A Review of Machine Learning Applications on Jet Tagging

Konkuk University Daohan Wang 2023.6.8

Brief motivation of ML in Jet Tagging



Simulation Details



We can use Deep learning to analyze low-level LHC data without constructing high-level observables !

The image-based representation is based on the reconstruction of jets with calorimeters. A calorimeter measures the energy deposition of a jet on fine-grained spatial cells. Treating the energy deposition on each cell as the pixel intensity naturally creates an image for a jet.

When jets are formed by particles reconstructed with the full detector information a jet image can be constructed by mapping each particle onto the corresponding calorimeter cell and sum up the energy if more than one particle is mapped to the same cell.

We can construct different channels to characterize more features. For example, based on energy flow algorithm in Delphes, we can construct 3 channels for EflowPhoton, EflowNeutralHadron and EflowTracks, respectively.

Introduction to CNN



CNN Architecture



Average Jet Images



Two shortcomings of image-based representation for jet tagging:

1, Treating jets as images leads to a very sparse representation.

Without considering the pileup effects, a typical jet has O(10) to O(100) particles, while a jet image typically needs O(1000) pixels (eg.32×32) in order to fully contain the jet, more then 90% of the pixels are blank.

2, How to incorporate additional information of the particles is unclear.

As it involves combining nonadditive quantities (e.g., the particle type) of multiple particles entering the same cell.

We organize a jet's constituent particles into an ordered structure (sequence or tree) based on the p_T .

Dense Network

We start with N p_T -sorted particles per jet, the arguably simplest deep network architecture is a dense network taking all (p_T, η, ϕ) information as a fixed set. We again improve the training through physics-motivated pre-processing.





- Sigmoid Function:
- The sigmoid function is a commonly used activation function defined as $\frac{1}{e^{-x}+1}$, where x is the input. It produces values between 0 and 1, which can be interpreted as probabilities.
- Tanh (tanch):
- The tanh function is similar to the sigmoid function but produces values between -1 and 1.
- ReLU (Rectified Linear Unit):
- The ReLU (Rectified Linear Unit) function is a popular activation function in deep learning. It is defined as f(x) = max(0, x), which means that the output is 0 for any negative input and equal to the input for any positive input.
- LReLU (Leaky ReLU)
- The LeakyReLU function is a modification of the ReLU function that solves the "dying ReLU" problem, where ReLU units can get stuck in the zero state during training. It is defined as $f(x) = max(\alpha x, x)$, where α is a small positive constant that determines the slope of the function for negative inputs.

Each activation function has its advantages and disadvantages, and the choice of function depends on the specific task and network architecture.



Introduction to Particle-based Representation



P-CNN

The particle-level convolutional neural network (P-CNN) is a customized 1-dimensional CNN for jet tagging. Each input jet is represented as a sequence of constituents with a fixed length of N, organized in descending order of p_T . For each constituent, several input features are computed from the 4-momenta of the constituent and used as inputs to the network. The P-CNN is similar to CNN, but only uses a 1-dimensional convolution instead of 2-dimensional convolutions.

Two shortcomings of particle-based representation for jet tagging:

1, The jet length are variable.

As each jet may contain a different number of particles.

2, The particles are needed to be sorted in some way.

The constituent particles in a jet have no intrinsic order; thus, the manually imposed order may turn out to be suboptimal and impair the performance.

We consider a jet as an unordered set of its constituent particles. more natural ! **Permutation Invariant** IRC-safe Energy Flow network Deep Sets Framework IRC-not-safe Particle Flow network ParticleNet Typical Architectures ABCNet LorentzNet Transformer

The Deep Sets framework was adapted and specialized to particle physics in 1810.05165

The Deep Sets framework for point clouds demonstrates how permutation-invariant functions of variable-length inputs can be parametrized in a fully general way and it enables a natural visualization of the learned latent space, providing insights as to what exactly the NN is learning.

Observable Decomposition. An observable \mathcal{O} can be approximated arbitrarily well as:

$$\mathcal{O}(\{p_1,\ldots,p_M\}) = F\left(\sum_{i=1}^M \Phi(p_i)\right),\tag{1.1}$$

where $\Phi : \mathbb{R}^d \to \mathbb{R}^\ell$ is a per-particle mapping and $F : \mathbb{R}^\ell \to \mathbb{R}$ is a continuous function.

IRC safety corresponds to robustness of the observable under collinear splittings of a particle or additions of soft particles.

