

1 **A systematic assessment of the impact of rare canonical splice site variants on splicing**
2 **using functional and *in silico* methods**

3
4 **Authors:** Rachel Y. Oh^{1,2,*}, Ali AlMail^{2,3,*}, David Cheerie^{3,4}, George Guirguis^{3,4}, Huayun Hou³,
5 Kyoko E. Yuki^{3,5}, Bushra Haque^{3,4}, Bhooma Thiruvahindrapuram⁶, Christian R. Marshall^{5,7},
6 Roberto Mendoza-Londono^{1,3,8}, Adam Shlien^{3,4,5,7}, Lianna G Kyriakopoulou^{5,7}, Susan Walker⁶,
7 James J. Dowling^{3,4,8,9}, Michael D. Wilson^{3,4}, Gregory Costain^{1,3,4,8,#}

8
9 *Denotes equal contribution

10 #Author for correspondence

11
12 **Affiliations:**

13 ¹Division of Clinical and Metabolic Genetics, Hospital for Sick Children, Toronto, Canada

14 ²Temerty Faculty of Medicine, University of Toronto, Toronto, Canada

15 ³Program in Genetics and Genome Biology, SickKids Research Institute, Toronto, Canada

16 ⁴Department of Molecular Genetics, University of Toronto, Toronto, Canada

17 ⁵Division of Genome Diagnostics, Hospital for Sick Children, Toronto, Canada

18 ⁶The Centre for Applied Genomics, SickKids Research Institute, Toronto, Canada

19 ⁷Department of Laboratory Medicine and Pathobiology, University of Toronto, Toronto, Canada

20 ⁸Department of Paediatrics, University of Toronto, Toronto, Canada

21 ⁹Division of Neurology, Hospital for Sick Children, Toronto, Canada

22
23 **Contact information:**

24 Gregory Costain, MD, PhD

25 The Hospital for Sick Children 555 University Avenue Toronto, ON M5G1X8, CANADA

26 Phone: 416-813-7654 ext. 301480

27 E-mail: gregory.costain@sickkids.ca

28 **Main text word count: 2826**

29 **Abstract**

30 **Background/Objectives:** Canonical splice site variants (CSSVs) are often presumed to cause
31 loss-of-function (LoF) and are assigned very strong evidence of pathogenicity (according to
32 ACMG criterion PVS1). However, the exact nature and predictability of splicing effects of
33 unselected rare CSSVs in blood-expressed genes is poorly understood.

34 **Methods:** A total of 184 rare CSSVs in unselected blood-expressed genes were identified by
35 genome sequencing in 121 individuals, and their impact on splicing was interrogated manually in
36 RNA sequencing (RNA-seq) data. Blind to these RNA-seq data, we attempted to predict the
37 precise impact of CSSVs by applying *in silico* tools and the ClinGen Sequence Variant
38 Interpretation Working Group 2018 guidelines for applying PVS1 criterion.

39 **Results:** There was no evidence of a frameshift nor of reduced expression consistent with
40 nonsense-mediated decay (NMD) for 24% of CSSVs: 17% had wildtype splicing only and
41 normal junction depths, 3.25% resulted in cryptic splice site usage and in-frame indels, 3.25%
42 resulted in full exon skipping (in-frame), and 0.5% resulted in full intron inclusion (in-frame).
43 Misclassification rates for splicing outcome (frameshift/NMD vs. no frameshift/no NMD) using
44 (i) SpliceAI, (ii) MaxEntScan, and (iii) AutoPVS1 ranged from 30-41%, with none
45 outperforming a simple “zero rule” classifier.

46 **Conclusion:** Nearly 1 in 4 CSSVs may not cause LoF based on analysis of RNA-seq data.
47 Predictions from *in silico* methods were often discordant with findings from RNA-seq. More
48 caution may be warranted in applying PVS1-level evidence to CSSVs in the absence of
49 functional data.

50

51 **Introduction**

52 Canonical splice site variants (CSSVs) are DNA variants affecting splicing donor (+1 and +2)
53 and acceptor (-1 and -2) sites defining exon-intron boundaries^{1,2}. The consensus nucleotide
54 sequences at splicing donor and acceptor sites are GT and AG, respectively, and are essential in
55 interacting with the U2 spliceosome to result in normal splicing and generation of wildtype (WT)
56 transcripts²⁻⁵. CSSVs may modify the interactions between the precursor messenger RNA (pre-
57 mRNA) and spliceosome complex⁵⁻⁷. The resulting splice disruption events can include exon
58 skipping, full intron inclusion, and alternative use of nearby cryptic splice sites resulting in
59 insertions or deletions (indels) of nucleotides⁶⁻⁸. These effects may or may not induce a
60 frameshift and premature termination codon (PTC), which can then trigger nonsense-mediated
61 RNA decay (NMD) and result in a loss-of-function (LoF) of the gene^{9,10}.

62
63 Accurate variant interpretation is foundational to both genome diagnostics and screening¹¹⁻¹³.
64 Rare CSSVs are typically considered under the “null variant” code and assigned very strong
65 evidence for pathogenicity (PVS1)¹². The PVS1 guidelines were refined by the ClinGen
66 Sequence Variant Interpretation Working Group (ClinGen SVI) in 2018 and again recently in
67 2023^{11,14}. ClinGen SVI initially recommended assigning PVS1 at varying evidence strengths
68 (i.e., supporting, moderate, strong, very strong, or not at all) after directly inspecting the genomic
69 region to predict the impact of the CSSV on splicing and the overall reading frame¹¹. SpliceAI
70 subsequently emerged as a widely used and powerful method for annotating genetic variants
71 with their predicted effect on splicing^{15,16}. However, *in silico* tools continue to be recognized as a
72 valuable but imperfect adjunct method for predicting the impact of CSSVs¹⁴.

