

Genome analysis

SpliceAI-10k calculator for the prediction of pseudoexonization, intron retention, and exon deletion

Daffodil M. Canson ^{1,2,*}, Aimee L. Davidson¹, Miguel de la Hoya ³,
Michael T. Parsons¹, Dylan M. Glubb^{2,4}, Olga Kondrashova⁴, Amanda B. Spurdle^{1,2}

¹Population Health Program, QIMR Berghofer Medical Research Institute, Herston, QLD 4006, Australia

²Faculty of Medicine, The University of Queensland, Herston, QLD 4006, Australia

³Molecular Oncology Laboratory CIBERONC, IdISSC (Instituto de Investigación Sanitaria del Hospital Clínico San Carlos), Hospital Clínico San Carlos, 28040 Madrid, Spain

⁴Cancer Research Program, QIMR Berghofer Medical Research Institute, Herston, QLD 4006, Australia

*Corresponding author. Population Health Program, QIMR Berghofer Medical Research Institute, Herston, QLD 4006, Australia.

E-mail: daffodil.canson@qimrberghofer.edu.au

Associate Editor: Christina Kendziorski

Received 15 December 2022; revised 27 February 2023; accepted 30 March 2023

Abstract

Summary: SpliceAI is a widely used splicing prediction tool and its most common application relies on the maximum delta score to assign variant impact on splicing. We developed the SpliceAI-10k calculator (SAI-10k-calc) to extend use of this tool to predict: the splicing aberration type including pseudoexonization, intron retention, partial exon deletion, and (multi)exon skipping using a 10 kb analysis window; the size of inserted or deleted sequence; the effect on reading frame; and the altered amino acid sequence. SAI-10k-calc has 95% sensitivity and 96% specificity for predicting variants that impact splicing, computed from a control dataset of 1212 single-nucleotide variants (SNVs) with curated splicing assay results. Notably, it has high performance ($\geq 84\%$ accuracy) for predicting pseudoexon and partial intron retention. The automated amino acid sequence prediction allows for efficient identification of variants that are expected to result in mRNA nonsense-mediated decay or translation of truncated proteins.

Availability and implementation: SAI-10k-calc is implemented in R (<https://github.com/adavi4/SAI-10k-calc>) and also available as a Microsoft Excel spreadsheet. Users can adjust the default thresholds to suit their target performance values.

1 Introduction

SpliceAI is a neural network that predicts splicing from a pre-mRNA sequence (Jaganathan et al. 2019). Previous evaluations (Ha et al. 2021; Moles-Fernández et al. 2021; Riepe et al. 2021; Rowlands et al. 2021; Wai et al. 2020) have identified SpliceAI as the best predictor of variants that impact splicing, here termed spliceogenic variants. These studies assessed single-nucleotide variants (SNVs) and small indels across multiple locations (i.e. splice site motifs, deep intronic regions >20 bp from the acceptor and >6 bp from the donor site, and exonic). They used the maximum delta score (of the four possible output scores) that passed the respective study-designated thresholds to predict variant spliceogenicity, but did not assess the splicing aberration type. SpliceAI sensitivity for detecting spliceogenic intronic variants >50 bp from exons was originally reported to be 41% using a 0.5 maximum delta score threshold (Jaganathan et al. 2019), but an improved sensitivity of 94% was observed for variants >20 bp from exons by lowering the

threshold to 0.05 (Moles-Fernández et al. 2021). Paired donor–acceptor splice site scores were observed for validated pseudoexonization events (Moles-Fernández et al. 2021). Moreover, manual checking of donor–acceptor splice site pairing was incorporated into a scheme to prioritize likely spliceogenic deep intronic variants (Qian et al. 2021).

We developed the SpliceAI-10k calculator (SAI-10k-calc) to systematically predict different SNV-associated splicing aberrations, altered transcript sizes, and consequent amino acid sequences, with a focus on accurate prediction of aberration sizes due to deep intronic variation.

2 Methods

SAI-10k-calc was designed to predict specific types of splicing aberrations, namely: pseudoexonization, partial intron retention, partial exon deletion, (multi)exon skipping, and whole intron retention.

Its features were derived from the application of all four raw delta scores and their corresponding delta positions generated by the SpliceAI tool (Jaganathan et al. 2019) using the maximum distance of ± 4999 bp flanking the variant of interest. SAI-10k-calc can process SpliceAI scores resulting from SNVs at any exonic or intronic position, but not scores resulting from indels due to the complexity of distance interpretations for such variants. The decision flowchart is shown in [Supplementary File S1](#).

We established default thresholds for SpliceAI delta scores (0.02–0.2 for exon skipping or whole intron retention and 0.02–0.05 for pseudoexon gain) and the gained exon size range of 25–500 bp based on two training sets derived from published splicing data: (i) SNVs in *BRCA1*, *BRCA2*, *MLH1*, *MSH2*, *MSH6*, and *PMS2* from Shamsani et al. (2019); and (ii) deep intronic SNVs in various Mendelian disease genes from Moles-Fernández et al. (2021) ([Supplementary Table S1](#)). The 0.2 upper threshold for exon skipping is based on the lower limit set by SpliceAI developers (Jaganathan et al. 2019). For deep intronic variants, the 0.05 upper threshold for pseudoexon gain is also supported by previous findings (Moles-Fernández et al. 2021). The 25–500 bp exon size range encompasses the optimal size for efficient splicing that is between 50 and 250 bp (Movassat et al. 2019) and is expected to capture most gained pseudoexons.

