

Streamlined jet tagging network assisted by jet prong structure

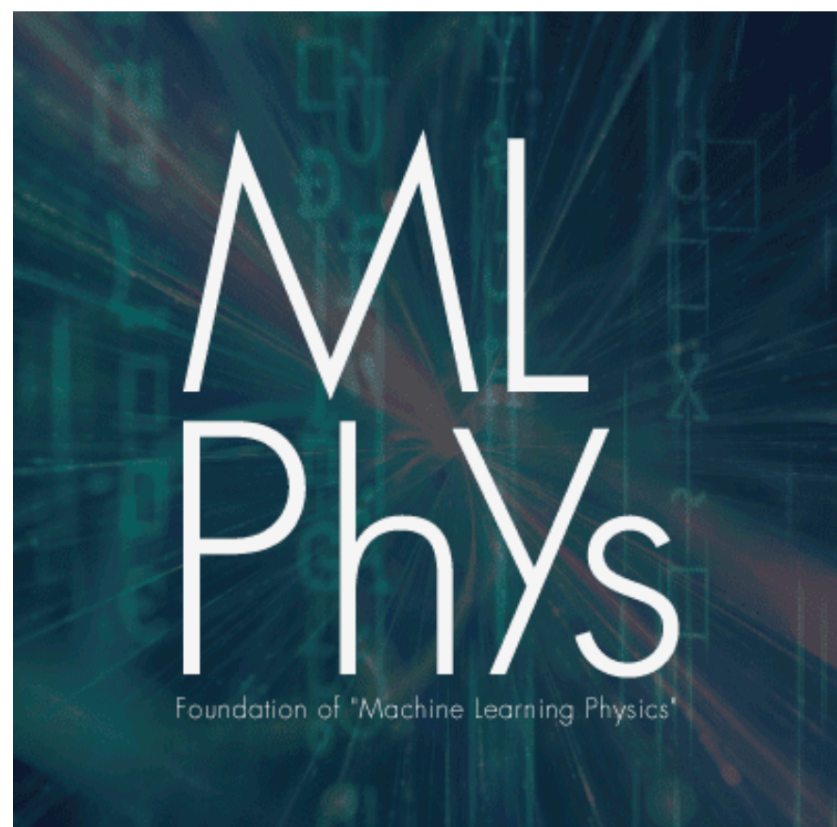
Speaker: Ahmed Hammad with Mihoko Nojiri

Theory center, KEK, Japan

LCWS workshop 2024 @ Tokyo

9th July 2024

MLPhYs 学術変革領域研究(A) 学習物理学の創成
Foundation of "Machine Learning Physics"



Introduction

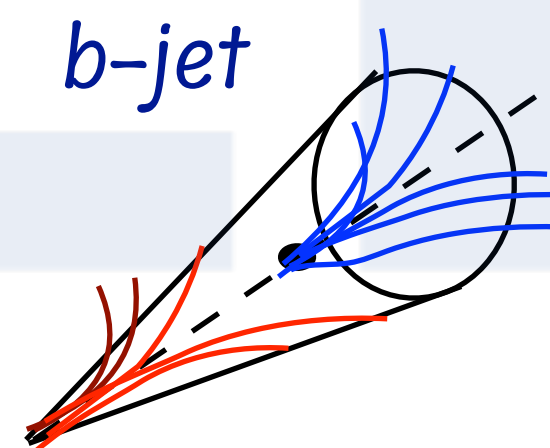
Jet tagging:

Identify the hard scattering particle that initiates the jet.

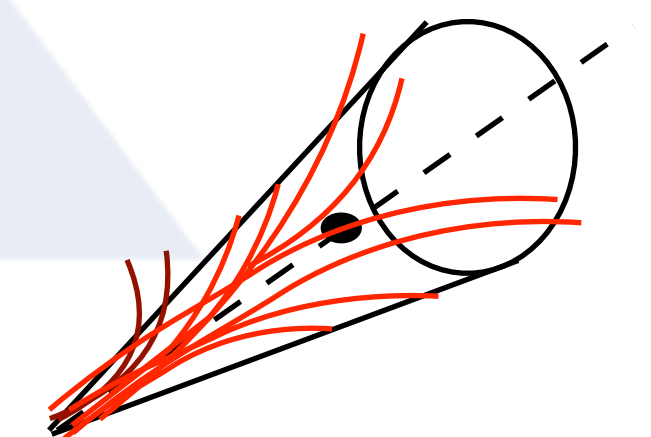
*Jet identification is important to improve our understanding for the QCD processes.
Moreover, it improves the significance for new physics searches.*

Examples:

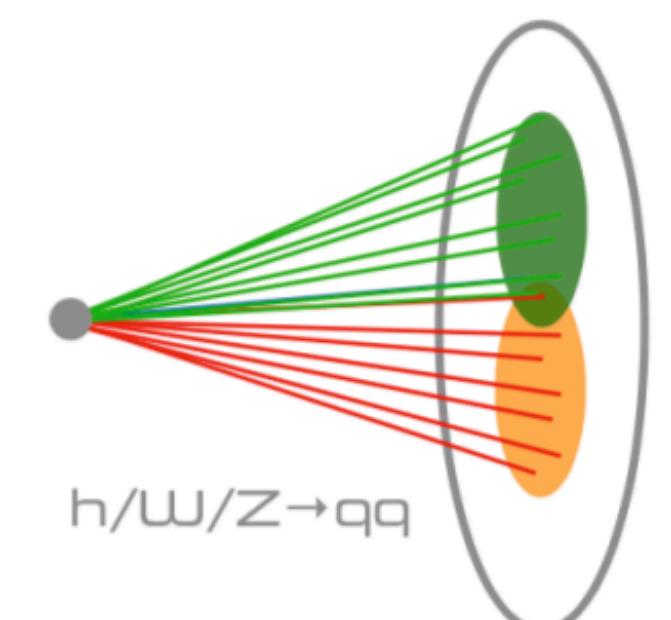
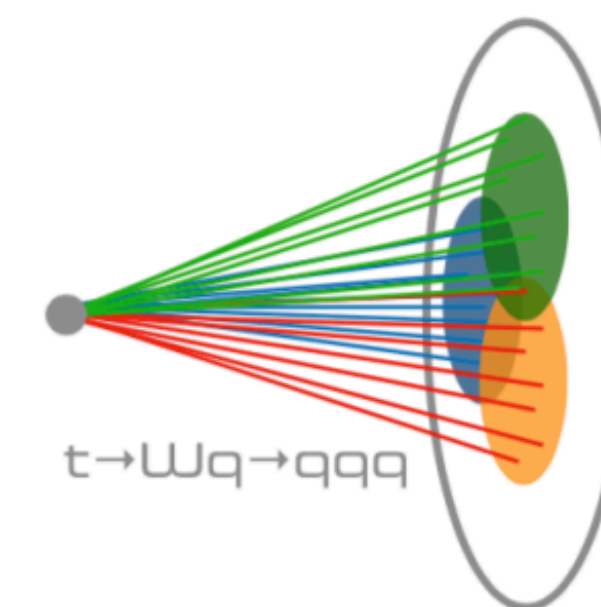
- Heavy flavor tagging, **bottom and charm tagging**



