# DESIGNING TRANSFORMER FOR PARTICLE PHYSICS MIHOKO NOJIRI (KEK)

with Waleed Esmail(Munster), Ahmed Hammad(KEK)

with Amon Furuihchi, Sung Hak Lim(IBS)

#### MLPHYSICS GRANT IN JAPAN"MACHINE LEARNING PHYSICS "

Overview Organization Events Acheivements Outreach

O y e r y i e w



The research area "Machine Learning Physics" will begin with the aim of discovering new laws and pioneering new materials

B01 Math and application of DL

B02 Statistical data and ML

B03 Topology and Geometry of ML

A01 Lattice

A02 Mihoko Nojiri HEP

Junichi Tanaka (ICEPP Tokyo, ATLAS)

Masako lawasaki (Osaka Metropolitan Belle II )

Noriko Takemura and Hajime Nagahara (Data Science)

A03 Condensed Matter

A04 Quantum and Gravity

PD. Ahmed Hammad

2017-2020: Ph.D Basel University,

**Basel Switzerland** 

2020-2023: SeoulTech, Korea

2023- KEK

# Deep Learning is changing particle physics

# Flavor Tagging Performance

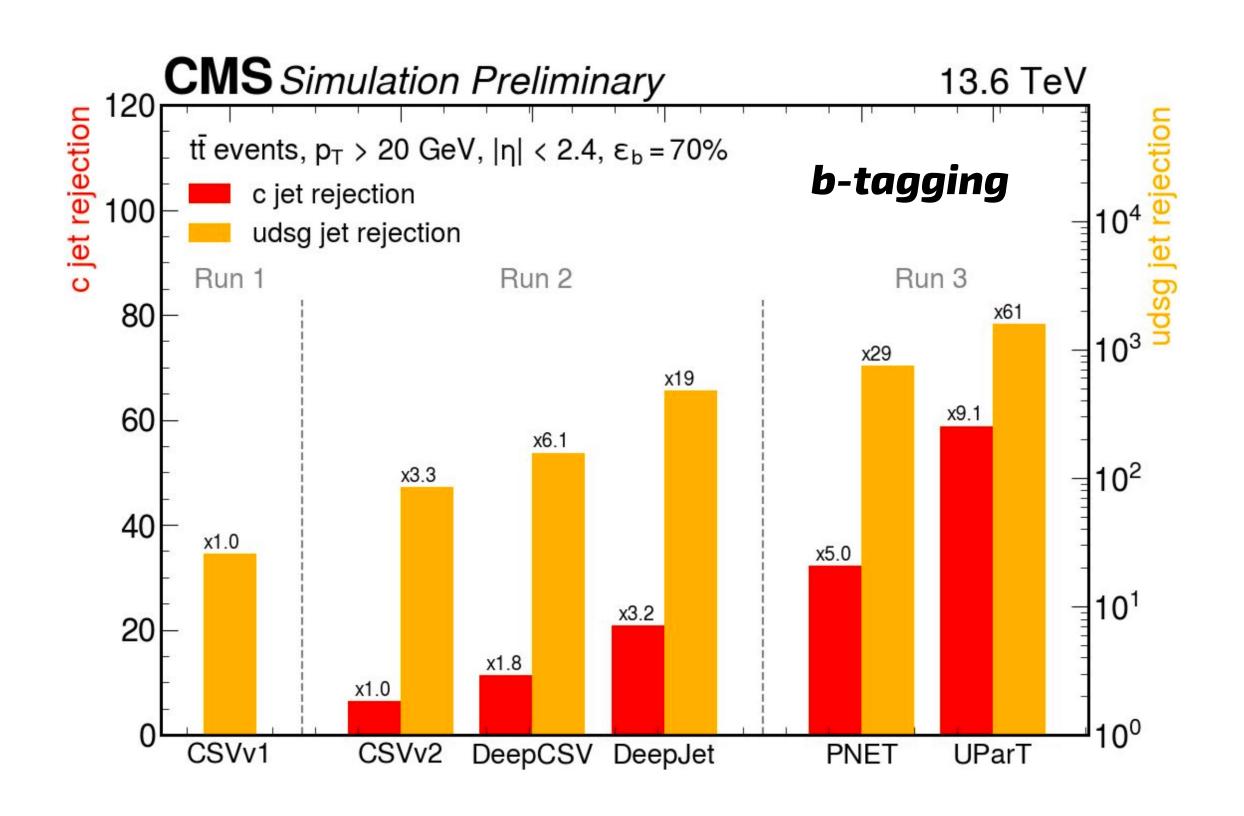


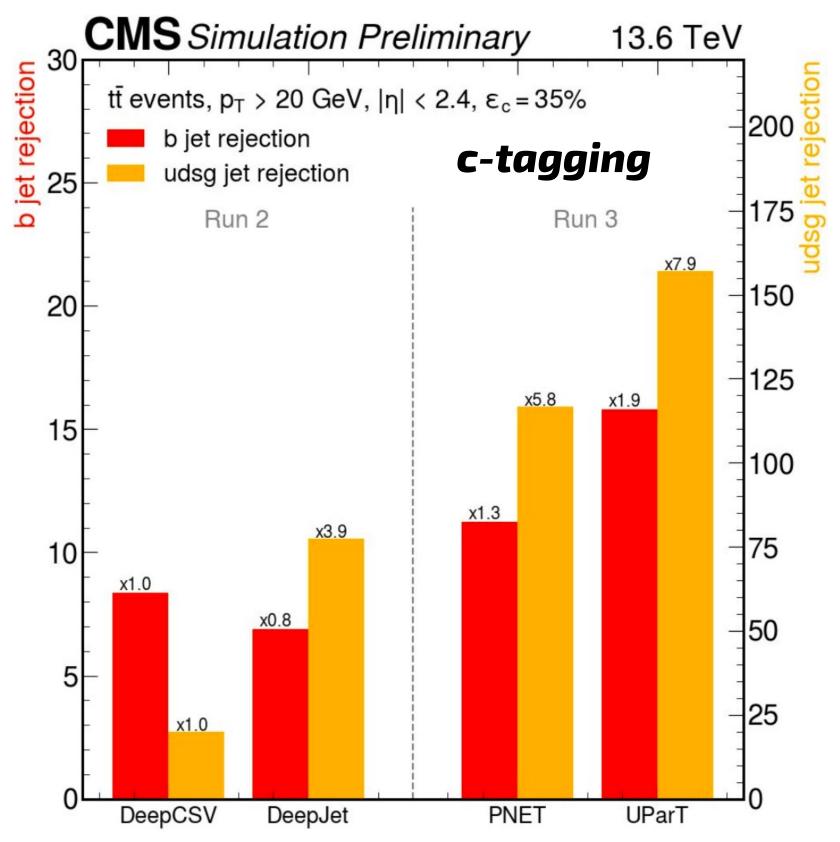




#### b/c-tagging performance

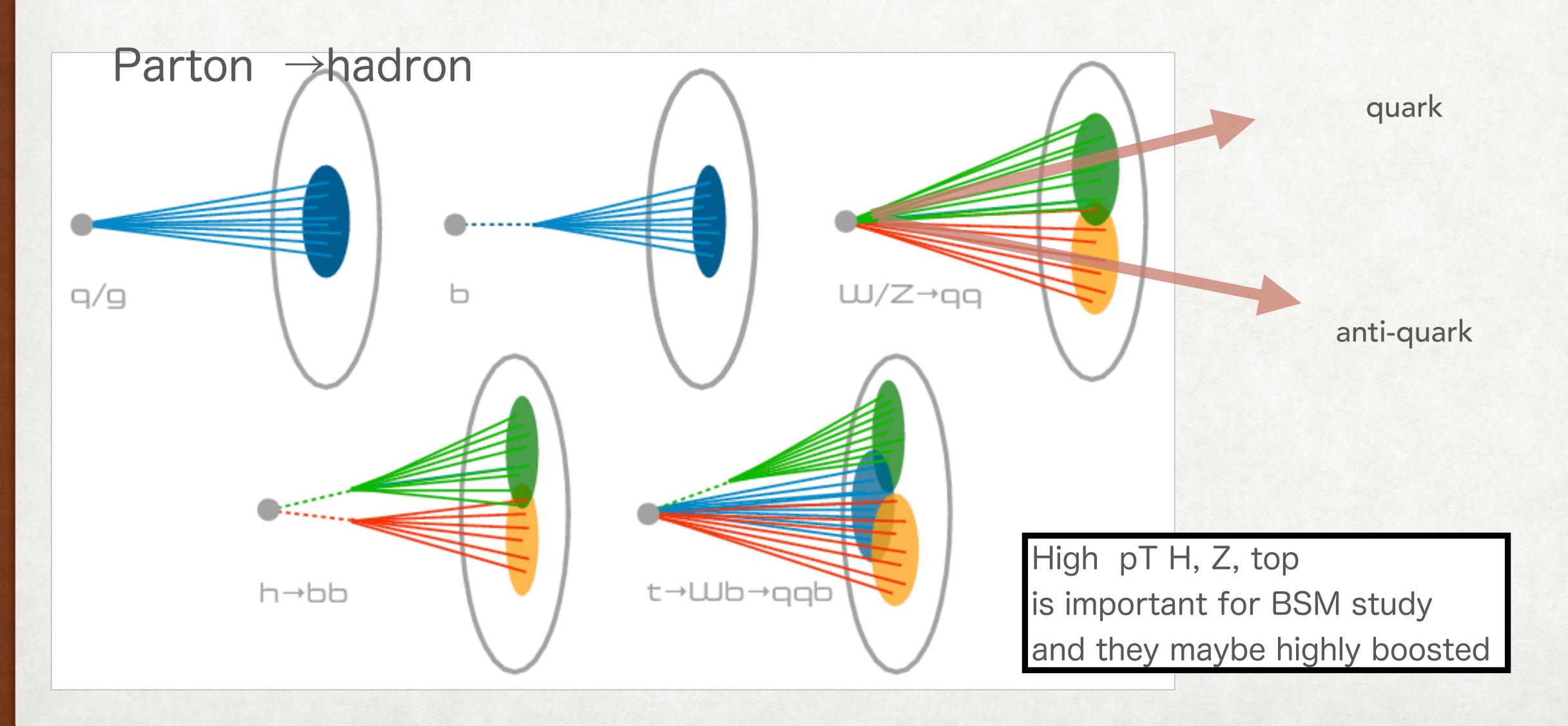
- Promising performance compared to previous taggers
  - ×3 better light jet rejection (at b-jet eff 70%) than DeepJet
  - ×2 better light rejection + ×2 better b-jet rejection (at c-jet eff 35 %) than DeepJet





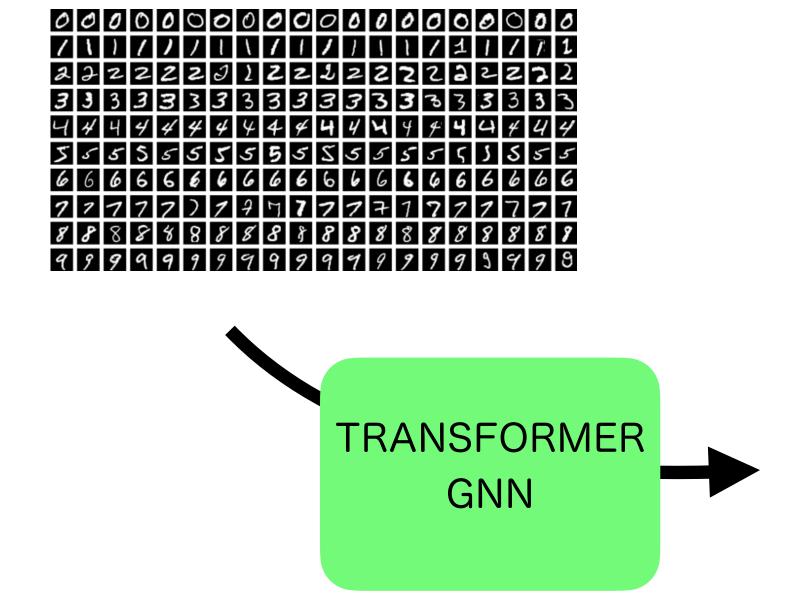
## JET TAGGING: WINDOW TO THE NEW PHYISICS

