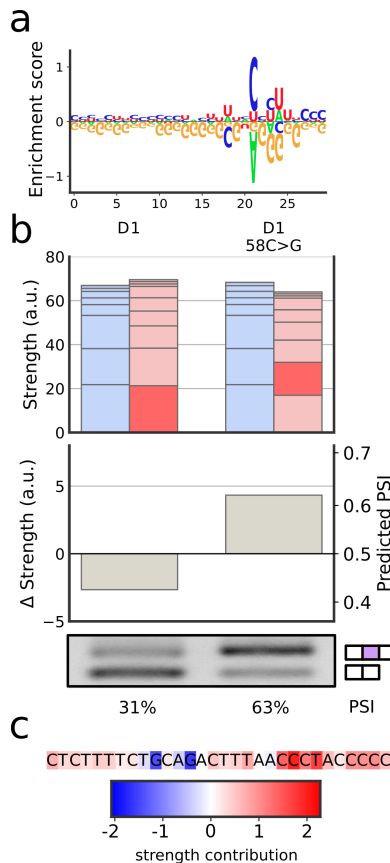


Figure 4 | Validation of the G-poor feature. **a**, The G-poor filter, represented by its enrichment-depletion \log^{31} . **b**, The machine learning model identifies a G-poor stretch in an exon (D1) as having a strong skipping strength (dark red bar, top), leading to skipping prediction (middle). A single nucleotide C>G mutation is predicted to disrupt the G-poor stretch and restore exon inclusion (right bars). RT-PCR validation (bottom) confirms the machine learning predictions. **c**, The G-poor stretch identified in D1, with the individual contributions to its strength by each nucleotide.

151 To the best of our knowledge, a long G-poor sequence has not been described in the literature.

152 Collectively, these experiments confirm that stem loops and G-poor sequences identified by the
 153 model reflect bona fide splicing features.

154 Discussion

155 In this study, we demonstrate that an interpretable-by-design model advanced scientific discovery.
 156 Our model accurately predicts splicing outcomes on both our assay and on previously published
 157 assays, demonstrating that interpretability need not come at the expense of accuracy or generalizabil-
 158 ity. Model interpretability enabled a systematic understanding of RNA splicing logic, including the
 159 identification of two candidate novel exon skipping features which were subsequently experimentally
 160 validated. The model's ability to quantify contributions of specific features to splicing outcomes for
 161 individual exons has considerable potential for a range of medical and biotechnology applications,
 162 including genome- or RNA-editing of target exons to correct splicing behavior or guiding rational
 163 design of RNA-based therapeutics like antisense oligonucleotides³².

164 In addition, model-identified features hint at novel biochemical mechanisms that warrant further
 165 study. For example, the fact that splicing decisions are modeled well by an additive quantity (Δ
 166 strength) supports a biochemical mechanism involving the nuclear spatial organization of SR and
 167 hnRNP proteins³³. Furthermore, the novel skipping-promoting G-poor feature may point to an

168 uncharacterized RNA binding protein or complex. These open questions further underscore how
169 interpretable-by-design models can advance scientific discovery by aiding hypothesis generation.

170 Our model performs well on synthetic datasets from immortalized cell lines, yet further work
171 is needed to capture the dynamics of developmentally regulated splicing logic^{34–36}. Importantly,
172 splicing outcomes change depending on the expression level of cell type-specific RNA binding
173 proteins³⁷. These questions can be addressed by generation of additional synthetic splicing datasets
174 in developmentally relevant cell types paired with interpretable-by-design models that capture cell
175 type-specific regulatory features.

176 Beyond the context of splicing, the interpretable-by-design framework can be used to decipher the
177 multiple, complex, and overlapping codes that dictate biomolecular processing. Importantly, many
178 rich synthetic datasets that address RNA untranslated 5'³⁸ and 3'³⁹ region regulation, methylation⁴⁰,
179 and small RNA biogenesis⁴¹, have already been generated. We expect that additional data generation
180 efforts paired with the interpretable-by-design framework will stimulate advances in understanding
181 biological codes more broadly.

182 **Data availability** Sequence data that support the findings of this study have been deposited in the
183 NCBI's Gene Expression Omnibus under accession number GSE200096.

184 **Code availability** Custom code, preprocessed datasets, and trained model are available on GitHub
185 (<https://github.com/regev-lab/interpretable-splicing-model>).

186 **Acknowledgements** We thank members of the Regev laboratory and Lawrence Chasin for feedback;
187 and Megan S. Hogan and Matthew T. Maurano (NYU Institute for Systems Genetics) for technical
188 assistance with amplicon sequencing. We thank Georg Seelig, Jef Boeke, and Brenton Graveley
189 for plasmids used to construct the reporter assay. This work was partly supported by a PhRMA
190 Fellowship (M.S.), Lalor Foundation Fellowship (S.E.L.), Life Sciences Research Foundation Fellowship
191 from Additional Ventures (S.E.L.), Simons Investigator Award, and NSF MCB-2226731 (O.R.), and
192 NYU IT High Performance Computing resources, services, and staff expertise.

193 **Author contributions** The study was initiated by S.E.L. and O.R. Experiments and machine learning
194 analysis were performed by S.E.L. and M.S. All authors wrote, reviewed, and provided feedback on
195 the manuscript.