Light flavor-jet



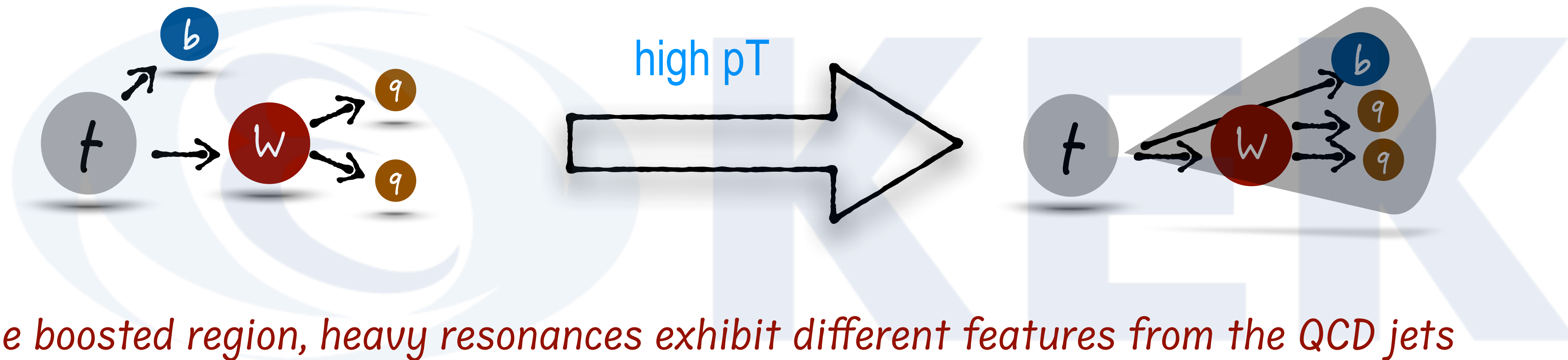
- Heavy resonance tagging, **Top-jet, W-jet, Higgs-jet**



Introduction

○ Boosted Jet tagging:

At high p_T , the decay products from heavy particles, **Higgs, W, Z, top**, become collimated and can be contained in a single **large-R jet**



○ Machine learning

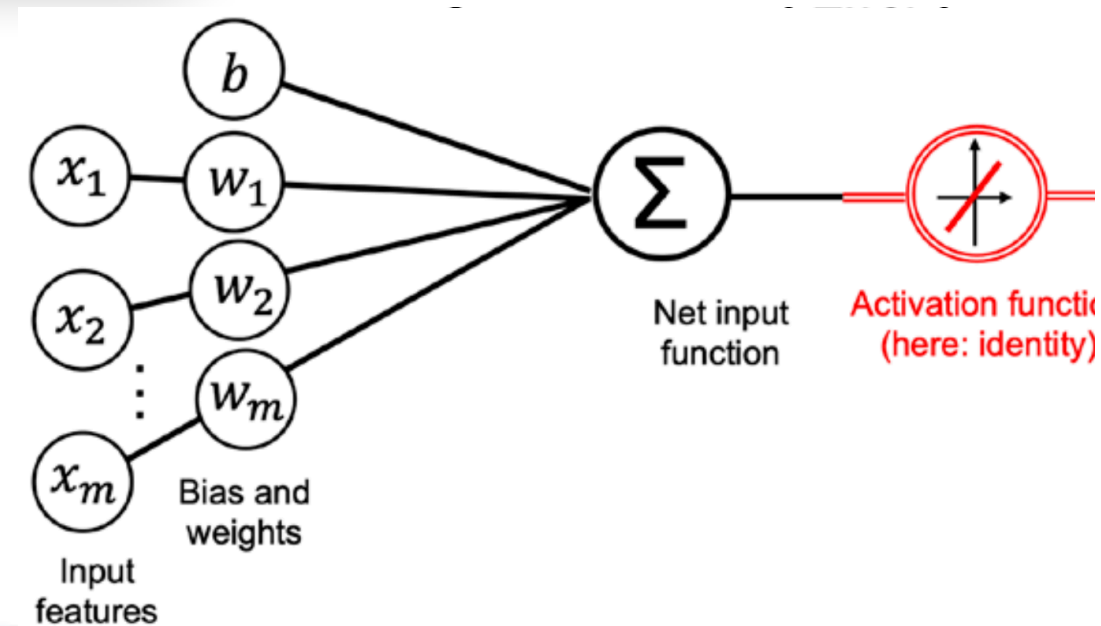
Existence of many new approaches has been proposed in the past few years, leading to significant improvement in performance and deeper insights into jet physics.

What is the difference between the different ML models? Which is the best one?

Introduction

Multi-Layers perceptron:

$$Z_i^l = \sigma \left(\sum_{j=1}^l W_{i,j}^l \cdot x_j^{l-1} + b_i^l \right)$$



The sum runs over each neuron in the layer fixed pattern of the latent neurons

MLP is **not** invariant under translation, rotation and permutation.
Suffer from the input sparsity.

Convolution NeuralNetwork

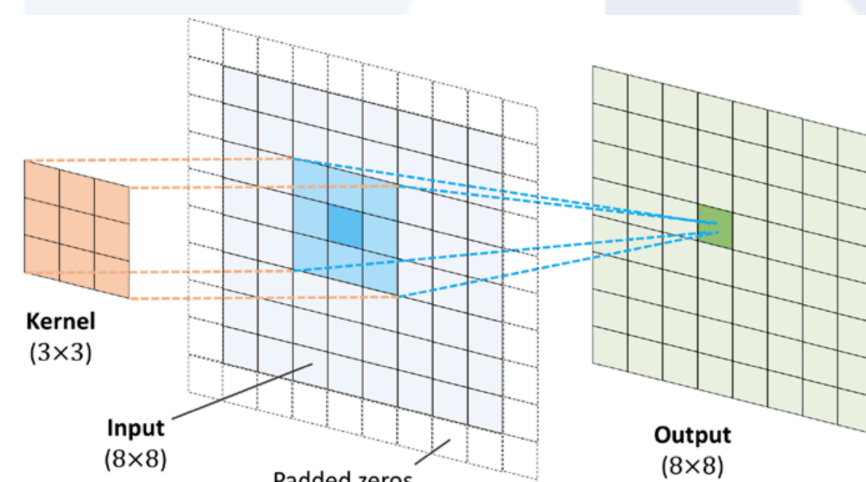
$$Z_{i,j}^l = \sigma \left(\sum_m \sum_n X(i+m, j+n) \times W(m,n) \right)$$

i, j run over spatial dimensions of the image. Local weights captured by each kernel is shared to allow for Translation invariance

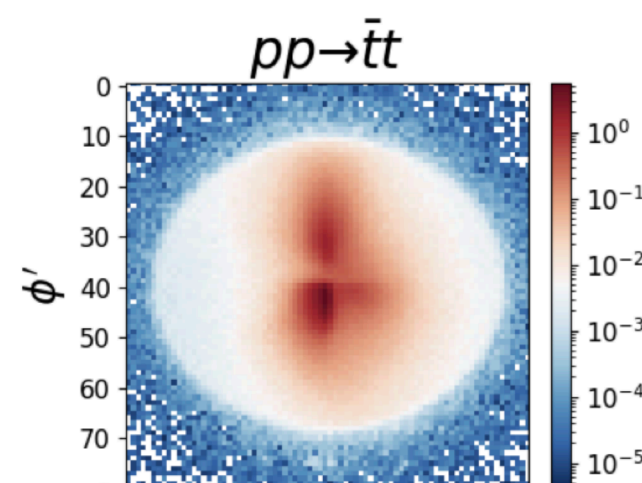
Only the important information are kept via pooling

$$P_{i,j}^l = \text{Max}_{(m,n)} \left[Z_{(i \times M + m, j \times M + n)}^l \right]$$

CNN is invariant **only** under translation. Suffer from the input sparsity.



