## Improving ATLAS Hadronic Object Performance with ML/AI Algorithms



#### Ben Hodkinson

University of Cambridge

On behalf of the ATLAS Collaboration

March 29, 2023





#### Introduction

- Jets and p<sub>r</sub><sup>miss</sup> are complex objects
  - Promising setting for ML applications
- ML applications:
  - > Regress truth-level quantities from detector-level information (eg. true jet  $p_{T}$ )
  - Classify type of object (eg. top jet vs. gluon jet)
- Highlights covered here:
  - > Missing transverse momentum (  $p_{T}^{miss}$  )
    - → Regressing truth-level p<sub>T</sub><sup>miss</sup> (<u>ATL-PHYS-PUB-2021-025</u>)
  - > Pion reconstruction
    - → Classification & calibration (<u>ATL-PHYS-PUB-2022-040</u>)
  - Boosted jet taggers:
    - → Top jets
      - (ATL-PHYS-PUB-2021-028, ATL-PHYS-PUB-2022-039)
    - → W/Z jets

(ATL-PHYS-PUB-2021-029, JETM-2022-006)

- > Jet calibration
  - → (See <u>Naseem Bouchhar's talk</u>)

# 1) Missing transverse momentum (p<sub>T</sub><sup>miss</sup>) performance



Optimal working point for a given event depends on the topology and pile-up conditions

- Tighter working points reject more pile-up but risk rejecting hard-scatter jets
- > METNet = Neural network to pick and combine working points into a single  $p_{\tau}^{miss}$  estimate
  - → Regresses truth-level p<sub>T</sub><sup>miss</sup>
  - Trained on ttbar and di-boson MC events.

#### **METNet performance**

• Improved resolution compared to individual working points.



#### **METNet performance**

- Improved resolution compared to individual working points.
- Negative bias

ATL-PHYS-PUB-2021-025

 Reduced using sample weights and sinkhorn contribution to loss function



#### **METNet performance**

- Improved resolution compared to individual working points.
- Negative bias

ATL-PHYS-PUB-2021-025

- Reduced using sample weights and sinkhorn contribution to loss function
- METNet can generalize to a number of event topologies
- Autotunes to event conditions





- $p_T^{miss}$  significance: separates real and fake  $p_T^{miss}$
- Use Gaussian Negative Log Likelihood loss to

Neutrinos / bSM

produce 'confidence'  $\sigma$  as well as central prediction:

$$\mathcal{L}_{\text{GNLL}} = \log \sigma + 0.5 \left(\frac{y - \hat{y}}{\sigma}\right)^2$$

$$\text{METNetSig} = p_{\text{T}}^{\text{miss, NN}} / \sigma$$

Mismeasured objects + finite

acceptance/resolution

- Similar behaviour and performance to ATLAS object-based p<sub>T</sub><sup>miss</sup> significance
  - Despite being constructed from entirely different inputs



# 2) Pion reconstruction

#### **Pion reconstruction**

- Goal:
  - > Classify pions as charged ( $\pi^{+}$ ) or neutral ( $\pi^{0}$ )
  - Regress true pion energy
- Why?
  - ATLAS has non-compensating calorimetry<sup>1</sup>
    - → Deposits from charged and neutral pions need to be restored to different scales
  - First step towards larger goal of using ML in jet energy reconstruction (eg. for Particle Flow)
- ML methods:
  - Convolutional Neural Network (CNN)
  - > Deep Neural Network (DNN) regression only
  - DeepSets
  - Graph Neural Network (GNN)
  - Transformer regression only

#### Input representations:





# Point clouds → 3D vectors of individual topocluster cells and tracks → Target truth particle energy instead of truth cluster energy

#### ATL-PHYS-PUB-2022-040

#### **Pion classification**

- All methods outperform baseline
- Excellent performance for **GNN** in particular



### **Pion energy calibration**

#### Without tracking information



- All ML models significantly outperform EM and LCW baselines
- Including tracking information significantly improves performance
- Up to ~50 GeV, the results approximate tracker energy resolution
  - Beyond this, calorimetry energy resolution dominates
  - ML giving the best of both
  - Advantage of point-cloud methods



# 3) Boosted jet tagging

## **Boosted jet tagging**

- Taggers distinguish large-radius jets from massive particles (W/Z/top) from light quark/gluon-initiated jets
  - Use jet substructure information
  - Enhances performance of bSM searches and precision SM measurements
- Latest taggers use jets reconstructed from Unified Flow Objects<sup>1</sup>
  - Combine topocluster and tracking information
  - Improved pile-up resilience and jet mass resolution
- ML classification techniques can improve on previous cut-based taggers



Latest taggers use kinematic properties of **jet constituents** rather than **high-level** jet substructure quantities

- Impressive performance on simplified Delphes simulated data-sets
- How do they perform with realistic GEANT4-simulated samples?



