Study for jet flavor tagging with <del>recent Neural Network ightarrow Transformer</del>



Masahiro Morinaga

The University of Tokyo ♥ (ICEPP ♥,Beyond Al ∧;)

Boost23 @Laurence Berkeley National Lab

券 🏵 🎮 Masahiro Morinaga

# What is Flavor Tagging?

### **Jet Flavor**

- QCD jets originated from quark or gluon decay products.
- Tagging jet flavor is an important technology to improve analysis sensitivity
- Quark and Gluon classification
- W/Z boson, Higgs, Top tagging :
  - $\circ$  Using a large R-jet can identify the jet if a particle is boosted
  - $\circ~$  Only higher  $p_{
    m T}$

## **Heavy Flavor Tagging**

- Rely on secondary vertex from B-meson or D-meson decay
  - $\circ \ b$ -quark : b o B o D o K
  - $\circ \ c\operatorname{-quark} : c \to D \to K$
- Machine Learning : BDT, CNN, RNN, DeepSets, Graph NN

## Quark/Gluon Tagging

- The separation between light flavor and gluon
- Rely on #of tracks or similar variables
- ML : Graph or similar specialized model(<u>ParticleNet</u>, <u>Lorentz Group NN</u>)



Strategy

#### Problem

- Models that are too specialized for HEP are hard to use in implementation and application.
- Standard models (ResNet, Transformer, etc.) are easy to implement and apply.
- Flavor tagging, such as b-tagging, uses only track information, but essentially clusters are also considered to have useful information
- Flavor tagging and H/Z/W/top-tagging are often constructed separately.

## Aim of this study

- Use models widely used in the ML community  $\rightarrow$  Transformer
- Create an all-in-one tagger that includes all flavors.
  - Classify light-flavor, gluon-jet, *c*-jet, *b*-jet simultaniously
  - $\circ~$  Same model can be applied for the large-R jet (qar q)
- Adding cluster information
  - $\circ~$  Expect better performance at higher  $p_{
    m T}$  by using calorimeter information

## Model Architechture : Transformer



### Transformer

- Simple Transformer encoder
- Taking input sequences into a tokenizer
- Repeat two residual blocks *L* times:
  - Multihead self-attention:
  - Feed Forward Network:
- Pool by Global Average Pooling
  - There is a Linear layer after the pooling layer as a head
- No jets physics acknowledgments
- No Lorentz summation

# Model Architechture : Tokenizer



### Tokenizer

- Tokenizer make tokens from track inputs
  - Feature : 4-vector or other features for each tracks
  - $\circ~$  Position : Relative position difference to parent jet,  $\Delta \phi, \Delta \eta~$
- Sum up X and position, not concat!

Training Setup

#### **Training Sample**

- <u>top-data</u>: Open access data with Top and QCD dijet samples
  - 14 TeV, Delphes ATLAS detector card with Pythia8
  - Clustering of particle-flow with anti-kT 0.8 jets in the pT range [550,650] GeV
  - 1.2M for training, 0.4M for test/validation samples

### **Model Setup**

- Total #of Parameters : 101,794
  - Depth: 4
  - Latent dimension : 32 , factor x4 for internal
  - Attention head : 4
  - Dropout: 0.1
  - Stochastic path : 0.1

## **Training setup**

- Trained by <u>JAX/Haiku</u>
- Optimizer: Adam with wieght decay 1.0e-4
- Learning rate: CosineAnnealingWarmupRestart
  - 1.0e-4 to 1.0e-7 ,
  - Wramup: 10% of total update steps
- Using NVidia A100 GPU x 8
  - data/model/operator parallel by <u>alpa</u>
- Using mixed precision bfloat16
- Batch size : default is 4096 = 4k .
  - Gradient accumulation : #of micro batch 8
- Target loss: <u>Class Balanced Loss</u>
  - The number of each signal is not equal in size
  - To correct the imbalance in training

## Loss Curve





- training loss curve and top-1 accuracy
  - 60000 steps, roughly ~200 epochs
  - Total training time ~2.5h

## Performance : Transformer





- ROC curve
- Confusion matrix:
  - $\circ$  91% of the true top can be labeled as the predicted top.