IRC-Safe Observable Decomposition. An IRC-safe observable \mathcal{O} can be approximated arbitrarily well as:

$$z_i = p_{T,i} / \sum_j p_{T,j}$$
 $\mathcal{O}(\{p_1, \dots, p_M\}) = F\left(\sum_{i=1}^M z_i \Phi(\hat{p}_i)\right),$ (1.2)

where z_i is the energy (or p_T) and \hat{p}_i the angular information of particle *i*.

EFN:
$$F\left(\sum_{i=1}^{M} z_i \Phi(\hat{p}_i)\right)$$
, PFN: $F\left(\sum_{i=1}^{M} \Phi(p_i)\right)$

Permutation-invariant & Variable lengths

Many common observables are naturally constructed by simple choices of Φ and F. Furthermore, Φ and F can be parametrized by NN layers, capable of learning essentially any function, in order to explore more complicated observables.



Network Implementation





Figure 4. The particular dense networks used here to parametrize (a) the per-particle mapping Φ and (b) the function F, shown for the case of a latent space of dimension $\ell = 8$. For the EFN, the latent observable is $\mathcal{O}_a = \sum_i z_i \Phi_a(y_i, \phi_i)$. For the PFN family, the latent observable is $\mathcal{O}_a = \sum_i \Phi_a(y_i, \phi_i, z_i, \text{PID}_i)$, with different levels of particle-ID (PID) information. The output of F is a softmaxed signal (S) versus background (B) discriminant.

(not necessary) Preprocessing

 $p_{T,i} = \frac{p_{T,i}}{\sum_j p_{T,j}}, y_i = y_i - y_j, \phi_i = \phi_i - \phi_j$



Advantages:

Easily stacked. A deep network can be built with many EdgeConv operations to learn features of point cloud hierarchically.

The graph describing the point clouds are **dynamically updated** to reflect the changes in the edges, i.e., the neighbors of each point.

Introduction to ParticleNet (DGCNN)



FIG. 1. The structure of the EdgeConv block.

FIG. 2. The architectures of the (a) ParticleNet and the (b) ParticleNet-Lite networks.

The "edge features" are constructed from the "features" input using the indices of k nearest neighboring particles. The EdgeConv operation is implemented as a 3-layer MLP.

After the EdgeConv blocks, a channel wise global average pooling operation is applied to aggregate the learned features over all particles in the cloud.

The attention-based cloud net (ABCNet) treats each collider event as an unordered set of points that defines a point cloud. To enhance the extraction of local information, an attention mechanism is used. The main difference between ABCNet and ParticleNet is that ABCNet takes advantage of attention mechanisms to enhance the local feature extraction, allowing for a more compact and efficient architecture. To capture the global information, direct connections for global input features can be directly added.

Key part: GAPLayer (Graph Attention Pooling Layer)

The point cloud is first represented as a graph with vertices represented by the points themselves. The edges are constructed by connecting the points to their k-nearest neighbors, while the edge features, $y_{ij} = (x_i - x_{ij})$, are taken as the difference between features of each point x_i and its k-neighbors x_{ij} . A GAPLayer is constructed by first encoding each point and edge to a higher-level feature space of dimension F using a single-layer neural network (NN), with learnable parameters θ, in the following form:

$$\begin{array}{ll} \mbox{Point Feature} & x_i' = h(x_i, \theta_i, F) \\ \mbox{Edge Feature} & y_{ij}' = h(y_{ij}, \theta_{ij}, F) & y_{ij} = (x_i - x_{ij}) \end{array}$$

$$\label{eq:expectation}$$

$$\begin{array}{ll} \mbox{Attention coefficient} & c_{ij} = {\rm LeakyRelu}(h(x_i', \theta_i', 1) + h(y_{ij}', \theta_{ij}', 1)) \end{array}$$

A single attention feature for each point is
$$\hat{x}_i = ext{Relu}\left(\sum_j c_{ij}y_{ij}'
ight)$$

The outputs of each GAPLayer consist of attention features (\hat{x}_i) and graph features (y'_{ij}) . The graph features are further aggregated in the form:

$$y_{ij}^{max} = max(y_{ij}')$$

A M-head process repeats the same procedure described above M times, differing only on the random weight initialisation. The M results are combined by taking the maximum.



Fig. 1. ABCNet architecture used for quark-gluon tagging. Fully connected layer and encoding node sizes are denoted inside "{}". For each GAPLayer, the number of k-nearest neighbours (k) and heads (H) are given.