#### **Top taggers**

- 1) Baseline **DNN** trained on **high-level** input features (from <u>ATL-PHYS-PUB-2021-028</u>)
- 2) **DNN** trained on **constituent-level** input features
- 3) ParticleNet (arXiv:1902.08570)
  - Graph neural network (GNN)
  - Jets represented as graph (nodes = constituents)
- 4) Energy Flow Network (EFN) (arXiv:1810.05165)
  - DeepSets structure
  - > Can only consider constituent-level quantities linear in  $p_{T}$
- 5) Particle Flow Network (PFN) (arXiv:1810.05165)
  - Similar to EFN but permits all constituent-level quantities
- 6) ResNet50 (arXiv:1512.03385)
  - Large-scale convolutional neural network (CNN)
  - > Trained on "images" of constituent  $p_T$  in  $(\eta, \phi)$  plane

16

#### **Top taggers**

- 3 / 5 constituent-based taggers outperform hIDNN baseline
- EFN and ResNet50 underperform relative to Delphes-based studies
  - Highlights the need to develop taggers in a realistic context
- Data-set publicly available for future development: <u>https://gitlab.cern.ch/atlas/ATLAS-top-tagging-open-data</u>
- Work to assess systematic uncertainties ongoing



#### **Top Taggers: Model-dependence**

- **PFN** and **ParticleNet** show increased model dependence
  - More dependent on QCD modelling than high-level tagger
- EFN shows lowest model dependence due to IRC safety constraint



## W/Z taggers

#### ATL-PHYS-PUB-2021-029 JETM-2022-006

- DNN has improved performance but high mass correlation
  - Complicates background estimation strategies
- Adversarial NN used to de-correlate jet mass
  - Corresponding decrease in performance could be partially recovered with analysis-specific mass-window cuts
  - Allows use of better side-band regions



GeV

Fraction of jets / 5

 $10^{-2}$ 

 $10^{-3}$ 

10

ATLAS Simulation Preliminary

QCD jets

QCD jets after ZNN cut

QCD jets after D cut

W jets

 $\sqrt{s} = 13 \text{ TeV}, W$  jet tagging anti- $k_1 \text{ R}=1.0 \text{ UFO Soft-Drop CS+SK jets}$ Cuts at  $\varepsilon_{eq}^{eej} = 50\%, p_{\tau} \in [500, 1000] \text{ GeV}$ 



- Hadronic object reconstruction is ripe for ML applications
- Calibration (predicting truth-level energy/pT with detector-level quantities)
  - > METNet
  - Pion energy calibration
  - Jet calibration (See <u>Naseem Bouchhar's talk</u>)
- Classification
  - Boosted top and W/Z taggers
  - $\succ$   $\pi^{+}$  vs.  $\pi^{0}$
- Development of all of these applications is ongoing and promises to enhance the performance of bSM searches and precision SM measurements

# BACKUP



#### **METNet set-up**

- Dense neural network regresses truth-level p,<sup>miss</sup>
  - Trained on **ttbar** and **di-boson** MC events.  $\succ$
- **Sample weights** flatten the  $p_{\tau}^{miss}$  distribution of the training data  $\rightarrow$  reduces bias
- Loss functions: Huber + Sinkhorn **Pre-processing**



Standardize

inputs



## $p_{T}^{miss}$ working points

			Selections	
$p_{\rm T}$ [GeV] for jets with:			fJVT for jets with	
Working point	$ \eta  < 2.4$	$2.4 <  \eta  < 4.5$	JVT for jets with $ \eta  < 2.4$	$2.5 <  \eta  < 4.5$ and $p_{\rm T} < 120~{\rm GeV}$
Loose	> 20	> 20	$>0.5$ for $p_{\rm T}<60{\rm GeV}$ jets	-
Tight	> 20	> 30	$> 0.5$ for $p_{\rm T} < 60 {\rm GeV}$ jets	< 0.4
Tighter	> 20	> 35	$> 0.5$ for $p_{\rm T} < 60 {\rm GeV}$ jets	-
Tenacious	> 20	> 35	$> 0.91$ for $20 < p_{\rm T} < 40 {\rm GeV}$ jets	< 0.5
			$> 0.59$ for $40 < p_{\rm T} < 60 {\rm GeV}$ jets	
			$> 0.11$ for $60 < p_{\rm T} < 120 {\rm GeV}$ jets	

#### **METNet set-up**

Layer	Nodes	Activation	Parameters
Input	60	None	0
Hidden 1	100	SELU	6100
LayerNorm 1	100	SELU	200
Hidden 2	100	SELU	10100
LayerNorm 2	100	SELU	200
Hidden 3	100	SELU	10100
LayerNorm 3	100	SELU	200
Output	2	Linear	202
Total			27,102

Hyperparameter	Value
Optimiser	Adam [67]
Weight initialization	Kaiming He [68]
Learning rate	0.001
Batch size	256
Huber loss $\delta$	1.5

#### **Top Taggers: Model-dependence**



#### W/Z taggers: Model-dependence

- Sensitive to modelling differences between MC generators
- Also sensitive to W boson polarization
- Further studies and careful calibration required to avoid large systematic uncertainties



## Top tagger: ML set-up

Model	Hyper-parameters
PerNet 50	Bottom Layer: 7x7 2D convolution with strides (2, 2) and zero
	padding
	Number of Stages: 4
	Blocks per Stage: $(3, 4, 6, 3)$
	Block Type: bottleneck
	Block Output Filters: (64, 128, 256, 512)
Residet 50	Activation Functions: ReLU
	Kernel Initialization: he uniform
	Batch Normalization Momentum: 0.1
	Global Pooling: average
	Initial Learning Rate: $1 \times 10^{-2}$
	Scheduler: decrease learning rate by factor of 0.1 every 10 epochs
	Batch Size: 256
	$\Phi$ Number of Stages: 3
	Blocks per Stage: $(3, 3, 3)$
	Block Output Features: $(64, 224, 384)$
	k Nearest Neighbors: 18
ParticleNet	Top Layer Nodes: 125
	Activation Functions: ReLU
	Kernel Initialization: glorot normal
	Batch Normalization Momentum: 0.7
	Global Pooling: max
	Learning Rate: $4.2 \times 10^{-4}$
	Batch Size: 250

