

Update of SmartBKG -- Improved Selective Background Monte Carlo Simulation at Belle II with Graph Attention Networks and Weighted Events

Boyang Yu¹, Nikolai Hartmann¹, Thomas Kuhr¹

¹ *Ludwig-Maximilians-Universität München*

KISS B2, Sept 28, 2023

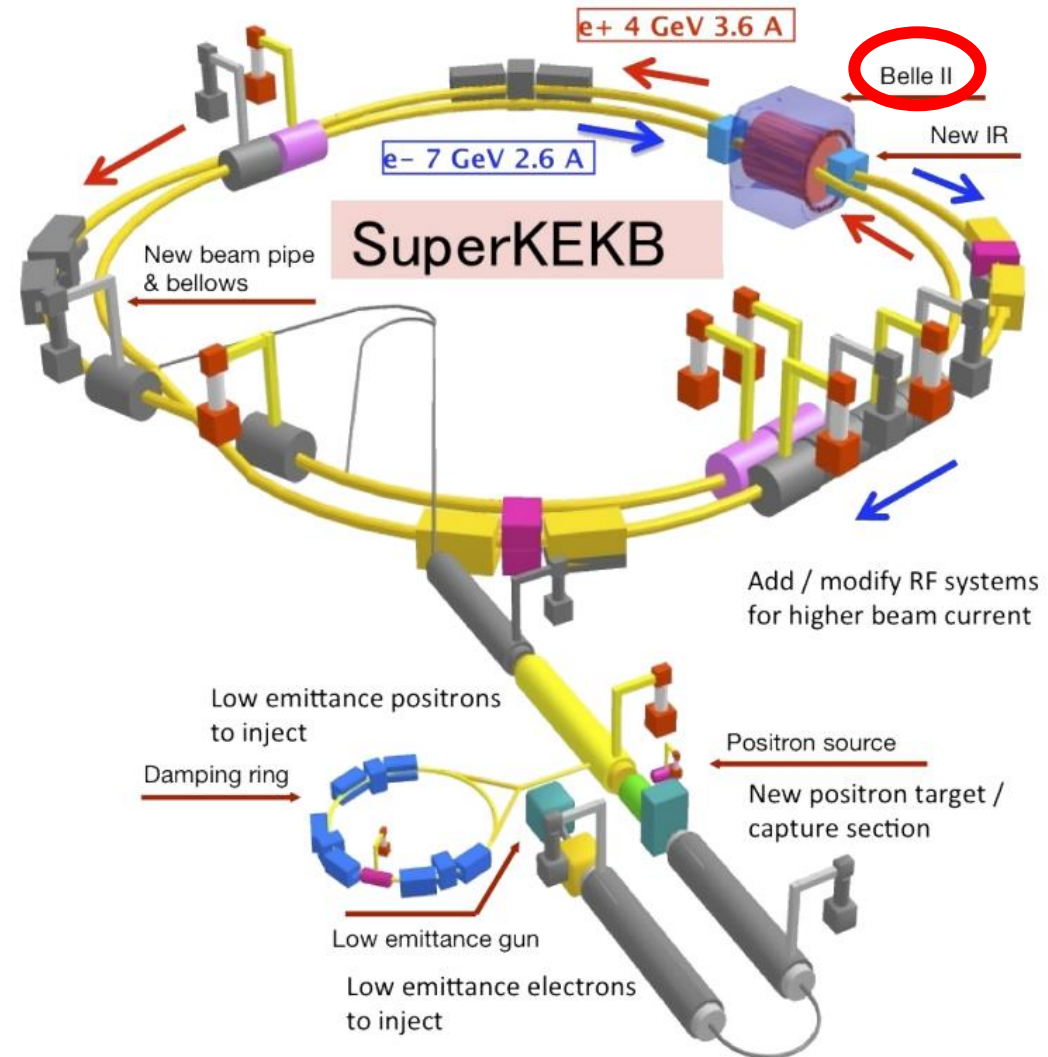


Bundesministerium
für Bildung
und Forschung



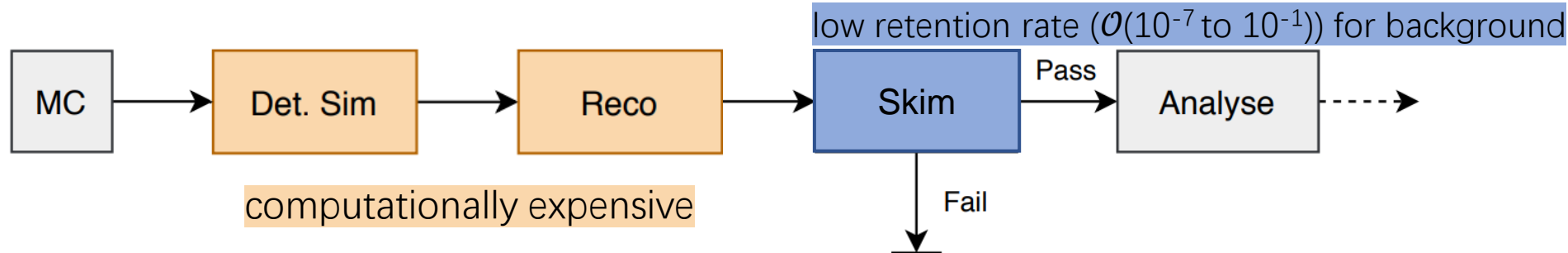
Belle II Experiment:

- At SuperKEKB
 - Electron-positron collider
 - Centre-of-momentum energy close to the mass of $Y(4S)$ resonance to mainly produce B mesons
 - Located in Tsukuba, Japan
- Detector for reconstruction and identification of charged and neutral particles
- Search for new physics
- World's highest luminosity
- Huge MC dataset for analysis

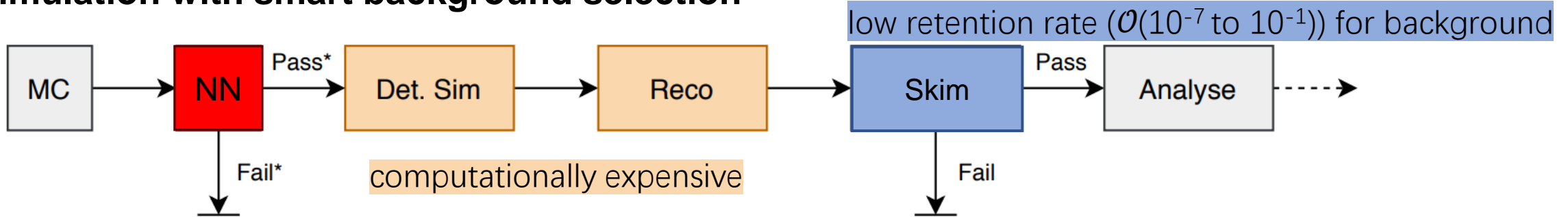




Normal Monte Carlo Simulation data flow



Simulation with smart background selection



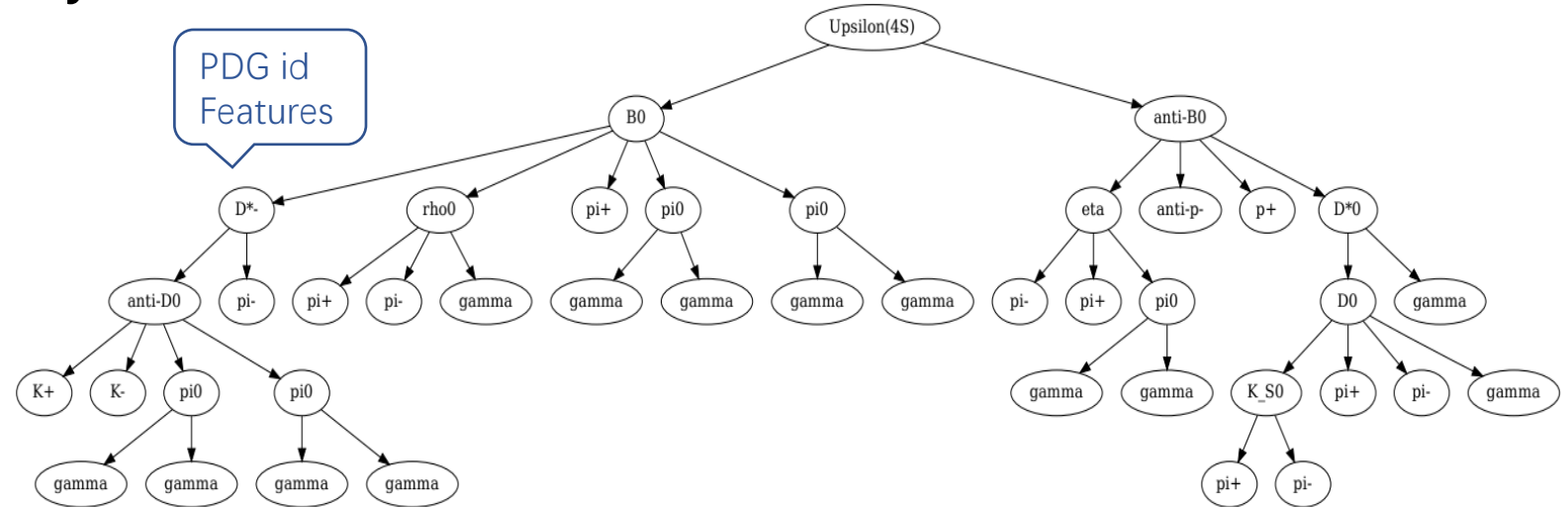
Previous Works:

- **PhD Thesis:** *Hadronic Tag Sensitivity Study of $B \rightarrow K^{(*)} \nu \bar{\nu}$ and Selective Background Monte Carlo Simulation at Belle II*, James Kahn, 2019
- **Talk:** *Selective background Monte Carlo simulation at Belle II*, James Kahn, CHEP 2019

Tree Structures of Particle Decay



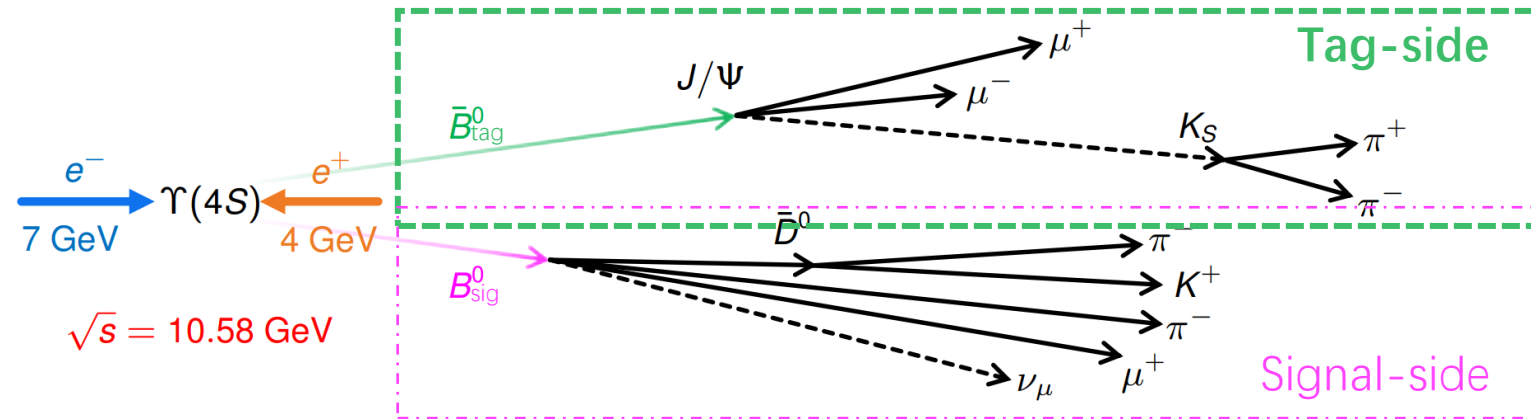
Graph Neural Network



Dataset:

- Each event (each **Graph**):
 - Decay of $\Upsilon(4S) \rightarrow B^0 \bar{B}^0$
 - Particles (**Nodes**)
 - Mother/Daughter relations (two way **Edges**) + self loops
- ❑ Each particle (each **Node**)
 - **PDG id**
 - **8 Features: Production time, Energy, Position (3d), Momentum (3d)**
- **Label** per event: Pass/Fail after the skims
 - * FEI Hadronic B0, retention rate 4.25%
- Other event level **attributions** for further analysis: e.g. M_{bc} , ΔE etc.