196 **Competing interests** The authors declare no competing interests.

197 REFERENCES

- 198 1. Geirhos, R. *et al.* Shortcut learning in deep neural networks. *Nature Machine Intelligence* **2**, 665–673
199 (2020).
- 200 2. Kashima, T. & Manley, J. L. A negative element in SMN2 exon 7 inhibits splicing in spinal
201 muscular atrophy. *Nature genetics* **34**, 460–3 (2003).
- 202 3. Cheung, R. *et al.* A Multiplexed Assay for Exon Recognition Reveals that an Unappreciated
203 Fraction of Rare Genetic Variants Cause Large-Effect Splicing Disruptions. *Molecular cell* **73**,
204 183–194.e8 (2019).
- 205 4. Scalzitti, N. *et al.* Spliceator: multi-species splice site prediction using convolutional neural
206 networks. *BMC bioinformatics* **22**, 561 (2021).

- 207 5. Jaganathan, K. *et al.* Predicting Splicing from Primary Sequence with Deep Learning. *Cell* **176**,
208 535–548.e24 (2019).
- 209 6. Julien, P., Miñana, B., Baeza-Centurion, P., Valcárcel, J. & Lehner, B. The complete local genotype-
210 phenotype landscape for the alternative splicing of a human exon. *Nature communications* **7**,
211 11558 (2016).
- 212 7. Ke, S. *et al.* Saturation mutagenesis reveals manifold determinants of exon definition. *Genome*
213 *research* **28**, 11–24 (2018).
- 214 8. Rosenberg, A. B., Patwardhan, R. P., Shendure, J. & Seelig, G. Learning the sequence determi-
215 nants of alternative splicing from millions of random sequences. *Cell* **163**, 698–711 (2015).
- 216 9. De Boer, C. G. *et al.* Deciphering eukaryotic gene-regulatory logic with 100 million random
217 promoters. *Nature biotechnology* **38**, 56–65 (2020).
- 218 10. Hertel, K. J. Combinatorial control of exon recognition. *The Journal of biological chemistry* **283**,
219 1211–5 (2008).
- 220 11. Cho, K., van Merriënboer, B., Bahdanau, D. & Bengio, Y. On the Properties of Neural Ma-
221 chine Translation: Encoder–Decoder Approaches. *Syntax, Semantics and Structure in Statistical*
222 *Translation*, 103 (2014).
- 223 12. Vaswani, A. *et al.* Attention is all you need. *Advances in neural information processing systems* **30**
224 (2017).
- 225 13. Singh, N. N., Androphy, E. J. & Singh, R. N. An extended inhibitory context causes skipping of
226 exon 7 of SMN2 in spinal muscular atrophy. *Biochemical and biophysical research communications*
227 **315**, 381–8 (2004).
- 228 14. Singh, N. N., Androphy, E. J. & Singh, R. N. In vivo selection reveals combinatorial controls that
229 define a critical exon in the spinal muscular atrophy genes. *RNA (New York, N.Y.)* **10**, 1291–305
230 (2004).
- 231 15. Singh, N. N., Singh, R. N. & Androphy, E. J. Modulating role of RNA structure in alternative
232 splicing of a critical exon in the spinal muscular atrophy genes. *Nucleic acids research* **35**, 371–89
233 (2007).
- 234 16. Baeza-Centurion, P., Miñana, B., Schmiedel, J. M., Valcárcel, J. & Lehner, B. Combinatorial
235 Genetics Reveals a Scaling Law for the Effects of Mutations on Splicing. *Cell* **176**, 549–563.e23
236 (2019).
- 237 17. Fairbrother, W. G., Yeh, R.-F., Sharp, P. A. & Burge, C. B. Predictive identification of exonic
238 splicing enhancers in human genes. *Science (New York, N.Y.)* **297**, 1007–13 (2002).
- 239 18. Ke, S. *et al.* Quantitative evaluation of all hexamers as exonic splicing elements. *Genome research*
240 **21**, 1360–74 (2011).
- 241 19. Graveley, B. R. Mutually exclusive splicing of the insect Dscam pre-mRNA directed by competing
242 intronic RNA secondary structures. *Cell* **123**, 65–73 (2005).
- 243 20. Lorenz, R. *et al.* ViennaRNA Package 2.0. *Algorithms for Molecular Biology* **6** (Nov. 2011).
- 244 21. Zhu, J., Mayeda, A. & Krainer, A. R. Exon identity established through differential antagonism
245 between exonic splicing silencer-bound hnRNP A1 and enhancer-bound SR proteins. *Molecular*
246 *cell* **8**, 1351–61 (2001).
- 247 22. Schneider, T. D. & Stephens, R. M. Sequence logos: a new way to display consensus sequences.
248 *Nucleic acids research* **18**, 6097–100 (1990).