$ \begin{array}{c c} \mbox{Hidden Layers: 5} \\ \mbox{Nodes per Layer: 180} \\ \mbox{Activation Functions: ReLU} \\ \mbox{Activation Functions: ReLU} \\ \mbox{Kernel Initialization: glorot uniform} \\ \mbox{Learning Rate: 4 \times 10^{-5}} \\ \mbox{Batch Normalization: not used} \\ \mbox{Hidden Layers: 5} \\ \mbox{Nodes per Layer: 400} \\ \mbox{Activation Functions: ReLU} \\ \mbox{Kernel Initialization: glorot uniform} \\ \mbox{L1 Regularization: 2 \times 10^{-4}, applied to all layers} \\ \mbox{Learning Rate: 1.2 \times 10^{-5}} \\ \mbox{Batch Normalization: applied before activation function for all} \\ \mbox{layers except output layer} \\ \mbox{\Phi Hidden Layers: 5} \\ \mbox{\Phi Nodes per Layer: 350} \\ \mbox{Latent Dropout: 0.084} \\ \mbox{F Hidden Layers: 5} \\ \mbox{F Nodes per Layer: 300} \\ \mbox{F Dropout: 0.036} \\ \mbox{Activation Functions: ReLU} \\ \mbox{Kernel Initialization: glorot normal} \\ \mbox{Learning Rate: 6.3 \times 10^{-5}} \\ \mbox{Batch Nize: 350} \\ \mbox{\Phi Nodes per Layer: 250} \\ \mbox{Latent Dropout: 0.072} \\ \mbox{F Hidden Layers: 5} \\ \mbox{\Phi Nodes per Layer: 500} \\ \mbox{F Dropout: 0.022} \\ \mbox{Activation Functions: ReLU} \\ \mbox{Kernel Initialization: glorot normal} \\ \mbox{Learning Rate: 6.3 \times 10^{-5}} \\ \mbox{Batch Nize: 350} \\ \mbox{PFN} \end{array}$	Model	Hyper-parameters		
$ \begin{array}{c c} \mbox{Nodes per Layer: 180} \\ \mbox{Activation Functions: ReLU} \\ \mbox{Activation Functions: glorot uniform} \\ \mbox{Learning Rate: } 4 \times 10^{-5} \\ \mbox{Batch Size: 250} \\ \mbox{Batch Normalization: not used} \\ \hline \\ \mbox{Hidden Layers: 5} \\ \mbox{Nodes per Layer: 400} \\ \mbox{Activation Functions: ReLU} \\ \mbox{Kernel Initialization: glorot uniform} \\ \mbox{L1 Regularization: } 2 \times 10^{-4}, applied to all layers \\ \mbox{Learning Rate: } 1.2 \times 10^{-5} \\ \mbox{Batch Size: } 250 \\ \mbox{Batch Normalization: applied before activation function for all \\ \mbox{layers except output layer} \\ \hline \\ \mbox{\Phi Hidden Layers: 5} \\ \mbox{\Phi Nodes per Layer: } 350 \\ \mbox{Latent Dropout: } 0.084 \\ F & \mbox{Hidden Layers: 5} \\ \hline \\ \mbox{F Nodes per Layer: } 300 \\ F & \mbox{Dropout: } 0.036 \\ \mbox{Activation Functions: ReLU} \\ \mbox{Kernel Initialization: glorot normal} \\ \mbox{Learning Rate: } 6.3 \times 10^{-5} \\ \mbox{Batch Size: } 350 \\ \hline \\ \mbox{\Phi Hidden Layers: 5} \\ \hline \\ \mbox{PFN} \\ \mbox{F Hidden Layers: 5} \\ F & \mbox{Nodes per Layer: } 250 \\ \mbox{Latent Dropout: } 0.072 \\ F & \mbox{Hidden Layers: 5} \\ \hline \\ \mbox{PFN} \\ \mbox{F Dropout: } 0.022 \\ \mbox{Activation Functions: ReLU} \\ \mbox{Kernel Initialization: glorot normal} \\ \mbox{Learning Rate: } 7.9 \times 10^{-5} \\ \hline \\ \mbox{PFN} \\ \mbox{F Nodes per Layer: } 500 \\ F & \mbox{Dropout: } 0.022 \\ \mbox{Activation Functions: ReLU} \\ \mbox{Kernel Initialization: glorot normal} \\ \mbox{Learning Rate: } 7.9 \times 10^{-5} \\ \hline \\ \mbox{PFN} \\ \mbox{F Hidden Layers: 5} \\ \mbox{F Nodes per Layer: 500} \\ \mbox{F Dropout: } 0.022 \\ \mbox{Activation Functions: ReLU} \\ \mbox{Kernel Initialization: glorot normal} \\ \mbox{Learning Rate: } 7.9 \times 10^{-5} \\ \hline \\ \mbox{F Nodes per Layer: } 500 \\ \mbox{F Dropout: } 0.022 \\ \mbox{Activation Functions: ReLU} \\ \mbox{Kernel Initialization: glorot normal} \\ \mbox{Learning Rate: } 7.9 \times 10^{-5} \\ \hline \\ \mbox{F Nodes per Layer: } 50 \\ \mbox{F Notes per Layer: } 50 \\ \mbox{F Notes per Layer: } 500 \\ \mbox{F Notes per Layer: } 500 \\ \hline \\ F Notes pe$		Hidden Layers: 5		
