

# Intron detection notes

Anh Vu Le

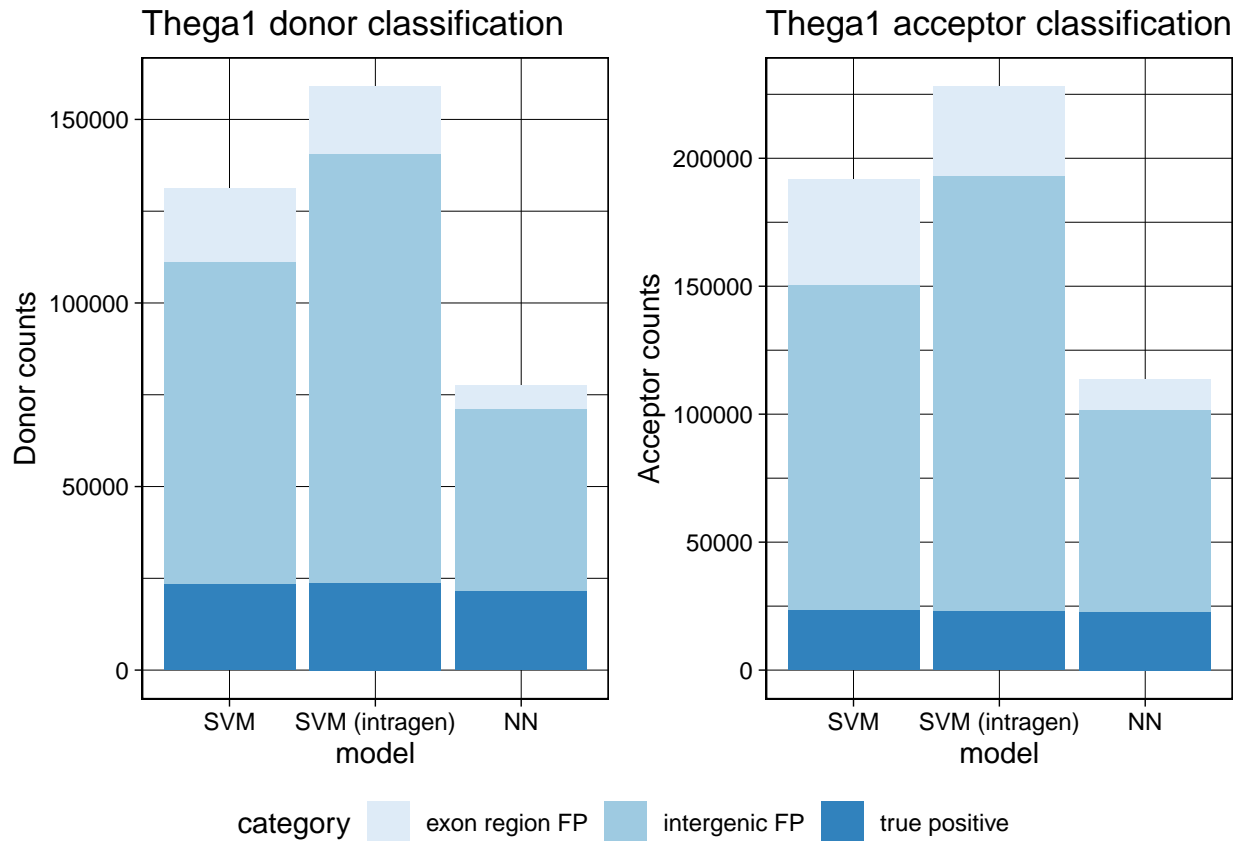
## Comparison of splice site classification models (on *Thega1* phylum)

This section will be devoted to comparing splice site classification models. In total 3 models were developed - two SVM models and NN model (technically there 6 of them as donors/acceptor site require each its own predictor). The SVM models differ in data, which they were trained on. First one (let's call it "general SVM") was given negative splice site candidates sampled from the entire genome with no distinction, whether they come from an intergenic or an intragenic (i.e. exon) region. Second SVM was given only negative candidates from within exons ("intragenic SVM"). NN models were also trained ignoring intergenic sequences. The motivation behind this decision comes from the assumption, that the accuracy in intergenic regions is irrelevant for the task at hand, which is basically about recognizing intron sequences from exon sequences.

In general validation, the models ended up toe to toe in sense of both precision and recall. The tests however do not capture the real nature of the task, which is intron cutting for the purposes of enhancing BLAST search results. Good results here are measured in a slightly more complex manner, as we will see shortly.

