



# **Xudong 2021 Summer Research Placement Summary**

# Technologies

- Python 3 standard libraries: numpy, pandas, sklearn, matplotlib
- Machine learning libraries: tensorflow and xgboost
- Also used: shap (<https://github.com/slundberg/shap>)
  - uproot (<https://uproot.readthedocs.io/en/latest/>)
  - hh4b\_utils (<https://gitlab.cern.ch/hh4b/hh4b-python-utils>)
  - mplhep (<https://github.com/scikit-hep/mplhep>)

# Files used

- mc files from /eos/user/l/lborgna/bbbb/NNT/cryptotuples-MAY21/mc/sm\_hh/mm\_hh\_pythia\_mc16a.root, mm\_hh\_pythia\_mc16d.root, mm\_hh\_pythia\_mc16e.root  
masks: pass\_vbf\_sel==False, X\_wt\_tag>=1.5, ntag>=4
- Data files from /eos/user/s/sgasioro/public/ggF\_push/quad\_yrOHE\_bkt\_Xhh/Xhh\_45/data16\_Xhh\_45\_NN\_100\_bootstraps.root, data17\_Xhh\_45\_NN\_100\_bootstraps.root, data18\_Xhh\_45\_NN\_100\_bootstraps.root  
masks: pass\_vbf\_sel==False, X\_wt\_tag>=1.5, ntag==2, rw\_to\_4b==True

# Some initial results

At the beginning, the more features, the higher the AUC (and accuracy)

- 'X\_hh','dEta\_hh','njets','m\_hh\_cor' , 'X\_wt\_tag', 'eta\_i','year'

AUC ~ 0.70-0.78

- 'm\_hh','X\_hh','dEta\_hh','njets','X\_wt\_tag','cos\_theta\_star','X\_wt\_notag','bkt\_HT','pt\_hh','pT\_2','pT\_4','eta\_i','dRjj\_1','dRjj\_2','year'

AUC ~ 0.85-0.90

- Basically all features (55)

AUC ~ 0.95-0.96

# Some initial results

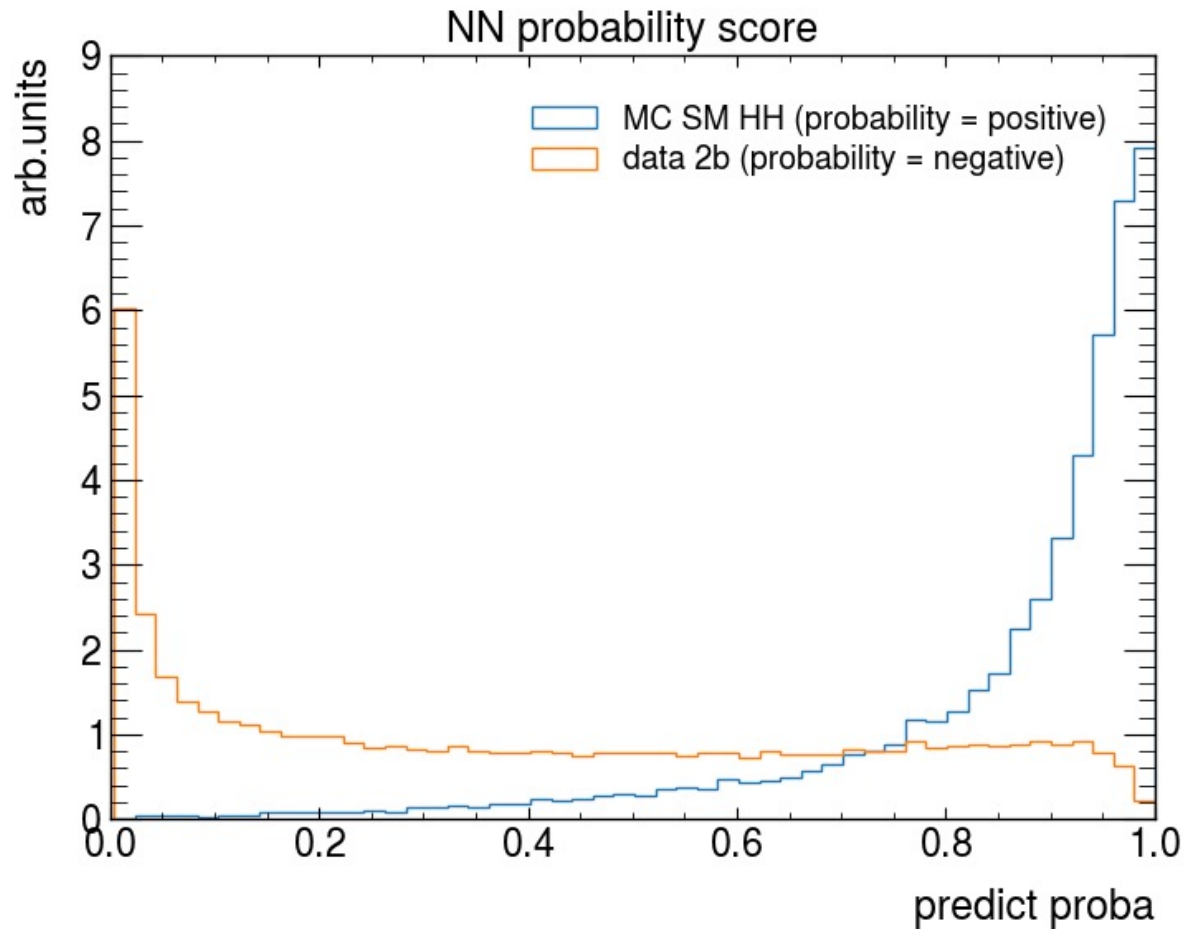
The 55 features:

'm\_hh', 'X\_hh', 'dEta\_hh', 'njets', 'X\_wt\_tag', 'cos\_theta\_star', 'X\_wt\_notag',  
'bkt\_HT', 'pt\_hh', 'pT\_2', 'pT\_4', 'eta\_i', 'dRjj\_1', 'dRjj\_2', 'm\_min\_dj',  
'm\_max\_dj', 'pairing\_score', 'bkt\_lead\_jet\_pt', 'bkt\_third\_lead\_jet\_pt', 'year',  
'm\_h1', 'E\_h1', 'pT\_h1', 'eta\_h1', 'phi\_h1', 'm\_h1\_j1', 'E\_h1\_j1', 'pT\_h1\_j1',  
'eta\_h1\_j1', 'phi\_h1\_j1', 'angle\_h1\_j1'

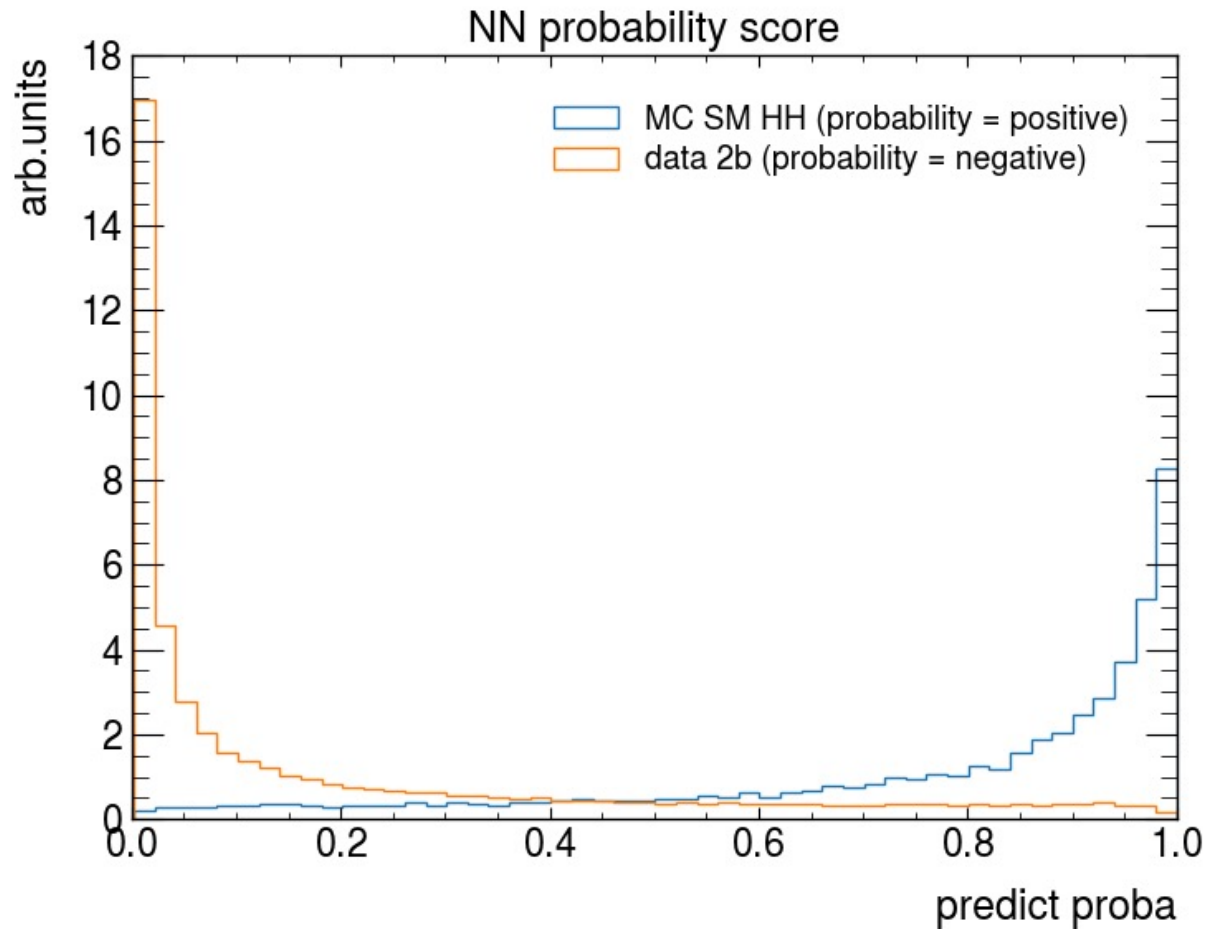
Note that for m\_h1, there is also m\_h2, and so on

