

# OpenSpliceAl and Splam: Enhancing Splice Site Prediction Across Species to Improve Transcriptome Assembly

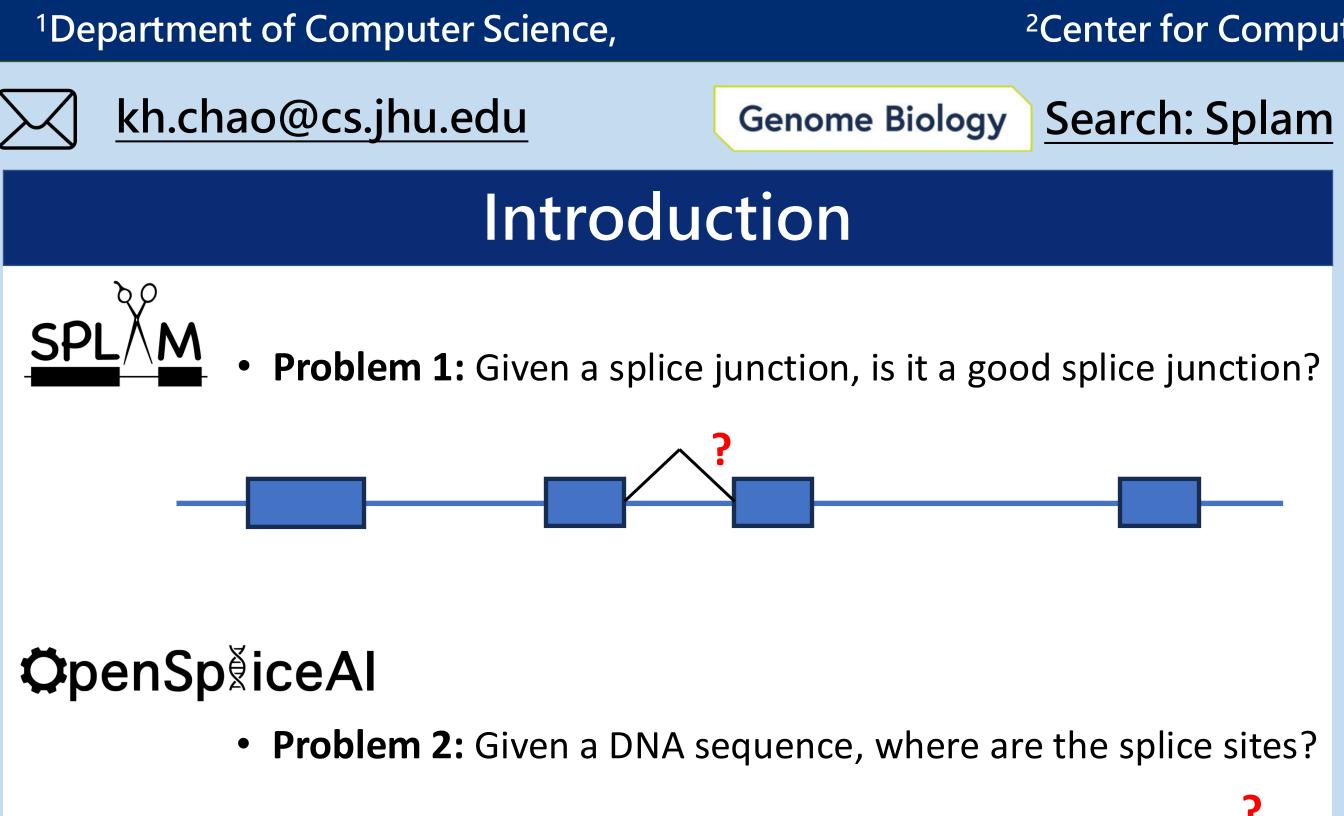


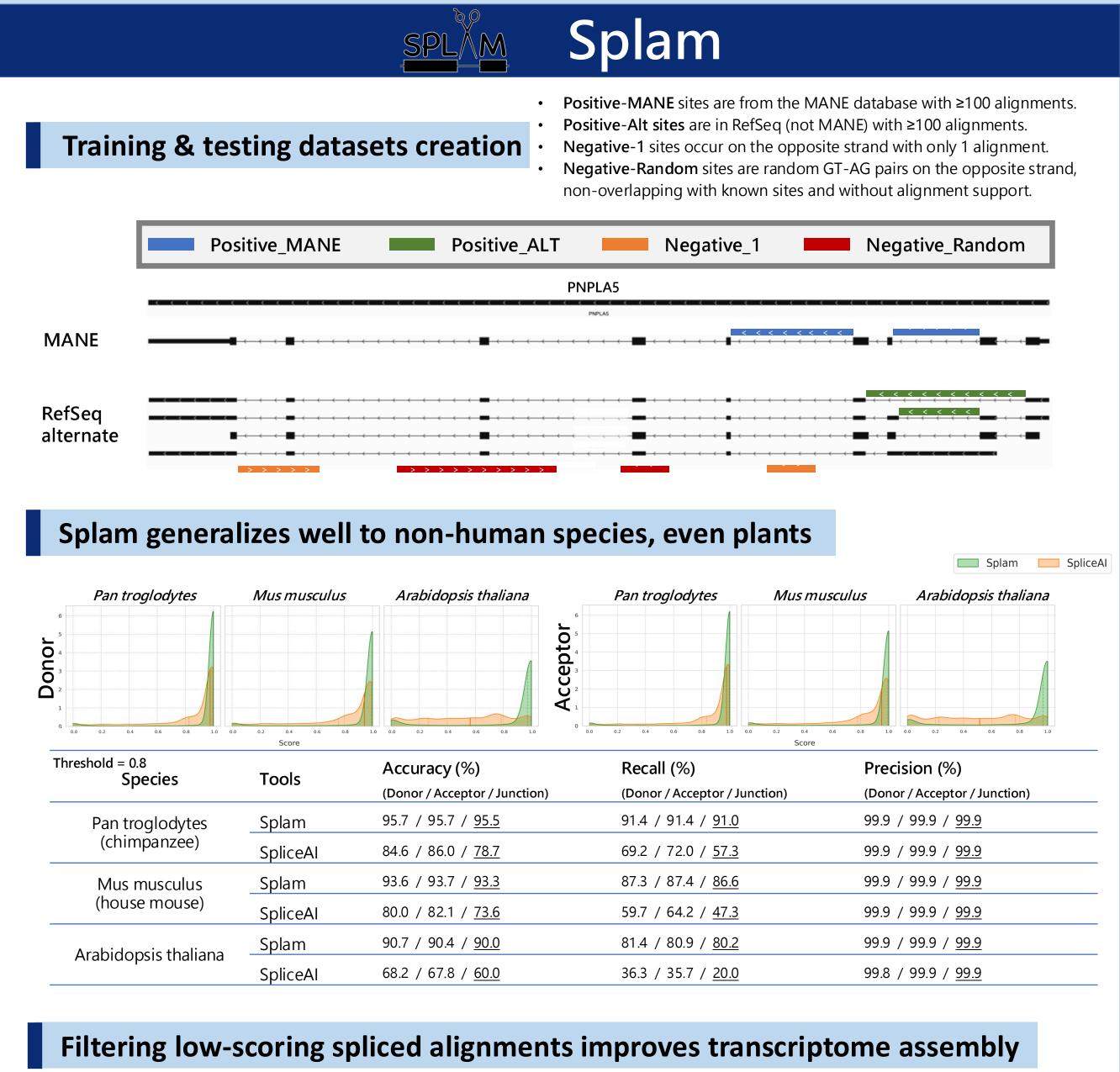
Kuan-Hao Chao<sup>1, 2, \*, †</sup>, Alan Mao<sup>1, 2, 3, †</sup>, Anqi Liu<sup>1</sup>, Steven L Salzberg<sup>1, 2, 3, 4, \*</sup>, and Mihaela Pertea<sup>1, 2, 3</sup>