-
- 249 23. Cavaloc, Y., Bourgeois, C. F., Kister, L. & Stévenin, J. The splicing factors 9G8 and SRp20
250 transactivate splicing through different and specific enhancers. *RNA (New York, N.Y.)* **5**, 468–83
251 (1999).
- 252 24. Ray, D. *et al.* A compendium of RNA-binding motifs for decoding gene regulation. *Nature* **499**,
253 172–7 (2013).
- 254 25. Dominguez, D. *et al.* Sequence, Structure, and Context Preferences of Human RNA Binding
255 Proteins. *Molecular cell* **70**, 854–867.e9 (2018).
- 256 26. Schaal, T. D. & Maniatis, T. Selection and characterization of pre-mRNA splicing enhancers:
257 identification of novel SR protein-specific enhancer sequences. *Molecular and cellular biology* **19**,
258 1705–19 (1999).
- 259 27. Chen, C. D., Kobayashi, R. & Helfman, D. M. Binding of hnRNP H to an exonic splicing silencer
260 is involved in the regulation of alternative splicing of the rat beta-tropomyosin gene. *Genes &
261 development* **13**, 593–606 (1999).
- 262 28. Zhang, X. H.-F., Arias, M. A., Ke, S. & Chasin, L. A. Splicing of designer exons reveals unexpected
263 complexity in pre-mRNA splicing. *RNA (New York, N.Y.)* **15**, 367–76 (2009).
- 264 29. Wang, Z. *et al.* Systematic identification and analysis of exonic splicing silencers. *Cell* **119**, 831–45
265 (2004).
- 266 30. Williamson, C. L., Desai, N. M. & Burke, J. M. Compensatory mutations demonstrate that P8
267 and P6 are RNA secondary structure elements important for processing of a group I intron.
268 *Nucleic acids research* **17**, 675–689 (1989).
- 269 31. Dey, K. K., Xie, D. & Stephens, M. A new sequence logo plot to highlight enrichment and
270 depletion. *BMC bioinformatics* **19**, 473 (2018).
- 271 32. Pitout, I., Flynn, L. L., Wilton, S. D. & Fletcher, S. Antisense-mediated splice intervention to
272 treat human disease: the odyssey continues. *F1000Research* **8**, F1000 Faculty Rev–710 (2019).
- 273 33. Liao, S. E. & Regev, O. Splicing at the phase-separated nuclear speckle interface: a model. *Nucleic
274 acids research* **49**, 636–645 (2021).
- 275 34. Yeo, G., Holste, D., Kreiman, G. & Burge, C. B. Variation in alternative splicing across human
276 tissues. *Genome biology* **5**, R74 (2004).
- 277 35. Wang, E. T. *et al.* Alternative isoform regulation in human tissue transcriptomes. *Nature* **456**,
278 470–6 (2008).
- 279 36. Mazin, P. V., Khaitovich, P., Cardoso-Moreira, M. & Kaessmann, H. Alternative splicing during
280 mammalian organ development. *Nature genetics* **53**, 925–934 (2021).
- 281 37. Calarco, J. A. *et al.* Regulation of vertebrate nervous system alternative splicing and development
282 by an SR-related protein. *Cell* **138**, 898–910 (2009).
- 283 38. Jia, L. *et al.* Decoding mRNA translatability and stability from the 5' UTR. *Nature structural &
284 molecular biology* **27**, 814–821 (2020).
- 285 39. Rabani, M., Pieper, L., Chew, G.-L. & Schier, A. F. A Massively Parallel Reporter Assay of 3'
286 UTR Sequences Identifies In Vivo Rules for mRNA Degradation. *Molecular cell* **68**, 1083–1094.e5
287 (2017).
- 288 40. Luo, Z., Zhang, J., Fei, J. & Ke, S. Deep learning modeling m⁶A deposition reveals the importance
289 of downstream cis-element sequences. *Nature communications* **13**, 2720 (2022).

-
- 290 41. Fukunaga, R. Loquacious-PD removes phosphate inhibition of Dicer-2 processing of hairpin
291 RNAs into siRNAs. *Biochemical and biophysical research communications* **498**, 1022–1027 (2018).
- 292 42. Smith, S. A. & Lynch, K. W. Cell-based splicing of minigenes. *Methods in molecular biology (Clifton,*
293 *N.J.)* **1126**, 243–55 (2014).
- 294 43. Adamson, S. I., Zhan, L. & Graveley, B. R. Vex-seq: high-throughput identification of the impact
295 of genetic variation on pre-mRNA splicing efficiency. *Genome biology* **19**, 71 (2018).
- 296 44. Yeo, G. & Burge, C. B. Maximum entropy modeling of short sequence motifs with applications
297 to RNA splicing signals. *Journal of computational biology* **11**, 377–394 (2004).
- 298 45. König, J. *et al.* iCLIP reveals the function of hnRNP particles in splicing at individual nucleotide
299 resolution. *Nature structural & molecular biology* **17**, 909–915 (2010).
- 300 46. Kivioja, T. *et al.* Counting absolute numbers of molecules using unique molecular identifiers.
301 *Nature methods* **9**, 72–74 (2012).
- 302 47. Van Rossum, G. & Drake, F. L. *Python 3 Reference Manual* (CreateSpace, Scotts Valley, CA, 2009).
- 303 48. Martin Abadi *et al.* *TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems* Software
304 available from tensorflow.org. 2015.
- 305 49. Harris, C. R. *et al.* Array programming with NumPy. *Nature* **585**, 357–362 (Sept. 2020).
- 306 50. Virtanen, P. *et al.* SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature*
307 *Methods* **17**, 261–272 (2020).
- 308 51. Bach, S. *et al.* On Pixel-Wise Explanations for Non-Linear Classifier Decisions by Layer-Wise
309 Relevance Propagation. *PloS one* **10**, e0130140 (2015).

