

Deep indel mutagenesis reveals the regulatory and modulatory architecture of alternative exon splicing

Received: 23 April 2024

Accepted: 4 August 2025

Published online: 30 August 2025

 Check for updates

Pablo Baeza-Centurión ^{1,2,3,7}, Belén Miñana ^{1,7}, Andre J. Faure ¹, Mike Thompson¹, Sophie Bonnal ¹, Gioia Quarantani ^{1,2}, Joseph Clarke ^{4,5}, Ben Lehner ^{1,2,4,6,8} ✉ & Juan Valcárcel ^{1,2,6,8} ✉

While altered pre-mRNA splicing is a frequent mechanism by which genetic variants cause disease, the regulatory architecture of human exons remains poorly understood. Antisense oligonucleotides (AONs) that target pre-mRNA splicing have been approved as therapeutics for various pathologies including patient-customised treatments for rare diseases, but AON discovery is currently slow and expensive, limiting the wider adoption of the approach. Here we show that deep indel mutagenesis (DIM) – which can be made experimentally at very low cost – provides an efficient strategy to chart the regulatory landscape of human exons and rapidly identify candidate splicing-modulating oligonucleotides. DIM reveals autonomous effects of insertions, while systematic deletion scans delineate the checkerboard architecture of sequential enhancers and silencers in a model alternative exon. The results also suggest a mechanism for repression of transmembrane domain-encoding exons and for the generation of microexons. Leveraging deep learning tools, we provide a resource, DANGO, that predicts the splicing regulatory landscape of all human exons and can help to identify effective splicing-modulating antisense oligonucleotides.

Pre-mRNA splicing is the process by which introns are removed from transcripts and exons are joined together to form mature mRNAs, which are then exported to the cytoplasm and translated¹. Altered splicing is an important mechanism by which genetic variants cause disease^{2–4}. Multiplex assays of variant effects (MAVEs) have revealed that random nucleotide substitutions in exons frequently affect splicing, with 60–70% of substitutions in over 90% of positions in alternatively-spliced exons^{5–7} and 5% of substitutions in constitutively-spliced exons⁸ altering exon inclusion. Comprehensive testing has also shown that ~10% of disease-causing missense variants, as well as 3% of common exonic substitutions, affect splicing^{9,10}.

The impact of genetic variation beyond substitutions on splicing has been far less studied. Insertions and deletions (indels^{11,12}) are abundant variants evident in 24% of Mendelian diseases¹³, and disease-causing indels are enriched close to splice sites¹⁴. Despite this, the effects of indels on splicing have not been systematically tested.

The frequent disruption of splicing in human disease has led to extensive efforts to therapeutically modulate splicing¹⁵. In particular, antisense oligonucleotides (AONs) that modulate splicing have been approved as therapies for spinal muscular atrophy¹⁶ and Duchenne muscular dystrophy¹⁷. Indeed AONs – because of their programmable sequence specificity – may represent a general strategy to

¹Centre for Genomic Regulation (CRG), The Barcelona Institute of Science and Technology, Barcelona, Spain. ²Universitat Pompeu Fabra (UPF), Barcelona, Spain. ³German Cancer Research Centre (DKFZ), Heidelberg, Germany. ⁴Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton, UK. ⁵University of Cambridge, Cambridge, UK. ⁶Institució Catalana de Recerca i Estudis Avançats (ICREA), Barcelona, Spain. ⁷These authors contributed equally: Pablo Baeza-Centurión, Belén Miñana. ⁸These authors jointly supervised this work: Ben Lehner, Juan Valcárcel. ✉ e-mail: ben.lehner@crg.eu; juan.valcarcel@crg.eu

therapeutically modulate many splicing changes¹⁸. However, predicting the impact of AONs on splicing is very challenging: AONs need to be long to achieve sequence specificity (typically 18 to 21 nts) whereas splicing regulatory elements are short and very poorly mapped genome-wide. Parameters that influence AON efficacy include length, proximity to splice sites and binding energy^{19–23}. In practice, however, identifying effective AONs requires laborious testing of many different designs.

Here, we use deep indel mutagenesis (DIM) to comprehensively quantify for the first time the impact of all possible deletions and all small insertions on the splicing of a human exon, FAS exon 6. This exon encodes a transmembrane helix of the FAS/CD95 receptor. Inclusion of this exon generates a pro-apoptotic receptor whereas skipping of this exon produces an anti-apoptotic soluble inhibitor^{24–27}. A variant in exon 6 that causes exon 6 skipping causes autoimmune lymphoproliferative syndrome (ALPS)²⁸. Our data show that insertions and deletions frequently disrupt splicing, with deletions in particular representing an efficient experimental strategy to map the regulatory architecture of the exon. The deletion scan also defines the minimum length of an exon, and reveals a novel repressive regulatory mechanism whereby cryptic sequences within an exon are recognised by a core component of the spliceosome (U2 snRNP) that normally recognises intronic branch points. We argue that this recognition of cryptic exonic elements might be a mechanism for the evolutionary birth of unusually short microexons²⁹. Using our data, we show that deep learning methods accurately predict the effects of indels on splicing and that predicted deletion effects accurately predict the splice-altering effects of AONs in a model alternatively spliced exon. Finally, we provide a genome-wide resource, DANGO, that charts splicing regulatory landscapes and can help to identify splice-modifying AONs across the human exome.

Results

Deep indel mutagenesis of a human alternatively spliced exon

To quantify in parallel how single nucleotide (nt) substitutions, deletions and insertions affect alternative splicing of a human exon, we designed a library containing all single-nt substitutions in the 63 nt-long FAS exon 6 ($n = 189$), all deletions ranging in length from 1 to 60 nts ($n = 2010$), all possible 1-, 2-, and 3-nt insertions ($n = 5208$), and a random selection of 585 4-nt-long insertions (Fig. 1A). We cloned this library into a plasmid minigene vector spanning FAS exons 5–7, transfected it into HEK293 cells, isolated RNA 48 h post-transfection and used deep sequencing of reverse transcription (RT)-PCR products to quantify the inclusion of each variant by counting how often it is present in the final exon inclusion product compared to each of the other variants in the library (Fig. 1B). The resulting enrichment scores allow the percent-spliced-in (PSI) value of each variant to be measured by reference to the PSI value of the wild type sequence (49.1%). Our results were highly reproducible across nine experimental replicates (Pearson's r between 0.97 and 0.98 for all pairs of replicates, Supplementary Fig. 1) and our estimated PSI values were very well correlated with PSI values determined by RT-PCR for 40 individual, independently transfected, mutant minigenes, which included indel and substitution mutants (Spearman's $\rho = 0.95$, Fig. 1C). As the minigenes do not contain natural AUG start codons, and other AUGs are neither in a good Kozak context for translation, nor are followed by long ORFs, Nonsense Mediated Decay is unlikely to contribute to the effects of mutations. We note, however, that other effects on RNA stability cannot be ruled out.

The effects of single-nt substitutions, insertions and deletions on exon inclusion

We quantified the effects of all single-nt substitutions ($n = 189$), insertions ($n = 187$) and deletions ($n = 63$) on the inclusion of FAS exon 6 (Fig. 1D, E). The results revealed that nearly the entire range of

inclusion values can be obtained as a consequence of at least one of the three classes of variant. Regardless of the type of mutation, just under two thirds affected splicing by more than 10 PSI units (Fig. 1D, E): 61.9% of substitutions changed the inclusion of FAS exon 6 by more than 10 PSI units (mean absolute Δ PSI = 17.5 PSI units; consistent with previous data⁵), similar to 58.7% of single-nt deletions (mean absolute Δ PSI = 15.5 PSI units) and 61.3% of single-nt insertions (mean absolute Δ PSI = 19.7 PSI units).

Substitutions have a mode near the wild type (WT) PSI (49.1%) and they more often promote skipping than inclusion, with 47.6% decreasing and 14.3% increasing inclusion by more than 10 PSI units. Similarly, deletions tend to promote skipping, with 15.9% of single-nt deletions promoting inclusion and 42.9% promoting skipping by more than 10 PSI units (Fig. 1E, see also Fig. 2A). In contrast, insertions more frequently promote inclusion, with 41.1% increasing the inclusion of FAS exon 6 and only 20.2% decreasing inclusion by more than 10 PSI units (Fig. 1E).

We next compared the effects of different mutation types in the same position along the exon. The effects of single-nt substitutions and deletions in the same positions correlate moderately well (Spearman's $\rho = 0.46$, Fig. 1F), consistent with some of these variants disrupting existing regulatory elements and with positive regulatory elements covering a larger proportion of the exon than negative elements (Fig. 1D, E). However, the effects of substitutions and deletions correlate poorly with the effects of insertions before or after the substituted/deleted position (Fig. 1G, H, Supplementary Figs. 2, 3). This suggests that insertions often affect splicing by a different mechanism, for example, by the creation of new regulatory sequences (see below).

The effects of short multi-nt deletions and insertions on exon inclusion

We next analysed the influence of short multi-nt deletions (ranging from 2 to 9 nts) and insertions (either 2,3 or 4 nts) on exon inclusion. Short multi-nt deletions had effects on exon inclusion which were similar to those of single-nt deletions. Thus, 64.5% of 2-nt deletions affected splicing by more than 10 PSI units (mean absolute Δ PSI = 17.0), as well as 62.3% of 3-nt deletions (mean absolute Δ PSI = 17.9) and 73.3% of 4-nt deletions (mean absolute Δ PSI = 20.8). Considering all short deletions up to 10 nts long, 66.0% altered inclusion by more than 10 PSI units (mean absolute Δ PSI = 18.7). Short multi-nt deletions longer than 2nt more frequently promote skipping, with 30.6% of 2-nt deletions, 39.3% of 3-nt deletions, 43.3% of 4-nt deletions and 45.8% of all deletions spanning 10 or fewer nts decreasing exon inclusion by more than 10 PSI units. In contrast, 33.9% of 2-nt deletions, 23.0% of 3-nt deletions, 30% of 4-nt deletions and 20.2% of all deletions up to 10 nts long increase inclusion by more than 10 PSI units.

Longer multi-nt deletions are particularly effective at revealing splicing regulatory elements (Fig. 2A). Namely, we found that while 1–4 nt deletions at the 5' end of the exon showed multiple, seemingly contradictory effects, longer deletions delineated -8 5' terminal nts whose collective deletion strongly promotes exon skipping. Similarly, deletions that cover exon positions 9 to 18 consistently increase exon inclusion. Interestingly, deletions that partially overlap this region (on either side) have the inverse effect and promote exon skipping. Systematic deletion scans thus delineate consistent discrete regulatory elements (enhancers and silencers), often adjacent to each other, and reveal patterns such as alternating elements with antagonistic effects (Fig. 2A).

Short multi-nt insertions had a stronger effect on FAS exon 6 inclusion compared to single-nt insertions: 74.5% of 2-nt insertions and 77.5% of 3-nt insertions changed inclusion by more than 10 PSI units (mean absolute Δ PSI = 23.9 and Δ PSI = 24.2, respectively). Our library also contained a random selection of 585 4-nt insertions, of which 88.6% changed splicing by more than 10 PSI units (mean absolute Δ PSI = 31.2). Similar to single-nt insertions, but in contrast with short

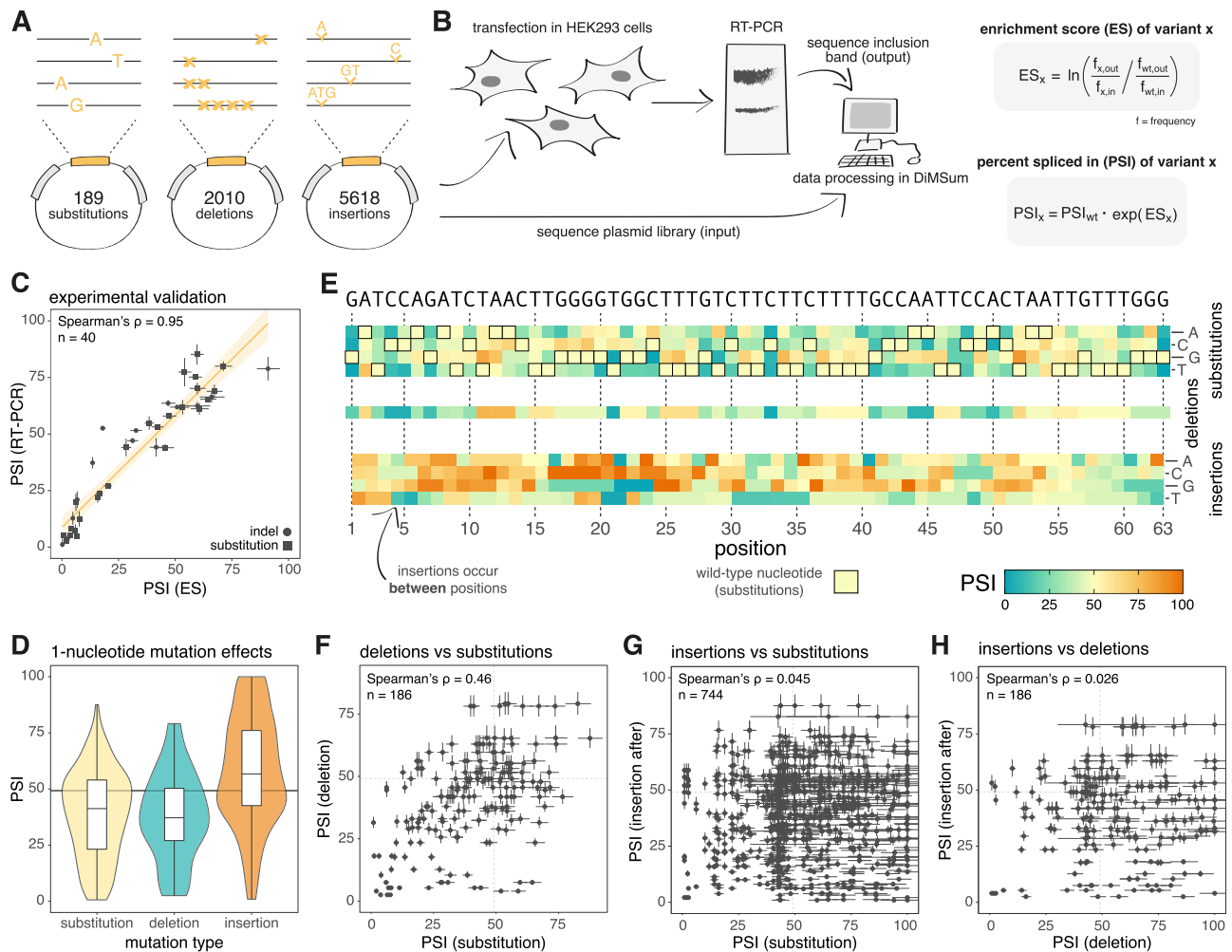


Fig. 1 | Deep indel mutagenesis of FAS exon 6. A Plasmid library design. **B** Experimental protocol of our massively-parallel splicing assay (MPSA). **C** Correlation between PSI values from 40 individual transfections and MPSA-derived PSI values. Error bars = 95% confidence interval for the mean. **D** Distribution of PSI values for variants with 1-nucleotide (nt) mutations. The horizontal line shows the PSI value of the reference wild type clone. Substitutions $n = 189$. Deletions $n = 63$. Insertions $n = 62$. Boxplot definition: centre line = median (50th percentile); box bounds = 25th–75th percentile; whiskers = show minima and maxima within 1.5