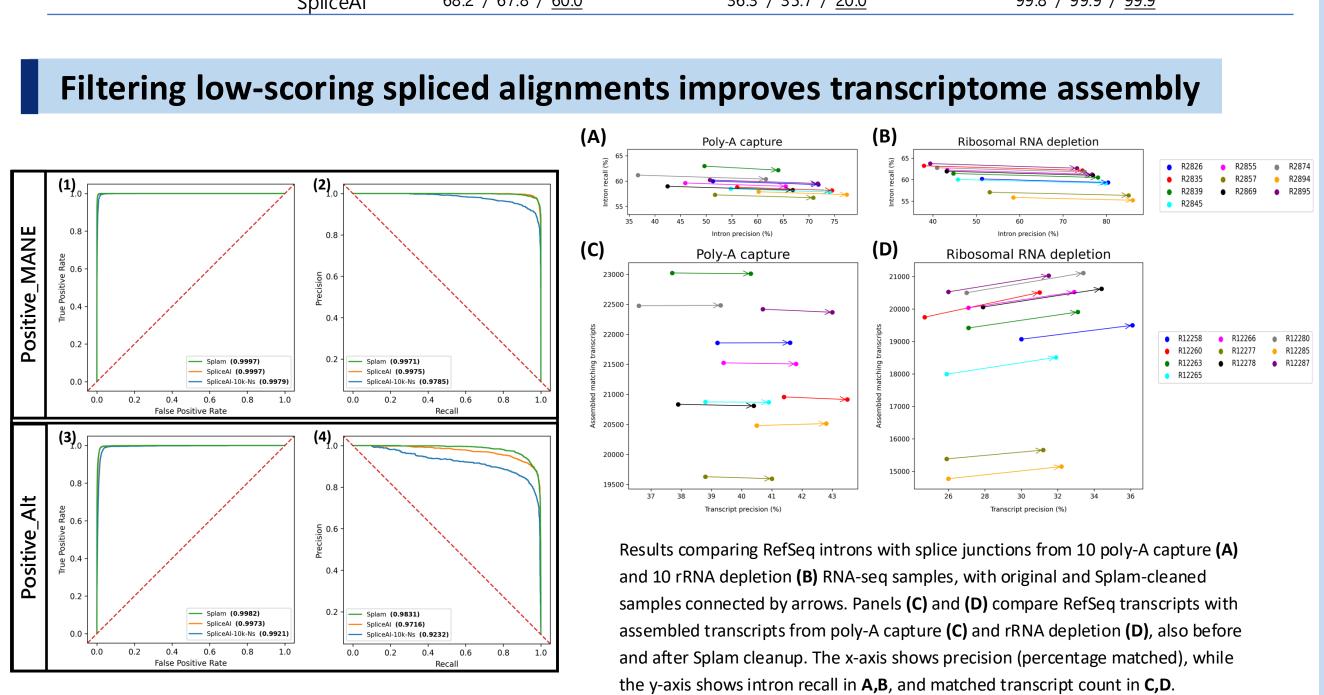
<sup>2</sup>Center for Computational Biology,