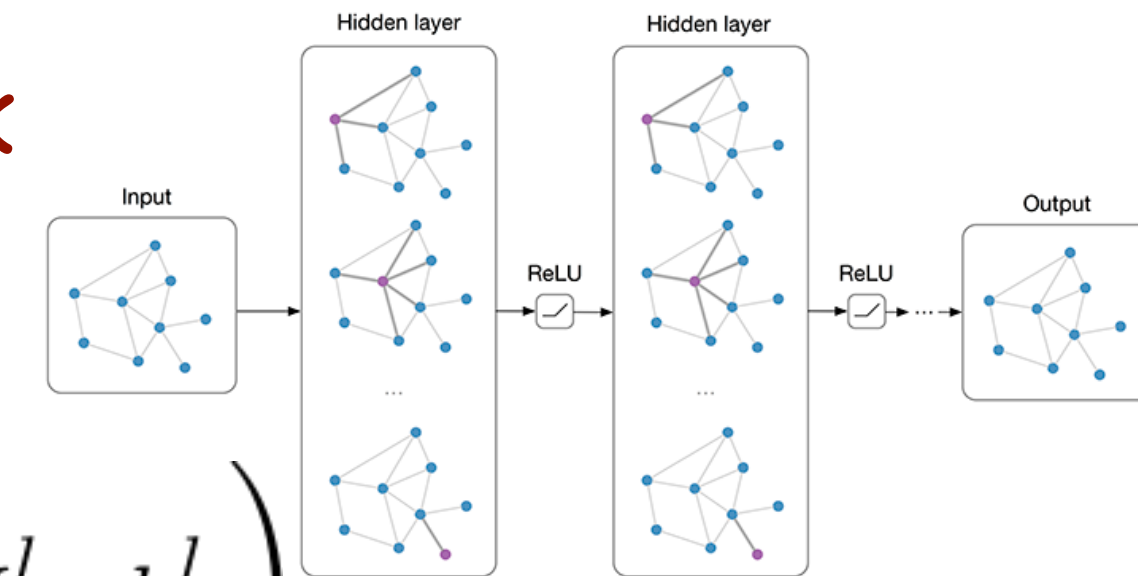
Low computation cost



Graph NeuralNetwork

High computation cost

$$h_v^{l+1} = \sigma \left(W^{(l)} h^l + \sum_n W_n^l \cdot h_n^l \right)$$



For a given graph the nodes for the next layer are updated by aggregating the information from the near by nodes from the previous layer

The node update runs sequentially over the graph nodes and thus GNN is **not** invariant under permutation but no sparsity issue.

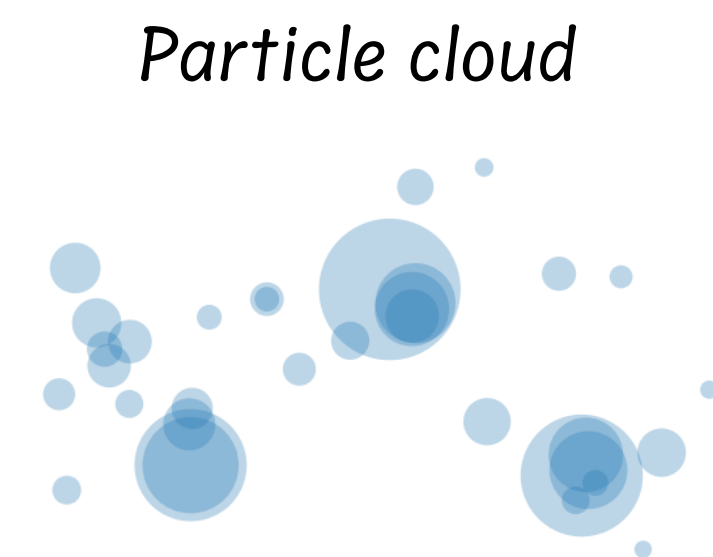
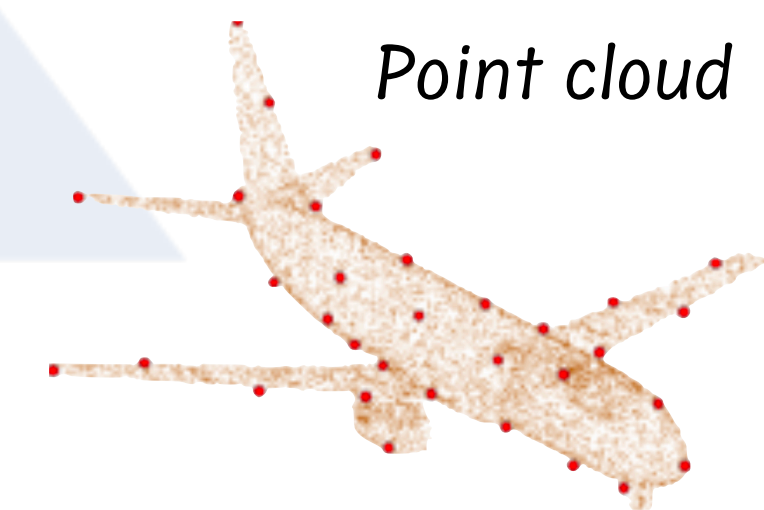
Transformer Network

High computation cost

$$\text{Attention}_{(i,j)} = \text{softmax} \left(\frac{Q \cdot K^T}{\sqrt{d}} \cdot V \right)$$

Query, key and value matrices mix up all the particle and feature tokens

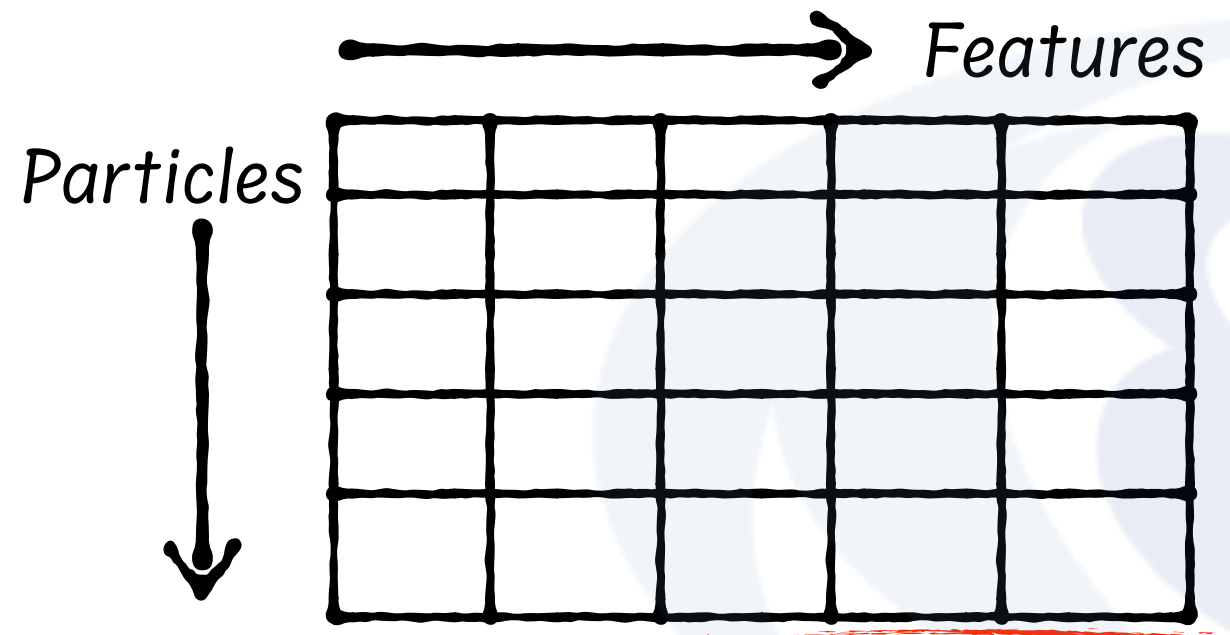
Transformer networks are **permutation invariant with no sparsity issue.**