Mostly talking about top vs QCD classification in this talk

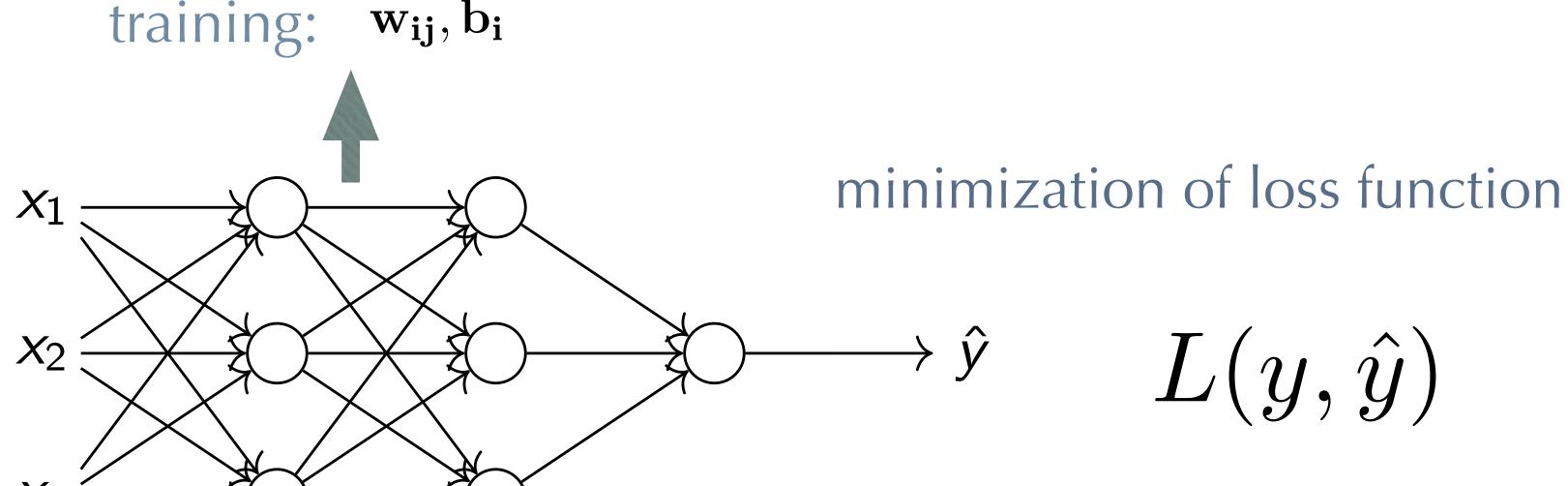


# Training for classification

data pool of images



(28x28) の画像データをn 個



y: truth

output

$$\hat{y} = (\hat{y}_1, \hat{y}_2, \dots, \hat{y}_{10}), \sum_{i} \hat{y}_i = 1$$

 $L(y, \hat{y})$ 

$$\hat{y}_i = \frac{\exp(x_i)}{\sum_j \exp(x_j)}$$

 $\hat{y}$ 's represent likeliness to be y

✓表現力 expressive power

✓データを学習 learn from data

✓ 微分可能Simple linear algebra + activation

# 1. "TRANSFORMER": SELF ATTENTION

X: n( #particle in the jet ) x (feature of particle)

output size = input size

Transformation before MLP

$$X' = X + \delta X, \ \delta X = \alpha \cdot V \equiv \operatorname{softmax}\left(\frac{QK^T}{\sqrt{d}}\right) \cdot V$$
  $Q = XW_Q, K = XW_K, V = XW_V$ 

Pro

- 1.Attention Matrix  $\alpha = QK^T$  All particle-particle attention calculated at once, while GNN takes care nearby information only
- 2. All feature can be included.
- 3. Each layer produce  $X^3$  effect, while MLP need three steps to take into acctound the effect.
- 4. W can be chosen so that X and  $\delta X$  can be added. Freedom choose attention pair.
- 5.skip connection-no information loss X → X'→X''

Con

Tough for memory, calculation. Physics may allow more compact descriptions.

# Features in the context of jet classification

			Definition	JET
Particle momentum	Kinematics	$\Delta\eta$ $\Delta\phi$ $\log p_{ m T}$ $\log E$ $\log rac{p_{ m T}}{p_{ m T}({ m jet})}$ $\log rac{E}{E({ m jet})}$	difference in pseudorapidity $\eta$ between the particle and the jet axis difference in azimuthal angle $\phi$ between the particle and the jet axis logarithm of the particle's transverse momentum $p_{\rm T}$ logarithm of the particle's energy logarithm of the particle's $p_{\rm T}$ relative to the jet $p_{\rm T}$ logarithm of the particle's energy relative to the jet energy angular separation between the particle and the jet axis $(\sqrt{(\Delta \eta)^2 + (\Delta \phi)^2})$	
charge,particle ID	Particle identification	charge Electron Muon Photon CH NH	electric charge of the particle if the particle is an electron ( pid ==11) if the particle is an muon ( pid ==13) if the particle is an photon (pid==22) if the particle is an charged hadron ( pid ==211 or 321 or 2212) if the particle is an neutral hadron ( pid ==130 or 2112 or 0)	
displaced vertex	Trajectory displacement	$ anh d_0 \  anh d_z \ \sigma_{d_0} \ \sigma_{d_z}$	hyperbolic tangent of the transverse impact parameter value hyperbolic tangent of the longitudinal impact parameter value error of the measured transverse impact parameter error of the measured longitudinal impact parameter	

# THIS TALK: IMPROVEMENT OF THE PARTICLE TRANSFORMER

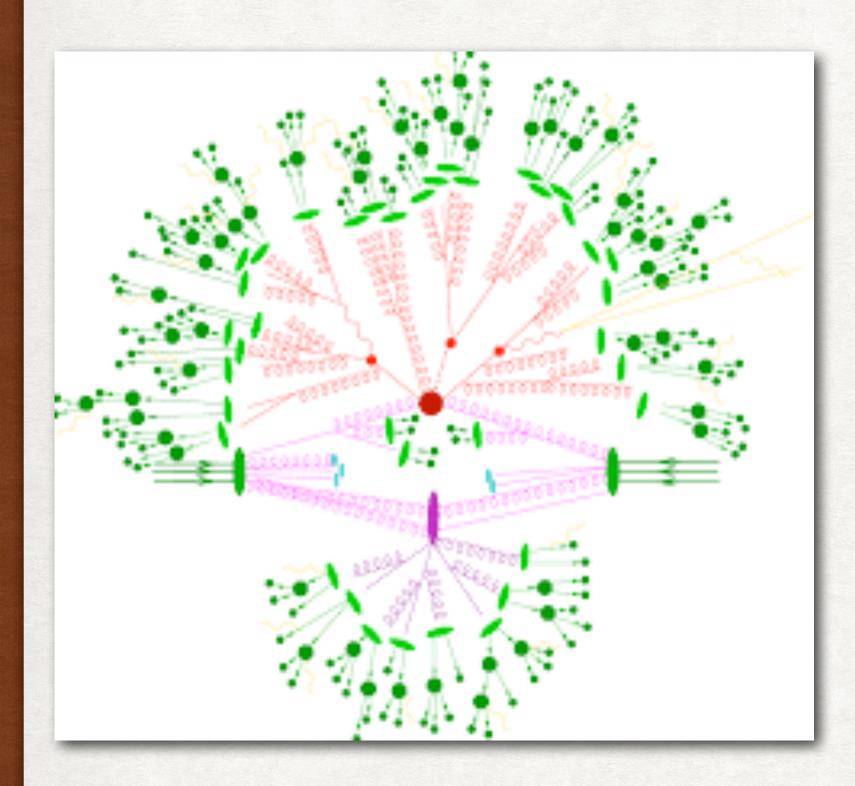
	Accuracy	AUC	$1/\epsilon_B(\epsilon_s = 0.5)$	$1/\epsilon_B(\epsilon_s=0.3)$	Parameters
Lorentz invariance based networks					
PELICAN[35]	0.9426	0.987		$2250\pm75$	208K
LorentzNet[70]	0.942	0.9868	$498 \pm 18$	$2195 \pm 173$	224K
L-GATr[71]	0.942	0.9870	$540 \pm 20$	$2240\pm70$	
Attention based networks					
ParT[49]	0.940	0.9858	$413 \pm 6$	$1602 \pm 81$	2.14M
MIParT[50]	0.942	0.9868	$505 \pm 8$	$2010 \pm 97$	720.9K
Mixer[21]	0.940	0.9859	$416 \pm 5$	<del></del>	86.03K
OmniLearn[72]	0.942	0.9872	$568 \pm 9$	$2647 \pm 192$	1.6M
Plain Transformer*	0.927	0.979	$362 \pm 7$	$780 \pm 73$	1.7M
IAFormer*	0.942	0.987	$510\pm 6$	$2012 \pm 30$	211K

Yellow bands highlight our works!

