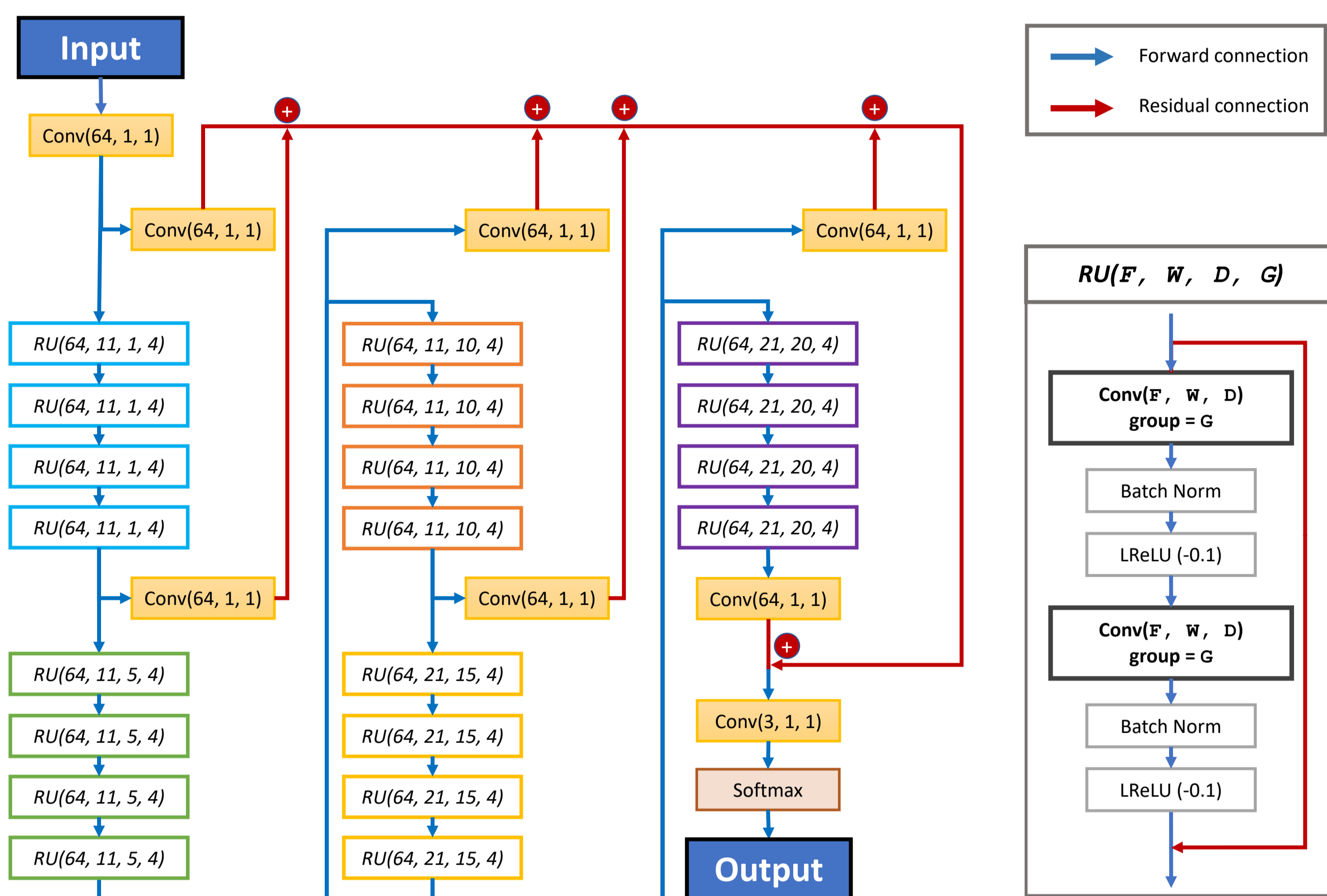


Introduction

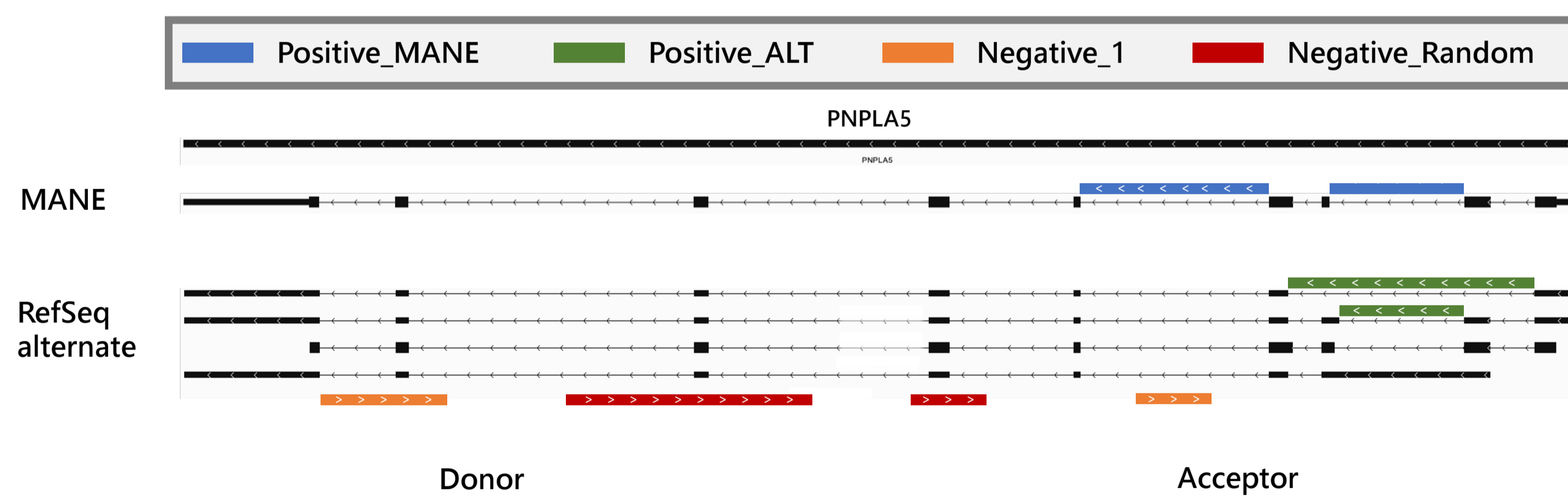
- Splam introduces the idea of training a neural network on donor and acceptor pairs together, inspired by the splicing machinery itself, which recognizes both ends of each intron at the same time.
- Splam uses a limited window of 400 bp flanking each splice site, again motivated by the biological process of splicing, which relies primarily on signals within this window
- Splam recognizes splice sites from genomic sequence alone more accurately than existing methods.
- Splam can improve the accuracy of transcript assemblies by removing spurious alignments produced by spliced aligners.