Intron cutting is a classification task. And as in all such tasks, we measure how successful we are at solving such a task by means of metrics such as precision and recall. In this case however, we need to adjust (or rather extend) the notion of these metrics, especially precision. In intron classification and cutting, not all false positives "hurt" the same - in fact, some false positives are even beneficial. This is because we deal with two types of false positives introns here - those which are in intergenic regions and those that are inside exons. The former do not concern us at all. If we cut into the intergenic region, it is only for good, as BLAST will receive smaller data at the cost of nothing. Intra-exon cuts are on the other hand the "real" false positives, as they will lower the quality of BLAST hits.

One way to compare the two models is therefore to look at the composition of examples, to which the models say they are introns (or splice sites). How many of them are true positives? And if they are false positives, how many of them will actually cause damage? If we compare SVM and NN splice site models in this fashion, we will reveal, that they are not quite the same, as they may seem after the first round testing.

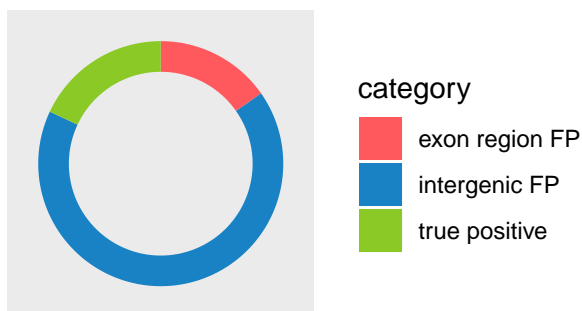


The chart captures several interesting facts about the models - first of all, all models recognized very similar number of true donors and acceptors. This might indicate, that there is set of introns which cannot be detected by neither of the methods. We will verify this hypothesis later on. Second, NN models produce substantially less false positives, than both SVM models, which carries forward to intragenic false positives, the harmful type of false positives. We can also see, that the intragenic SVM, even though producing more raw FP, than the general SVM, has actually better results in exon regions. Afterall, the model was built for such purpose.

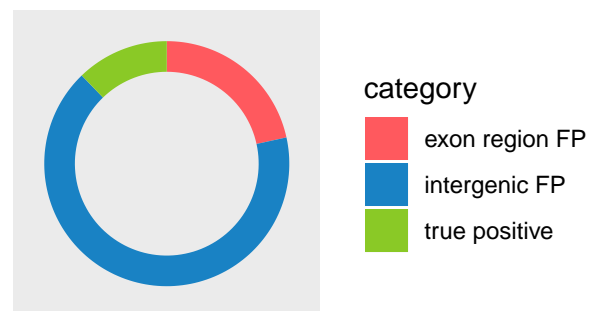
Let's visualize solely the proportions of the types of classified true candidates. We will be particularly interested in ratios between correctly classified donors/acceptors (green) and falsely classified candidates, that lie inside an exon (red).

## SVM models

### donor candidates

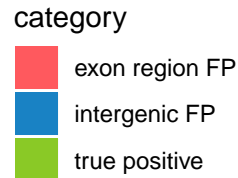
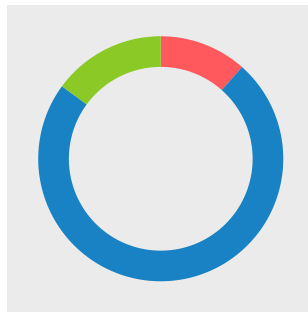


### acceptor candidates

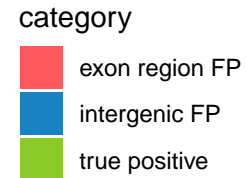
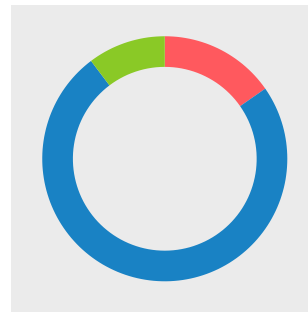


## SVM intragen models

### donor candidates

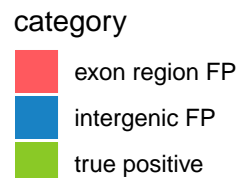
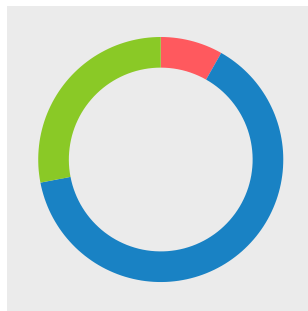