hIDNNActivation Functions: ReLU Kernel Initialization: glorot uniform Learning Rate: $4 \times 10^{-5}$ Batch Normalization: not usedHidden Layers: 5 Nodes per Layer: 400 Activation Functions: ReLU Kernel Initialization: glorot uniform L1 Regularization: $2 \times 10^{-4}$ , applied to all layers Learning Rate: $1.2 \times 10^{-5}$ Batch Normalization: applied before activation function for all layers except output layer $\Phi$ Hidden Layers: 5 $\Phi$ Nodes per Layer: 350 Latent Dropout: 0.084 $F$ Hidden Layers: 5 $F$ Nodes per Layer: 300 $F$ Dropout: 0.036 Activation Functions: ReLU Kernel Initialization: glorot normal Learning Rate: $6.3 \times 10^{-5}$ Batch Size: 350EFN $\Phi$ Hidden Layers: 5 $F$ Nodes per Layer: 300 $F$ Dropout: 0.036 Activation functions: ReLU Kernel Initialization: glorot normal Learning Rate: $6.3 \times 10^{-5}$ Batch Size: 350PFN $F$ Hidden Layers: 5 $F$ Nodes per Layer: 250 Latent Dropout: 0.072 $F$ Hidden Layers: 5 $F$ Nodes per Layer: 500 $F$ Dropout: 0.022 Activation Functions: ReLU Kernel Initialization: glorot normal Learning Rate: $7.9 \times 10^{-5}$ $F$ Nodes per Layer: 500 $F$ Dropout: 0.022 $Acitvation Functions: ReLUKernel Initialization: glorot normalLearning Rate: 7.9 \times 10^{-5}$	HDNN	Nodes per Layer: 180		
$\begin{array}{c c} \text{InDNA} & Kernel Initialization: glorot uniform \\ Learning Rate: 4 \times 10^{-5} \\ \text{Batch Size: 250} \\ \text{Batch Normalization: not used} \\ & \text{Hidden Layers: 5} \\ \text{Nodes per Layer: 400} \\ & Activation Functions: ReLU \\ Kernel Initialization: glorot uniform \\ L1 Regularization: 2 \times 10^{-4}, applied to all layers \\ Learning Rate: 1.2 \times 10^{-5} \\ \text{Batch Normalization: applied before activation function for all layers except output layer \\ \hline \Phi \text{Hidden Layers: 5} \\ & \Phi \text{Nodes per Layer: 350} \\ \text{Latent Dropout: 0.084} \\ F \text{ Hidden Layers: 5} \\ F \text{ Nodes per Layer: 300} \\ F \text{ Dropout: 0.036} \\ Activation Functions: ReLU \\ Kernel Initialization: glorot normal \\ Learning Rate: 6.3 \times 10^{-5} \\ \text{Batch Size: 350} \\ \hline \Phi \text{ Hidden Layers: 5} \\ F \text{ Nodes per Layer: 500} \\ Latent Dropout: 0.072 \\ F \text{ Hidden Layers: 5} \\ F \text{ Nodes per Layer: 500} \\ Latent Dropout: 0.072 \\ F \text{ Hidden Layers: 5} \\ F \text{ Nodes per Layer: 500} \\ Latent Dropout: 0.072 \\ F \text{ Hidden Layers: 5} \\ F \text{ Nodes per Layer: 500} \\ Latent Dropout: 0.072 \\ F \text{ Hidden Layers: 5} \\ F \text{ Nodes per Layer: 500} \\ Latent Dropout: 0.072 \\ F \text{ Hidden Layers: 5} \\ F \text{ Nodes per Layer: 500} \\ Latent Dropout: 0.072 \\ F \text{ Hidden Layers: 5} \\ F \text{ Nodes per Layer: 500} \\ F \text{ Dropout: 0.022} \\ Activation Functions: ReLU \\ Kernel Initialization: glorot normal \\ Learning Rate: 7.9 \times 10^{-5} \\ \hline \end{array}$		Activation Functions: ReLU		
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	mbrit	Kernel Initialization: glorot uniform		
$\begin{array}{ c c c c c c c c c c c c c c c c c c c$		Learning Rate: $4 \times 10^{-5}$		
$\begin{array}{ c c c c c c } & \text{Batch Normalization: not used} \\ & \text{Hidden Layers: 5} \\ & \text{Nodes per Layer: 400} \\ & Activation Functions: ReLU \\ & Kernel Initialization: glorot uniform \\ & \text{L1 Regularization: } 2 \times 10^{-4}, applied to all layers \\ & \text{Learning Rate: } 1.2 \times 10^{-5} \\ & \text{Batch Size: } 250 \\ & \text{Batch Normalization: applied before activation function for all } \\ & \text{layers except output layer} \\ \hline & \Phi \text{ Hidden Layers: 5} \\ & \Phi \text{ Nodes per Layer: 350} \\ & \text{Latent Dropout: } 0.084 \\ & F \text{ Hidden Layers: 5} \\ & F \text{ Nodes per Layer: 300} \\ & F \text{ Dropout: } 0.036 \\ & Activation Functions: ReLU \\ & Kernel Initialization: glorot normal \\ & \text{Learning Rate: } 6.3 \times 10^{-5} \\ & \text{Batch Size: 350} \\ \hline & \Phi \text{ Hidden Layers: 5} \\ & F \text{ Nodes per Layer: 5} \\ & \Phi \text{ Nodes per Layer: 5} \\ & F \text{ Hidden Layers: 5} \\ & F \text{ Nodes per Layer: 5} \\ & F \text{ Nodes per Layer: 500} \\ & \text{Latent Dropout: 0.072} \\ & F \text{ Hidden Layers: 5} \\ & F \text{ Nodes per Layer: 500} \\ & F \text{ Dropout: 0.022} \\ & Acitvation Functions: ReLU \\ & Kernel Initialization: glorot normal \\ & \text{Learning Rate: 7.9 \times 10^{-5} \\ \hline & \text{Det Midden Layers: 5} \\ & \text{Dropout: 0.022} \\ & Acitvation Functions: ReLU \\ & Kernel Initialization: glorot normal \\ & \text{Learning Rate: 7.9 \times 10^{-5} \\ \hline & \text{Det Midden Layers: 5} \\ & \text{Dropout: 0.022} \\ & Acitvation Functions: ReLU \\ & Kernel Initialization: glorot normal \\ & \text{Learning Rate: 7.9 \times 10^{-5} \\ \hline & \text{Det Midden Layers: 5} \\ & \text{Dropout: 0.021} \\ & \text{Acitvation Functions: ReLU} \\ & \text{Kernel Initialization: glorot normal} \\ & \text{Learning Rate: 7.9 \times 10^{-5} \\ \hline & \text{Dropout: 0.021} \\ & \text{Acitvation Functions: ReLU} \\ & \text{Kernel Initialization: glorot normal} \\ & \text{Learning Rate: 7.9 \times 10^{-5} \\ \hline & \text{Dropout: 0.022} \\ & \text{Acitvation Functions: ReLU} \\ & \text{Kernel Initialization: glorot normal} \\ & \text{Learning Rate: 7.9 \times 10^{-5} \\ \hline & \text{Dropout: 0.021} \\ & \text{Acitvation Functions: ReLU} \\ & Ac$		Batch Size: 250		
$\begin{array}{c c} \mbox{Hidden Layers: 5}\\ \mbox{Nodes per Layer: 400}\\ \mbox{Activation Functions: ReLU}\\ \mbox{Kernel Initialization: glorot uniform}\\ \mbox{L1 Regularization: 2 \times 10^{-4}, applied to all layers}\\ \mbox{Learning Rate: } 1.2 \times 10^{-5}\\ \mbox{Batch Size: 250}\\ \mbox{Batch Normalization: applied before activation function for all}\\ \mbox{layers except output layer}\\ \hline \mbox{\Phi Hidden Layers: 5}\\ \mbox{\Phi Nodes per Layer: 350}\\ \mbox{Latent Dropout: 0.084}\\ F \mbox{Hidden Layers: 5}\\ F \mbox{Nodes per Layer: 300}\\ F \mbox{Dropout: 0.036}\\ \mbox{Activation Functions: ReLU}\\ \mbox{Kernel Initialization: glorot normal}\\ \mbox{Learning Rate: } 6.3 \times 10^{-5}\\ \mbox{Batch Size: 350}\\ \hline \mbox{\Phi Hidden Layers: 5}\\ F \mbox{Nodes per Layer: 250}\\ \mbox{Latent Dropout: 0.072}\\ F \mbox{Hidden Layers: 5}\\ F \mbox{Nodes per Layer: 500}\\ \mbox{Latent Dropout: 0.072}\\ F \mbox{Hidden Layers: 5}\\ F \mbox{Nodes per Layer: 500}\\ F \mbox{Dropout: 0.022}\\ \mbox{Activation Functions: ReLU}\\ \mbox{Kernel Initialization: glorot normal}\\ \mbox{Learning Rate: 7.9 \times 10^{-5}}\\ \hline \mbox{PFN} \end{tabular}$		Batch Normalization: not used		
$\begin{array}{c c} & \operatorname{Nodes \ per \ Layer: \ 400} \\ & Activation \ Functions: \ ReLU \\ & Kernel \ Initialization: \ glorot \ uniform \\ & L1 \ Regularization: \ 2 \times 10^{-4}, \ applied \ to \ all \ layers \\ & Learning \ Rate: \ 1.2 \times 10^{-5} \\ & Batch \ Size: \ 250 \\ & Batch \ Normalization: \ applied \ before \ activation \ function \ for \ all \\ & layers \ except \ output \ layer \\ \hline & \Phi \ Hidden \ Layers: \ 5 \\ & \Phi \ Nodes \ per \ Layer: \ 350 \\ & Latent \ Dropout: \ 0.084 \\ & F \ Hidden \ Layers: \ 5 \\ \hline & F \ Nodes \ per \ Layer: \ 300 \\ & F \ Dropout: \ 0.036 \\ & Activation \ Functions: \ ReLU \\ & Kernel \ Initialization: \ glorot \ normal \\ & Learning \ Rate: \ 6.3 \times 10^{-5} \\ \hline & Batch \ Size: \ 350 \\ \hline & \Phi \ Hidden \ Layers: \ 5 \\ \hline & \Phi \ Nodes \ per \ Layer: \ 250 \\ & Latent \ Dropout: \ 0.072 \\ & F \ Hidden \ Layers: \ 5 \\ \hline & \Phi \ Nodes \ per \ Layer: \ 250 \\ & Latent \ Dropout: \ 0.072 \\ & F \ Hidden \ Layers: \ 5 \\ \hline & F \ Nodes \ per \ Layer: \ 250 \\ & Latent \ Dropout: \ 0.072 \\ & F \ Hidden \ Layers: \ 5 \\ \hline & F \ Nodes \ per \ Layer: \ 250 \\ & Latent \ Dropout: \ 0.072 \\ & F \ Hidden \ Layers: \ 5 \\ \hline & F \ Nodes \ per \ Layer: \ 250 \\ & Latent \ Dropout: \ 0.072 \\ & F \ Hidden \ Layers: \ 5 \\ \hline & F \ Nodes \ per \ Layer: \ 250 \\ & Latent \ Dropout: \ 0.072 \\ & F \ Hidden \ Layers: \ 5 \\ \hline & F \ Nodes \ per \ Layer: \ 500 \\ & F \ Dropout: \ 0.022 \\ & Acitvation \ Functions: \ ReLU \\ & Kernel \ Initialization: \ glorot \ normal \\ & Learning \ Rate: \ 7.9 \times 10^{-5} \\ \hline & P \ Hidden \ Layers: \ 5 \\ \hline & P \ Hidden \ Layers: \ 5 \\ \hline & P \ Layers: \ 5 \\ \hline & F \ Nodes \ P \ Layers: \ 5 \\ \hline & F \ Nodes \ P \ Layers: \ 5 \\ \hline & F \ Nodes \ P \ Layers: \ 5 \\ \hline & F \ Nodes \ P \ Layers: \ 5 \\ \hline & F \ Nodes \ P \ Layers: \ 5 \\ \hline & F \ Nodes \ P \ Layers: \ 5 \\ \hline & F \ Nodes \ P \ Layers: \ 5 \\ \hline & F \ Nodes \ P \ Layers: \ 5 \\ \hline & F \ Nodes \ P \ Layers: \ 5 \\ \hline & F \ Nodes \ P \ Layers: \ 5 \\ \hline & F \ Nodes \ P \ Layers: \ 5 \\ \hline & F \ Nodes \ Layers: \ 5 \\ \hline & F \ Nodes \ P \ Layers: \ 5 \\ \hline & F \ Nodes \ P \ Layers: $		Hidden Layers: 5		
$\begin{array}{c c} Activation Functions: ReLU \\ Kernel Initialization: glorot uniform \\ L1 Regularization: 2 \times 10^{-4}, applied to all layersLearning Rate: 1.2 \times 10^{-5} Batch Size: 250Batch Normalization: applied before activation function for alllayers except output layer$		Nodes per Layer: 400		
$\begin{array}{c c} \text{DNN} & \begin{matrix} \text{Kernel Initialization: glorot uniform} \\ \text{L1 Regularization: } 2 \times 10^{-4}, \text{ applied to all layers} \\ \text{Learning Rate: } 1.2 \times 10^{-5} \\ \text{Batch Size: } 250 \\ \text{Batch Normalization: applied before activation function for all} \\ \text{layers except output layer} \\ & \Psi \text{ Hidden Layers: } 5 \\ & \Psi \text{ Nodes per Layer: } 350 \\ \text{Latent Dropout: } 0.084 \\ & F \text{ Hidden Layers: } 5 \\ & F \text{ Nodes per Layer: } 300 \\ & F \text{ Dropout: } 0.036 \\ & Activation Functions: ReLU \\ & Kernel Initialization: glorot normal \\ & \text{Learning Rate: } 6.3 \times 10^{-5} \\ & \text{Batch Size: } 350 \\ \hline \\ & \Psi \text{ Hidden Layers: } 5 \\ & \Psi \text{ Nodes per Layer: } 250 \\ & \text{Latent Dropout: } 0.072 \\ & F \text{ Hidden Layers: } 5 \\ & F \text{ Nodes per Layer: } 500 \\ & F \text{ Dropout: } 0.022 \\ & Acitvation Functions: ReLU \\ & Kernel Initialization: glorot normal \\ & \text{Learning Rate: } 7.9 \times 10^{-5} \\ & \text{Peth}  H  Of methods of m$		Activation Functions: ReLU		
DINNL1 Regularization: $2 \times 10^{-4}$ , applied to all layers Learning Rate: $1.2 \times 10^{-5}$ Batch Size: 250 Batch Normalization: applied before activation function for all layers except output layer $\Phi$ Hidden Layers: 5 $\Phi$ Nodes per Layer: 350 Latent Dropout: 0.084 F Hidden Layers: 5 F Nodes per Layer: 300 F Dropout: 0.036 Activation Functions: ReLU Kernel Initialization: glorot normal Learning Rate: $6.3 \times 10^{-5}$ 	DNN	Kernel Initialization: glorot uniform		
$\begin{array}{c c} & \mbox{Learning Rate: } 1.2 \times 10^{-5} \\ & \mbox{Batch Size: } 250 \\ & \mbox{Batch Normalization: applied before activation function for all layers except output layer } \\ & $\Psi$ Hidden Layers: 5 $ $ $ $ $ $ $ $ $ $ $ $ $ $ $ $ $ $ $	DIVIN	L1 Regularization: $2 \times 10^{-4}$ , applied to all layers		
$\begin{array}{c c c c c c c c c c c c c c c c c c c $		Learning Rate: $1.2 \times 10^{-5}$		
$\begin{array}{c c c c c c c c c c c c c c c c c c c $		Batch Size: 250		
$\begin{array}{c c c c c c c c c c c c c c c c c c c $		Batch Normalization: applied before activation function for all		
$ \begin{array}{c} \Phi \mbox{ Hidden Layers: 5} \\ \Phi \mbox{ Nodes per Layer: 350} \\ \mbox{ Latent Dropout: 0.084} \\ F \mbox{ Hidden Layers: 5} \\ F \mbox{ Nodes per Layer: 300} \\ F \mbox{ Dropout: 0.036} \\ Activation Functions: ReLU \\ Kernel Initialization: glorot normal \\ \mbox{ Learning Rate: } 6.3 \times 10^{-5} \\ \mbox{ Batch Size: 350} \\ \hline \end{tabular} \\ \Phi \mbox{ Hidden Layers: 5} \\ \Phi \mbox{ Nodes per Layer: 250} \\ \mbox{ Latent Dropout: 0.072} \\ F \mbox{ Hidden Layers: 5} \\ F \mbox{ Nodes per Layer: 500} \\ F \mbox{ Dropout: 0.022} \\ Acitvation Functions: ReLU \\ Kernel Initialization: glorot normal \\ \mbox{ Learning Rate: } 7.9 \times 10^{-5} \\ \mbox{ Dropout: 0.076} \\ \hline \end{array} $		layers except output layer		