I. METHODS

I.1 Reporter assay design and cloning

The splicing reporter is based on a three-exon beta globin minigene⁴² under the control of a truncated mammalian CAG promoter. The massively parallel splicing assay allows for high-throughput characterization of exon variants on splicing outcomes⁴³ using Gibson assembly and ligation cloning. The assay replaces the middle beta globin exon with 70nt random sequences flanked by weak splice sites (MaxEnt scores⁴⁴: 3'ss 9.41, 5'ss 5.06). Each 70nt exon is coupled with a 20nt barcode downstream of the third exon, allowing for identification of middle exon identity in exon skipping products. Briefly, two pools of oligonucleotides (random exons and barcodes with complementary linker regions and flanking overlap sequences for Gibson assembly) were synthesized as single-stranded oligonucleotides (IDT DNA Technologies) and joined using an anneal-extend procedure. 200nM of each oligonucleotide were joined in a 100µL reaction (Phusion[®] Hot Start Flex 2X Master Mix, NEB). Oligonucleotides were denatured at 98°C for 10min, cooled slowly to 60°C (0.1°C/sec), annealed at 60°C for 5min, and extended at 72°C for 60min. Single-stranded products were removed from pooled double-stranded exon-barcode using 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 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).

I.2 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, 200nM of each oligonucleotide were joined in a 100µL reaction with 5U DNA polymerase (NEB Klenow). Oligonucleotides were denatured at 98°C for 10min, annealed after cooling slowly to 37°C (1°C/sec), and extended at 37°C for at least 2 hours. Reactions were heat inactivated at 75°C for 20min and used directly for Gibson assembly into a digested receiving plasmid with a fixed barcode.

I.3 Cell culture

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

354 **I.4 Transfection, RNA extraction, and reverse transcription**

355 Cells were transfected at 60-80% confluence with FuGENE HD[®] according to the manufacturer's
356 protocol at a 3:1 FuGENE HD[®] to DNA ratio. For high-throughput measurements of splicing
357 outcomes, 10 µg pooled reporter assay DNA was transfected in three 100 mm plates. For biochemical
358 analysis of individual reporters, 1 µg or 2.5 µg individual reporter DNA was transfected into each
359 well of a 12- or 6-well plate (respectively). 24 hours after transfection, total RNA was isolated from
360 detached cells (Accutase[®], ThermoFisher). For amplicon sequencing, total RNA was isolated using
361 phenol-chloroform (Ambion) extraction (5PRIME Phase Lock Gel, Quantabio) followed by DNase
362 treatment (TURBO DNase). For biochemical analysis, RNA was isolated using a silica column
363 (illustra[™] RNAspin Mini RNA Isolation Kit, GE Healthcare) with on-column DNase digestion
364 following manufacturer's automated protocol. DNase-treated RNA was reverse transcribed using a
365 reporter-specific primer following manufacturer's specifications (SuperScript IV Reverse Transcrip-
366 tase, Thermo Fisher) with RNase H treatment. Reverse transcription primers included degenerate
367 nucleotides to serve as unique molecular identifiers (UMIs) during amplicon sequencing^{45,46}. cDNA
368 products were used for amplicon sequencing or biochemical analysis.

369 **I.5 Amplicon sequencing**

370 Amplicon sequencing was used to identify exon-barcode pairings in Lib1 and to quantify splicing
371 products from reverse-transcribed cDNA. Second-strand synthesis added additional UMIs in a single
372 anneal-extend cycle of 98°C for 10min, cooled slowly to 60°C (0.1°C/sec), annealed at 60°C for
373 5min, and extended at 72°C for 5min (Phusion[®] Hot Start Flex 2X Master Mix, NEB). Resulting
374 double-stranded amplicons were amplified using a two-stage procedure. In the first stage, targets
375 were amplified by PCR primers. PCR was performed using the following protocol: 98°C for 30sec
376 initial denaturation, then 16 cycles of 98°C denaturation for 10sec, 60°C annealing for 15sec, 72°C
377 extension for 1min 45sec, and a final extension step at 72°C for 5min (Phusion[®] Hot Start Flex
378 2X Master Mix, NEB). Longer extension times and minimal number of PCR cycles were used to
379 avoid recombination across exons and barcodes. The number of cycles was determined for each
380 sample by first running 10µL qPCR reactions (LightCycler[®] 480 SYBR Green I Master, Roche). In
381 the second stage, index primers were added using 5 PCR cycles. PCR was performed using the
382 following protocol: 98°C for 30 s initial denaturation, then 5 cycles of 98°C denaturation for 10sec,
383 71°C annealing for 15sec, 72°C extension for 1min 45sec, and a final extension step at 72°C for 5min
384 (Phusion[®] Hot Start Flex 2X Master Mix, NEB). Final DNA concentrations were measured using
385 fluorometric measurements (Qubit 1X dsDNA HS Assay, Thermo Fisher) on a Qubit 3 fluorometer.
386 Paired-end sequencing was carried out on an Illumina NextSeq 550 with 10% PhiX spiked in, with 54
387 cycles in read 1 (reverse) and 106 in read 2 (forward). About 4M paired-end reads (> 10X coverage)
388 were acquired for Lib1 exon-barcode sequencing and an average of 22M paired-end reads (> 50X
389 coverage) for each PSI quantification replicate.