#### 1. TRANSFORMER FOR PARTON SHOWER+HADRONIZATION

"Ahmed Hammad, & MN

arXiv 2404 14677 JHEP 06 (2024) 176

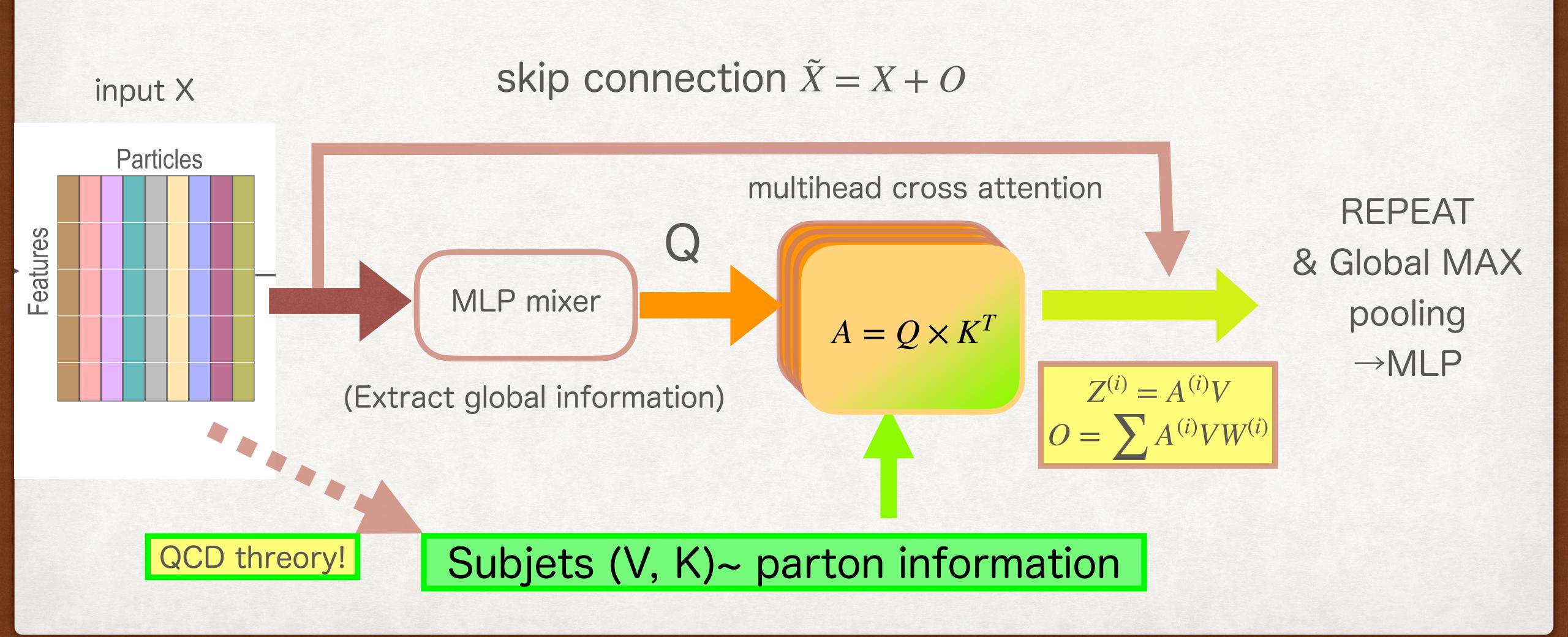


- · Hard Process = Partons(quarks and gluons) {y}
- a jet: P(hadrons in jets | parton ~ jet) =  $P(\{x_i\} | \{y\})$
- . jet with substructure  $P(\{x_i\} | \{y_\alpha\})$
- · Maybe several fatjets in an event (factorization)

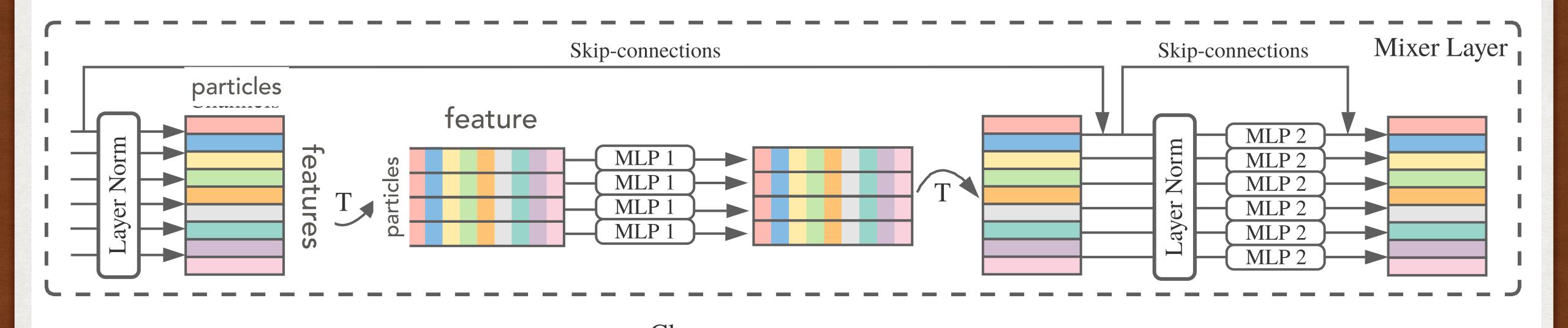
$$P(\{x_i\}, \{x_j'\}, \{y_\alpha\}, \{y_\beta'\}) \sim P(\{x_i\} \mid \{y_\alpha\}) P(\{x_i'\} \mid \{y_\beta'\}) P(\{y_\alpha, y_\beta'\})$$