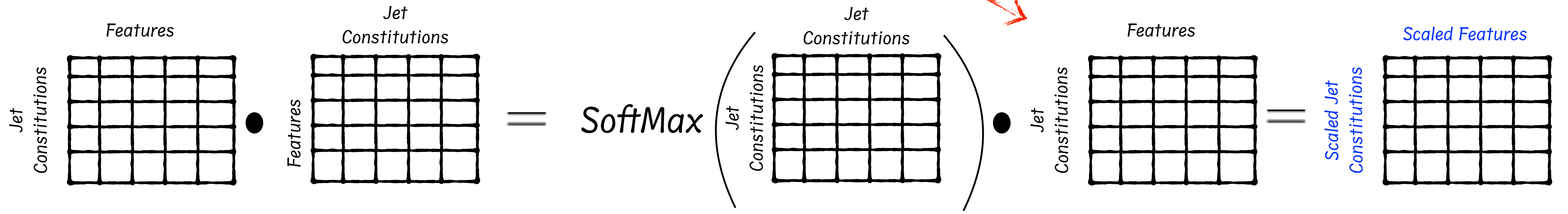
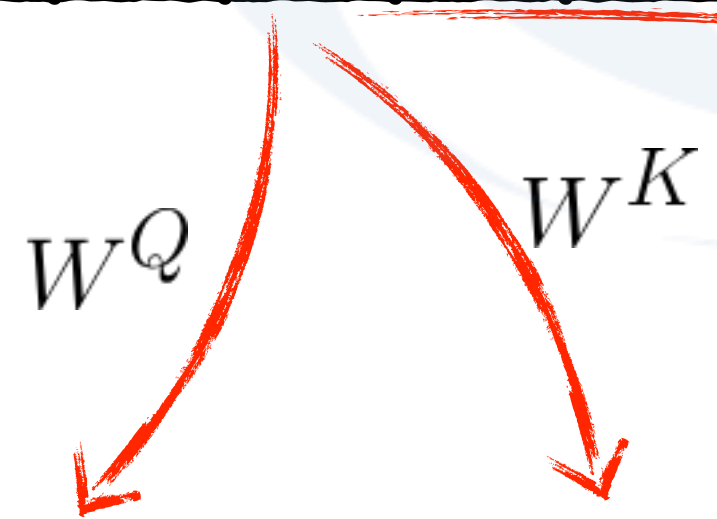
Transformer network

Transformer network mixes particle and feature tokens to highlight the most important tokens for the model decision making. It allows for global and local information extraction.

Input data structured as a fixed size unordered grid



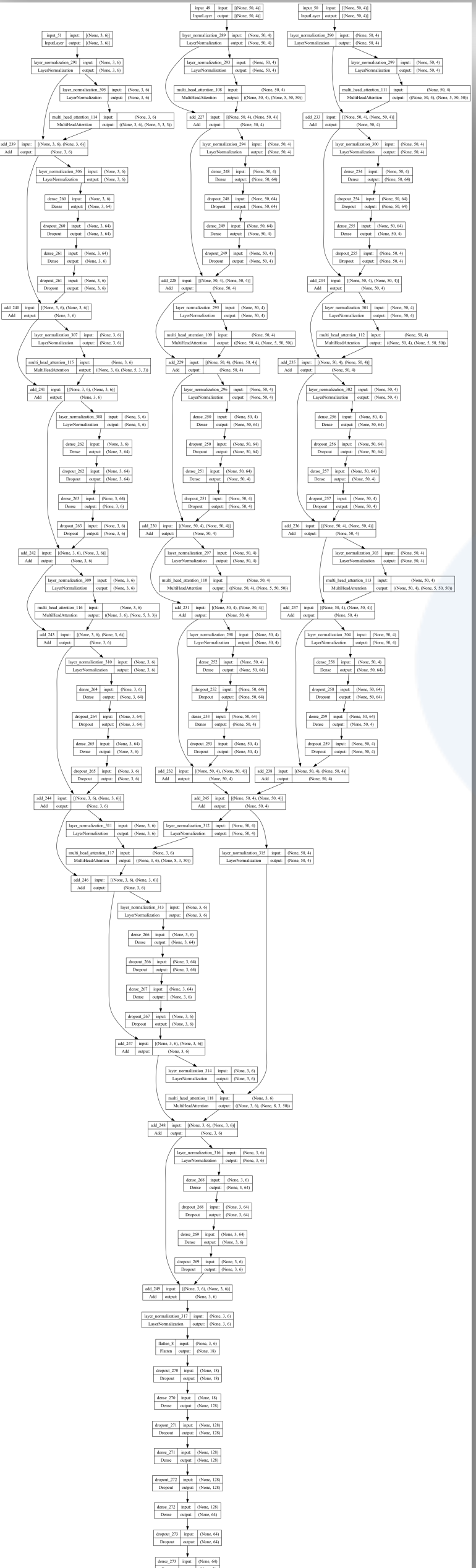
$$\text{Attention}_{(i,j)} = \text{softmax} \left(\frac{Q \cdot K^T}{\sqrt{d}} \cdot V \right)$$



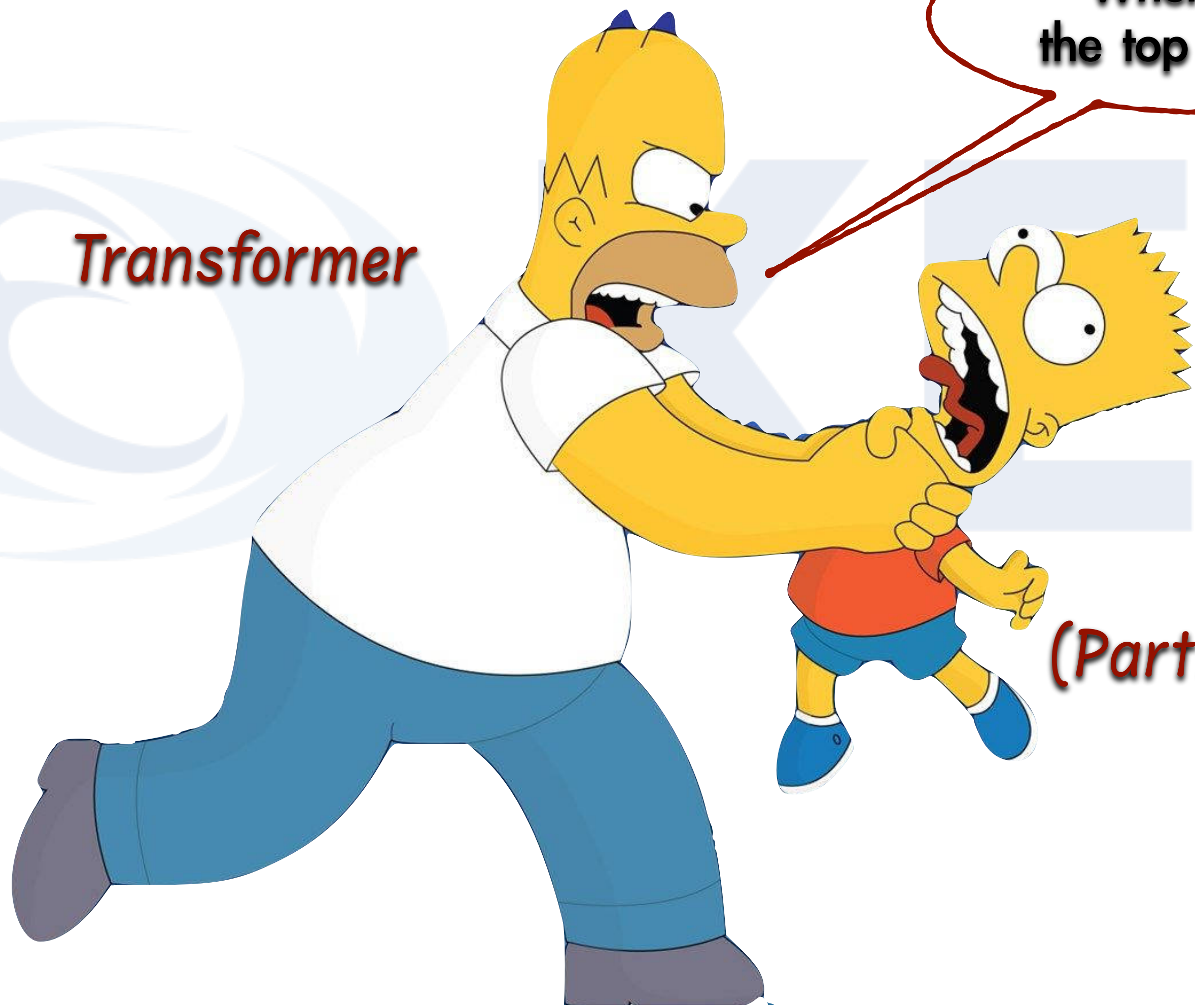
Transformer network

Transformers have the best performance but with high computational complexity.