# Some initial results

TensorFlow

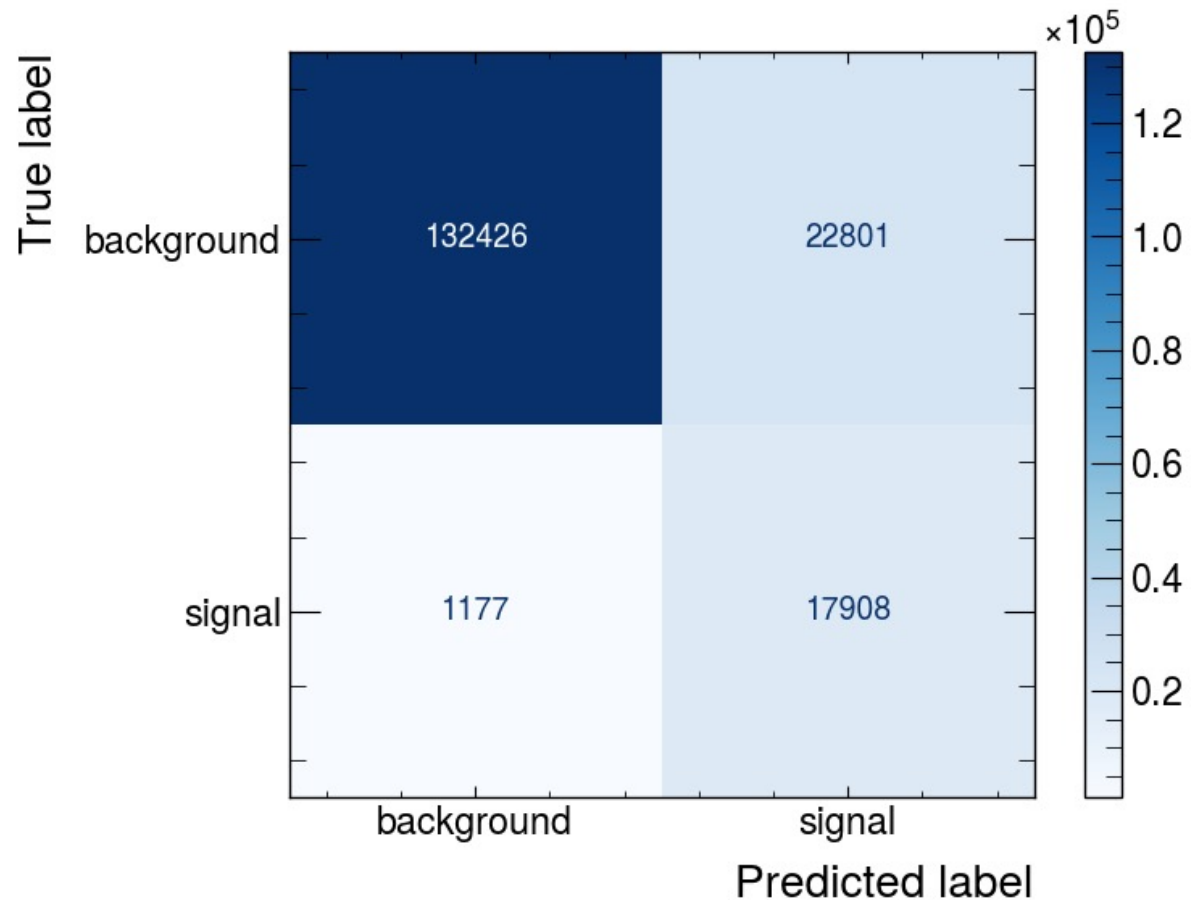


XGBoost

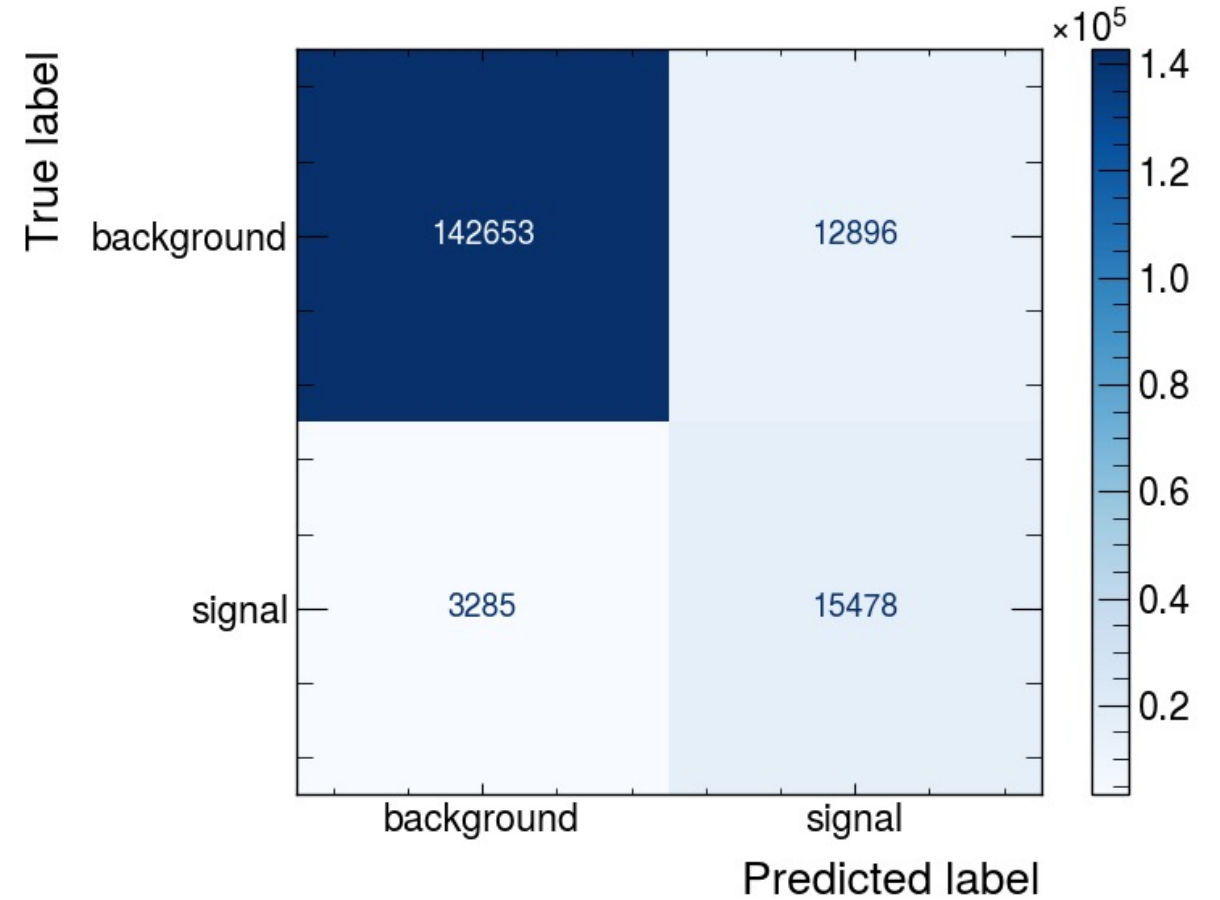


# Some initial results

TensorFlow



XGBoost



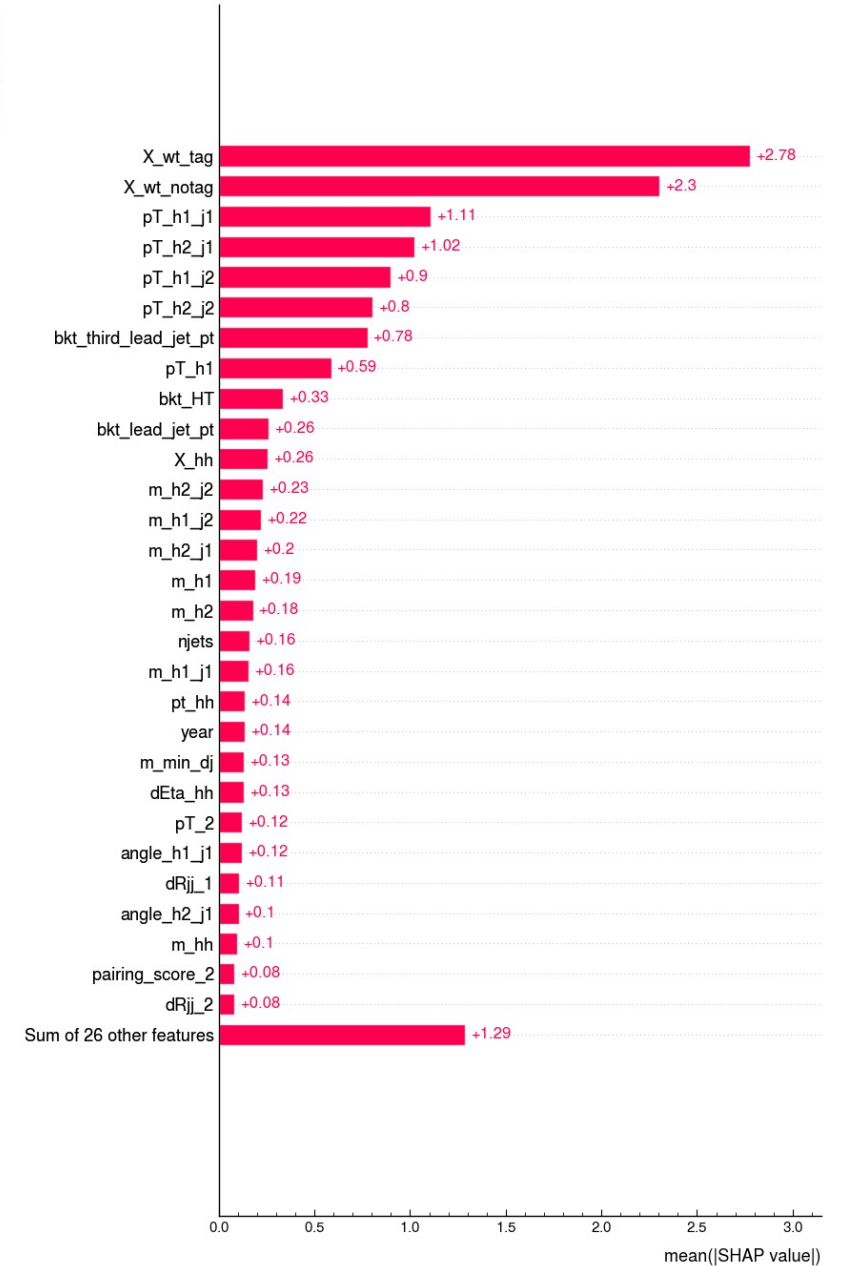
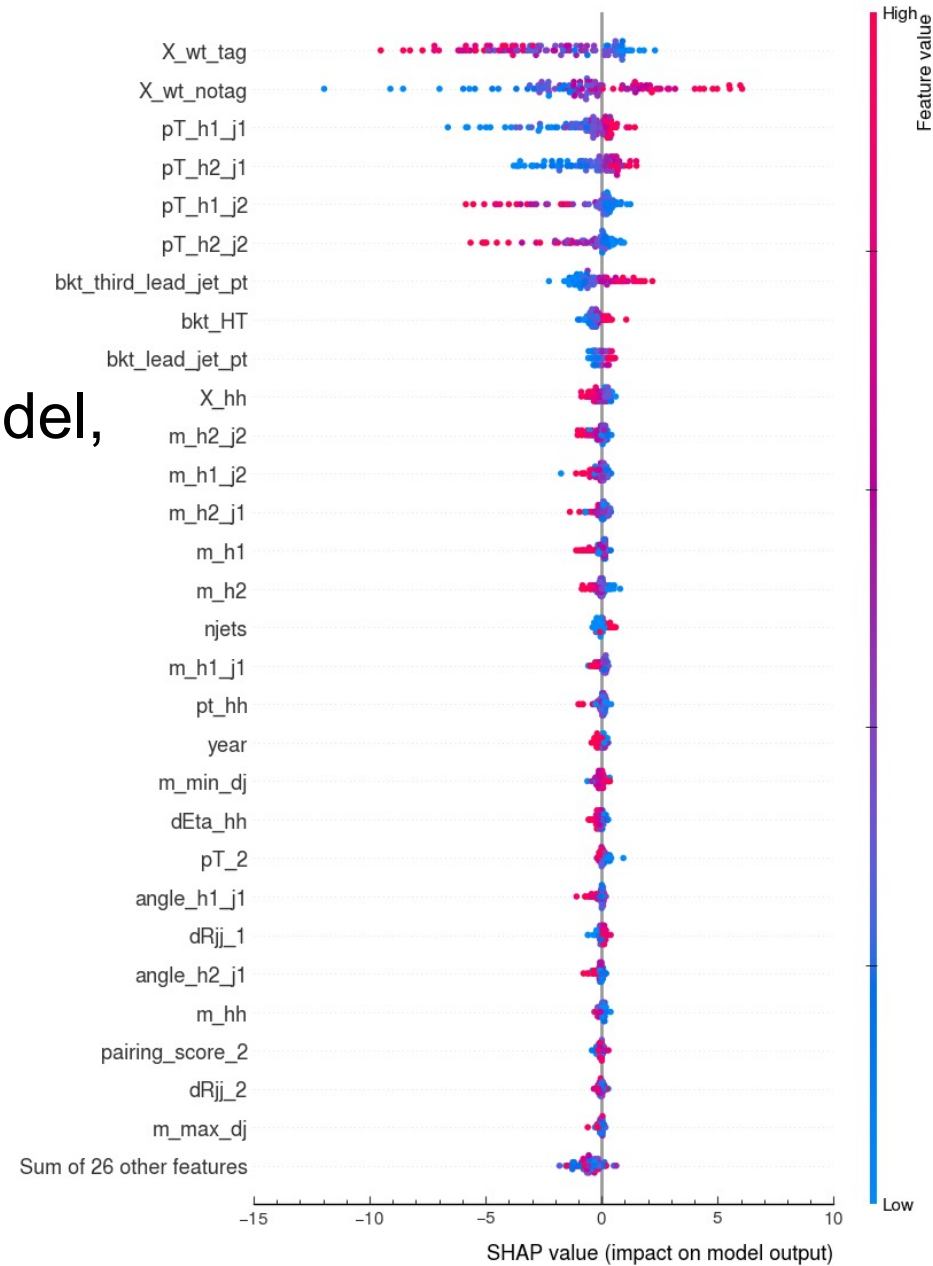
# Some initial results

- Predict signal if only very confident → higher precision
- Avoid missing too many cases of signal → higher recall