We need the network forcusing on partons(subjets/jets) vs hadrons

# ATTENTION → CROSS Attention for P(h| subjets) estimation



## MLP MIXER



MLP 1: operate on features

MLP 2: operate on particles

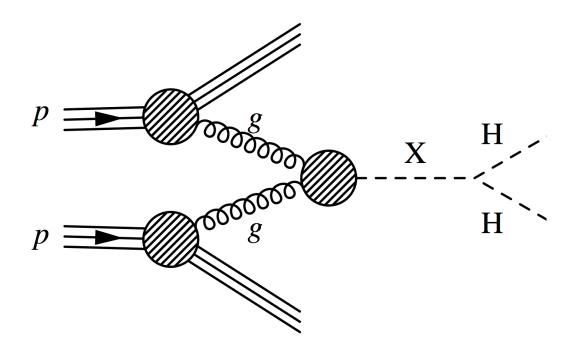
# THE PERFORMANCE FOR TOP VS QCD CLASSIFICATION

	Accuracy	AUC	$1/\epsilon_B(\epsilon_s = 0.5)$	$1/\epsilon_B(\epsilon_s=0.3)$	Parameters
Lorentz invariance based networks					
PELICAN[35]	0.9426	0.987		$2250\pm75$	208K
LorentzNet[70]	0.942	0.9868	$498 \pm 18$	$2195 \pm 173$	224K
L-GATr[71]	0.942	0.9870	$540 \pm 20$	$2240\pm70$	
Attention based networks					
ParT[49]	0.940	0.9858	$413 \pm 6$	$1602 \pm 81$	2.14M
MIParT[50]	0.942	0.9868	$505 \pm 8$	$2010 \pm 97$	720.9K
Mixer[21]	0.940	0.9859	$416 \pm 5$	<del></del>	86.03K
OmniLearn[72]	0.942	0.9872	$568 \pm 9$	$2647 \pm 192$	1.6M
Plain Transformer*	0.927	0.979	$362 \pm 7$	$780 \pm 73$	1.7M
IAFormer*	0.942	0.987	$510\pm 6$	$2012 \pm 30$	211K

Yellow bands highlight our works!

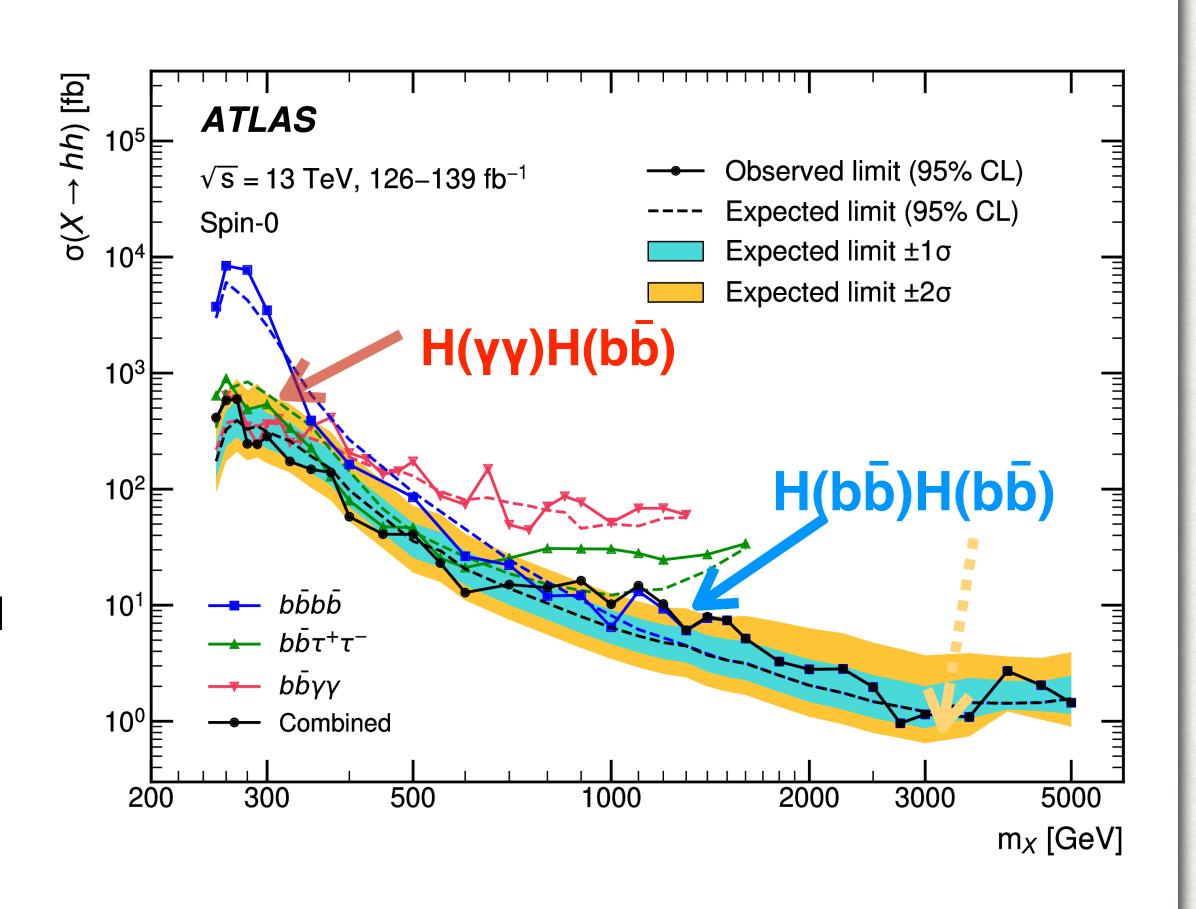
# 2. GLOBAL EVENT ANALYSIS

# $X \rightarrow HH$



H(bb)H(bb) most sensitive channel for  $m_X > 400/500 \text{ GeV}$ H(yy)H(bb) complement in the low mass

#### Phys. Rev. Lett. 132 (2024) 231801





# cross attention for 2 fatjet events

MLP & softmax

**CROSS ATTENTION** 

jet constituent information relevant gives extra weight to the corresponding jets though backward propagation

Hammad, Moretti, MN JHEP 03, 2024

step 2:multihead cross attention transform jet kin by cross Att. [substracture]x [jet kin]

ADD

step 1 : multihead self attention [substructure] x[substructure] [jet kin]x [jet kin]

We can replace transformer to "mixer+subjet" network

Transformer

1st Leading jet Transformer

2nd leading jet

Transformer

jet kinematics

# 3. IAFormer (=InterAction transFormer)

3-1. Improvement of attention matrix.