### acceptor candidates

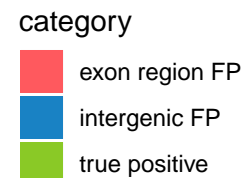
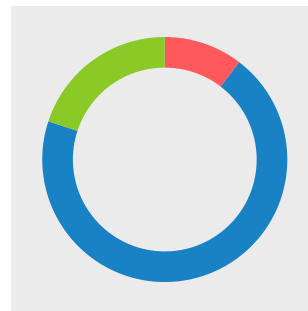


## NN models

### donor candidates



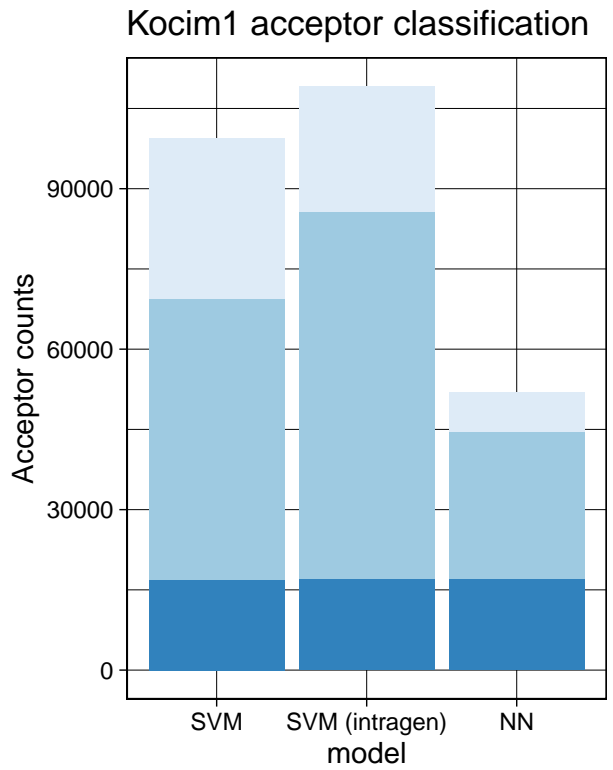
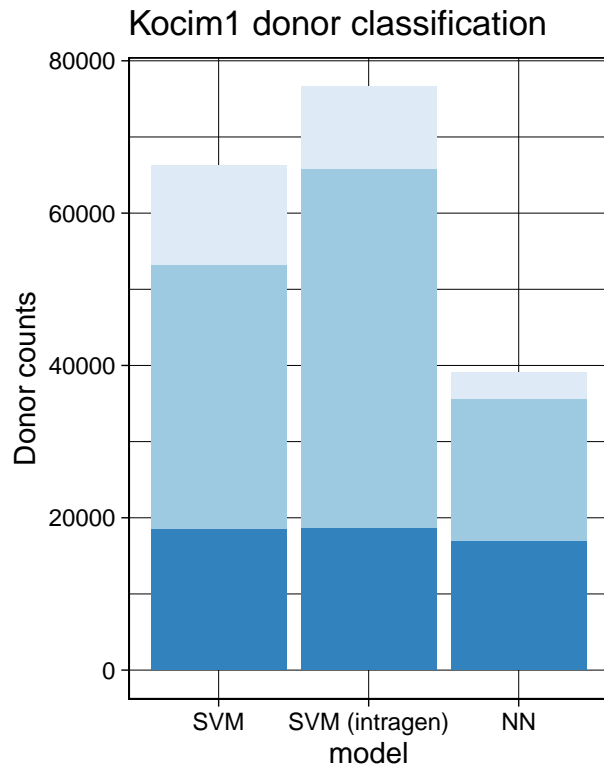
### acceptor candidates



Looking at the three pairs of graphs and the proportions of green-red regions we can clearly see, that NN models have significantly better ratio between correctly classified splice site candidates and intra-exon false positives. Note, that the green regions capture in all cases very similar number of examples. In case of SVM models, they appear to produce one false intragenic splice site for each correctly detected one. This ratio seems unfavorable and this is why a subsequent intron model is needed to improve those odds. Next, let's look at the the same analysis, this time for *Kocim1*.

## Comparison of splice site classification models (on *Kocim1* phylum)

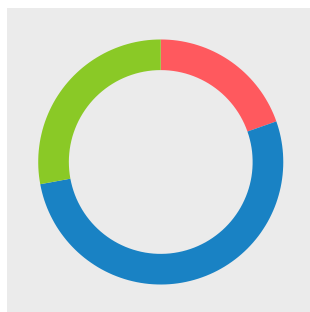
The graphs tell a similar story as in the *Thega1* case. Models find almost the same number of correct splice sites, NN however produce less false positives, both intergenic and - more importantly - intragenic. Intragenic SVM model produces more FP, than the general SVM, but slightly less intragenic FP.



