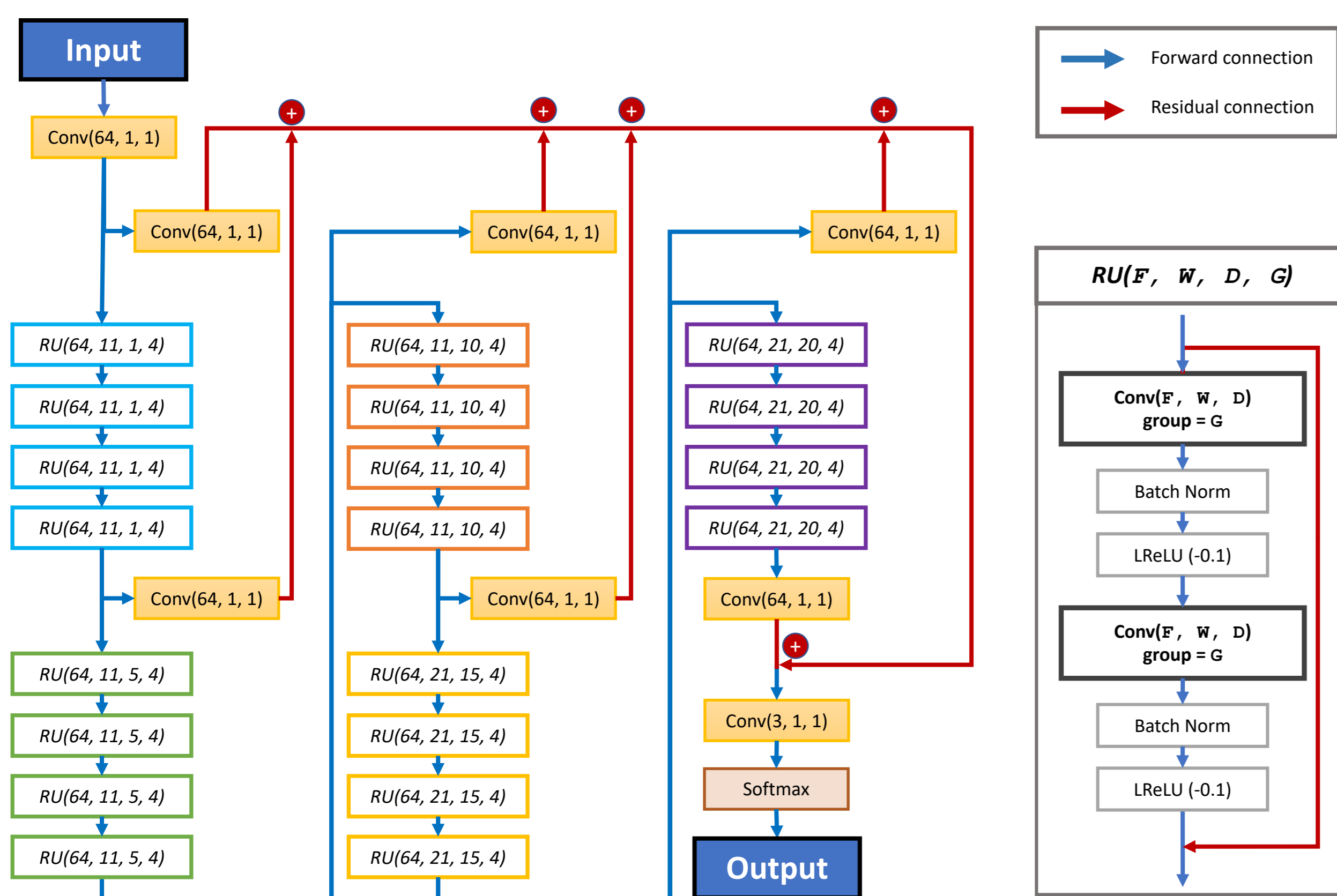


Introduction

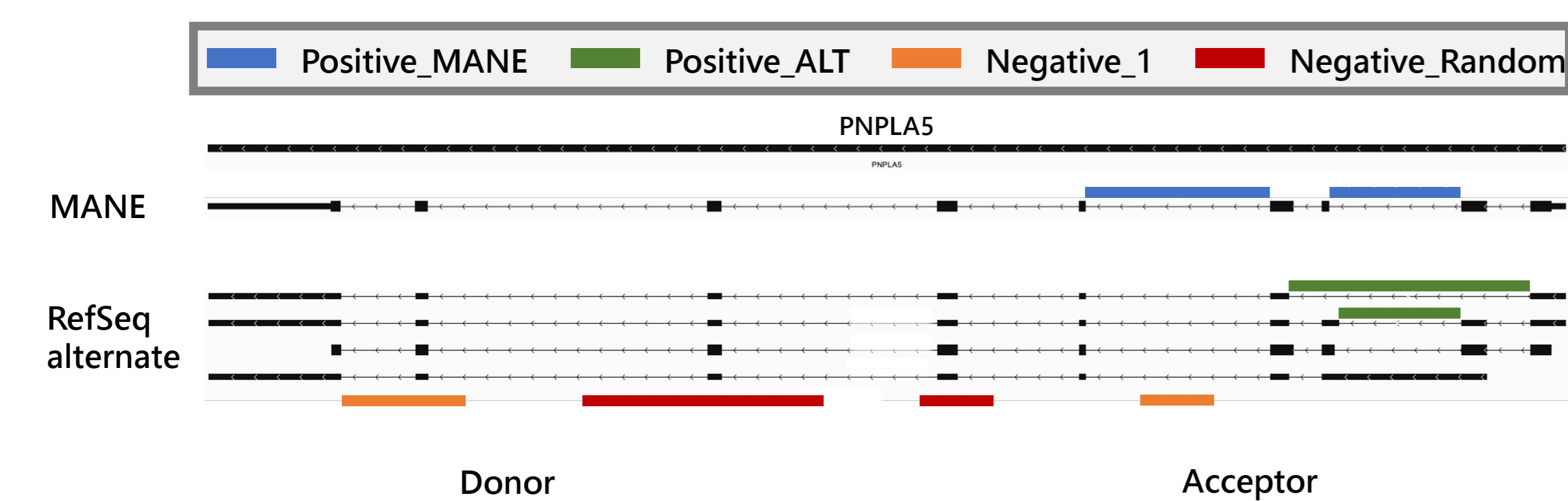
- Splam introduces the idea of training a neural network on donor and acceptor pairs together, inspired by the splicing machinery itself, which recognizes both ends of each intron at the same time.
- Splam uses a limited window of 400 bp flanking each splice site, again motivated by the biological process 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