Esmail, Hammad, Nojiri 2025.03258

original input for attention  $\alpha = softmax(QK^T)$ 

particle information

-  $P_4 = (p_x, p_y, p_z, E)$  : particle 4-momentum

-  $\Delta \eta = \eta - \eta_{\rm jet}$  : pseudorapidity difference

-  $\Delta \phi = \phi - \phi_{\rm jet}$  : azimuthal angle difference

-  $\Delta R = \sqrt{(\Delta \eta)^2 + (\Delta \phi)^2}$  : angular distance from jet axis

-  $\log(p_T)$  : transverse momentum (GeV)

 $-\log(E)$  : energy (GeV)

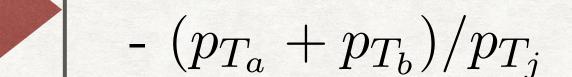
-  $\log\left(\frac{p_T}{p_{T_{\rm jet}}}\right)$  : normalized  $p_T$  (GeV)

 $-\log\left(\frac{E}{E_{\text{int}}}\right)$  : normalized energy (GeV)

**IAFormer attention**  $\alpha = \operatorname{softmax}(\mathcal{F}_{ij})$ 

$$\mathcal{F}_{ij} = W \cdot I_{ij}$$

 $I_{ij}$  pairwize and boost invariant quantity



$$-(E_a+E_b)/E_j$$

$$-\Delta = \sqrt{(\eta_a - \eta_b)^2 + (\phi_a - \phi_b)^2}$$

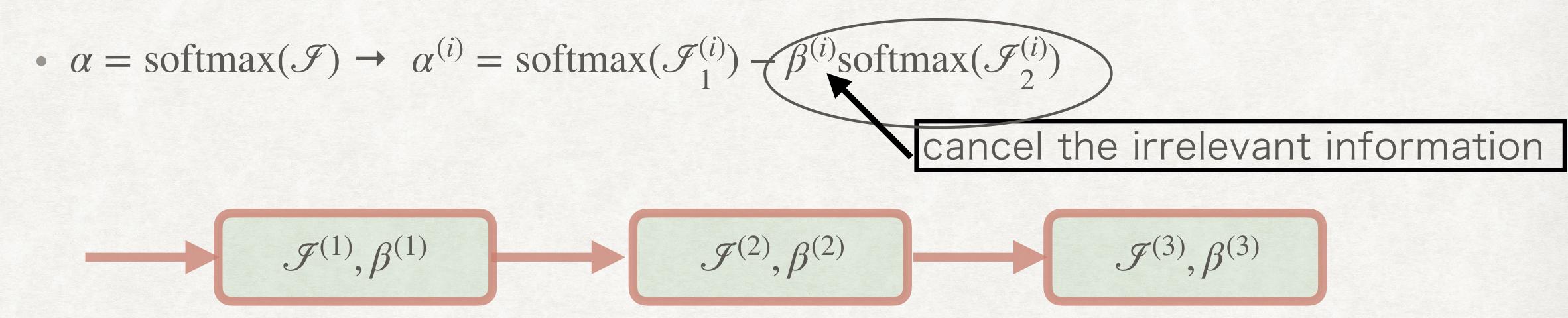
$$-k_T = \min(p_{T_a}, p_{T_b}) \cdot \Delta$$

$$-z = \min(p_{T_a}, p_{T_b})/p_{T_a} + p_{T_b}$$

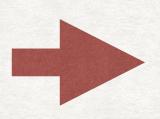
$$-m^2 = (E_a + E_b)^2 - |\mathbf{p}_a + \mathbf{p}_b|^2$$

#### 3-2 Introduction of Differential attention

We have introduced a new dynamic attention called "differential attention" to the network. (see arXiv:2410.05258)



fixed sparse attention

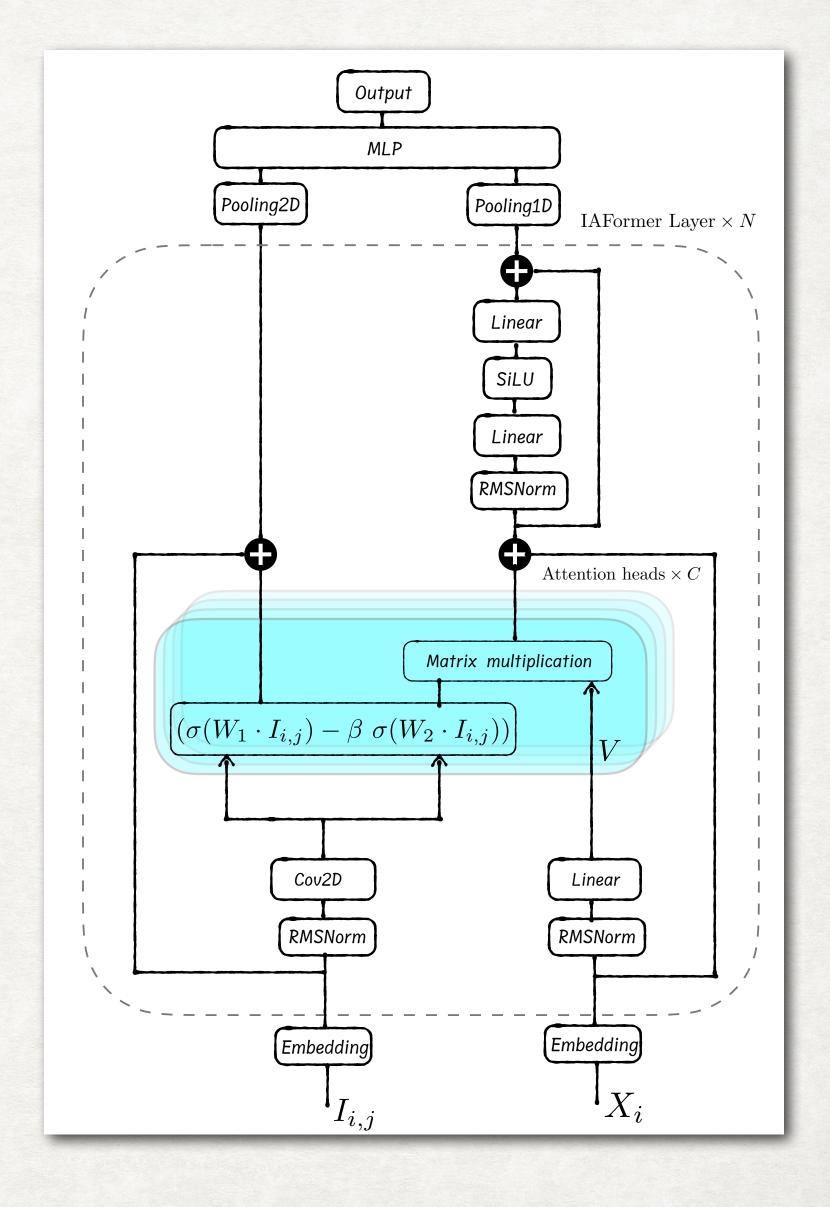


Each layer built different filters dynamically

 $I_{ij}$ ,  $X_i$  are both updated by skip connection(not like PART)

For particle transformer, Attention build from transformed X, but in our network, they develop independently.