* interquartile range. **E** Heatmaps displaying inclusion levels of 1-nt substitutions, deletions, and insertions. **F** Correlation between PSI values for 1-nt deletions and substitutions at the same position. Error bars = 95% confidence interval for the mean. **G** Correlation between PSI values for 1-nt insertions after a given position and substitutions at that position. Error bars = 95% confidence interval for the mean. **H** Correlation between PSI values for 1-nt insertions after a given position and 1-nt deletions at the same position. Error bars = 95% confidence interval for the mean. Source data are provided as a Source Data file.

multi-nt deletions, short multi-nt insertions tended to promote inclusion rather than skipping: 44.6% of 2-nt insertions increased exon inclusion by more than 10 PSI units (compared to 30.0% which decreased inclusion by this same amount), as well as 42.4% of 3-nt insertions (35.1% for skipping) and 57.5% of 4-nt insertions (31.2% for skipping).

In contrast to the regulatory landscape emerging from deletion analyses (Fig. 2A), double- or triple-nt insertions tended to show autonomous effects that were strongly influenced by the nature of the inserted nts (Fig. 2B, C, Supplementary Fig. 4). Specifically, insertion of GC, CG or GA dinucleotides resulted in general increases in exon inclusion, almost independently of the position of the insertion, with the prominent exception of the exon 3' end. In contrast, insertion of GG or CC showed markedly different effects depending on the site of insertion.

The effects of triple nt insertions were even more idiosyncratic (Fig. 2C). For example, all CG-containing triplets enhanced exon inclusion in nearly all positions (as was the case of GC, CG or GA dinucleotides), with the notable exception of the five 3' nts of the exon.

Almost any insertion in this region (except those positioning an A at the 3' end of the exon) led to enhanced exon skipping (black vertical rectangles on the right of panel 2B–C), suggesting that this region harbours a strong enhancer sequence very sensitive to insertions or deletions that may be important for activation of the adjacent 5' splice site. Effects similar to those of CG-containing triplets were observed upon insertion of a variety of GC- / GA- / GG-containing triplets (lower rows of Fig. 2C); some of these effects might be related to enrichment in purine residues which, together with other purines present in the insertion site, could function as purine-rich exonic enhancers, a well-known class of exonic regulatory elements³⁰. Indeed, these effects are also observed for other purine-rich triplets such as GAG or AAG, albeit not for all (e.g., AGG or GGG). Triplets containing AU/UA dinucleotides (e.g., UAG, UAA, UUA, AUU, CUA) promote skipping when inserted at most exonic positions (cluster of blue colour in mid-low rows of panel 2C), which might be explained by enhanced binding of hnRNP proteins such as hnRNP A1³¹, known to mediate effects of exonic silencers³². Pyrimidine-rich triplets (e.g., UCC, UUC, CUU), which could provide/reinforce binding sites for other repressive hnRNPs such as PTB/

forming G-quadruplexes recognised by hnRNP F/H factors³⁴, is disrupted by C-containing triplets. The latter effects could be in part linked to the creation of CG dinucleotides, which as discussed above display strong enhancing effects, or to effects of other G-rich enhancer sequences^{35,36}.

We next evaluated the extent to which the effects of 2- and 3-nt insertions in FAS exon 6 predict the association between 2- and 3-nt kmer content and exon inclusion transcriptome-wide. Strikingly, there is a strong positive correlation (Spearman's rho between 0.72 and 0.87 for 2-mers, between 0.69 and 0.79 for 3-mers) between the Δ PSI induced by kmer insertions in our library and the PSI of exons containing at least 20 (for 2-mers) or 10 (for 3-mers) such kmers in the GTEx database, either in adipose tissue (Fig. 2D) or in all GTEx tissues (Supplementary Fig. 5). Furthermore, the relationship between the content of each individual nucleotide in the inserted 3-mers and the PSI of the mutated exon (Fig. 2E) resembles the relationship between the content of each individual nucleotide in cassette exons and their PSI (Fig. 2F). For example, increasing the number of uridines in an inserted triplet correlates with increased exon skipping (Fig. 2E), similarly to how a larger uridine content is associated with increased skipping in exons throughout the genome (results for GTEx adipose tissue are shown in Fig. 2F, results for all GTEx tissues are shown in Supplementary Figs. 6–9). These results argue that the effects learned from systematic analysis of insertion mutations observed in our experiment for one individual model exon have captured sequence features relevant for the inclusion of alternatively-spliced exons genome-wide.

Small insertions create novel microexons by activating cryptic splice sites within FAS exon 6

We next analysed the relationship between exon length and inclusion in our library, examining the effects on inclusion of all possible deletions from size 1 to 60 nts. Specifically, we sought to elucidate, at 1 nt resolution, how short an exon can be while still being recognised by the splicing machinery. There is no clear dependence of exon inclusion on length for exons longer than 50 nts (up to 13 nt-long deletions). At exon lengths less than 50, inclusion gradually decreases with increasing deletion length (i.e., shorter exons are less included), with almost no exons shorter than 30 nts showing detectable levels of inclusion in our library (Fig. 3A). Consistent with this and previous large deletions in constitutive exons^{37–39}, exons shorter than 30 nts are more likely to be skipped genome-wide compared to longer exons (results for GTEx adipose tissue are shown in Fig. 3B, results for all GTEx tissues are shown in Supplementary Fig. 10).

Exons shorter than 27 nts are detected in multicellular animals but are categorised as a special class - microexons - whose recognition requires a dedicated set of regulatory sequences and factors (such as SRRM3/4) that enable their inclusion in specific tissues (e.g., the brain or endocrine pancreas^{40,41}). However, we observed that a group of very short exons (microexons, length-wise) from our library were detectably included (Fig. 3A). These microexons might correspond to large deletions at the 5' or 3' ends of FAS exon 6 (Supplementary Fig. 11A). Nonetheless, these deletion mutants showed no evidence of exon inclusion when tested individually (Supplementary Fig. 11B). We realised that these apparently contradictory results could be explained if the clones detected by deep sequencing of the exon inclusion amplicon product were not the result of splicing of very short exons flanked by FAS exon 6 splice sites, but rather the result of the use of cryptic splice sites within exon 6 that have been activated by another mutation that is no longer present in the spliced-in exon sequence. The central part of FAS exon 6 (positions 24–40) contains a pyrimidine-rich tract that resembles the polypyrimidine (Py)-tracts that precede 3' splice sites, and nts in positions 10–15 contain a sequence that, strikingly, matches a branch point sequence⁴² (Fig. 3C). We reasoned that, if an AG-containing kmer

were to be introduced after the pyrimidine (Py)-rich segment, this would result in a 3' splice site-like sequence arrangement that, if recognised by the spliceosome, could create a novel microexon spanning from this new 3' splice site to the 3' end of exon 6 (Fig. 3D). To test this possibility, we inserted AG-creating triplets (like CAG or CUA -as the next nt is a G-) after the Py-tract in the FAS exon 6 minigene and observed that exons corresponding to the final part of the exon were included to some extent in the mature RNA (Fig. 3E). In contrast, insertion of GAC, which does not create a 3' splice site, does not result in activation of a shorter exon, but rather enhances the inclusion of its full-length version (Fig. 3E).

Interestingly, the exonic Py-tract is longer and more uridine-rich than the Py-tract associated with the natural 3' splice site of intron 5 (Fig. 3C). To assess whether the interplay between these 3' splice sites plays a role in regulation, we strengthened the 3' splice site by replacing the branch point and polypyrimidine tract of intron 5 by a consensus branch point followed by a stretch of 16 uninterrupted pyrimidines (Fig. 3F, see also Methods). In the presence of this mutation, inclusion of the full-length exon was enhanced in the wild type minigene and activation of the cryptic 3' splice site in the AG-containing construct was greatly reduced (Fig. 3G).

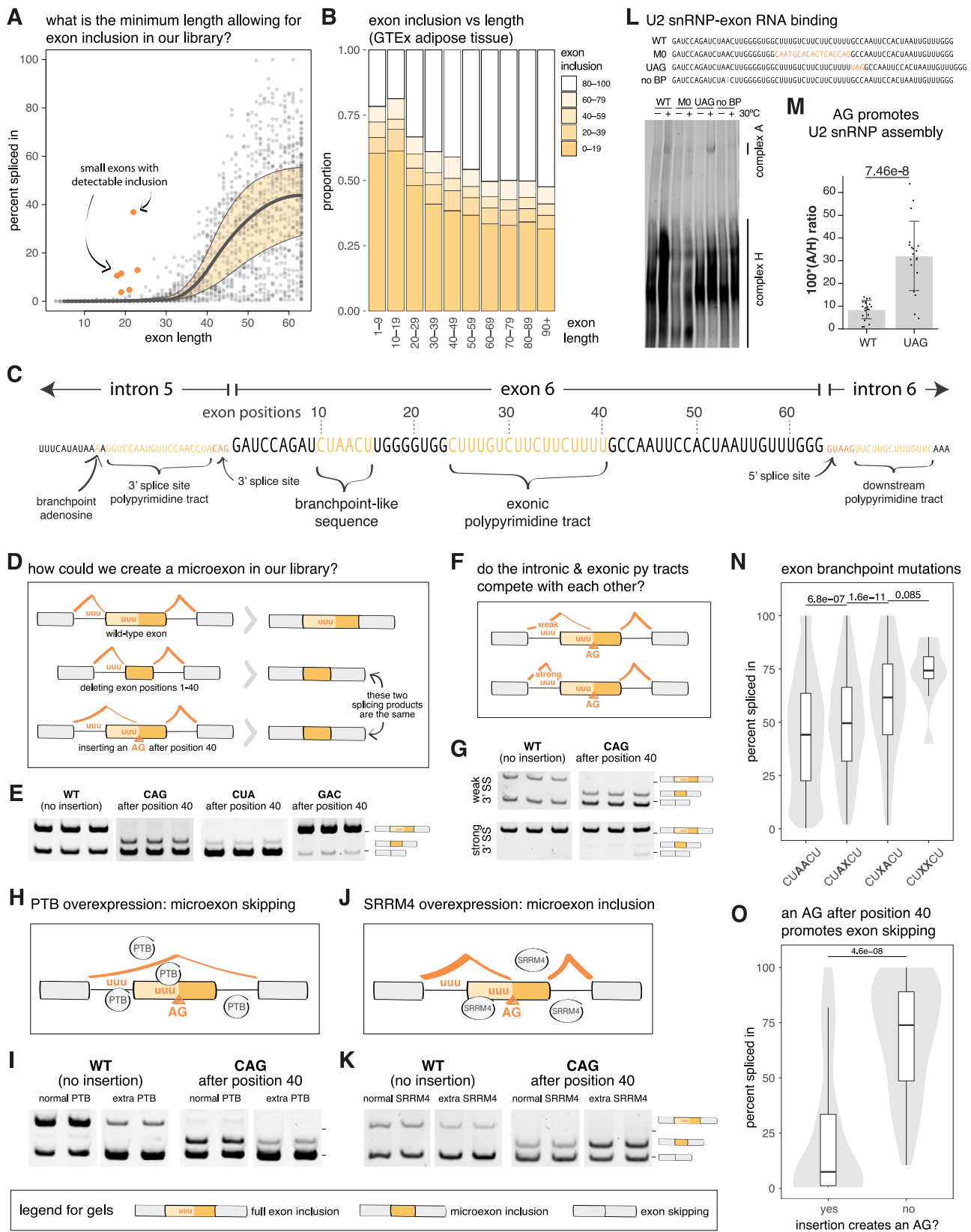
Previous work showed that the pyrimidine-rich sequence within FAS exon 6 functions as a silencer when bound by PTB³³ (Fig. 3H). As expected, overexpression of PTB led to skipping of the wild-type exon (Fig. 3I, left panel) and also to reduced inclusion of the shorter exon in the AG-containing construct (Fig. 3I, right panel), most likely due to direct competition between PTB and the Py-tract-binding splicing factor U2AF^{33,43}. Finally, we tested whether over-expression of SRRM4, which triggers inclusion of microexons in neurons⁴⁰, has any effect on inclusion of the shorter version of FAS exon 6 (Fig. 3J). Surprisingly, SRRM4 overexpression enhanced inclusion of the shorter exon (Fig. 3K) despite this exon not being flanked by *cis*-acting sequences typically required for the inclusion of microexons (e.g., intronic UGC motifs⁴⁴). Interestingly, SRRM4 reduced, rather than enhanced, inclusion of the wild-type full length exon 6 (Fig. 3K). These results show that the short exon activated by a cryptic 3' splice site in exon 6 is not only recognised by the splicing machinery but can be subject to splicing regulation by mechanisms similar to those operating on natural microexons.