73

74 To our knowledge, there have been no systematic attempts to catalogue the precise consequences
75 of rare CSSVs on splicing using functional data. The degree to which the impact(s) of CSSVs are
76 predictable via inspection of the genomic region and application of *in silico* tools is also unclear.
77 Here, we analyzed all rare CSSVs identified by genome sequencing (GS) from children and
78 parents, in blood-expressed but otherwise unselected genes, using genome-wide RNA
79 sequencing (RNA-seq). We hypothesized that a significant minority of CSSVs would show no
80 evidence of a frameshift nor of reduced expression consistent with nonsense-mediated decay,
81 recognizing that $\sim 1/3$ of inserted or deleted DNA segments would be expected to have a size
82 divisible by 3¹⁷. We determined the proportions of various splicing outcomes, and we compared
83 the results to outcomes predicted in a blinded fashion by three *in silico* tools: SpliceAI,
84 MaxEntScan, and AutoPVS1^{12,13,15}. Our results revealed a previously underappreciated
85 complexity in CSSV impact prediction, and underscore the value of functional data in the
86 interpretation of rare CSSVs.

87 **Materials and Methods:**

88 *Identification of rare CSSVs expressed in blood*

89 In this study, we analyzed blood-derived GS and RNA-seq data from a convenience sample of
90 121 total individuals^{18,19}. The GS and RNA-seq studies were approved by the Research Ethics
91 Board at The Hospital for Sick Children. Parents and/or guardians provided written consent on
92 their child's behalf. Where appropriate, children provided written and/or oral assent.

93 Demographic details for the cohort are described in Table S1. Detailed GS methods, and a subset
94 of the GS data, were published previously¹⁸⁻²². DNA variants from GS files were filtered
95 according to the following criteria: (I) single nucleotide substitution in a canonical splice site
96 identified in a MANE or Ensembl Canonical transcript, (II) allele frequency of <0.05 (per 1000
97 Genomes, NHLBI-ESP, and ExAC), (III) $\geq 99\%$ Genotype Quality Score, and (IV) gene
98 possibly detected in whole blood (according to GTEx, V8, transcripts per million (TPM) >0.05)
99 and detected in our internal cohorts^{20,23}. Variants in untranslated regions (UTRs) were
100 excluded²³. All 184 CSSVs included in this study are listed in Table S2. Two of the identified
101 CSSVs were considered diagnostic for the phenotype(s) that prompted GS, and all CSSVs were
102 heterozygous^{18-22,25}.

103

104 *Analysis of splicing impact of CSSVs using RNA-seq data*

105 We analyzed the impact of CSSVs on splicing in the canonical transcript using the
106 accompanying short-read blood-derived RNA-seq data. RNA extraction, sequencing, and data
107 processing methods were previously described in full^{20,26}. Each splicing junction was manually
108 inspected using the Integrated Genome Visualizer (IGV) by two independent evaluators (R.Y.O.
109 and A.A.). For every CSSV, an average of five random, age-range matched "controls" (i.e., with

110 normal DNA sequence at the affected CSS) from this cohort were used to identify the WT
111 splicing event(s) and provide a reference on junction read depths to account for any possible
112 fluctuations. Sex-matched “controls” were specifically used for one CSSV located on the X
113 chromosome. Only junctions with >5 uniquely mapped reads, which is a low cut-off, were
114 considered in the analysis. The junction with the highest read depth was considered the
115 predominant splicing outcome. Splicing outcome categories are as follows:

- 116 (A) presumed NMD (if the raw WT splicing junction coverage was $\geq 20\%$ less than
117 control individuals but no aberrant splicing events were captured)²⁶,
- 118 (B) NMD not detected [if there was comparable ($\leq 20\%$ difference) junction depth
119 between individuals and no aberrant splicing events were captured],
- 120 (C) exon skipping leading to frameshift deletion,
- 121 (D) exon skipping leading to in-frame deletion,
- 122 (E) full intron inclusion leading to frameshift insertion,
- 123 (F) full intron inclusion leading to in-frame insertion,
- 124 (G) activation of a cryptic splice site leading to frameshift indel,
- 125 (H) activation of a cryptic splice site leading to in-frame indel.

126

127 *Prediction of splicing outcomes using in silico tools*

128 Blind to the RNA-seq data, two independent assessors (D.C. and G.G.) attempted to predict the
129 precise impact of CSSVs using *in silico* tools and ClinGen SVI recommendations¹¹. Standard
130 score thresholds were used for MaxEntScan (>3) and for SpliceAI (delta score >0.2). For each
131 CSSV, the splice junction was manually inspected using Alamut Visual Plus (version v1.7,
132 ©2022 SOPHiA GENETICS) with a ± 20 bp window for a cryptic splice site. In the case of two

133 or more cryptic splice sites in the same region with comparable *in silico* scores, if the use of any
134 one site was predicted to result in an in-frame indel then NMD was not predicted as the outcome.
135 In the absence of a cryptic splice site in the neighboring region, for acceptor losses we predicted
136 exon skipping and for donor losses we predicted full intron inclusion. The length of the exon or
137 intron, respectively, was then used to predict whether the impact would be in-frame or out-of-
138 frame. While the ClinGen SVI recommendations specify use of a +/-20 bp window, we
139 conducted additional exploratory analyses using an expanded SpliceAI window of +/-5000 bp
140 (“SpliceAI_expanded”). Concordance of outcomes was then calculated if the frameshift/non-
141 frameshift outcome predicted by the *in silico* tools matched the results from RNA-seq.
142
143 AutoPVS1, an automatic classification tool for PVS1 interpretation of null variants, was also
144 used to predict the impact of CSSVs on splicing²⁷. This algorithm uses Variant Effect Predictor
145 for annotation of variants and MaxEntScan to predict the use of cryptic splice sites and resulting
146 impact on splicing²⁷. The three possible outcomes of AutoPVS1 are (I) Exon skipping or cryptic
147 splice site usage that leads to a frameshift and NMD, (II) Exon skipping or cryptic splice site
148 usage that leads to a frameshift without NMD, and (III) Exon skipping or cryptic splice site
149 usage that preserves the reading frame.