$ \begin{array}{c} \Phi \ {\rm Nodes \ per \ Layer: \ 350} \\ {\rm Latent \ Dropout: \ 0.084} \\ F \ {\rm Hidden \ Layers: \ 5} \\ F \ {\rm Nodes \ per \ Layer: \ 300} \\ F \ {\rm Dropout: \ 0.036} \\ Activation \ Functions: \ ReLU \\ Kernel \ Initialization: \ glorot \ normal \\ {\rm Learning \ Rate: \ 6.3 \times 10^{-5}} \\ {\rm Batch \ Size: \ 350} \\ \hline \\ \Phi \ {\rm Hidden \ Layers: \ 5} \\ \Phi \ {\rm Nodes \ per \ Layer: \ 250} \\ {\rm Latent \ Dropout: \ 0.072} \\ F \ {\rm Hidden \ Layers: \ 5} \\ F \ {\rm Nodes \ per \ Layer: \ 500} \\ F \ {\rm Dropout: \ 0.022} \\ Acitvation \ Functions: \ ReLU \\ Kernel \ Initialization: \ glorot \ normal \\ {\rm Learning \ Rate: \ 7.9 \times 10^{-5}} \\ \end{array} $		Φ Hidden Layers: 5		
$ \begin{array}{c} \mbox{Latent Dropout: } 0.084 \\ F \mbox{Hidden Layers: 5} \\ F \mbox{Nodes per Layer: 300} \\ F \mbox{Dropout: } 0.036 \\ Activation Functions: ReLU \\ Kernel Initialization: glorot normal \\ Learning Rate: 6.3 \times 10^{-5} \\ \mbox{Batch Size: } 350 \\ \hline \\ \Psi \mbox{Hidden Layers: 5} \\ \Phi \mbox{Nodes per Layer: } 250 \\ \mbox{Latent Dropout: } 0.072 \\ F \mbox{Hidden Layers: 5} \\ F \mbox{Nodes per Layer: 500} \\ F \mbox{Dropout: } 0.022 \\ Acitvation Functions: ReLU \\ Kernel Initialization: glorot normal \\ \mbox{Learning Rate: } 7.9 \times 10^{-5} \\ \hline \end{array} $		$\Phi$ Nodes per Layer: 350		
$ \begin{array}{c} F \mbox{ Hidden Layers: 5} \\ F \mbox{ Nodes per Layer: 300} \\ F \mbox{ Dropout: 0.036} \\ Activation Functions: ReLU \\ Kernel Initialization: glorot normal \\ Learning Rate: 6.3 \times 10^{-5} \\ \mbox{ Batch Size: 350} \\ \hline \end{tabular} \\ \hline \end{tabular} \\ \hline \end{tabular} \\ \hline \end{tabular} \\ F \mbox{ Hidden Layers: 5} \\ \Phi \mbox{ Nodes per Layer: 250} \\ \mbox{ Latent Dropout: 0.072} \\ F \mbox{ Hidden Layers: 5} \\ F \mbox{ Nodes per Layer: 500} \\ F \mbox{ Dropout: 0.022} \\ Acitvation Functions: ReLU \\ Kernel Initialization: glorot normal \\ \mbox{ Learning Rate: 7.9 \times 10^{-5} \\ \hline \end{tabular} \\ \hline \end{tabular} \\ \hline \end{tabular} $		Latent Dropout: 0.084		
$ \begin{array}{c} {\rm EFN} & F \ {\rm Nodes \ per \ Layer: \ 300} \\ F \ {\rm Dropout: \ 0.036} \\ Activation \ Functions: \ ReLU \\ Kernel \ Initialization: \ glorot \ normal \\ {\rm Learning \ Rate: \ 6.3 \times 10^{-5}} \\ {\rm Batch \ Size: \ 350} \\ \hline \\ \hline \\ \Phi \ {\rm Hidden \ Layers: \ 5} \\ \Phi \ {\rm Nodes \ per \ Layer: \ 250} \\ {\rm Latent \ Dropout: \ 0.072} \\ F \ {\rm Hidden \ Layers: \ 5} \\ F \ {\rm Nodes \ per \ Layer: \ 500} \\ F \ {\rm Dropout: \ 0.022} \\ Acitvation \ Functions: \ ReLU \\ Kernel \ Initialization: \ glorot \ normal \\ {\rm Learning \ Rate: \ 7.9 \times 10^{-5}} \\ \hline \end{array} $		F Hidden Layers: 5		
EFN F Dropout: 0.036 Activation Functions: ReLU Kernel Initialization: glorot normal Learning Rate: $6.3 \times 10^{-5}$ Batch Size: 350 $\Phi$ Hidden Layers: 5 $\Phi$ Nodes per Layer: 250 Latent Dropout: 0.072 F Hidden Layers: 5 F Nodes per Layer: 500 F Dropout: 0.022 Activation Functions: ReLU Kernel Initialization: glorot normal Learning Rate: $7.9 \times 10^{-5}$	FFN	F Nodes per Layer: 300		
Activation Functions: $ReLU$ Kernel Initialization: glorot normal Learning Rate: $6.3 \times 10^{-5}$ Batch Size: $350$ $\Phi$ Hidden Layers: 5 $\Phi$ Nodes per Layer: 250 Latent Dropout: $0.072$ $F$ Hidden Layers: 5 $F$ Nodes per Layer: 500 $F$ Dropout: $0.022$ Acitvation Functions: $ReLU$ Kernel Initialization: glorot normal Learning Rate: $7.9 \times 10^{-5}$ De the Gime Gime State St	EFIN	F Dropout: 0.036		