Our results thus far do not account for the detection of short exons spanning the first third of FAS exon 6 (Fig. 3C). However, inserting sequences (e.g., UAA) after exon positions 18 or 19, mimicking a 5' splice site (G/GUAAGG), induced the accumulation of spliced products containing these sequences (Supplementary Fig. 11C). Interestingly, the Py-tract in the central region of FAS exon 6 is similar to a Py-tract found downstream of FAS exon 6 5' splice site (Fig. 3C), which is recognised by the protein TIA1 and enhances 5' splice site recognition by U1 snRNP^{45,46}. The Py-tract in the central region of FAS exon 6 could therefore enhance recognition of upstream 5' splice sites generated by exonic mutations.

These findings reveal that FAS exon 6 contains 3' and 5'-like sequences that can be activated by simple mutations to function as bona fide splice sites, promoting the inclusion of very short exons even in the absence of regulatory sequence elements known to be involved in the activation of microexons. The evolutionary birth of new microexons is therefore likely to be simpler and more frequent than previously appreciated.

Exonic binding of U2 snRNP promotes FAS exon 6 skipping

The presence of a relatively strong Py-tract preceded by a near-consensus branch point sequence within exon 6 (Fig. 3C), as well as the consequential inclusion of a shorter version of the exon when a mutation creates a functional 3' splice site AG downstream of the Py-tract, opened the possibility that 3' splice site recognising factors assemble on FAS exon 6 (Fig. 3L, M).



To directly assess whether U2 snRNP, the key ribonucleoprotein complex involved in 3' splice site recognition, can assemble on FAS exon 6 sequences, we incubated *in vitro* transcribed FAS exon 6 (wild type and mutants, all lacking the flanking splice sites) with HeLa nuclear extracts and measured the interaction by native gel electrophoresis. The results indicated that U2 snRNP can indeed assemble (complex A)

on FAS exon 6, an interaction that was decreased upon mutation of the Py-tract or branch point sequences and enhanced upon introduction of an AG dinucleotide (Fig. 3L, M and Supplementary Fig. 12).

It is conceivable that U2 snRNP assembly on the wild-type exon, in the absence of a 3' splice site, competes with recognition of the 3' splice site of intron 5 by the splicing machinery and that this

Fig. 3 | Deep indel mutagenesis of FAS exon 6 reveals the origin of novel microexons. **A** Relationship between exon length and inclusion (black curve represents a constrained B-spline fit to rolling median PSI values, and the yellow shaded area indicates the rolling interquartile range of PSI values). **B** Distribution of cassette exon inclusion values versus exon length (GTEX adipose tissue). **C** Sequence of FAS exon 6 and surrounding intronic sequences. **D** Hypothetical mechanisms explaining how our experimental assay could result in the detection of a microexon in our mutant library. **E** RT-PCR analysis of FAS exon 6 inclusion for the WT exonic sequence, two variants with insertions that introduce an AG after exon position 40 and one variant with an insertion after exon position 40 that does not introduce an AG dinucleotide. The figure shows the results of 3 technical replicates. Two additional biological replicates were performed with consistent results. The position of molecular weight markers (150 and 200 bp) is indicated for this and subsequent panels. **F** Illustration showing the impact of a weak or strong 3' splice site before FAS exon 6 on the recognition, by the splicing machinery, of the exonic 3' splice-site-like sequence in the central region of the exon (followed by an AG insertion). **G** RT-PCR analysis of FAS exon 6 inclusion in the presence of a weak or strong 3' splice site, for the WT exon and an exon with an insertion introducing an AG dinucleotide after exon position 40. The Figure shows results of 3 technical replicates. Two additional biological replicates were performed with consistent results. **H, I** RT-PCR analysis of FAS exon 6 inclusion upon overexpression of PTB (WT sequence and variant with an insertion introducing an AG dinucleotide after exon position 40). The Figure shows results of 2 technical replicates. Two additional biological replicates were performed with consistent results. **J, K** RT-PCR analysis of FAS exon 6 inclusion upon overexpression of SRRM4 (WT sequence and

variant with an insertion introducing an AG dinucleotide after exon position 40). The Figure shows results of 2 technical replicates. Two additional biological replicates were performed with consistent results. **L** Spliceosome assembly assays using the indicated fluorescently-labelled RNAs (wild type or mutant Fas exon 6 sequences) and HeLa nuclear extracts. The position of complexes assembling U2 snRNP (A complex) and hnRNP proteins (H complex) are indicated. **M** Quantification of the ratio between A and H complexes for wild type (WT) and 3' ss-containing (AUG) mutant RNAs as in **L**. *P* value corresponds to a 2-tailed t-test from 20 replicates with WT RNA, and 19 replicates with UAG-containing RNA. Error bars represent ± 1 standard deviation of the data. **N** Distributions of PSI values for double-nucleotide substitutions targeting (i) neither of the putative exonic branchpoint adenines; (ii) the second putative exonic branchpoint adenine as well as another nucleotide in the exon; (iii) the first putative exonic branchpoint adenine as well as another nucleotide in the exon; (iv) both putative exonic branchpoint adenines. *P* values correspond to 2-tailed Wilcoxon tests with 16230 data points in the CTAAC group, 544 data points in the CTAXCT group, 549 data points in the CTXXCT group, and 9 data points in the CTXXCT group. Boxplot definition: centre line = median (50th percentile); box bounds = 25th-75th percentile; whiskers = show minima and maxima within 1.5 * interquartile range. **O** Distribution of PSI values for insertions that either introduce (left) or do not introduce (right) an AG dinucleotide after exonic position 40. *P* value corresponds to a 2-tailed Wilcoxon test with 60 data points in the "no" group and 29 data points in the "yes" group. Boxplot definition: centre line = median (50th percentile); box bounds = 25th-75th percentile; whiskers = show minima and maxima within 1.5 * interquartile range. Source data are provided as a Source Data file.

contributes to modulate the levels of exon 6 inclusion. To test this possible mechanism, we took advantage of our saturation mutagenesis results. We observed that exonic variants with an intact branchpoint-like sequence at positions 10–15 (CUAACU) displayed an average inclusion of 45%. Variants harbouring mutations at either of the two adenosines that could serve as branch sites in this sequence, however, increased the levels of exon inclusion, and mutation of both adenosines further increased exon inclusion to an average of 75% (Fig. 3N). Also consistent with our model, insertion of AG-containing sequences after position 40 reduced full length exon inclusion, compared to insertion of non-AG-containing sequences, from an average of 75% to 25% inclusion (Fig. 3O).

Collectively, our results reveal a novel mechanism of exon skipping based upon assembly of U2 snRNP on exonic sequences that resemble (but cannot be active as) 3' splice sites. This illustrates the value of saturation mutagenesis approaches to discover and test mechanistic hypotheses.

Cryptic 3' splice sites regulate alternative exons encoding one-pass transmembrane helices

It has previously been reported that the Py-tract binding protein U2AF2 binds to an exonic polypyrimidine tract in IL7R exon 6⁴⁷, promoting exon skipping. Like FAS exon 6, IL7R exon 6 encodes a one-pass transmembrane helix. Interestingly, transmembrane helices are enriched in nonpolar amino acid residues that are encoded by codons with the highest number of pyrimidines (Fig. 4A). Therefore, transmembrane-encoding exons are expected to be rich in pyrimidines, allowing regulation by mechanisms similar to those described for IL7R exon 6 or FAS exon 6 (Fig. 3).

To investigate whether alternative exons, and particularly those coding for transmembrane domains, are generally regulated by this type of mechanism, we first used SVM-BPfinder⁴⁸ to identify throughout the genome exons containing features of 3' splice site regions, including a branchpoint motif followed by a polypyrimidine tract. For each input sequence, this tool returns a score ('SVM-BP score') that reflects the predicted strength of a 3' splice site. Interestingly, 23% of all exons had an SVM-BP score greater than that of FAS exon 6 (1.19), suggesting that a significant proportion of exons across the genome may have cryptic 3' splice sites or at least sequence elements that resemble 3' splice site regions. We found that shorter exons (<100 nts)

with SVM-BP scores greater than 1.19 encoded the most hydrophobic amino acid sequences (Fig. 4B), consistent with transmembrane domains. These exons had a lower average PSI compared to other exons (Fig. 4C), which could be attributed mainly to the effect of their polypyrimidine tract (Supplementary Fig. 13). These results are compatible with the existence of a category of (relatively) short exons containing 3' splice site-like sequences and, in particular, those encoding individual transmembrane helices (Fig. 4D), whose inclusion is decreased by exonic 3' splice site-like sequences.

To experimentally validate this hypothesis, we built minigenes containing one-pass transmembrane domain-encoding exon 5 of CHODL⁴⁹ and exon 6 of CXADR⁵⁰, and mutated their exonic putative branchpoint adenosines (Fig. 4E, F). In both cases, the mutations reduced exon skipping, suggesting that the levels of exon skipping were regulated by the recognition of 3' splice site-like sequences. Cryptic 3' splice sites in exons may therefore be a widespread mechanism to regulate the inclusion of alternative exons encoding transmembrane helices, and thus modulate the balance between soluble and membrane-bound protein isoforms.

Systematic deletion scans reveal the checkerboard regulatory landscape of a human exon

Deletions are more likely to cause loss-of-function molecular effects than insertions, which are more likely to create new regulatory elements. Plotting the effects on inclusion of all possible deletions from size 1 to 60 nts provides a comprehensive deletion map of FAS exon 6 (Fig. 5B, bottom triangle). The map reveals several regions whose deletion leads to enhanced inclusion (exonic silencers, red areas) or enhanced skipping (exonic enhancers, more green areas), arranged in an alternating 'checkerboard' pattern and covering most of the exon length. These regulatory elements can also be visualised by plotting the effects on exon inclusion of 1 to 6 nt-long deletions versus the position of the deletion along the exon (Fig. 5C, left panel) and using local polynomial regression (LOESS) to identify sequence 'blocks' where deletions promote more inclusion or more skipping. These regulatory blocks recapitulate well-known regulatory elements in FAS exon 6, namely EWS binding exonic enhancers at positions 15–23 and 55–63³⁵, a PTB-binding silencer in positions 25–40³³, and an SRSF6 enhancer at positions 40–45³⁶. Further, regulatory blocks suggest other elements that have not

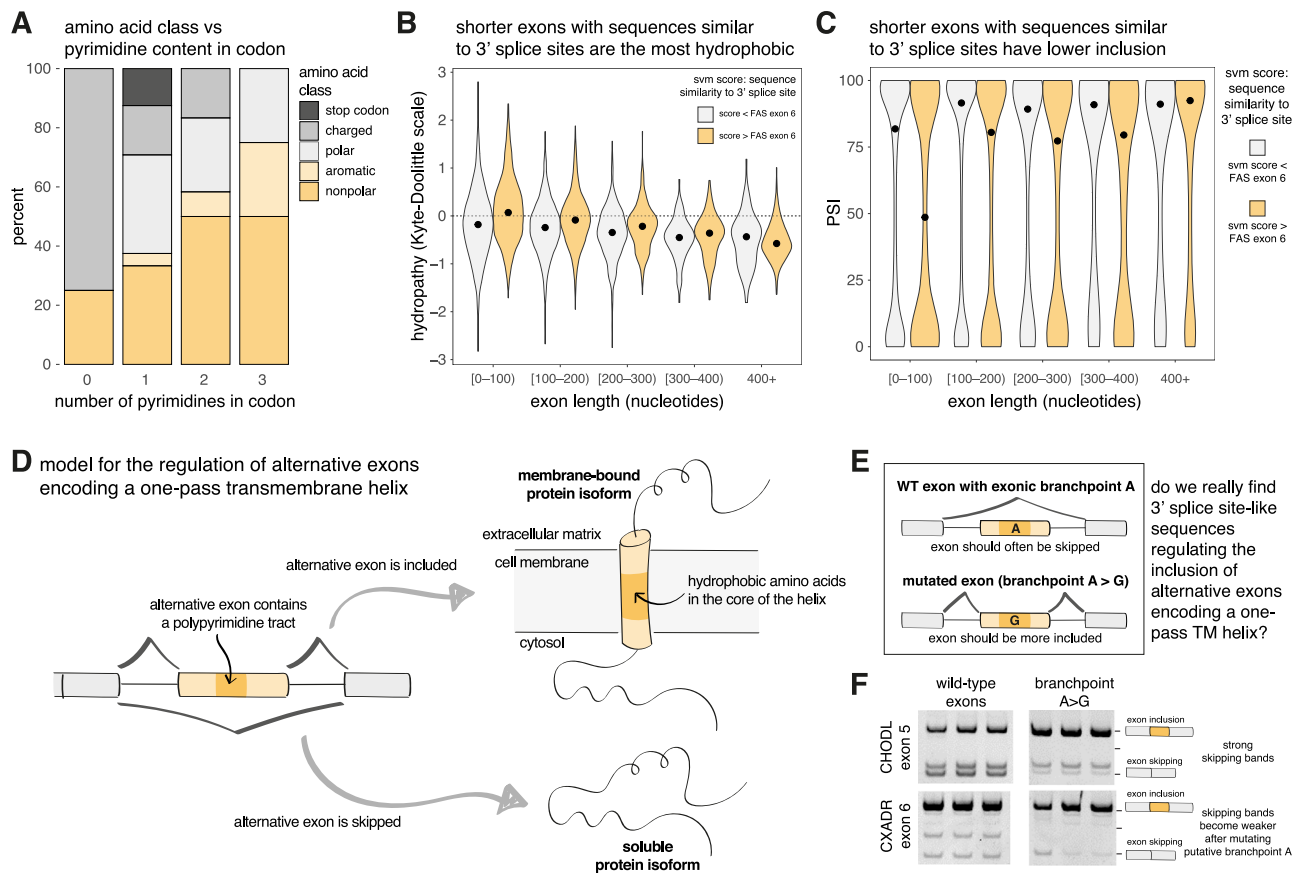


Fig. 4 | The inclusion of short alternative exons encoding one-pass transmembrane helices is regulated by exonic 3' splice site-like sequences. **A** Amino acid category encoded by each codon, categorised by the number of pyrimidines in the codon. **B** Hydropathy score of exons categorised into different length groups, divided by whether they have an exonic sequence more or less similar to a 3' splice site than FAS exon 6. **C** Inclusion of exons categorised into different length groups, divided by whether they have an exonic sequence more or less similar to a 3' splice site than FAS exon 6. **D** Model for the regulation of alternative exons encoding a one-pass transmembrane helix. **E** Hypothesis suggesting that 3' splice site-like

sequences in alternative exons encoding transmembrane helices promote exon skipping. **F** RT-PCR analysis of the inclusion of two alternative exons (CHODL exon 5 and CXADR exon 6) that each encode a one-pass transmembrane helix, including wild-type (WT) sequences and sequence variants with the putative branchpoint adenine mutated. The Figure shows results of 3 technical replicates. Two additional biological replicates were performed with consistent results. The position of molecular weight markers (150, 200 and 250 bp) is indicated. Source data are provided as a Source Data file.

been fully characterised to date, in particular a very active silencer located between positions 8 and 13. Interestingly, deletions affecting nt C33 show different effects compared to deletions affecting only its flanking nts, and these differential effects are consistent across a wide range of deletion mutant lengths (blue square in the lower middle part of panel 5B). Mutation of C33 has been reported to cause FAS exon 6 skipping in patients with ALPS²⁸ and our deletion scans further emphasise the particular effects of deletions containing this nucleotide, demarcating a splicing enhancer embedded within two silencers previously associated with PTB-mediated repression³³. Finally, a general observation is that deletions above a certain length tend to induce strong skipping (lower left triangle in panel 5B), defining the exon length below which exon definition is likely failing (40–50 nts^{37,38}; but see below).