Not sure about boost invariance is kept along with layers but leading term is OK....



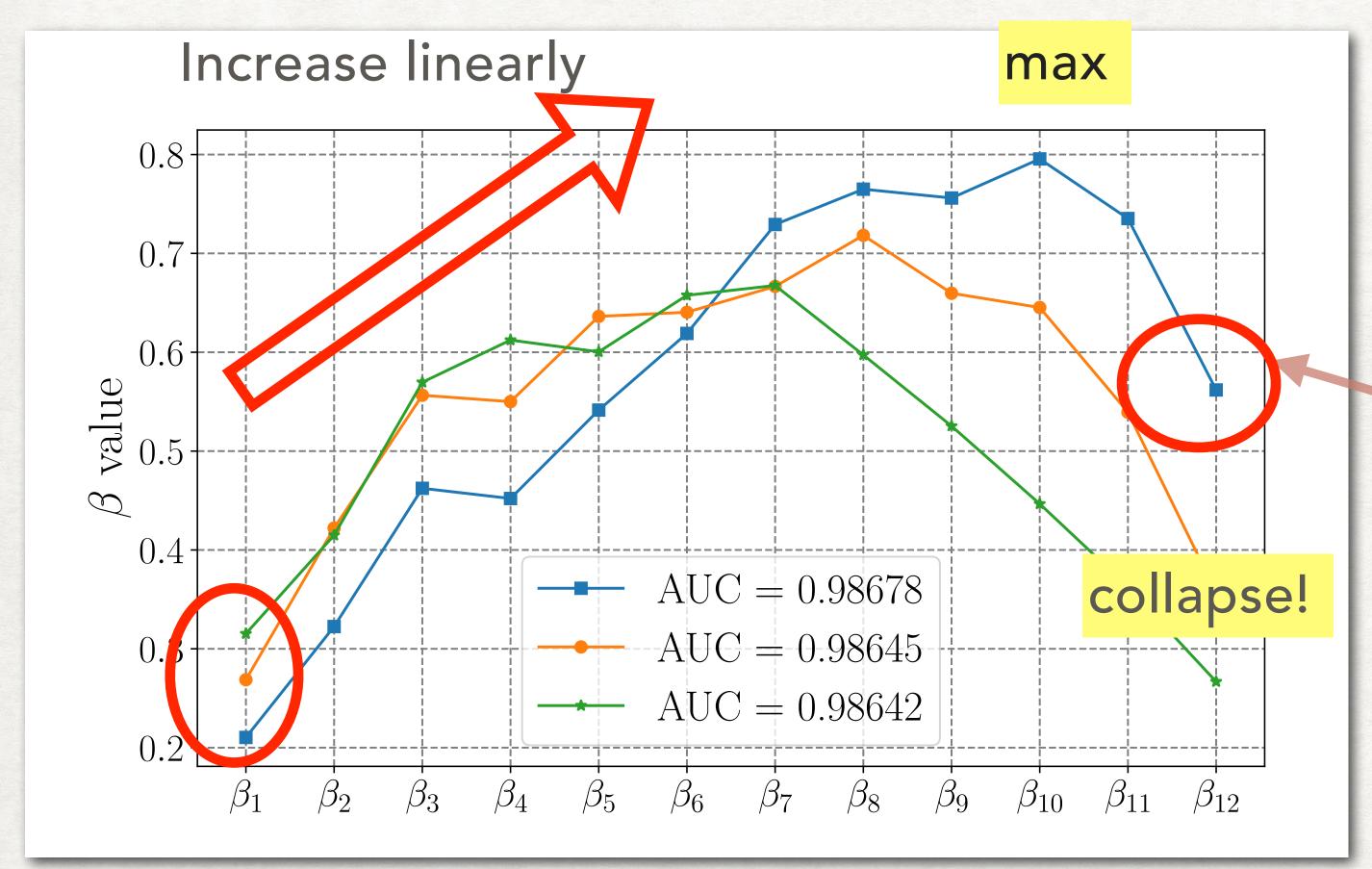
# THE PERFORMANCE FOR TOP VS QCD CLASSIFICATION

	Accuracy	AUC	$1/\epsilon_B(\epsilon_s=0.5)$	$1/\epsilon_B(\epsilon_s=0.3)$	Parameters
Lorentz invariance based networks					
PELICAN[35]	0.9426	0.987		$2250\pm75$	208K
LorentzNet[70]	0.942	0.9868	$498 \pm 18$	$2195 \pm 173$	224K
L-GATr[71]	0.942	0.9870	$540 \pm 20$	$2240 \pm 70$	
Attention based networks					
ParT[49]	0.940	0.9858	$413 \pm 6$	$1602 \pm 81$	2.14M
MIParT[50]	0.942	0.9868	$505 \pm 8$	$2010 \pm 97$	720.9K
Mixer[21]	0.940	0.9859	$416 \pm 5$	<u>—</u>	86.03K
OmniLearn[72]	0.942	0.9872	$568 \pm 9$	$2647 \pm 192$	1.6M
Plain Transformer*	0.927	0.979	$362 \pm 7$	$780 \pm 73$	1.7M
IAFormer*	0.942	0.987	$510 \pm 6$	$2012 \pm 30$	211K

Yellow bands highlight our works!