A.H, S.M, Nojiri
JHEP 03 (2024),2401.00452



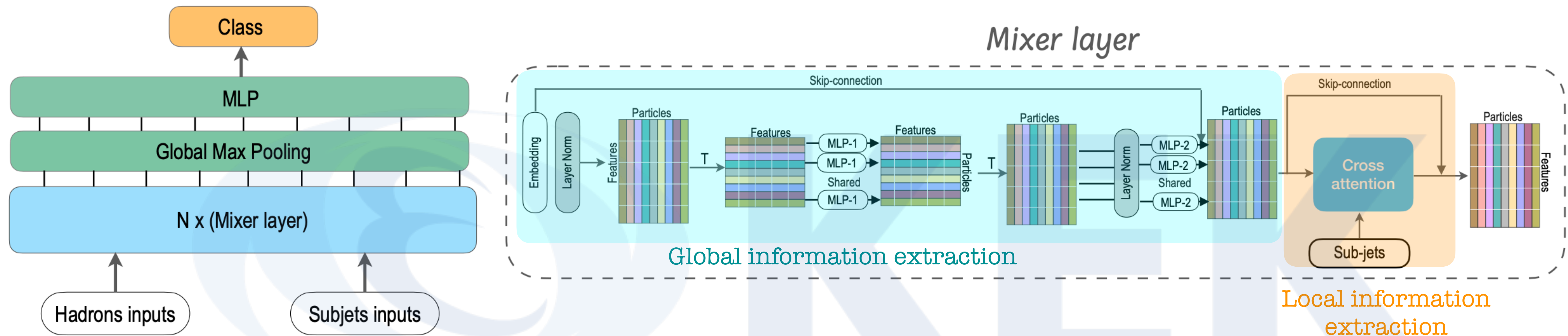
Transformer



Data
(Particle cloud)

Mixer network

Mixer layer has two MLP that mix both features and Particle tokens (similar to the transformer) which allow for fast extraction of the global features of the event. Local information is extracted from the subjects via Cross-attention layer.



Jet tagging task can be divided into two main parts:

- Global information extraction

- Local information extraction

The network learns how important each jet constituent to all other constituents via two MLPs.

The network learns how important each jet constituent to the sub-jet it belongs to

$$Y_{i,j} = X_{i,j} + \left[\left(W_2 \sigma W_1 (\text{LayerNorm}(\mathbf{X})^T) \right)^T \right]_{i,j},$$

$$\tilde{X}_{i,j} = Y_{i,j} + (W_4 \sigma W_3 (\text{LayerNorm}(Y_{i,j})))$$

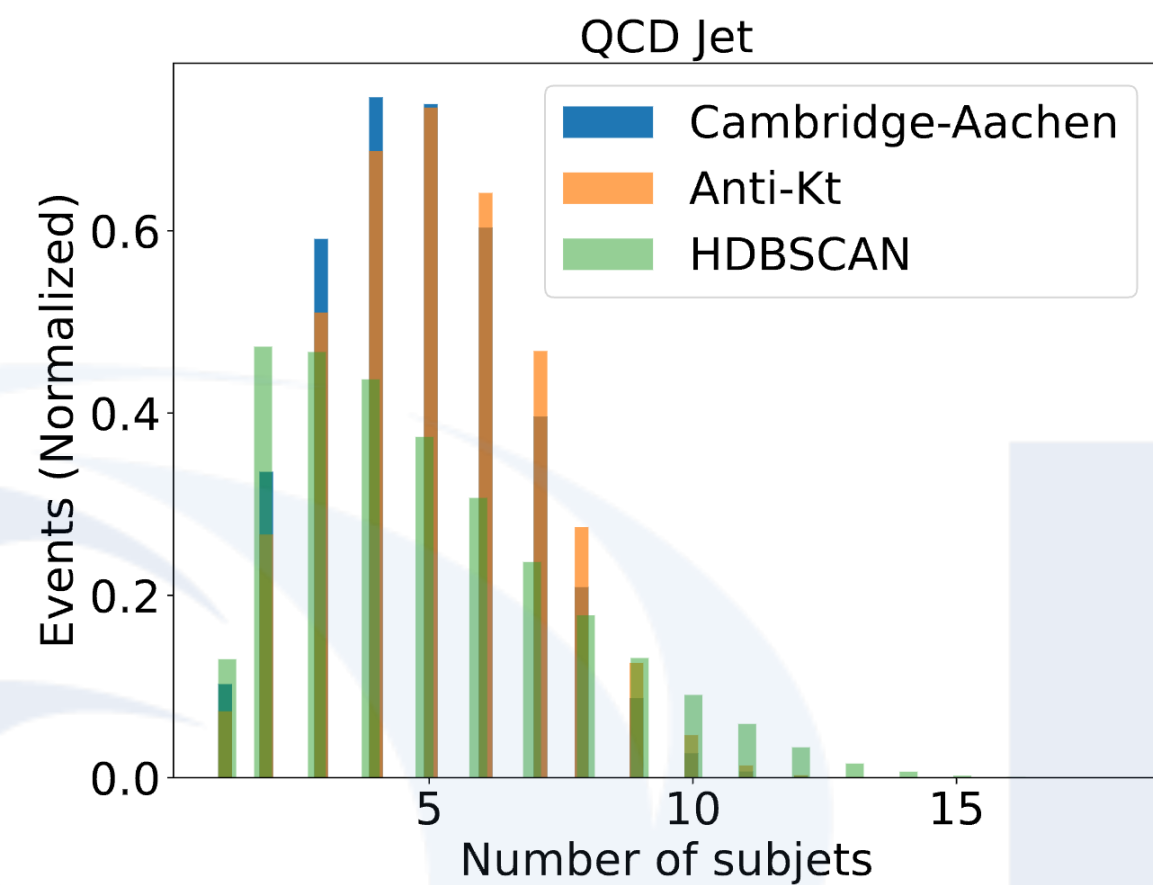
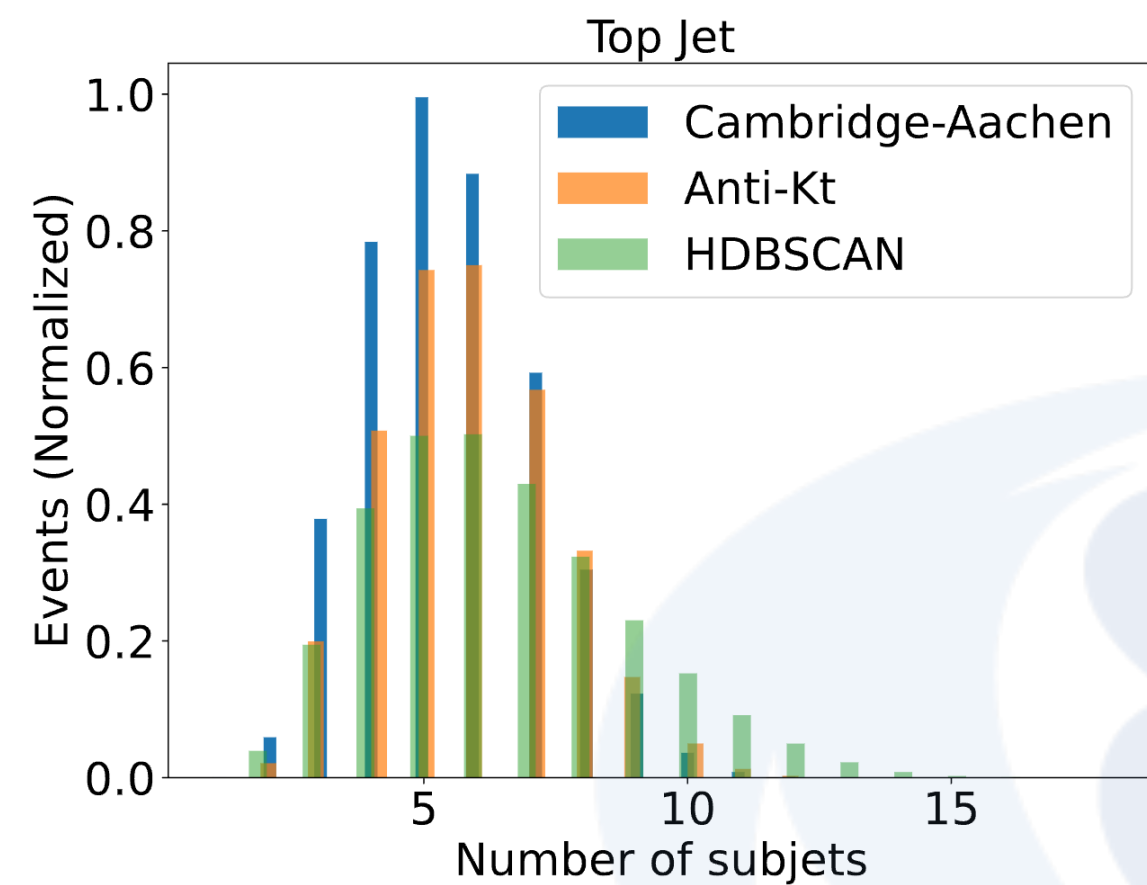
$$P(\text{Hadrons inside jet} \mid \text{subject cluster}) = P(x_i \mid y_\alpha)$$