Collectively, our results suggest that systematic deletion mutagenesis is a particularly informative experimental design to rapidly identify splicing regulatory elements throughout an exon.

Deep learning variant effect predictors accurately predict the effects of different types of genetic variation on alternative splicing

Multiple computational models have been used to predict the effects of genetic variation on splicing. Our deep indel mutagenesis dataset provides a unique opportunity to test the performance of these

models for indel mutations, as well as independent evaluation of their accuracy for predicting the effects of substitutions.

We evaluated the performance of five different models: SMS score (an additive model using 7-mer sequences as input features with parameters learnt in a saturation mutagenesis assay⁷), HAL (an additive model using hexamers as input features with parameters learned from millions of random 50-nt-long exonic and intronic sequences⁵¹), MMSplice (a modular neural network where modules were trained to predict the effects of mutations on different splicing-relevant sequence regions⁵²), SpliceAI (a deep learning model trained on human sequencing data⁵³) and Pangolin (a deep learning model based on the SpliceAI architecture but also trained with data from three additional mammalian species⁵⁴).

All models predicted the effects of single-nt substitutions at least moderately well (Fig. 5A), with Pangolin showing the best performance ($\rho = 0.82$), followed by SpliceAI ($\rho = 0.79$), MMSplice ($\rho = 0.74$), HAL ($\rho = 0.69$) and SMS scores ($\rho = 0.50$). These models followed a similar range and order of performance when predicting the effects of all insertions (1-, 2-, 3- and 4-nts long) in our library ($\rho = 0.80$ for Pangolin, $\rho = 0.78$ for SpliceAI, $\rho = 0.66$ for MMSplice, $\rho = 0.67$ for HAL and $\rho = 0.54$ for SMS scores).

Interestingly, the predictive performance for deletions was substantially worse for most methods ($\rho = 0.18$ for MMSplice, $\rho = -0.41$ for HAL, and $\rho = 0.17$ for SMS scores) apart from

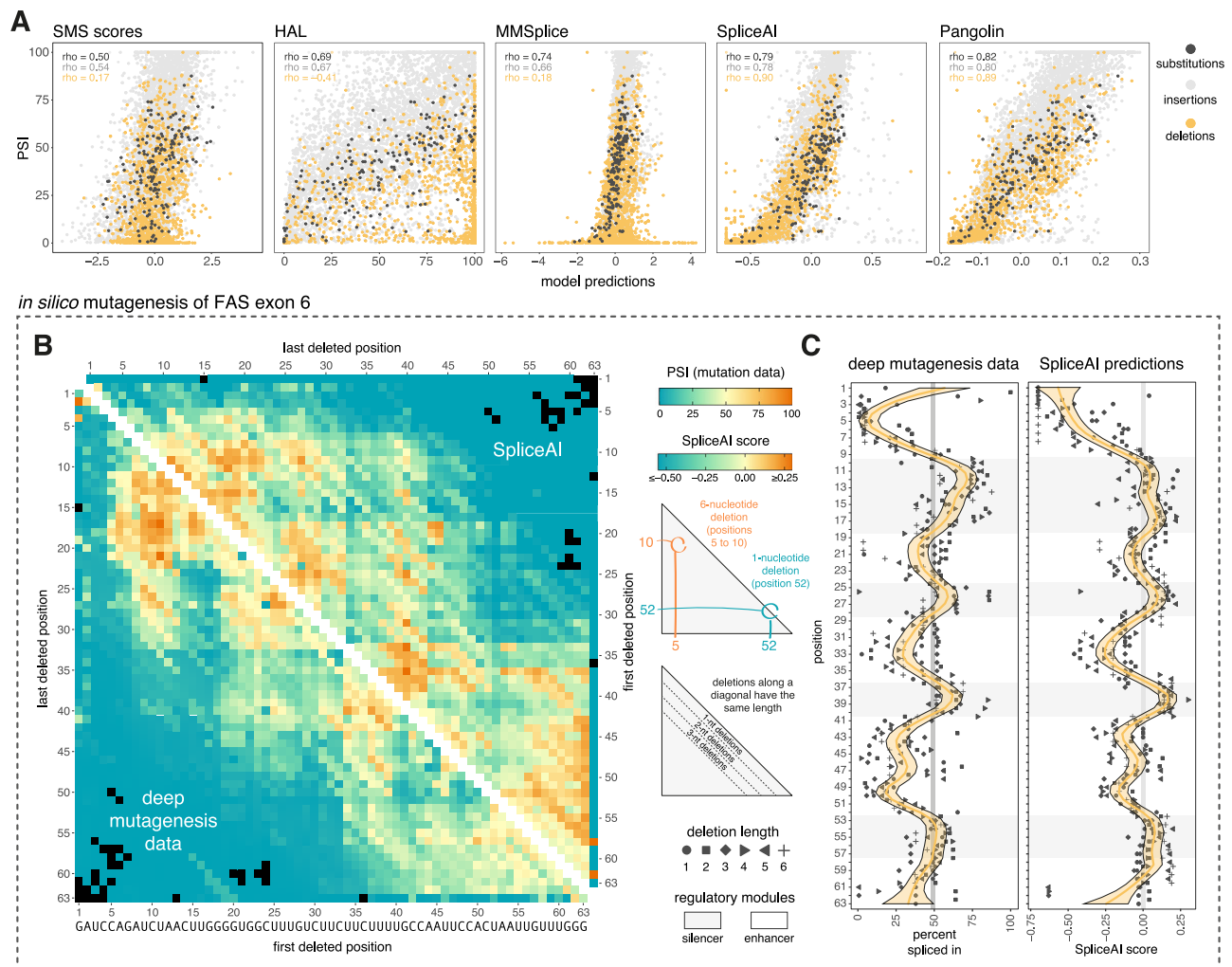


Fig. 5 | Deep learning predicts the inclusion of variants in our deep indel mutagenesis library. **A** Correlations between the inclusion levels of all variants in our library and their inclusion levels according to five different predictors (Substitutions $n = 189$, Deletion $n = 1985$, Insertion $n = 5744$). **B** Lower triangle: Heatmap displaying inclusion levels of all deletion variants in our mutant library. Upper

triangle: Heatmap showing inclusion levels of all deletion variants in our mutant library as predicted by SpliceAI. **C** Left: Inclusion levels of all 1–6 nt long deletion variants along the sequence of FAS exon 6. The yellow line represents a LOESS fit with a 95% confidence band. Right: Inclusion levels of the same variants as predicted by SpliceAI. Source data are provided as a Source Data file.

SpliceAI ($\rho = 0.90$) and Pangolin ($\rho = 0.89$) which remained highly predictive.

Considering all variants of all types in the library, SpliceAI performed best, with a Spearman correlation of 0.84. Indeed, SpliceAI predictions closely replicated our comprehensive deletion maps (Fig. 5B, C), recapitulating a similar regulatory architecture as that uncovered by our experimental dataset. This opens the possibility to use this type of model to perform *in silico* deletion mutagenesis experiments and build regulatory maps for other exons as well.

In silico deletion mutagenesis reveals the regulatory architecture of exons genome-wide

To test the utility of *in silico* deletion mutagenesis, we used SpliceAI to predict the effects of all 4-nt deletions in 18,551 exons expressed in at least 80% of GTEx tissues with a length between 50 and 200 nts. The impact of a 4-nt deletion is predicted to depend on each individual exon, although highly-included exons ($\text{PSI} > 90\%$) are predicted to be more robust to PSI changes compared to exons included at lower levels (Fig. 6A), an expected consequence of the scaling law that the effects of splicing mutations follow⁵⁵.

To systematically analyse the regulatory architecture of exons throughout the genome, we used SpliceAI predictions to train a hidden

Markov model with 3 states (Fig. 6B): E (enhancer – corresponding to regions of the exon that promote skipping upon deletion), S (silencer – regions that promote inclusion upon deletion) and N (neutral – which have no consistent effect upon deletion). Our model captured the regulatory architecture of FAS exon 6 as uncovered by our deep indel mutagenesis experiment: regions of the exon corresponding to inferred enhancers were predicted to be in state E, and regions corresponding to inferred silencers in state S (Fig. 6C, matching experimental results of 4-mer deletions summarised in the heatmap below). This suggests that the model can accurately detect splicing regulatory elements along an exon sequence.

We first used our model to study the distribution of splicing regulatory element (SRE) lengths (i.e., stretches of nts in the same E or S states within an exon) throughout the collection of 18,551 exons. This revealed that most exonic SREs are short, with a median length of 5 nts and a mean length of 8.57 nts (compatible with the average binding site of various RNA-binding protein domains^{56–58}) and similar to what we find in FAS exon 6. Enhancers were predicted to be slightly shorter than silencers (median lengths = 4 vs 6 nts; mean lengths = 6.15 vs 10.70 nts, respectively; Fig. 6D). The model also correctly interpreted exonic sequences that are part of splice site sequences and their immediate neighbourhood as inclusion-promoting (i.e., belonging to

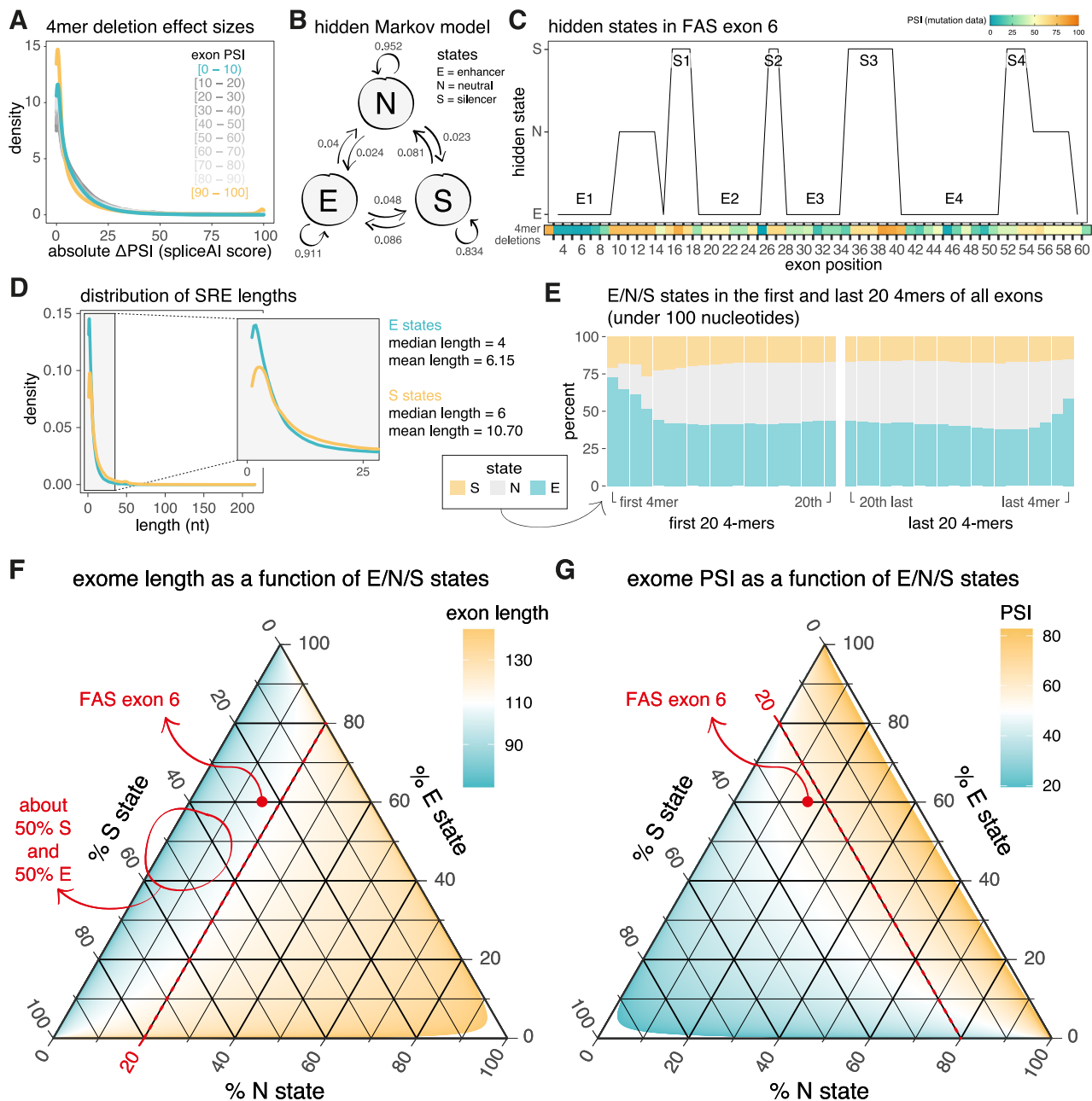


Fig. 6 | Architecture of regulatory elements in alternatively spliced exons across the transcriptome. **A** Distribution of absolute effect sizes of all 4mer deletions in the exome, as predicted by SpliceAI and split by exon PSI groups. **B** Hidden Markov model (HMM) with three states (enhancer, silencer, neutral) used to model the splicing regulatory architecture of exons across the genome. **C** Predicted regulatory architecture of FAS exon 6 based on HMM (above), compared with experimental results of 4-mer deletions (below; deletion effects are centred around the position that the deletions share in common, e.g., deletions 1–4, 2–5, 3–6, 4–7 all cover position 4, therefore they are centred around position 4). **D** Distribution of exonic splicing enhancer and silencer lengths across the exome,

as predicted by our HMM. **E** Distribution of the three states of our model in the first and last 20 4mers of all exons under 100 nucleotides long. **F** Ternary plot illustrating exon length as a function of the relative composition of E/N/S states along the sequences of all 18,551 exons in our dataset. To help to visualise trends in the data, a local regression model (LOESS) was fitted over the raw data (shown in Fig. S15). **G** Ternary plot illustrating exon PSI as a function of the relative composition of E/N/S states along the sequences of all 18,551 exons in our dataset. To help to visualise trends in the data, a local regression model (LOESS) was fitted over the raw data (shown in Supplementary Fig. 15). Source data are provided as a Source Data file.

the E state, Fig. 6E for exons below 100 nt, Supplementary Fig. 14 for longer exons).