150 **Results**

151 *Comparison of frameshift outcomes of CSSVs using RNA-seq vs. in silico predictions*

152 We assessed a total of 184 rare, blood expressed CSSVs in 180 otherwise unselected genes. By
153 RNA-seq, 24% of these CSSVs did not result in a frameshift nor in reduced expression
154 consistent with NMD (Figure 1a). There was no apparent difference in the patterns of variant
155 location and specific nucleotide substitution of CSSVs by frameshift/non-frameshift outcome
156 (Figures S1 and S2). Considering the n=18 CSSVs that showed only WT splicing and with
157 comparable read depth to controls (outcome category B), 14 CSSVs were in the donor splice site
158 including only 1 GT>GC variant²⁸ (Figure S3).

159

160 For the 184 CSSVs, blinded predictions from *in silico* methods for non-frameshift outcomes are
161 28% (SpliceAI), 29% (SpliceAI_expanded), to 43% (AutoPVS1) of all CSSVs (Figure 1a). For
162 CSSVs resulting in a frameshift per RNA-seq, SpliceAI_expanded had the highest pairwise
163 concordance (77%), followed by SpliceAI (74%), MaxENTScan (69%), and AutoPVS1 (61%;
164 Figure 1b). For CSSVs not resulting in a frameshift per RNA-seq, SpliceAI was concordant in
165 36%, MaxENTScan in 41%, SpliceAI_expanded in 48%, and AutoPVS1 in 55% (Figure 1b).
166 Across all 184 variants, SpliceAI_expanded had the highest pairwise concordance to RNA-seq
167 with respect to the frameshift vs. no-frameshift outcome, at 71%. In order to assess the
168 performance of each *in silico* method, we calculated misclassification rates from each technique,
169 which ranged from 30-41% (Figure 1c). In comparison, a zero-rule classifier (which would
170 predict that every CSSV causes frameshift/NMD, the predominant outcome category in this
171 study) had a misclassification rate of only 24% (Figure 1c).

172

173 *Comparison of specific splicing outcomes of CSSVs using RNA-seq vs. in silico predictions*
174 Next, we compared the specific splicing outcome of CSSVs (cryptic splice site usage, exon
175 skipping, or intron inclusion; categories C-H above) between RNA-seq and *in silico* predictions,
176 for the n=23 variants where this could be determined from RNA-seq (Figure 2a). Total pairwise
177 concordance was 70% (16/23) for SpliceAI_expanded (improved from 35% (8/23) for SpliceAI
178 and 26% (6/23) for MaxENTScan; AutoPVS1 does not provide such predictions. The
179 performance of both *in silico* methods appeared to vary by outcome category, e.g., with
180 SpliceAI_expanded correctly predicting all usage of cyptic splice sites (10/10), some of the exon
181 skipping events (5/10) and none of the intron inclusion events (0/3) (Figure 2b). Two selected
182 donor CSSVs are used to illustrate that *in silico* were often correct in predicting frameshift vs. no
183 frameshift outcomes, however, for incorrect and/or incomplete mechanisms of abnormal splicing
184 (Figure 3a-d).

185 **Discussion**

186 *Challenging common assumptions and interpretation heuristics for CSSVs*

187 As demonstrated in this study cohort, most human genomes contain one or more rare CSSVs in
188 blood-expressed genes. Although there is widespread recognition that a rare CSSV need not
189 necessarily result in LoF, experiences from decades of targeted diagnostics-focused clinical
190 genetic testing (with a resulting ascertainment bias) may have contributed to a misconception
191 that CSSVs are comparable to nonsense and frameshift variants. We present the first systematic
192 assessment of the impact of rare “unselected” CSSVs in blood-expressed genes, using RNA-seq.
193 We found that nearly 1 in 4 CSSVs may not cause LoF, and that *in silico* predictions using
194 established tools and published guidelines were often discordant with RNA-seq data.

195
196 In recent years, there have been numerous computational tools developed to predict the location
197 of novel splice sites and thus the impact of DNA variants on splicing^{12,15}. These tools were
198 validated using few if any CSSVs (e.g., near intronic variants at least 3 nucleotides from a
199 canonical exon boundary in SpliceAI-10k; n=55 in 300K-RNA Top-4) and were not created to
200 predict specific CSSV outcomes like exon skipping or intron inclusion^{15,16,29,30}. A prior report
201 found that intron inclusion was poorly predicted using SpliceAI when compared with
202 transcriptome sequencing data³¹. AutoPVS1 does not list intron inclusion as a specific outcome,
203 nor does it distinguish which variants result in exon skipping versus cryptic splice site usage²⁷.
204 For now, no *in silico* methods appear to predict the precise impact of CSSVs with the sensitivity
205 and specificity needed for clinical diagnostics. A recent study has proposed using the
206 ACMG/AMP PM4 (moderate evidence of pathogenicity) criterion for CSSVs that are predicted
207 to result in intron inclusion and 3 or more in-frame events as predicted by 300K-RNA Top-4¹⁶.

208 Relative weighting of biological function or evolutionary conservation of the affected gene
209 region(s) may need to be considered in addition to the length of the in-frame disruption¹⁶. We
210 propose, in addition, that future updates to published guidelines on the use of PVS1 should
211 consider the use of SpliceAI with an expanded window of +/-5000 bp.