Validation on Top dataset

The network is validated on TopLandscape community dataset

<https://zenodo.org/records/2603256>

$R = 0.3$ for CA and anti-KT



For subjects clustering, three methods are used:

- Cambridge-Aachen
- Anti-Kt
- Hierarchical Density-Based Spatial Clustering of Applications with Noise (HDBSCAN)

IRC safe, Non parametric clustering with dynamic R

[A.H. M.Nojiri](#)
[arXiv:2404.14677, JHEP 06 \(2024\) 176](#)

TABLE I. performance of the Mixer network for top quark tagging compared with other models. Results for PFN [36], ParticleNET [37], and ParT [38] are quoted from their published results. Transformer(subjet) model is trained from scratch using the CA subjects dataset only. Training time is per epoch with a batch size of 1024. The GPU training time is measured on an NVIDIA RTX A6000 card.

	AUC	Rej _{50%}	Parameters	Time (GPU) [s]
PFN	0.9819	247	86.1K	30
ParticleNET	0.9858	397	370K	729
ParT	0.9858	413	2.14M	612
Transformer(subjets)	0.9640	186	398k	129
Mixer(Anti-kt)	0.9854	375	86.03k	33
Mixer(CA)	0.9856	392	86.03k	33
Mixer(HDBSCAN)	0.9859	416	86.03k	33

Mixer network can achieve high performance as the transformer network with *low computational cost*

Do the two components of the mixer layer capture different information?

Central kernel alignment

Analysis of the hidden Layers representation

CKA computes the the similarity of the hidden layers representations independent on the size of each layer.

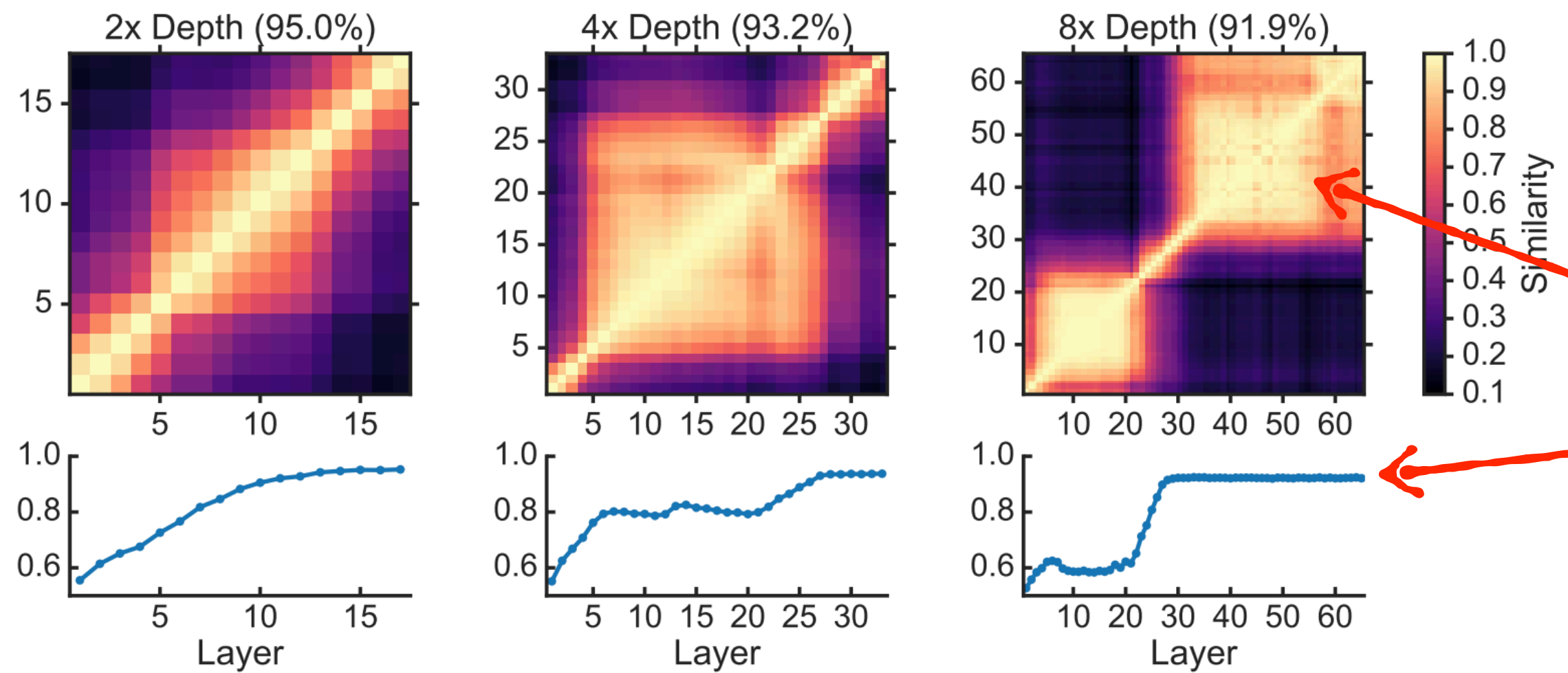
Two inputs from two hidden layers with the form $X \in \mathbb{R}^{d \times P_1}$ $Y \in \mathbb{R}^{d \times P_2}$