We next used our model to gain a comprehensive overview of the regulatory architecture of the entire exome under analysis. To do this, we used ternary plots to visualise all 18,551 exons in our dataset based on their predicted E/N/S states. This revealed two strong trends. First, for exons shorter than 100 nts, the percentage of nts in the *N* state tends to be below 20%, irrespective of the proportion of nts in the *S* or

E states, while the percentage of nts in the *N* state increases with exon length (Fig. 6F, Supplementary Figs. 15A and 16). The percentage of nts in the *N* state tends to be higher, however, in highly (90–100%) included exons, likely due to stronger flanking splice sites characteristic of highly included/constitutive exons, in line with previous reports^{59–61}. This suggests that the inclusion of short alternative exons, such as FAS exon 6, may require a higher density of SREs compared to longer exons. Indeed, the percentage of nts in the *E* state tends to

decrease as exon length increases (Fig. 6F and Supplementary Fig. 16). Second, across highly-included exons, consistently fewer than 20% of nts are in the S state, regardless of the proportion of nts in the N or E states (Fig. 6G, Supplementary Figs. 15B and 16). More generally, the majority of exons examined show fewer than 25% of nts in the S state, and levels of inclusion tend to inversely correlate with this percentage (Supplementary Fig. 16). These observations suggest that maintaining exon inclusion relies more on a low proportion of silencers than on a high proportion of enhancers. Interestingly, FAS exon 6, which has intermediate inclusion levels, approximately aligns with this boundary, with 23% of its nts predicted to be in the S state.

Our deletion analysis revealed not only that nearly the entire sequence of FAS exon 6 is covered by SREs (as is apparently typical of short exons across the genome), but also that its enhancers and silencers alternate in a ‘checkerboard’ pattern along the exon. We used our genome-wide *in silico* deletion mutagenesis to evaluate if this ‘checkerboard’ pattern is likely to be common in additional exons

We first hypothesised that exons encompassed entirely by SREs arranged in an alternating checkerboard pattern would distribute approximately 50% of their nts in the S state and the remaining 50% in the E state. Our ternary plots suggest that exons meeting this criterion tend to be shorter than 100 nts (Fig. 6F and Supplementary Fig. 15A) and display relatively low levels of inclusion (Fig. 6G and Supplementary Fig. 15B). We next explicitly evaluated this hypothesis by counting the number of times the sequence of each exon in our dataset transitions from the E to the S state directly (i.e., without passing through the N state), and vice versa. For example, in the case of FAS exon 6, we counted 7 such transitions (Fig. 6C), equivalent to 9.5 E/S state transitions per 100 nts of exon sequence. The sequences of short (≤ 100 nt long) highly included ($\geq 90\%$) exons were predicted to have very few E/S state transitions, with an average of 1.5 E/S state transitions per 100 nts (Supplementary Fig. 17A, results for longer exons shown in Supplementary Fig. 17B, C). In contrast, short alternatively spliced exons included at lower levels ($< 90\%$) had many more state changes (two-tailed Wilcoxon rank sum test $p < 2.2e-16$), with an average of 3.7 E/S state transitions per 100 nts (Supplementary Fig. 17A, results for longer exons shown in Supplementary Fig. 17B, C).

Repeating this analysis for E/N state transitions (i.e., where the sequence transitions from the E to the N state and vice versa, without passing through the S state) reveals that short highly-included exons have significantly more E/N state transitions per 100 nts compared to short exons with a PSI below 90% (median 2.4 vs 1.4, two-tailed Wilcoxon rank sum test $p < 2.2e-16$, Supplementary Fig. 18A), in agreement with previous findings that constitutive exons are sustained by strong enhancers^{62,63}. Interestingly, this result did not hold true for exons longer than 100 nts (median 1.7 vs 1.8, Wilcoxon rank sum test p value 0.57, Supplementary Fig. 18B, C).

Alternative exons are therefore predicted to have a high density of splicing regulatory elements, which suggests that their precise inclusion levels are tightly regulated and therefore sensitive to mutation. The alternating pattern of enhancers and silencers further suggests that some of these regulatory domains likely act by modulating the function of a neighbouring domain (e.g., a silencer protein binding to its site might sterically prevent a neighbouring enhancer from being bound by an enhancer protein). On the other hand, the lower density of enhancer-silencer alternations in constitutive exons suggests that their high inclusion levels have not been achieved by fine-tuning the binding levels of the splicing regulatory machinery. Their high inclusion levels might therefore simply be a function of their stronger splice sites^{59–61}.

Deletion scans accurately predict the effects of antisense oligonucleotides on exon inclusion

Antisense oligonucleotides (AONs) are an increasingly appealing therapeutic strategy to clinically modulate alternative splicing, as

recently illustrated by the clinical success of Nusinersen for the treatment of Spinal Muscular Atrophy¹⁶ and Eteplirsen for Duchenne muscular dystrophy¹⁷. AONs base-pair to splice sites or regulatory sequences, competing with the binding of splicing factors and regulators⁶⁴.

Clinically used AONs are typically longer than individual regulatory elements (18–21 nt AONs versus 5–10 nt regulatory motifs), making the prediction of AON effects challenging. We reasoned that deletion mutagenesis might provide a rapid method to predict the effects of AONs binding to different regions of a transcript, since deleting (a set of) regulatory motif(s) will inhibit the assembly of cognate *trans*-acting regulatory factors, which is also the mechanism of action of AONs, as they typically compete with the binding of *trans*-acting factors to the same sequences (Fig. 7A).

We compared the effects on splicing of an array of partially overlapping AONs collectively covering the entire length of FAS exon 6 (AON walk) with the effects of deletions of the same length (deletion walk). Changes in exon inclusion correlated well for 21-nt AONs and 21-nt deletions spaced every 5nt along the exon (Spearman $\rho = 0.75$, $n = 9$, Fig. 7B). These AONs modulate exon 6 inclusion over a wide dynamic range (from 20% to 50% inclusion, compared with the approximately 50% inclusion level of the WT exon). The correlation between the AON effects and the SpliceAI-predicted effects of 21-nt deletions (Dango score, see next section) was similarly strong ($\rho = 0.74$, Fig. 7C), suggesting that *in silico* deletion mutagenesis could be an efficient, affordable and high resolution strategy to identify regions in an exon that can be targeted by AONs to achieve a range of desired splicing outcomes for therapeutic or biotechnological applications. Additional analyses using previously published AON walks on two other exons^{55,66} confirm that SpliceAI-based predictions of the effects of exonic deletions can be used to identify AONs that efficiently modulate splicing of alternative exons, although the extent to which it provides a comprehensive picture of regulatory architectures depends on the particular exon analyzed (Supplementary Fig. 19).

DANGO: a genome-wide resource for AON discovery

We used SpliceAI to predict the effects of all possible 21-nt deletions across the exome to generate a resource we refer to as DANGO (Deletion/ANtisense oliGO), where each DANGO score corresponds to the SpliceAI predictions for a particular 21-nt deletion. 12.4% of all 21-nt deletions had an absolute DANGO score greater than 0.1 (mean absolute DANGO score across the exome = 0.05). Short exons (≤ 100 nt) were most vulnerable to these deletions (Fig. 7D), with 44.8% of all 21nt deletions in these exons having an absolute DANGO score > 0.1 , compared to 10.2% in longer exons. Varying the DANGO score threshold (0.025, 0.05, 0.2, or 0.4) altered the proportion of deletions considered impactful, but it did not change the fundamental finding that short exons are more sensitive to the effects of 21nt deletions.

Regardless of the length of the exon, the proportion of negative DANGO scores is greater than the proportion of positive scores (Fig. 7E). This suggests that AONs targeting exonic regions are more likely to reduce recognition of the exon, rather than increase inclusion.

To visualise our results on an exon-by-exon basis, we generated a custom genome browser track (Supplementary Data 1) displaying the DANGO scores for all exons in the genome. This track allows users to interactively explore exonic regions of interest for sequences that may be susceptible to splicing changes upon 21-nt deletions. Visualising FAS exon 6 at single nucleotide resolution reveals that DANGO scores cluster around the identified regulatory domains of this exon (Fig. 7F), suggesting that these scores can accurately reflect the regulatory architecture of exonic sequences. Interestingly, since 21-nt deletions push FAS exon 6 below the length threshold for exon definition (Fig. 3A), nearly all 21mer deletions in this exon are predicted to promote skipping (Fig. 7F). Indeed, these deletions promote skipping as demonstrated in our experimental assay (Fig. 5B), however deletions

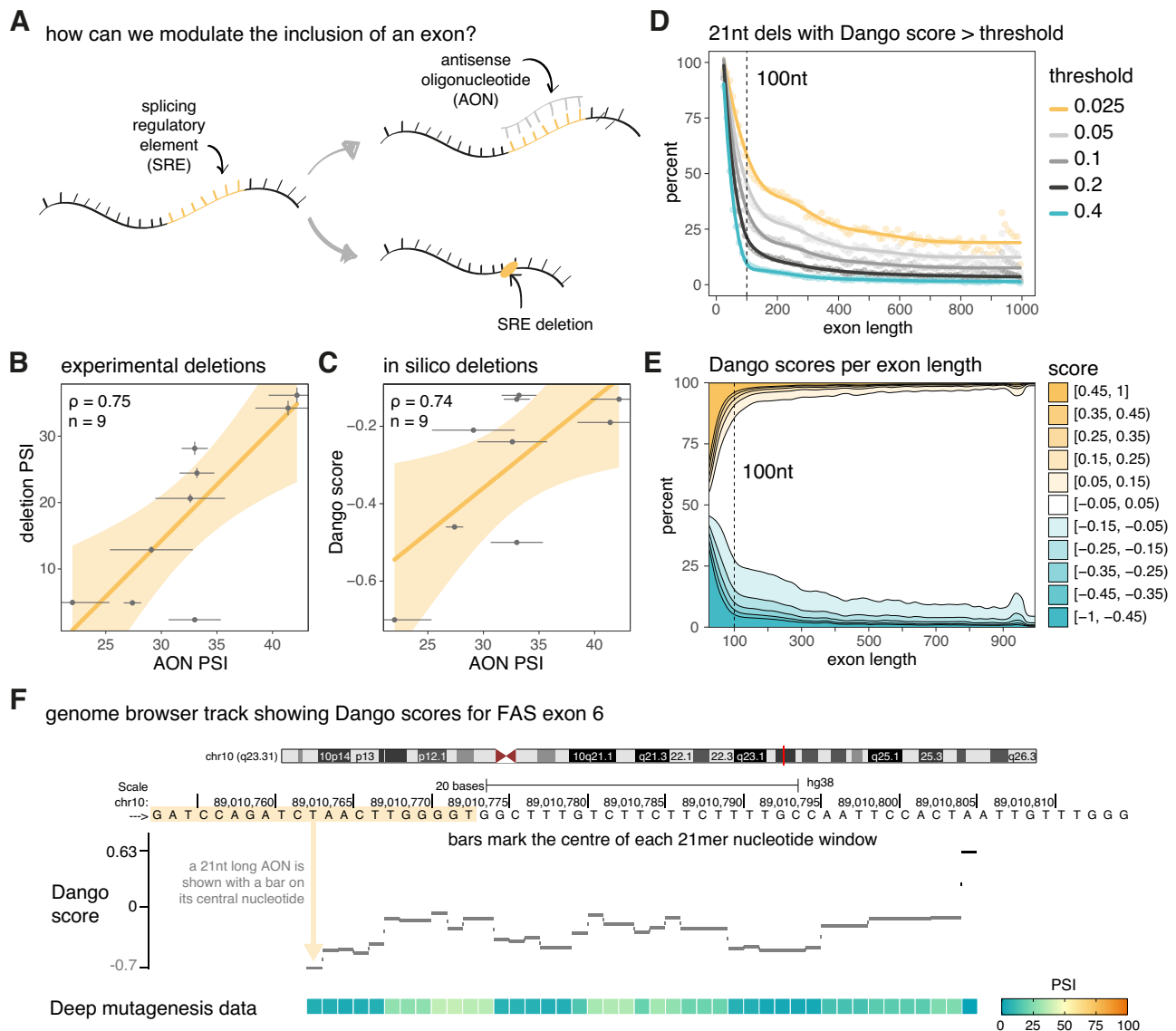


Fig. 7 | Design of splicing-modulating antisense oligonucleotides. **A** The activity of a splicing regulatory element (SRE) could be modulated by using an antisense oligonucleotide (AON) to base pair with this region (therefore sterically blocking any proteins that might bind to the SRE) or alternatively by deleting the SRE altogether. **B** Correlation between the PSI values of nine FAS exon 6 variants with 21-nt deletions and the PSI values of WT FAS exon 6 with a 21-nt AON base pairing to the corresponding regions. Horizontal error bars represent the standard deviation of three replicates. Vertical error bars represent the standard error of the mean in our deep insertion mutagenesis library. Error bars = 95% confidence interval for the mean. **C** Correlation between the Dango scores of the same nine FAS exon 6

variants as in **B**, and the PSI values of the WT FAS exon 6 with a 21-nt AON base pairing to the corresponding regions. Horizontal error bars represent the standard deviation of three replicates. Error bars = 95% confidence interval for the mean. **D** Percentage of 21-nt deletions with Dango scores above the indicated thresholds as a function of exon length. **E** Distribution of Dango scores as a function of exon length. **F** Custom genome browser track displaying the Dango scores for FAS exon 6. The corresponding PSI values as measured experimentally in our deep mutagenesis assay are shown in the heatmap below, using the same colour scale as Fig. 5B. Source data are provided as a Source Data file.

spanning silencer regions promote less skipping than deletions spanning enhancer regions.