Input features:
$$\Delta \eta$$
, $\Delta \phi$, log p_T , log E, log $\frac{p_T}{p_{Tj}}$, log $\frac{E}{E_j}$, ΔR , PID Global features: m_j , p_{Tj}

LorentzNet, a new symmetry-preserving deep learning model for jet tagging. The message passing of LorentzNet relies on an efficient Minkowski dot product attention.

Preprocessing cannot achieve full invariance to arbitrary Lorentz transformations. LorentzNet directly scalarizes the input 4-vectors to realize Lorentz symmetry. Specifically, they design Minkowski dot product attention, which aggregates the 4-vectors with weights learned from Lorentz-invariant geometric quantities under the Minkowski metric.

They regard the constituent particles as a point cloud, which is an unordered, permutation invariant set of particles $V = (v_1, ..., v_N) \in \mathbb{R}^{N \times 4}$

InputVector
$$v_i = (E^i, p_x^i, p_y^i, p_z^i)$$
Scalar $s_i = (s_1^i, s_2^i, \cdots, s_{\alpha}^i)$ Graph $G = (V, E)$ where V is the set of nodes and E is the set of edges.

Let Q be the Lorentz transformation, the Lorentz equivariance of $\phi(\cdot)$ means:

$$Q\phi(v) = \phi(Qv), \quad \text{for } \phi(v) \in \mathbb{R}^4; \tag{2.2}$$

$$\phi(v) = \phi(Qv), \text{ for } \phi(v) \in \mathbb{R}.$$
 (2.3)

The construction of the LorentzNet is based on the following universal approximation theorem for the Lorentz group equivariant continuous function.

Proposition 3.1. [64] A continuous function $\phi : (\mathbb{R}^{N \times 4}) \to \mathbb{R}^4$ is Lorentz-equivariant if and only if

$$\phi(v_1, v_2, \cdots, v_N) = \sum_{i=1}^N g_i(\langle v_i, v_j \rangle_{i,j=1}^N) v_i, \qquad (3.1)$$

where g_i are continuous Lorentz-invariant scalar functions, and $\langle \cdot, \cdot \rangle$ is the Minkowski inner product.



Figure 1. (*left*): the structure of the Lorentz Group Equivariant Block (LGEB). (*right*): the network architecture of the LorentzNet.

The inputs include the PID, mass and the 4-momenta of each particles.

They use $h^l = (h_1^l, h_2^l, ..., h_N^l)$ to denote the node embedding scalars and $x^l = (x_1^l, x_2^l, ..., x_N^l)$ to denote the coordinate embedding vectors in the l-th LGEB layer.

Edge message: $m_{ij}^l = \phi_e\left(h_i^l, h_j^l, \psi(\|x_i^l - x_j^l\|^2), \psi(\langle x_i^l, x_j^l \rangle)\right)$ $\psi(\cdot) = \text{sgn}(\cdot) \log(|\cdot|+1)$

Minkowski dot product attention: $x_i^{l+1} = x_i^l + c \sum_{j \in [N]} \phi_x(m_{ij}^l) \cdot x_j^l$ $h_i^{l+1} = h_i^l + \phi_h \left(h_i^l, \sum_{j \in [N]} w_{ij} m_{ij}^l\right)$ $w_{ij} = \phi_m(m_{ij}^l) \in [0, 1]$

 ϕ_x , ϕ_h , ϕ_e , ϕ_m are all scalar functions modeled by NN (linear layers).

Introduction to Transformer



Introduction to Point Cloud Transformer

Point Cloud Transformer incorporates the advantages of the Transformer architecture to an unordered set of particles resulting from collision events.



To complete the general architecture, the SA layers are combined through a simple concatenation over the feature dimension, followed by a mean aggregation, resulting in the overall means of each feature across all the particles.

Introduction to Point Cloud Transformer



EdgeConv uses a k-nearest neighbors approach to define a vicinity for each point in the point cloud.

> Off-set Attention: $F_{sa} - F_e$ results in a superior classification performance



The output of the feature extractor F_e is used as the input of the first SA layer.

