



## The human genome contains over a million autonomous exons

Nicholas Stepankiw, Ally W H Yang and Timothy R Hughes

*Genome Res.* published online November 9, 2023  
Access the most recent version at doi:[10.1101/gr.277792.123](https://doi.org/10.1101/gr.277792.123)

---

<b>P&lt;P</b>	Published online November 9, 2023 in advance of the print journal.
<b>Accepted Manuscript</b>	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
<b>Open Access</b>	Freely available online through the <i>Genome Research</i> Open Access option.
<b>Creative Commons License</b>	This manuscript is Open Access. This article, published in <i>Genome Research</i> , is available under a Creative Commons License (Attribution 4.0 International license), as described at <a href="http://creativecommons.org/licenses/by/4.0/">http://creativecommons.org/licenses/by/4.0/</a> .
<b>Email Alerting Service</b>	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or <a href="#">click here</a> .

---



---

To subscribe to *Genome Research* go to:  
<https://genome.cshlp.org/subscriptions>

---

Published by Cold Spring Harbor Laboratory Press

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20

# The human genome contains over a million autonomous exons

Nicholas Stepankiw <sup>1</sup>, Ally W.H. Yang <sup>1</sup>, and Timothy R. Hughes <sup>1,2,\*</sup>

- 1. Donnelly Centre, University of Toronto, Toronto ON CANADA
  - 2. Department of Molecular Genetics, University of Toronto, Toronto, ON CANADA
- \* To whom correspondence should be addressed:  
[t.hughes@utoronto.ca](mailto:t.hughes@utoronto.ca); 416-859-1492

21

22 **ABSTRACT**

23 Eukaryotic mRNAs and lncRNA exons are often small compared to introns. The exon  
24 definition model predicts that exons splice autonomously, dependent on proximal exon  
25 sequence features, explaining their delineation within large introns. This model has not  
26 been examined on a genome-wide scale, however, leaving open the question of how  
27 often mRNA and lncRNA exons are autonomous. It is also unknown how frequently  
28 such exons can arise by chance. Here, we directly assayed large fragments (500-1000  
29 bp) of the human genome by exon trapping, which detects exons spliced into a  
30 heterologous transgene, here designed with a large intron context. We define these  
31 exons as “autonomous”. We obtained ~1.25 million exons, including most known mRNA  
32 and well-annotated lncRNA internal exons, demonstrating that human exons are  
33 predominantly autonomous. mRNA exons are trapped with highest efficiency. Nearly a  
34 million of the trapped exons are unannotated, most located in intergenic regions and  
35 antisense to mRNA, with depletion from the forward strand of introns. These exons are  
36 not conserved, indicating they are non-functional and likely arose from random  
37 mutations. They are nonetheless highly enriched with known splicing promoting  
38 sequence features delineating known exons. Novel autonomous exons are more  
39 abundant than annotated lncRNA exons, and computational models also indicate they  
40 will occur with similar frequency in any randomly generated sequence. These results  
41 show that most human coding exons splice autonomously, and provide an explanation  
42 for the existence of many unconserved lncRNAs, as well as a new annotation and  
43 inclusion levels of spliceable loci in the human genome.

44

45 Keywords:

46 Splicing, exons, exon trapping

47

48

49

50

51

52 **INTRODUCTION**

53 Eukaryotic transcripts, including most human mRNAs, are often composed of  
54 alternating exons and introns. In human, and most vertebrates, the introns are generally  
55 much larger than the exons (Consortium 2001). The intron/exon boundaries consist of  
56 relatively short and degenerate 5' and 3' splice site sequences (hereafter, 5' and 3'SS,  
57 respectively), and due to random chance, large introns will contain many sequences  
58 that resemble 5' and 3'SS (Sun and Chasin 2000; Cote et al. 2001; Yeo and Burge  
59 2004). Precise removal of introns is thought to be facilitated mainly by a mechanism  
60 known as exon definition (Robberson et al. 1990; Piovesan et al. 2019), in which the  
61 recognition of adjacent flanking 3' and 5'SS are facilitated by bridging of the splicing  
62 complexes across the exon. The specificity gained by the characteristic 80-220 base  
63 spacing between the 3' and 5'SS is insufficient to precisely specify human exons,  
64 however, as many exons are outside this range. Presumably as a consequence, human  
65 exons are often associated with additional sequences that promote inclusion (known as  
66 splicing enhancers) (Liu et al. 1998; Reed 2000; Cote et al. 2001; Wang et al. 2004b;  
67 Zhang et al. 2005; Ke et al. 2011; Wang et al. 2012). Exon splicing therefore depends  
68 on a variety of sequence features, with the 3'SS and 5'SS being essential. The  
69 sequence features that delineate exons remain incompletely known, however, and as a  
70 result, the exon definition model has not been explicitly confirmed on a genomic scale.  
71 Thus, it remains unknown what proportion of human exons are autonomous – i.e.  
72 containing sequences that are sufficient to enable splicing into a mature transcript.

73 A variety of computational approaches have been taken to predict exon identity and  
74 inclusion level from primary sequence. These algorithms would presumably learn or  
75 incorporate features that are employed by cells to delineate exons, but they also tend to  
76 include additional correlated information that are not relevant to mechanistic  
77 understanding of exon recognition. Thus, they do not explicitly predict exon autonomy,  
78 nor do they reveal the required sequence features. For example, gene-finding programs  
79 perform this task (e.g. GenScan (Burge and Karlin 1997)), but these typically  
80 incorporate coding potential, sequence conservation, and other factors (see (Scalzitti et  
81 al. 2020) for a recent overview)). Much of the literature has focused on predicting  
82 inclusion levels of alternative exons (Cartegni et al. 2003; Fairbrother et al. 2004;  
83 Barash et al. 2010; Rosenberg et al. 2015; Xiong et al. 2015) but these methods  
84 assume exon boundaries are known. Most coding exons are constitutive, in any case  
85 (Pan et al. 2008; Wang et al. 2008); it is not clear that alternative splicing signals would  
86 be the same signals that define the exons. SpliceAI, a well-known predictor of splice  
87 sites, operates directly from primary sequence, using Convolutional Neural Networks  
88 (CNNs) trained on splice site locations within known mRNA genes, thus presumably  
89 incorporating local context (Jaganathan et al. 2019). SpliceAI captures splice sites of  
90 both constitutive and alternative exons, but it does not directly report the identity of full  
91 exons. In addition, as a CNN with ~700,000 parameters, it is inherently challenging to  
92 interpret. Moreover, exon-associated sequence features (e.g. those that would  
93 contribute to protein-coding ability or transcript stability, and not splicing *per se*) may  
94 contribute to computational discrimination of annotated exons vs. other sequences,  
95 without necessarily being mechanistic drivers of splicing itself.

96 A related fundamental question is how many exons (broadly defined as sequences that  
97 can splice into a surrounding transcript) exist in the human genome. There are  
98 ~181,000 annotated internal exons within the ~20,000 known protein-coding genes;  
99 these constitute roughly 1% of the human genome. An even larger fraction may  
100 comprise the enigmatic lncRNAs (long noncoding RNAs), however. In aggregate,  
101 upwards of 800,000 lncRNA exons have been catalogued, of which at least 250,000 are  
102 internal exons (Lagarde et al. 2017; The et al. 2017; Pertea et al. 2018; Volders et al.  
103 2019; Zhao et al. 2021). In contrast, only ~25,000 internal exons annotated as part of a  
104 lncRNA in the curated ENCODE collection (GENCODE v37 (Frankish et al. 2019)).  
105 Many lncRNAs appear to be extremely rare, as they are found at low levels, or in only  
106 one dataset. The vast majority of lncRNAs have no known function (Ponting and Haerty  
107 2022), and many display weaker splicing signals than protein coding genes (Deveson et  
108 al. 2018). It has been proposed that many arise as transcriptional “noise” (Ponting and  
109 Haerty 2022), which could arise as a consequence of transcription from enhancers  
110 (Engreitz et al. 2016). While eRNAs are typically unstable, well-annotated lncRNAs  
111 typically contain multiple exons (Orom et al. 2010), presumably stabilizing the lncRNAs,  
112 as splicing signals are known to enhance both transcription and RNA stability (Le Hir et  
113 al. 2003; Core et al. 2014).

114 The classical exon definition model predicts that exons would be largely self-determined  
115 by local sequence features, but the fact that splicing does not always occur in a strictly  
116 linear fashion along primary transcripts (e.g. (Drexler et al. 2020)) suggests the  
117 possibility of distant interactions and dependency among exons and flanking genomic  
118 features. We therefore sought to survey the human genome to ask which exons are

119 “autonomous” (i.e. self-defined, as they will splice into a heterologous transgene with  
120 relatively large introns). We employed a classical “exon trapping” assay to survey the  
121 human genome for autonomous exons (Duyk et al. 1990) whereby genomic fragments  
122 are assayed outside of their normal contextual setting, e.g. flanking exons, promoter,  
123 transcription level and distal intronic sequences. We reasoned that this survey would  
124 allow us to query whether protein-coding exons are generally autonomous, whether  
125 exons exist elsewhere in the genome, what sequence features they possess, and  
126 whether exons arise at random, which would partly explain the existence of lncRNAs.  
127 The results clarify several aspects of the human exon complement, identify a large  
128 number of previously undocumented exons, and indicate that multi-exon lncRNAs are  
129 an expected feature of large genomes.

130

## 131 **RESULTS**

### 132 **Genome-wide exon trapping**

133 We employed a classical exon trapping assay, in which a query sequence is cloned into  
134 the middle of a 1.6 kb intron, to survey the human genome (see **Fig. 1** for a schematic  
135 overview and example data). If the query sequence contains an autonomous exon, that  
136 exon (or more than one exon) will be included in the resulting spliced transcript. We  
137 employed five vectors that differ in either reading frame relative to the splice sites  
138 (0,+1,+2), removal of a predicted downstream intronic splicing enhancer sequence (a  
139 binding site for RBFOX1 within a hairpin), or a 5'SS mutation that we engineered to  
140 weaken the predicted splice site strength (see **Methods, Supplemental Fig. S1A**).  
141 These variations in the splicing context were included initially to ask whether gross

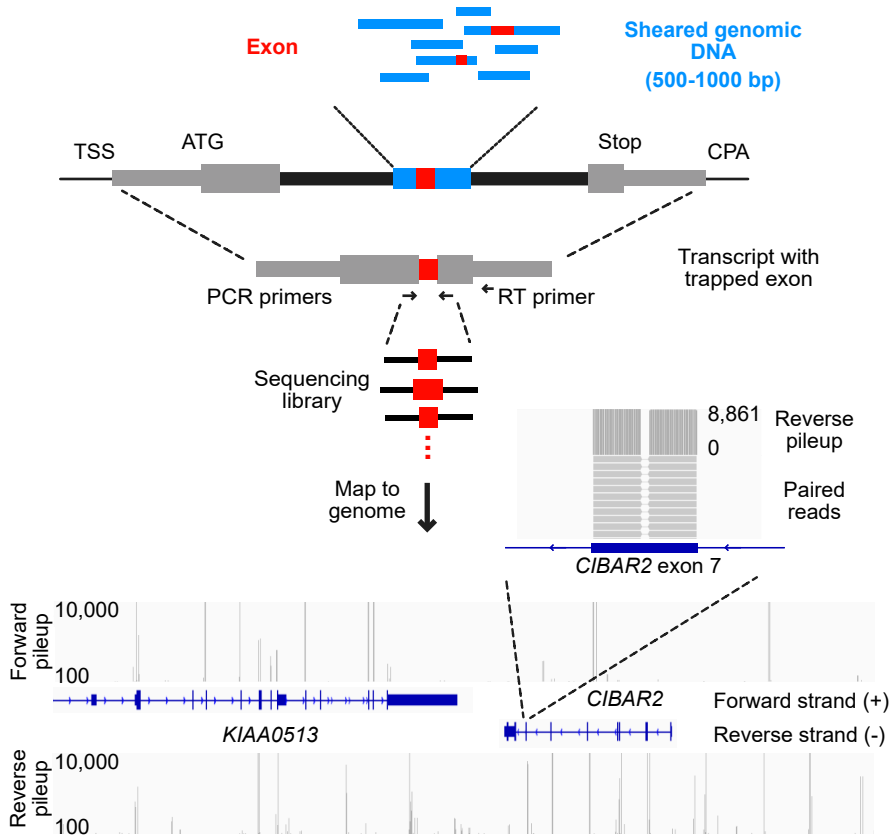


Figure 1

## Figure 1. Overview of genome exon trapping method

Diagram depicting exon trapping approach, sequencing library construction and example sequencing read maps. Sheared genomic DNA library fragments (blue boxes) 500-1000 bp in length are cloned into the middle of the 6<sup>th</sup> intron from *TRA2B* (black boxes), in a pcDNA 3.1 vector backbone. First and terminal exons (grey boxes) are labeled with the transcriptional start site (TSS), start codon (ATG), stop codon (Stop), and the cleavage and polyadenylation site (CPA). Internal exons (red boxes) are amplified by RT-PCR, using indicated primers, then sequenced and mapped to the human genome (hg38). Bottom panel shows mapped sequencing read counts (separated into forward and reverse strand pileups) for regions containing KIAA0513 and a portion of *CIBAR2* (display region coordinates: Chr16:85,062,938-85,134,585). The zoom-in region corresponds to exon 7 of *CIBAR2*.

142 systematic differences would result, but even those specifically sought (i.e. impact of  
143 reading frame) appeared relatively minor. We therefore pooled all of the reads, to  
144 increase numbers and provide redundancy.

145 Query sequences consisted of 23 libraries (4 or 5 libraries for each vector) which were  
146 generated from sheared human genomic DNA fragments (500-1000 bp). Each of the 23  
147 plasmid libraries consisted of 2 to 20 million bacterial clones, and was transfected into  
148 ~2 million HEK293 cells. The number of plasmids per transfected cell was not  
149 assessed, but 2 million unique plasmids would represent roughly a quarter of the  
150 stranded genome. We therefore anticipated up to 1-fold genomic coverage per library,  
151 and up to 5-fold sampling of the genome over all 23 libraries. Following transient  
152 transfection of reporter construct libraries, RNA extraction, and poly(A) selection, we  
153 generated RT-PCR products containing the trapped exon. We then sequenced the  
154 resulting PCR products and mapped the reads to the genome. As with other Massively  
155 Parallel Reporter Assays, we measure the splicing frequency of exons across the  
156 genome using the sequencing read counts, which we take as a proxy for the inclusion  
157 rate of the exon.