DANGO is therefore a genome-wide resource that can help to predict the effects of 21-nt deletions and inform AON selection.

Discussion

Here, by comprehensively quantifying the effects of substitutions, insertions and deletions on splicing, we have not only shown that experimental and in silico deletion scanning are particularly effective strategies for revealing the splicing regulatory landscape of exons, but we have also provided a first overview of the exonic splicing regulatory landscape across the human genome. Moreover, we have shown that deletion scanning accurately predicts the

effects of AONs on exon inclusion and have consequently provided a resource, DANGO, to facilitate AON design genome-wide. Finally, by analysing the effects of insertions, we have discovered a novel regulatory mechanism whereby the inclusion of short exons encoding one-pass transmembrane domains is repressed by the binding of a core component of the spliceosome that normally recognises introns, U2 snRNP, to exons. This recognition of cryptic intron-like sequences in exons provides a simple mechanism for the evolutionary birth of short microexons. Indeed, given the special cis-acting sequences and trans-acting factor requirements for microexon regulation, we were surprised to find that an *accidentally* generated microexon was included and even regulated by SRRM4⁶⁷. This may reflect a general activity for SRRM4 to overcome steric

constraints for the simultaneous engagement of U1 and U2 snRNPs on short exons.

The FAS exon 6 regulatory landscape has a checkerboard organisation of alternating silencers and enhancers that cover the entire exon. *In silico* deletion scans indicate that this architecture is not unusual, but typical of short alternatively spliced exons. The dense splicing regulatory landscape of human exons (Fig. 6) is remarkable and a valuable resource to be further explored in future work. A first insight is that while splicing enhancer elements contribute to maintain some level of inclusion of alternative exons, the density of the (generally less frequent) splicing silencer elements appears to be more decisive in determining the levels of exon skipping.

Collectively, our results highlight the power of deep indel mutagenesis for charting regulatory landscapes, generating novel mechanistic hypotheses, and deriving biological insights. One shortcoming of our study is that we have only considered exonic sequences. In future work, it will be important to extend deletion scanning to introns, and to quantify the effects of inhibiting combinations of intronic and exonic regulatory sequences, which would potentially enable finer control of desired splicing changes with AONs. The good correlation between *in silico* and experimental deletion scans, as well as their high resolution and relatively low cost compared to affordable experimental AON walks, make them interesting tools to prioritise AONs to be tested for eventual clinical or biotechnological applications. Our benchmarking efforts indicate, however, that the predictive power of *in silico* deletion scans varies among different alternatively spliced events: while predicted AONs generally display effects on splicing, *in silico* deletion scans failed to predict the effects of some AONs targeting certain transcript regions or splicing events (Fig. 7C and Supplementary Fig. 19). Nevertheless, we envision that both *in silico* and high-throughput experimental deletion scans will play an increasingly important role in accelerating the discovery of AONs to effectively modulate splicing for many different therapeutic goals.

Methods

Indel library construction

A sequence library was designed to include: 63-nt-long wild-type sequence of FAS exon 6, all possible 189 single-nt substitutions, all 2010 possible deletions ranging in length from 1 to 60 nts, all 5208 possible 1-, 2-, and 3-nt-long insertions as well as 400 randomly-selected 4-nt insertions (full library design along with measured PSI values available in Supplementary Data 1). The library was synthesised and purified by Twist Bioscience.

Indel library amplification

Oligo libraries were resuspended in 10 mM Tris buffer, pH 8.0 to a concentration of 20 ng/ul. 20 ng of template ssDNA Twist library were PCR amplified with Pfx Accuprime polymerase (Thermo scientific, 12344024) in a total reaction volume of 50 ul in quadruplicate, for 12 cycles (as recommended by Twist Bioscience for a 100–150 nt oligo pool) using the following flanking intronic primers: FAS_i5_GC_F (5'-tgtccaatgtccaacctacag-3') and FAS_i6_GC_R (5'-ctacttccaagtatttcaatcg-3'). PCR reactions were combined and cleaned-up with the Qiaquick PCR purification kit, eluted with 50 ul elution buffer and dsDNA measured with a NanoDrop spectrophotometer.

Indel library subcloning

The amplified library was recombined with pCMV FAS wt minigene exon 5-6⁷⁴⁵. We used a vector:insert ratio of 1:8, using 150 ng of vector backbone and 20 ng of dsDNA amplified libraries and incubated at 50 °C for 2 h for DNA assembly, using a Gibson master mix developed at the CRG Protein Technologies Unit, which contains a mix of T5 exonuclease (T5E4111K 1000U from Epicentre Biotech-Ecogen), Phusion polymerase (F530s 100U from VITRO and a Taq DNA ligase (Protein Technologies Unit CRG, homemade). After transformation

into Stellar competent cells (Clontech, 636766), combining five replicates in order to maximise the number of individual transformants amplified, cells were grown for 18 h in LB medium containing ampicillin. We obtained approximately 4.29 million clones. After bacterial transformation, the final plasmid library was purified using the Qiagen plasmid maxi kit (50912163, Qiagen) and quantified with a NanoDrop spectrophotometer.

Transfection of Hek293 cell line to generate output libraries

750,000 Hek293 cells (purchased from ATCC, Cat. # CRL-1573, authenticated by SPR) were plated on 100 × 20 mm petri dishes and transfected with 80 ng cloned libraries in 8 ml OPTIMEM Reduced Serum Medium with no phenol Red (Life technologies, 11058021) using Lipofectamine 2000 (Life technologies, 11668019) in nine biological replicates. 48 h post transfection, cells were collected and RNA was prepared using Maxwell simplyRNA Tissue Kit (Promega, AS1280). cDNA was prepared with 400 ng total RNA using specific vector backbone PT2 primer (5'-AAGCTTGCATCGAATCAGTAG -3') and Superscript III reverse transcriptase (Thermo Fisher, 18080085). PCR amplification of cDNA samples was performed with GoTaq flexi (PROMEGA, M7806) using a distinct 8-mer barcoded oligos to distinguish the nine experimental replicates. PCR products were run on a 2% agarose gel and the smear corresponding to sizes of the amplification product expected from exon inclusion (full length and insertion and deletion mutants) was excised, purified using the Qiaquick Gel extraction kit (Qiagen, 50928704) and quantified with a NanoDrop spectrophotometer.

Input indel library

20 ng of the plasmid library were amplified in triplicates using GoTaq flexi DNA polymerase (M7806, Promega) for 25 cycles with three different pairs of barcoded intronic primers FAS_i5_TR_F and PT2 ("Indel library amplification primers" in Supplementary Table 1). Since the insertions and deletions result in a library with exons of different length, this resulted in a PCR smear (corresponding to exons ranging in length from 3 to 67 nts), which was gel-purified and sequenced. Each pair of primers had a distinct 8-mer barcode sequence to discriminate between technical replicates ("Primers used for amplifying technical replicates (input library)" in Supplementary Table 1).

Sequencing

Equimolar quantities of three independent amplifications of the input library and equimolar quantities of the purified inclusion smear (output library) of each of the nine replicates were pooled and sequenced at the CRG Genomics Core Facility, where Illumina Ampliseq PCR-free libraries were prepared and run on a single lane of an Illumina HiSeq2500. In total, 424 million paired-end reads were obtained (188 and 236 million for input and output, respectively). The median sequencing coverage for all exon variants in the input was 2114 reads. In the output, the median sequencing coverage was between 278 and 468 reads. Raw sequencing data has been submitted to GEO with accession number GSE244179.

Data processing and calculation of PSI values

FastQ files from paired-end sequencing were processed with DiMSum v1.2.7⁶⁸ using default settings with minor adjustments (<https://github.com/lehner-lab/DiMSum>). First, DiMSum was run in default paired-end mode to demultiplex reads into input and replicate output samples (Stage 0 only). Second, DiMSum Stages 1–5 were run in single-end mode ('--paired' = F) using only demultiplexed forward reads that have full coverage of the exon sequence. Reverse reads, originally intended to cover a unique molecular identifier (UMI) and a 3' portion of the exon sequence, were discarded. The final stage estimates an enrichment score (ES) and associated error for each mutant variant based on its frequency in the input and output libraries, and relative to the wild

type sequence in both libraries. Experimental design files and command-line options required for running DiMSum on this dataset are available on GitHub (<https://github.com/lehner-lab/fas-indel-library>).

The PSI of the wild type FAS exon 6 sequence has been previously experimentally shown to be 49.1%⁵. Therefore, the PSI of a variant of interest is estimated as follows:

$$\frac{PSI_{variant}}{PSI_{WT}} = \frac{\exp(ES_{variant})}{\exp(ES_{WT})} \quad (1)$$

$$PSI_{variant} = PSI_{WT} \cdot \frac{\exp(ES_{variant})}{\exp(ES_{WT})} \quad (2)$$

$$PSI_{variant} = 49.1 \cdot \frac{\exp(ES_{variant})}{\exp(ES_{WT})} \quad (3)$$

Because ES_{wt} is 0, this is equivalent to:

$$PSI_{variant} = 49.1 \cdot \exp(ES_{variant}) \quad (4)$$

Experimental validation of estimated PSI values

To confirm the accuracy of our PSI estimates on single-nt substitutions, we took the previously experimentally-determined values of 25 exon variants⁵ and plotted them against the experimentally-determined values (Fig. 1C). To validate our estimates of indel PSI values, 15 individual clones from the indel library were Sanger sequenced. They covered a wide range of estimated PSI and were therefore good for validation and checking correlation.

Specific individual mutants were transfected into Hek293 cells in triplicates to quantify the ratio between exon 6 inclusion and skipping. For RT-PCR, minigene-specific primers were used (“Primers used for amplifying biological replicates (Output library)” in Supplementary Table 1). To avoid amplification of endogenous FAS RNAs, these primers (PT1 and PT2) are complementary to a plasmid backbone sequence distinct from endogenous DNA. RT-PCR products were fractionated by electrophoresis using 6% polyacrylamide gels in 1 x TBE and Sybr safe staining (ThermoFisher Scientific, S33102). The bands corresponding to exon inclusion or skipping were quantified using ImageJ v1.47 (NIH, USA). PSI measurements are shown in Supplementary Table 2.

Under the particular experimental conditions in which these indel mutants were tested, the wild-type exon was included with a PSI of 49.9% (compared to 49.1% in the experiment done to validate the single nt substitutions⁵). To visualise these results in the same plot as the single-nt substitutions (Fig. 1C), we used the splicing scaling law⁵⁵ to adjust all experimentally-determined PSI values to what their values should have been if the wild type had a PSI of 49.1%.

Estimating PSI values in the GTEx dataset

We estimated the PSI of exons in the GTEx dataset (GTEx Consortium, 2017) from the proportion of reads supporting exon inclusion in the GTEx junction read counts file (GTEx_Analysis_2016-01-15_v7_STAR-v2.4.2a_junctions.gct.gz; available for download at <https://www.gtexportal.org/home/datasets>). To do this, we used the *quantifySplicing* function from the Psychomics package in R⁶⁹. The *minReads* argument was set to 10 (such that a splicing event requires at least 10 reads for it to be quantified) and the *eventType* argument was set to ‘SE’ (instructing the *quantifySplicing* function to quantify alternative exon events). All estimates were based on the Psychomics hg19/GRCh37 alternative splicing annotations.

Experimentally validating microexon inclusion

We initially cloned the microexon sequences observed in the output (i.e., those sequences corresponding to large deletions, Supplementary Fig. 11A) into our plasmid vector backbone (pCMV_FAS_exon4_exon6) and transfected them. No inclusion band was found in the polyacrylamide gels (i.e., these sequences were 100% skipped, Supplementary Fig. 11B).

Since the nt composition of the PTB binding domain in the central region of FAS exon 6 is very similar to the polypyrimidine tract of a 3' splice site, we reasoned that an AG-containing insertion right after nt 40 (e.g., CAG, AG, TA, A, CTA) could create a new 3' splice site in this region of the exon. Such a splice site would produce the microexons detected in the deep mutagenesis experiment. We introduced these insertions (as well as the non-AG-containing GAC insertion as a negative control) into the vector backbone using site-directed mutagenesis (Agilent, 200523) using the relevant mutagenesis primers. These minigene constructs were transfected into Hek293 cells in triplicate. RT-PCR products were fractionated by electrophoresis on 6% native acrylamide gels.

Experimentally testing exon inclusion with different splice sites

To test the effects of mutations in the presence of different 3' splice site strengths, we used partially complementary oligonucleotides in combination with TaqPlus precision (Agilent, 600212) and PCR “around the world” (primers pointing in opposite directions from the mutagenesis site to amplify the full length of the plasmid) to replace the naturally weak 3' splice site of FAS exon 6 (5'UUUCAUAUAAAAU-GUCCAAUGUCCAACCUACAG3') with a strong 3' splice site sequence (5'UACUACGGCUUUUUUUUCCUUUUUCAG3').

PTB/SRRM4 overexpression experiments

We overexpressed the SRRM4 (or PTBP1) protein by co-transfecting minigenes containing a CAG or AG insertion after position 40 along with 1000 ng of pcDNA5_SRRM4_flag (or pcDNA5_PTBP1_T7) in lipofectamine 2000 for 24 h. RT-PCR was then used to analyse the splicing ratios as described above.

In vitro transcription and spliceosome assembly experiments

T7 promoter-containing transcription templates were generated by PCR using Gotaq flexi enzyme (Promega, M7806): FAS exon 6 WT/noBP templates were generated from FAS WT minigene, Fas M0 template was generated from Fas M0 minigene³³, FAS exon 6 UAG mutant template was generated from a ssDNA oligonucleotide. PCR products were purified on agarose gel.