Tagging method:



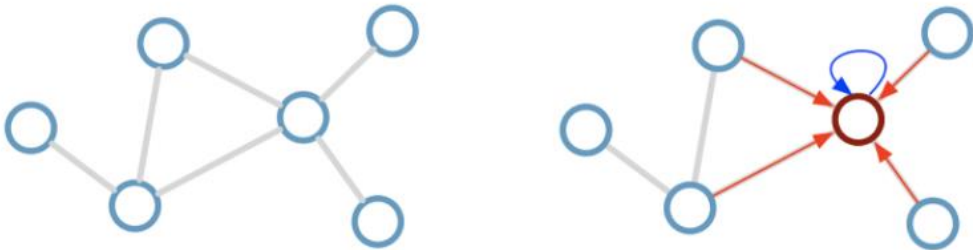
Retention rate after reconstruction and selection of tag-side B candidate:

	Hadronic B^+	Hadronic B^0
Mixed ($\Upsilon(4s) \rightarrow B^0 \bar{B}^0$)	5.62%	4.25%



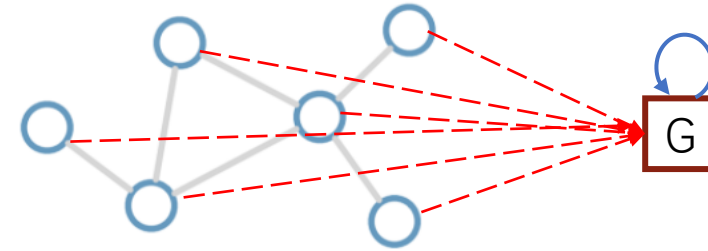
Updating node features:

Graph Convolutional Networks (GCN)
-> Graph structure remains



Updating global features:

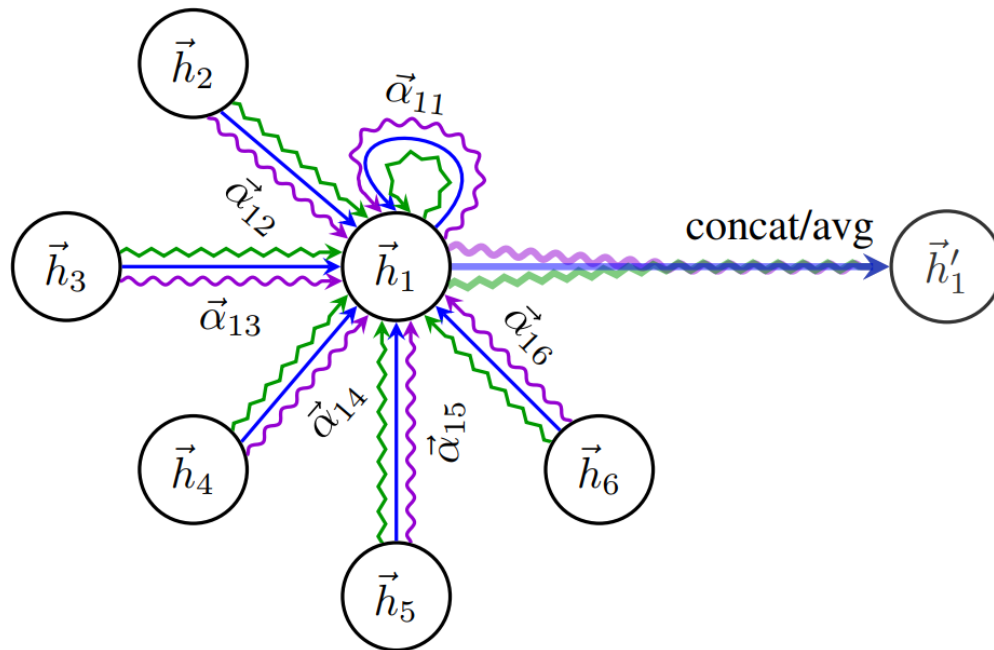
Global Average Pooling
-> Graph structure degenerated



Improvement with attention mechanism

Graph Attention Networks (GAT)

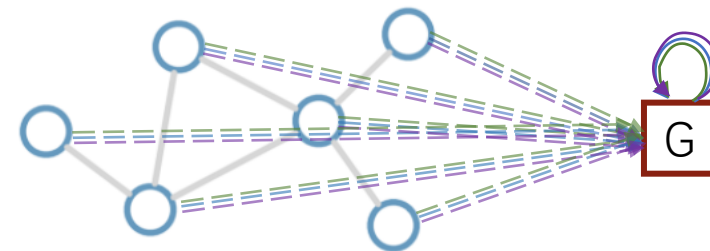
-> Graph structure remains



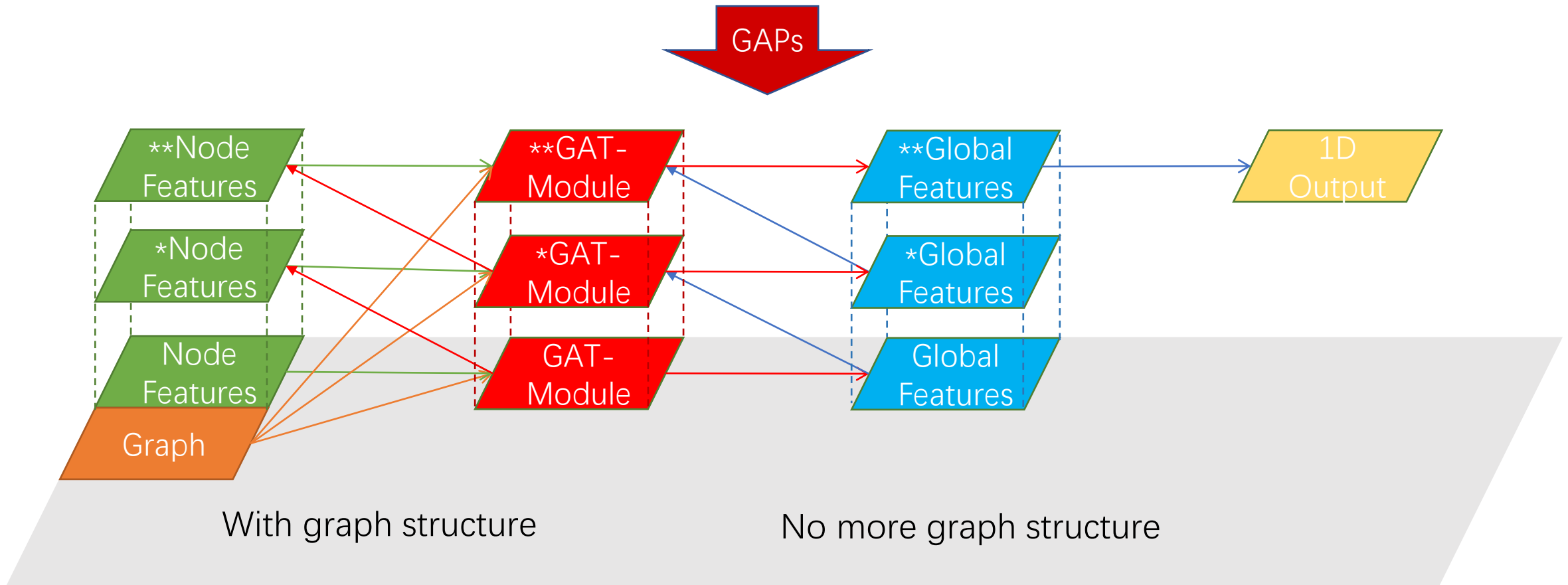
Each head (color) represents a different set of attention weights

Global Attention Pooling (GAP)