158 In total, we obtained ~4.2 billion paired end reads that mapped uniquely to the human  
159 genome (an average of ~200 million reads per library). These reads mapped to ~6  
160 million clusters with identical or nearly identical ends (see **Methods** for details). These  
161 initial clusters encompassed ~9% of the genome. The majority of exon clusters were  
162 supported by 10 or fewer reads, however (**Fig. 2A, Supplemental Fig. S2A**), and many  
163 were found in only one library (**Supplemental Fig. S2B**). Internal exons from protein-

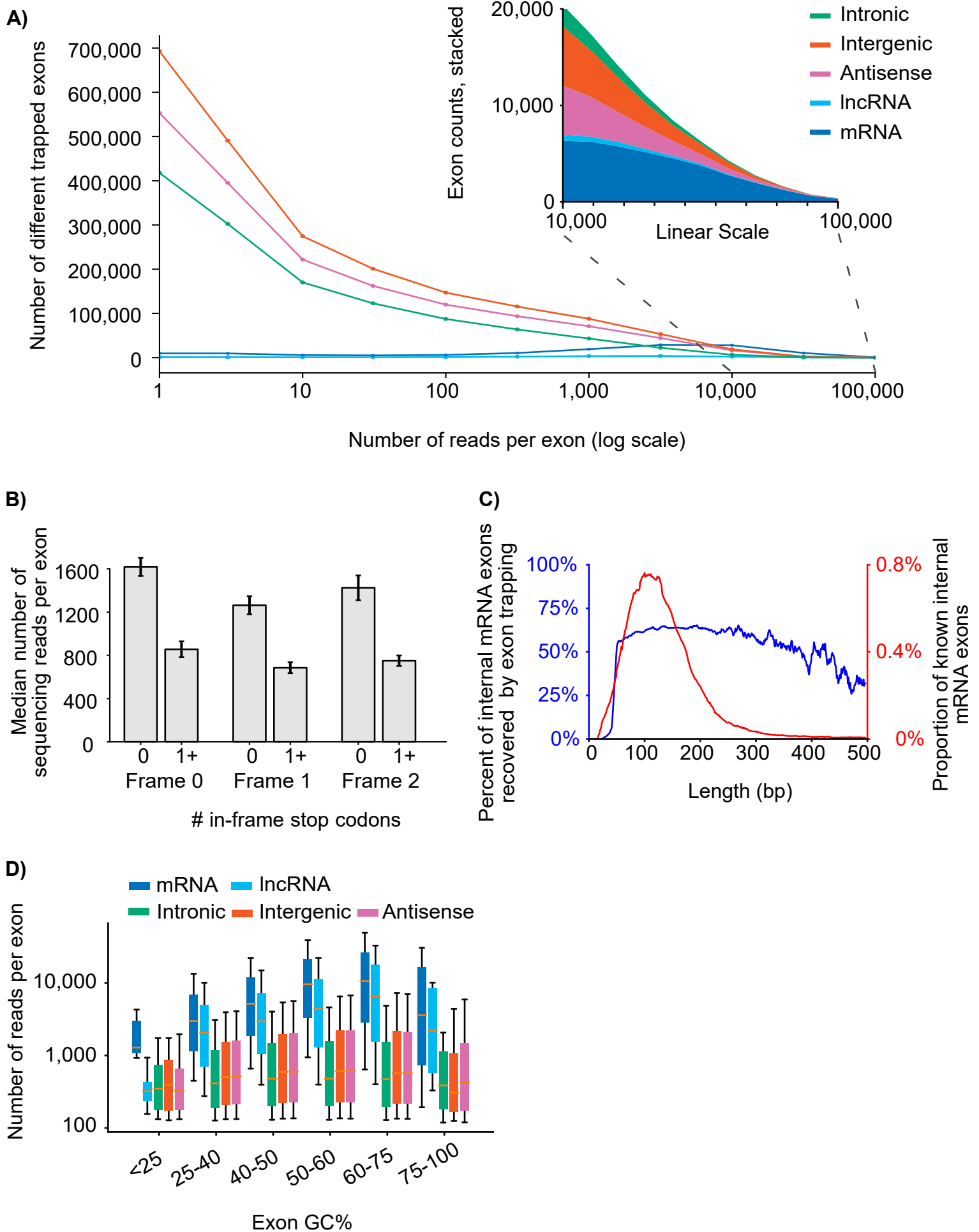


Figure 2

## Figure 2. Properties of trapped exons.

- A) Histograms of sequencing read counts for trapped internal exons within different genomic regions. Outset plot shows logarithmic exon counts and inset shows zoomed linear exon counts. Logarithmic bin boundaries indicated by dots corresponding to  $10^x$  for  $x$  from 0 to 5 with step size 0.5. Linear bins boundaries range from 100 to 100,000 with a step size of 1,000.
- B) Bar plots depicting sequencing read counts of internal exons containing zero or at least one in-frame stop codon. Results for reading frames in the first (“Frame 0”), second (“Frame 1”), and third (“Frame 2”) positions are shown.
- C) Line plots depicting the distribution of mRNA exon lengths and the percent of mRNA exons at each length recovered by exon trapping. Plots have a 9 bp smoothing window applied.
- D) Boxplots depicting GC content of trapped internal exons from different genomic regions. Y-axis indicates sequencing read counts of trapped exons, within indicated GC content ranges (x-axis). Whiskers indicate 10<sup>th</sup> and 90<sup>th</sup> percentiles.

164 coding genes, in contrast, had an average of 10,636 reads, and were typically found in  
165 multiple libraries (**Supplemental Fig. S2C**), together encompassing 32% of all reads.

166 We limited most of our subsequent analyses to exon clusters with at least 100 reads  
167 across all libraries (hereafter, we refer to “exon clusters” as “exons”), with the goal of  
168 cataloguing exons with very high confidence. When we randomized the positions of  
169 reads across the forward strand of a chromosome (Chr17, chosen for its relatively small  
170 size) and identified exons using the same process described above, none exceeded ten  
171 reads. In the original data set, however, 27,715 of the trapped exons exceeded 100  
172 reads. Thus, 100 reads is a very conservative threshold, which we anticipate will be  
173 robust to alternative statistical tests. Across the entire genome in the real exon trapping  
174 data, a threshold of 100 reads captures 1,245,947 exons in total, encompassing 3.2% of  
175 the stranded genome (i.e. 6.2 Gb) (hereafter referred to as “1.25 million exons”). This  
176 figure is much higher than we had initially expected, and is thus the focus of this paper,  
177 beyond this section.

178 We retrospectively examined the detection of the 1.25 million exons in the different  
179 vectors and libraries, in order to estimate coverage. The individual vectors each  
180 captured between 41% and 75% of all sequence present in any vector (i.e. there is not  
181 even a single read for the remaining 59% and 25% of the 1.25 million exons,  
182 respectively), roughly consistent with the 60% breadth that would be expected from 1×  
183 coverage. Overlap between the vectors is consistent with random sampling (i.e. the  
184 intersection is similar to expectation if both vectors sampled from 1.25 million exons)  
185 (**Supplemental Fig. S1B**). The distribution of individual exons across libraries also  
186 appears random, with most of the 1.25 million exons present in multiple libraries

187 **(Supplemental Fig. S2D,E)**. Manual inspection of read counts on genome browser  
188 displays typically show an “all or nothing” pattern, in which an exon is detected in a  
189 library either hundreds of times, or not at all, again consistent with incomplete sampling  
190 in individual libraries. Most of the 1.25 million exons were detected in at least half of the  
191 libraries (examples in **Supplemental Fig. S2F**), showing that they are not an artifact of  
192 a single library or vector. Importantly, these results do not demonstrate coverage of all  
193 possible exons in the genome; below, we describe several exon attributes that are  
194 depleted. They do, however, indicate that the vast majority of exons that would be  
195 detected in this experimental system are present in the dataset. Our estimated 5-fold  
196 coverage would correspond to ~99% breadth; even 3-fold would correspond to over  
197 90%.

198 The five different libraries did vary systematically in the inclusion of individual exons to  
199 some degree, but not as greatly as we had anticipated. The impact of Nonsense  
200 Mediated Decay (Maquat 2004) was clearly observed among the vectors in three  
201 different reading frames, but with only a 2-fold decrease, on average, associated with  
202 stop codons in the trapped exon (**Fig. 2B**). Read counts of identical exons, compared  
203 between two libraries – with a stop codon in one library but not the other, due to reading  
204 frame – also showed a median decrease of just over 2-fold (**Supplemental Fig.**  
205 **S3A,B**). Thus, NMD has a quantitative, but generally not qualitative effect in the splicing  
206 reporter assay. As our main goal in this manuscript is to provide an overall picture of the  
207 unexpectedly high number of exons observed, we pooled the reads from all libraries for  
208 subsequent analyses.

209 We note that the size of the genomic fragments in the libraries is large enough to  
210 accommodate two closely spaced exons, which could be captured in the assay  
211 simultaneously. If the exons are short, such cases should be detectable in single reads.  
212 Indeed, we observed 14,830 trapped exons corresponding to such mRNA “doublets”,  
213 associated with 8.2% of trapped mRNA exons. 9,590 of these, however, are also  
214 associated with respective single exons. Anecdotally, the same appears to be true for  
215 non-mRNA exons among the 1.25 million, but due to uncertainties in mapping from the  
216 ends of reads we did not further examine this phenomenon; extrapolating from mRNA  
217 exons, we expect that roughly 3% of new exons described here may in fact represent  
218 two adjacent exons.

### 219 **Exon inclusion rate varies among types of RNA, and with exon properties**

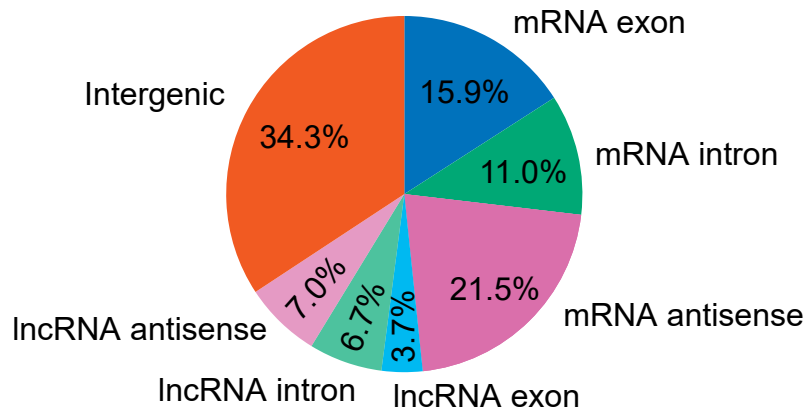
220 We next determined the proportion of known exons captured (from mRNAs and  
221 lncRNAs), and where these trapped exons are found relative to known transcript  
222 structures (Lagarde et al. 2017; The et al. 2017; Pertea et al. 2018; Volders et al. 2019;  
223 Zhao et al. 2021). As noted above, at least a subset of known exons are very well  
224 captured: **Fig. 2A** shows that exon clusters with very high read counts (higher than  
225 20,000) largely correspond to internal exons from protein-coding genes, despite these  
226 exons representing less than 1% of the genome.

227 The cutoff of 100 reads retains the majority of GENCODE mRNA (61%) (and lncRNA  
228 (53%)) internal exons in the trapped exons. Thus, even though only a single cell line  
229 was employed, the assay clearly demonstrates that the majority of human protein-  
230 coding exons are autonomous (78% and 68% are detected at a threshold of one read  
231 for GENCODE mRNA and lncRNA, respectively).

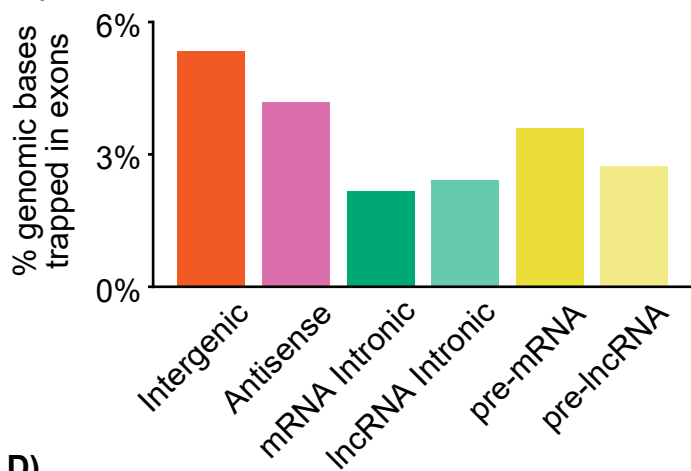
232 We examined what properties of exons may control detection in exon trapping. Internal  
233 exons with high inclusion rates tend to have stronger splice sites; this expected finding  
234 is explored in more detail below. The capture rate of GENCODE mRNA internal exons  
235 also depends on exon length, and is highest between 50 and 250 bases (blue line in  
236 **Fig. 2C**). This range encompasses the ~140 bp that is typical of internal exons, which  
237 presumably facilitates U2/U1 bridging (Robberson et al. 1990) (red line in **Fig. 2C**). The  
238 capture rate also depends strongly on GC content (**Fig. 2D**). We cannot rule out a  
239 technical origin for this phenomenon, but we note that base content has the potential to  
240 impact RBP binding site frequency, and could influence nucleosome occupancy (Tillo  
241 and Hughes 2009). There are many indications that nucleosomes promote splicing  
242 (Hollander et al. 2016), and indeed, exons with 50-75% GC are captured at the highest  
243 rates (**Fig. 2D**). The effect of exon length is also influenced by base content, in that high  
244 GC content is associated with recovery of longer exons (**Supplemental Fig. S3C,D**),  
245 consistent with high GC content helping to overcome lack of U2/U1 bridging.

246 The internal exons of housekeeping genes (Houkpe et al. 2021), genes expressed  
247 highly in HEK293 cells (Nieborak et al. 2023), and the exons of all other coding genes  
248 displayed similar read counts in the exon trapping data (**Supplemental Fig. S4A**). Exon  
249 sequencing read counts are, however, dependent on the PSI measured in HEK293 cells  
250 (Ellis et al. 2023) (**Supplemental Fig. S4B**), and alternative HEK293 mRNA exons are  
251 frequently missing in the exon trapping data (**Supplemental Fig. S4C**). Alternative  
252 exons of all types have lower read counts on average (see below). Altogether,  
253 sequence properties of exons themselves have a strong impact on exon trapping, but  
254 expression level of the corresponding gene has almost none.