Cy5-CTP/Cy5-UTP labelled RNA were transcribed directly from the PCR templates using Megascript T7 Transcription kit (Ambion) according to the manufacturer's instructions. A complex formation 15 ng/ul fluorescently labelled RNA were incubated with 3 ul of HeLa cell nuclear extracts (CILBIOTECH, CC-01-20-50) supplemented with 3 mM MgCl₂, 24.9 mM KCl, 3.33% PVA, 13.3 mM HEPES pH 8, 0.13 mM EDTA, 13.3 % glycerol, 0.03 % NP-40, 0.66 mM DTT, 2 mM ATP and 22 mM creatine phosphate in a final volume of 9 ul. The mixture was incubated for 18 min at 30 °C. 1 microliter of heparin (10 ug/ul stock) was added and incubated for 10 min at room temperature. 3 ul of 50% glycerol were added and 9 ul loaded on a composite gel (4% acrylamide, 0.05% bis-acrylamide, 0.5% agarose, 50 mM Tris, 50 mM glycine).

The gel was run for 6 h at 200 Volts in 50 mM Tris / 50 mM glycine buffer. After electrophoresis, fluorescence was detected using a Typhoon PhosphorImager. The inactivation of U1 snRNP and U2 snRNP was performed as described in Dönmez et al.⁷⁰ using 2'-O-methylated oligoribonucleotide complementary to U1 snRNA (5'-CUGCCAG-GUAGUAU-3') or U2 snRNA (5'-CAGAUACUACACUUG-3').

Scanning the exome for sequences similar to 3' splice sites

To identify exonic sequences similar to 3' splice sites in our GTEx dataset (see Estimating PSI values in the GTEx dataset section), we used

SVM-BPfinder, a support vector machine that scores how closely a nt sequence resembles a 3' splice site preceded by a branchpoint⁴⁸. SVM-BPfinder was run at each position of each exon with the *--species* argument set to Hsap, the *--max-len* argument set to 1000, and the *--min-dist* argument set to 15. The final score assigned to each exon was the maximum score across all of its positions. This corresponds to the sequence within the exon that most closely resembles a 3' splice site.

Measuring hydropathy in exons throughout the genome

To measure hydropathy, the *hydrophobicity* function from the *Peptides* package⁷¹ in R was used with the *scale* argument set to “KylDoolittle”.

Preparation of minigene constructs carrying a transmembrane domain (CHODL exon 5 and CXADR exon 6)

Genomic DNA sequences were amplified from commercial genomic DNA (PROMEGA, G304A), and branch site mutations were produced with the help of Taq Plus precision (Agilent, 600212). Sequences of genomic regions were cloned into the pCMV_FAS567 minigene replacing FAS exon 6. The amplified genomic sequences were:

CHODL exon 5 (103 bp) GRCh37/hg19 chr21:19635108-19635210
GTATAATCCCAATCTAATTTATGTTGTTATACCAACAA-
TACCCTGCTCTTACTGA-
TACTGGTTGCTTTTGAACCTGTTGTTCCAGATGCTGCATAAAAG

CXADR exon 6 (139 bp) GRCh37/hg19 chr21:18933656-18933794
CTTCAAATAAAGCTGGACTAATTGCAGGAGCCATTA-
TAGGAACTTTGCTTGTCTAGCGCTCATTGGTCTTAT-
CATCTTTTGTGCTGCGTAAAAAGCGCAGAGAAGAAAAA-
TATGAAAAGGAAGTTCATCACGATATCAG

Minigene constructs with the wild type exons contained the amplified genomic sequences above. In the case of exons whose putative branchpoint adenines were mutated and substituted with guanines, the minigene constructs contained the following exonic sequences:

CHODL exon 5

GTATAATCCCAATCTAATTTATGTTGTTATACCAACAA-
TACCCTGCTCTTGCTGGTGCTGGTTGCTTTTGAACCTGTTGTTTC-
CAGATGCTGCATAAAAG

CXADR exon 6

CTTCAAATAAAGCTGGACTAATTGCAGGAGCCATTA-
TAGGAACTTTGCTTGTCT-
TAGCGCTCGTTGGTCTTGCTCGTCTTTTGTGCTGCGTAAAAAGCGCAGA-
GAAGAAAAATATGAAAAGGAAGTTCATCACGATATCAG

Sequences were all confirmed by Sanger sequencing of acrylamide purified PCR bands by crush and soak method. All minigene constructs were confirmed by Sanger sequencing.

Predicting FAS exon 6 mutation effects using SMS Scores

We downloaded Supplementary Table 7 from Ke et al.⁷, which lists the SMS scores for all possible 7-mers. To calculate the total SMS score for each exon in our indel library, we performed a sliding window analysis by adding the SMS scores of consecutive 7-mers along its sequence. The final SMS score for each variant was obtained by subtracting the total SMS score for the wild type exon from that of the variant:

$$SMS_{final(variant)} = SMS_{total(variant)} - SMS_{total(wild\ type)} \quad (5)$$

Predicting FAS exon 6 mutation effects using HAL

To predict the effects of exon variants on inclusion with HAL⁵¹, we uploaded a file containing the sequences in our library to <http://splicing.cs.washington.edu/SE> using 49.1% as the wild type levels of

inclusion. An output file was returned that contains the predicted PSI values for each sequence in the input file.

Predicting FAS exon 6 mutation effects using MMSplice

We converted our indel library design file to VCF format, and used this new file as input for MMSplice⁵². We ran the algorithm online, on the Google Colab notebook provided for this purpose (available at <https://colab.research.google.com/drive/1Kw5rHMxaxXXsmE3WecxbXyGQJma80Eq6>). This returned a CSV file with multiple columns containing different metrics for each exon variant. We selected *delta_logit_psi* as the predictor for the mutation effects.

Predicting FAS exon 6 mutation effects using SpliceAI

We converted our indel library design file to VCF format, and used this new file as input for SpliceAI⁵³. We ran SpliceAI using GRCh38 as both the genome reference and gene annotation files. All other parameters were set to the default configuration (parameter “D”, the maximum distance between the variant and gained/lost splice site was left to its default of 50, which means that the algorithm could be capturing information about splice site gain and loss beyond the boundaries of the exon). As SpliceAI outputs four scores for each mutant sequence (corresponding to splice site acceptor loss, splice site acceptor gain, splice site donor loss, and splice site donor gain), we selected the score associated with the highest absolute value in each case. If this score corresponded to a splice site loss, the score was multiplied by -1 .

Predicting FAS exon 6 mutation effects using Pangolin

We converted our indel library design file to VCF format, and used this new file as input for the Google Colab Notebook made available by the authors of Pangolin⁵⁴ at <https://colab.research.google.com/github/tkzeng/Pangolin/blob/main/PangolinColab.ipynb>. Pangolin was used with the default options chosen for the Colab Notebook, including GRCh37 as the genome reference. Like SpliceAI, Pangolin outputs four scores for each mutant sequence (corresponding to splice site acceptor loss, splice site acceptor gain, splice site donor loss, and splice site donor gain). The Pangolin score selected for each mutation corresponded to that with the highest absolute value out of these four. If this score corresponded to a splice site loss, the score was multiplied by -1 .

In silico 4mer deletions in exons genome-wide

The SpliceAI developers created a file with annotations for all possible substitutions, 1 base insertions, and 1–4 base deletions across the genome. This file is available for download at <https://basespace.illumina.com/s/otSPW8hnhZR>. We downloaded the file and extracted 4mer deletion data for all exons we calculated PSI values for (see Estimating PSI values in the GTeX dataset section). For each 4mer deletion in these exons, we computed its SpliceAI score by taking the maximum value among the acceptor gain, acceptor loss, donor gain, and donor loss scores. If the maximum value was the acceptor loss or the donor loss score, we multiplied the value by -1 .

Hidden Markov model

We used the *depmixS4* package in R to build a hidden Markov model that predicts the locations of exonic splicing enhancers (which promote exon inclusion) and silencers (which promote skipping) in each exon based on SpliceAI scores for 4mer deletions (see section above titled in silico 4mer deletions in exons genome-wide). Exons with lengths between 50 and 200 nts were used as input for the model, with each exon as a separate time series during training. The model has three hidden states: E (enhancer), N (neutral), and S (silencer), with mean scores of -0.15 , 0 , and 0.15 , respectively. The standard deviations for the states were fixed at 0.1 , 0.025 , and 0.1 to account for the variability of positive and negative values in the dataset. For example, although

breaking a splice site has a much stronger effect than breaking a weak enhancer (resulting in a much more negative spliceAI score), both sequence elements should be classified together in the E state.

AON walk

100,000 HEK293 cells in a 6-well plate were transfected with Antisense oligonucleotide harbouring 2'-O Me phosphorothioate modifications at each nucleotide position (Integrated technologies) using 3 μ l of Lipofectamine 2000 (11668027, ThermoFisher Scientific) in one ml OPTIMEM I Reduced Serum Medium with no phenol red (11058021, ThermoFisher Scientific) to a final concentration of 2.5 nMolar (exact sequences shown in Supplementary Table 3). Six hours post-transfection, the cell culture medium was replaced with DMEM Glutamax (61965059, ThermoFisher Scientific) containing 10% FBS and Pen/Strep antibiotics. 24 h post-transfection, total RNA was isolated using the automated Maxwell LEV 16 simplyRNA tissue kit (AS1280, Promega). cDNA was synthesised with 400 ng total RNA using Superscript III (18080085, Life Technologies) with a mix of random primers and oligodT. Effects on endogenous FAS exon 6 inclusion were determined by PCR using GoTaq flexi DNA polymerase (M7806, Promega) and the following primers:

FAS_e5_for 5'-TGTGAACATGGAATCATCAAGG-3'
FAS_e7_endo_R 5'-AAAGTTGGAGATTCATGAGAACC-3'

Exome-wide 21mer deletion scan

After validating the ability of SpliceAI to predict AON effects, we characterised the AON targetability across the genome by performing an exon-wide scan of deletions using SpliceAI scores as a proxy for targetability. Specifically, for all exons in the genome, we produced SpliceAI scores for each length-21 deletion within each of the exons' boundaries. To compile the list and sequences of each exon, we used R package *biomaRt*⁷² with the *hsapiens_gene_ensembl* dataset. We limited the analyses to canonical exons only (referring to the *transcript_is_canonical* attribute of each exon). After compiling the list of all exons and each of their length-21 deletions, we passed each sequence to SpliceAI using the parameters as described under Predicting FAS exon 6 mutation effects using SpliceAI. The resulting scores are referred to as the "DANGO scores" of each exon.

Statistical tests

All statistical tests were performed in R 3.6.2 using custom code (see Code availability section).

Reporting summary

Further information on research design is available in the Nature Portfolio Reporting Summary linked to this article.

Data availability

Dango scores have been uploaded to the GitHub repository. DNA sequencing data are available in NCBI's Gene Expression Omnibus under the GEO Series accession number [GSE244179](https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE244179). All reagents used in this study, including deep mutagenesis libraries are available upon request from the corresponding authors. Source data are provided with this paper.

Code availability

All scripts used in this study have been made available at the following GitHub repository: <https://github.com/lehner-lab/fas-indel-library> (<https://doi.org/10.5281/zenodo.15705337>).

References

- Braunschweig, U., Guerousov, S., Plocik, A. M., Graveley, B. R. & Blencowe, B. J. Dynamic integration of splicing within gene regulatory pathways. *Cell* **152**, 1252–1269 (2013).
- Anna, A. & Monika, G. Splicing mutations in human genetic disorders: examples, detection, and confirmation. *J. Appl. Genet.* **59**, 253–268 (2018).
- Scotti, M. M. & Swanson, M. S. RNA mis-splicing in disease. *Nat. Rev. Genet.* **17**, 19–32 (2016).
- Sterne-Weiler, T. & Sanford, J. R. Exon identity crisis: disease-causing mutations that disrupt the splicing code. *Genome Biol.* **15**, 201 (2014).
- Julien, P., Miñana, B., Baeza-Centurion, P., Valcárcel, J. & Lehner, B. The complete local genotype-phenotype landscape for the alternative splicing of a human exon. *Nat. Commun.* **7**, 11558 (2016).
- Braun, S. et al. Decoding a cancer-relevant splicing decision in the RON proto-oncogene using high-throughput mutagenesis. *Nat. Commun.* **9**, 1–18 (2018).
- Ke, S. et al. Saturation mutagenesis reveals manifold determinants of exon definition. *Genome Res.* **28**, 11–24 (2018).
- Baeza-Centurion, P., Miñana, B., Valcárcel, J. & Lehner, B. Mutations primarily alter the inclusion of alternatively spliced exons. *Elife* **9**, e59959 (2020).
- Soemedi, R. et al. Pathogenic variants that alter protein code often disrupt splicing. *Nat. Genet.* **49**, 848–855 (2017).
- Cheung, R. 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).
- Montgomery, S. B. et al. The origin, evolution, and functional impact of short insertion-deletion variants identified in 179 human genomes. *Genome Res.* **23**, 749–761 (2013).
- Mullaney, J. M., Mills, R. E., Pittard, W. S. & Devine, S. E. Small insertions and deletions (INDELs) in human genomes. *Hum. Mol. Genet.* **19**, R131–R136 (2010).
- Stenson, P. D. et al. The Human Gene Mutation Database: building a comprehensive mutation repository for clinical and molecular genetics, diagnostic testing and personalized genomic medicine. *Hum. Genet.* **133**, 1–9 (2014).
- Zhang, X. et al. Impact of human pathogenic micro-insertions and micro-deletions on post-transcriptional regulation. *Hum. Mol. Genet.* **23**, 3024–3034 (2014).
- Rogalska, M. E., Vivori, C. & Valcárcel, J. Regulation of pre-mRNA splicing: roles in physiology and disease, and therapeutic prospects. *Nat. Rev. Genet.* **24**, 251–269 (2023).
- Qiu, J. et al. History of development of the life-saving drug 'Nusinersen' in spinal muscular atrophy. *Front. Cell. Neurosci.* **16**, 942976 (2022).
- Sheikh, O. & Yokota, T. Pharmacology and toxicology of eteplirsen and SRP-5051 for DMD exon 51 skipping: an update. *Arch. Toxicol.* **96**, 1–9 (2022).
- Kim, J. et al. A framework for individualized splice-switching oligonucleotide therapy. *Nature* **619**, 828–836 (2023).
- Popplewell, L. J., Trollet, C., Dickson, G. & Graham, I. R. Design of phosphorodiamidate morpholino oligomers (PMOs) for the induction of exon skipping of the human DMD gene. *Mol. Ther.* **17**, 554–561 (2009).
- Harding, P. L., Fall, A. M., Honeyman, K., Fletcher, S. & Wilton, S. D. The influence of antisense oligonucleotide length on dystrophin exon skipping. *Mol. Ther.* **15**, 157–166 (2007).
- Pramono, Z. A. D. et al. A prospective study in the rational design of efficient antisense oligonucleotides for exon skipping in the DMD gene. *Hum. Gene Ther.* **23**, 781–790 (2012).
- Wee, K. B. et al. Dynamics of co-transcriptional pre-mRNA folding influences the induction of dystrophin exon skipping by antisense oligonucleotides. *PLoS One* **3**, e1844 (2008).
- Aartsma-Rus, A., Houlleberghs, H., van Deutekom, J. C. T., van Ommen, G.-J. B. & 't Hoen, P. A. C. Exonic sequences provide better targets for antisense oligonucleotides than splice site sequences in