category ■ exon region FP ■ intergenic FP ■ true positive

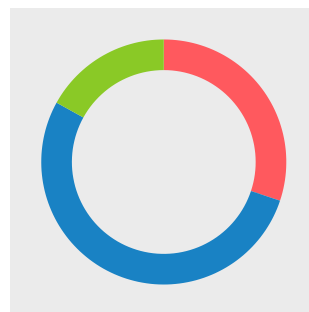
## SVM models

### donor candidates



category  
■ exon region FP  
■ intergenic FP  
■ true positive

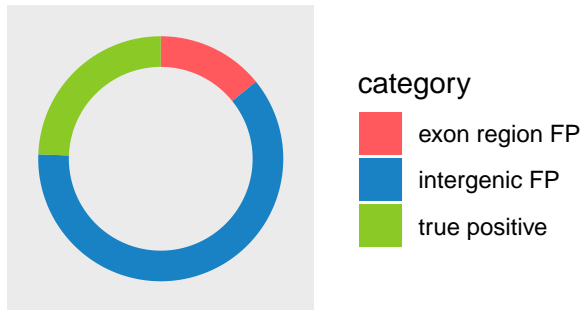
### acceptor candidates



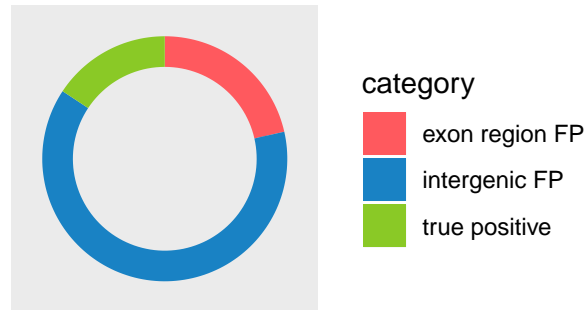
category  
■ exon region FP  
■ intergenic FP  
■ true positive