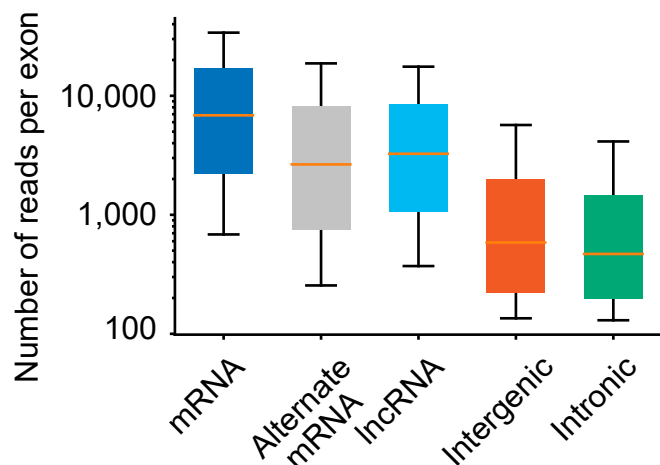
**A)**



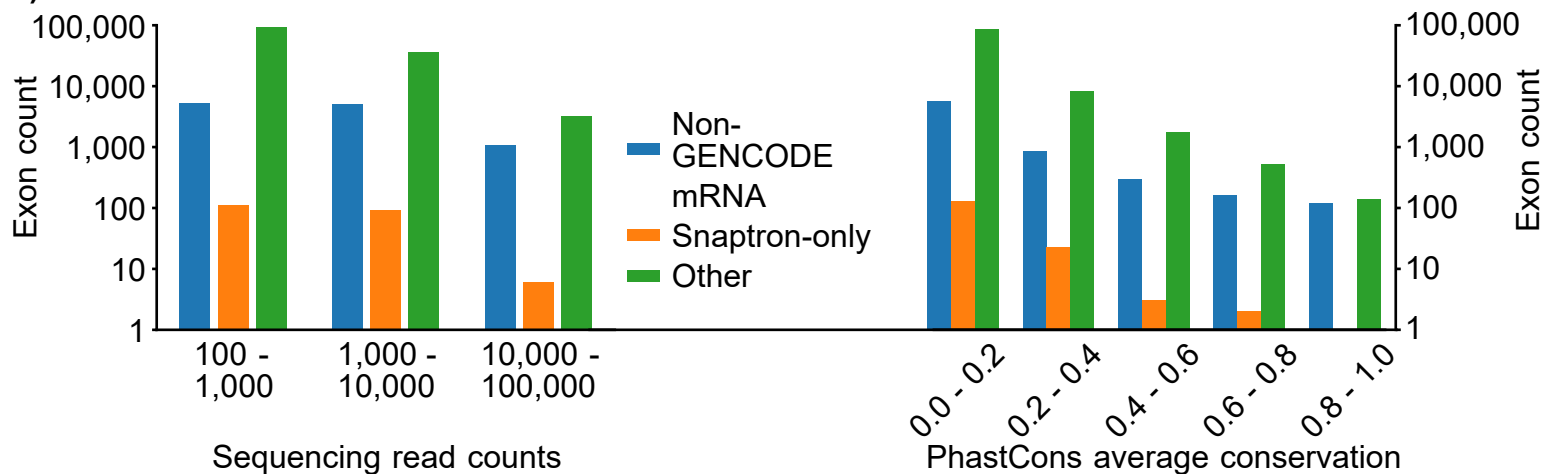
**B)**



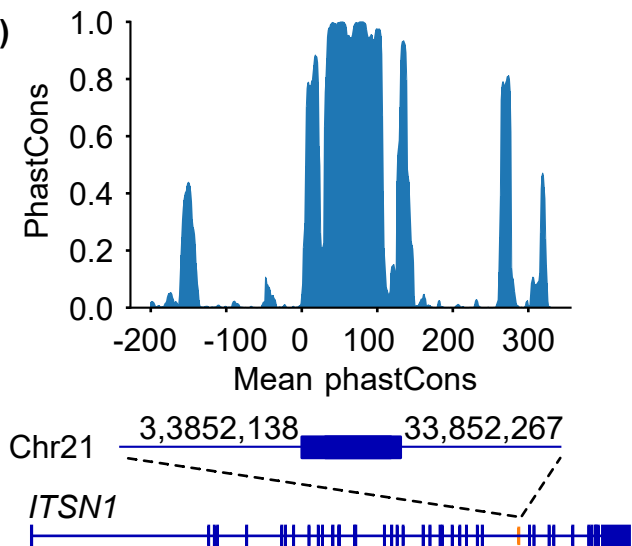
**C)**



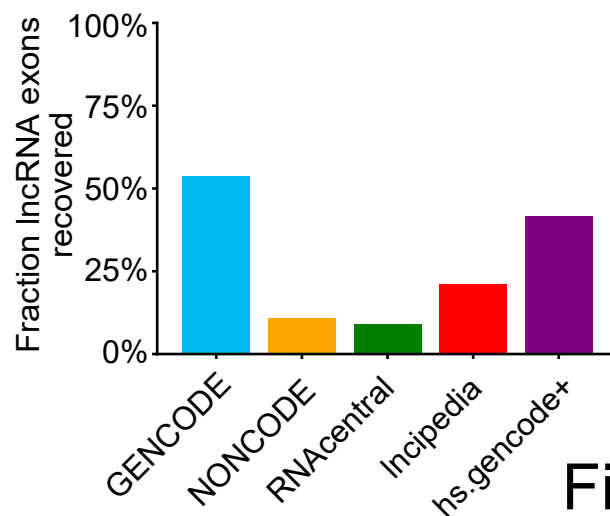
**D)**



**E)**



**F)**



**Figure 3**

### Figure 3. Trapped exons found in different categories of genomic region

- A) Pie chart depicting the proportions of exons from different genomic regions. The percentage of exons for indicated genomic regions relative to the total read count is indicated.
- B) Bar plots showing the percent of genomic bases in a trapped exon for different categories of genomic region.
- C) Boxplot depicting trapped exon sequencing read counts for different genomic regions. Whiskers depict 10<sup>th</sup> and 90<sup>th</sup> percentiles.
- D) Bar plots depicting the exon counts at different read count bins (left) or average PhastCons score bins (right) for different non-GENCODE v37 exon annotations. Exon counts are shown with logarithmic scale. Non-GENCODE v37 exon annotations, Snaptron database, and exons found by this exon trapping study are displayed. Average PhastCons scores were calculated using the sequence of the exons.
- E) Bar plot depicting PhastCons (30-way) scores for +/- 200 bp around an unannotated sense intronic exon found by exon trapping in gene *ITSN1*.
- F) Proportion of annotated exons recovered from various databases. Blue bars indicate trapped exons that are annotated in GENCODE mRNA/lncRNAs. Trapped exons annotated in other lncRNA databases are shown, with annotated GENCODE lncRNAs removed.

255 We next categorized the ~1.25 million trapped exons relative to known gene features.  
256 **Fig. 3A** shows the proportion of exons that overlap major categories of genomic  
257 sequence annotation. The largest fraction of trapped exons is “intergenic”, followed by  
258 mRNA antisense, likely due to the fact that these sequences represent most of the  
259 genome: per base, only a small fraction of each is trapped (**Fig. 3B**). For example, 5.3%  
260 of all intergenic region bases (i.e. excluding any kind of mRNA or lncRNA, and their  
261 antisense sequence) are part of a trapped exon, corresponding to an exon every 3,311  
262 bases on average. Because there is a large amount of intergenic sequence, the  
263 absolute number of “intergenic” exons is high (424,632). Similar proportions, and  
264 corresponding numbers of exons, are obtained for antisense strands (**Fig. 3B**).

265 A large proportion of the trapped exons (11.0%) were found within mRNA introns, in the  
266 forward strand (**Fig. 3A**), but the fraction of intronic sequence encompassed is lower  
267 than it is for other regions, particularly the mRNA antisense strand (**Fig. 3B**). This  
268 outcome is consistent with selection against fortuitous exon-like sequences within  
269 introns; i.e. exons that arise by chance in the antisense orientation are inconsequential,  
270 while exons that arise by chance in the sense orientation (i.e. within introns of coding  
271 pre-mRNAs) will be deleterious, and thus removed over time. These sequences  
272 nonetheless exist, and could represent potential alternative exons, or regulatory exons  
273 that trigger NMD. It is also conceivable that they are excluded by context-specific  
274 mechanisms (e.g., the sequences of neighboring exons) that are not present in our  
275 library plasmids. The “intronic” trapped exons displayed lower overall inclusion rates  
276 than any other category, including “intergenic” exons (**Fig. 3C**). Many of the “intronic”  
277 exons, especially those with higher inclusion levels, are found in other mRNA databases

278 (but not GENCODE) (6,985); an additional 240 are found in mRNA-seq data (from the  
279 Snaptron database (Wilks et al. 2018)), indicating that they are utilized in their genomic  
280 context (**Fig. 3D, left**). In addition, a subset of the “intronic” exons display primary  
281 sequence conservation (**Fig. 3D, right**). **Fig. 3E** shows an unannotated region from the  
282 *ITSN1* gene which is both conserved and trapped at high levels (10,917 read counts).  
283 Known alternative cassette exons displayed, on average, 2.5-fold lower inclusion rates  
284 when compared to all internal mRNA exons (**Fig. 3C**). First and last (i.e. terminal)  
285 mRNA exons, however, which would be expected to lack either the 3' or the 5'SSs,  
286 respectively, were rarely captured: only ~4% of these are present among the trapped  
287 exons, consistent with the rate of fortuitous splice sites across the genome.

288 Notably, 52% of GENCODE lncRNA internal exons were trapped with at least 100 reads  
289 (**Fig. 3F**), a figure comparable to that of mRNA internal exons (61%). The recovery of  
290 lncRNA internal exons that are only present in the other lncRNA databases (and not  
291 GENCODE), however, averages only 9.6% for the four databases interrogated (**Fig. 3F**,  
292 **Supplemental Fig. S5**), suggesting that these exons may have lower splicing  
293 efficiency. We assume that the well-curated GENCODE dataset is enriched for lncRNAs  
294 that splice efficiently, relative to lncRNAs that are only present in the other databases,  
295 because the impact on RNA abundance leads to a higher curation rate.

296

### 297 **Sequence features of trapped exons**

298 We next examined whether unannotated exons in the exon trapping dataset contained  
299 known sequence features of exons. We first considered the splice site scores, using

300 MaxEntScan, which outputs a maximum entropy-based score indicating whether a  
301 given base location is a 5' or 3' splice site (Yeo and Burge 2004). Read counts per exon  
302 displayed a positive overall correlation with MaxEntScan scores, for mRNA, lncRNA,  
303 and “intergenic” exons (**Fig. 4A,B**). The MaxEntScan scores of the intergenic exons  
304 display a wider spread, however, and a lower (albeit overlapping) score distribution to  
305 the annotated mRNA and lncRNA exons, for identical read counts.

306 We also investigated the prevalence of known splicing enhancing sequences within  
307 exonic regions, focusing on potential general splicing enhancer (ESE) hexamers from  
308 (Ke et al. 2011). We observed that the frequency of ESEs increases with the exon read  
309 count and that this relationship is more prominent when exons are subset by their splice  
310 site scores (**Fig. 4C, Supplemental Fig. S6A**). We reasoned that strong ESEs may  
311 offset weak splice sites, and consistent with this notion, the relationship between ESE  
312 count and read count becomes more prominent when the splice site scores are subset  
313 to narrow ranges (e.g. **Fig. 4A, Supplemental Fig. S6B**). The number of ESEs within  
314 exons also contrasts with surrounding sequence (**Fig. 4D**), illustrating that the ESE  
315 enrichment is not simply a feature of local genomic sequence, and instead contributes  
316 to exon inclusion. Moreover, exons with lower MaxEntScan scores have higher ESE  
317 density, on average (**Fig. 4E, Supplemental Fig. S6B**), further indicating that ESEs  
318 play a role in exon identity. There is a particularly strong trend for the individual SR  
319 protein-binding ESE hexamer GAAGAA (**Fig. 4F**) (Fairbrother et al. 2002), which is  
320 present more than twice as often in intergenic exons with the weakest splice sites vs.  
321 strongest splice sites (20% vs. 8%).

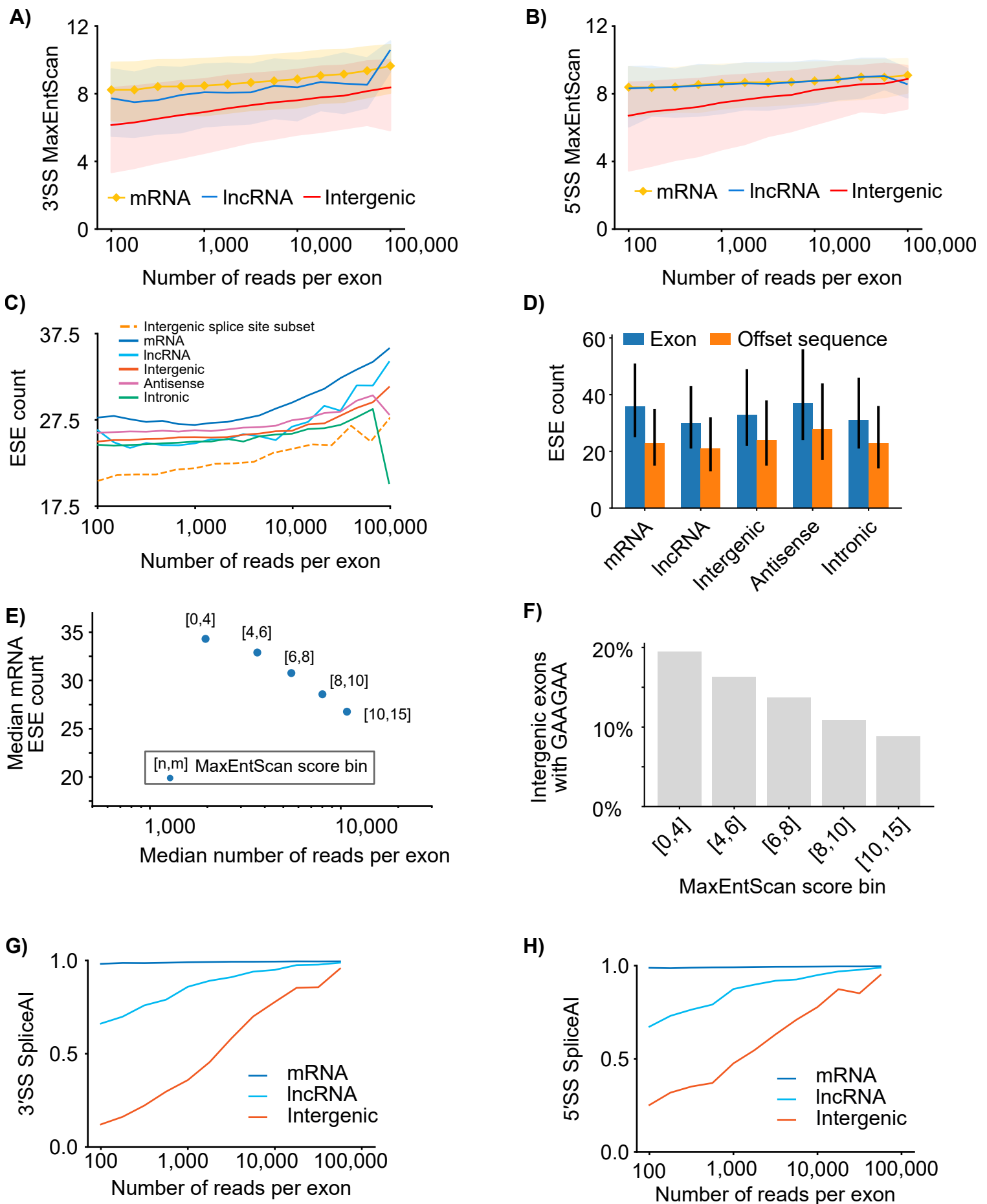


Figure 4

#### Figure 4. Known splicing signals correlations with exon read counts

- A) Line plots depicting 3' MaxEntScan scores for trapped exons from different genomic regions. Exons from mRNA, lncRNA and Intergenic regions are indicated and are binned by logarithmic read counts. Median values are displayed as lines with shaded region corresponding to 25<sup>th</sup>-75<sup>th</sup> percentiles.
- B) Same as 4A, above, but depicting 5' MaxEntScan scores.
- C) Line plots representing Splicing Enhancer (ESE) counts for trapped exon sequences from different genomic regions. ESE median values are displayed, and exons are binned by their logarithmic sequencing read counts using logarithmic bins ranging from 100 to 10,000.
- D) Bar plots depicting the median ESE counts for trapped exons (blue bars) and nearby sequence of the same length offset by 250 bp (orange bars) for different genomic regions. Offset sequences are the same length as the associated exon and corresponds to coordinates 250 bp upstream for reverse strand exons and 250 bp downstream for forward strand exons. For forward strand exons this is downstream the exon and for reverse strand exons this is upstream the exon. Range lines indicate 25<sup>th</sup>-75<sup>th</sup> percentiles.
- E) Scatter plot representing ESE count (y-axis) vs median sequencing read count (x-axis) for trapped exons, subdivided by MaxEntScan scores into groups with weaker to stronger splice sites based on splice site score bin (point label). Splice site bins indicate that contained exons have both their 3'SS and 5'SS splice sites within the labeled MaxEntScan score boundaries, between values indicate by [n,m], where n=lower score and m=upper score.
- F) Bar plot depicting fraction of intergenic exons that contain the ESE GAAGAA nucleotide sequence. Individual bars correspond to exons with both 3'SS and 5'SS MaxEntScan scores (see Figure 4E, above) within the range given in the bar label (e.g. [n < splice site MaxEntScan score < m] for both 3' and 5'SS MaxEntScan scores]).
- G) Line plots depicting 3'SS SpliceAI scores for trapped exons in different genomic regions. Values in the x-axis are logarithmic sequencing read counts using bins from 100 to 10,000 with 25 steps. For intergenic exons, the Spearman correlation between SpliceAI scores and read counts is 0.31.
- H) Same as Figure 4G, above, except for 5'SS SpliceAI scores. For intergenic exons, the Spearman correlation between SpliceAI scores and read counts is 0.12.