Gram matrix can be constructed which is independent on the dimension of the hidden layer $M = XX^T$ and $N = YY^T$

$$CKA(M, N) = \frac{HSIC(M, N)}{\sqrt{HSIC(M, M)HSIC(N, N)}}$$

With $HSIC(M, N) = \frac{1}{(d-1)^2} tr(MHNH)$

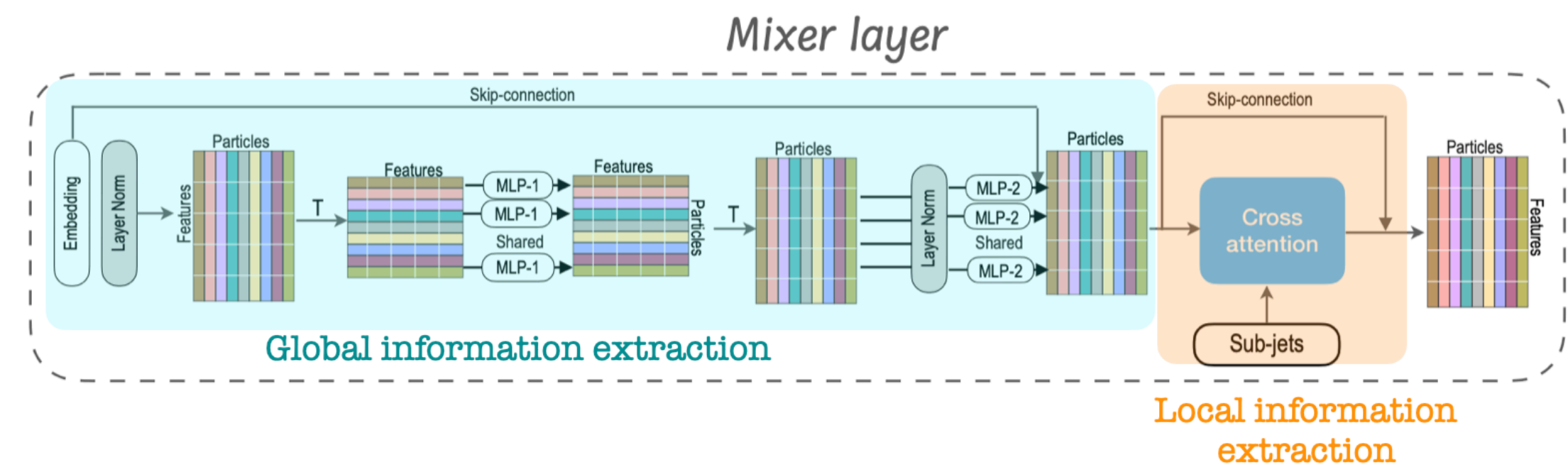
[arXiv:1905.00414](https://arxiv.org/abs/1905.00414)



Layers with similar CKA value, learns the similar information

Central kernel alignment

CKA similarity for 5000 test events
for the Top and QCD jets

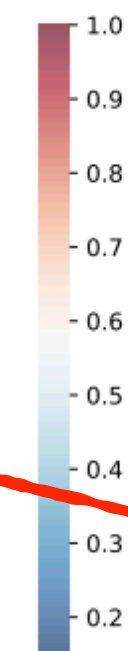
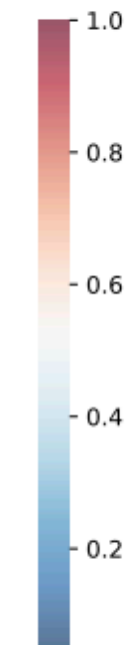


CKA-Top jets

Embedding	1	0.66	0.56	0.31	0.4	0.53	0.17	0.12
FC ¹ (MLP ₁)	0.66	1	0.74	0.29	0.48	0.57	0.2	0.16
FC ² (MLP ₁)	0.56	0.74	1	0.53	0.58	0.46	0.28	0.15
FC ¹ (MLP ₂)	0.31	0.29	0.53	1	0.39	0.17	0.12	0.047
FC ² (MLP ₂)	0.4	0.48	0.58	0.39	1	0.56	0.3	0.07
SkipLayer	0.53	0.57	0.46	0.17	0.56	1	0.17	0.071
Attention	0.17	0.2	0.28	0.12	0.3	0.17	1	0.18
FC	0.12	0.16	0.15	0.047	0.07	0.071	0.18	1

CKA-QCD jets

Embedding	1	0.61	0.61	0.48	0.49	0.46	0.19	0.39
FC ¹ (MLP ₁)	0.61	1	0.8	0.65	0.64	0.28	0.22	0.55
FC ² (MLP ₁)	0.61	0.8	1	0.76	0.62	0.61	0.25	0.68
FC ¹ (MLP ₂)	0.48	0.65	0.76	1	0.57	0.43	0.15	0.48
FC ² (MLP ₂)	0.49	0.64	0.62	0.57	1	0.34	0.57	0.72
SkipLayer	0.46	0.28	0.61	0.43	0.34	1	0.21	0.34
Attention	0.19	0.22	0.25	0.15	0.57	0.21	1	0.58
FC	0.39	0.55	0.68	0.48	0.72	0.34	0.58	1

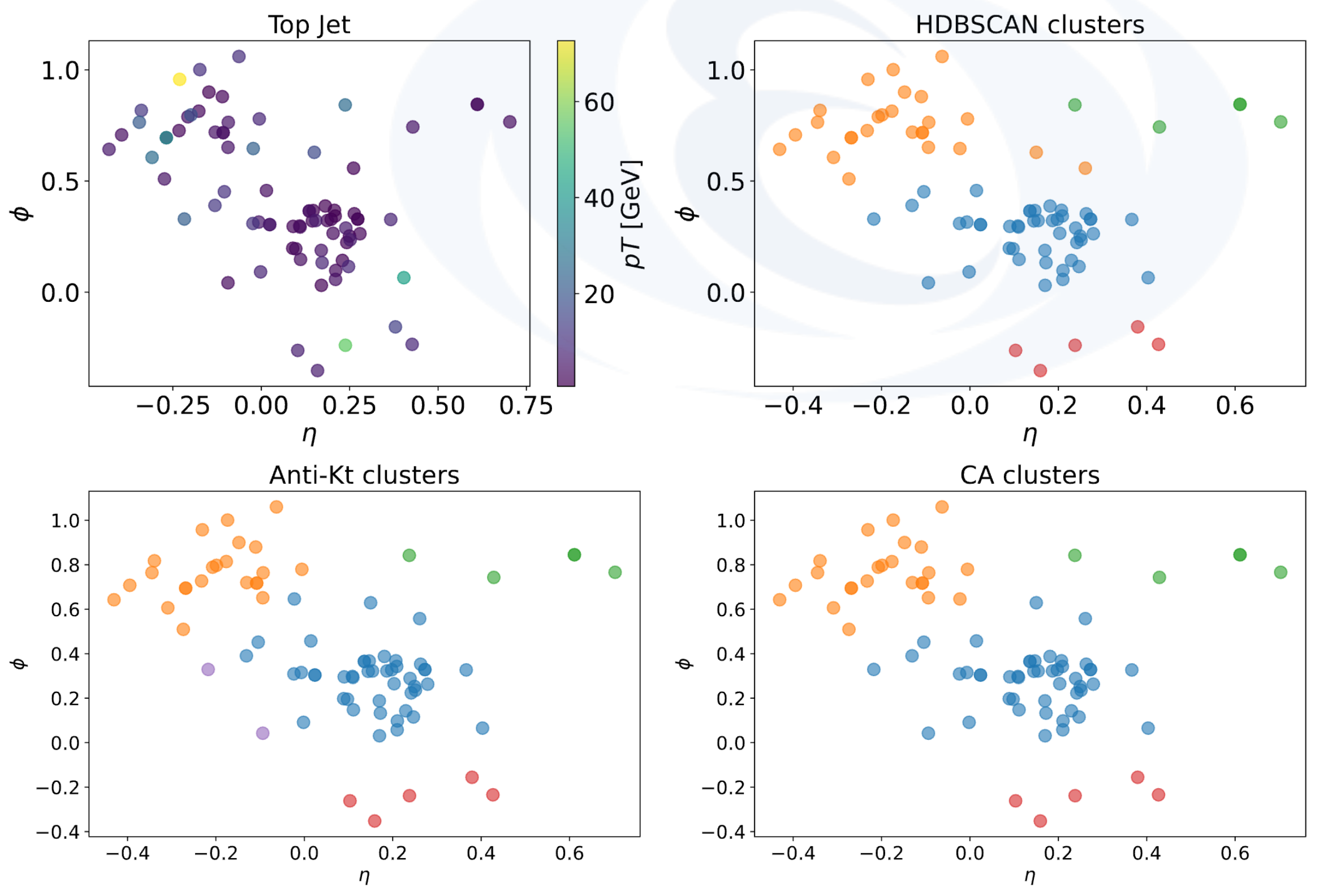


For both data, top and QCD, the cross-attention layer has different CKA value from the two MLPs