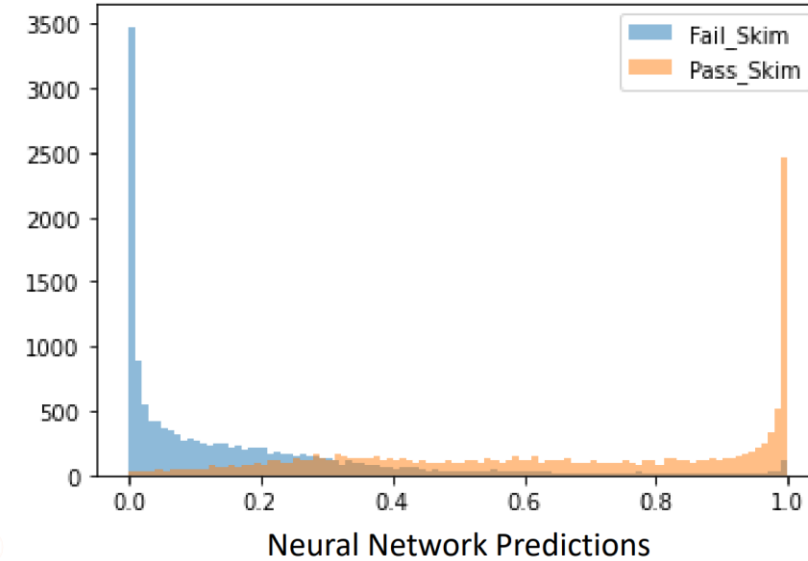
-> Graph structure degenerated



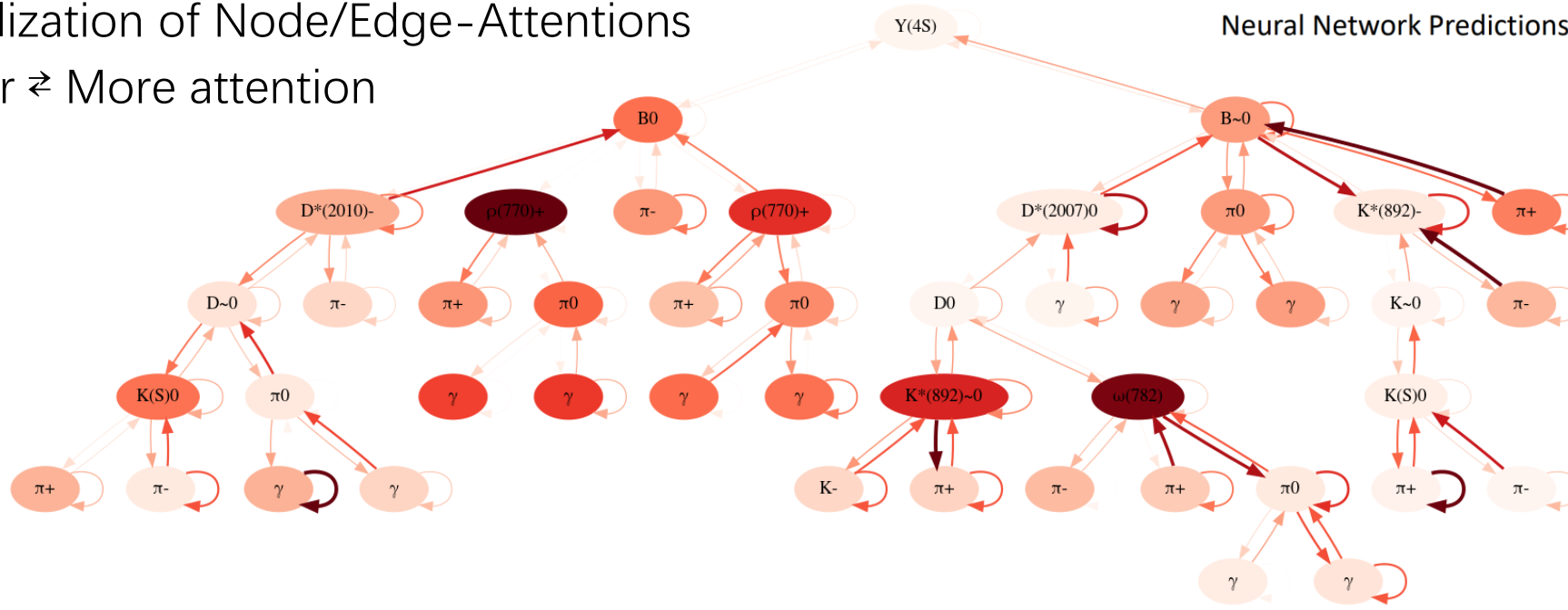
Final Architecture: GAT+GAP



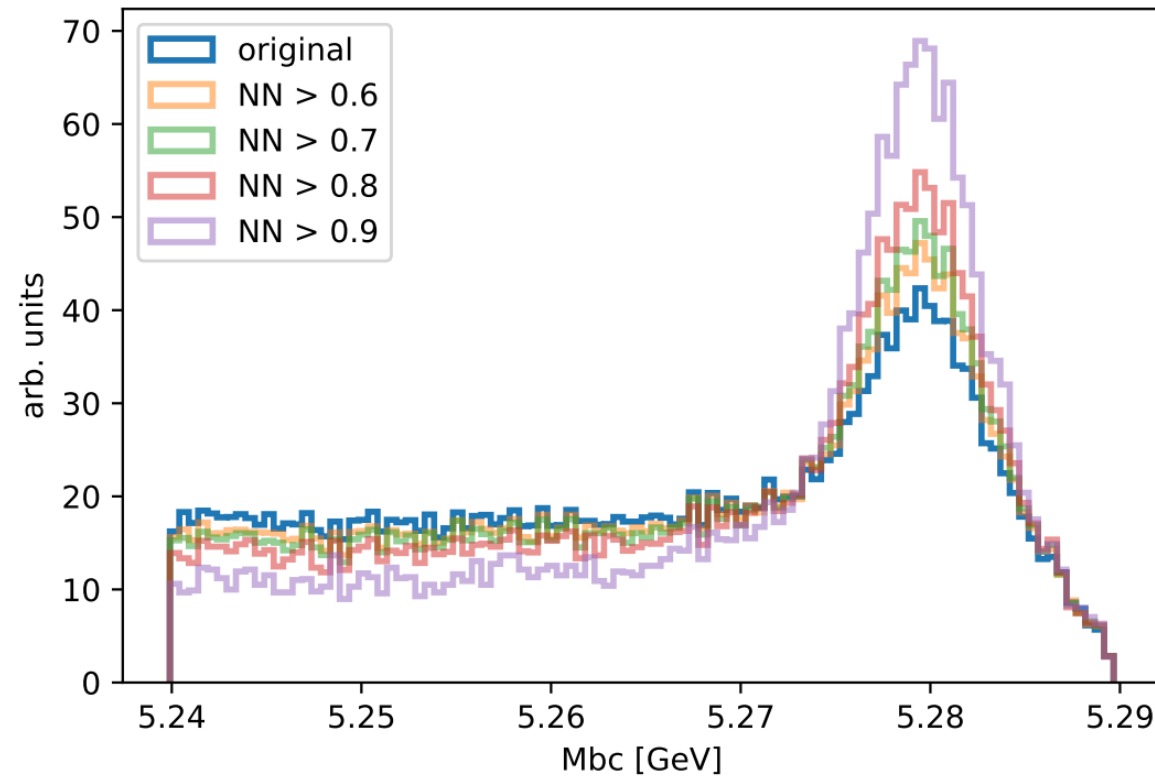
- Best AUC* improved from 0.9083 to 0.9122
- * Area under the Curve of ROC
(The closer to 1 the better)



- Visualization of Node/Edge-Attentions
Darker \Rightarrow More attention



Bias due to False-Negatives with Naive Filtering



Skim \ NN	Positive	Negative
	True-Positive (TP)	False-Negative (FN)
Pass	True-Positive (TP)	False-Negative (FN)
Fail	False-Positive (FP)	True-Negative (TN)

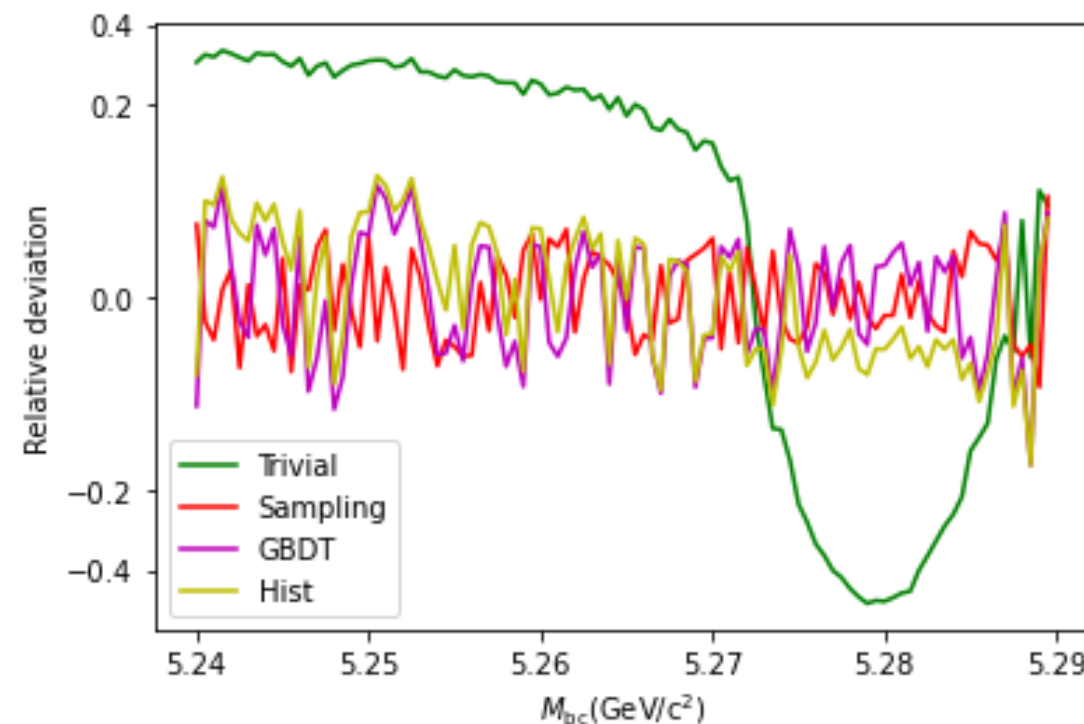
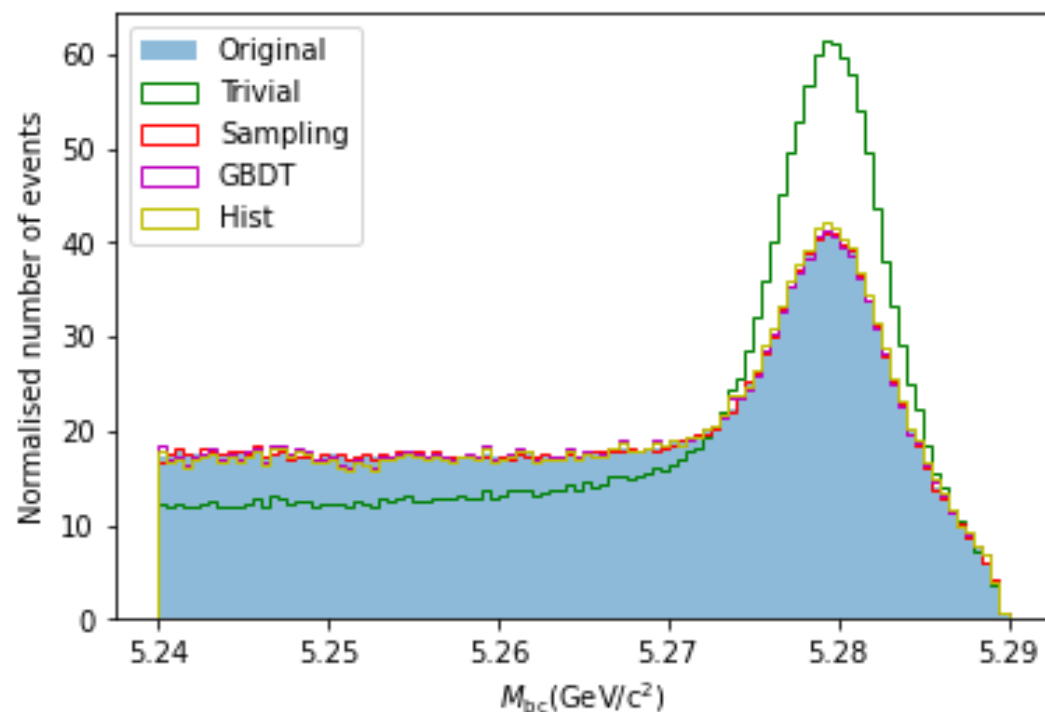
	Sampling Method	Reweighting Method
Use of NN output	As probability to keep event randomly	As score for selection according to fix threshold
Weight	Inverse of NN output	Decided with the help of another classifier
Loss to train NN	Speedup	Binary cross entropy

Metric: Speedup

--Improvement of computation time to produce the same effective number of events with the help of NN filter:

$$\text{Speedup: } s = \frac{t_{no_filter}}{t_{filter}}$$

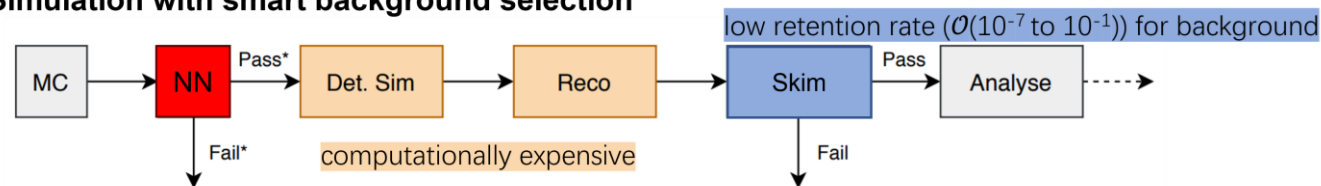
$$\text{Effective Sample Size: } N_{eff} = \frac{(\sum \omega_i)^2}{\sum \omega_i^2}$$



	Sampling	Reweighting
Maximum speedup	2.0	6.5
Bias	No bias	Small bias on some of the variables