- F1 score: weighted average of precision and recall

	Precision	Recall	F1 Score
TensorFlow	0.440	0.938	0.599
XGBoost	0.545	0.825	0.657



From the XGboost model,  
using SHAP



# Some initial results in CR

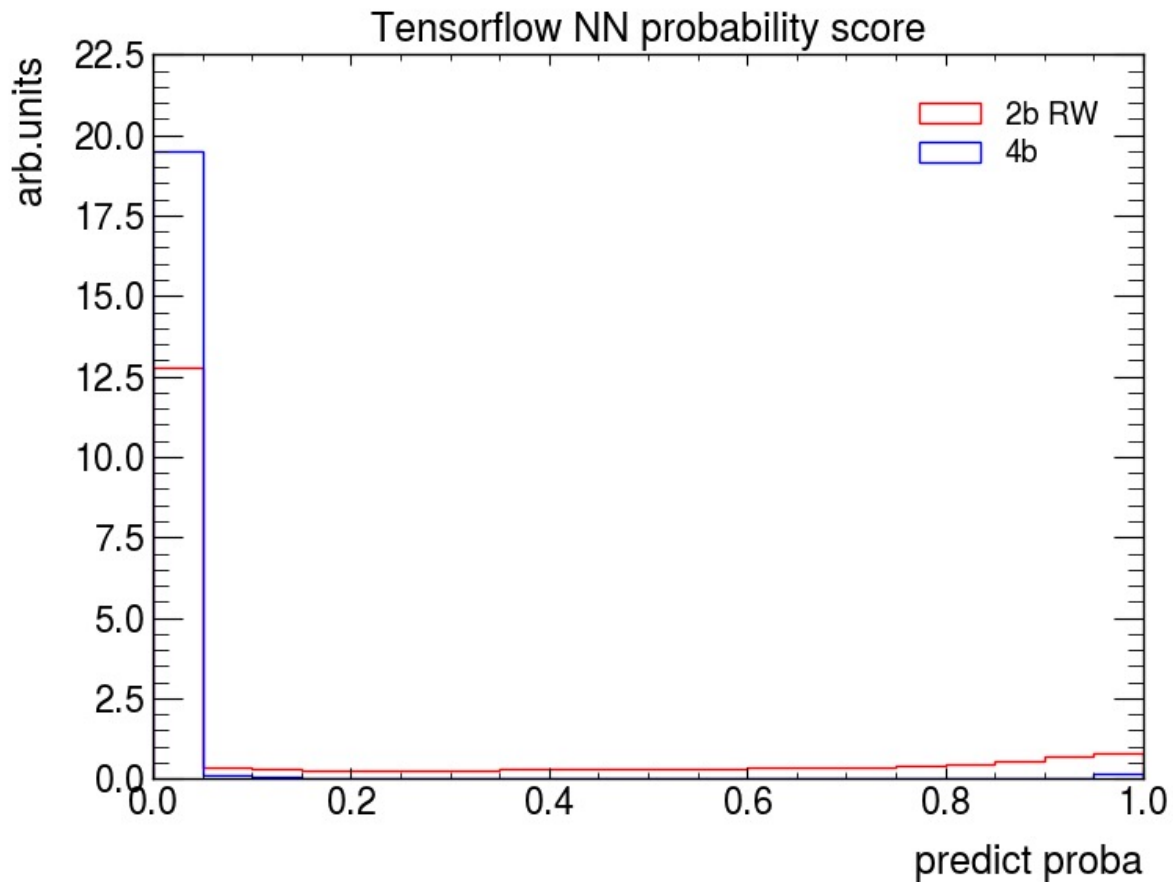
Masks:

- 2b RW: `pass_vbf_sel==False`, `X_wt_tag>=1.5`, `ntag==2`, `rw_to_4b==True`
- 4b: `pass_vbf_sel==False`, `X_wt_tag>=1.5`, `ntag>=4`