390 **I.6 Biochemical analysis**

391 PCR amplification reactions to determine splicing products were carried out in 20µL reactions
392 containing 10µL OneTaq[®] 2X Master Mix with Standard Buffer (NEB), 200 nM each forward and
393 reverse primers (IDT), and 1µL cDNA. PCR was performed using the following protocol: 94°C for
394 30 s initial denaturation, then 25 cycles of 94°C denaturation for 10 s, 62°C annealing for 15 s, 68°C
395 extension for 20 s, and a final extension step at 68°C for 1 min. 5µL final PCR product was run
396 out on 2.0% agarose (Denville Scientific) Tris-acetate-EDTA (TAE) gel with ethidium bromide and
397 visualized on a Bio-Rad imager. Densitometry measurements to calculate PSI were measured using
398 Bio-Rad Image Lab (Windows v6.1).

399 I.7 Reporter assay preprocessing

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

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

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

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

414 where $n_{\text{inclusion}}$ is the total number of exon inclusion reads, and n_{skipping} is the total number of
 415 exon skipping reads. In addition to the measured PSI, the dataset includes for each barcode: (1)
 416 a 90 nucleotide sequence, containing the 70 nucleotide variable exon sequence plus the 10 fixed
 417 flanking nucleotides on each side; (2) structure in dot-bracket notation predicted by RNAFold (Vienna
 418 RNA²⁰, version 2.4.17), using default parameters; (3) an indicator vector indicating which nucleotide
 419 participates in a predicted G-U wobble base pair. The dataset was split randomly into a training set
 420 and a test set in an 80/20 split, using a fixed seed for reproducibility.

421 I.8 Model design

The model's input is a triple of vectors $(x_{\text{seq}}, x_{\text{struct}}, x_{\text{wobble}})$,

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

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

$$\begin{aligned} f_a^b : x &\mapsto \text{Sum}(\text{Softplus}(\text{Position-Bias}(\text{Convolution}(x; \alpha_a^b); \beta_a^b))) && \text{(SCM)} \\ \alpha_a^b &\in \mathbb{R}^{w_a^b \times c_a^b \times k_a^b}, && \beta_a^b \in \mathbb{R}^{(d-w_a^b+1) \times k_a^b} \end{aligned}$$

422 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
 423 $x = [x_{\text{seq}}, x_{\text{struct}}, x_{\text{wobble}}]$ (structure SCM). The 1D convolutional layer contains $k_a^b = 20$ convolutional
 424 filters of width $w_a^b = 6$ for each sequence SCM ($b = \text{seq}$), and $k_a^b = 8$ convolutional filters of
 425 width $w_a^b = 30$ for each structure SCM ($b = \text{struct}$). The number of input channels is $c_a^b = 4$ for
 426 sequence SCM (corresponding to the one-hot encoded four nucleotides) and $c_a^b = 8$ for structure
 427 SCM (corresponding to sequence, structure, and wobble indicator). The output of the convolution
 428 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
 429 $z + \beta_a^b$, 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

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

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 .

I.9 Model training

The model was implemented in Python 3.8⁴⁷ using Tensorflow 2.6⁴⁸ and Numpy 1.20⁴⁹. 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 hours on a mid-range 4-core with 16GB 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 (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 (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.

I.10 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. Sequences were padded to 70 nucleotides, and predicted PSI was then computed using our model. To account for differences in splice sites, flanking sequences, and cell types, one correction term was introduced per assay, as described previously¹⁶.

466 **I.11 Filter visualization**

467 To avoid reporting redundant sequence filters, hierarchical clustering using Scipy⁵⁰ was applied.
468 Each sequence filter was represented by a vector containing its total strength for each of the exons in
469 the dataset. The strongest filter in each cluster was then used to generate a sequence logo²². The
470 logo represents the set of 6-mers that lead to positive filter activation.

471 The structure filters included one G-poor filter and three stem loop filters. Since enumerating all
472 30-mers is not tractable, the G-poor sequence logo was computed by evaluating the filter on a subset
473 of sequences from our dataset. As the three stem loop filters differed in the length of the loop (short,
474 medium, long) but were otherwise very similar, they were considered as one cluster. Layer-wise
475 relevance propagation was used to visualize individual nucleotide contributions to filter strength⁵¹.

476 **I.12 Ruling out sequencing artifacts**

477 The reduction in measured PSI associated with the presence of the stem loop and G-poor stretch
478 features could potentially be a technical artifact due to decreased amplification or sequencing
479 efficiency of exon inclusion products. To rule this out, we verified that the presence of these features
480 was not accompanied by a decrease in the total number of sequencing reads, and was instead
481 accompanied by an increase in the number of exon skipping reads (Extended Data Fig. 6).