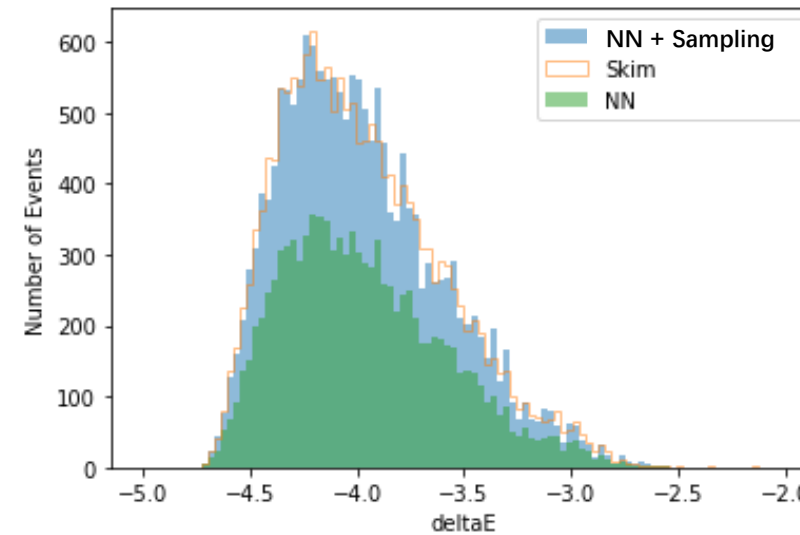
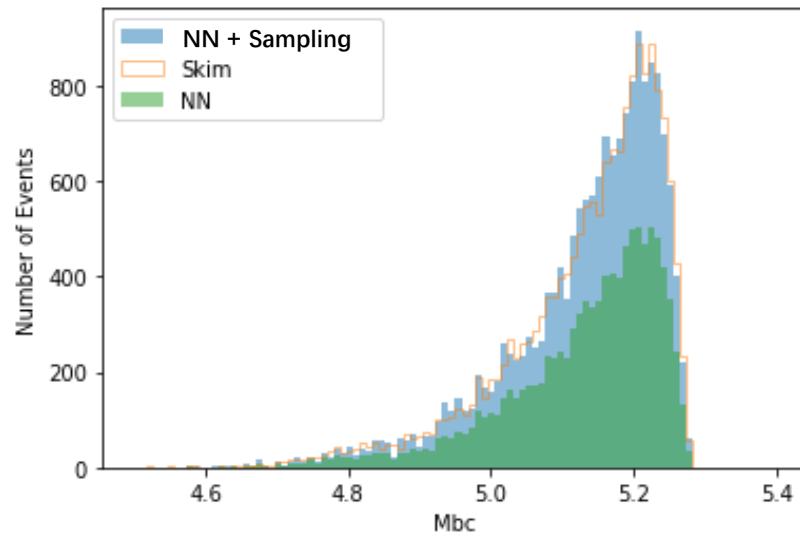
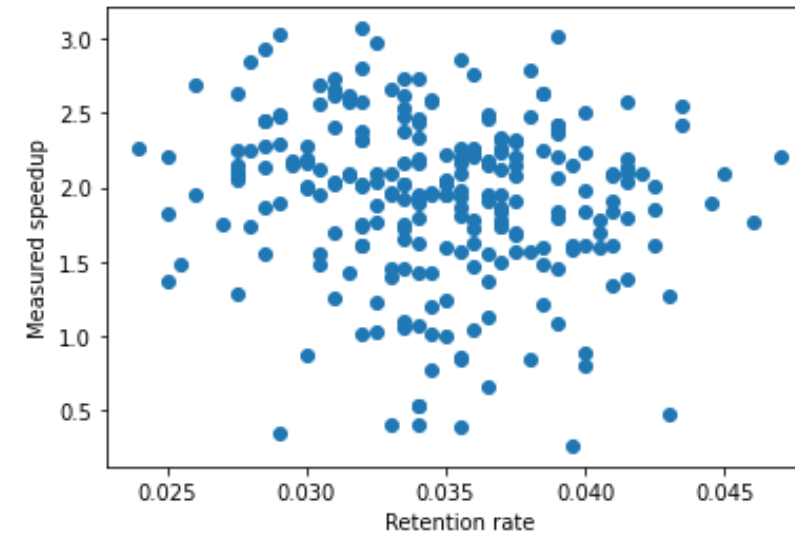
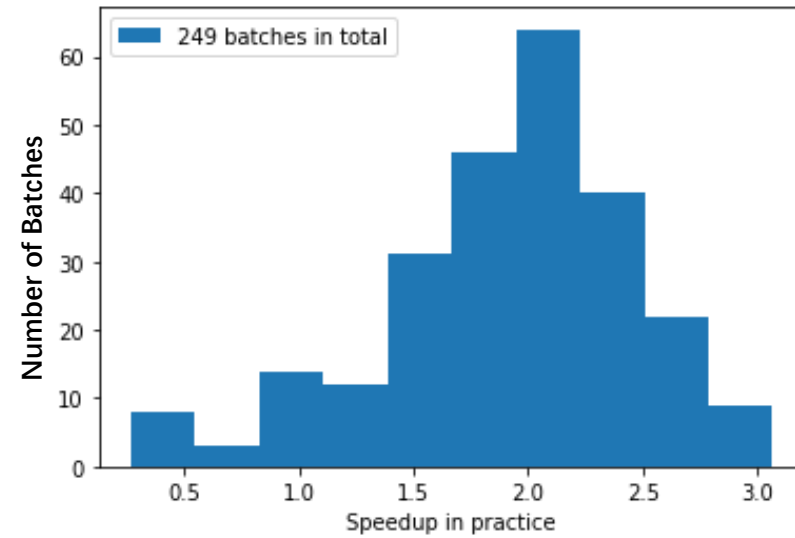
Test the module using $B^+ \rightarrow K^+ \nu \nu$ inclusive reconstruction:

Simulation with smart background selection



$B^+ \rightarrow K^+ \nu \nu$ inclusive	skim.WGs.ewp	SmartBKG
Datasets	Run full chain with charged generic MC	Train: Charged generic MC14 Test: Run full chain with charged generic MC and SmartBKG
Process (Time measurement)	<ul style="list-style-type: none"> • DetSim & Rec • Skim • ROE • Y(4S) Reconstruction 	<ul style="list-style-type: none"> • NN Prediction & Sampling • DetSim & Rec (Test only) • Skim • ROE • Y(4S) Reconstruction
Sample sizes	0.5M	Train: 1.7M Test: 0.5M
Retention rate	3.68%	16.1% (True-Positive-Rate: 60.4%)
Speedup	-	Theoretical during training: 2.09 Measured in practice: 1.92

Test the module using $B^+ \rightarrow K^+ \nu \nu$ inclusive reconstruction





Conclusion:

- Attention mechanism can improve NN performance for selective background monte carlo simulation
- Bias is avoided with sampling method while a speedup of factor 2 can still be maintained
- Reweighting method can reach much higher speedup up to 6.5 but will still have some bias in the variables that are not used in the training of the extra classifier

Plans:

- Further improvements of the NN and its training
- Further improvements of weighting methods

Thank You for your Attention

Boyang Yu¹, Nikolai Hartmann¹, Thomas Kuhr¹

¹ *Ludwig-Maximilians-Universität München*

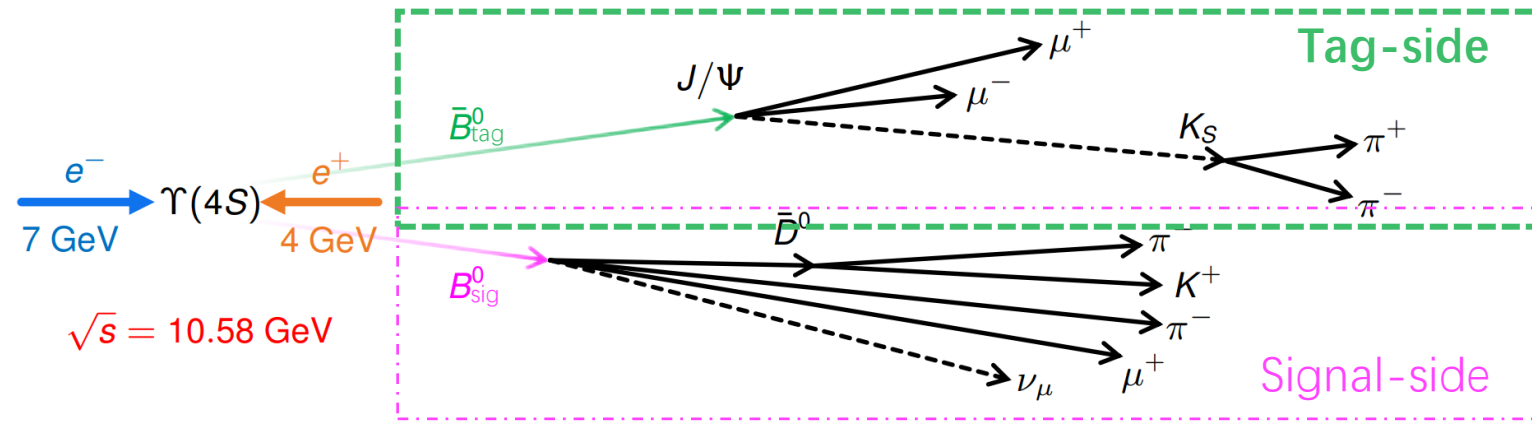
KISS B2, Sept 28, 2023



Backup



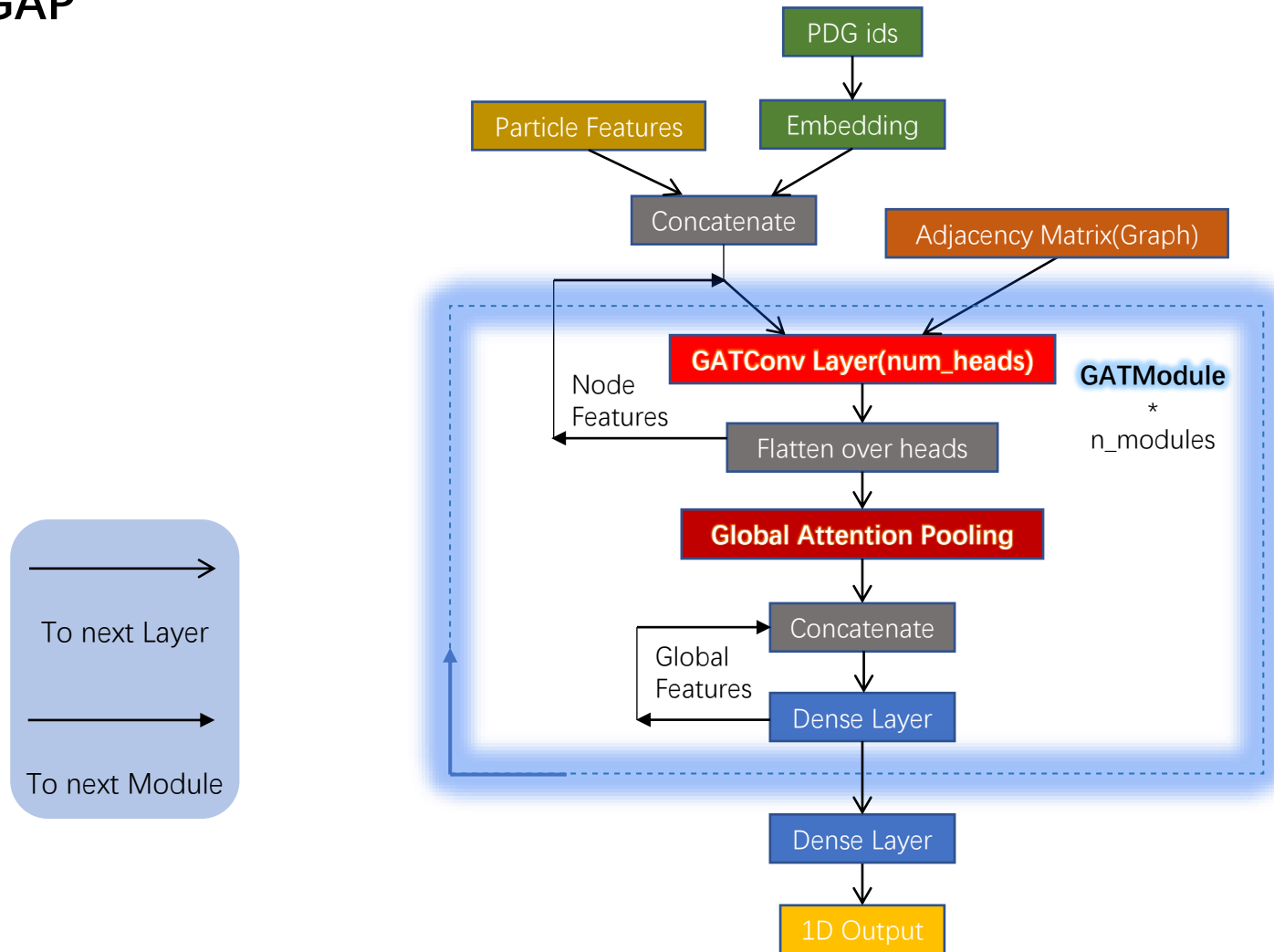
Tagging method:

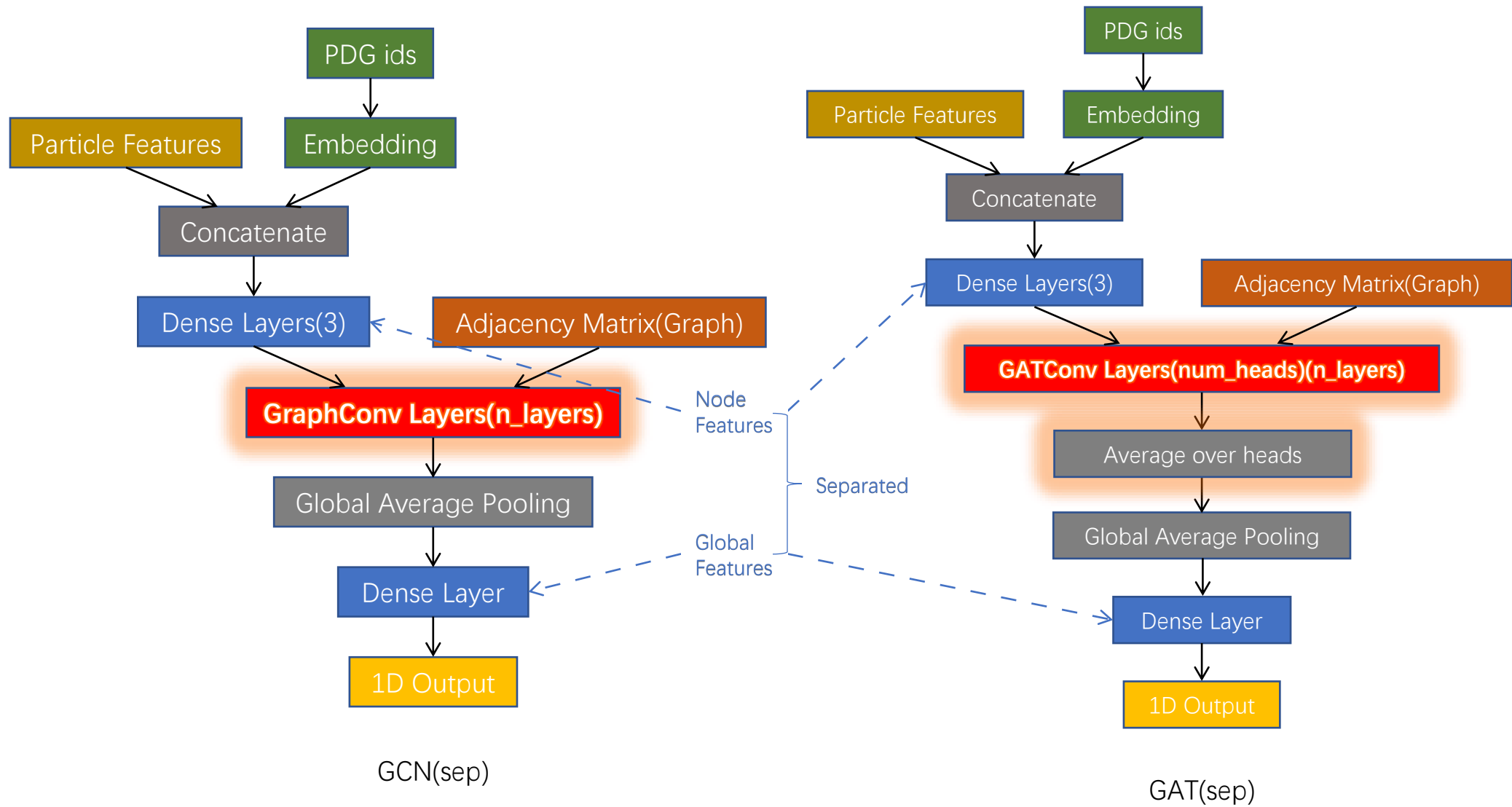


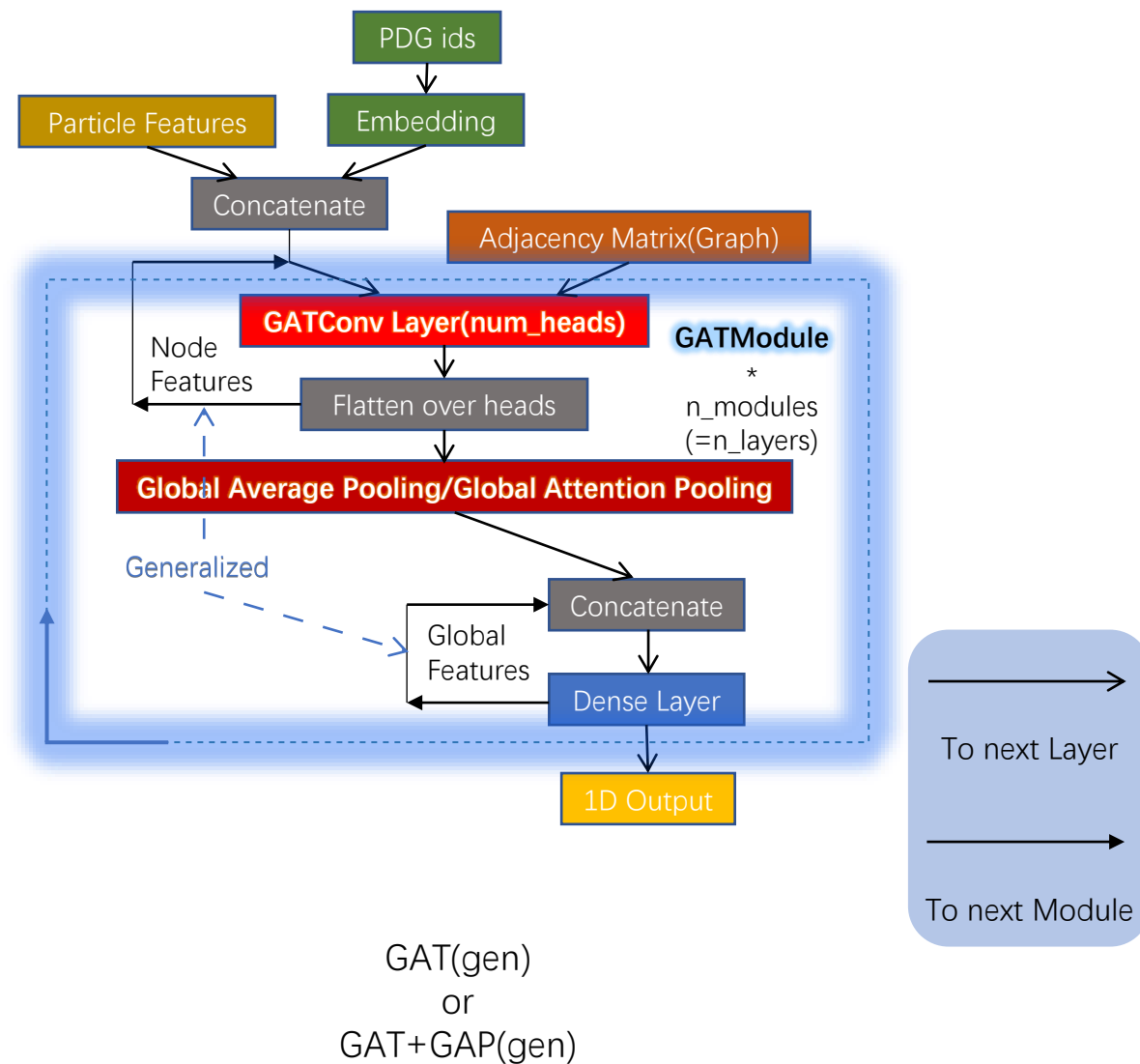
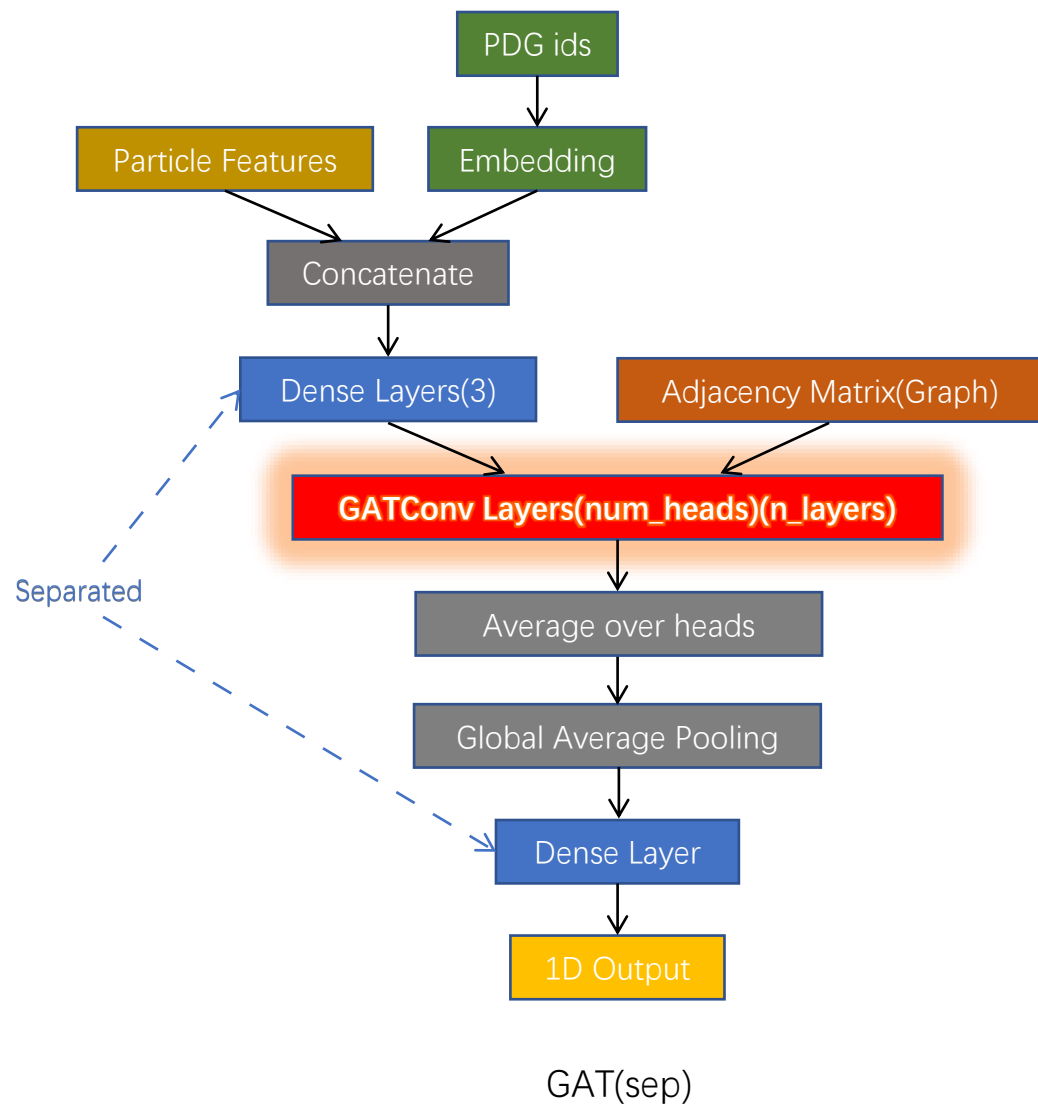
Retention rate after reconstruction and selection of tag-side B candidate:

FEI Skim	Hadronic B^+	Hadronic B^0
Mixed ($\Upsilon(4s) \rightarrow B^0 \bar{B}^0$)	5.62%	4.25%