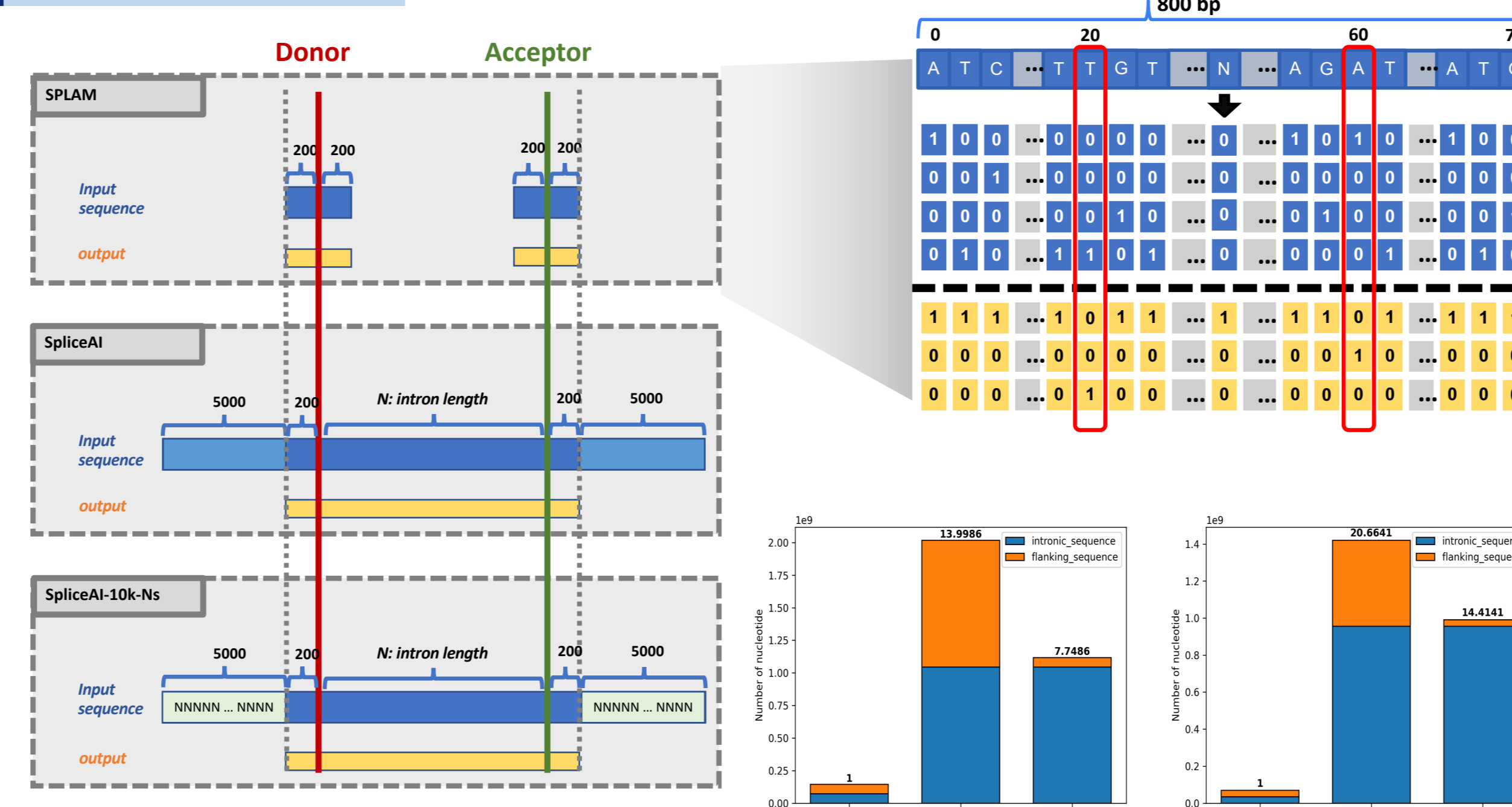


Training & testing dataset creation



- Positive-MANE sites are selected from the MANE database with ≥ 100 alignments.
- Positive-Alt sites are in RefSeq but not in MANE, also with ≥ 100 alignments.
- Negative-1 sites occur on the opposite gene strand, supported by only 1 alignment.
- Negative-Random sites are random GT-AG pairs on the opposite strand, not overlapping with known sites and lacking alignment support.