<sup>3</sup>Department of Biomedical Engineering,

<sup>4</sup>Department of Biostatistics, Johns Hopkins University













https://github.com/Kuanhao-Chao/splam



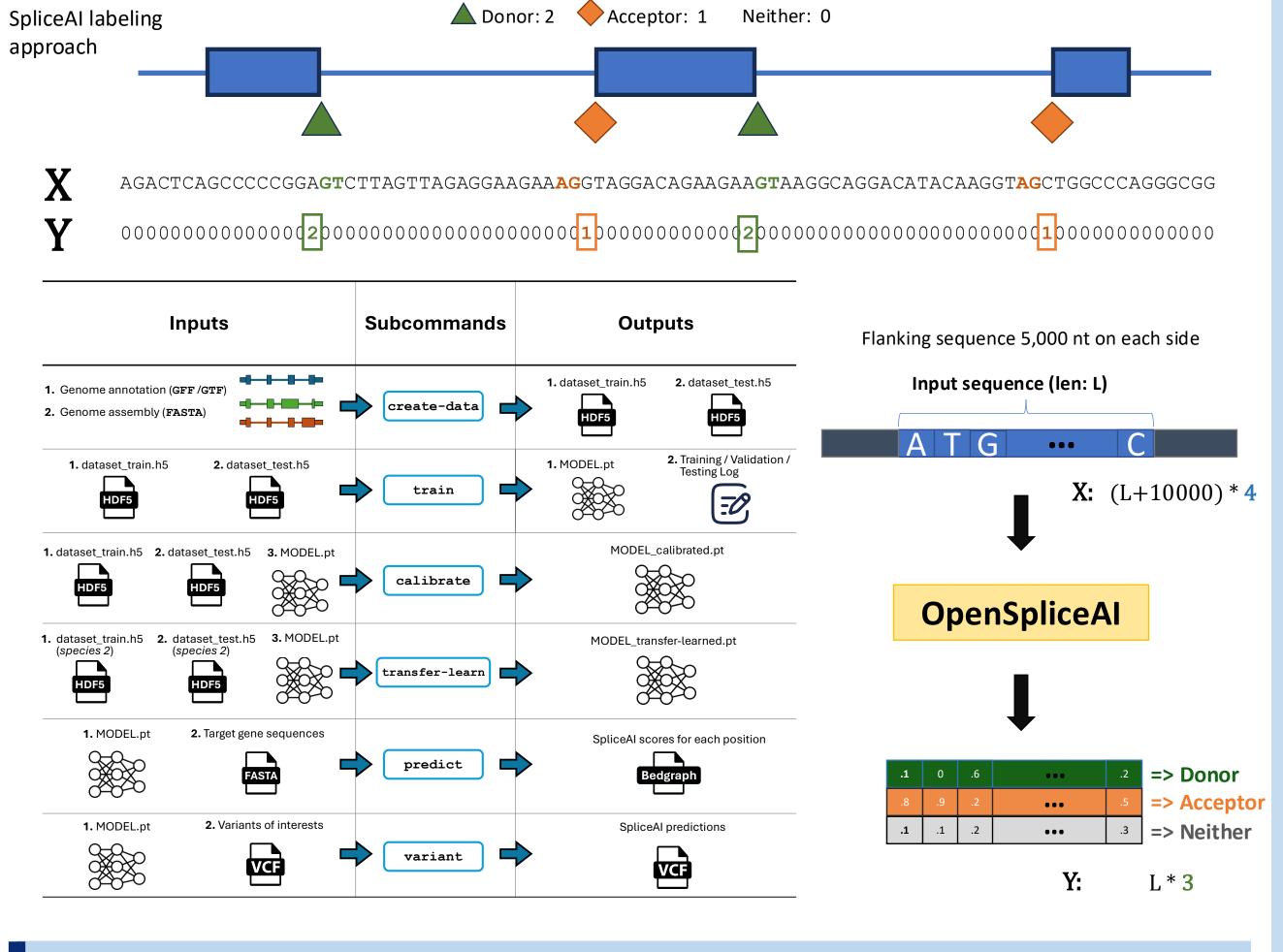
## **OpenSp**iceAl

OpenSpliceAl is an open-source version of the SpliceAl program<sup>1</sup>, a highly accurate splice site prediction system. OpenSpliceAI uses the newer PyTorch ML package instead of the older, slower Keras package, but otherwise uses the same code and is intended to replicate SpliceAl and allow users to re-train on their own species of interest.