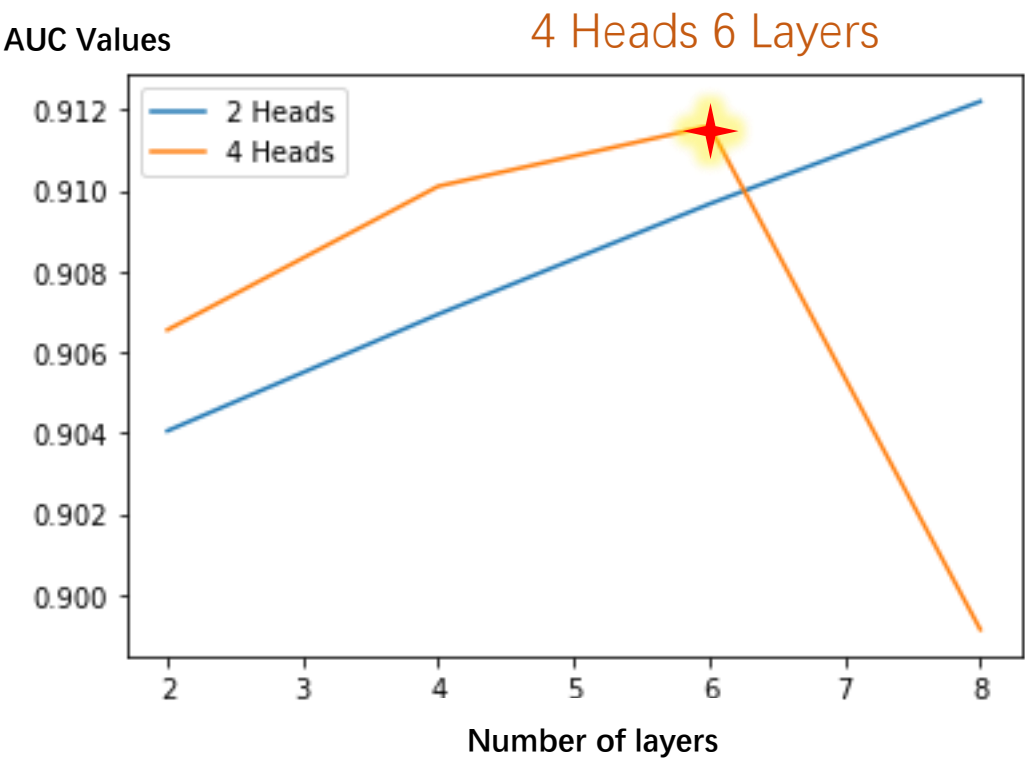
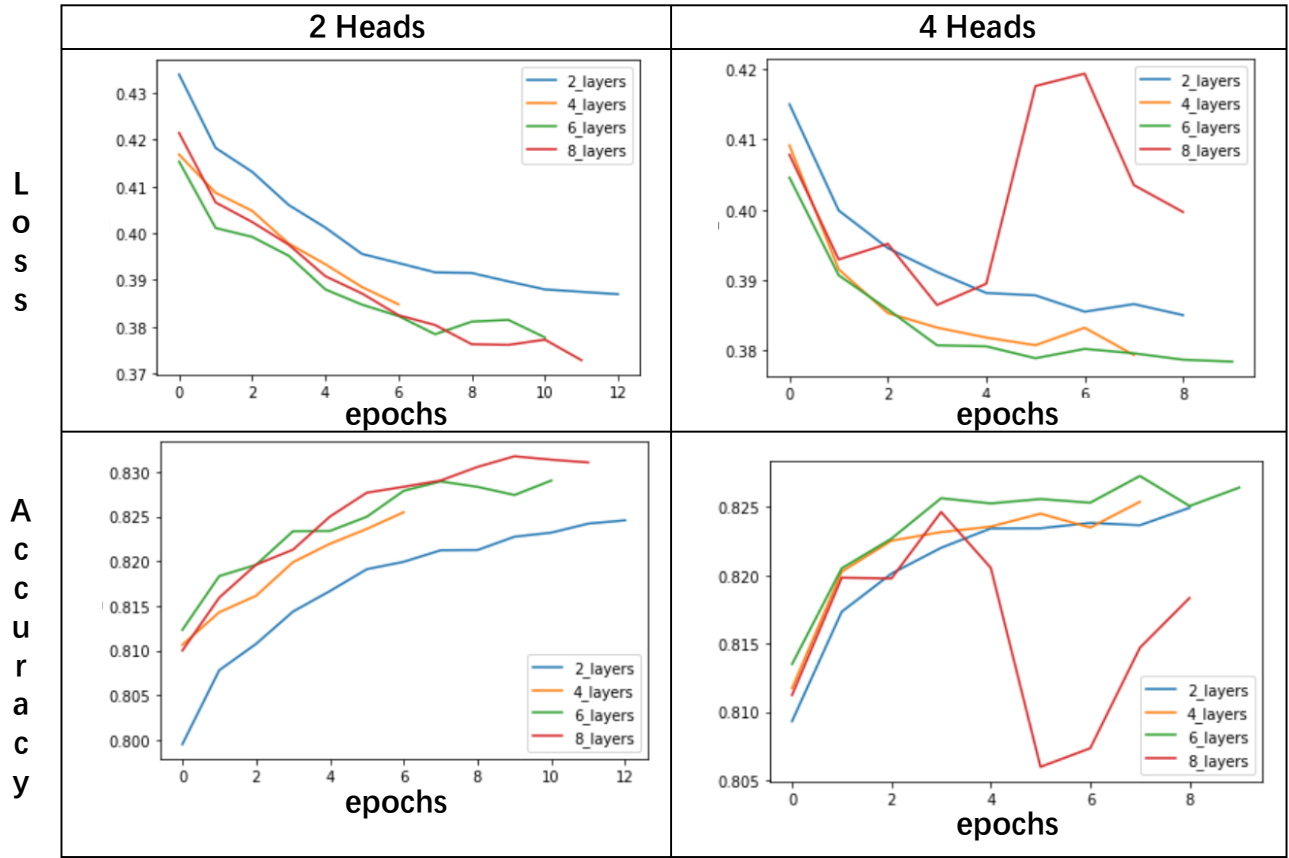
Final Architecture: GAT+GAP







Quantitative Studies



Comparison

Parameters:

- $n_heads = 4$
- $n_layers = 6$
- $n_units = 128$
- $batch_size = 128$
- $n_train = 0.9M$
- $n_val = 0.1M$
- $n_test = 0.5M$

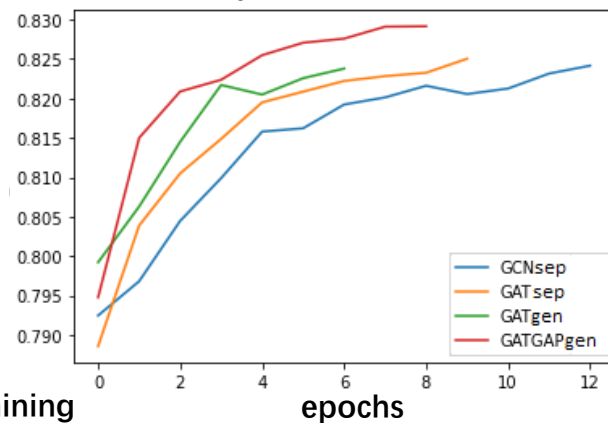
Loss:

- Entropy

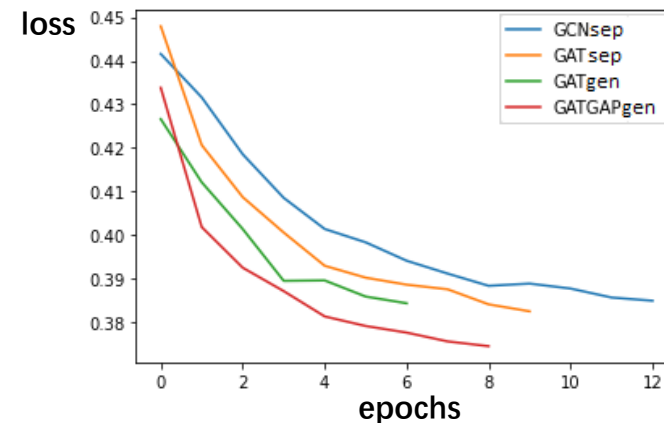
EarlyStopping:

- $patience = 3$
- $\delta = 1e-5$

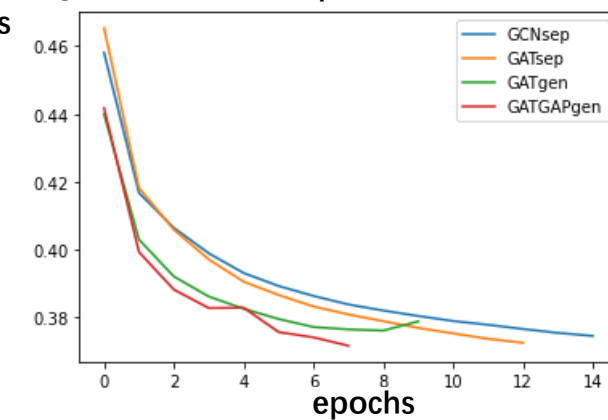
Validation accuracy



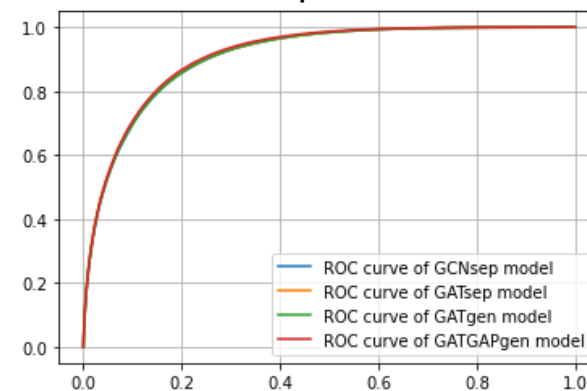
Validation loss



Training loss



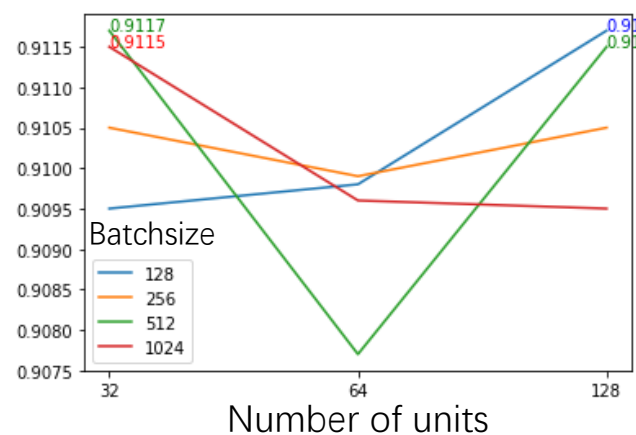
ROC curve



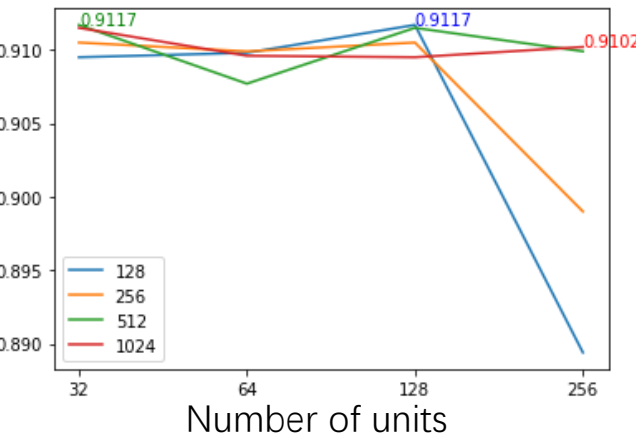
	GCN(sep)	GAT(sep)	GAT(gen)	GAT+GAP(gen)
TrainingTime	3619.46s	4047.47s	3471.48s	5049.81s
AUCValues	0.90831	0.90937	0.90891	0.91216

Grid Search

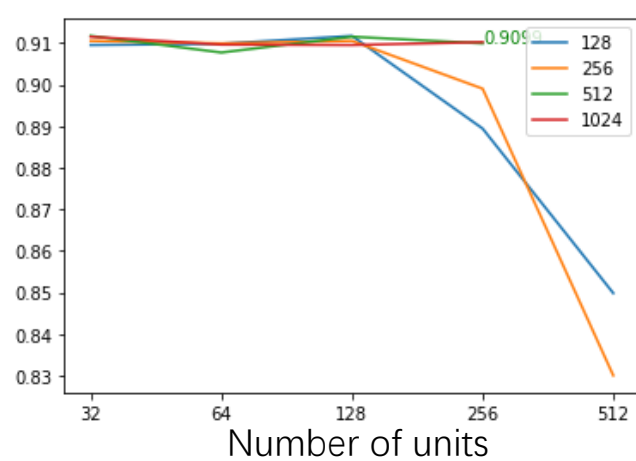
AUC



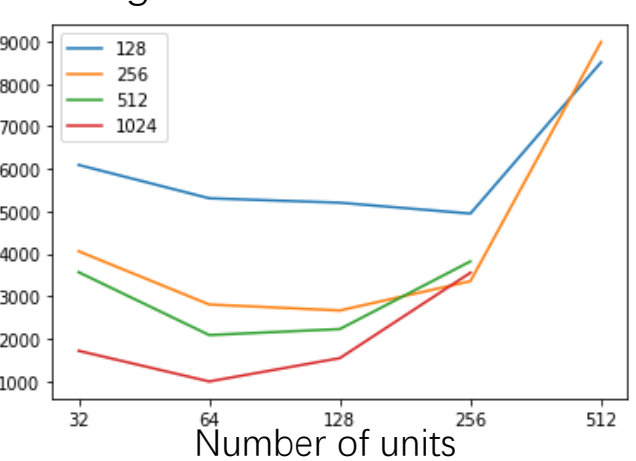
AUC



AUC



Training Time



Best Combinations

Batch-size	Number of units	AUC	Training Time
128	128	0.9117	5205
256	32	0.9105	4061
256	128	0.9105	2666
512	32	0.9117	3568
512	128	0.9115	2228
1024	32	0.9115	1716
1024	256	0.9102	3556