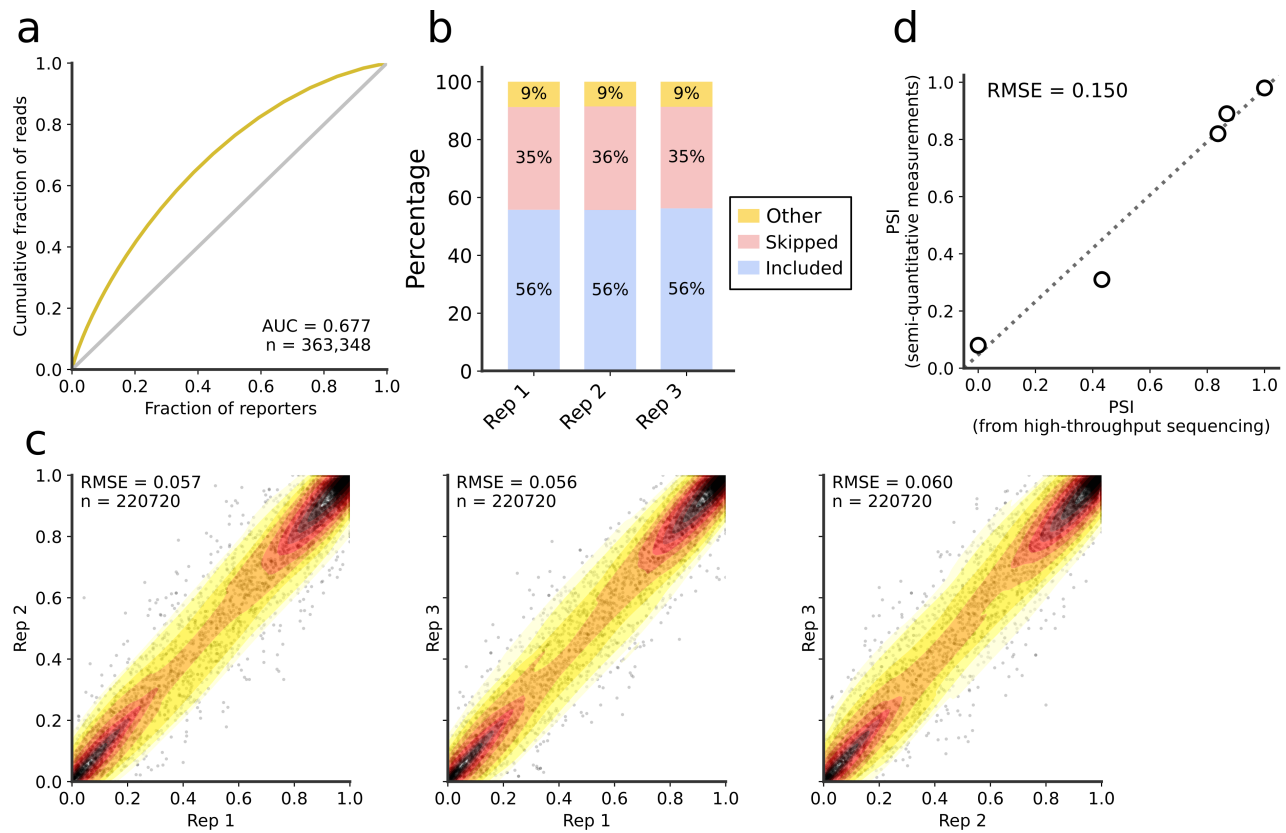
482 **I.13 Design of mutant constructs**

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

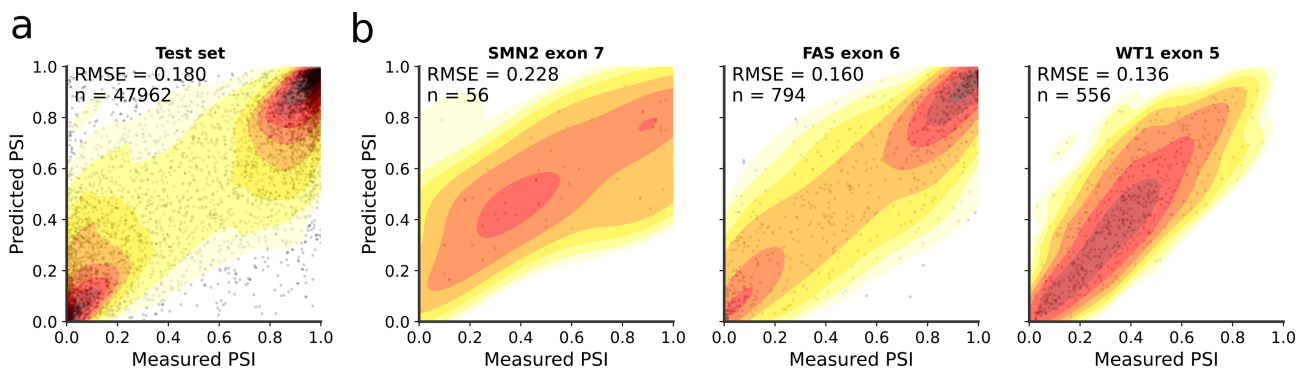
488 To validate the G-poor stretch feature, candidate exons that strongly activate the G-poor filter
489 exactly once along the exon were selected. For each candidate exon, a C-to-G mutation in the middle
490 of the activated filter's window was introduced. As before, to ensure this does not introduce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