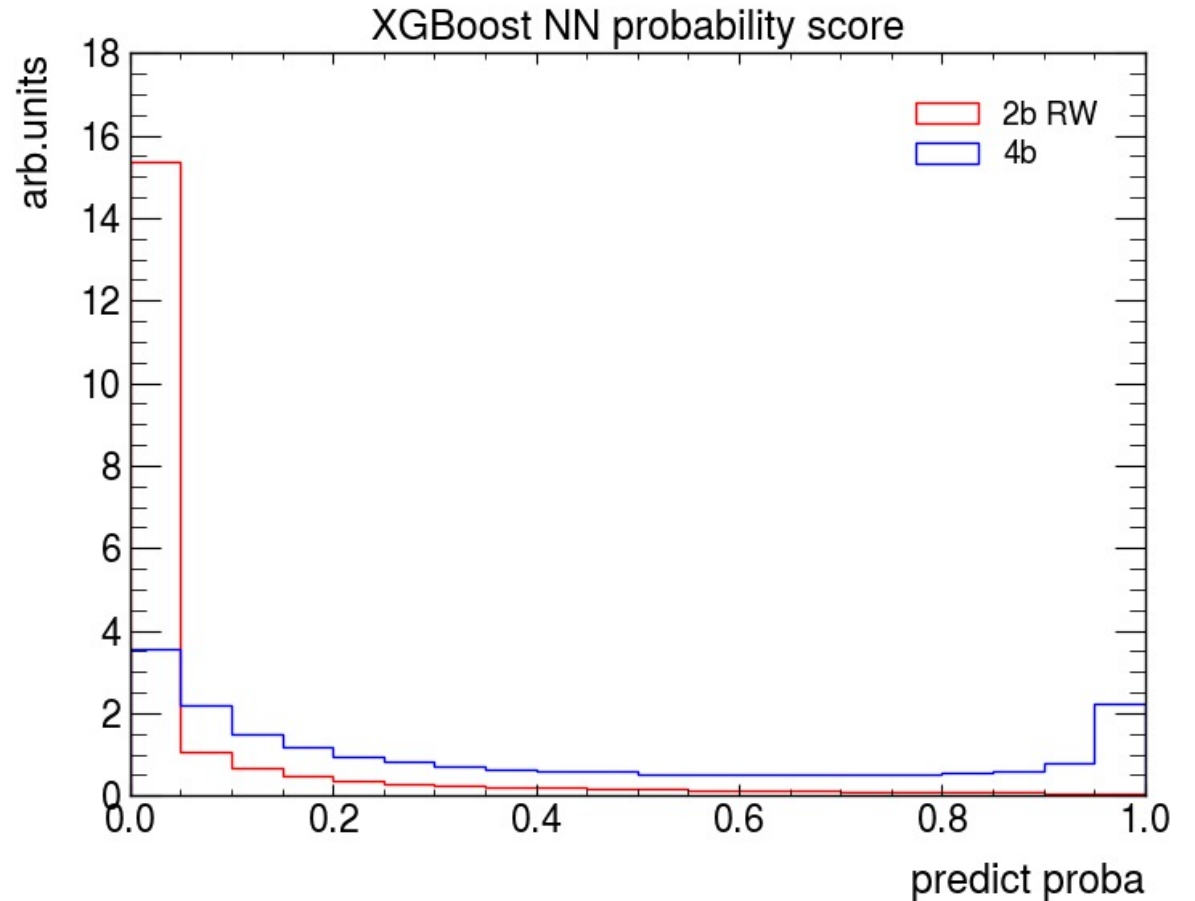
However, 2bRW and 4b did not agree well

# Some initial results in CR

TensorFlow

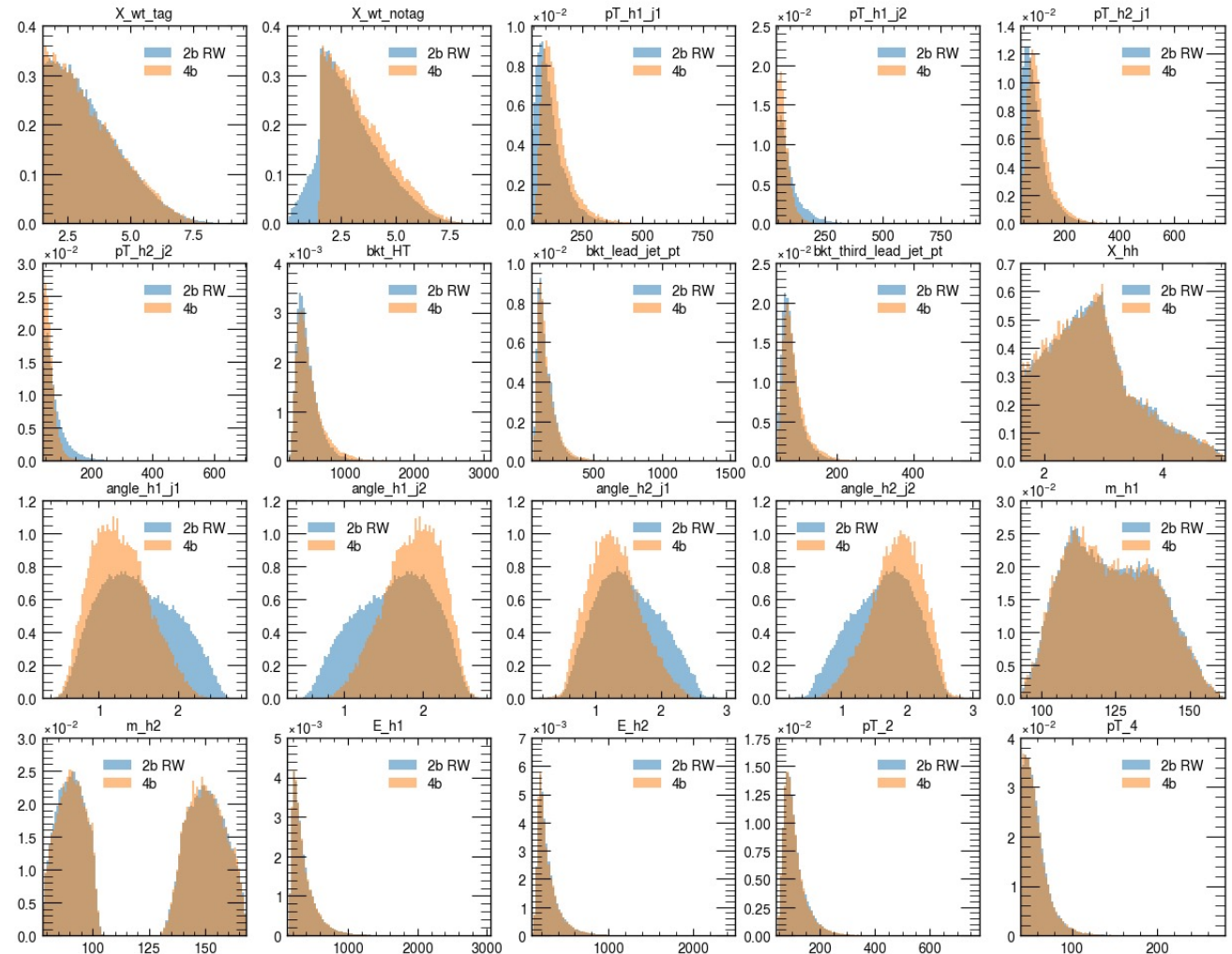


XGBoost



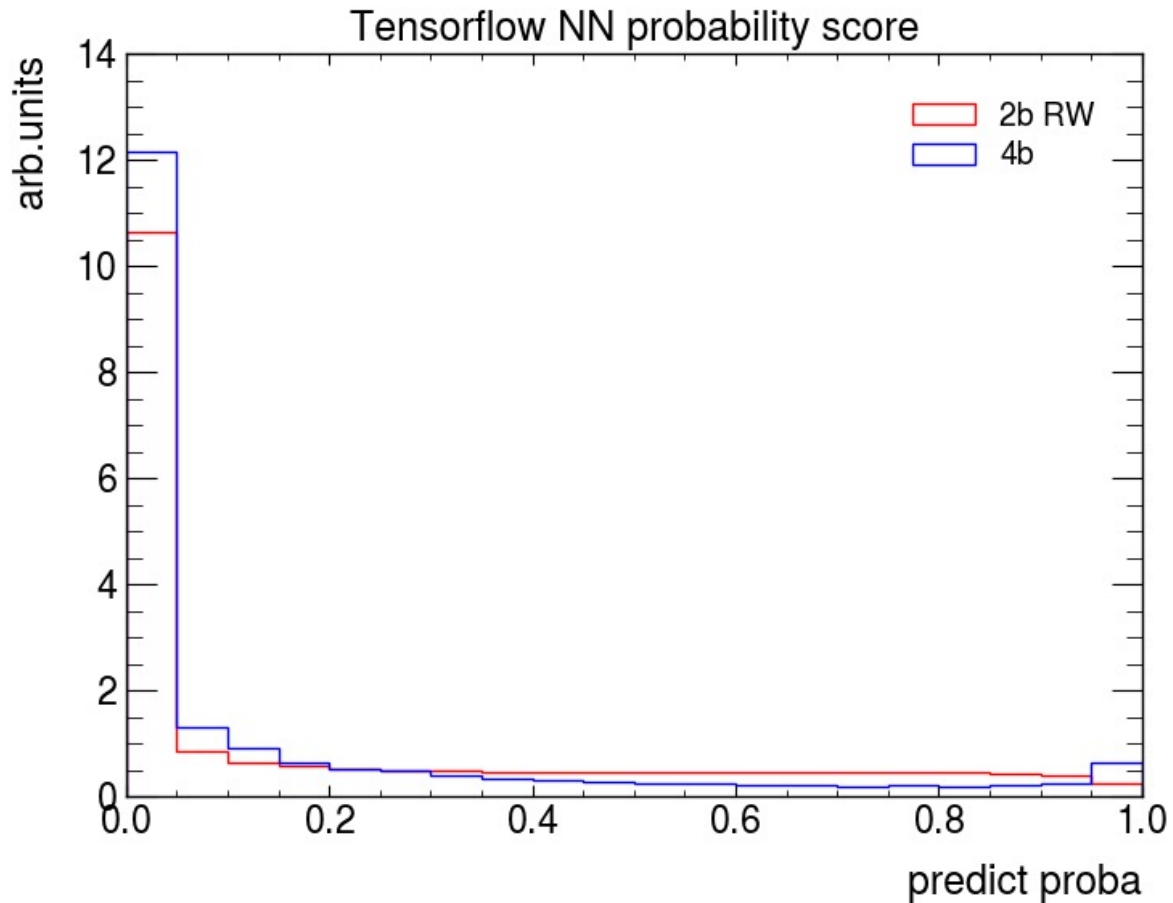
# CR features: 2bRW vs 4b

- The most important features according to SHAP values
- Some features like `X_wt_notag` or the angles are **very different**, so these were removed