# Comparison : ParticleNet





- <u>ParticleNet</u> is a good candidate to compare with this work.
  - Graph Neural Network like style
  - $^{\circ}~$  Edge convolution with k-NN
  - Total # of paramters : 365,162
    - Three times larger than Transformer in this study.

## Performance : ParticleNet



- Slightly worse than Transformer
- Note that : difference from <u>the paper</u>
  - We used a bigger batch size and a different learning rate scheduler. The original learning scheduler cannot handle training
  - The number of total steps is much more significant due to different batch size compared with Transformer training

# Performance : Comparison



- Slightly worse(roughly 1%) than Transformer, *both results are worse than the paper!!*
- Even if my implementation/training of the ParticleNet is wrong, the Transformer performs similarly.
- At least the Transformer has good ability for the top tagger?

Sample with Clusters

### Original sample

- To test model performance with more realistic situations, add cluster information
- Making new samples by <u>Delphes</u> with more labels:
  - $\circ$  Lables: total 8 labels,  $l,s,c,b,g, lar{l},sar{s},car{c},bar{b}$  where l is u or d
  - $\circ~~Z o qar{q}, H o sar{s}/car{c}/bar{b}, g o uar{u}/dar{d}/sar{s}/car{c}/bar{b}$ ,other
  - Label composition:

Label	uu/dd	SS	сс	bb	u/d	S	С	b	gluon
	2.9%	9.0%	8.8%	8.0%	24.3%	8.24%	9.3%	15.4%	14.1%

- Using Delphes with ATLAS geometry
  - MG5\_aMC of v3.2.0 w/ <u>Pythia8</u>
  - Delphes with ATLAS no pileup card using a partcle flow based jets
- Generate pp
  - $\,\circ\,$  Selection :  $p_{
    m T}>250$  GeV and  $|\eta|$  < 2
  - At least one jet constituent(track or cluster)

# Model Architechture : Tokenizer with Clusters



#### Tokenizer

- Tokenize tracks and clusters as tokens
- Absorb the feature differences between tracks and clusters
- one hot vector are additionally taken as input
- Affine transformation is applied

### Feature

- input feature
  - $\circ~~$  Tracks:  $p_{\mathrm{T}}, \Delta\eta, \Delta\phi, d_{0}$
  - $\circ~$  Cluster:  $p_{
    m T}, \Delta\eta, \Delta\phi, E_{
    m EM}/E_{
    m total}, E_{
    m Had}/E_{
    m total}$
- Position :  $\Delta\eta, \Delta\phi$  for both tracks and clusters
- Tracks and clusters are input into the same linear layer

Orignal Samples



- Transformer with 4M samples for 200 epochs.
  - Difference between train/test comes from difference of sample composition between train/test samples.
- ParticleNet cannot be converged, or poor separation power

Orignal Samples



		Efficiency								
	ss pair	-52.7%	4.1%	1.0%	19.1%	10.8%	9.2%	1.9%	1.1%	
	cc pair	- 2.9% -	57.8%	8.5%	3.5%	1.4%	0.7%	18.6%	6.6%	
	bb pair	- 0.3% -	9.1%	70.3%	1.0%	0.1%	0.1%	2.6%	16.6%	
label	gluon	- 8.4% -	2.6%	1.4%	71.3%	7.8%	4.6%	2.5%	1.6%	
True	ıd-quark	- - <b>1</b> 1.6% -	1.6%	0.6%	22.0%	40.2%	20.4%	2.6%	1.1%	
	s-quark	- -13.3% -	1.7%	0.8%	21.6%	29.0%	29.0%	3.5%	1.4%	
	c-quark	2.2%	10.4%	2.0%	6.5%	3.8%	2.4%	60.4%	12.3%	
	b-quark	- 0.8% 	5.2%	12.6%	3.3%	0.5%	0.4%	13.8%	63.3%	
55 pair c pair oluon quart quart quart quart quart										
	Predicted label									

- Transformer with 4M samples for 200 epochs.
- The most difficult pair is strange and up/down separation.