Data encoding

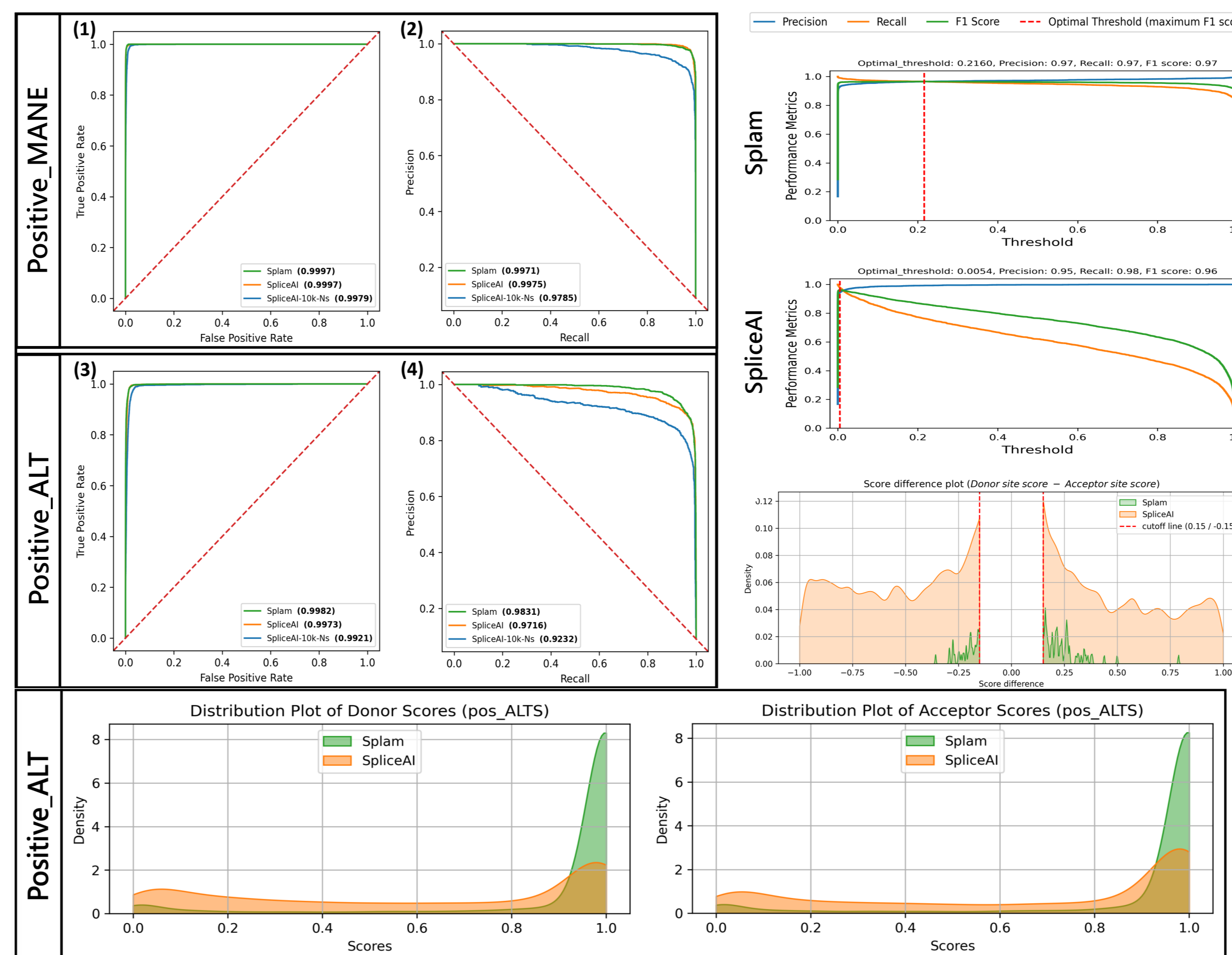


- Splam's input (top row) uses 400nt flanking the donor site and another 400nt flanking the acceptor site. The output is labels for the 800nt region (shown in yellow).
- SpliceAI's input (second row) follows its standard configuration, using 200nt upstream and downstream of the donor and acceptor sites, the entire intron, and 10Kb of flanking sequences.