493

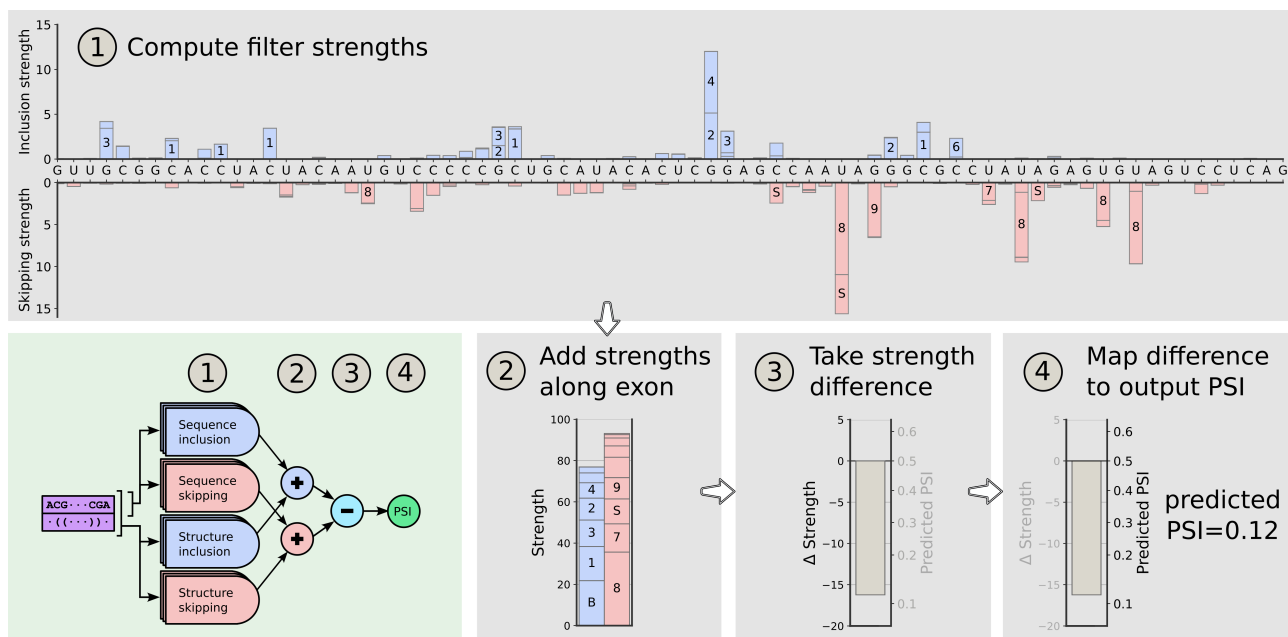
II. EXTENDED DATA FIGURES



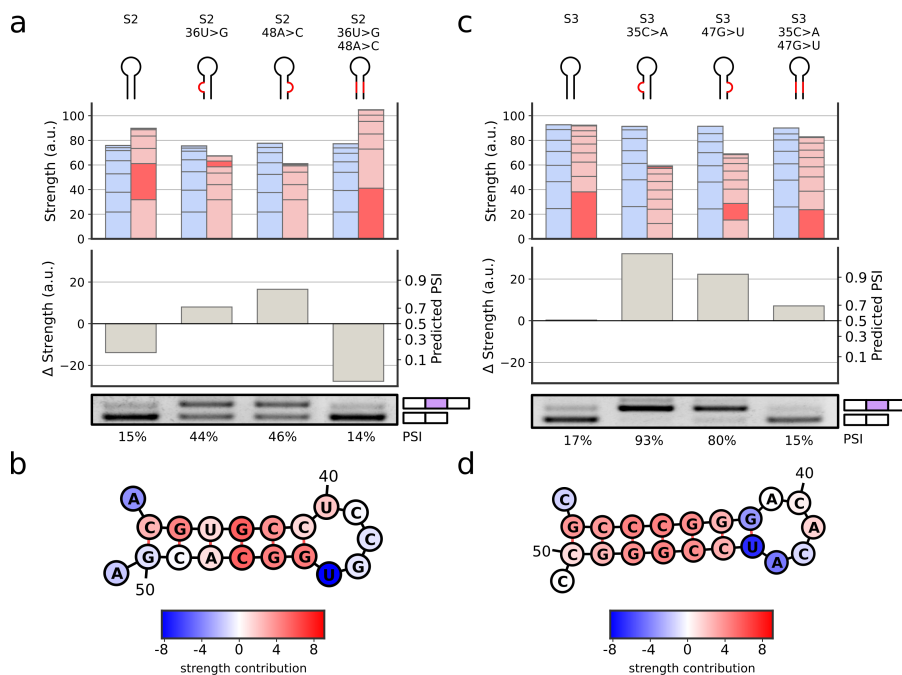
Extended Data Figure 1 | Assay quality checks. **a**, Lorenz plot showing the distribution of reads across the reporters in a high-throughput sequencing of the library DNA (gold). A perfectly even library would be on the diagonal (gray). **b**, Over 90% of splicing products corresponded to exon inclusion or exon skipping products, as measured by high-throughput RNA sequencing. **c**, Comparison of PSI measurements across the three biological replicates. **d**, Comparison of PSI measurements from high-throughput sequencing and semi-quantitative measurements for five individual reporters (V1-V5). AUC: area under curve; RMSE: root-mean-square error.