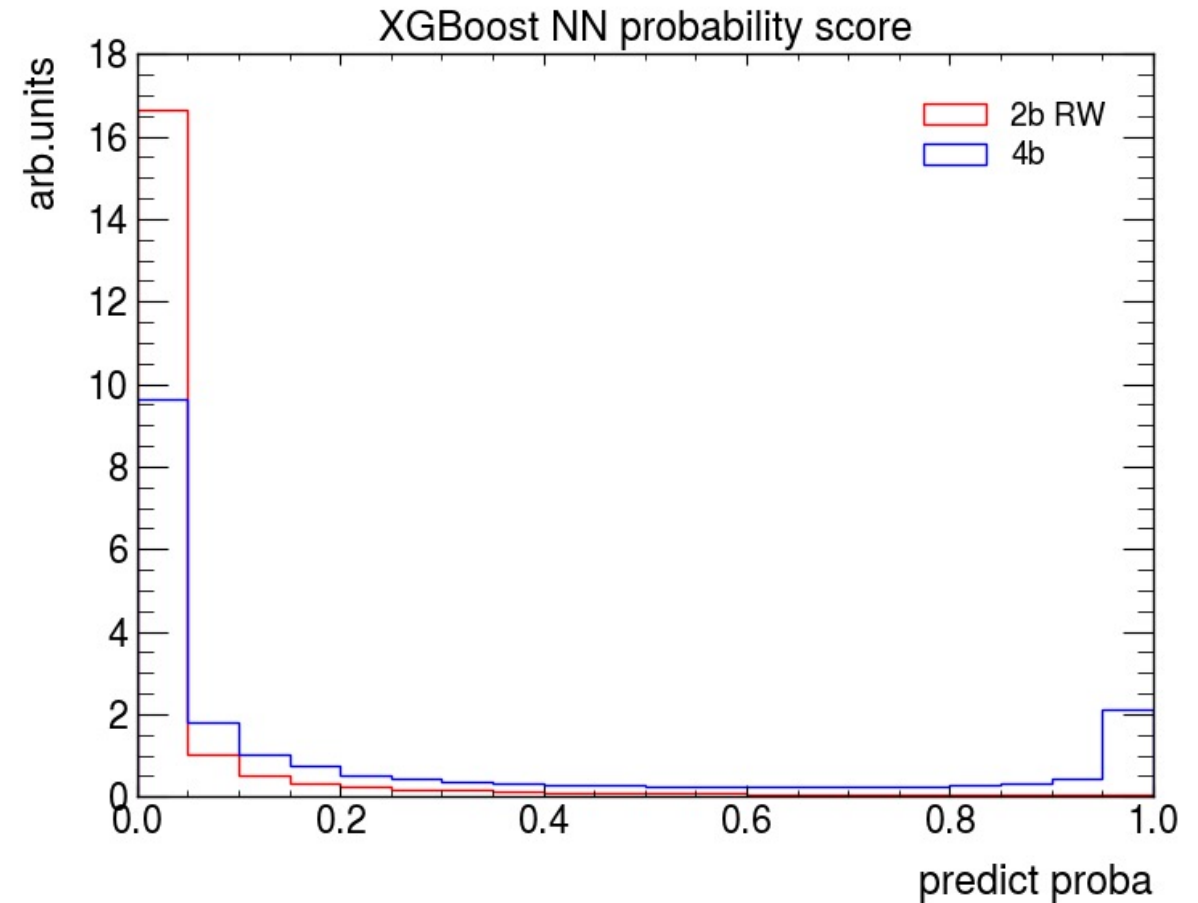


# Removed very different features in CR

TensorFlow



XGBoost

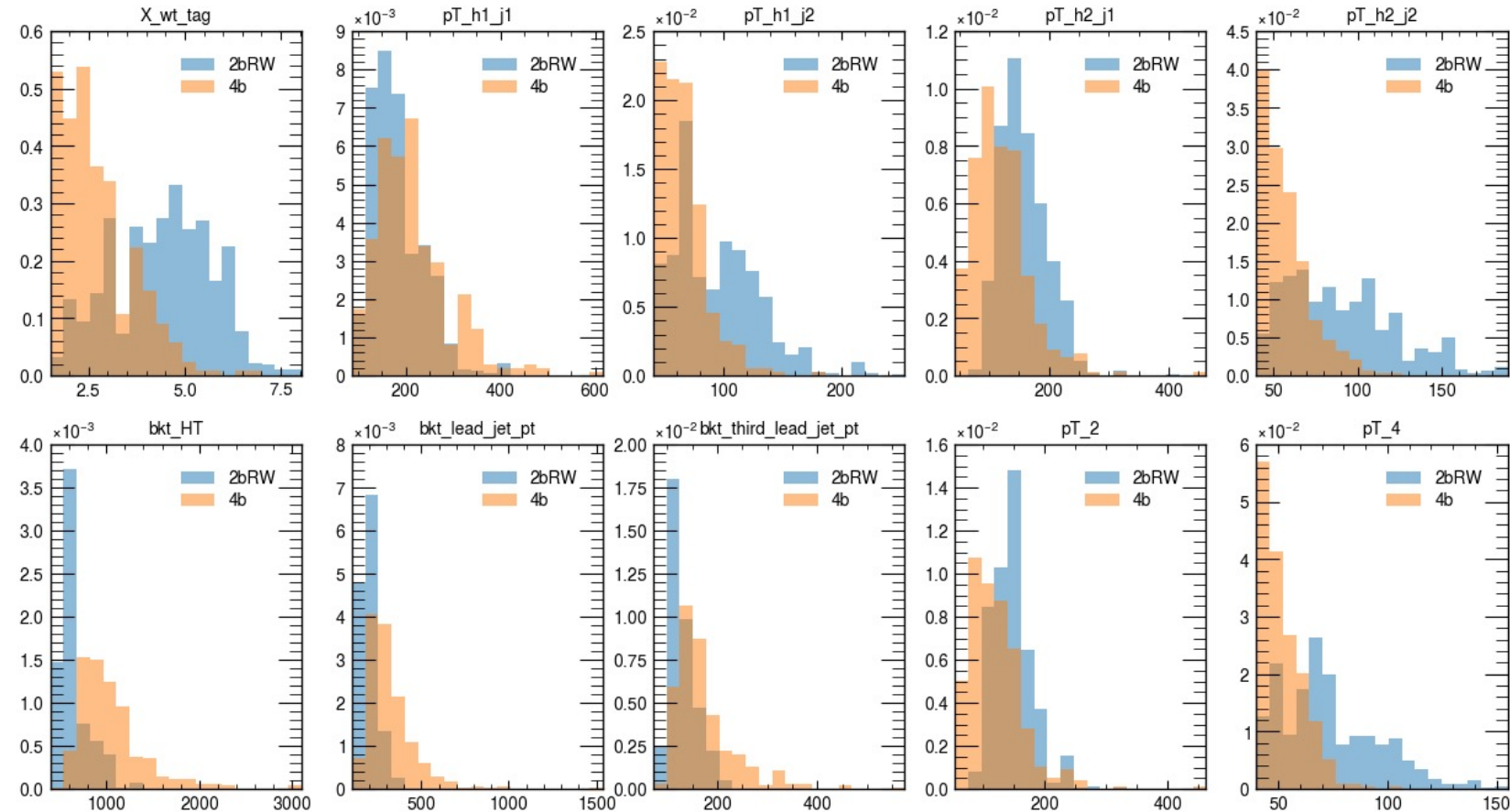


# Removed very different features in CR

- After removing ['angle\_h2\_j2', 'angle\_h1\_j2', 'angle\_h2\_j1', 'angle\_h1\_j1', 'X\_wt\_notag'], 2bRW and 4b NN scores agree more
- However, XGboost models are not agreeing as well as TensorFlow models
- It is also more difficult to tune hyperparameters with XGBoost models, it is easier for TensorFlow models
- Decided to use **TensorFlow models only**
- Also found **small peak for 4b at signal (1.0)**

# Removed very different features in CR

- Compared events from TensorFlow models where predictions were greater than 0.99
- Decided to **remove more variables** bkt\_HT and pT\_h1\_j1 type features, and potentially X\_wt\_tag