Attention weights are created by matching Q and K through matrix multiplication. These attention weights, after normalization, represent the weighted importance between each pair of particles. The self-attention is then the result of the weighted elements of V, defined as the result of the matrix multiplication between the attention weights and the value matrix.

$$\begin{aligned} \mathbf{Q}, \mathbf{K}, \mathbf{V} &= \mathbf{F}_{e}.(\mathbf{W}_{\mathbf{q}}, \mathbf{W}_{\mathbf{k}}, \mathbf{W}_{\mathbf{v}}) \\ \mathbf{Q}, \mathbf{K} &\in \mathbb{R}^{N \times d_{\mathbf{a}}}, \mathbf{V} \in \mathbb{R}^{N \times d_{\mathrm{out}}} \\ \mathbf{A} &= \mathrm{Softmax}(\mathbf{Q}.\mathbf{K}^{\mathrm{T}})/N, \mathbf{A} \in \mathbb{R}^{N \times N} \\ \mathbf{F}_{sa} &= \mathbf{A} \cdot \mathbf{V}, \mathbf{F}_{sa} \in \mathbb{R}^{N \times d_{\mathrm{out}}} \end{aligned}$$

Particle Transformer: Incorporating pairwise interactions in the attention mechanism



Particles form an array of (N, C), each particle has C features. Interactions form an array of (N, N, C'), each pair of particles has C' features.

 $\begin{array}{ll} \ln \Delta & \Delta = \sqrt{(y_a - y_b)^2 + (\phi_a - \phi_b)^2}, \\ lnk_T & k_{\rm T} = \min(p_{{\rm T},a}, p_{{\rm T},b})\Delta, \\ lnz & z = \min(p_{{\rm T},a}, p_{{\rm T},b})/(p_{{\rm T},a} + p_{{\rm T},b}), \\ m^2 = (E_a + E_b)^2 - \|{\bf p}_a + {\bf p}_b\|^2, \\ lnm^2 \end{array}$

Introduction to Particle Transformer





P-MHA(Q, K, V) = SoftMax $(QK^T/\sqrt{d_k} + \mathbf{U})V$, They take U as the attention mask matrix

 $\begin{array}{ll} \mathsf{Q} \in R^{N \times d_1} & \mathsf{Scale:} \ d_k = d_1 \ \mathsf{embedding \ dimension} \\ \mathsf{K} \in R^{N \times d_1} & Q_i = QW_i^Q, K_i = KW_i^K, V_i = VW_i^V, \quad i = 1, \cdots, 8 \\ & \mathsf{head}_i = \mathsf{Attention}(Q_i, K_i, V_i), \quad i = 1, \cdots, 8 \\ \mathsf{V} \in R^{N \times d_2} & \mathsf{MultiHead}(Q, K, V) = \mathsf{Concat}(\mathsf{head}_1, \cdots, \mathsf{head}_8) W^O \\ \mathsf{U} \in R^{N \times N} & We \ set \ d_1 = d_2 = d_0 / h \end{array}$

Typical evaluation metrics for performance

Accuracy

AUC (Area under ROC curve)