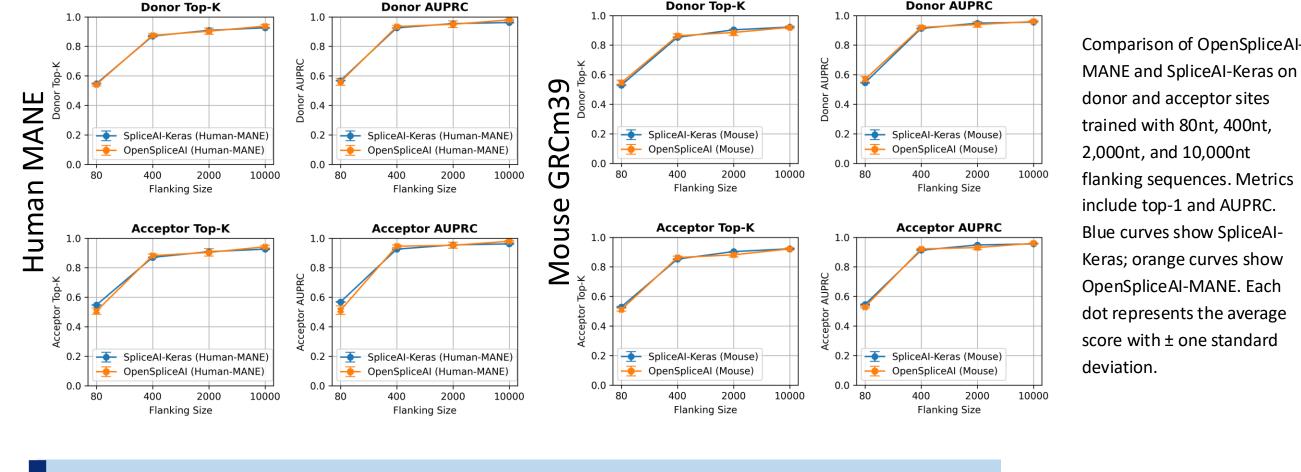
1. Jaganathan, Kishore, et al. "Predicting splicing from primary sequence with deep learning." Cell 176.3 (2019): 535-548.



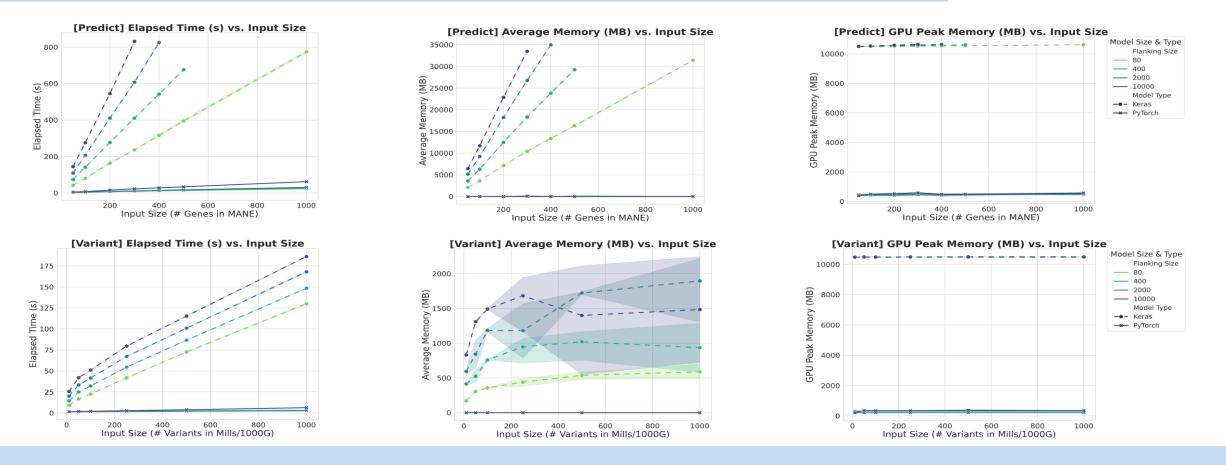
#### OpenSpliceAl design & model training