## SVM intragen models

### donor candidates

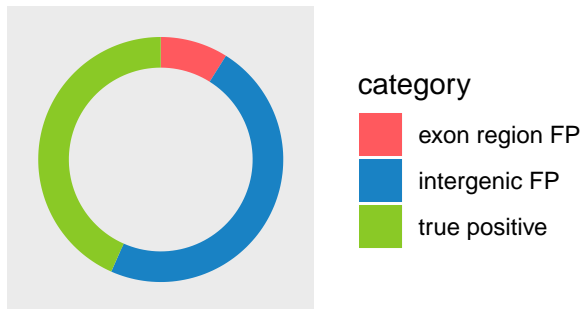


### acceptor candidates

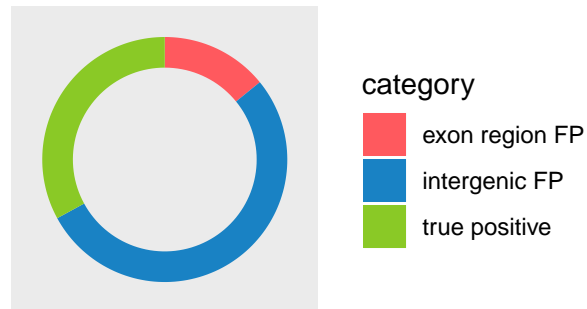


## NN models

### donor candidates



### acceptor candidates

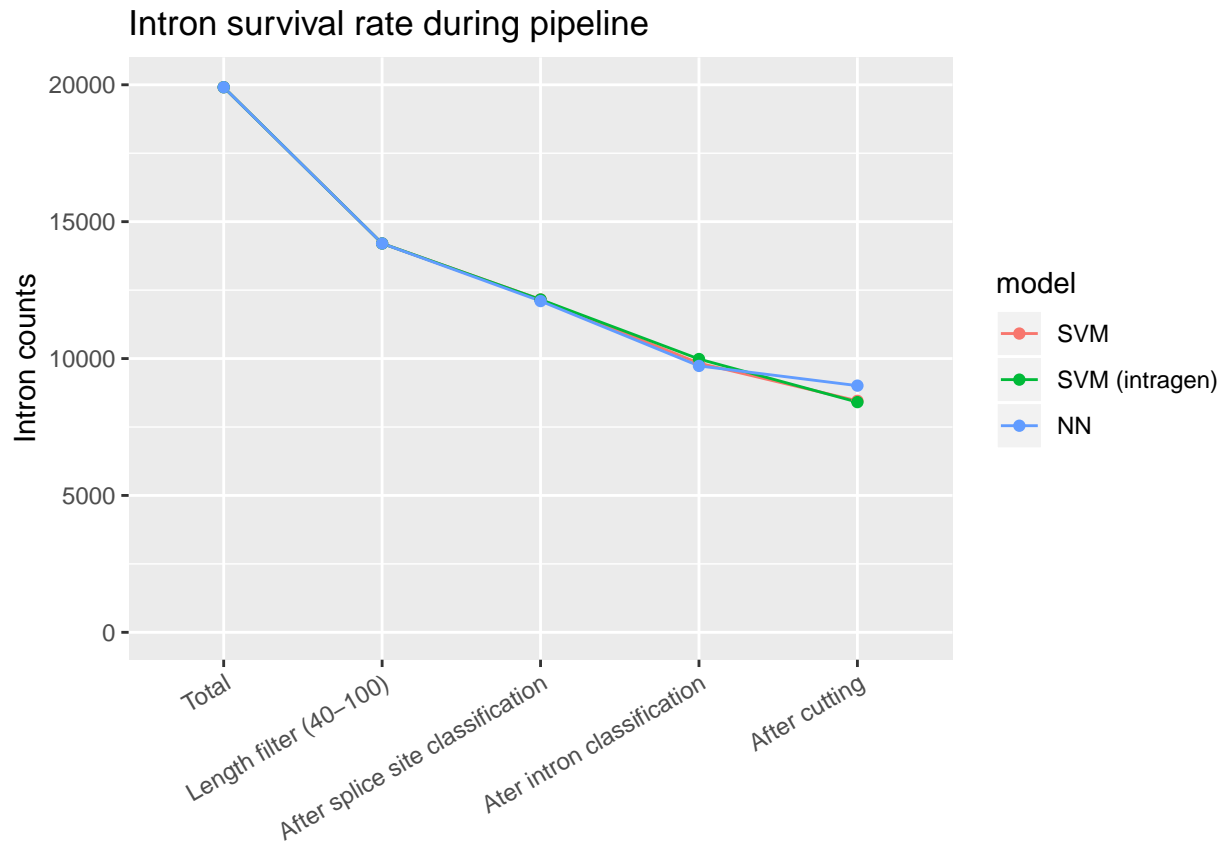