Background rejection $(\frac{1}{\epsilon_h})$ at a certain signal efficiency ϵ_s

Top Tagging	Ac	c AUG	$C = 1/\epsilon_B \ (\epsilon_S = 0.5)$	$1/\epsilon_B \ (\epsilon_S = 0.3)$	
ResNeXt-50 [17]	0.93	36 0.983	37302 ± 5	1147 ± 58	
P-CNN [17]	0.93	30 0.980	201 ± 4	$759{\pm}24$	
PFN [33]	-	0.981	19 247 ± 3	$888{\pm}17$	
ParticleNet-Lite [17]	0.93	37 0.984	325 ± 5	$1262{\pm}49$	
ParticleNet [17]	0.94	40 0.98	$58 \qquad 397{\pm}7$	$1615{\pm}93$	
JEDI-net 21	0.92	0.978	- 36	590.4	
JEDI-net with $\sum O$ [21] 0.93	0.980)7 -	774.6	
SPCT	0.95	28 0.979	$99 201 \pm 9$	725 ± 54	
PCT	0.94	40 0.985	392 ± 7	$1533{\pm}101$	
LorentzNet	0.94	42 0.98	68 498 <u>+</u> 18	2195 <u>+</u> 173	
Church (Quarth					
(-IIIOn/()IIark					
		ATTO	1/ (0.5)	1/ (0.0)	
Discrimination	Acc	AUC	$1/\epsilon_B \ (\epsilon_S = 0.5)$	$1/\epsilon_B \ (\epsilon_S = 0.3)$	
Discrimination ResNeXt-50 [17]	Acc 0.821	AUC 0.9060	$\frac{1/\epsilon_B \ (\epsilon_S = 0.5)}{30.9}$	$\frac{1/\epsilon_B \ (\epsilon_S = 0.3)}{80.8}$	
Discrimination ResNeXt-50 [17] P-CNN [17]	Acc 0.821 0.827	AUC 0.9060 0.9002	$\frac{1/\epsilon_B \ (\epsilon_S = 0.5)}{30.9}$ $\frac{34.7}{34.7}$	$\frac{1/\epsilon_B \ (\epsilon_S = 0.3)}{80.8}$ 91.0	
Discrimination ResNeXt-50 [17] P-CNN [17] PFN [33]	Acc 0.821 0.827	AUC 0.9060 0.9002 0.9005	$1/\epsilon_B \ (\epsilon_S = 0.5)$ 30.9 34.7 34.7 ± 0.4	$\frac{1/\epsilon_B \ (\epsilon_S = 0.3)}{80.8}$ 91.0	
Discrimination ResNeXt-50 [17] P-CNN [17] PFN [33] ParticleNet-Lite [17]	Acc 0.821 0.827 - 0.835	AUC 0.9060 0.9002 0.9005 0.9079	$\frac{1/\epsilon_B \ (\epsilon_S = 0.5)}{30.9} \\ 34.7 \\ 34.7 \pm 0.4 \\ 37.1$	$\frac{1/\epsilon_B \ (\epsilon_S = 0.3)}{80.8}$ 91.0 - 94.5	
Discrimination ResNeXt-50 [17] P-CNN [17] PFN [33] ParticleNet-Lite [17] ParticleNet [17]	Acc 0.821 0.827 - 0.835 0.840	AUC 0.9060 0.9002 0.9005 0.9079 0.9116	$\frac{1/\epsilon_B \ (\epsilon_S = 0.5)}{30.9} \\ 34.7 \\ 34.7 \pm 0.4 \\ 37.1 \\ 39.8 \pm 0.2$	$\frac{1/\epsilon_B \ (\epsilon_S = 0.3)}{80.8}$ 91.0 - 94.5 98.6±1.3	
Discrimination ResNeXt-50 [17] P-CNN [17] PFN [33] ParticleNet-Lite [17] ParticleNet [17] ABCNet [18]	Acc 0.821 0.827 - 0.835 0.840 0.840	AUC 0.9060 0.9002 0.9005 0.9079 0.9116 0.9126	$\frac{1/\epsilon_B \ (\epsilon_S = 0.5)}{30.9} \\ 34.7 \\ 34.7 \pm 0.4 \\ 37.1 \\ 39.8 \pm 0.2 \\ 42.6 \pm 0.4$	$\frac{1/\epsilon_B \ (\epsilon_S = 0.3)}{80.8}$ 91.0 - 94.5 98.6±1.3 118.4±1.5	
Discrimination ResNeXt-50 [17] P-CNN [17] PFN [33] ParticleNet-Lite [17] ParticleNet [17] ABCNet [18] SPCT	Acc 0.821 0.827 - 0.835 0.840 0.840 0.840	AUC 0.9060 0.9002 0.9005 0.9079 0.9116 0.9126 0.8910	$\frac{1/\epsilon_B \ (\epsilon_S = 0.5)}{30.9}$ $\frac{34.7}{34.7 \pm 0.4}$ $\frac{37.1}{39.8 \pm 0.2}$ $\frac{42.6 \pm 0.4}{31.6 \pm 0.3}$	$\frac{1/\epsilon_B \ (\epsilon_S = 0.3)}{80.8}$ 91.0 - 94.5 98.6±1.3 118.4±1.5 93.0±1.2	
Discrimination ResNeXt-50 [17] P-CNN [17] PFN [33] ParticleNet-Lite [17] ParticleNet [17] ABCNet [18] SPCT PCT	Acc 0.821 0.827 - 0.835 0.840 0.840 0.815 0.841	AUC 0.9060 0.9002 0.9005 0.9079 0.9116 0.9126 0.8910 0.9140	$\frac{1/\epsilon_B \ (\epsilon_S = 0.5)}{30.9}$ $\frac{34.7}{34.7 \pm 0.4}$ $\frac{37.1}{39.8 \pm 0.2}$ $\frac{42.6 \pm 0.4}{31.6 \pm 0.3}$ $\mathbf{43.2 \pm 0.7}$	$\frac{1/\epsilon_B \ (\epsilon_S = 0.3)}{80.8}$ 91.0 - 94.5 98.6±1.3 118.4±1.5 93.0±1.2 118.0±2.2	



Taken from 2102.05073

2206.xxxx Minxuan He & Daohan Wang

> Spatial tokens: Particle numbers (Different particles) Local

Self-Attention

Channel tokens: Feature numbers (Different features) Global



(a) Spatial Window Multihead Self-attention



With spatial tokens, the spatial dimension defines the token scope, and the channel dimension defines the token feature dimension.