212

213 The impact of CSSVs on splicing can be complex. Multiple and/or partial effects on splicing
214 have been observed in the past (e.g., with some aberrantly spliced transcripts resulting in LoF
215 and others showing no apparent impact or producing a functional transcript)²⁴. Surprisingly,
216 some CSSVs in our data showed no direct (aberrant/non-WT splice junctions) or indirect
217 (reduced read depth, compatible with NMD) impact on splicing. The ACMG/AMP rule, BS3,
218 may be applied to CSSVs with evidence of normal splicing patterns demonstrated by RNA-seq
219 fulfilling specific criteria^{24,32}. A recent study using cell culture-based full-length gene splicing
220 assays demonstrated that specific nucleotide substitutions (GT>GC in the donor splice site) can
221 generate WT transcript levels in an estimated 15-18% of cases; no other nucleotide substitutions
222 in the +2 donor splice site were able to generate WT transcripts²⁸. Moreover, WT splicing as a
223 result of the 5' splice site GT>GC substitutions was not accurately predicted by *in silico* tools²⁸.
224 Our results showed that WT transcripts can be generated with diverse nucleotide substitutions, in
225 no consistent ranking order in the 5' donor splice site as recently described (though this study
226 performed RNA-seq in fibroblast samples and also included common variants), affecting +/- 1
227 and 2 canonical splice sites, highlighting *in vitro* assays' inability to capture the full complexity
228 of splicing in human cells (Figure S3)³³. Of note, sensitivity of our RNA-seq methods for
229 detecting evidence of NMD will be <100%. Long-read RNA-seq might have facilitated more

230 robust quantitation of transcript isoforms and revealed additional splicing outcomes. Testing
231 additional tissues beyond blood was also outside the scope of this initial study.

232

233 Our study has several additional limitations. First, we recognize that in-frame insertions and
234 deletions can still result in non-functional proteins (e.g., through disruption of an essential
235 protein domain) and that protein function cannot be inferred completely from RNA-seq. In the
236 absence of any evidence-based guidance, we assumed for all *in silico* predictions that exon
237 skipping would be the typical impact of an acceptor site loss and that intron retention/inclusion
238 would be the typical impact of a donor site loss, whenever an alternative/cryptic splice site was
239 not present. There is growing appreciation that this reasoning is overly simplistic^{31,34}. The
240 landscape of impacts of rare CSSVs may change based on tissue and/or age. Last, we were
241 underpowered in this study to identify substitution- and location/motif-specific predictors of
242 splicing outcomes.

243

244 *Implications for diagnostics, predictive testing, and screening*

245 These data reinforce prior expert consensus recommendations that cautioned against applying
246 PVS1 to CSSVs in the absence of additional supportive evidence²⁴. Recent ClinGen SVI
247 recommendations emphasize the importance of annotating physiologically occurring alternative
248 splicing events (or “leaky” splicing events) as candidate rescue transcripts that may not result in
249 LoF, consideration of gene-specific details (e.g., the impact on critical regions), and generating
250 more data to guide the integration of various lines of evidence (i.e., computational and *in vitro*
251 assays to assess impact on splicing), with the end goal of improved navigation of equivocal
252 clinical scenarios requiring interpretation of rare CSSVs³⁴. The advantages and limitations of

253 RNA-seq, which can be done high-throughput, should be weighed against a targeted approach
254 like RT-PCR; the latter may have greater sensitivity in detecting mis-spliced reads of low read
255 depth as a result of low gene expression in whole blood, and being able to distinguish partial
256 from complete abnormal splicing³³. Although a majority of CSSVs result in LoF, this assumption
257 should be questioned when genome-wide sequencing identifies novel rare CSSVs in genes
258 associated with ultra-rare, poorly characterized conditions with non-specific phenotypes (e.g.,
259 developmental delays)³⁵. The nuances of CSSV interpretation take on added importance when
260 the pre-test probability is low and phenotypes are absent, as is the case with most secondary
261 findings and in newborn genomic screening programs. Last, confirmation of genetic diagnosis
262 and pathomechanism will be foundational to the coming wave of clinical trials of genetic
263 therapies³⁶.

264 **Declaration of interest**

265 S.W. is currently an employee of Genomics England Limited. The other authors declare no
266 competing interests.

267

268 **Acknowledgements**

269 We gratefully acknowledge all the individuals and their families who participated in this study.

270 We thank the many healthcare providers involved in the diagnosis and care of the study

271 participants. Special thanks to all staff affiliated with the Complex Care Program and The Centre

272 for Applied Genomics. M.D.W. was supported by the Canada Research Chairs Program.

273 Funding was provided by Genome Canada (OGI-158; M.D.W., A.S., and J.J.D.), the SickKids

274 Centre for Genetic Medicine and Translational Genomics Node, the Sickkids Research Institute,

275 the Canadian Institutes of Health Research (Funding Reference Number: PJT186240), and the

276 University of Toronto McLaughlin Centre.

277

278 **Web resources**

279 SpliceAI: <https://spliceailookup.broadinstitute.org>

280

281 **Data availability**

282 GS and RNA-seq data were published and shared previously. The *in silico* predictions can be

283 shared upon request to the corresponding author.

284

285 **Author contributions**

286 Conceptualization: R.Y.O., A.A., S.W., G.C.; Data curation: R.Y.O., A.A., D.C., G.G., B.H.,
287 G.C.; Formal analysis: R.Y.O., A.A., D.C., G.G., G.C.; Funding acquisition: C.R.M., R.M-L.,
288 A.S., L.G.K., J.J.D., M.D.W., G.C.; Investigation: R.Y.O., A.A., D.C., G.C.; Methodology:
289 R.Y.O., A.A., D.C., G.G., H.H., K.Y., S.W., G.C.; Project administration: R.Y.O., A.A., T.K.,
290 B.T., G.C.; Resources: C.R.M., R.M-L., A.S., L.G.K., J.J.D., M.D.W., G.C.; Software: R.Y.O.,
291 A.A., D.C., G.G., H.H., K.Y.; Visualization: R.Y.O., A.A., D.C., H.H.; Writing-original draft:
292 R.Y.O., A.A., G.C.; Writing-review & editing: all other authors