Methods

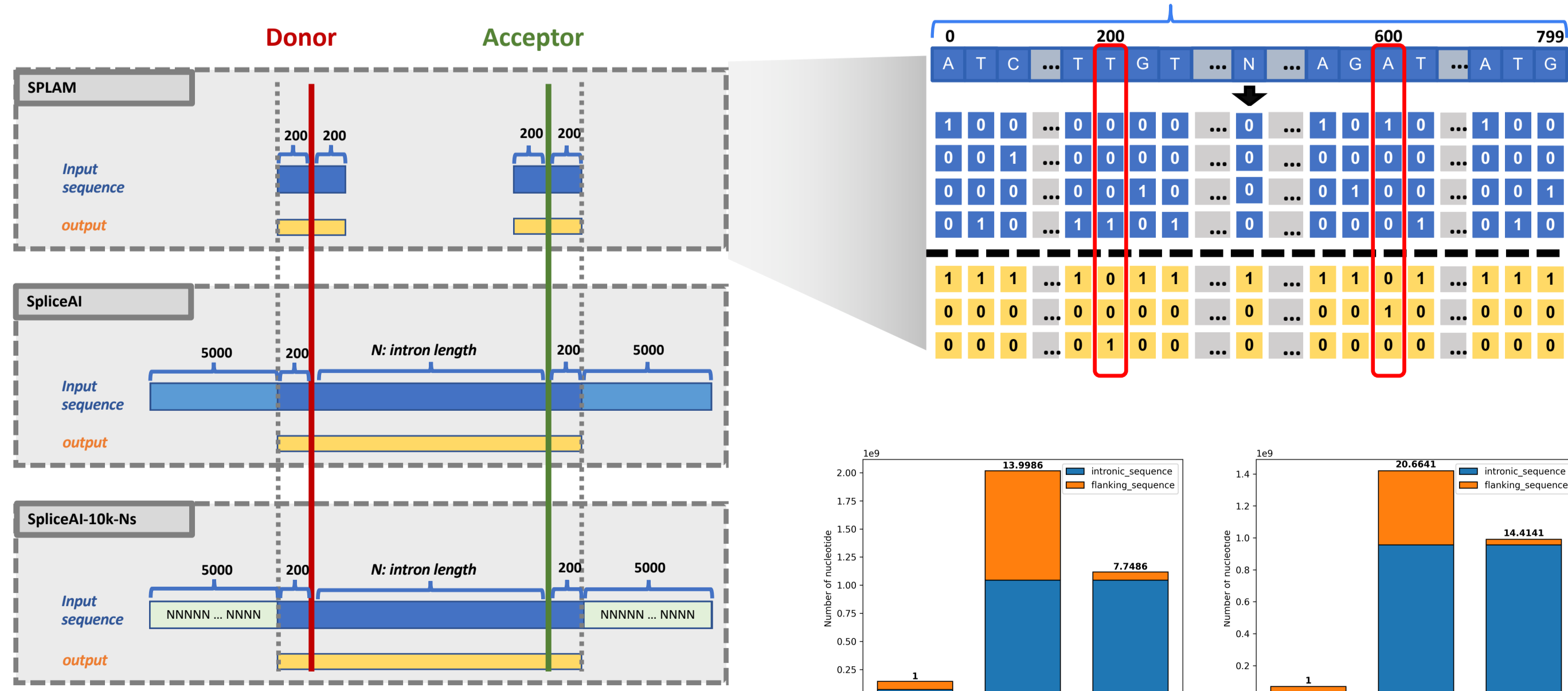
Model Architecture



Training & testing dataset creation



Data encoding



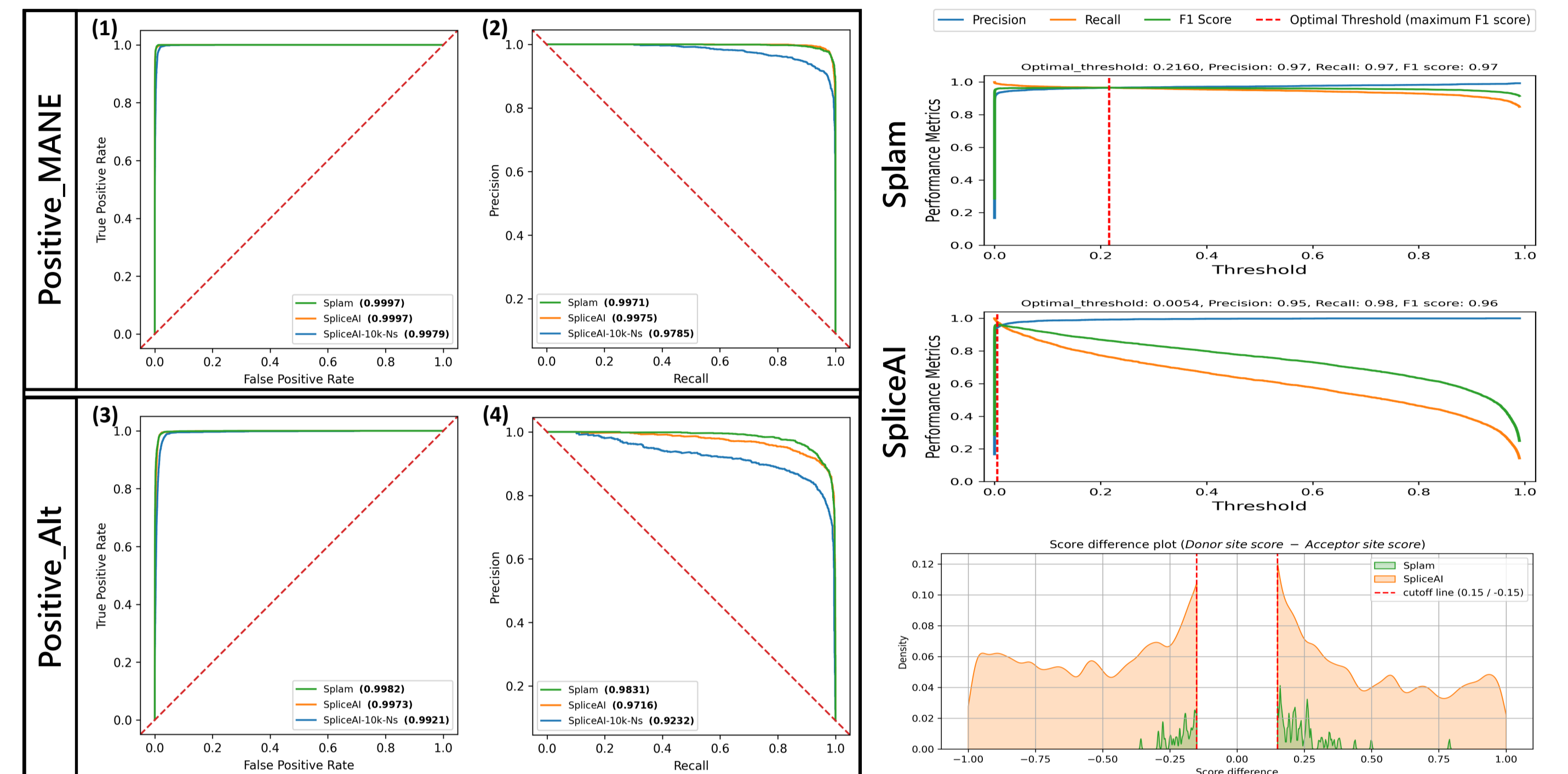
Splam's input (top row) uses 400nt flanking the donor site and another 400nt flanking the acceptor site. The output is labels for the 800nt region (shown in yellow).

SpliceAI's input (second row) follows its standard configuration, using 200nt upstream and downstream of the donor and acceptor sites, the entire intron, and 10Kb of flanking sequences.

SpliceAI-10k-Ns (third row) uses a modified input, replacing the 5Kb flanking sequences with Ns.