3 Usage and features

The R code (<https://www.R-project.org/>) version of SAI-10k-calc requires two input files: a SpliceAI output VCF file and a tab-separated file with gene names and RefSeq transcript IDs (to match transcripts used in SpliceAI calculations). SAI-10k-calc was developed using human genome reference GRCh37, but is compatible with GRCh38.

The SAI-10k-calc output is a tab-separated file with summary of splicing predictions indicating the type of aberration, possible combinations of aberrations (e.g. one SNV resulting in both exon skipping and partial intron retention), the exact size of inserted and/or deleted sequences, and the effect on reading frame and translation (Fig. 1). The latter is critical to predict the pathogenicity of the splicing alteration, and to design and interpret laboratory validation experiments. Amino acid sequence predictions could also be useful for additional applications, for example cancer neoantigen predictions (Yarchoan et al. 2017).

We also provide a lightweight Microsoft Excel spreadsheet ([Supplementary File S2](#), processing up to 1000 SNVs) that predicts the types and sizes of aberrations. In this version, users need to provide the raw scores generated by either SpliceAI Lookup (<https://spliceailookup.broadinstitute.org/>) or SpliceAI run from the command line. However, the predicted aberration sizes for partial intron retentions or partial exon deletions for this lightweight version are less accurate than the R code implemented version. Specifically, the R code uses native splice site positions derived from the given RefSeq transcript, whereas the lightweight version uses the SpliceAI-predicted acceptor and donor site positions. For example, the Excel version incorrectly predicted NM_007294.4(*BRCA1*):c.4868C>G to result in a 125-bp partial exon deletion, while the R code gave the correct size of 119 bp.

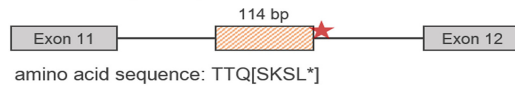
We note that, due to SpliceAI limitations, the calculator cannot be designed to predict three specific combinations of aberrant transcripts: (i) exon skipping and multi-exon skipping; (ii) exon skipping and whole intron retention; and (iii) partial exon deletion and partial intron retention. Multi-exon skipping and whole intron retention can only be predicted if the positions of donor and acceptor losses are within the analysis window, i.e. <4999 bp from the variant.

4 Performance

Using our training set data, SAI-10k-calc (R code version) has an overall sensitivity of 95% (441/464 confirmed spliceogenic SNVs) and specificity of 96% (715/748 non-spliceogenic SNVs) using our

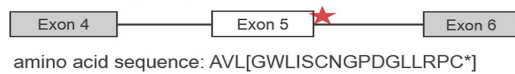
A. Pseudoexon activation

NM_007294.4(*BRCA1*):c.4185+4105C>T



B. Exon skipping

NM_000179.3(*MSH6*):c.3438+1G>A



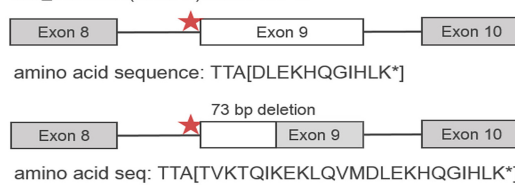
C. Partial exon deletion

NM_000251.3(*MSH2*):c.2635–1G>A



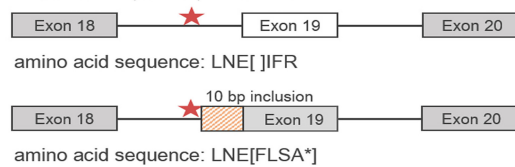
D. Exon skipping and partial exon deletion

NM_000059.4(*BRCA2*):c.682–2A>C



E. Exon skipping and partial intron retention

NM_007294.4(*BRCA1*):c.5194–12G>A



F. Whole and partial intron retention

NM_000249.4(*MLH1*):c.1667G>C

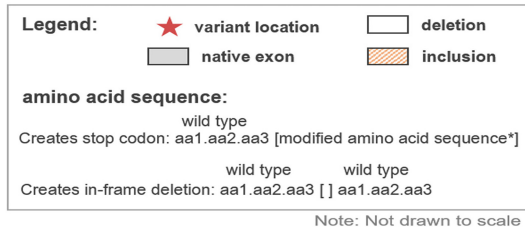


Figure 1 Types of splicing aberrations predicted by the SpliceAI-10k calculator. Six SNVs that were experimentally confirmed to alter splicing are correctly predicted by SAI-10k-calc (A–F). Of these, two represent correct prediction of combinations of splicing aberration types (D, E). Amino acid sequence predictions include three amino acids preceding the first variant amino acid, followed by the modified sequence inside square brackets. In-frame deletions that do not introduce any missense amino acid change (example in panel E) are indicated by blank square brackets flanked by three wild type amino acids preceding and following the deleted sequence. In all cases, deleted amino acids are not shown

thresholds. Furthermore, SAI-10k-calc demonstrates high accuracy for prediction of pseudoexonization (85%), partial intron retention (84%), and exon skipping (81%), highlighting its applicability for prioritization of variants through clinical or research sequencing. R code output data from the training set variants are shown in