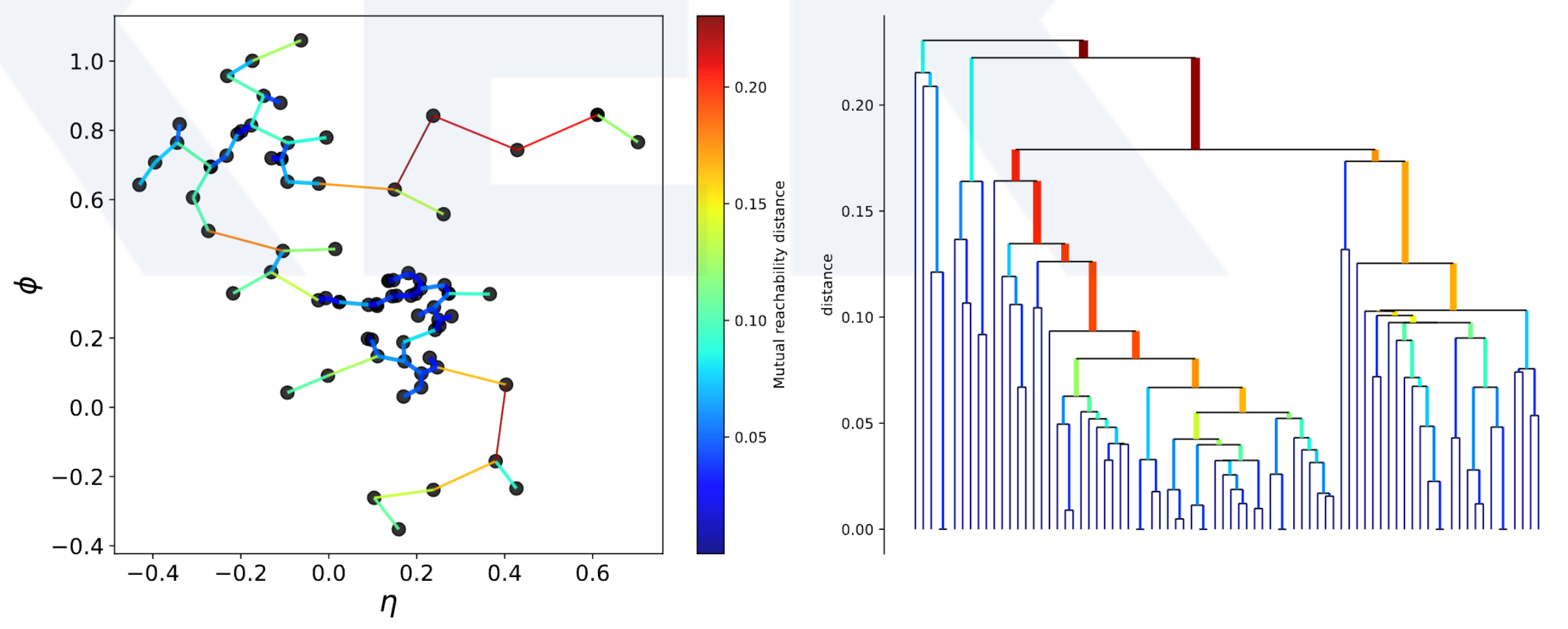
For the QCD, there is no clear prong structure. This is why attention layer capture global information similar to the MLPs

*Thank you
for your attention*

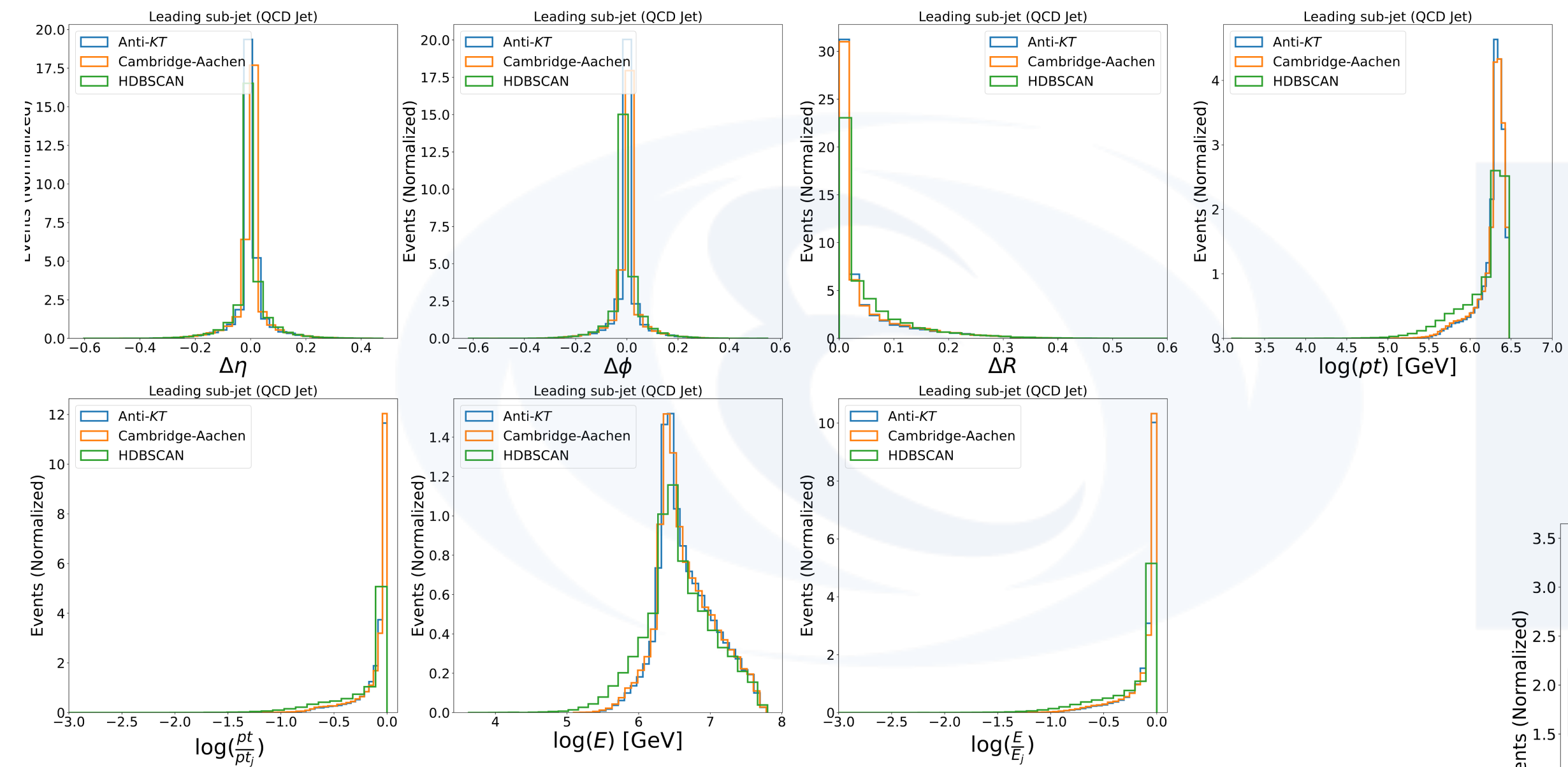
Example of top jet clustering by different methods



Minimum spanning tree by the HDBSCAN



Subjets kinematics for the QCD events



Subjets kinematics for the top events