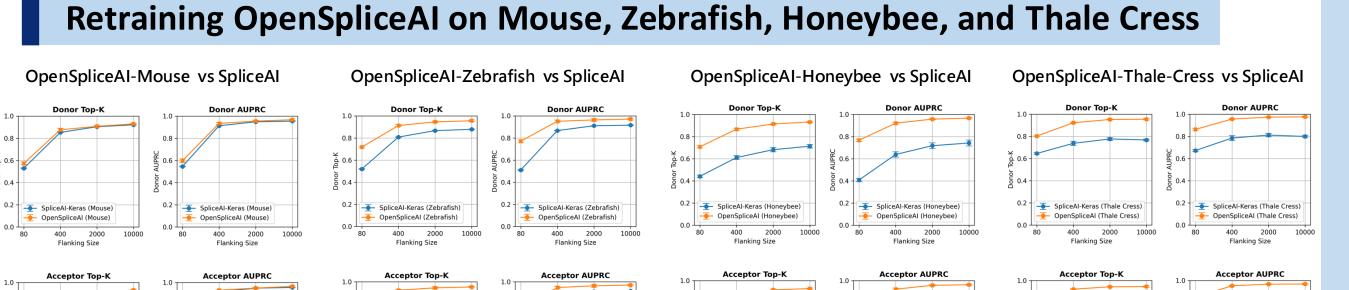


#### Testing OpenSpliceAI-MANE on Human MANE and Mouse GRCm39 splice sites

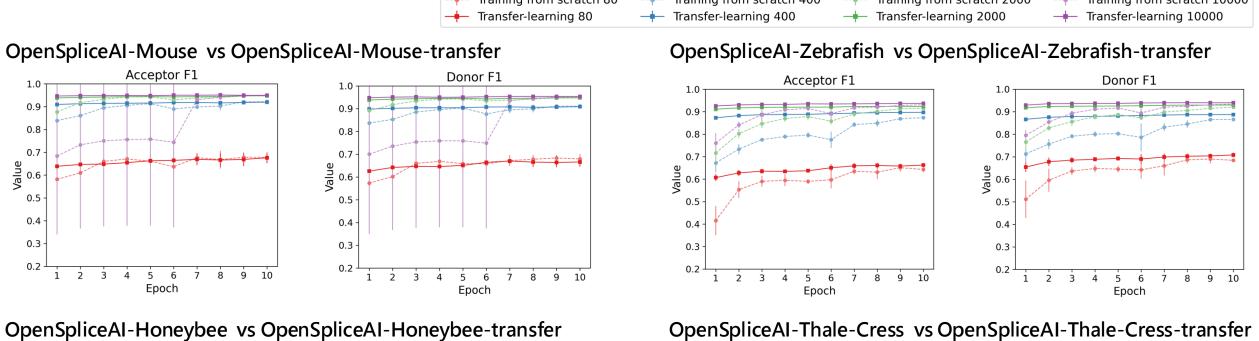


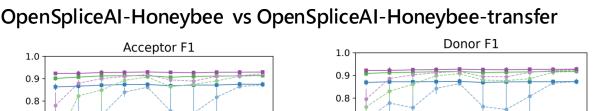
#### Elasped time, GPU Memory, and main memory benchmark