293 **References**

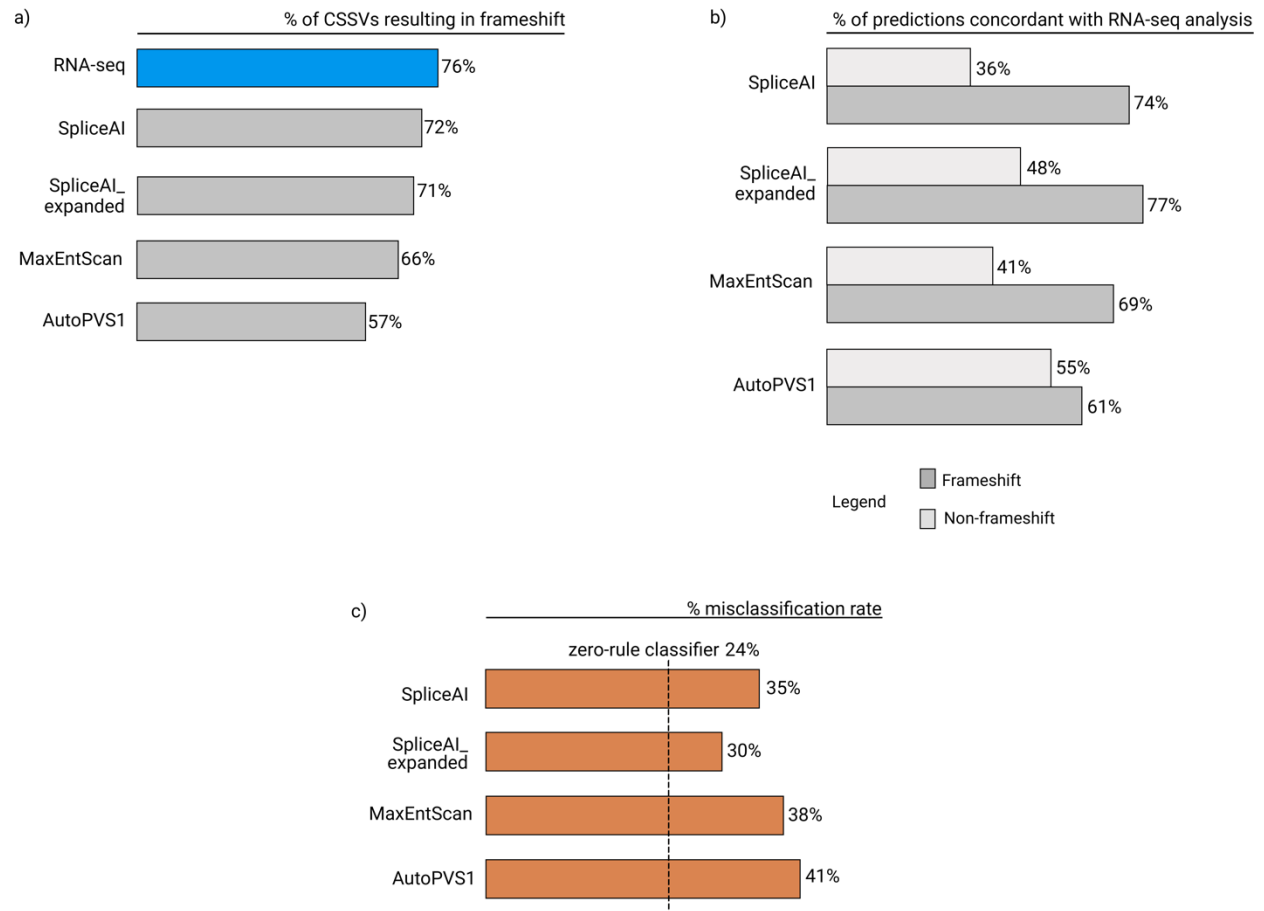
- 294 1. Krawczak M, Reiss J, Cooper DN. The mutational spectrum of single base-pair substitutions
295 in mRNA splice junctions of human genes: causes and consequences. *Hum Genet.* 1992;90(1-
296 2):41-54. doi:10.1007/BF00210743
- 297 2. Burset M, Seledtsov IA, Solovyev VV. Analysis of canonical and non-canonical splice sites in
298 mammalian genomes. *Nucleic Acids Res.* 2000;28(21):4364-4375.
299 doi:10.1093/nar/28.21.4364
- 300 3. Hastings ML, Krainer AR. Pre-mRNA splicing in the new millennium. *Curr Opin Cell Biol.*
301 2001;13(3):302-309. doi:10.1016/s0955-0674(00)00212-x
- 302 4. Matera AG, Wang Z. A day in the life of the spliceosome. *Nat Rev Mol Cell Biol.*
303 2014;15(2):108-121. doi:10.1038/nrm3742
- 304 5. Rogalska ME, Vivori C, Valcárcel J. Regulation of pre-mRNA splicing: roles in physiology
305 and disease, and therapeutic prospects. *Nat Rev Genet.* 2023;24(4):251-269.
306 doi:10.1038/s41576-022-00556-8
- 307 6. Krawczak M, Thomas NST, Hundrieser B, et al. Single base-pair substitutions in exon–intron
308 junctions of human genes: nature, distribution, and consequences for mRNA splicing. *Human*
309 *Mutation.* 2007;28(2):150-158. doi:10.1002/humu.20400
- 310 7. Ward AJ, Cooper TA. The pathobiology of splicing. *J Pathol.* 2010;220(2):152-163.
311 doi:10.1002/path.2649
- 312 8. Dufner-Almeida LG, do Carmo RT, Masotti C, Haddad LA. Chapter Two - Understanding
313 human DNA variants affecting pre-mRNA splicing in the NGS era. In: Kumar D, ed.
314 *Advances in Genetics.* Vol 103. Academic Press; 2019:39-90.
315 doi:10.1016/bs.adgen.2018.09.002
- 316 9. Fatscher T, Boehm V, Gehring NH. Mechanism, factors, and physiological role of nonsense-
317 mediated mRNA decay. *Cell Mol Life Sci.* 2015;72(23):4523-4544. doi:10.1007/s00018-015-
318 2017-9
- 319 10. Hug N, Longman D, Cáceres JF. Mechanism and regulation of the nonsense-mediated
320 decay pathway. *Nucleic Acids Res.* 2016;44(4):1483-1495. doi:10.1093/nar/gkw010
- 321 11. Abou Tayoun AN, Pesaran T, DiStefano MT, et al. Recommendations for interpreting the
322 loss of function PVS1 ACMG/AMP variant criterion. *Hum Mutat.* 2018;39(11):1517-1524.
323 doi:10.1002/humu.23626
- 324 12. Richards S, Aziz N, Bale S, et al. Standards and Guidelines for the Interpretation of
325 Sequence Variants: A Joint Consensus Recommendation of the American College of Medical
326 Genetics and Genomics and the Association for Molecular Pathology. *Genet Med.*
327 2015;17(5):405-424. doi:10.1038/gim.2015.30