$\begin{array}{c c} & Kernel \ Initialization: \ glorot \ normal \\ \ Learning \ Rate: \ 6.3 \times 10^{-5} \\ & \ Batch \ Size: \ 350 \\ \hline \\ & \Phi \ Hidden \ Layers: \ 5 \\ & \Phi \ Nodes \ per \ Layer: \ 250 \\ & \ Latent \ Dropout: \ 0.072 \\ & F \ Hidden \ Layers: \ 5 \\ & F \ Nodes \ per \ Layer: \ 500 \\ & F \ Dropout: \ 0.022 \\ & \ Acitvation \ Functions: \ ReLU \\ & \ Kernel \ Initialization: \ glorot \ normal \\ & \ Learning \ Rate: \ 7.9 \times 10^{-5} \\ \hline \end{array}$		Activation Functions: ReLU		
$\begin{array}{c c} & \text{Learning Rate: } 6.3 \times 10^{-5} \\ & \text{Batch Size: } 350 \\ \hline & \Phi \text{ Hidden Layers: } 5 \\ & \Phi \text{ Nodes per Layer: } 250 \\ & \text{Latent Dropout: } 0.072 \\ & F \text{ Hidden Layers: } 5 \\ & F \text{ Nodes per Layer: } 500 \\ & F \text{ Dropout: } 0.022 \\ & Acitvation Functions: ReLU \\ & Kernel Initialization: glorot normal \\ & \text{Learning Rate: } 7.9 \times 10^{-5} \\ & \text{Petric} = 1000 \\ \hline & \text{Rection Functions} = 1000 \\ \hline & \text{Rection Functions} = 1000 \\ \hline & \text{Rection Functions} = 1000 \\ \hline & \text{Rection Function} = 10000 \\ \hline & \text{Rection Function} = 1000 \\ \hline & \text{Rection Function Function} = 1000 \\ \hline & \text{Rection Function Function} = 1000 \\ \hline & Rection Function Functi$		Kernel Initialization: glorot normal		
Batch Size: 350 $\Phi$ Hidden Layers: 5 $\Phi$ Nodes per Layer: 250Latent Dropout: 0.072 $F$ Hidden Layers: 5 $F$ Nodes per Layer: 500 $F$ Dropout: 0.022 $Acitvation Functions: ReLU$ $Kernel Initialization: glorot normal$ Learning Rate: $7.9 \times 10^{-5}$ Performed for the formed for the formed for		Learning Rate: $6.3 \times 10^{-5}$		
$ \begin{array}{l} \Phi \mbox{ Hidden Layers: 5} \\ \Phi \mbox{ Nodes per Layer: 250} \\ \mbox{ Latent Dropout: 0.072} \\ F \mbox{ Hidden Layers: 5} \\ F \mbox{ Nodes per Layer: 500} \\ F \mbox{ Dropout: 0.022} \\ Acitvation Functions: ReLU \\ Kernel Initialization: glorot normal \\ \mbox{ Learning Rate: } 7.9 \times 10^{-5} \\ \end{array} $		Batch Size: 350		
$ \begin{array}{l} \Phi \text{ Nodes per Layer: } 250 \\ \text{Latent Dropout: } 0.072 \\ F \text{ Hidden Layers: } 5 \\ F \text{ Nodes per Layer: } 500 \\ F \text{ Dropout: } 0.022 \\ Acitvation Functions: ReLU \\ Kernel Initialization: glorot normal \\ \text{Learning Rate: } 7.9 \times 10^{-5} \\ P  Prime Prim$		$\Phi$ Hidden Layers: 5		
PFN Latent Dropout: 0.072 F Hidden Layers: 5 F Nodes per Layer: 500 F Dropout: 0.022 Acitvation Functions: ReLU Kernel Initialization: glorot normal Learning Rate: $7.9 \times 10^{-5}$		$\Phi$ Nodes per Layer: 250		
$\begin{array}{c} F \mbox{ Hidden Layers: 5} \\ F \mbox{ Nodes per Layer: 500} \\ F \mbox{ Dropout: 0.022} \\ Acitvation Functions: ReLU \\ Kernel Initialization: glorot normal \\ Learning Rate: 7.9 \times 10^{-5} \\ P \mbox{ Product} \end{array}$	PFN	Latent Dropout: 0.072		
PFN F Nodes per Layer: 500 F Dropout: 0.022 Acitvation Functions: ReLU Kernel Initialization: glorot normal Learning Rate: $7.9 \times 10^{-5}$ Production Functions		F Hidden Layers: 5		
F Dropout: 0.022 Acitvation Functions: ReLU Kernel Initialization: glorot normal Learning Rate: $7.9 \times 10^{-5}$		F Nodes per Layer: 500		
Acitvation Functions: ReLU Kernel Initialization: glorot normal Learning Rate: $7.9 \times 10^{-5}$		F Dropout: 0.022		
Kernel Initialization: glorot normal Learning Rate: $7.9 \times 10^{-5}$		Acitvation Functions: ReLU		
Learning Rate: $7.9 \times 10^{-5}$		Kernel Initialization: glorot normal		
D I I GL OFO		Learning Rate: $7.9 \times 10^{-5}$		
Batch Size: 250		Batch Size: 250		

#### **Top-tagger: Baseline DNN inputs**

Quantity Type	Symbols	References
N-subjettiness	$ au_1,   au_2,   au_3,   au_4$	[66] $[67]$
$k_t$ Splitting Scales	$\sqrt{d_{12}}, \sqrt{d_{23}}$	[15]
Generalized Energy Correlation Functions	$ECF_1, ECF_2, ECF_3, C_2, D_2, L_2, L_3$	[68] $[69]$ $[70]$
Minimum Pair-wise Invariant Mass	$Q_w$	[15]
Thrust Major	$T_m$	[71]

Table 1: A listing of the 15 quantities used to train the baseline high level quantity tagger