## Follow-up intron classification (*Kocim1*)

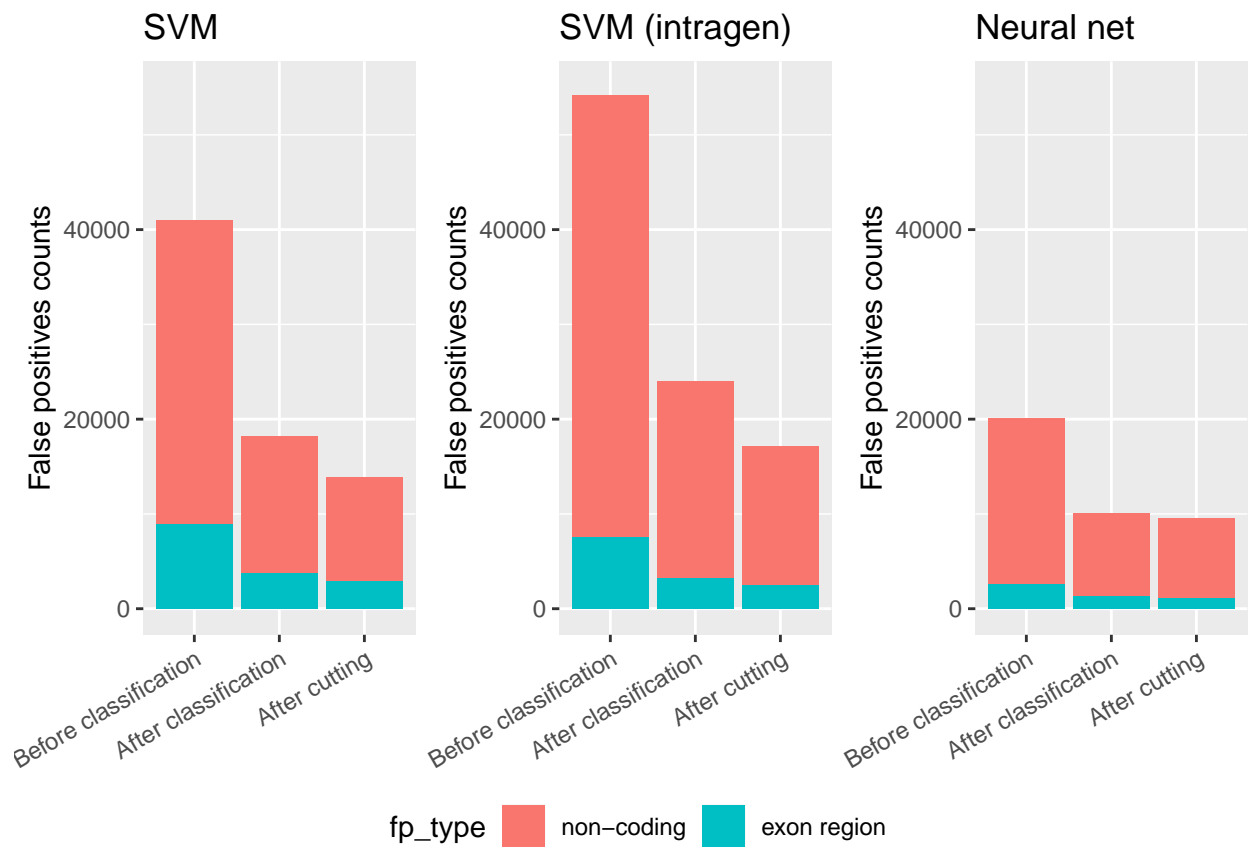
Next step after splice site classification is pairing positively classified donors/acceptor to form intron candidates. Intron model is then trained and candidates classified. Positively classified introns are then cut out of a scaffold.