- 328 13. Costain G, Cohn RD, Scherer SW, Marshall CR. Genome sequencing as a diagnostic test.
329 *CMAJ*. 2021;193(42):E1626-E1629. doi:10.1503/cmaj.210549
- 330 14. Walker LC, Hoya M de la, Wiggins GAR, Lindy A, Vincent LM, Parsons MT, et al.
331 Using the ACMG/AMP framework to capture evidence related to predicted and observed
332 impact on splicing: Recommendations from the ClinGen SVI Splicing Subgroup. *The*
333 *American Journal of Human Genetics* [Internet]. <https://doi.org/10.1016/j.ajhg.2023.06.002>
- 334
- 335 15. Jaganathan K, Kyriazopoulou Panagiotopoulou S, McRae JF, et al. Predicting Splicing
336 from Primary Sequence with Deep Learning. *Cell*. 2019;176(3):535-548.e24.
337 doi:10.1016/j.cell.2018.12.015
- 338 16. Dawes R, Bournazos AM, Bryen SJ, et al. SpliceVault predicts the precise nature of
339 variant-associated mis-splicing. *Nat Genet*. 2023;55(2):324-332. doi:10.1038/s41588-022-
340 01293-8
- 341 17. Resch A, Xing Y, Alekseyenko A, Modrek B, Lee C. Evidence for a subpopulation of
342 conserved alternative splicing events under selection pressure for protein reading frame
343 preservation. *Nucleic Acids Res*. 2004;32(4):1261-1269. doi:10.1093/nar/gkh284
- 344 18. Costain G, Walker S, Marano M, et al. Genome Sequencing as a Diagnostic Test in
345 Children With Unexplained Medical Complexity. *JAMA Network Open*. 2020;3(9):e2018109.
346 doi:10.1001/jamanetworkopen.2020.18109
- 347 19. Stavropoulos DJ, Merico D, Jobling R, et al. Whole Genome Sequencing Expands
348 Diagnostic Utility and Improves Clinical Management in Pediatric Medicine. *NPJ Genom*
349 *Med*. 2016;1:15012-. doi:10.1038/npjgenmed.2015.12
- 350 20. Deshwar A, Yuki K, Hou H, et al. Trio RNA sequencing in a cohort of medically
351 complex children. *The American Journal of Human Genetics*. Published online March 1,
352 2023. doi:10.1016/j.ajhg.2023.03.006
- 353 21. Lionel AC, Costain G, Monfared N, et al. Improved diagnostic yield compared with
354 targeted gene sequencing panels suggests a role for whole-genome sequencing as a first-tier
355 genetic test. *Genetics in Medicine*. 2018;20(4):435-443. doi:10.1038/gim.2017.119
- 356 22. Walker S, Lamoureux S, Khan T, et al. Genome sequencing for detection of pathogenic
357 deep intronic variation: A clinical case report illustrating opportunities and challenges. *Am J*
358 *Med Genet A*. 2021;185(10):3129-3135. doi:10.1002/ajmg.a.62389
- 359 23. GTEx Portal. Accessed May 14, 2023. <https://gtexportal.org/home/>
- 360 24. Ellingford JM, Ahn JW, Bagnall RD, et al. Recommendations for clinical interpretation
361 of variants found in non-coding regions of the genome. *Genome Medicine*. 2022;14(1):73.
362 doi:10.1186/s13073-022-01073-3