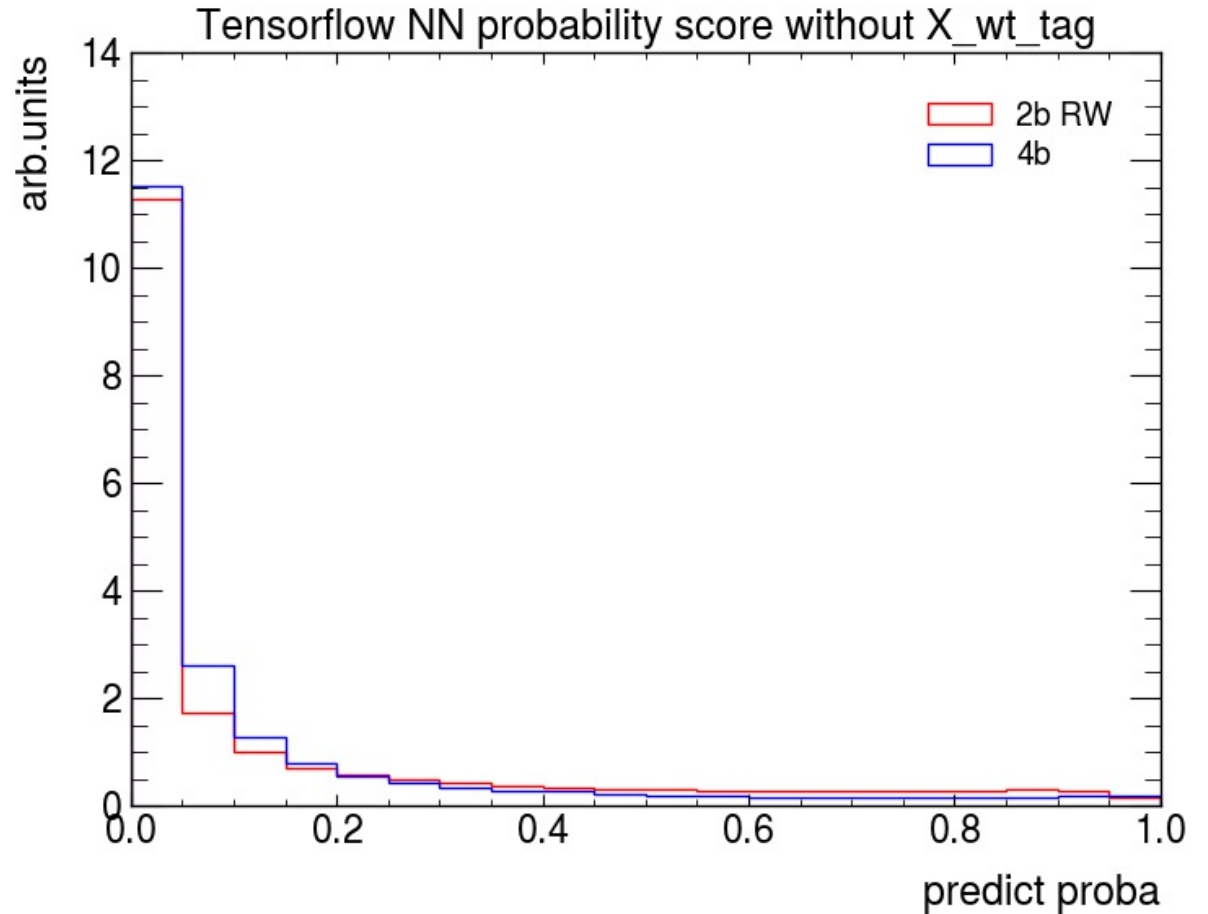
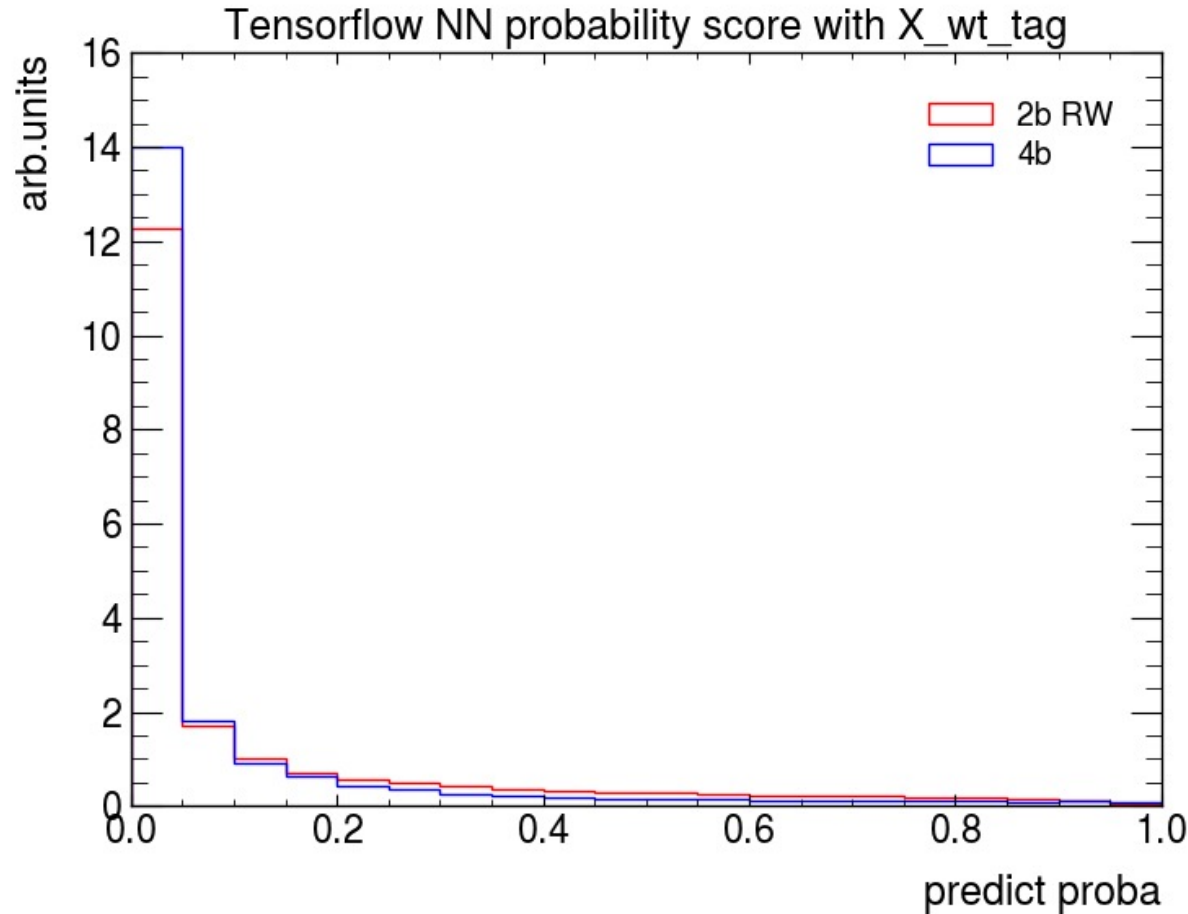


# Hyperparameter tuning

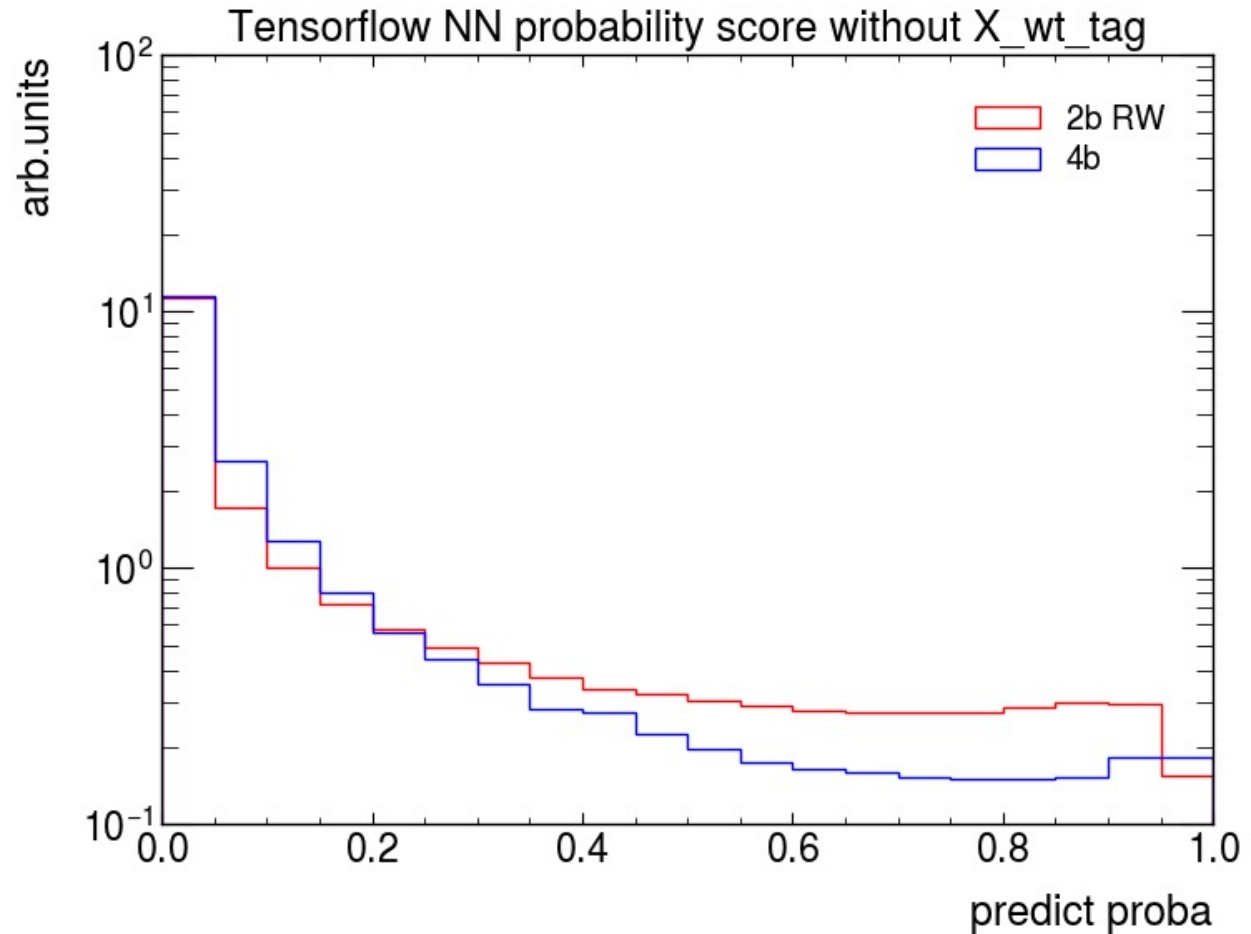
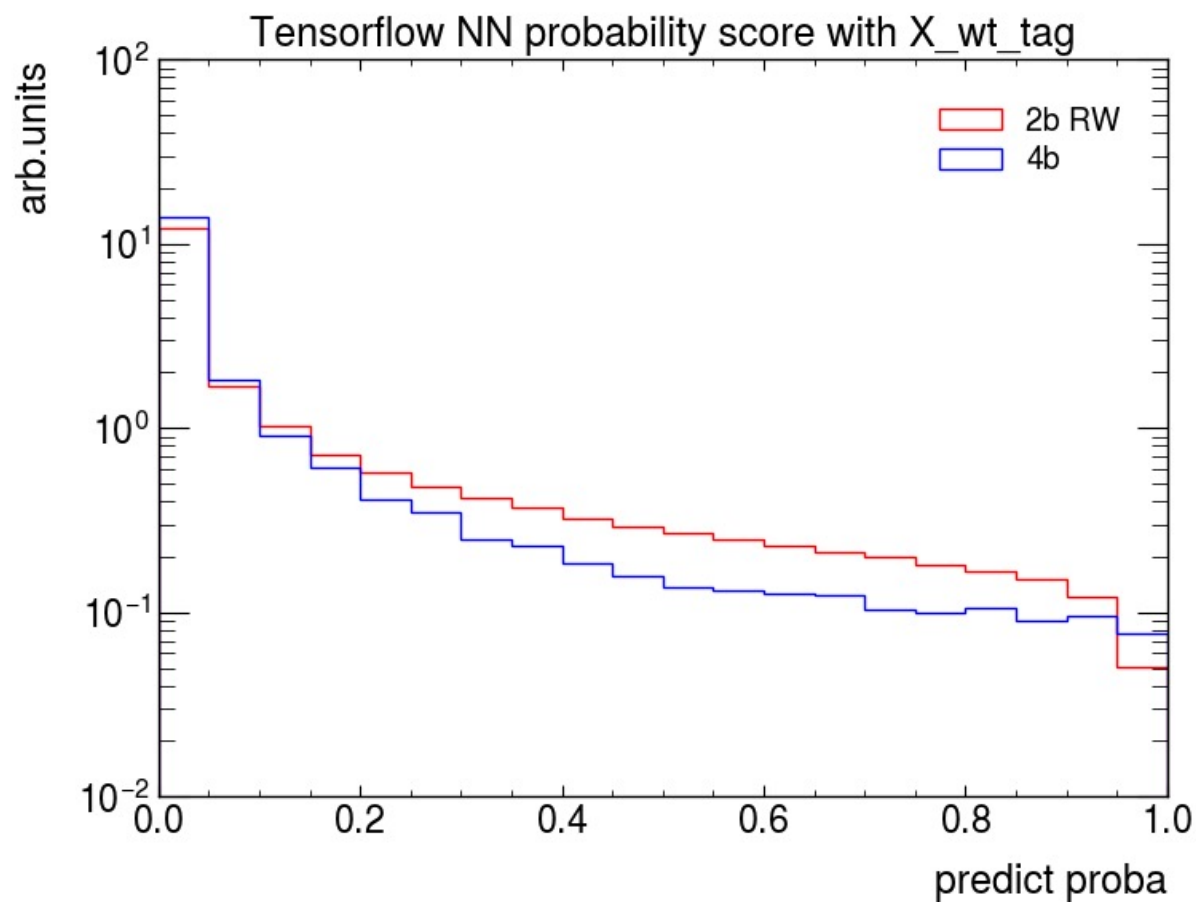
- TensorFlow models' hyperparameters can be tuned quite easily using **KerasTuner**. I decided to use **Bayesian Optimisation**.
- The model of the network is: Dense\_0 + Dropout\_0 + Dense\_1 + Dropout\_1 + ... + Dense\_output
- The number of hidden layers (Dense+Dropout) is a tunable parameter
- The number of units, the activation functions, the rate of dropouts and the learning rate for the optimiser (adam) are tunable
- Look for 100 trials (100 different combinations) and 3 executions per trial to reduce variance in the results.