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

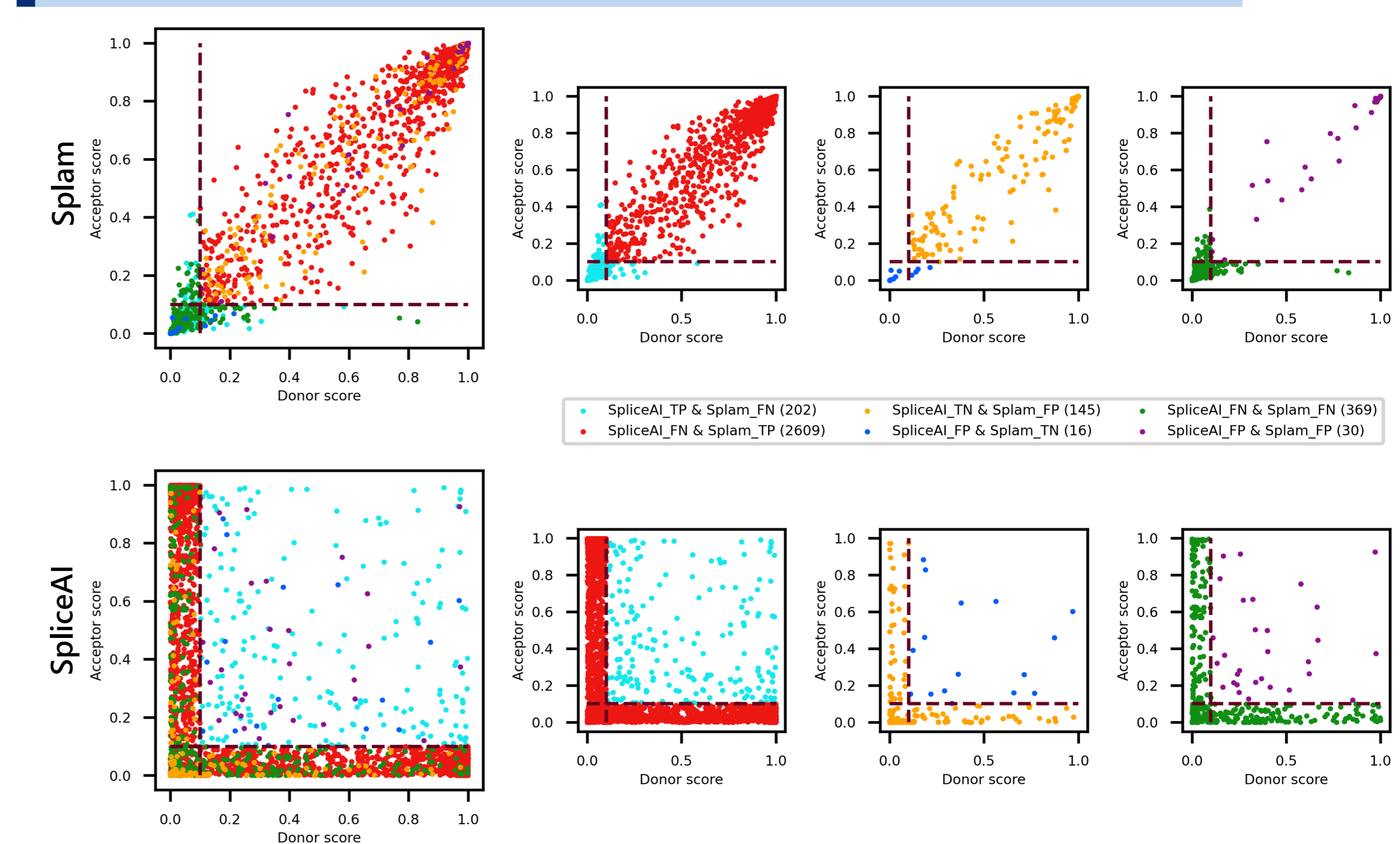
Results

Comparisons on human splice junctions in test dataset



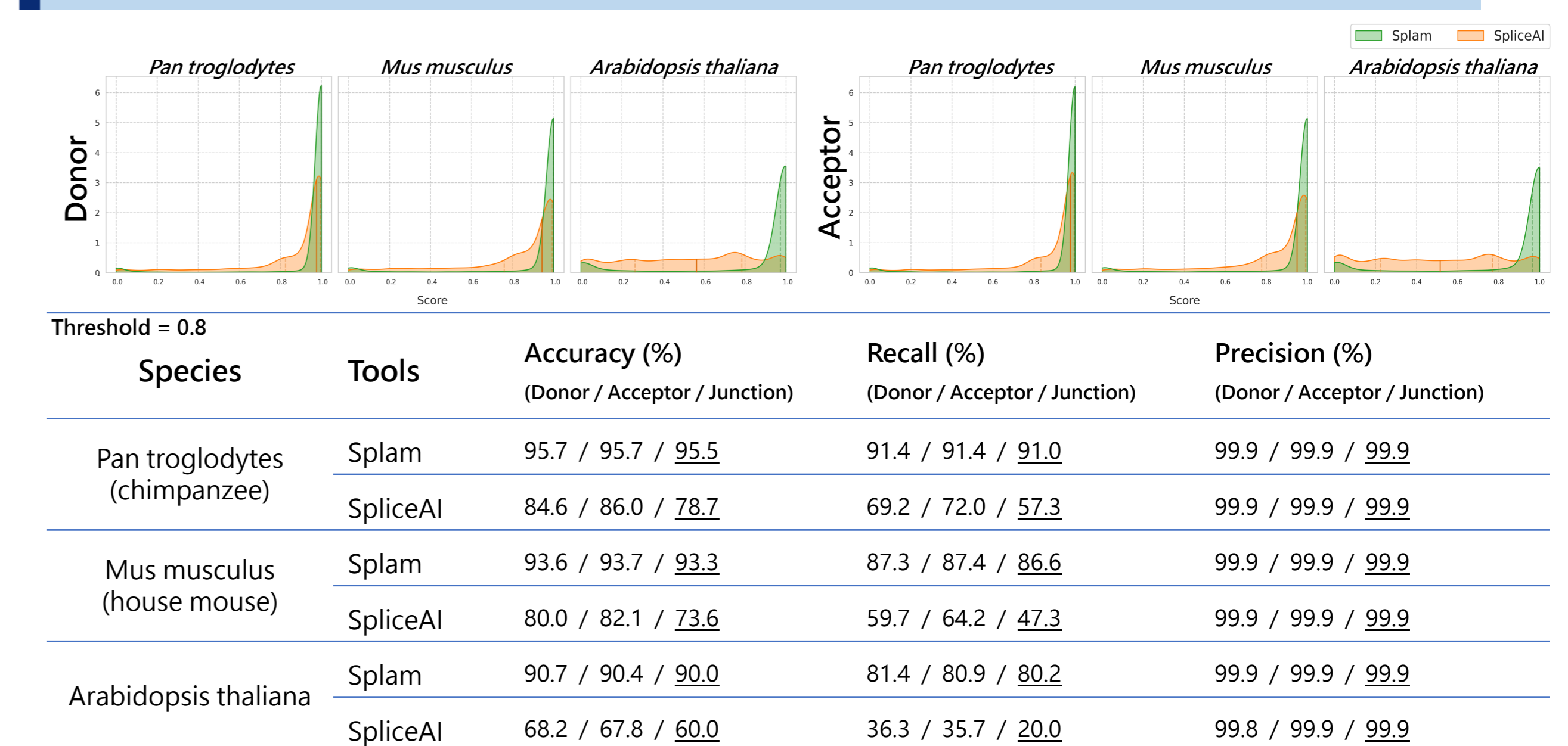
- ROC / PR curves results are shown at the junction level, where the junction score is determined by the minimum of its donor and acceptor scores.
- Discrimination threshold plots show the precision (blue curve), recall (orange curve), and F1 score (green curve) calculated at different thresholds. The optimal threshold (maximum F1 score) is indicated by a red dashed line.
- Kernel density plot shows the differences between donor and acceptor scores (donor score - acceptor score).
- Distribution of Splam and SpliceAI scores on alternative splice junctions.

Comparisons score distributions of Splam and SpliceAI



- Each dot represents a splice junction. The red dashed lines are 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 assembly