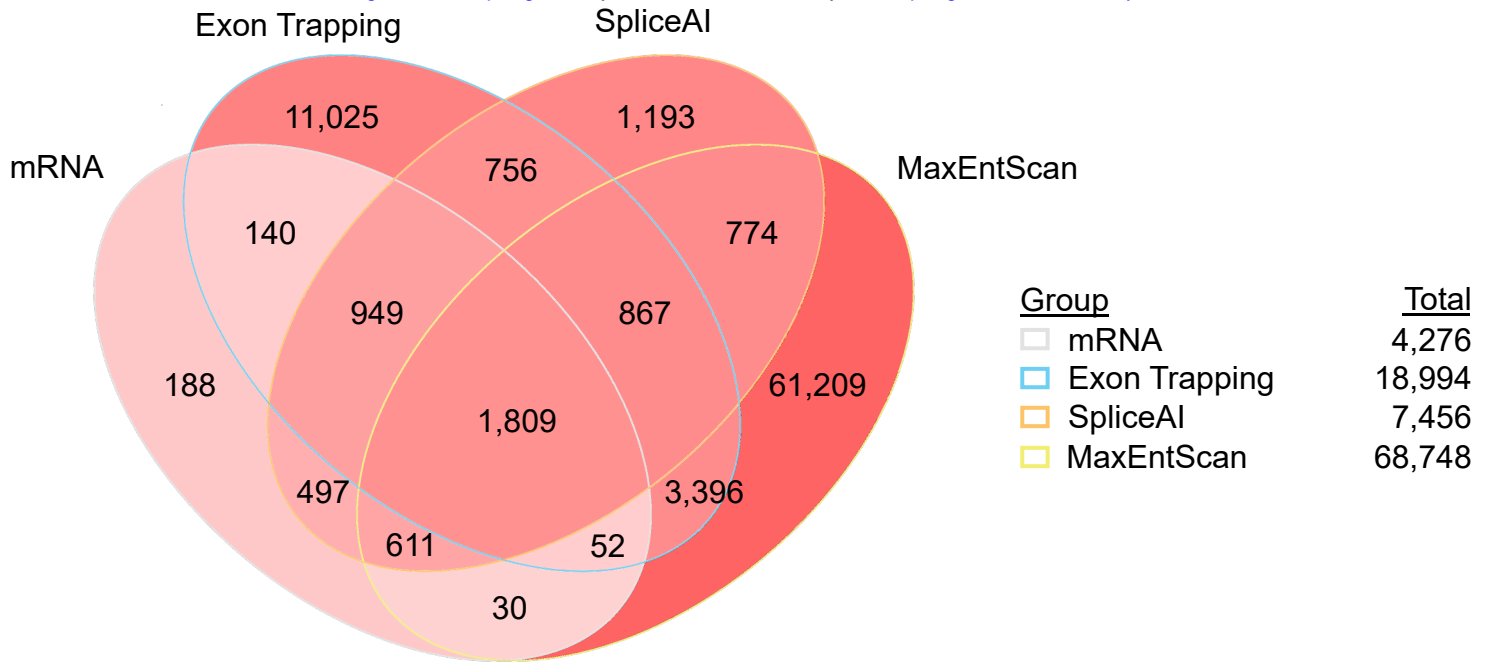
322 We also considered splice site scores from SpliceAI (Jaganathan et al. 2019), which  
323 should recognize both splice site strength and the presence of other sequences that  
324 impact splicing. The median SpliceAI prediction scores correlate most strongly with  
325 exon trapping inclusion rate for "intergenic" exons, which were not part of the SpliceAI  
326 training data (**Fig. 4G,H**). The mRNA exons, used in training SpliceAI, typically have  
327 maximal SpliceAI scores, irrespective of their autonomous splicing potential (i.e. read  
328 count), suggesting that SpliceAI has learned features of mRNA exons that are distinct  
329 from their autonomous splicing potential. Average SpliceAI scores for GENCODE  
330 lncRNA exons are higher than those of intergenic exons, but lower than those of  
331 mRNAs exons, for equivalent read counts (**Fig. 4G,H**), consistent with the notion the  
332 additional features learned by SpliceAI do not relate only to coding potential, and may  
333 encompass productive elongation or transcript stability (Jaganathan et al. 2019).

334

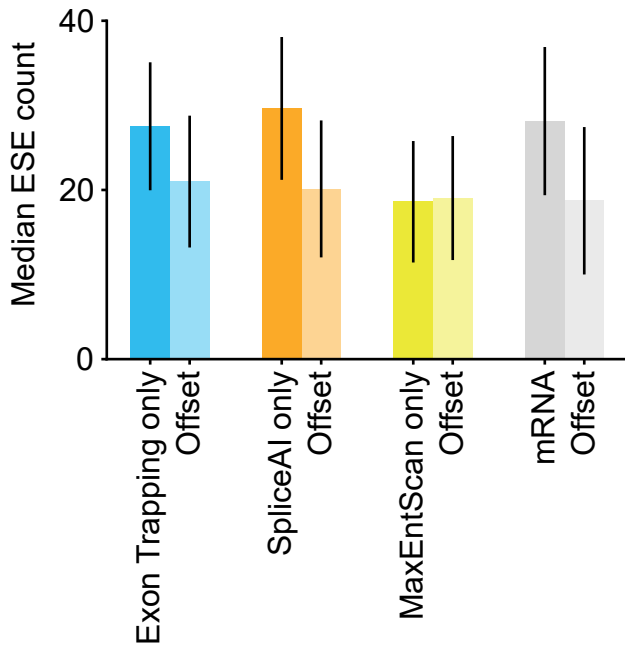
### 335 **Many trapped exons lack known exon-associated sequence features**

336 We next asked whether the sequence features above could completely account for the  
337 trapped exons. To do this, we required exon predictions that are based only on these  
338 sequence features. Neither SpliceAI nor MaxEntScan explicitly predict exons, but it is  
339 possible to derive exon predictions by simply associating strong predicted 3' splice sites  
340 with proximal strong predicted 5' splice sites (here, using a SpliceAI score cutoff of 0.2,  
341 a MaxEntScan score cutoff of 6, separated by 63-222 bases (10-90% percentile of  
342 mRNA internal exons). We examined the overlap of such exons predicted across  
343 Chromosome 17 (to reduce computation time) from either the SpliceAI or MaxEntScan  
344 outputs, and compared to exonic regions we obtained from exon trapping, and

A)



B)



C)

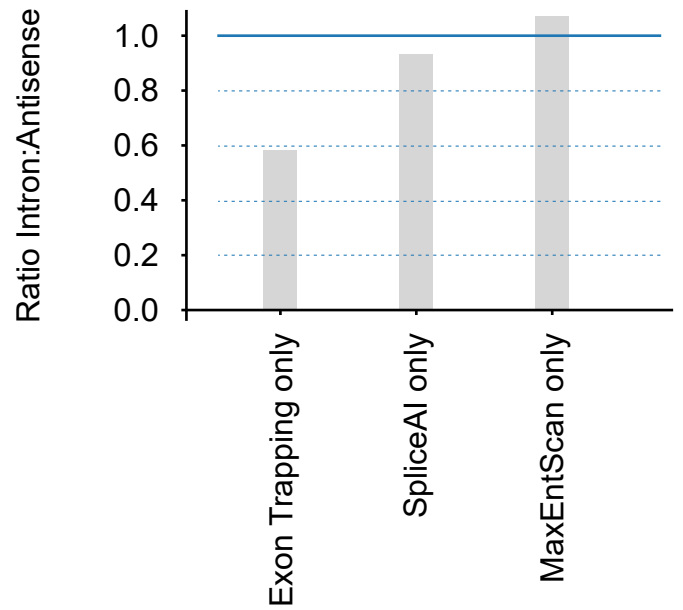


Figure 5

### Figure 5. Overlaps between exons detected using different approaches

- A) Venn diagram depicting Chromosome 17 forward strand exons found by different exon calling approaches. Exons are labeled as mRNA (annotated mRNA & lncRNA internal exons), ET (exons found by exon trapping), MaxEntScan based on MaxEntScan scoring, and SpliceAI based on SpliceAI scoring. Exon counts corresponding to overlapping regions are indicated and are colored red linearly with intensity determined by  $(\log_{10} \#exons)$ .
- B) Bar plot of median ESE counts for exons and offset sequences identified using approaches listed in 5A. Offset sequences are the same length as the associated exon and corresponds to coordinates 500 bp upstream for reverse strand exons and 500 bp downstream for forward strand exons. For forward strand exons this is downstream the exon and for reverse strand exons this is upstream the exon. Range lines indicate 25<sup>th</sup>-75<sup>th</sup> percentiles. Refer to **Figure 5A** for x-axis labels.
- C) Bar plots showing the ratio of intronic to antisense exon counts found for the different exon finder approaches.

345 annotated exons. We separately tallied exons that overlapped perfectly (**Fig. 5A**) from  
346 those that overlapped partially (**Supplemental Fig. S7A**).

347 As described above, most of the annotated exons were captured by exon trapping, but  
348 a large proportion of the trapped exons did not overlap with annotated exons. The  
349 majority of trapped exons also did not overlap with exons predicted by either SpliceAI or  
350 MaxEntScan: more than half of the trapped exons - 11,025/18,994 - did not overlap with  
351 any of the other exon sets. We asked whether these 11,025 trapped exons have  
352 characteristic properties beyond their low 5' or 3' SS scores (which we infer because  
353 they were not detected by MaxEntScan) (**Supplemental Fig. S7B,C**). Their lengths and  
354 base content are not unusual (**Supplemental Fig. S7D,E**). They overlap with repetitive  
355 elements at roughly the frequency (~50%) that repetitive elements occur in the genome.  
356 The trapped exons are enriched for ESEs, however, relative to adjacent sequence, to a  
357 degree that is similar to known exons and SpliceAI-predicted exons (**Fig. 5B**). The  
358 difference in ESE frequency appears insufficient to explain why these exons are  
359 included, however: the range largely overlaps with that of the tens of thousands of  
360 MaxEntScan-predicted exons, which have stronger splice sites overall, and yet are not  
361 trapped. We also asked whether the detection of these exons might be due to their  
362 genomic context. To do this, we queried the density of ISS sequences (Wen et al. 2010)  
363 for mRNA and trapped intergenic exons, reasoning that there may be ISS sequences in  
364 distal intergenic genomic sequence (and not mRNA introns) that would be absent from  
365 the reporter inserts. We observed little difference however (**Supplemental Fig. S8A,B**).

366 An additional and intriguing observation is that, as noted above, the trapped exons are  
367 very significantly depleted from introns ( $P < 7.1 \times 10^{-293}$ , Wilcoxon rank sum test), but not

368 the mRNA antisense strand, and trapped exons that are detected in introns tend to have  
369 low read counts (median 448 vs. 6214 for mRNA). Among exon detected only by exon  
370 trapping there is a 2-fold bias towards introns (**Fig. 5C**). In contrast, exons predicted  
371 only by SpliceAI or MaxEntScan have roughly similar numbers of exons within  
372 annotated introns, relative to exons predicted in the antisense strand (**Fig. 5C**). Thus,  
373 the exon trapping data must contain some biologically meaningful information not  
374 captured by the splicing predictors.

375

### 376 **Transposons as a source of novel exons**

377 Transposable elements (TEs) make up roughly half of the human genome, and are a  
378 prevalent source of new genetic material, including exons (Sorek 2007). *Alu* elements,  
379 for example, are a source of alternative mRNA exons due to the presence of sequences  
380 consistent with splice sites near the repeat's 5' end and an upstream polypyrimidine  
381 tract arising as a product of retrotransposition (Makalowski et al. 1994; Dewannieux and  
382 Heidmann 2005; Sorek 2007). We asked whether specific classes of TEs are enriched  
383 or depleted among the trapped exons, on a base-by-base level, and found many cases  
384 of both enrichment and depletion (**Fig. 6A**). We reasoned that enrichment of specific TE  
385 classes might be due to splice sites within the TE, and indeed such cases are readily  
386 identified. Examination of compiled instances across the genome, and inspection of the  
387 consensus models on Dfam 3.4 (Storer et al. 2021), revealed that the trapped exon  
388 often corresponds to a common segment of the transposon, beginning and/or ending at  
389 a location where the ancestral element contained sequences resembling 5' and/ or 3'  
390 splice sites. An example (DF0000317.4, the 5' end of L1 retrotransposon L1P2) is