Network Sizes

# Units	# Parameters
32	120,527
64	459,951
128	1,808,495
256	7,184,367
512	28,651,247

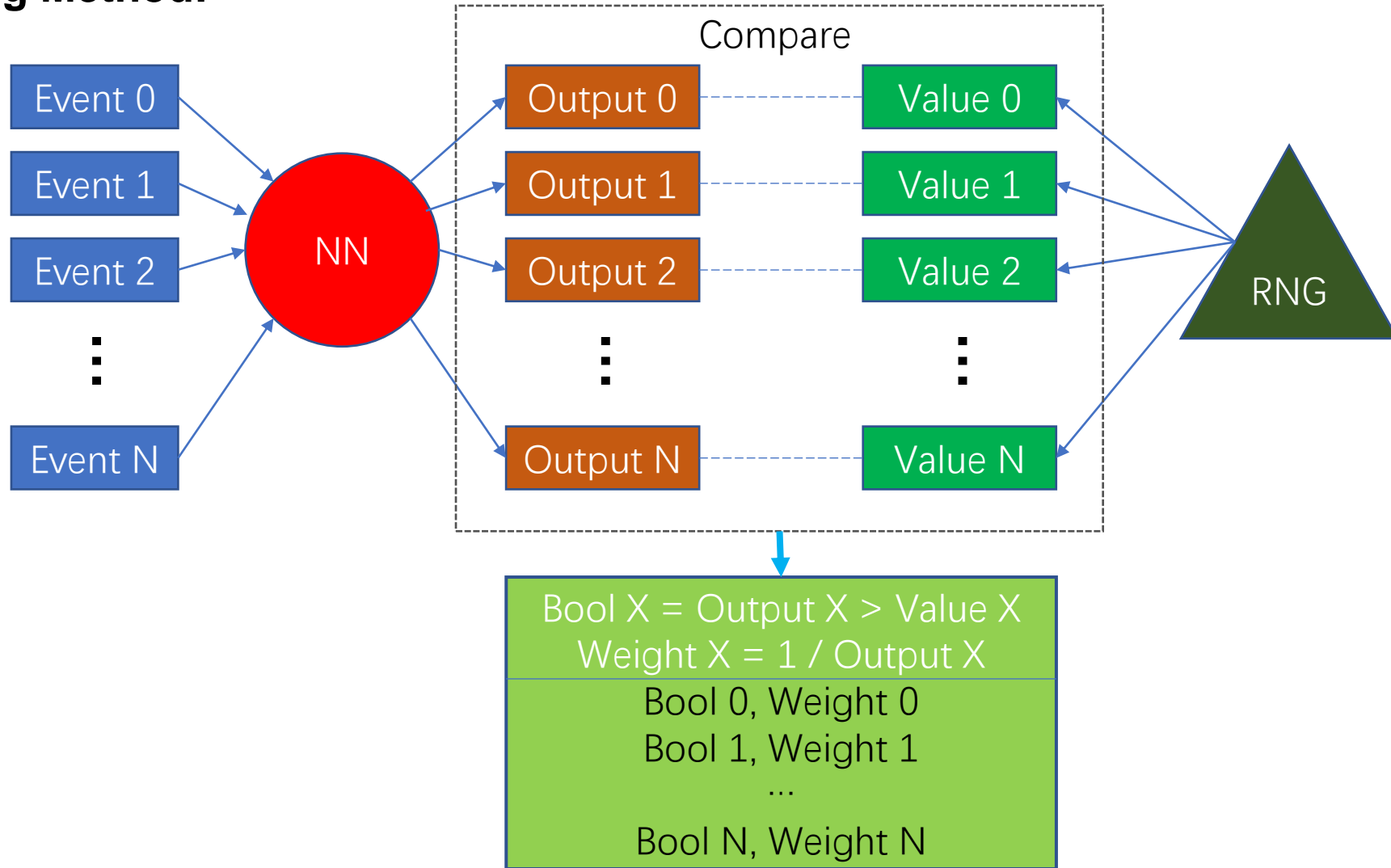
Hyperparameter Optimization

Model	AUC	Batch Size	Number of Units	AUC	Training Time in s	Number of Units	Number of Parameters
GCN(sep)	0.908	128	16	0.9131	10940		
GAT(sep)	0.909	512	32	0.9117	3568		
GAT(gen)	0.909	128	128	0.9117	5205	16	34,911
GATGAP(gen)	0.912	1024	32	0.9115	1716	32	120,527
		512	128	0.9115	2228	64	459,951
		256	128	0.9115	2666	128	1,808,495
		256	32	0.9115	4061		

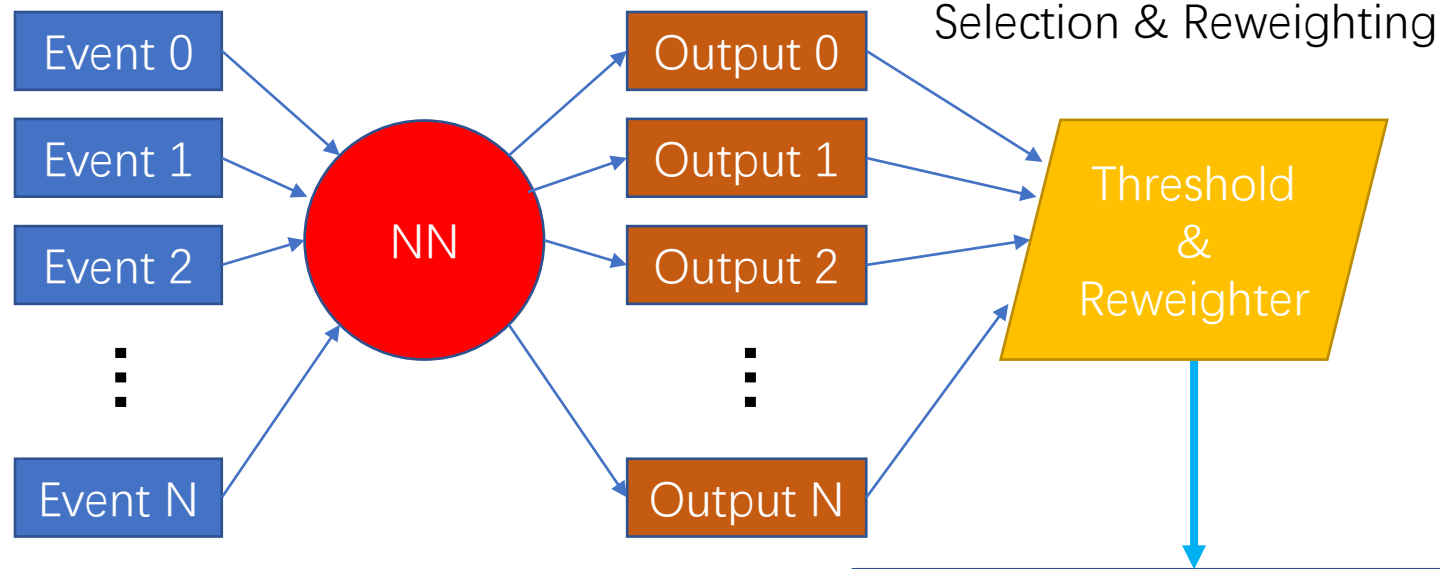
Final Configuration:

- GATGAP Model using PyTorch + Deep Graph Library (DGL)
- 6 layers with 4 attention heads each and 32 units for GAT output & global features
 -> \approx 120k parameters
- Batch size 1024 (GPU training)

Sampling Method:



Reweighting Method:



Studied reweighters:

- GBDT Reweighting
- Histogram Reweighting

Bool X = Output X > Threshold
If Bool X:
Weight X = f(Output X)

Bool 0, Weight 0
Bool 1, Weight 1
...
Bool N, Weight N

Reweighting Method:

- Train a Gradient Boosting Decision Tree (GBDT) classifier with some event level variables to distinguish between True-Positive events and False-Negative events
- GBDT Reweighting: use the outputs of the classifier directly:

$$w = \frac{1}{p_{clf}} = \frac{1}{p_{TP}/p_{TP+FN}} = \frac{p_{pass_skim}}{p_{TP}}$$

- Histogram Reweighting: compare the score histogram of all the events that can pass the skim (True-Positive + False-Negative) with the score histogram of True-Positives to give each bin of score a scaling factor:

$$w = w_{bin_i|p_{clf} \in bin_i} = \frac{H_{pass_skim,i}}{H_{TP,i}} \Big|_{p_{clf} \in bin_i}$$

Skim \ NN	Positive	Negative
	Pass	Fail
Pass	True-Positive (TP)	False-Negative (FN)
Fail	False-Positive (FP)	True-Negative (TN)

Relative statistical uncertainty and effective sample size

Variable	Formula	Remark
NN outputs / Probabilities to pass	$\{p_i\}$	'i' refers to each event in the whole sample (batch)
Weights	$\{\omega_i\} = \left\{ \frac{1}{p_i} \right\}$	Infinites (at $p_i = 0$) are excluded and set to 0 Avoid the bias by construction
Relative statistical uncertainty	$S = \frac{\sqrt{\sum \omega_i^2 p_i}}{\sum \omega_i p_i}$	$\sum \omega_i^2 p_i = \sum \omega_i$ $\sum \omega_i p_i = N$ Here consider only passed events (label = 1)
Effective sample size	$N_{eff} = \frac{1}{S^2}$	Number of events needed to reach the same statistical uncertainty without sampling

Speedup rate

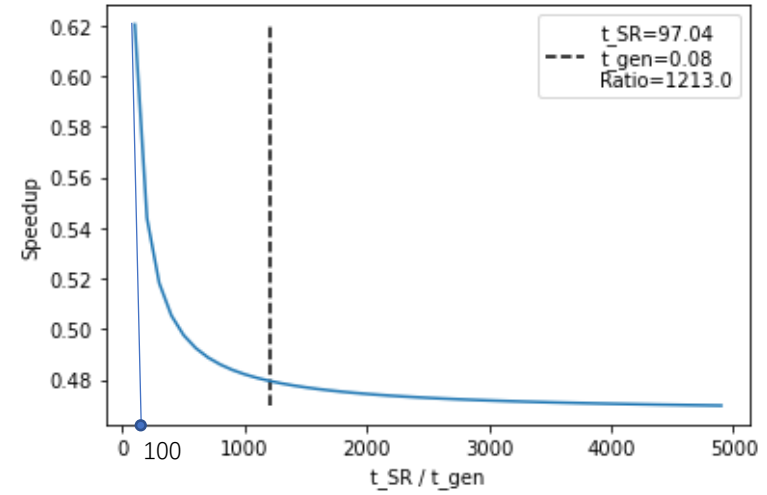
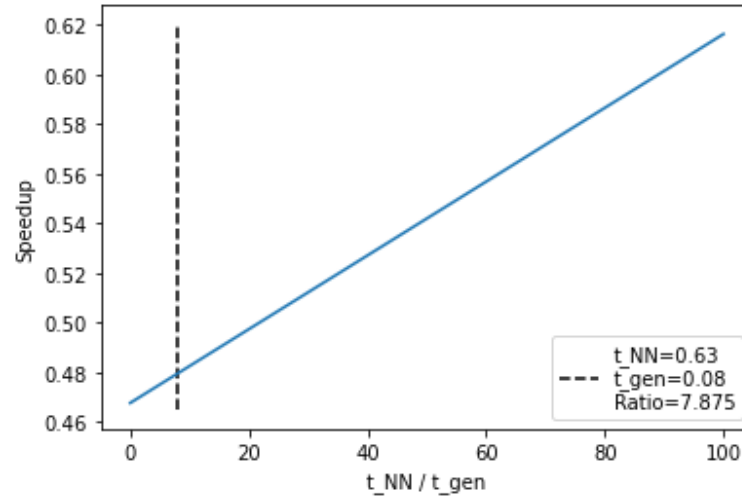
Variable	Formula	Remark
Skim retention rate	$r = 0.05$	Probability to pass the skim process
Times of different phases in ms	$t_{gen} = 0.08$ $t_{NN} = 0.63$ $t_{SR} = 97.04$	Taken from previous studies
Effective number of events after sampling	$n_+ = \sum p_i$ $n_- = \sum (1 - p_i)$	$\{p_i\}$ will be divided into two subsets where the events will/won't pass the skim process
Time consuming with NN filter	$t_+ = [n_{TP}r + n_{FP}(1 - r)](t_{gen} + t_{NN} + t_{SR})$ $t_- = [n_{FN}r + n_{TN}(1 - r)](t_{gen} + t_{NN})$	Positive/Negative: Result of sampling True/False: Result of sampling == skim process
Time consuming without NN	$t_0 = N_{eff}(t_{gen} + t_{NN})$	To reach the same statistical uncertainty
(Inverse) Speedup rate	$R = \frac{t_+ + t_-}{t_0}$	The lower the better