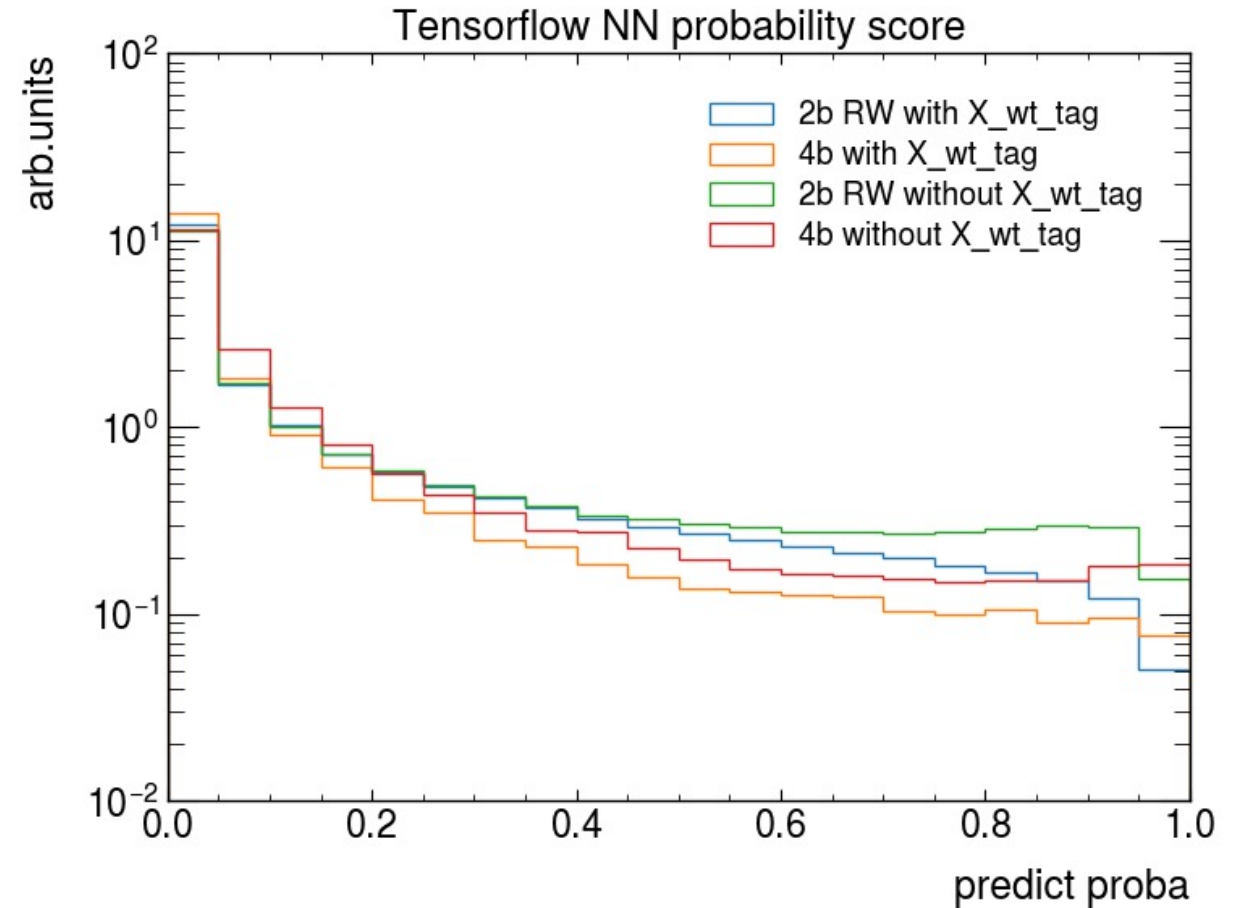
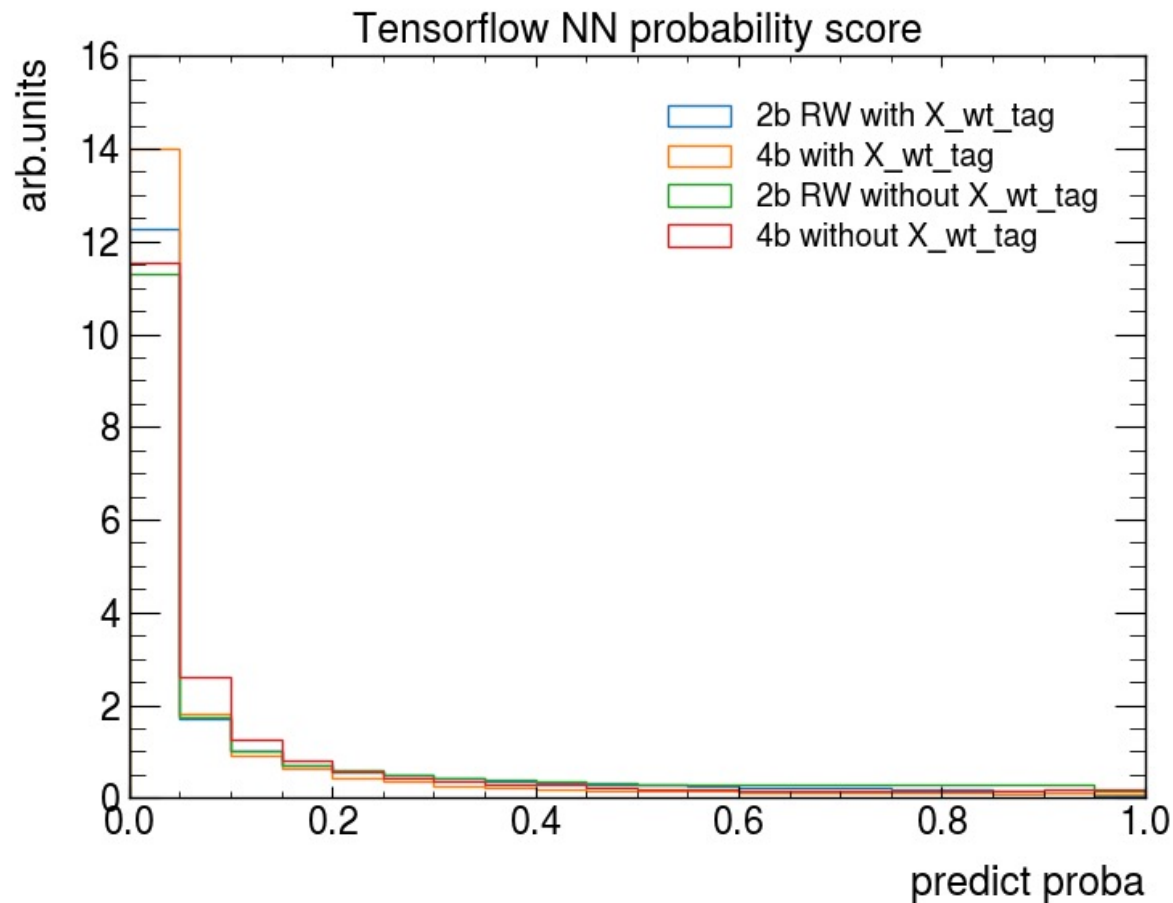
# In the control region updated