- 363 25. Costain G, Jobling R, Walker S, et al. Periodic reanalysis of whole-genome sequencing
364 data enhances the diagnostic advantage over standard clinical genetic testing. *Eur J Hum*
365 *Genet.* 2018;26(5):740-744. doi:10.1038/s41431-018-0114-6
- 366 26. Abstracts from the 54th European Society of Human Genetics (ESHG) Conference: Oral
367 Presentations. *Eur J Hum Genet.* 2022;30(1):3-87. doi:10.1038/s41431-021-01025-2
- 368 27. Xiang J, Peng J, Baxter S, Peng Z. AutoPVS1: An automatic classification tool for PVS1
369 interpretation of null variants. *Hum Mutat.* 2020;41(9):1488-1498. doi:10.1002/humu.24051
- 370 28. Lin JH, Tang XY, Boulling A, et al. First estimate of the scale of canonical 5' splice site
371 GT>GC variants capable of generating wild-type transcripts. *Hum Mutat.* 2019;40(10):1856-
372 1873. doi:10.1002/humu.23821
- 373 29. de Sainte Agathe JM, Filser M, Isidor B, et al. SpliceAI-visual: a free online tool to
374 improve SpliceAI splicing variant interpretation. *Human Genomics.* 2023;17(1):7.
375 doi:10.1186/s40246-023-00451-1
- 376 30. Palmisano A, Vural S, Zhao Y, Sonkin D. MutSpliceDB: A database of splice sites
377 variants with RNA-seq based evidence on effects on splicing. *Hum Mutat.* 2021;42(4):342-
378 345. doi:10.1002/humu.24185
- 379 31. Shiraishi Y, Okada A, Chiba K, et al. Systematic identification of intron retention
380 associated variants from massive publicly available transcriptome sequencing data. *Nat*
381 *Commun.* 2022;13(1):5357. doi:10.1038/s41467-022-32887-9
- 382 32. Bournazos AM, Riley LG, Bommireddipalli S, et al. Standardized practices for RNA
383 diagnostics using clinically accessible specimens reclassifies 75% of putative splicing
384 variants. *Genet Med.* 2022;24(1):130-145. doi:10.1016/j.gim.2021.09.001
- 385 33. Erkelenz S, Theiss S, Kaisers W, et al. Ranking noncanonical 5' splice site usage by
386 genome-wide RNA-seq analysis and splicing reporter assays. *Genome Res.* 2018;28(12):1826-
387 1840. doi:10.1101/gr.235861.118
- 388 34. Rivas MA, Pirinen M, Conrad DF, et al. Impact of predicted protein-truncating genetic
389 variants on the human transcriptome. *Science.* 2015;348(6235):666-669.
390 doi:10.1126/science.1261877
- 391 35. Sanders SJ, Schwartz GB, Farh KKH. Clinical impact of splicing in neurodevelopmental
392 disorders. *Genome Medicine.* 2020;12(1):36. doi:10.1186/s13073-020-00737-2

- 394 36. Peña LDM, et. al. Contributions from medical geneticists in clinical trials of genetic
395 therapies: A points to consider statement of the American College of Medical Genetics and
396 Genomics (ACMG). *Genetics in Medicine*. 2023 Jun;25(6):100831. doi:
397 10.1016/j.gim.2023.100831