It is prominent, that data intron model gets is determined by what splice site models produce. It is therefore reasonable, that each donor/acceptor pair of models should have a dedicated intron model. At this stage however, our main focus are still splice site models and how they influence the intron one. Therefore, for benchmarking purposes, we use an intron model trained with general SVM data. Even though this will give the general SVM models a slight edge, the comparison should still be possible. As we will see later, this is thanks to the fact there is a considerable overlap in the splice site models outputs.

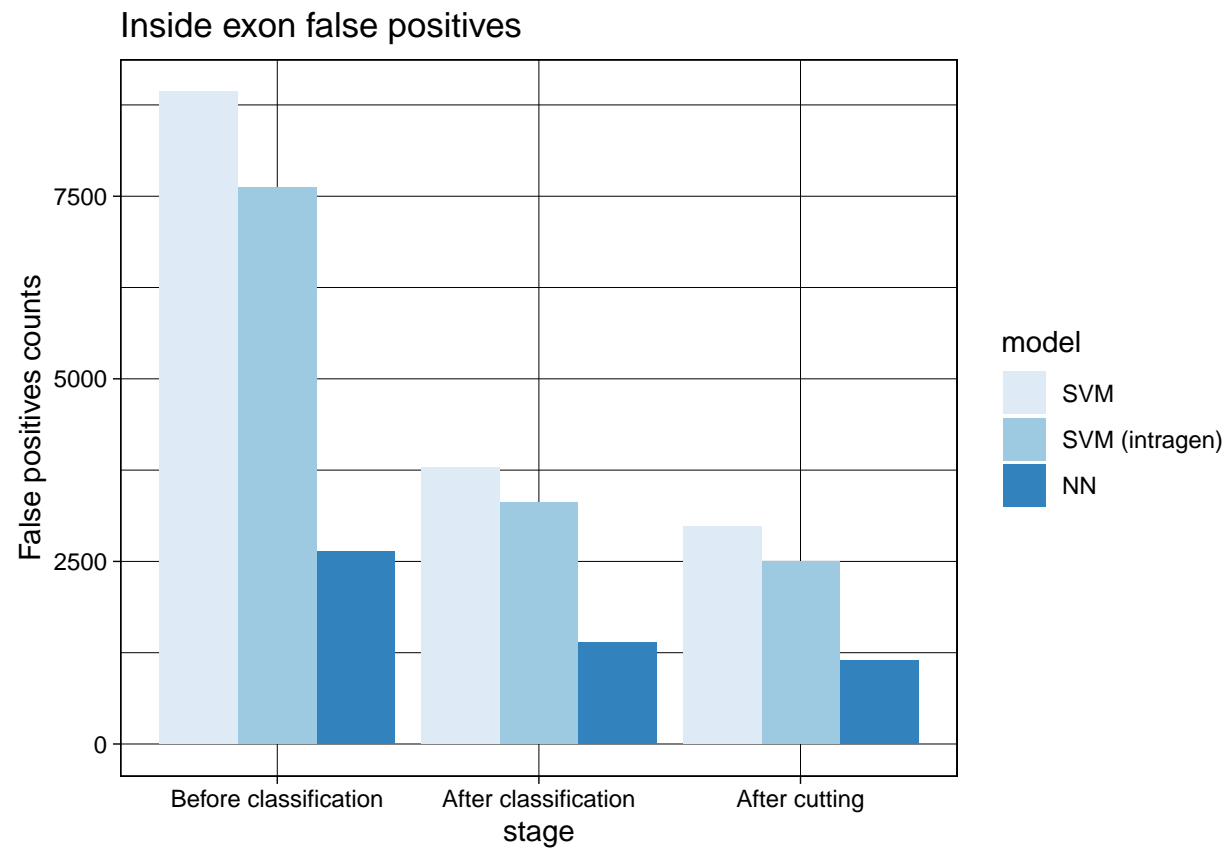
As mentioned before, intron cutting and purging is done in steps. Each steps loses a portion of introns due to some filtration or non-perfect recall. Let's visualize, how many introns are lost in each step. The first stage is the number of introns we start with - all introns on the positive strand. Next, we will lose those introns, which are not in the length range of 40-100. Undetected donors and acceptor will also cause some introns to be lost. Intron classification is an apparent source of losses and some are lost during cutting as well. This final decrease is due to incorrect decisions when dealing with candidate overlaps - two (or more) positively classified introns overlap and only one can be cut. And it may happen, that the wrong candidate is chosen over the correct one.