With channel tokens, the channel dimension defines the token scope, and the spatial dimension defines the token feature dimension.





$$\left\{ E, p_x, p_y, p_z, p_T, \sum p_{Tf}, \sum E_f, \overline{\Delta \eta}, \overline{\Delta \phi}, \overline{\Delta R}, PID \right\}.$$
(3.3)

Ι	Е	p_x	p_y	p_z	p_T	$\sum p_{Tf}$	$\sum E_f$	$\overline{\Delta\eta}$	$\overline{\Delta \phi}$	$\overline{\Delta R}$	PID
Е	1	$\frac{p_x}{E}$	$\frac{p_y}{E}$	$\frac{p_z}{E}$	$\frac{p_T}{E}$	0	1	0	0	0	$\frac{E_{PID}}{E}$
p_x	$\frac{p_x}{E}$	1	0	0	$rac{p_x}{p_T}$	0	0	0	0	0	$rac{p_{xPID}}{p_{x}}$
p_y	$\frac{p_y}{E}$	0	1	0	$rac{p_y}{p_T}$	0	0	0	0	0	$rac{p_{yPID}}{p_y}$
p_z	$\frac{p_z}{E}$	0	0	1	0	0	0	0	0	0	$rac{p_{zPID}}{p_{z}}$
p_T	$\frac{p_T}{E}$	$rac{p_x}{p_T}$	$rac{p_y}{p_T}$	0	1	1	0	0	0	0	$rac{p_{TPID}}{p_T}$
$\sum p_{Tf}$	0	0	0	0	1	1	0	0	0	0	p_{TfPID}
$\sum E_f$	1	0	0	0	0	0	1	0	0	0	E_{fPID}
$\overline{\Delta\eta}$	0	0	0	0	0	0	0	1	0	$rac{\overline{\Delta\eta}}{\overline{\Delta R}}$	$\overline{\Delta \eta}_{PID}$
$\overline{\Delta\phi}$	0	0	0	0	0	0	0	0	1	$rac{\overline{\Delta \phi}}{\overline{\Delta R}}$	$\overline{\Delta \phi}_{PID}$
$\overline{\Delta R}$	0	0	0	0	0	0	0	$\frac{\overline{\Delta \eta}}{\overline{\Delta R}}$	$rac{\overline{\Delta \phi}}{\overline{\Delta R}}$	1	$\overline{\Delta R}_{PID}$
PID	$\frac{E_{PID}}{E}$	$rac{p_{xPID}}{p_x}$	$rac{p_{yPID}}{p_y}$	$\frac{p_{zPID}}{p_z}$	$rac{p_{TPID}}{p_T}$	p_{TfPID}	E_{fPID}	$\overline{\Delta\eta}_{PID}$	$\overline{\Delta \phi}_{PID}$	$\overline{\Delta R}_{PID}$	1

Table 1. The jet feature pairwise interaction matrix used as the inputs for the P-DAT.

Constraints in Pairwise Interaction Matrix

1, Limited by computational time and memory consumption, it is difficult to use the full pairwise interaction matrix.

2, We solved the memory usage problem by importing and deleting data during training.

3, The computational time problem regarding of using the full pairwise interaction matrix is still unresolved. We need novel approaches for particle embeddings and self-attentions in order to fully utilize the physics principles.

- Point Cloud based representation and Transformer show excellent performance. We can explore more novel methods for input embedding and self attentions that improve the classification performance and alleviate the computational cost.
- 2, Point Cloud based representation can be used for other interesting physics problems such as pileup subtraction, jet grooming and jet energy calibration.
- 3, Based on the Point Cloud representation and EFN/PFN, we can explore different variant architectures and try to obtain a natural visualization of the learned latent space, providing insights as to what exactly the NN is learning.
- 4, We can try to incorporate more priors or constraints from physics principles in the architecture designs.
- 5, How to evaluate the statistical uncertainties of deep learning in jet tagging?

Thanks for your attention !