Results

Accuracy comparisons on human splice junctions in test dataset

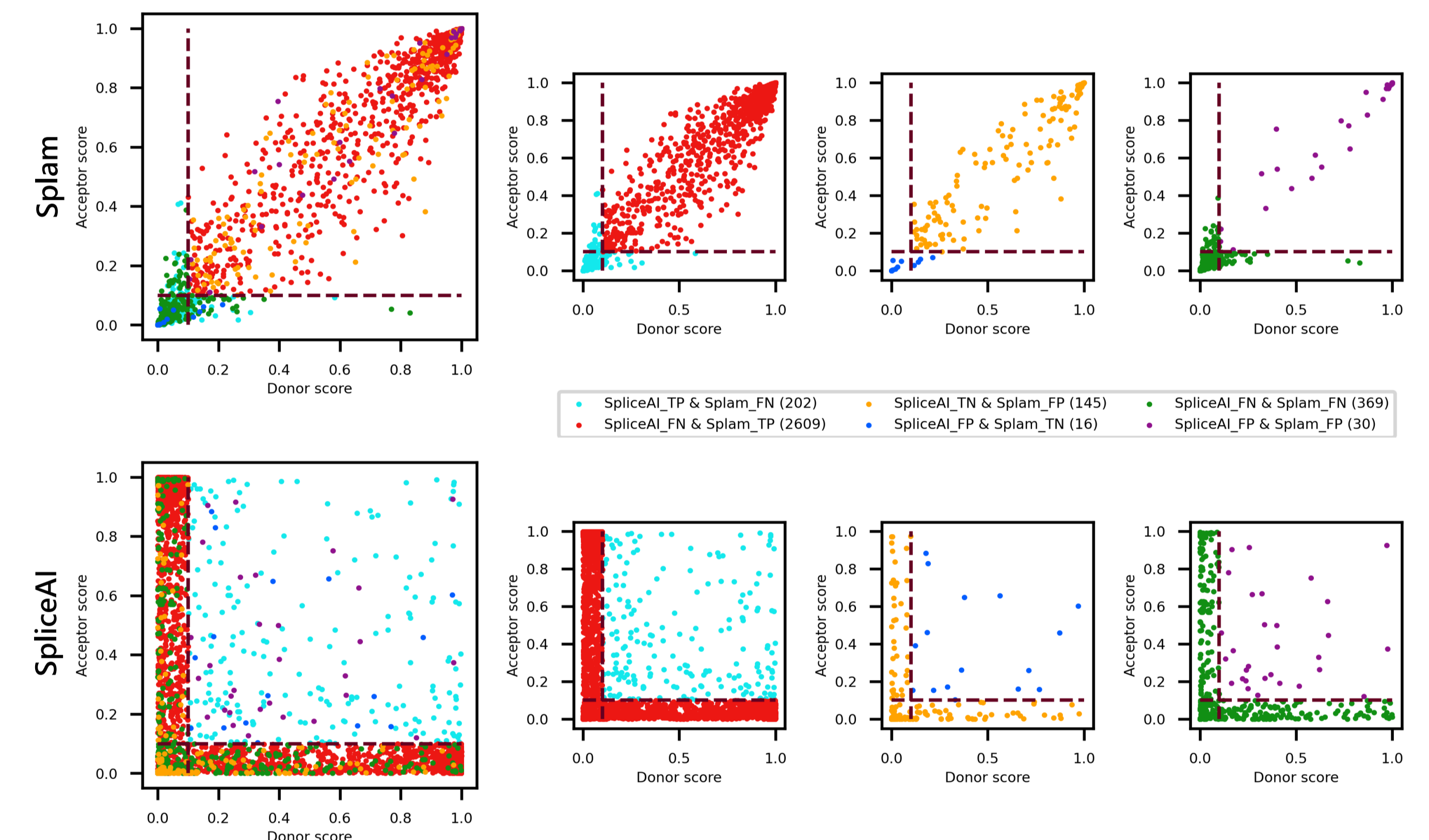


ROC / PR curves results are shown at the junction level, where the junction score is determined by the minimum of its donor and acceptor scores.

Discrimination threshold plots show the precision (blue curve), recall (orange curve), and F1 score (green curve) calculated at different thresholds. The optimal threshold (maximum F1 score) is indicated by a red dashed line.