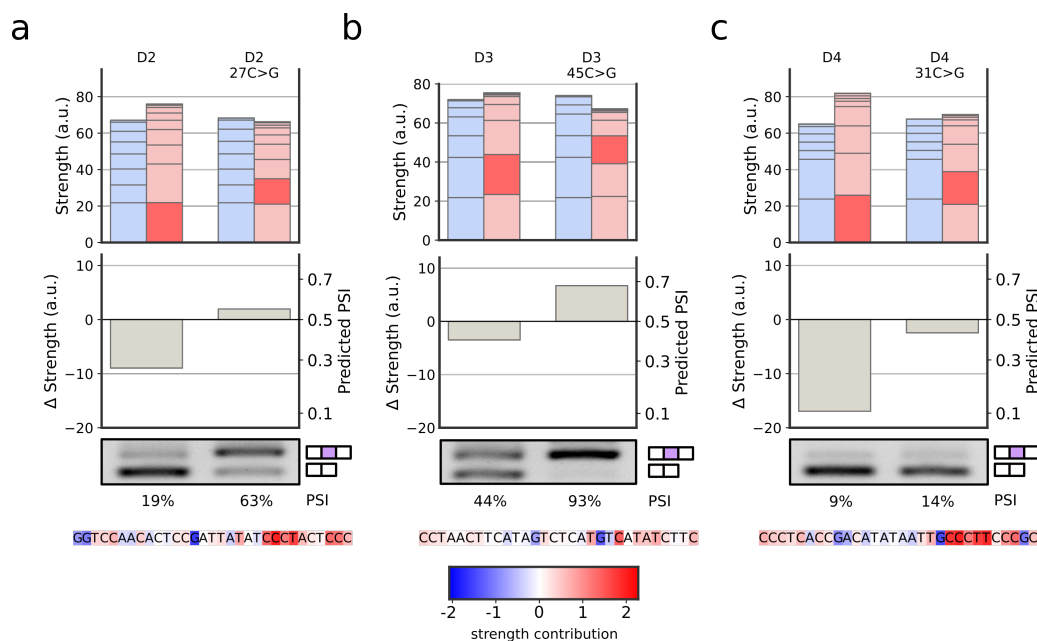
Extended Data Figure 2 | Model predictive accuracy. **a**, Predictions on the held-out experimental data. **b**, Predictions on previously-published assays: SMN2 exon 7 (C33a cells)^{8,13-15}, FAS exon 6 (HEK293 cells)¹⁶, and WT1 exon 5 (HEK293 cells)⁷.



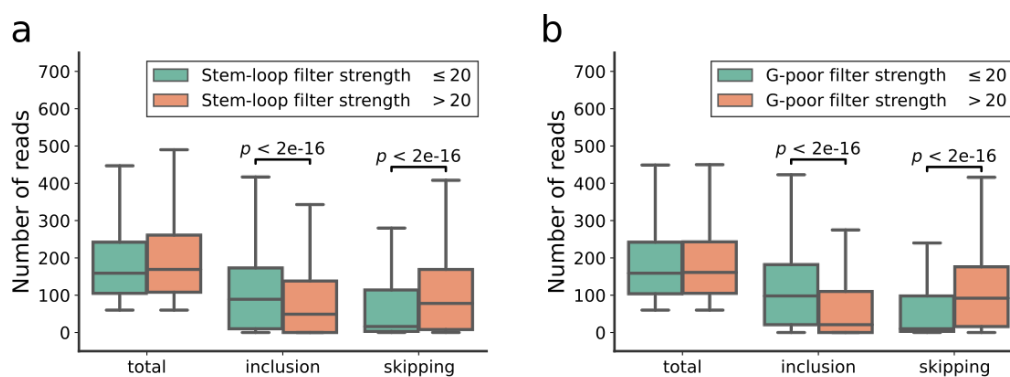
Extended Data Figure 3 | Tracing the neural network computation from exon sequence to predicted PSI. Our neural network (bottom left) arrives at its predictions in four steps. First, using exon sequence and predicted structure, it computes filter strengths for each position along the exon (1). Next, it adds all inclusion strengths together and all skipping strengths together (2). Then, the difference between these two strengths (Δ strength) is computed (3). Finally, it maps that difference to a predicted output PSI (4).



Extended Data Figure 4 | Additional validation of novel stem loop feature. Two more exons (S2, S3) were chosen and validated as in Figure 3.



Extended Data Figure 5 | Additional validations of the G-poor feature. Three more exons (D2, D3, D4) were chosen and validated as in Figure 4.



Extended Data Figure 6 | Effect of novel features on absolute read counts. **a**, Box plot showing the distribution of the total number of sequencing reads, the number of exon inclusion reads, and the number of exon skipping reads, for exons with stem loop strength at most 20 and greater than 20. **b**, As in panel **a** for G-poor strengths. Center line: median; box limits: upper and lower quartiles; whiskers: 1.5x interquartile range. p values: Student's t -test.