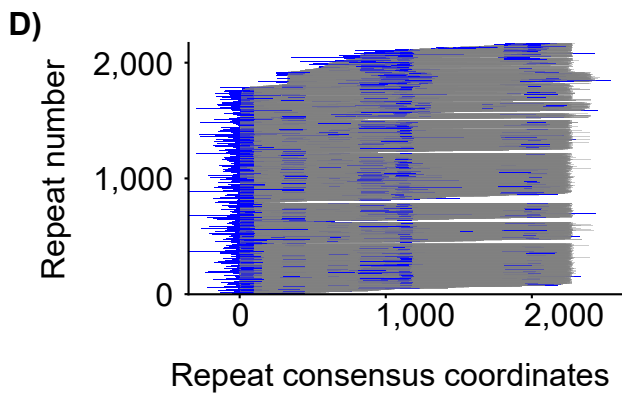
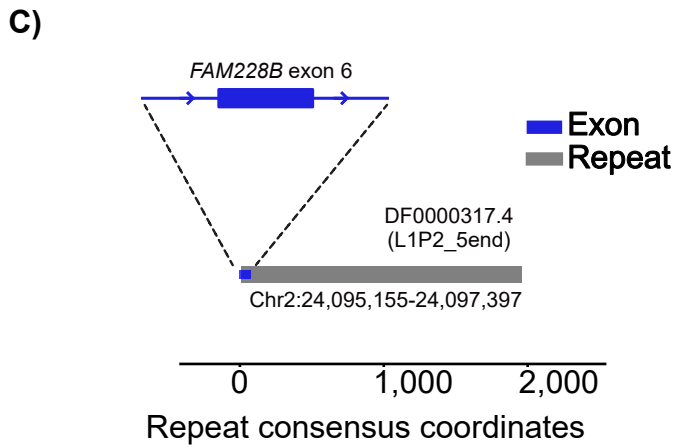
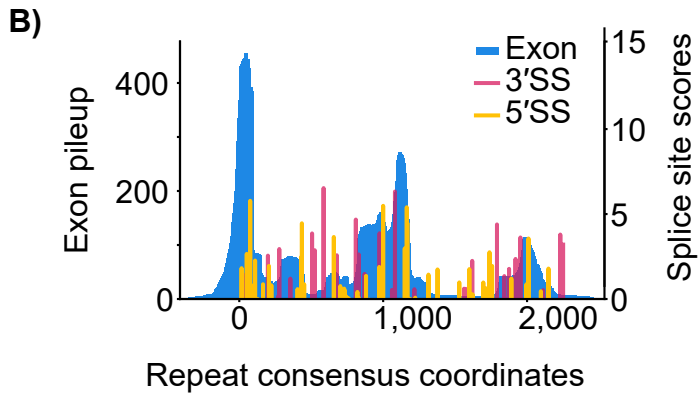
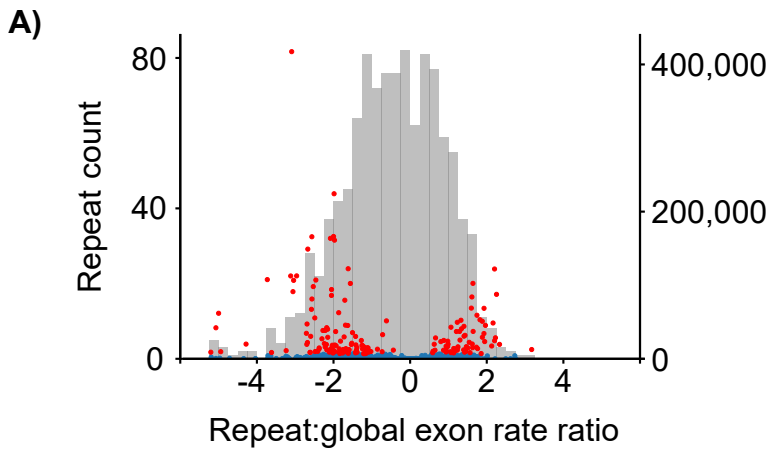


Figure 6

## Figure 6. Exons overlapping repetitive elements

- A) Histogram depicts ratio of exon overlaps for different repeat families, relative to the global genomic exon rate. Volcano plot shows  $\log_{10}$  p-values of repeat family exon bases enrichment (hypergeometric test) vs the repeat exon enrichment relative to the global exon rate. Red dots indicate top 5% of p-values. Blue dots indicate the bottom 95% of repeat-values.
- B) Histogram pileup depicting sequencing reads overlapping repeat instances of DF0000317.4 (5' end of L1 retrotransposon L1P2) in the human genome. Histogram maps sequencing reads in the genome to the Dfam repeat consensus model. MaxEntScan 3'SS and 5'SS scores Dfam are also shown across the Dfam repeat consensus model with scores above 0 shown with colored bars. Repeat consensus coordinates start at 0.
- C) Diagram depicting overlap between trapped exon 6 from *FAM228B* (Chr2:24,095,141-24,095,230) and a genome instance of L1P2 transposon DF0000317.4 (Chr2:24,095,155-24,097,397).
- D) Diagram depicting overlap between all trapped exons and associated genomic repeat instances for L1P2 transposon DF0000317.4. Rows are sorted by repeat start and end coordinates for the Dfam repeat consensus model.

391 shown in **Fig. 6B**. This repeat overlaps four annotated mRNA exons; one example is  
392 depicted in **Fig. 6C**. The ~2,000 additional trapped exons overlapping this repeat largely  
393 favor internal exons arising inside near full-length repeat sequences (**Fig. 6D**). We  
394 assume that these sequences are fortuitous, since splice signals are short and  
395 degenerate, but they nonetheless represent a ready means by which these transposons  
396 can contribute to the evolution of existing genes.

397

### 398 **Conservation of splice sites and trapped exons**

399 Finally, we examined trends in sequence constraint of various types of exons: coding  
400 exons, trapped exons found in annotated introns (“intronic”), lncRNA exons, and novel  
401 trapped exons found in intergenic regions (“intergenic”). **Fig. 7A,B** shows that the  
402 splicing sequences of known coding exons are highly constrained (by phyloP (Pollard et  
403 al. 2010)), and also shows a three-base periodicity within the exons, presumably due to  
404 codon bias and wobble. In contrast, none of the other classes displayed strong  
405 conservation, on average. For lncRNAs this phenomenon is well-documented (Wang et  
406 al. 2004a; Kutter et al. 2012; Li and Yang 2017). The intergenic regions detected by  
407 exon trapping therefore behave much like lncRNA exons in this aspect.

408

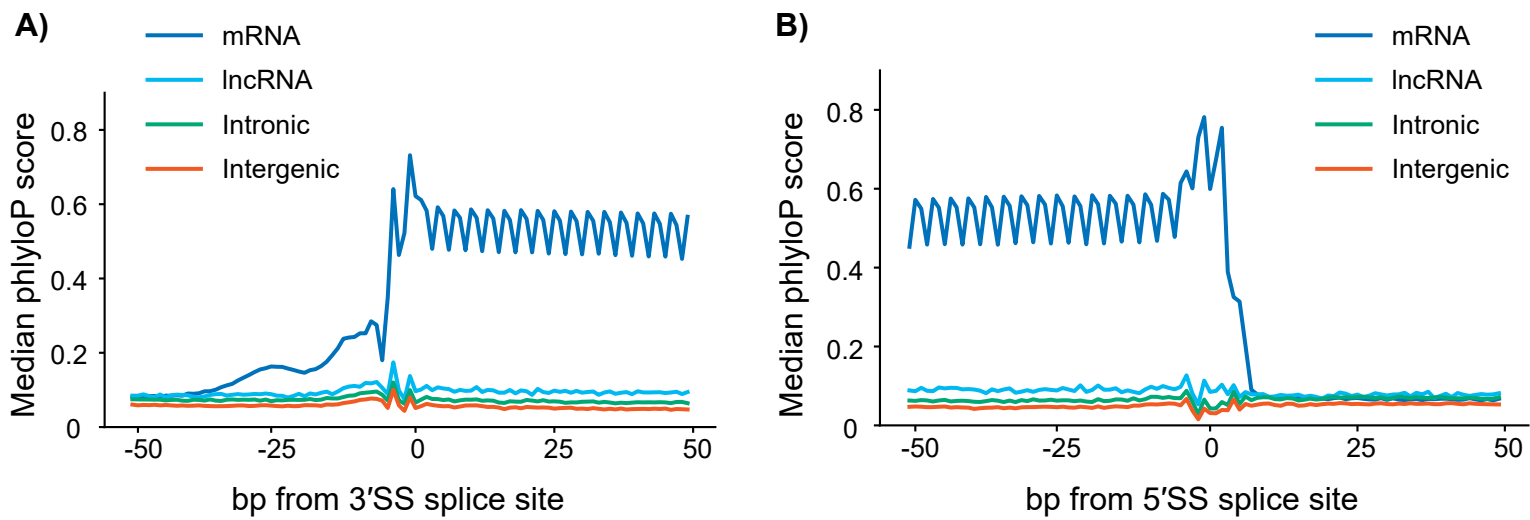


Figure 7

### **Figure 7. Conservation around 3' and 5' splice sites of trapped exons**

- A) Diagram depicting sequence conservation of trapped exons around 3'SS for different genomic regions. Shown are phyloP scores for the 100 bp region centered around the 3'SS for trapped exons from mRNA, lncRNA, and intergenic regions, as indicated.
- B) Same as A) except phyloP scores were calculated for the 100 bp region centered around the 5'SS for trapped exons.

409

410 **DISCUSSION**

411

412 An exon trapping screen of the human genome in a single cell type captured >1 million  
413 sequences that splice as internal exons. Most known mRNA exons are trapped,  
414 demonstrating that exons of coding genes are predominantly autonomous, consistent  
415 with the exon definition model. Moreover, splicing-competent sequences are common in  
416 non-genic locations of the human genome. This finding presents a ready explanation for  
417 the large number of poorly-conserved lncRNAs.

418 A central finding of these exon trapping assays is that a majority of the known exons in  
419 human protein-coding genes are autonomous. This observation does not rule out a role  
420 for context (i.e. flanking sequence) in exon recognition. A clear limitation is that only five  
421 different splicing vectors were employed, and they are highly related to each other. The  
422 assay in its current form was not intended as a universal exon discovery tool, only as an  
423 initial survey. The fact that most coding exons were trapped, however, indicates that a  
424 specific context is not required for splicing of most individual exons. Another limitation of  
425 this study is that the screens were performed in a single cell type. Thus, the impact of  
426 cell-type-specific splicing factors would not have been captured. To our knowledge,  
427 such factors would mainly be expected to influence alternative exons, which were  
428 indeed trapped at lower rates.

429 The data also demonstrate that there are many sequences in intergenic regions, and  
430 antisense to known genes, that can be spliced into a heterologous transcript. Such  
431 sequences also exist in introns, but at a significantly reduced frequency. These “intronic

432 exons” bear a resemblance to pseudoexons – intronic sequences that are flanked by  
433 splice sites but that are not observed in spliced mRNA. Aberrant inclusion of  
434 pseudoexons is thought to represent an underreported disease mechanism (Petersen et  
435 al. 2022). These exons appear to be determined in part by the simple presence of  
436 strong splice sites in close proximity, but there are many which have weak splice sites.  
437 Frequency of ESE sequences provides only a partial explanation for their inclusion as  
438 exons. Given the large number of variables to explore (e.g. positioning and relative  
439 weights among ESEs; presence of intronic cues; splicing silencers, and combinations or  
440 conditional relationships among features), we did not attempt to develop new exon  
441 predictors as part of this study. Nonetheless, the data presented here, which includes  
442 not only many new autonomous exons, but also associated read counts, provides a new  
443 resource for analysis of both exon identity and exon inclusion levels spanning a large  
444 dynamic range.

445 It is possible that the hundreds of thousands of exons detected by exon trapping in  
446 intergenic and antisense regions are not already annotated as exons mainly because  
447 they are not expressed, and/or do not form part of stable transcripts. Moreover, the  
448 observation that autonomous exons with little or no sequence constraint occur  
449 frequently in intergenic space suggests that the required sequence features arise  
450 frequently by random mutations. SpliceAI, which was trained to predict splice sites in  
451 coding exons, also readily predicts splice sites in intergenic sequence (**Fig. 5**), SpliceAI  
452 also predicts splice sites in random sequence at similar rates to what is observed in the  
453 exon trapping data: using the thresholds described earlier, it identifies 1,856 exons in a  
454 dinucleotide-permuted positive strand of Chromosome 17. This number, and the range

455 of splice site scores, is similar to the number of unannotated exons found on the real  
456 positive strand Chromosome 17 (**Supplemental Fig. S7F**). This outcome is consistent  
457 with the low information content of splice sites, and the fact that nearly 25% of 6mers  
458 show ESE activity (Ke et al. 2011).

459 We speculate that lncRNAs, whose promoters often overlap with endogenous retroviral  
460 LTRs and enhancers (Kelley and Rinn 2012; Engreitz et al. 2016), may result from  
461 coupling promoter-like sequences to randomly arising autonomous exons in the  
462 surrounding DNA. This notion is consistent with the fact that there is little overall  
463 evidence of selection on lncRNA primary sequence (Wang et al. 2004a). In addition, the  
464 read counts we obtained for lncRNA exons are only about half of those obtained for  
465 coding exons, and similar to the read count abundances of alternative exons. Lower  
466 splicing efficiency could partly explain the lower expression levels and lower exon  
467 numbers of lncRNAs relative to mRNAs (Orom et al. 2010) and is consistent with  
468 observations that lncRNA exons are often chaotically spliced (Deveson et al. 2018).

469 These findings reinforce speculation that much of the “dark matter” transcriptome may  
470 be a by-product, or even an expected component, of the regulation of known genes, as  
471 well as a source of novel genetic entities. The exon trapping data presented here  
472 identify regions that would be incorporated into such transcripts, if expressed. These  
473 data also offer an orthogonal set of genomic exons appropriate for understanding  
474 splicing sequences not under mRNA selection, providing sequences across a large  
475 dynamic range that should enable additional insights into the splicing code.

476

477

478 **METHODS**479 ***Reporter intron***

480 We gene synthesized the 6<sup>th</sup> intron (Chr3:185919497-185921103, '-' strand) of the gene  
481 *Tra2B* as the reporter intron. This intron has strong splice sites and a native length of  
482 1.6 kb. We changed several bases to remove possibly cryptic 5ss splice sites and  
483 added sequence to the middle of the intron to facilitate genomic fragment insertion by  
484 restriction digest and Gibson Assembly.

485

486 ***Plasmid preparation***

487 The five different plasmid backbones (**Supplemental\_details.xlsx** sheet "Vectors",  
488 **Supplemental\_vector\_details.docx**) were generated as follows. The three plasmids  
489 that vary in reading frame position of the splice site were designed so that the splice  
490 sites are compatible with the coding sequence; these were ordered by gene synthesis  
491 and cloned by Gateway into a pcDNA3.1 based plasmid. The two variants with  
492 mutations in the 5'SS and putative RBFOX1 binding site, respectively, were made using  
493 Gibson assembly and primers containing variant sequences (see  
494 **Supplemental\_details.xlsx** and **Figure S1A**).

495 To generate libraries, each of the five backbones was inoculated in *E. coli* (ElectroMAX  
496 Stbl4 competent cells) and plasmid DNA was extracted using Qiagen HiSpeed Midiprep  
497 kit. The plasmids were digested with restriction enzymes as follows. 10 µg plasmid was  
498 incubated in 250 µl 1X CutSmart buffer with 10 µl AgeI-HF enzyme and 10 µl NotI-HF  
499 enzyme, at 37°C for 20 minutes. DNA was then purified using Zymo Clean and  
500 Concentrator 5 at 2:1 binding buffer:restriction digest. DNA was eluted twice, incubating  
501 for 1 minute with 70°C 8 µl of NEB DNA elution buffer. (For reagents, see  
502 **Supplemental\_details.xlsx** sheet "Products").

503 We prepared two different adapter duplexes (age\_25/age\_common and  
504 not\_25/not\_common, **Supplemental\_details.xlsx** sheet "Oligos"), one for the AgeI site  
505 and one for NotI by heating 10 mM of the oligo pair in 25 mM NaCl + 0.5 mM EDTA to  
506 90°C for 2 minutes moving to room temp for 5 minutes.