Now, let's explore the number of FP introns and their composition (in terms of location) for each model. The IG SVM has by far the most number of false positives. It is however better, than the general SVM, if we consider only intragenic false positives. In the setting, where the cutting only prepares data for BLAST, high number of false positives might even be beneficial. Considering the number of IG FP (the only adversary FPs) is in both types of SVMs similar, it can be beneficial to employ the IG version as it will shorten the input for the BLAST search. Still, the difference is marginal (after intron classification) and cannot make up for considerably better performance of NN.



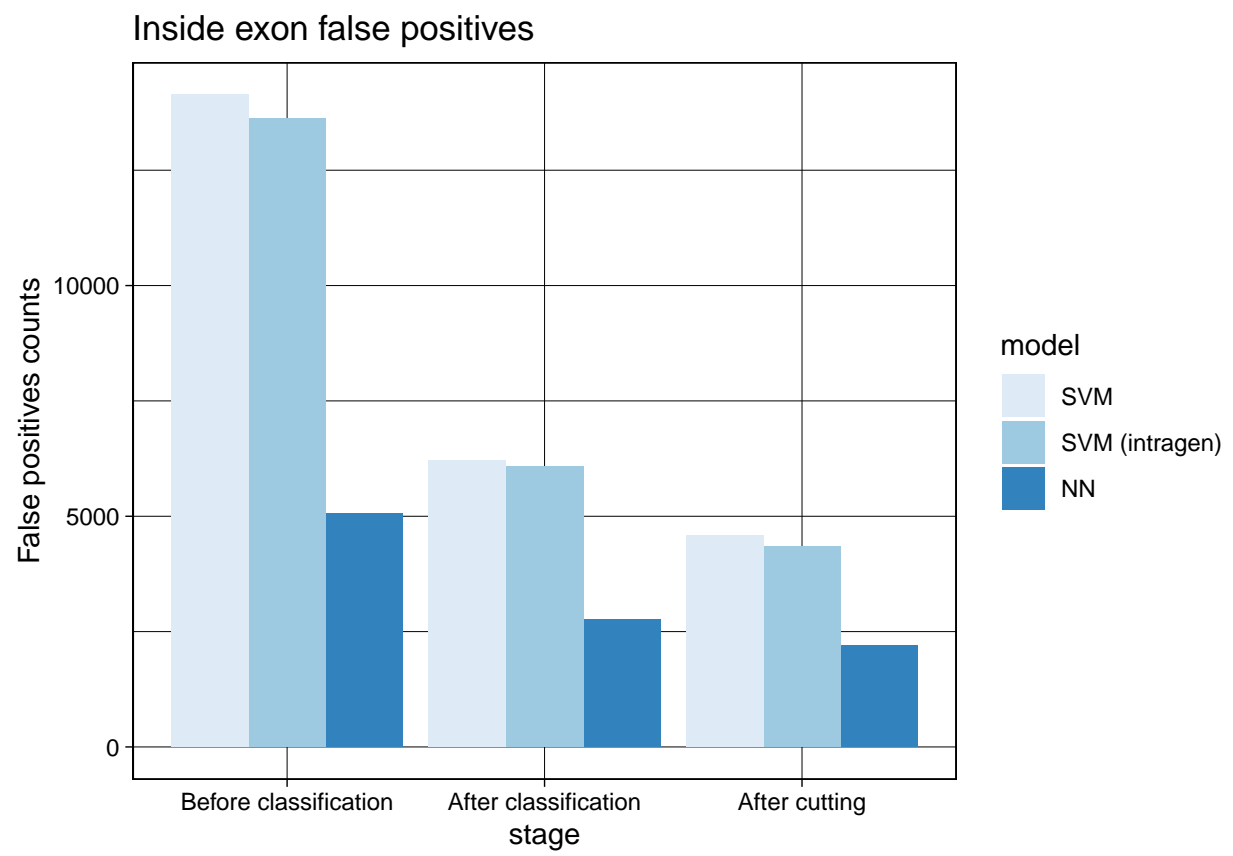
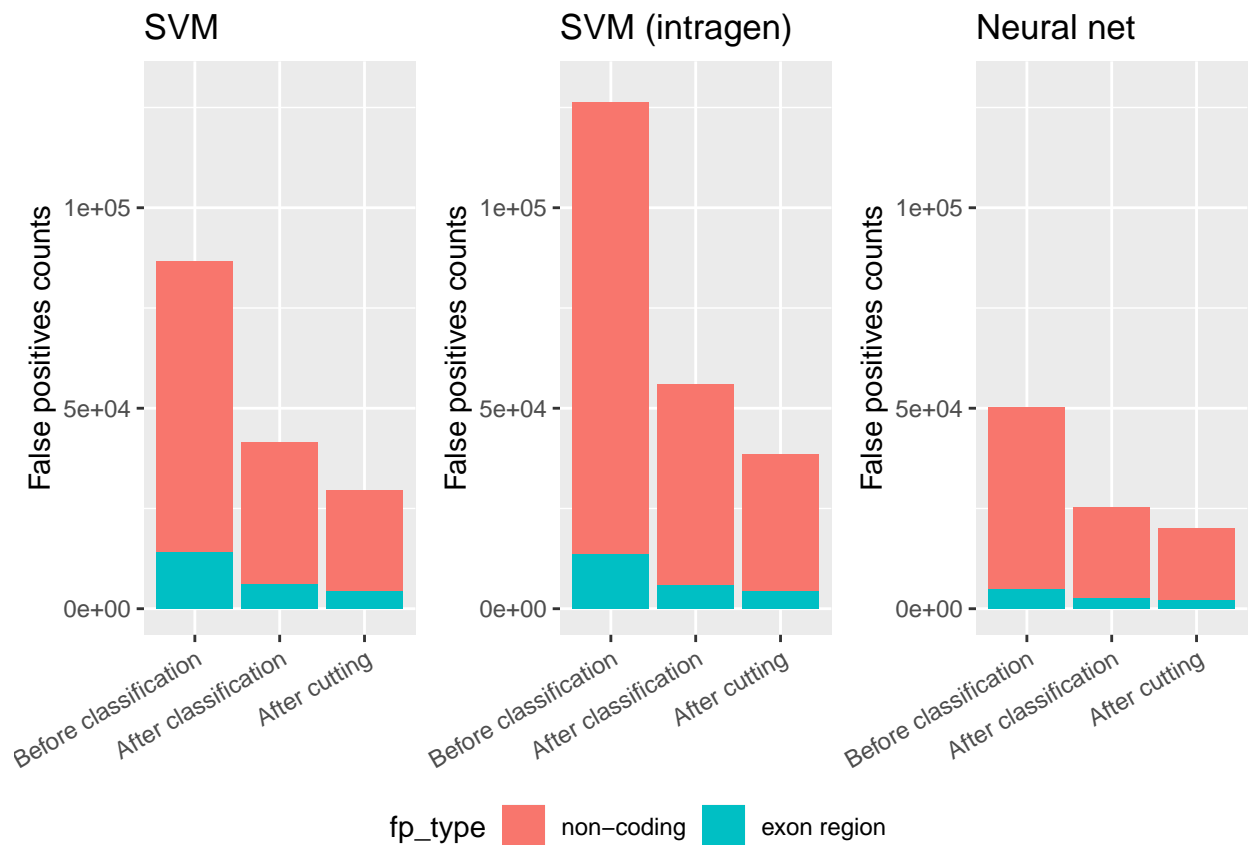
Zooming in and considering the intragenic FP reveals slight advantages of IG SVM. NN models are however still vastly superior to both of them.



## Follow-up intron classification (*Thega1*)

Similar results as in splice site classification.





## Misc

### Intron model recall

This chart shows how many introns we would be able to detect, if we had an intron model with perfect recall. I.e. we see here, how many true introns survived the splice site classification step (and pairing of GT-AG candidates). The orange line represents the relative 100% recall level, as we only pair GT-AG candidates in a range of 40-100nt. Longer or shorter introns are therefore not even included into the classification dataset, not speaking of being detected.

#### After splice site classification – intron survival rate