# In the control region updated



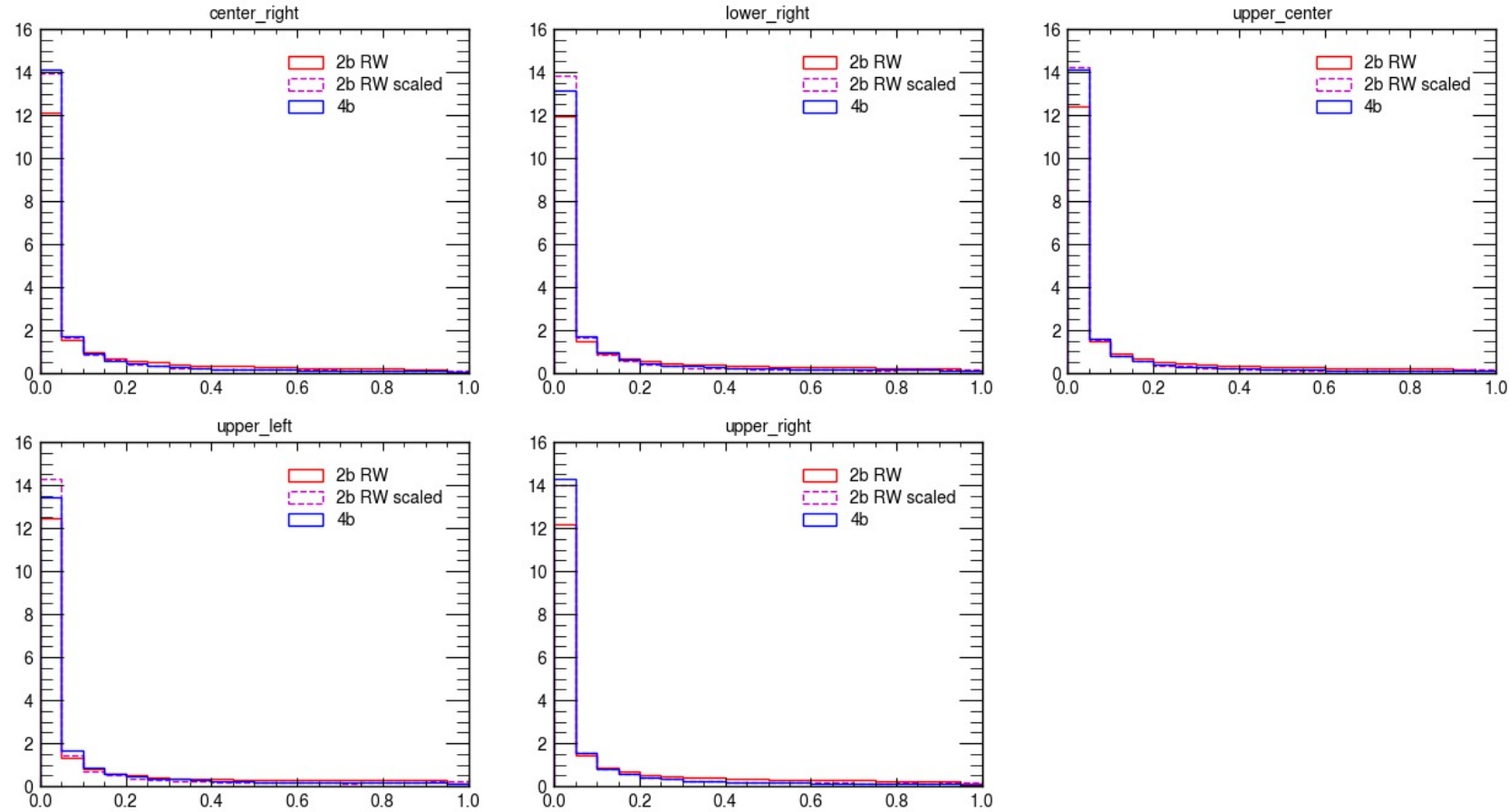
# In the control region updated



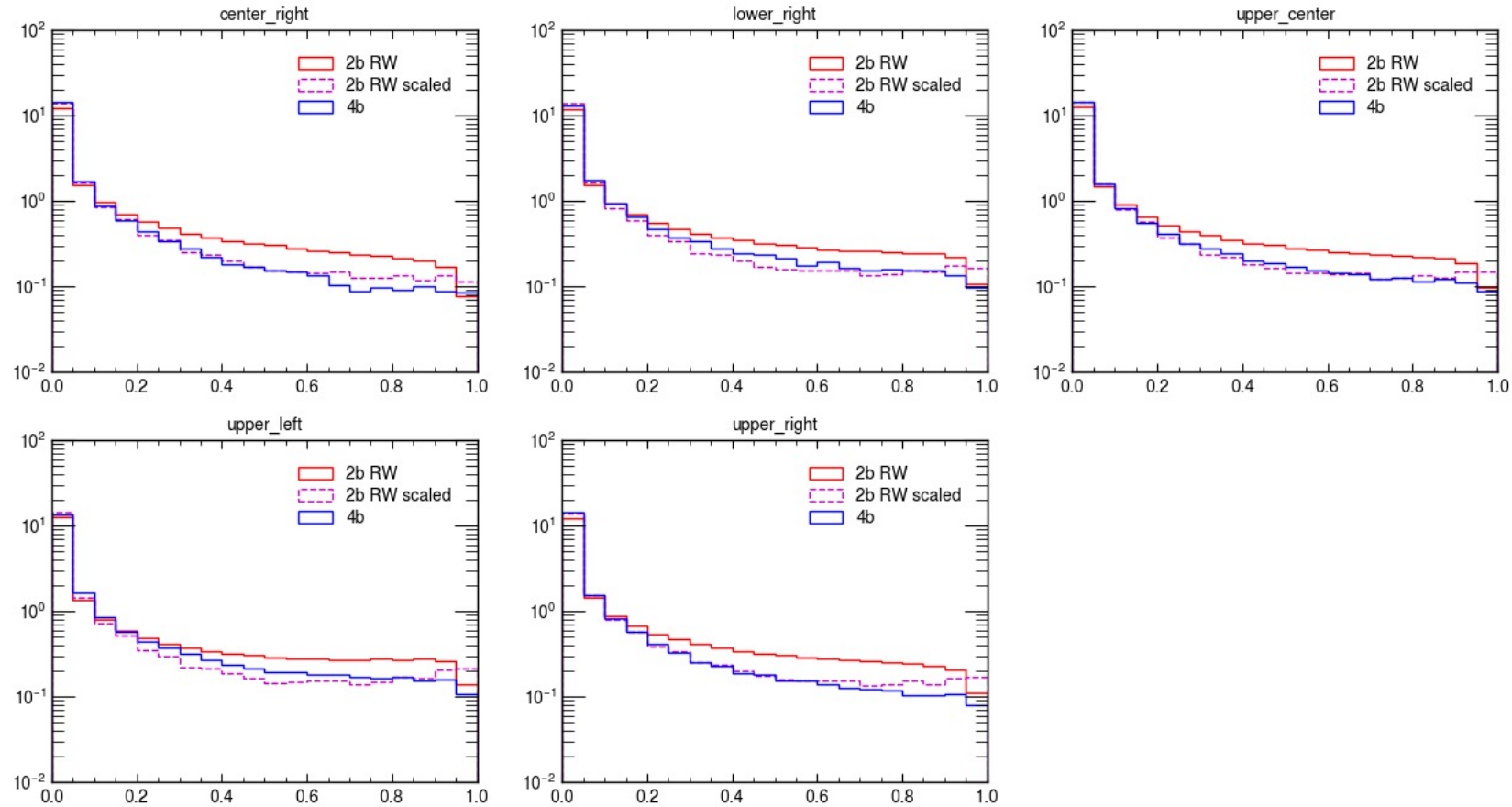
# In the control region updated

- The best model using  $X_{wt\_tag}$  outperformed and the 4b signal peak was also solved
- Found corrections for 2bRW by dividing the histogram weights in the CR and tested these in shifted regions from `/eos/user/s/shhayash/HH4b/HH4b-ShiftReg-Validation/shifted/Data/`

# In the shifted regions



# In the shifted regions



# Other tests

- I also tried including  $pT_{h1_j1}$  type variables, it sometimes worked well but not always, it depends how the model was actually trained. One could explore this further. It is difficult to say which features caused that 4b peak at 1.0.
- One thing that I tried but didn't work well was Adversarial trained Neural Networks. I followed some works from other people and tried to de-correlate the masses  $m_{h1}$  and  $m_{h2}$ . However, considering that normal Neural Networks already perform well in shifted regions, one could leave this for the future.

# Future steps

- Select the best models from KerasTuner and validate them in the control region. The second or third best model could work better than the best model. The metric used is not trivial.
- Include pT\_h1\_j1 type features (with X\_wt\_tag). More analysis could be done and again, the metric is not trivial.
- Explore more on XGboost models, they worked well at the beginning, but it is more difficult to tune the parameters.
- Use Adversarial Neural Networks to see if further improvement can be made in other regions.



# Resources

Code:

- <https://github.com/Xudong-Ke-Lin/HEP-HH4b-machine-learning>

Useful Adversarial Neural Networks code:

- CERN jet tagger trained adversarially <https://github.com/asogaard/adversarial/wiki>