507

508 ***Adapter ligated gDNA preparation***

509 We extracted gDNA from HEK293 cell cultures (six well plate) using the PureLink  
510 Genomic DNA Mini Kit, and fragmented 10 µg of the resulting DNA using NEB NEBNext  
511 dsDNA Fragmentase (1.6 mg of gDNA in 54 µl of 1X dsDNA Fragmentase buffer,  
512 incubate on ice for 5 minutes, add 6 µl dsDNA Fragmentase, mix by vortexing, incubate

513 at 37°C for 20 minutes, stop with 15 µl 0.5M EDTA mixed by pipetting). The resulting  
514 fragmented DNA was purified with Zymo Clean and Concentrator 5 µg and eluted with  
515 NEB DNA elution buffer. We then end repaired the DNA using NEB end repair kit  
516 (E6050L) and added dA tails using NEB dA-tailing. We then ligated the eluted gDNA (20  
517 ng/µl reaction volume) to dT-tailed DNA adapters (generated by annealing two oligos  
518 pairs (age\_25/age\_common and not\_25/not\_common, **Supplemental\_details.xlsx**  
519 sheet “Oligos”) for 30 minutes at room temperature using NEB Quick Ligase. We  
520 cleaned the reactions using 1.8X AMPure XP beads and eluted the DNA with NEB  
521 elution buffer supplemented with Tween20 to 0.1%.

522 DNA was then amplified by PCR using NEB Q5 enzyme and primers complimentary to  
523 the adapters (NEBb\_1, NEBb\_2, **Supplemental\_details.xlsx** sheet “Oligos”) (PCR  
524 program 98°C 30 sec; 98°C 10 sec, 57°C 15 sec, 72°C 60 sec (N times); 72°C 120  
525 sec). At this and subsequent PCR steps, we performed qPCR to determine the number  
526 of cycles to amplify the library. We employed the Q5 HF protocol scaled to 60 µl with 1  
527 µl of cDNA, dispensed as 3 replicates of 15 µl, 1X Evagreen dye, 0.5 mM primer. We  
528 selected the cycle number that is two cycles before 50% amplification is exceeded.

529 DNA was sized on 1% agarose gel cast with 1X SYBR Safe to a range of ~500-1100 bp  
530 using Quick-Load Purple 100 bp DNA Ladder (N0551L). DNA was gel extracted using  
531 Qiagen Qiaex ii gel extraction kit, and eluted with 25 µl of NEB DNA elution buffer.  
532 Eluted DNA was again amplified, with qPCR to avoid saturation, in preparation for  
533 Gibson Assembly.

534

### 535 ***Gibson Assembly and Library Amplification***

536 Gibson Assembly of gDNA into plasmids was performed using NEBuilder HiFi DNA  
537 Assembly Master Mix in 200 µl with 2 µg plasmid and 300 ng of gDNA. The five  
538 assembled plasmid libraries were phenol extracted and precipitated, and DNA was  
539 resuspended in 12 µl NEB DNA elution buffer. Resuspended DNA was electroporated  
540 using 2 µl plasmid and 20 µl of ElectroMAX Stbl4 competent cells according to kit  
541 protocol, with five replicates for each of the five plasmid libraries. *E. coli* transformants  
542 were grown overnight with shaking at 30°C in 110 ml of LB with Carbenicillin antibiotic  
543 (100 mg/ml).

544 Overnight cultures were pelleted were split into 2×50 ml and pelleted in a tabletop  
545 centrifuge set to 4°C, supernatant discarded, and pellets stored at -20°C. Cell pellets  
546 were extracted using Qiagen HiSpeed Midiprep kit and eluted using 500 µl of the  
547 supplied TE buffer.

548

### 549 ***Transfections***

550 We plated 1.2 million HEK293 cells into 10 cm plates and incubated overnight in DMEM  
551 to ~75% confluency at the time of transient transfection. We performed transient  
552 transfections with 10 µg of plasmid DNA using Promega Viafect according to reagent  
553 protocol. DNA was incubated with the Viafect transfection reagent and Opti-MEM for 15  
554 minutes at room temperature before transfection. Growth medium was changed after 24  
555 hours, and after a further 24 hours, the transfected cells were harvested from each plate  
556 with 10 ml of TRIzol, and stored at -80 °C.

557

### 558 ***RNA extraction***

559 We thawed the TRIzol-cell mixture by incubating on ice for 10 minutes. Two (2) ml of  
560 chloroform was added and mixed by inversion, then tubes were incubated at room  
561 temperature for 2 minutes.

562 We then centrifuged the tubes for 15 minutes at ~4,000 × g at 4°C. We split each  
563 sample's aqueous phase into two pre-spun phaselock tubes (pre-spun with 1 ml of  
564 BCP). We added an equal volume of acid phenol (~3 ml) and .2 volume of BCP (~.6 ml)  
565 to each tube, mixed by inversion, and incubated at room temp for 2-3 minutes before  
566 centrifuging for 15 minutes at ~4,000 × g at 4°C. We performed a final spin by adding 1  
567 ml of chloroform to each sample which we mixed by inversion, incubate at room temp  
568 for 2-3 minutes, and centrifuged for 15 minutes at ~4,000 × g at 4°C. We transferred the  
569 supernatant.

570 We precipitated these samples using sodium acetate ethanol precipitation along with 10  
571 µl of Glycoblue co-precipitant. Precipitated RNA was resuspended in 300 µl of H<sub>2</sub>O.

572

### 573 ***Poly(A) enrichment***

574 We isolated poly(A) RNA using the Invitrogen Poly(A) Purist MAG Kit. Each sample of  
575 extracted RNA was brought to 500 µl with 200 µl H<sub>2</sub>O. We then added 500 µl of 2X  
576 binding buffer. This mixture was added to 30 µl of prepared beads. This mixture was  
577 denatured at 75°C while shaking at 300 rpm. Samples were eluted twice, using elution  
578 buffer heated to 80°C. 2 µl glycogen was added to the poly(A) RNA and samples were  
579 precipitated per kit instructions. Precipitated RNA was resuspended in 40 µl of THE  
580 RNA storage solution.

581

### 582 ***Reverse transcription***

583 We reverse transcribed the RNA using superscript IV (Thermo Fisher Scientific  
584 #18090010) and a primer targeting reporter transcript sequence downstream of the 3'  
585 exon (pTH\_T5\_RT\_3). We followed the kit protocol scaled to 3.5X volume with 2 mM of  
586 pTH\_T5\_RT\_3 and 38.5 µl of Template RNA. cDNA synthesis was carried out for 11

587 minutes. We hydrolyzed RNA using the cDNA clean-up protocol from Zymo DNA Clean  
588 Concentrator-5 (Cat #D4004) with volumes scaled 1.4-fold. We performed the final  
589 elution with 20 µl of NEB DNA elution buffer warmed to 65°C.

590

### 591 ***Amplification with reporter-specific primers***

592 We PCR amplified the reporter transcript with NEB Q5 polymerase and primers  
593 targeting the flanking exons, tailed to add Illumina Nextera sequences (primer pairs  
594 NS601- NS602, NS603- NS604, NS605- NS606) (PCR program 98°C 30 sec; 98°C 10  
595 sec, 70°C 15 sec, 72°C 60 sec (N times); 72°C 120 sec). We determined the number of  
596 cycles to amplify using qPCR as described previously.

597

### 598 ***PCR 1: Sizing using native acrylamide gels***

599 We cleaned the PCR with Ampure XP beads. We gel-sized this PCR reaction using  
600 BioRad precast 5% Mini-PROTEAN TBE Gels run for 45 minutes at 75V in 1X TBE  
601 buffer. We stained the gels using 1X SYBR Gold and extracted all material longer than  
602 the empty reporter PCR product (i.e. ~ to the top of the well). We eluted and recovered  
603 the DNA from the gels slice using QIAEX II Gel Extraction Kit, using the kit protocol for  
604 native acrylamide gel extraction. We eluted the DNA in diffusion buffer at 50°C and  
605 1400 rpm for 1 hour, followed by overnight rotation at room temperature. During the  
606 subsequent recovery of DNA from the diffusion buffer, DNA was incubated with shaking  
607 (1400 rpm) for 10 minutes, instead of occasional vortexing. The recovered DNA was  
608 eluted using 25 µl of NEB elution buffer supplemented to 0.1% Tween20.

609

### 610 ***PCR 2: Amplify DNA with Illumina index primers***

611 We amplified the gel-sized DNA in order to attach index sequences, and the remaining  
612 sequences required for Illumina sequencing. This PCR used Illumina UDI primers  
613 (Illumina CAT #20027216, IDT for Illumina Nextera DNA Unique Dual Indexes Set D)  
614 (PCR program 98°C 30 sec; 98°C 10 sec, 70°C 15 sec, 72°C 60 sec (N times, N  
615 determined by qPCR); 72°C 120 sec). Index usage in Resources Table and  
616 **Supplemental\_details.xlsx** sheet "Indexes" and "Oligos". We then cleaned and sized  
617 the PCR again as described above. We determined the number of cycles to amplify  
618 using qPCR as described previously.

619

### 620 ***PCR 3: Amplify DNA to quantities for sequencing***

621 We PCR amplified the sized DNA using oligos corresponding to Illumina P5 and P7  
622 sequences (PCR program 98°C 30 sec; 98°C 10 sec, 70°C 15 sec, 72°C 60 sec (N  
623 times); 72°C 120 sec). We determined the number of cycles to amplify using qPCR as

624 described previously. We cleaned these PCRs using Ampure XP beads and eluted  
625 using 20 µl NEB DNA elution buffer supplemented to 0.1% with Tween20.

626

### 627 ***Illumina sequencing***

628 Sequencing was performed on an Illumina NovaSeq S4 flowcell, with 2×150 bp reads.  
629 The sequencing facility (The Centre for Applied Genomics and The Hospital for Sick  
630 Children, Toronto) performed Ampure XP clean up on some samples, to reduce Illumina  
631 dimers, and “dark cycled” (chemistry only) the first 18 bases of read 1, in order to  
632 bypass constant sequences of the reporter assay, and increase sequencing read base  
633 complexity used for Illumina flowcell deconvolution processing. The dark cycling  
634 encompasses part of the 3’ splice site; thus, 3 or 4 bases of each exon are removed  
635 (because different length primers were used, depending on backbone).

636

### 637 ***Illumina data processing and read mapping***

638 We truncated both read 1 and read 2 to the first 45 bases, and removed reads with  
639 cigar strings containing 'S'. We used BBDUK (Key Resources Table) to trim adapter  
640 sequences and remove empty reads. BBDUK was run with additional parameters  
641 “ktrim=r k=13 mink=5 hdist=1 tpe tbo threads=2”).

642 We aligned trimmed sequencing reads to the hg38.2bit genome sourced from UCSC  
643 (<https://hgdownload.cse.ucsc.edu/goldenpath/hg38/bigZips/hg38.2bit>). Alignments were  
644 performed using HISAT2 (Kim et al. 2015) with parameter “--max-intronlen 1500”.

645 To aggregate reads, we sorted the individual aligned libraries by their chromosome  
646 position (chromosome, 5’ coordinate and then 3’ coordinate, then strand), using linux  
647 sort with C locale. We retained mapped reads with mapping quality greater than or  
648 equal to 20. We adjusted the mapped read termini to account for the 3-4 bases of the 5’  
649 exon end that were clipped by the “dark cycle” sequencing (library adjustment in and  
650 **Supplemental\_details.xlsx** sheet “Indexes”). Reads with identical ends were counted  
651 and collapsed into one entry in a new file, which represents exon read counts. Most of  
652 the analyses in the paper were performed using data pooled from the different libraries,  
653 which were merged by collapsing the results from the five libraries using the same  
654 overall approach. To account for sequencing indels, we also collapsed overlapping  
655 exons (i.e. groups of overlapping reads) by first adjusting the termini of all exons up to 2  
656 bases, if the adjustment coincides with a higher scoring splice site, and then again  
657 collapsing these reads to arrive at the final exon dataset.

658

### 659 ***Removal of exons associated with the reporter intron***

660 Some of the most abundant exons align to the native reporter intron sequence which is  
661 present in every assayed plasmid. We removed exons that are fully contained in the  
662 Tra2B intron we used (Chr3:185919497-185921103) since these counts are presumably  
663 rarely splicing exons in the reporter present in every assayed plasmid.

664

### 665 ***GENCODE gene annotation***

666 We used GENCODE v37 'basic' to determine annotated exons. Transcripts with any  
667 'transcript support level' were included in analyses. Transcripts with 'level' 3 or less  
668 were included for analysis. The mRNA exons are GENCODE exons of  
669 gene\_type='protein\_coding'. The lncRNA exons are GENCODE exons of  
670 gene\_type='lncRNA'. First and last exons were obtained for each mRNA transcript. The  
671 remaining exons of each transcript define the internal mRNA exons. The GENCODE  
672 lncRNA exons were obtained in the same manner. Transcript boundaries were obtained  
673 from the first and last exon coordinates.

674 We aggregated all remaining exons into a total exon trapping dataset. We created a  
675 primary exon dataset by grouping overlapping exons sharing a 3'SS. Primary 3'SS  
676 clusters with greater than or equal to 100 sequencing reads were used to threshold  
677 exon trapping exons. The exon with the highest number of sequencing reads was  
678 selected as the primary exon to represent this exonic region. For exons sharing the  
679 same 3'SS we selected the exon with the majority of reads as the primary exon.

680

### 681 ***Intergenic, intronic, and antisense exons and region categories***

682 We defined the intergenic regions as those that do not overlap with sense  
683 mRNA/lncRNA transcripts or their antisense regions. Intronic regions were identified in  
684 transcripts regions with all annotated exons removed. We selected the antisense  
685 regions by starting with the antisense regions of transcripts and removing any annotated  
686 transcript (this, for example, excludes annotated antisense transcripts). We determined  
687 the exon trapping exons of each region type by taking those exons fully contained within  
688 the above-described regions.

689

### 690 ***Conservation***

691 We used hg38 phyloP 7way sourced from UCSC for phyloP conservation scores  
692 (<https://hgdownload.cse.ucsc.edu/goldenPath/hg38/phyloP7way/>). Phastcon scores  
693 were computed using the hg38 30way file  
694 (<https://hgdownload.soe.ucsc.edu/goldenPath/hg38/phastCons30way/>)

695

### 696 ***Alternative exons***

697 We used the database HEXEvent to select alternative exons. We downloaded the data  
698 using the associated website (<http://hexevent.mmg.uci.edu/HEXEvent/example.html>) for  
699 all chromosomes and those exons of type cassette or constitutive. We called an exon  
700 alternative if the constitLevel score was  $\leq 0.9$ . We called exons constitutive if the  
701 constitLevel score  $\geq 0.99$ .

702

### 703 ***Repetitive elements***

704 We used Dfam v3.4 repeat annotations. We used both the coordinates of the  
705 consensus models (Dfam.embl) and tallied exons overlapping the Dfam genome hits  
706 (hg38\_dfam.nrph.hits) along the coordinates of the consensus models. We analyzed  
707 repeats overlapping an exon with a read count of 100 or greater.

708 We calculated log p-values for the enrichment of exons in each repeat using the  
709 hypergeometric test, computed using the Python SciPy function logpmf (with  $M$ =number  
710 of stranded genome bases,  $n$ =number of exon bases,  $N$ =the number of genome bases  
711 for the repeat family tested,  $k$ =the number of exon bases overlapping the repeat family  
712 tested).