Kernel density plot visualizes the differences between donor and acceptor scores (donor score - acceptor score)

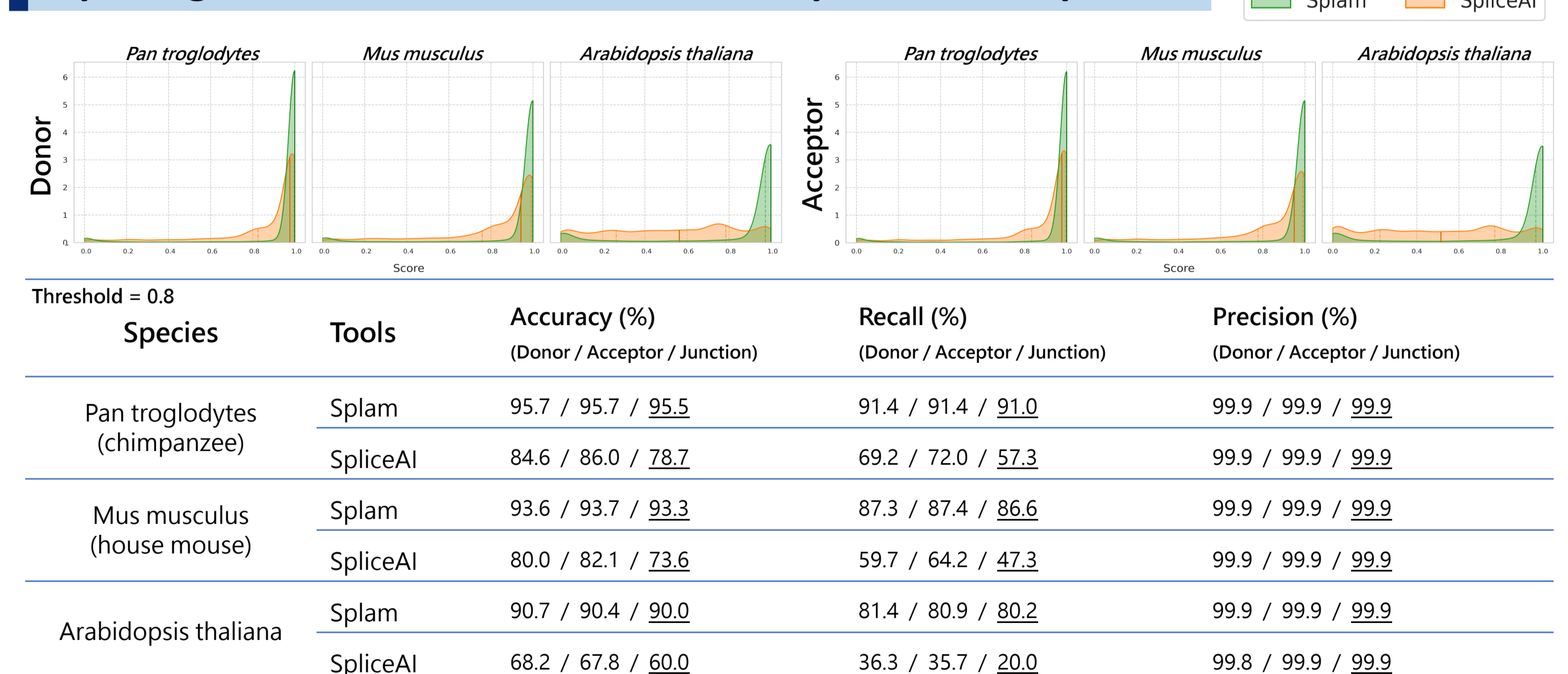
Comparing the score distributions of Splam and SpliceAI



Each dot represents a splice junction. The red dashed lines indicate the 0.1 cutoff threshold for labeling splice sites as true positives (TPs), true negatives (TNs), false positives (FPs), or false negatives (FNs).

Subplots in the second & third columns show cases where one program was correct while the other was incorrect (TP and FN, or FP and TN); the fourth column shows cases where both programs made incorrect predictions (FP and FN).

Splam generalizes well to non-human species, even plants



Filtering low-scoring spliced alignments improves transcriptome assembly