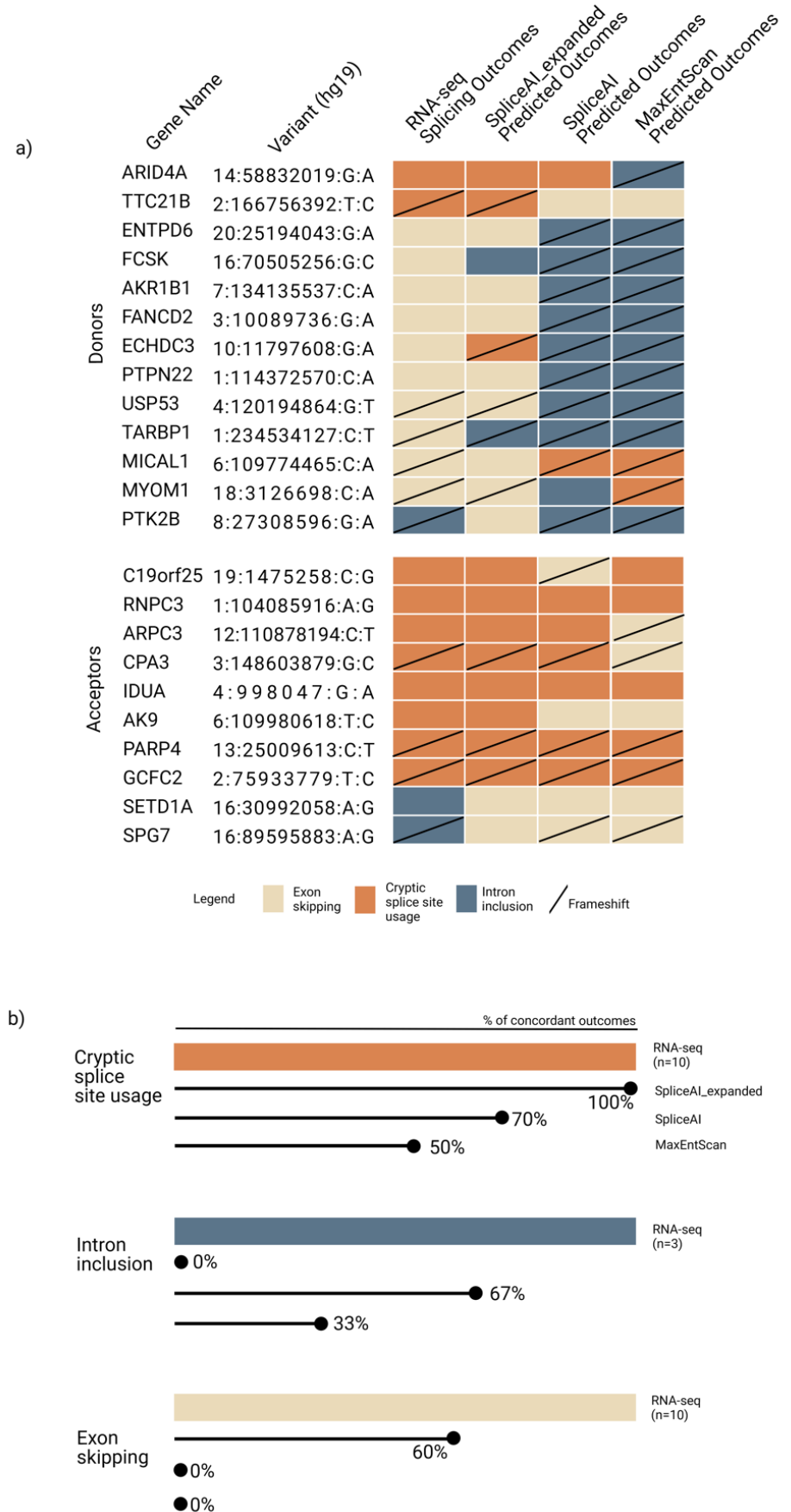
398 **Figures**



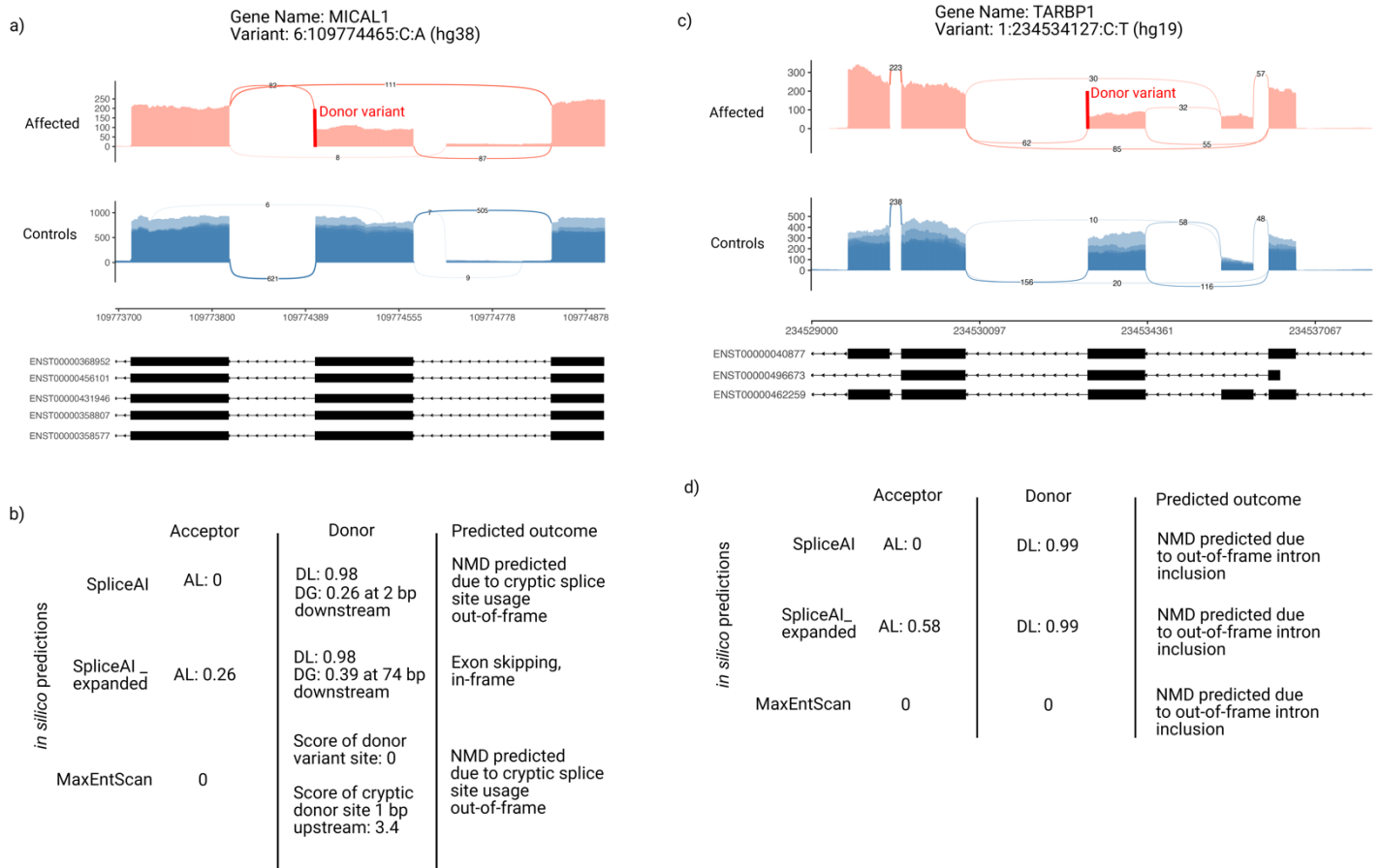
399 **Figure 1 a)** A comparison of the proportion of CSSVs resulting in a frameshift according to
400 RNA-seq analysis vs. *in silico* predictions. **b)** Concordance between RNA-seq and *in silico*
401 predictions of the impact of CSSVs on splicing. Concordant outcomes are defined as the RNA-
402 seq and respective *in silico* tool identifying the same outcome (frameshift or non-frameshift).
403 **c)** Misclassification rates of *in silico* tools compared to a zero-rule classifier.

Oh*, AlMail*, et al.

404 **Figure 2 a)** Selected variants
 405 (n=23) with specific splicing
 406 outcomes in RNA-seq (exon
 407 skipping, cryptic splice site
 408 usage, and intron inclusion)
 409 compared with splicing outcome
 410 predictions from *in silico* tools.
 411 Total pairwise concordance in
 412 specific splicing outcomes was
 413 70% (16/23) for
 414 SpliceAI_expanded, 35% (8/23)
 415 for SpliceAI and 26% (6/23) for
 416 MaxENTScan.
 417 **b)** Concordance of specific
 418 splicing outcomes in RNA-seq
 419 vs. *in silico* predictions.



420



425 **Figures 3)** Two examples of donor CSSVs showing discordant splicing events in blood RNA-
 426 seq vs. *in silico* predictions but overall correct outcome (frameshift vs. no frameshift). Splicing
 427 junctions with reads >5 are shown. First transcripts on the list are the canonical transcripts
 428 indicated by Ensembl. **a)** CSSV in *MICAL1* results in exon skipping leading to a frameshift in
 429 RNA-seq. **b)** SpliceAI, SpliceAI_expanded and MaxEntScan all predicted activation of a cryptic
 430 splice site, resulting in frameshift. **c)** CSSV in *TARBP1* results in exon skipping leading to a
 431 frameshift in RNA-seq. **d)** SpliceAI, SpliceAI_expanded and MaxEntScan all predicted intron
 432 inclusion, resulting in frameshift.