713

### 714 ***Exonic splicing enhancer database***

715 We used the hexmer sequences splicing enhancers (Ke et al. 2011), which yielded  
716 1,182 hexmer exonic splicing enhancers out of a possible 4,096 hexmers. These  
717 splicing enhancers sequences were used to count the occurrences of ESEs in the  
718 queried exon sequences. ESE counts were normalized to the number of ESE per 100  
719 bp of exon sequence.

720

### 721 ***SpliceAI scoring***

722 We obtained the SpliceAI models by installing the Python SpliceAI module and locating  
723 the trained model files. The models were loaded using the tensorflow Keras submodule.  
724 We obtained the average score of the 5 models for the query DNA sequence. The 5-  
725 model average score was used as the SpliceAI score for a given exon. When  
726 comparing mRNA, lncRNA, and intergenic exon SpliceAI scores, the scored sequence  
727 was the genomic sequence from 500 bp upstream and downstream of the exons 3'SS  
728 and 5'SS for up to 10,000 exons. This sequence was padded with 'N' characters for the  
729 remainder of the scoring window. Only exons that do not overlap repetitive elements  
730 were scored. Otherwise, the 11 kb sequence centered at the middle of the query exon  
731 was scored using SpliceAI when performing exon finding. Cutoffs were chosen to  
732 capture a majority of mRNA exons.

733

734 ***Maxentscan scoring***

735 We computed MaxEntScan scores using the Python module maxentpy, a Python  
736 wrapper for the original MaxEntScan implementation. Cutoffs were chosen to capture a  
737 majority of mRNA exons.

738

739 ***Exon Finding on Chromosome 17***

740 We used MaxEntScan and SpliceAI scores to call genome sequences as exons. We  
741 computed splice site scores across every base of Chromosome 17 using MaxEntScan  
742 and SpliceAI. We used a threshold and checked if the 5'SS score is above this  
743 threshold. We then checked if a 3'SS score is also above this threshold within an exon  
744 'length' of the selected 3'SS. An exon 'length' corresponds to between 63 and 222 bp,  
745 the 10% and 90% percentiles of GENCODE mRNA internal exons.

746 For MaxEntScan we required both splice sites to have a score of at least 6. For SpliceAI  
747 used a threshold of 0.2. We used this threshold since exons with a score delta greater  
748 than 0.8 were considered unlikely to splice in the publication associated with SpliceAI.  
749 We scored 1000 bp fragments padded with N characters to fit in the 10 kb SpliceAI  
750 model.

751

752 ***Shuffled Chromosome 17***

753 We created a shuffled chromosome preserving dinucleotide frequencies using the fasta-  
754 shuffle-letters package from the MEME suite with parameters -kmer 2 -seed 1337. We  
755 ran the command on 20 million bp, sequentially, (5 runs) across Chromosome 17 and  
756 joined the output into a single chromosome 'chr17\_shuffle' FASTA file.

757

758 ***Snaptron exon database***

759 We used the Snaptron SQL database SRAv2  
760 (<http://snaptron.cs.jhu.edu/srav2/snaptron?regions=ABCD3>) which contains junctions  
761 obtained from ~21 thousand samples from the Sequence Read Archive (SRA) using the  
762 reference hg38.

763

764 ***HEK293 database***

765 HEK293 gene expression data was obtained from supplemental data for (Nieborak et al.  
766 2023)  
767 ([https://ftp.ncbi.nlm.nih.gov/geo/series/GSE235nnn/GSE235387/suppl/GSE235387\\_HEK293.xlsx](https://ftp.ncbi.nlm.nih.gov/geo/series/GSE235nnn/GSE235387/suppl/GSE235387_HEK293.xlsx)) which contains data for HEK293 wild type control samples. HEK293  
768 transcript RPKM values are from the column 'HEK293-WT-S20190326-S20190326'.  
769

770

### 771 **HEK293 percent spliced in**

772 Percent spliced in numbers were obtained from the supplemental data of (Ellis et al.  
773 2023)

774 ([https://ftp.ncbi.nlm.nih.gov/geo/series/GSE221nnn/GSE221838/suppl/GSE221838\\_d0\\_v\\_d1000\\_SE.MATS.JCEC.txt.gz](https://ftp.ncbi.nlm.nih.gov/geo/series/GSE221nnn/GSE221838/suppl/GSE221838_d0_v_d1000_SE.MATS.JCEC.txt.gz)).

776 'IncLevel1'. Splicing inclusion levels containing a 'NA' were skipped and PSI values  
777 were calculated as the mean of the provided inclusion percents list. Splicing inclusion  
778 levels were associated with exon trapping data by finding exons that overlap the splicing  
779 inclusion exon list.

780

### 781 **Housekeeping gene list**

782 A list of housekeeping genes was obtained from (Hounkpe et al. 2021) supplemental file  
783 Supplementary\_Table1.xlsx.

784

### 785 **ISS frequency**

786 A list of sequencing displaying ISS activity was obtained from (Wen et al. 2010). The list  
787 of 5'SS and 3'SS ISS sequences were combined. The frequency of ISS sequences  
788 were computed around the 3'SS and 5'SS for the indicated exon dataset and the curves  
789 were smoothed with a 19 bp moving average. Exons overlapping repeats were not  
790 included datasets.

791

### 792 **Nonsense Mediated Decay**

793 Exon used for Nonsense Mediated Decay analysis consists of Internal mRNA exons  
794 with length divisible by 3 and nonzero read counts in the given reading frames. The ratio  
795 of NMD is calculated for an exon as follows for two given reading frames: the numerator  
796 lacks a stop codon in the first given reading frame vector and the denominator must  
797 have 1 or more stop codons in the second given reading frame vector. Positive values  
798 of the  $\log_2$  of this number indicate the fold increase of the exon abundance in the  
799 reading frame lacking a stop codon over the reading frame with a stop codon.

800

### 801 **DATA ACCESS**

802 All raw and processed sequencing data generated in this study have been submitted to  
803 the NCBI Gene Expression Omnibus (GEO; <https://www.ncbi.nlm.nih.gov/geo/>) under  
804 accession number GSE213006.

805 Supplemental data corresponding to figures is available as supplemental file  
806 **Supplemental\_data.zip** and at [https://hugheslab.cabr.utoronto.ca/supplementary-](https://hugheslab.cabr.utoronto.ca/supplementary-data/ExonTrappingGenome)  
807 [data/ExonTrappingGenome](https://hugheslab.cabr.utoronto.ca/supplementary-data/ExonTrappingGenome). Supplemental figures are available as supplemental file  
808 **Supplemental\_figures.pdf**. All code to analyze the data is available as supplemental  
809 file **Supplemental\_code.zip** and on Github at  
810 <https://github.com/nstep2/ExonTrapGenome>.

811

## 812 **COMPETING INTEREST STATEMENT**

813 The authors declare no competing interests.

814

## 815 **ACKNOWLEDGEMENTS**

816 We thank Debashish Ray and Ben Blencowe for helpful discussions, Mihai Albu for IT  
817 support, and the Donnelly Sequencing Centre (University of Toronto) for technical  
818 assistance. This work was supported by CIHR grants PJT-162255 and FDN-148403,  
819 and NIH grant R01HG008613, to TRH. TRH holds the Billes Chair of Medical Research  
820 at the University of Toronto, and a CIHR Canada Research Chair.

821

822 Author contributions: NS and TRH conceived of the study, analyzed the results, and  
823 wrote the paper. NS performed the experiments with assistance from AY. NS performed  
824 the computational data analyses, generated the figures, and orchestrated the database  
825 depositions.

826

827

828

829 **Figure 1. Overview of genome exon trapping method**

830 Diagram depicting exon trapping approach, sequencing library construction and  
 831 example sequencing read maps. Sheared genomic DNA library fragments (blue boxes)  
 832 500-1000 bp in length are cloned into the middle of the 6<sup>th</sup> intron from *TRA2B* (black  
 833 boxes), in a pcDNA 3.1 vector backbone. First and terminal exons (grey boxes) are  
 834 labeled with the transcriptional start site (TSS), start codon (ATG), stop codon (Stop),  
 835 and the cleavage and polyadenylation site (CPA). Internal exons (red boxes) are  
 836 amplified by RT-PCR, using indicated primers, then sequenced and mapped to the  
 837 human genome (hg38). Bottom panel shows mapped sequencing read counts  
 838 (separated into forward and reverse strand pileups) for regions containing KIAA0513  
 839 and a portion of *CIBAR2* (display region coordinates: Chr16:85,062,938-85,134,585).  
 840 The zoom-in region corresponds to exon 7 of *CIBAR2*.

841

842

843 **Figure 2. Properties of trapped exons.**

- 844 A) Histograms of sequencing read counts for trapped internal exons within different  
 845 genomic regions. Outset plot shows logarithmic exon counts and inset shows  
 846 zoomed linear exon counts. Logarithmic bin boundaries indicated by dots  
 847 corresponding to  $10^x$  for  $x$  from 0 to 5 with step size 0.5. Linear bins boundaries  
 848 range from 100 to 100,000 with a step size of 1,000.
- 849 B) Bar plots depicting sequencing read counts of internal exons containing zero or at  
 850 least one in-frame stop codon. Results for reading frames in the first ("Frame 0"),  
 851 second ("Frame 1"), and third ("Frame 2") positions are shown.
- 852 C) Line plots depicting the distribution of mRNA exon lengths and the percent of mRNA  
 853 exons at each length recovered by exon trapping. Plots have a 9 bp smoothing  
 854 window applied.
- 855 D) Boxplots depicting GC content of trapped internal exons from different genomic  
 856 regions. Y-axis indicates sequencing read counts of trapped exons, within indicated  
 857 GC content ranges (x-axis). Whiskers indicate 10<sup>th</sup> and 90<sup>th</sup> percentiles.

858

859

860 **Figure 3. Trapped exons found in different categories of genomic region**

- 861 A) Pie chart depicting the proportions of exons from different genomic regions. The  
 862 percentage of exons for indicated genomic regions relative to the total read count is  
 863 indicated.
- 864 B) Bar plots showing the percent of genomic bases in a trapped exon for different  
 865 categories of genomic region.

- 866 C) Boxplot depicting trapped exon sequencing read counts for different genomic  
867 regions. Whiskers depict 10<sup>th</sup> and 90<sup>th</sup> percentiles.
- 868 D) Bar plots depicting the exon counts at different read count bins (left) or average  
869 PhastCons score bins (right) for different non-GENCODE v37 exon annotations.  
870 Exon counts are shown with logarithmic scale. Non-GENCODE v37 exon  
871 annotations, Snaptron database, and exons found by this exon trapping study are  
872 displayed. Average PhastCons scores were calculated using the sequence of the  
873 exons.
- 874 E) Bar plot depicting PhastCons (30-way) scores for +/- 200 bp around an unannotated  
875 sense intronic exon found by exon trapping in gene *ITSN1*.
- 876 F) Proportion of annotated exons recovered from various databases. Blue bars indicate  
877 trapped exons that are annotated in GENCODE mRNA/lncRNAs. Trapped exons  
878 annotated in other lncRNA databases are shown, with annotated GENCODE  
879 lncRNAs removed.

880

881

#### 882 **Figure 4. Known splicing signals correlations with exon read counts**

- 883 A) Line plots depicting 3' MaxEntScan scores for trapped exons from different genomic  
884 regions. Exons from mRNA, lncRNA and Intergenic regions are indicated and are  
885 binned by logarithmic read counts. Median values are displayed as lines with shaded  
886 region corresponding to 25<sup>th</sup>-75<sup>th</sup> percentiles.
- 887 B) Same as 4A, above, but depicting 5' MaxEntScan scores.
- 888 C) Line plots representing Splicing Enhancer (ESE) counts for trapped exon sequences  
889 from different genomic regions. ESE median values are displayed, and exons are  
890 binned by their logarithmic sequencing read counts using logarithmic bins ranging  
891 from 100 to 10,000.
- 892 D) Bar plots depicting the median ESE counts for trapped exons (blue bars) and nearby  
893 sequence of the same length offset by 250 bp (orange bars) for different genomic  
894 regions. Offset sequences are the same length as the associated exon and  
895 corresponds to coordinates 250 bp upstream for reverse strand exons and 250 bp  
896 downstream for forward strand exons. For forward strand exons this is downstream  
897 the exon and for reverse strand exons this is upstream the exon. Range lines  
898 indicate 25<sup>th</sup>-75<sup>th</sup> percentiles.
- 899 E) Scatter plot representing ESE count (y-axis) vs median sequencing read count (x-  
900 axis) for trapped exons, subdivided by MaxEntScan scores into groups with weaker  
901 to stronger splice sites based on splice site score bin (point label). Splice site bins  
902 indicate that contained exons have both their 3'SS and 5'SS splice sites within the  
903 labeled MaxEntScan score boundaries, between values indicate by [n,m], where  
904 n=lower score and m=upper score.
- 905 F) Bar plot depicting fraction of intergenic exons that contain the ESE GAAGAA  
906 nucleotide sequence. Individual bars correspond to exons with both 3'SS and 5'SS

- 907 MaxEntScan scores (see Figure 4E, above) within the range given in the bar label  
 908 (e.g. [ $n < \text{splice site MaxEntScan score} < m$ ] for both 3' and 5'SS MaxEntScan  
 909 scores]).
- 910 G) Line plots depicting 3'SS SpliceAI scores for trapped exons in different genomic  
 911 regions. Values in the x-axis are logarithmic sequencing read counts using bins from  
 912 100 to 10,000 with 25 steps. For intergenic exons, the Spearman correlation  
 913 between SpliceAI scores and read counts is 0.31.
- 914 H) Same as Figure 4G, above, except for 5'SS SpliceAI scores. For intergenic exons,  
 915 the Spearman correlation between SpliceAI scores and read counts is 0.12.

916

917

### 918 **Figure 5. Overlaps between exons detected using different approaches**

- 919 A) Venn diagram depicting Chromosome 17 forward strand exons found by different  
 920 exon calling approaches. Exons are labeled as mRNA (annotated mRNA & lncRNA  
 921 internal exons), ET (exons found by exon trapping), MaxEntScan based on  
 922 MaxEntScan scoring, and SpliceAI based on SpliceAI scoring. Exon counts  
 923 corresponding to overlapping regions are indicated and are colored red linearly with  
 924 intensity determined by  $(\log_{10} \# \text{exons})$ .
- 925 B) Bar plot of median ESE counts for exons and offset sequences identified using  
 926 approaches listed in 5A. Offset sequences are the same length as the associated  
 927 exon and corresponds to coordinates 500 bp upstream for reverse strand exons and  
 928 500 bp downstream for forward strand exons. For forward strand exons this is  
 929 downstream the exon and for reverse strand exons this is upstream the exon. Range  
 930 lines indicate 25<sup>th</sup>-75<sup>th</sup> percentiles. Refer to **Figure 5A** for x-axis labels.
- 931 C) Bar plots showing the ratio of intronic to antisense exon counts found for the  
 932 different exon finder approaches.

933

934

### 935 **Figure 6. Exons overlapping repetitive elements**

- 936 A) Histogram depicts ratio of exon overlaps for different repeat families, relative to the  
 937 global genomic exon rate. Volcano plot shows  $\log_{10}$  p-values of repeat family exon  
 938 bases enrichment (hypergeometric test) vs the repeat exon enrichment relative to  
 939 the global exon rate. Red dots indicate top 5% of p-values. Blue dots indicate the  
 940 bottom 95% of repeat-values.
- 941 B) Histogram pileup depicting sequencing reads overlapping repeat instances of  
 942 DF000317.4 (5' end of L1 retrotransposon L1P2) in the human genome. Histogram  
 943 maps sequencing reads in the genome to the Dfam repeat consensus model.  
 944 MaxEntScan 3'SS and 5'SS scores Dfam are also shown across the Dfam repeat

945 consensus model with scores above 0 shown with colored bars. Repeat consensus  
946 coordinates start at 0.