Robustness:

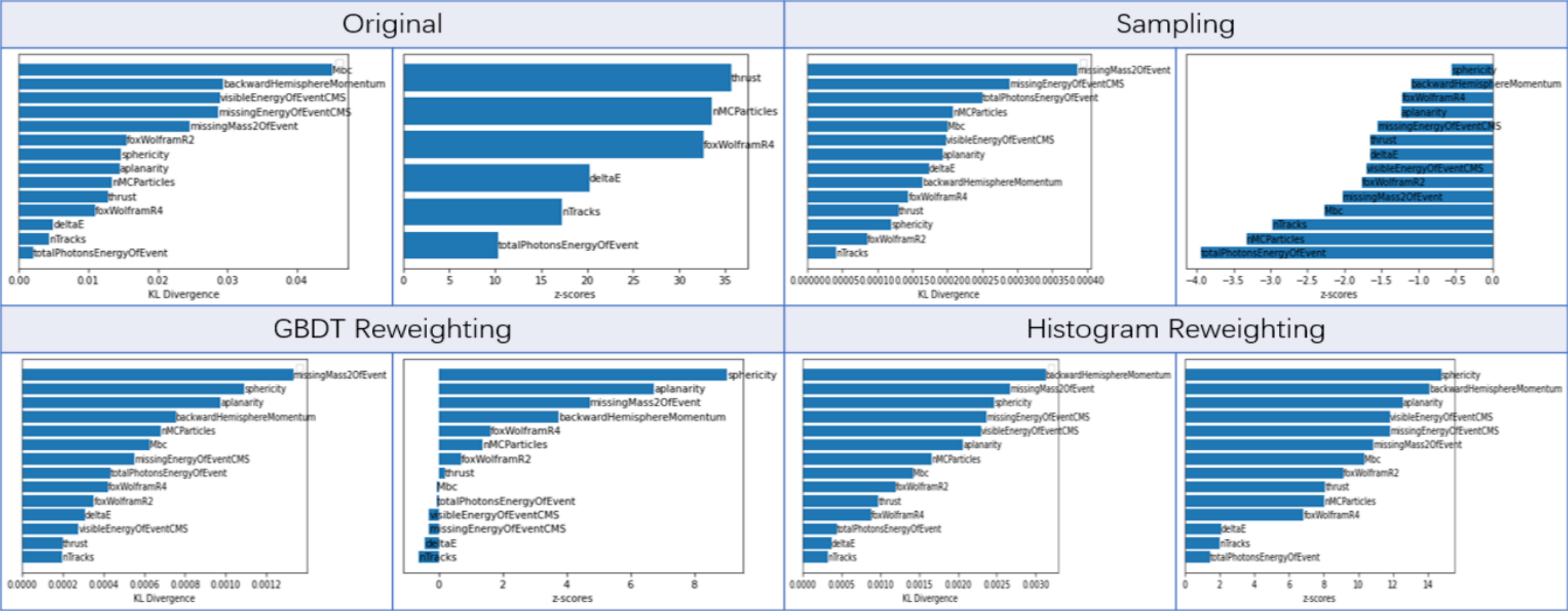
Weak dependency of
Speedup on t_{NN} and t_{SR}



Safe to generalize



KS-Test



Original

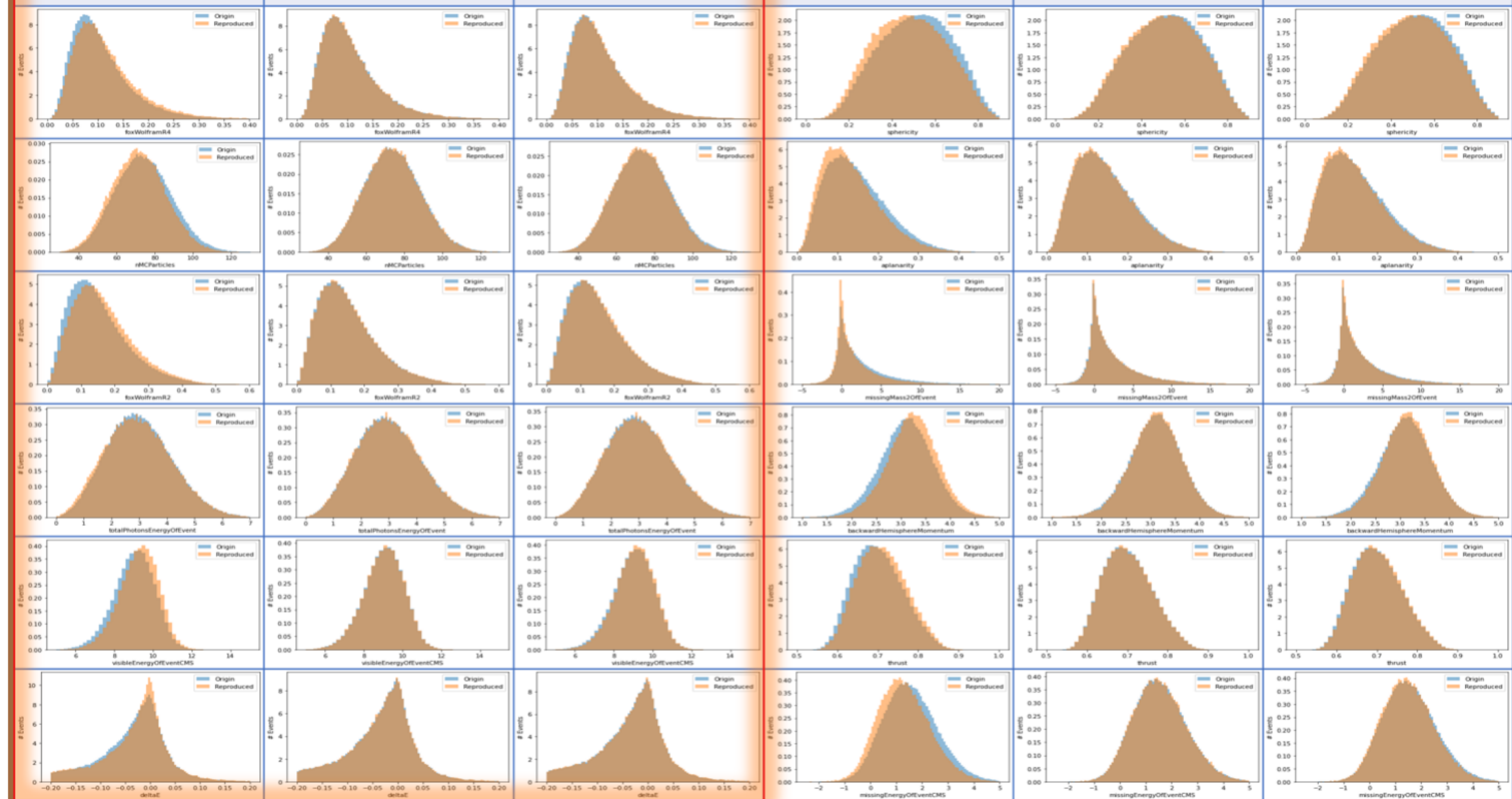
GBDT Reweighting

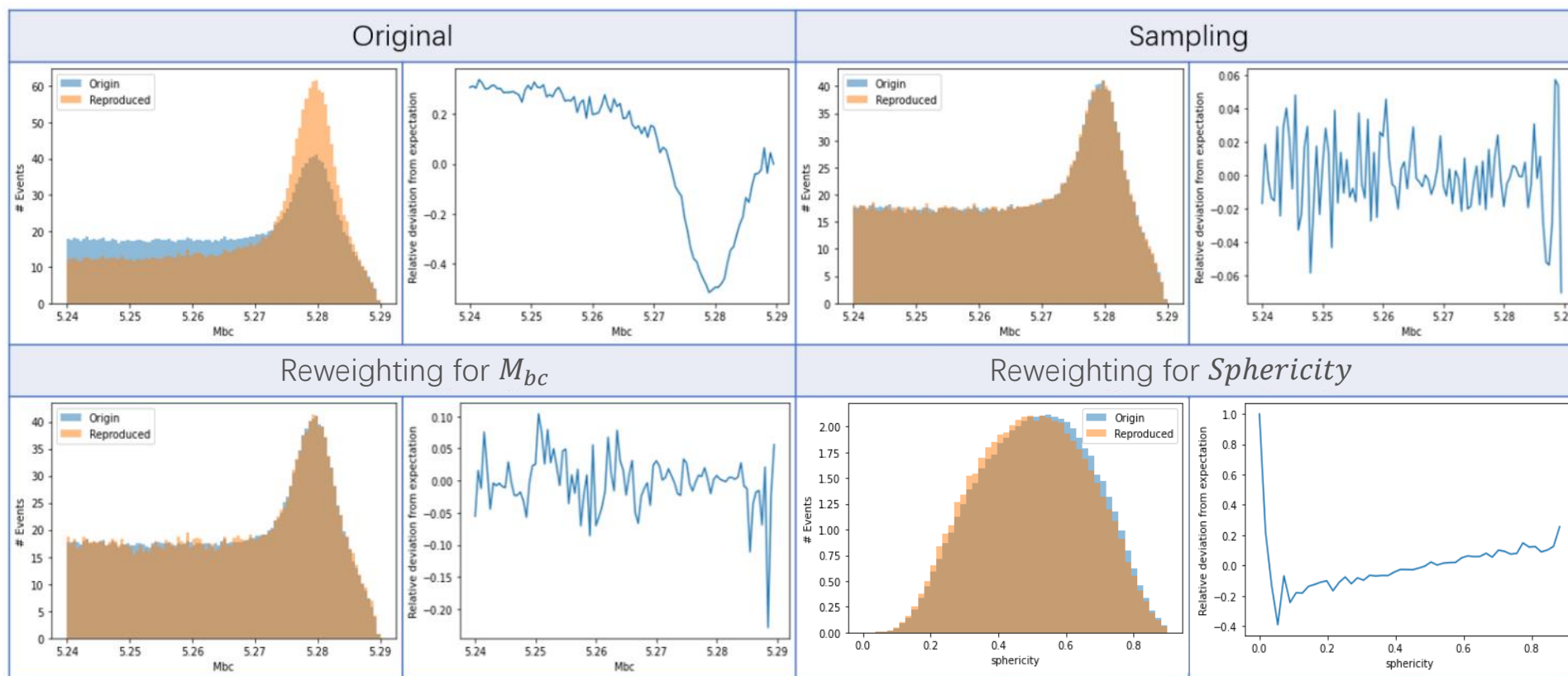
Hist Reweighting

Original

GBDT Reweighting

Hist Reweighting





skim.WGs.ewp.inclusiveBplusToKplusNuNu

- Track cleanup:
 - $p_t > 0.1$
 - $\theta_{\text{InCDCAcceptance}}$
 - $dr < 0.5$ and $\text{abs}(dz) < 3.0$
- Event cleanup:
 - $3 < \text{nCleanedTracks} < 11$
- Kaon pre-cuts:
 - $\text{track cleanup} + \text{event cleanup} + \text{nPXDHits} > 0$
- **K+ reconstruction**
- Kaon cuts:
 - $p_t \text{ rank}=1$
 - $\text{kaonID} > 0.01$
- **B+ reconstruction**
- B+ cut:
 - $\text{mva_identifier: MVAFastBDT_InclusiveBplusToKplusNuNu_Skim} > 0.5$