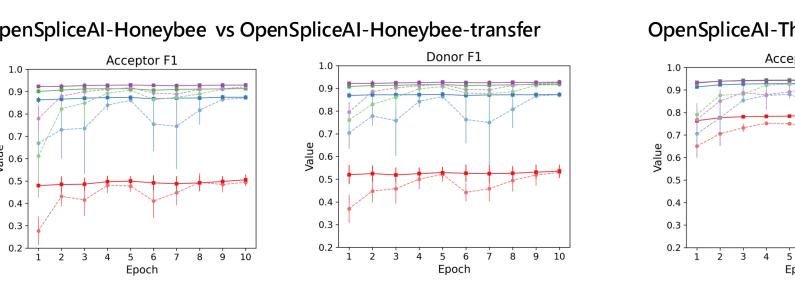


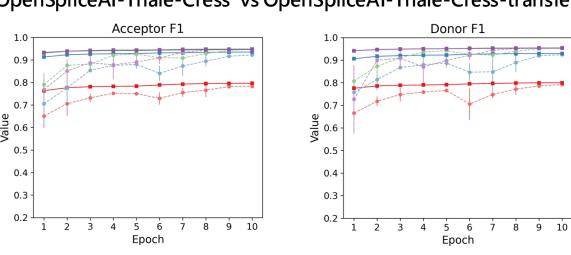


#### Transfer learning: Transfer OpenSpliceAI-MANE to four species





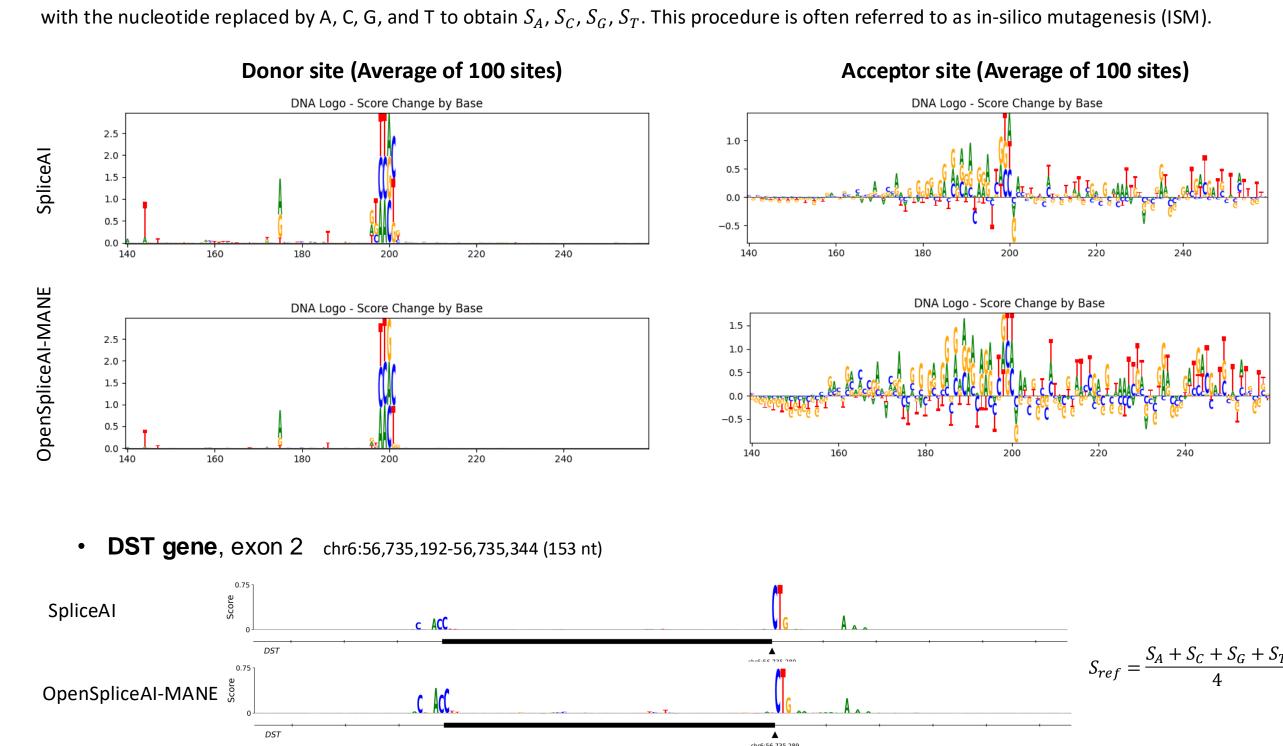




Transfer learning improves model performance by leveraging knowledge from related domains. We tested if OpenSpliceAI, trained on human splice annotations could adapt to predict splice sites in other species. Using five pre-trained OpenSpliceAI-MANE models, we fine-tuned species-specific models for mouse, honeybee, zebrafish, and thale cress, then compared them to scratch-trained models. Figures show F1 score for donor and acceptor sites in 80 nt, 400 nt, 2,000 nt, and 10,000 nt models over epochs 1–10 on the test dataset, comparing scratch-trained and transfer-trained variants.

### In Silico Mutagenesis (ISM) analysis

Perturbation-based forward propagation determines the relevance of a feature via mutating each nucleotide to all three possible alternatives The "importance score" of a nucleotide for a splice acceptor is calculated by taking the reference acceptor score,  $S_{ref}$ , and recalculating it with the nucleotide replaced by A, C, G, and T to obtain  $S_A$ ,  $S_C$ ,  $S_G$ ,  $S_T$ . This procedure is often referred to as in-silico mutagenesis (ISM).



• **U2SURP gene**, chr1:142,740,137–142,740,263 (127 nt) **SpliceAI** OpenSpliceAl-MANE

Acknowledgements: supported by NIH under grants R01-HG006677, R35-GM130151, and by NSF under grants DBI-2412449