Orignal Samples : More Samples



- Transformer with 256M samples for 64 epochs.
- Compared with 4M results, better loss values.
  - Training time about 1 week

Orignal Samples : Comparison



4M

• Training with 256M samples improved most of the values.

	Efficiency								
ss pair	57.4%	1.9%	0.5%	16.5%	8.1%	12.9%	1.6%	1.1%	
cc pair	- 1.8%	62.8%	7.0%	2.4%	0.8%	0.8%	19.6%	4.9%	
bb pair	- 0.3% -	6.0%	73.5%	0.9%	0.1%	0.2%	2.0%	17.0%	
gluon	- 7.2%	1.6%	0.9%	72.3%	7.3%	6.5%	2.4%	1.7%-	
ud-quark	- 9.9% -	0.9%	0.4%	20.7%	37.2%	28.0%	2.0%	1.0%-	
s-quark	- -11.0% -	0.9%	0.4%	20.6%	23.3%	40.0%	2.6%	1.1%_	
c-quark	- 1.5% -	6.7%	1.3%	4.6%	2.3%	2.8%	71.8%	9.0%	
b-quark	- -0.7%	2.6%	7.1%	2.5%	0.3%	0.5%	10.2%	75.9%-	
Predicted label									
256M									

Orignal Samples : Dataset size



- Better performance with a larger dataset.
- Training with 256M samples took about one week  $\rightarrow$  need to improve.

## Conclusion

### Conclusion

- A flavor tagger with all-in-one labels seems to be working.
- Transformer would be a good candidate as a jet flavor tagger.
- Tracks and clusters can be used simultaneously as inputs to the Transformer.
- Again, the amount of data is totally justice.

### Plans

- Check the scaling law for our data
  - More data
  - Bigger model
  - More training epochs
- The ultimate goal : Supervised training with real data

## Backup

# Original Samples with ParticleNet



- ParticleNet with original samples
- No hyperparameter optimization

Orignal Samples : Comparison



Efficiency

• Training with 256M samples improved most of the values.

	Purity								
ss pair	-63.9%	2.2%	0.6%	11.7%	10.1%	14.0%	1.5%	1.0%	
cc pair	- 2.0%	75.2%	7.7%	1.7%	1.0%	0.8%	17.5%	4.3%	
bb pair	- 0.3%	7.2%	80.6%	0.7%	0.1%	0.2%	1.8%	15.2%	
label gluon	- 8.0% -	2.0%	1.0%	51.4%	9.2%	7.1%	2.2%	1.5%_	
ບັບ ບັບ ບັບ ບັບ ບັບ	- -11.0% -	1.1%	0.4%	14.7%	46.8%	30.6%	1.8%	0.9%	
s-quark	- -12.2%	1.1%	0.4%	14.7%	29.4%	43.7%	2.3%	1.0%	
c-quark	- 1.6% -	8.0%	1.5%	3.3%	2.9%	3.0%	63.9%	8.1%	
b-quark	- -0.8%	3.1%	7.8%	1.8%	0.4%	0.6%	9.1%	68.0%	
55 pair c pair opon ou of ouar ouar couart ouart									
Predicted label									
Purerity									

# Orignal Samples : Dataset size





*c*-quark

# Orignal Samples : Dataset size



100  $10^{-1}$ False positive 10<sup>-2</sup> 10<sup>-3</sup> 10<sup>-2</sup> 4M:AUC=0.851  $10^{-5}$ 16M:AUC=0.863 64M:AUC=0.870 256M:AUC=0.878  $10^{-6}$ 0.0 0.2 0.8 0.4 1.0 0.6 True positive

u/d-quark

*c*-quark

# Orignal Samples : Dataset size



False positive  $10^{-3}$  $10^{-4}$ 10-5  $10^{-6}$ 0.0 0.2 0.8 1.0 0.4 0.6 True positive

10<sup>0</sup>

 $10^{-1}$ 

 $10^{-2}$ 

4M:AUC=0.853

16M:AUC=0.862 64M:AUC=0.874

256M:AUC=0.888

 $s\overline{s}$  pair