947 C) Diagram depicting overlap between trapped exon 6 from *FAM228B*  
948 (Chr2:24,095,141-24,095,230) and a genome instance of L1P2 transposon  
949 DF0000317.4 (Chr2:24,095,155-24,097,397).

950 D) Diagram depicting overlap between all trapped exons and associated genomic  
951 repeat instances for L1P2 transposon DF0000317.4. Rows are sorted by repeat start  
952 and end coordinates for the Dfam repeat consensus model.

953

954

### 955 **Figure 7. Conservation around 3' and 5' splice sites of trapped exons**

956 A) Diagram depicting sequence conservation of trapped exons around 3'SS for different  
957 genomic regions. Shown are phyloP scores for the 100 bp region centered around  
958 the 3'SS for trapped exons from mRNA, lncRNA, and intergenic regions, as  
959 indicated.

960 B) Same as A) except phyloP scores were calculated for the 100 bp region centered  
961 around the 5'SS for trapped exons.

962

963

964

## 965 REFERENCES

- 966 Barash Y, Calarco JA, Gao W, Pan Q, Wang X, Shai O, Blencowe BJ, Frey BJ. 2010.  
967 Deciphering the splicing code. *Nature* **465**: 53-59.
- 968 Burge C, Karlin S. 1997. Prediction of complete gene structures in human genomic DNA. *J Mol*  
969 *Biol* **268**: 78-94.
- 970 Cartegni L, Wang J, Zhu Z, Zhang MQ, Krainer AR. 2003. ESEfinder: A web resource to identify  
971 exonic splicing enhancers. *Nucleic Acids Res* **31**: 3568-3571.
- 972 Consortium IHGS. 2001. Initial sequencing and analysis of the human genome. *Nature* **409**:  
973 860-921.
- 974 Core LJ, Martins AL, Danko CG, Waters CT, Siepel A, Lis JT. 2014. Analysis of nascent RNA  
975 identifies a unified architecture of initiation regions at mammalian promoters and  
976 enhancers. *Nat Genet* **46**: 1311-1320.
- 977 Cote J, Dupuis S, Jiang Z, Wu JY. 2001. Caspase-2 pre-mRNA alternative splicing:  
978 Identification of an intronic element containing a decoy 3' acceptor site. *Proc Natl Acad*  
979 *Sci U S A* **98**: 938-943.
- 980 Deveson IW, Brunck ME, Blackburn J, Tseng E, Hon T, Clark TA, Clark MB, Crawford J, Dinger  
981 ME, Nielsen LK et al. 2018. Universal Alternative Splicing of Noncoding Exons. *Cell Syst*  
982 **6**: 245-255 e245.
- 983 Dewannieux M, Heidmann T. 2005. Role of poly(A) tail length in Alu retrotransposition.  
984 *Genomics* **86**: 378-381.
- 985 Drexler HL, Choquet K, Churchman LS. 2020. Splicing Kinetics and Coordination Revealed by  
986 Direct Nascent RNA Sequencing through Nanopores. *Mol Cell* **77**: 985-998 e988.
- 987 Duyk GM, Kim SW, Myers RM, Cox DR. 1990. Exon trapping: a genetic screen to identify  
988 candidate transcribed sequences in cloned mammalian genomic DNA. *Proc Natl Acad*  
989 *Sci U S A* **87**: 8995-8999.
- 990 Ellis JA, Hale MA, Cleary JD, Wang ET, Andrew Berglund J. 2023. Alternative Splicing  
991 Outcomes Across an RNA-Binding Protein Concentration Gradient. *J Mol Biol* **435**:  
992 168156.
- 993 Engreitz JM, Haines JE, Perez EM, Munson G, Chen J, Kane M, McDonel PE, Guttman M,  
994 Lander ES. 2016. Local regulation of gene expression by lncRNA promoters,  
995 transcription and splicing. *Nature* **539**: 452-455.
- 996 Fairbrother WG, Yeh RF, Sharp PA, Burge CB. 2002. Predictive identification of exonic splicing  
997 enhancers in human genes. *Science* **297**: 1007-1013.
- 998 Fairbrother WG, Yeo GW, Yeh R, Goldstein P, Mawson M, Sharp PA, Burge CB. 2004.  
999 RESCUE-ESE identifies candidate exonic splicing enhancers in vertebrate exons.  
1000 *Nucleic Acids Res* **32**: W187-190.
- 1001 Frankish A, Diekhans M, Ferreira AM, Johnson R, Jungreis I, Loveland J, Mudge JM, Sisu C,  
1002 Wright J, Armstrong J et al. 2019. GENCODE reference annotation for the human and  
1003 mouse genomes. *Nucleic Acids Res* **47**: D766-D773.
- 1004 Hollander D, Naftelberg S, Lev-Maor G, Kornblihtt AR, Ast G. 2016. How Are Short Exons  
1005 Flanked by Long Introns Defined and Committed to Splicing? *Trends Genet* **32**: 596-606.
- 1006 Hounkpe BW, Chenou F, de Lima F, De Paula EV. 2021. HRT Atlas v1.0 database: redefining  
1007 human and mouse housekeeping genes and candidate reference transcripts by mining  
1008 massive RNA-seq datasets. *Nucleic Acids Res* **49**: D947-D955.
- 1009 Jaganathan K, Kyriazopoulou Panagiotopoulou S, McRae JF, Darbandi SF, Knowles D, Li YI,  
1010 Kosmicki JA, Arbelaez J, Cui W, Schwartz GB et al. 2019. Predicting Splicing from  
1011 Primary Sequence with Deep Learning. *Cell* **176**: 535-548 e524.
- 1012 Ke S, Shang S, Kalachikov SM, Morozova I, Yu L, Russo JJ, Ju J, Chasin LA. 2011.  
1013 Quantitative evaluation of all hexamers as exonic splicing elements. *Genome Res* **21**:  
1014 1360-1374.

- 1015 Kelley D, Rinn J. 2012. Transposable elements reveal a stem cell-specific class of long  
 1016 noncoding RNAs. *Genome Biol* **13**: R107.
- 1017 Kim D, Langmead B, Salzberg SL. 2015. HISAT: a fast spliced aligner with low memory  
 1018 requirements. *Nat Methods* **12**: 357-360.
- 1019 Kutter C, Watt S, Stefflova K, Wilson MD, Goncalves A, Ponting CP, Odom DT, Marques AC.  
 1020 2012. Rapid turnover of long noncoding RNAs and the evolution of gene expression.  
 1021 *PLoS Genet* **8**: e1002841.
- 1022 Lagarde J, Uszczynska-Ratajczak B, Carbonell S, Perez-Lluch S, Abad A, Davis C, Gingeras  
 1023 TR, Frankish A, Harrow J, Guigo R et al. 2017. High-throughput annotation of full-length  
 1024 long noncoding RNAs with capture long-read sequencing. *Nat Genet* **49**: 1731-1740.
- 1025 Le Hir H, Nott A, Moore MJ. 2003. How introns influence and enhance eukaryotic gene  
 1026 expression. *Trends Biochem Sci* **28**: 215-220.
- 1027 Li D, Yang MQ. 2017. Identification and characterization of conserved lncRNAs in human and  
 1028 rat brain. *BMC Bioinformatics* **18**: 489.
- 1029 Liu HX, Zhang M, Krainer AR. 1998. Identification of functional exonic splicing enhancer motifs  
 1030 recognized by individual SR proteins. *Genes Dev* **12**: 1998-2012.
- 1031 Makalowski W, Mitchell GA, Labuda D. 1994. Alu sequences in the coding regions of mRNA: a  
 1032 source of protein variability. *Trends Genet* **10**: 188-193.
- 1033 Maquat LE. 2004. Nonsense-mediated mRNA decay: splicing, translation and mRNP dynamics.  
 1034 *Nat Rev Mol Cell Biol* **5**: 89-99.
- 1035 Nieborak A, Lukauskas S, Capellades J, Heyn P, Santos GS, Motzler K, Zeigerer A, Bester R,  
 1036 Protzer U, Schelter F et al. 2023. Depletion of pyruvate kinase (PK) activity causes  
 1037 glycolytic intermediate imbalances and reveals a PK-TXNIP regulatory axis. *Mol Metab*  
 1038 **74**: 101748.
- 1039 Orom UA, Derrien T, Guigo R, Shiekhattar R. 2010. Long noncoding RNAs as enhancers of  
 1040 gene expression. *Cold Spring Harbor symposia on quantitative biology* **75**: 325-331.
- 1041 Pan Q, Shai O, Lee LJ, Frey BJ, Blencowe BJ. 2008. Deep surveying of alternative splicing  
 1042 complexity in the human transcriptome by high-throughput sequencing. *Nat Genet* **40**:  
 1043 1413-1415.
- 1044 Pertea M, Shumate A, Pertea G, Varabyou A, Breitwieser FP, Chang YC, Madugundu AK,  
 1045 Pandey A, Salzberg SL. 2018. CHESSE: a new human gene catalog curated from  
 1046 thousands of large-scale RNA sequencing experiments reveals extensive transcriptional  
 1047 noise. *Genome Biol* **19**: 208.
- 1048 Petersen USS, Doktor TK, Andresen BS. 2022. Pseudoexon activation in disease by non-splice  
 1049 site deep intronic sequence variation - wild type pseudoexons constitute high-risk sites in  
 1050 the human genome. *Hum Mutat* **43**: 103-127.
- 1051 Piovesan A, Antonaros F, Vitale L, Strippoli P, Pelleri MC, Caracausi M. 2019. Human protein-  
 1052 coding genes and gene feature statistics in 2019. *BMC Res Notes* **12**: 315.
- 1053 Pollard KS, Hubisz MJ, Rosenbloom KR, Siepel A. 2010. Detection of nonneutral substitution  
 1054 rates on mammalian phylogenies. *Genome Res* **20**: 110-121.
- 1055 Ponting CP, Haerty W. 2022. Genome-Wide Analysis of Human Long Noncoding RNAs: A  
 1056 Provocative Review. *Annu Rev Genomics Hum Genet* doi:10.1146/annurev-genom-  
 1057 112921-123710.
- 1058 Reed R. 2000. Mechanisms of fidelity in pre-mRNA splicing. *Curr Opin Cell Biol* **12**: 340-345.
- 1059 Robberson BL, Cote GJ, Berget SM. 1990. Exon definition may facilitate splice site selection in  
 1060 RNAs with multiple exons. *Mol Cell Biol* **10**: 84-94.
- 1061 Rosenberg AB, Patwardhan RP, Shendure J, Seelig G. 2015. Learning the sequence  
 1062 determinants of alternative splicing from millions of random sequences. *Cell* **163**: 698-  
 1063 711.

- 1064 Scalzitti N, Jeannin-Girardon A, Collet P, Poch O, Thompson JD. 2020. A benchmark study of  
 1065 ab initio gene prediction methods in diverse eukaryotic organisms. *BMC Genomics* **21**:  
 1066 293.
- 1067 Sorek R. 2007. The birth of new exons: mechanisms and evolutionary consequences. *RNA* **13**:  
 1068 1603-1608.
- 1069 Storer J, Hubley R, Rosen J, Wheeler TJ, Smit AF. 2021. The Dfam community resource of  
 1070 transposable element families, sequence models, and genome annotations. *Mobile DNA*  
 1071 **12**: 2.
- 1072 Sun H, Chasin LA. 2000. Multiple splicing defects in an intronic false exon. *Mol Cell Biol* **20**:  
 1073 6414-6425.
- 1074 The RC, Petrov AI, Kay SJE, Kalvari I, Howe KL, Gray KA, Bruford EA, Kersey PJ, Cochrane G,  
 1075 Finn RD et al. 2017. RNAcentral: a comprehensive database of non-coding RNA  
 1076 sequences. *Nucleic Acids Res* **45**: D128-D134.
- 1077 Tillo D, Hughes TR. 2009. G+C content dominates intrinsic nucleosome occupancy. *BMC*  
 1078 *Bioinformatics* **10**: 442.
- 1079 Volders PJ, Anckaert J, Verheggen K, Nuytens J, Martens L, Mestdagh P, Vandesompele J.  
 1080 2019. LNCipedia 5: towards a reference set of human long non-coding RNAs. *Nucleic*  
 1081 *Acids Res* **47**: D135-D139.
- 1082 Wang ET, Sandberg R, Luo S, Khrebtkova I, Zhang L, Mayr C, Kingsmore SF, Schroth GP,  
 1083 Burge CB. 2008. Alternative isoform regulation in human tissue transcriptomes. *Nature*  
 1084 **456**: 470-476.
- 1085 Wang J, Zhang J, Zheng H, Li J, Liu D, Li H, Samudrala R, Yu J, Wong GK. 2004a. Mouse  
 1086 transcriptome: neutral evolution of 'non-coding' complementary DNAs. *Nature* **431**: 1 p  
 1087 following 757; discussion following 757.
- 1088 Wang Y, Ma M, Xiao X, Wang Z. 2012. Intronic splicing enhancers, cognate splicing factors and  
 1089 context-dependent regulation rules. *Nat Struct Mol Biol* **19**: 1044-1052.
- 1090 Wang Z, Rolish ME, Yeo G, Tung V, Mawson M, Burge CB. 2004b. Systematic identification  
 1091 and analysis of exonic splicing silencers. *Cell* **119**: 831-845.
- 1092 Wen J, Chiba A, Cai X. 2010. Computational identification of tissue-specific alternative splicing  
 1093 elements in mouse genes from RNA-Seq. *Nucleic Acids Res* **38**: 7895-7907.
- 1094 Wilks C, Gaddipati P, Nellore A, Langmead B. 2018. Snaptron: querying splicing patterns  
 1095 across tens of thousands of RNA-seq samples. *Bioinformatics* **34**: 114-116.
- 1096 Xiong HY, Alipanahi B, Lee LJ, Bretschneider H, Merico D, Yuen RK, Hua Y, Gueroussov S,  
 1097 Najafabadi HS, Hughes TR et al. 2015. RNA splicing. The human splicing code reveals  
 1098 new insights into the genetic determinants of disease. *Science* **347**: 1254806.
- 1099 Yeo G, Burge CB. 2004. Maximum entropy modeling of short sequence motifs with applications  
 1100 to RNA splicing signals. *J Comput Biol* **11**: 377-394.
- 1101 Zhang XH, Kangsamaksin T, Chao MS, Banerjee JK, Chasin LA. 2005. Exon inclusion is  
 1102 dependent on predictable exonic splicing enhancers. *Mol Cell Biol* **25**: 7323-7332.
- 1103 Zhao L, Wang J, Li Y, Song T, Wu Y, Fang S, Bu D, Li H, Sun L, Pei D et al. 2021.  
 1104 NONCODEV6: an updated database dedicated to long non-coding RNA annotation in  
 1105 both animals and plants. *Nucleic Acids Res* **49**: D165-D171.

1106