- the modulation of Duchenne muscular dystrophy splicing. *Oligonucleotides* **20**, 69–77 (2010).
24. Papoff, G. et al. An N-terminal domain shared by Fas/Apo-1 (CD95) soluble variants prevents cell death in vitro. *J. Immunol.* **156**, 4622–4630 (1996).
 25. Cascino, I., Fiucci, G., Papoff, G. & Ruberti, G. Three functional soluble forms of the human apoptosis-inducing Fas molecule are produced by alternative splicing. *J. Immunol.* **154**, 2706–2713 (1995).
 26. Liu, C., Cheng, J. & Mountz, J. D. Differential expression of human Fas mRNA species upon peripheral blood mononuclear cell activation. *Biochem. J.* **310**, 957–963 (1995).
 27. Cheng, J. et al. Protection from Fas-mediated apoptosis by a soluble form of the Fas molecule. *Science* **263**, 1759–1762 (1994).
 28. Agrebi, N. et al. Rare splicing defects of FAS underly severe recessive autoimmune lymphoproliferative syndrome. *Clin. Immunol.* **183**, 17–23 (2017).
 29. Ustianenko, D., Weyn-Vanhenhenryck, S. M. & Zhang, C. Micro-exons: discovery, regulation, and function. *Wiley Interdiscip. Rev. RNA* **8**, 10 (2017).
 30. Liu, H. X., Zhang, M. & Krainer, A. R. Identification of functional exonic splicing enhancer motifs recognized by individual SR proteins. *Genes Dev.* **12**, 1998–2012 (1998).
 31. Jain, N., Lin, H.-C., Morgan, C. E., Harris, M. E. & Tolbert, B. S. Rules of RNA specificity of hnRNP A1 revealed by global and quantitative analysis of its affinity distribution. *Proc. Natl. Acad. Sci. USA* **114**, 2206–2211 (2017).
 32. Yu, Y. et al. Dynamic regulation of alternative splicing by silencers that modulate 5' splice site competition. *Cell* **135**, 1224–1236 (2008).
 33. Izquierdo, J. M. et al. Regulation of Fas alternative splicing by antagonistic effects of TIA-1 and PTB on exon definition. *Mol. Cell* **19**, 475–484 (2005).
 34. Vo, T. et al. HNRNPH1 destabilizes the G-quadruplex structures formed by G-rich RNA sequences that regulate the alternative splicing of an oncogenic fusion transcript. *Nucleic Acids Res.* **50**, 6474–6496 (2022).
 35. Paronetto, M. P. et al. Regulation of FAS exon definition and apoptosis by the Ewing sarcoma protein. *Cell Rep.* **7**, 1211–1226 (2014).
 36. Choi, N. et al. SRSF6 regulates the alternative splicing of the apoptotic Fas gene by targeting a novel RNA sequence. *Cancers* **14**, 1990 (2022).
 37. Berget, S. M. Exon recognition in vertebrate splicing. *J. Biol. Chem.* **270**, 2411–2414 (1995).
 38. Ule, J. & Blencowe, B. J. Alternative splicing regulatory networks: functions, mechanisms, and evolution. *Mol. Cell* **76**, 329–345 (2019).
 39. Dominski, Z. & Kole, R. Selection of splice sites in pre-mRNAs with short internal exons. *Mol. Cell. Biol.* **11**, 6075–6083 (1991).
 40. Irimia, M. et al. A highly conserved program of neuronal microexons is misregulated in autistic brains. *Cell* **159**, 1511–1523 (2014).
 41. Juan-Mateu, J. et al. Pancreatic microexons regulate islet function and glucose homeostasis. *Nat. Metab.* **5**, 219–236 (2023).
 42. Gao, K., Masuda, A., Matsuura, T. & Ohno, K. Human branch point consensus sequence is γ UnAy. *Nucleic Acids Res.* **36**, 2257–2267 (2008).
 43. Singh, R., Valcárcel, J. & Green, M. R. Distinct binding specificities and functions of higher eukaryotic polypyrimidine tract-binding proteins. *Science* **268**, 1173–1176 (1995).
 44. Gonatopoulos-Pournatzis, T. et al. Genome-wide CRISPR-Cas9 interrogation of splicing networks reveals a mechanism for recognition of autism-misregulated neuronal microexons. *Mol. Cell* **72**, 510–524.e12 (2018).
 45. Förch, P. et al. The apoptosis-promoting factor TIA-1 is a regulator of alternative pre-mRNA splicing. *Mol. Cell* **6**, 1089–1098 (2000).
 46. Izquierdo, J. M. & Valcárcel, J. Fas-activated serine/threonine kinase (FAST K) synergizes with TIA-1/TIAR proteins to regulate Fas alternative splicing. *J. Biol. Chem.* **282**, 1539–1543 (2007).
 47. Schott, G. et al. U2AF2 binds IL7R exon 6 ectopically and represses its inclusion. *RNA* **27**, 571–583 (2021).
 48. Corvelo, A., Hallegger, M., Smith, C. W. J. & Eyras, E. Genome-wide association between branch point properties and alternative splicing. *PLoS Comput. Biol.* **6**, e1001016 (2010).
 49. Weng, L. et al. A novel alternative spliced chondrolectin isoform lacking the transmembrane domain is expressed during T cell maturation. *J. Biol. Chem.* **278**, 19164–19170 (2003).
 50. Excoffon, K. J. D. A., Bowers, J. R. & Sharma, P. 1. Alternative splicing of viral receptors: a review of the diverse morphologies and physiologies of adenoviral receptors. *Recent Res. Dev. Virol.* **9**, 1–24 (2014).
 51. Rosenberg, A. B., Patwardhan, R. P., Shendure, J. & Seelig, G. Learning the sequence determinants of alternative splicing from millions of random sequences. *Cell* **163**, 698–711 (2015).
 52. Cheng, J. et al. MMSplice: modular modeling improves the predictions of genetic variant effects on splicing. *Genome Biol.* **20**, 48 (2019).
 53. Jaganathan, K. et al. Predicting splicing from primary sequence with deep learning. *Cell* **176**, 535–548.e24 (2019).
 54. Zeng, T. & Li, Y. I. Predicting RNA splicing from DNA sequence using Pangolin. *Genome Biol.* **23**, 103 (2022).
 55. Baeza-Centurion, P., Miñana, B., Schmiedel, J. M., Valcárcel, J. & Lehner, B. Combinatorial genetics reveals a scaling law for the effects of mutations on splicing. *Cell* **176**, 549–563.e23 (2019).
 56. Hall, K. B. RNA-protein interactions. *Curr. Opin. Struct. Biol.* **12**, 283–288 (2002).
 57. Ray, D. et al. A compendium of RNA-binding motifs for decoding gene regulation. *Nature* **499**, 172–177 (2013).
 58. Pan, X., Rijnbeek, P., Yan, J. & Shen, H.-B. Prediction of RNA-protein sequence and structure binding preferences using deep convolutional and recurrent neural networks. *BMC Genom.* **19**, 511 (2018).
 59. Senapathy, P., Shapiro, M. B. & Harris, N. L. Splice junctions, branch point sites, and exons: sequence statistics, identification, and applications to genome project. *Methods Enzymol.* **183**, 252–278 (1990).
 60. Sorek, R. & Ast, G. Intronic sequences flanking alternatively spliced exons are conserved between human and mouse. *Genome Res.* **13**, 1631–1637 (2003).
 61. Zheng, C. L., Fu, X.-D. & Gribskov, M. Characteristics and regulatory elements defining constitutive splicing and different modes of alternative splicing in human and mouse. *RNA* **11**, 1777–1787 (2005).
 62. Schaal, T. D. & Maniatis, T. Multiple distinct splicing enhancers in the protein-coding sequences of a constitutively spliced pre-mRNA. *Mol. Cell. Biol.* **19**, 261–273 (1999).
 63. Fairbrother, W. G., Yeh, R.-F., Sharp, P. A. & Burge, C. B. Predictive identification of exonic splicing enhancers in human genes. *Science* **297**, 1007–1013 (2002).
 64. Havens, M. A. & Hastings, M. L. Splice-switching antisense oligonucleotides as therapeutic drugs. *Nucleic Acids Res.* **44**, 6549–6563 (2016).
 65. Hua, Y., Vickers, T. A., Baker, B. F., Bennett, C. F. & Krainer, A. R. Enhancement of SMN2 exon 7 inclusion by antisense oligonucleotides targeting the exon. *PLoS Biol.* **5**, e73 (2007).
 66. Mogilevsky, M. et al. Modulation of MKNK2 alternative splicing by splice-switching oligonucleotides as a novel approach for glioblastoma treatment. *Nucleic Acids Res.* **46**, 11396–11404 (2018).
 67. Gonatopoulos-Pournatzis, T. & Blencowe, B. J. Microexons: at the nexus of nervous system development, behaviour and autism spectrum disorder. *Curr. Opin. Genet. Dev.* **65**, 22–33 (2020).
 68. Faure, A. J., Schmiedel, J. M., Baeza-Centurion, P. & Lehner, B. DiMSum: an error model and pipeline for analyzing deep mutational

- scanning data and diagnosing common experimental pathologies. *Genome Biol.* **21**, 207 (2020).
69. Saraiva-Agostinho, N. & Barbosa-Morais, N. L. Psychomics: graphical application for alternative splicing quantification and analysis. *Nucleic Acids Res.* **47**, e7 (2019).
70. Dönmez, G., Hartmuth, K., Kastner, B., Will, C. L. & Lührmann, R. The 5' end of U2 snRNA is in close proximity to U1 and functional sites of the pre-mRNA in early spliceosomal complexes. *Mol. Cell* **25**, 399–411 (2007).
71. Osorio, D., Rondon-Villarreal, P. & Torres, R. Peptides: a package for data mining of antimicrobial peptides. *R. J.* **7**, 4–14 (2015).
72. Durinck, S., Spellman, P. T., Birney, E. & Huber, W. Mapping identifiers for the integration of genomic datasets with the R/Bioconductor package biomaRt. *Nat. Protoc.* **4**, 1184–1191 (2009).

Acknowledgements

The authors thank Manuel Irimia for insightful comments and Maxim Mogilevsky and Rotem Karni for providing primary data from their publication. This project has received funding from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (grant agreements 670146 and 883742) and from the European Union's Horizon Europe under the grant agreement No 101071936. We also acknowledge support of the Spanish Ministry of Science and Innovation through the Centro de Excelencia Severo Ochoa (CEX2020-001049-S, MCIN/AEI/10.13039/501100011033), and the Generalitat de Catalunya through the CERCA programme. We are grateful to the CRG Core Technologies Programme for their support and assistance in this work. We received funding from the Spanish State Research Agency (PID2020-114630GB-I00/AEI/10.13039/501100011033), LCF/PR/HR21/52410004, EMBL Partnership, the Bettencourt Schueller Foundation, the AXA Research Fund, and Agencia de Gestio d'Ajuts Universitaris i de Recerca (AGAUR, 2017 SGR 1322). GQ was supported by PRE2022-102744, financed by MCIN/AEI/10.13039/501100011033 and FSE + . J.C. is funded by the BBSRC DTP (Bio-technology and Biological Sciences Research Council, Biosciences Doctoral Training Programme, Cambridge, UK. The Genotype-Tissue Expression (GTEx) data used for the analyses described in this manuscript were obtained from the GTEx Portal on May 8, 2018 and dbGaP accession number phs000424.v7.p2 on May 8, 2018. The GTEx Project was supported by the Common Fund of the Office of the Director of the NIH and by NCI, NHGRI, NHLBI, NIDA, NIMH, and NINDS. Funded by the European Union. Views and opinions expressed are, however, those of the author(s) only and do not necessarily reflect those of the European Union. Neither the European Union nor the granting authority can be held responsible for them.

Author contributions

P.B.-C., A.J.F., M.T., G.Q. and J.C. performed computational analyses; B.M. and S.B. performed experiments; P.B.-C., B.L. and J.V. wrote the manuscript with input from all authors.

Competing interests

CRG has filed a patent (European Priority Application 24382126.0) for the use of deep indel mutagenesis as a method to identify and predict the effects of antisense oligonucleotides. P.B.-C., B.M., B.L. and J.V. are listed as co-inventors. J.V. is a member of the Scientific Advisory Boards of Remix Therapeutics, Stoke Therapeutics and IntronX. The remaining authors declare no competing interests.

Additional information

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s41467-025-62957-7>.

Correspondence and requests for materials should be addressed to Ben Lehner or Juan Valcárcel.

Peer review information *Nature Communications* thanks Ulrich Braunschweig, Yitzhak Pilpel and the other, anonymous, reviewer(s) for their contribution to the peer review of this work. A peer review file is available.

Reprints and permissions information is available at <http://www.nature.com/reprints>

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

Open Access This article is licensed under a Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License, which permits any non-commercial use, sharing, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if you modified the licensed material. You do not have permission under this licence to share adapted material derived from this article or parts of it. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by-nc-nd/4.0/>.

© The Author(s) 2025