Supplementary Table S2. General splicing prediction results and performance values are summarized in **Supplementary Tables S3 and S4.** Splicing aberration predictions (type and size) are detailed in **Supplementary Tables S5–S8.**

Author notes

The authors wish it to be known that Daffodil M. Canson and Aimee L. Davidson contributed equally as first authors. Olga Kondrashova and Amanda B. Spurdle shared chief supervision and last authorship.

Acknowledgements

Emma Tudini and Vaishnavi Nathan of the QIMR Berghofer Molecular Cancer Epidemiology Laboratory updated the *BRCA1* and *BRCA2* splicing table that partly provided curated splicing data for training set 1.

Author contributions

Daffodil M. Canson (Conceptualization-Lead, Data curation-Lead, Formal analysis-Lead, Methodology-Lead, Visualization-Lead, Writing—original draft-Lead, Writing—review & editing-Equal), Aimee L. Davidson (Conceptualization-Supporting, Formal analysis-Lead, Methodology-Lead, Software-Lead, Writing—original draft-Lead, Writing—review & editing-Equal), Miguel de la Hoya (Formal analysis-Supporting, Methodology-Supporting, Writing—review & editing-Equal), Michael T. Parsons (Data curation-Supporting, Formal analysis-Supporting, Writing—review & editing-Equal), Dylan M. Glubb (Supervision-Supporting, Writing—review & editing-Equal), Olga Kondrashova (Conceptualization-Supporting, Formal analysis-Supporting, Methodology-Lead, Software-Lead, Supervision-Lead, Writing—review & editing-Equal), and Amanda B. Spurdle (Funding acquisition-Lead, Supervision-Lead, Writing—review & editing-Equal).

Supplementary data

Supplementary data is available at *Bioinformatics* online.

Funding

D.M.C. was supported by QIMR Berghofer Ailsa Zinns PhD Scholarship, QIMR Berghofer HDC Top Up Scholarship, and The University of Queensland Research Training Tuition Fee Offset. The work of A.L.D. was

supported in part by National Institutes of Health grant R01 CA264971. M.d.l.H. is supported by a grant from the Spanish Ministry of Science and Innovation, Plan Nacional de I+D+I 2013–2016, ISCIII [PI20/00110] co-funded by FEDER from Regional Development European Funds (European Union). O.K. is supported by a National Health and Medical Research Council Emerging Leader 1 Investigator Grant [APP2008631]. A.B.S. is supported by National Health and Medical Research Council Investigator Fellowship funding [APP1177524].

Conflict of interest: None declared.

Data availability

The data underlying this article are available in its **Supplementary Material**.

References

- Ha C, Kim J-W, Jang J-H. Performance evaluation of SpliceAI for the prediction of splicing of NF1 variants. *Genes* 2021;12:1308.
- Jaganathan K, Kyriazopoulou Panagiotopoulou S, McRae JF *et al.* Predicting splicing from primary sequence with deep learning. *Cell* 2019;176:535–48.e24.
- Moles-Fernández A, Domènech-Vivó J, Tenés A *et al.* Role of splicing regulatory elements and in silico tools usage in the identification of deep intronic splicing variants in hereditary breast/ovarian cancer genes. *Cancers* 2021; 13:3341.
- Movassat M, Forouzmand E, Reese F *et al.* Exon size and sequence conservation improves identification of splice-altering nucleotides. *RNA* 2019;25:1793–805.
- Qian X, Wang J, Wang M *et al.* Identification of deep-intronic splice mutations in a large cohort of patients with inherited retinal diseases. *Front Genet* 2021;12:647400.
- Riepe TV, Khan M, Roosing S *et al.* Benchmarking deep learning splice prediction tools using functional splice assays. *Hum Mutat* 2021;42:799–810.
- Rowlands C, Thomas HB, Lord J *et al.*; Genomics England Research Consortium. Comparison of in silico strategies to prioritize rare genomic variants impacting RNA splicing for the diagnosis of genomic disorders. *Sci Rep* 2021;11:20607.
- Shamsani J, Kazakoff SH, Armean IM *et al.* A plugin for the ensembl variant effect predictor that uses MaxEntScan to predict variant spliceogenicity. *Bioinformatics* 2019;35:2315–7.
- Wai HA, Lord J, Lyon M *et al.*; Splicing and disease working group. Blood RNA analysis can increase clinical diagnostic rate and resolve variants of uncertain significance. *Genet Med* 2020;22:1005–14.
- Yarchoan M, Johnson BA, Lutz ER *et al.* Targeting neoantigens to augment antitumour immunity. *Nat Rev Cancer* 2017;17:209–22.