## 3-4 RESULTS: BEHAVIOR OF DYNAMIC FILTERS

filtered information



Networks minimize finite positive  $\beta$ . We need filters

last  $\beta$  is large

→ higher network performance

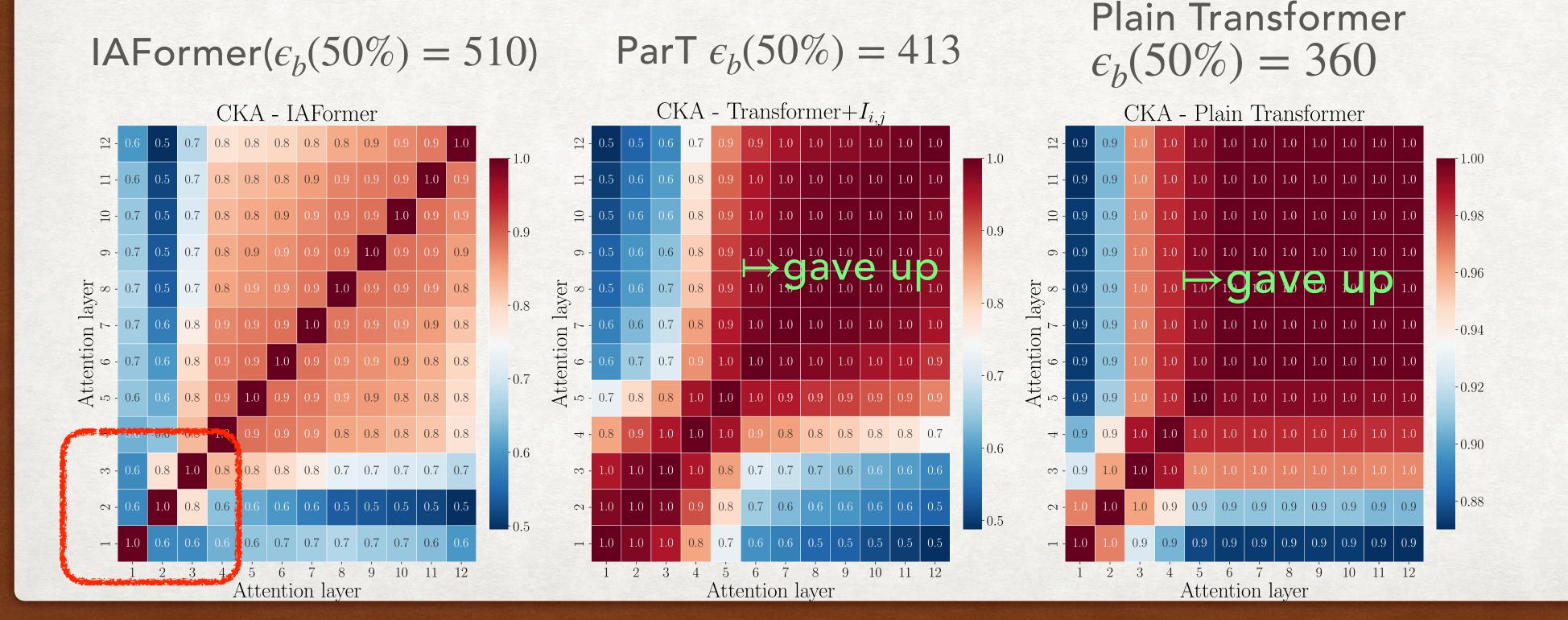
start from small beta

# 3-4 RESULT :Learning pattern analysis CKA similarity

d event-two layer output  $X_1(d \times h_1)$  and  $X_2(d \times h_1) \rightarrow \text{dxd}$  matrix  $M = X_1 X_1^T, N = X_2 X_2^T$ . Then

$$\mathrm{CKA}(\mathrm{M,N}) = \frac{\mathrm{HSIC}(\mathrm{M,N})}{\sqrt{\mathrm{HSIC}(\mathrm{M,M})\mathrm{HSIC}(\mathrm{N,N})}}\,, \qquad \mathrm{HSIC}(\mathrm{M,N}) = \frac{1}{(d-1)^2}\mathrm{Tr}(MHNH) \qquad \qquad H = \delta_{ij} - \frac{1}{d}$$

if CKA~1, two layers are equivalent—and not needed.



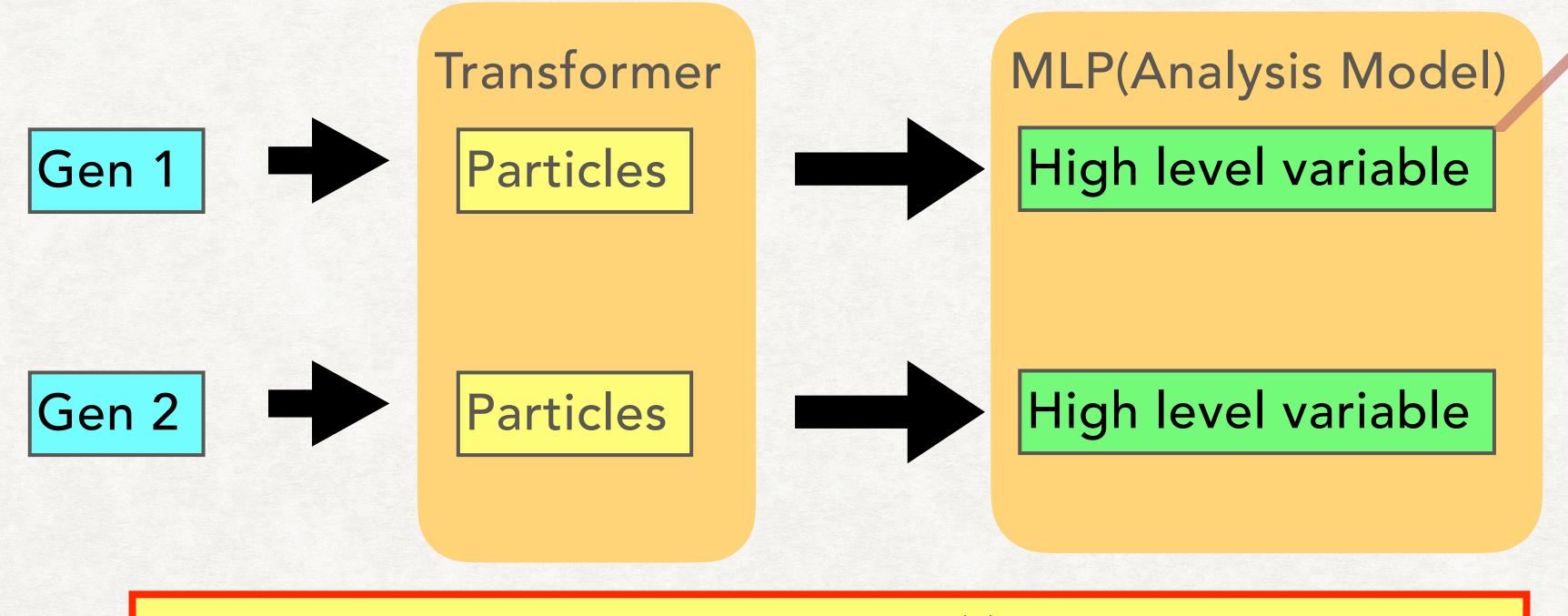
IAFormer is learning efficiently

## SUMMARY: OUR MODEL: CROSS ATTENTION

	Key	Query	Value	update
Particle transformer	particle	particle	particle	particle
"Mixer"	subjet	particle	subjet	particle
"Global analysis"	particle	jet	jet	jet
"IAFormer"	boost invaria	ant pairwise	particle	particle

## 4. ML FOR GENERATOR COMPARISON

What is this?



 $s_{gen}(x)$ : generator classifier output  $w = \frac{s_{gen}(x)}{1 - s_{gen}(x)}$  estimated probability ratio

Transformer: no human bias, poor training stablity

Analysis Model (AP):MLP using highlevel input: stable prediction, ideal for generator comparison

# HL feature for generator classification

Amon Furuichi, Sung Hak Lim, Mihoko M. Nojiri JHEP 07 (2024) 146 JHEP 07 (2025) 111

pt distribution of constituents

Jet spectrum

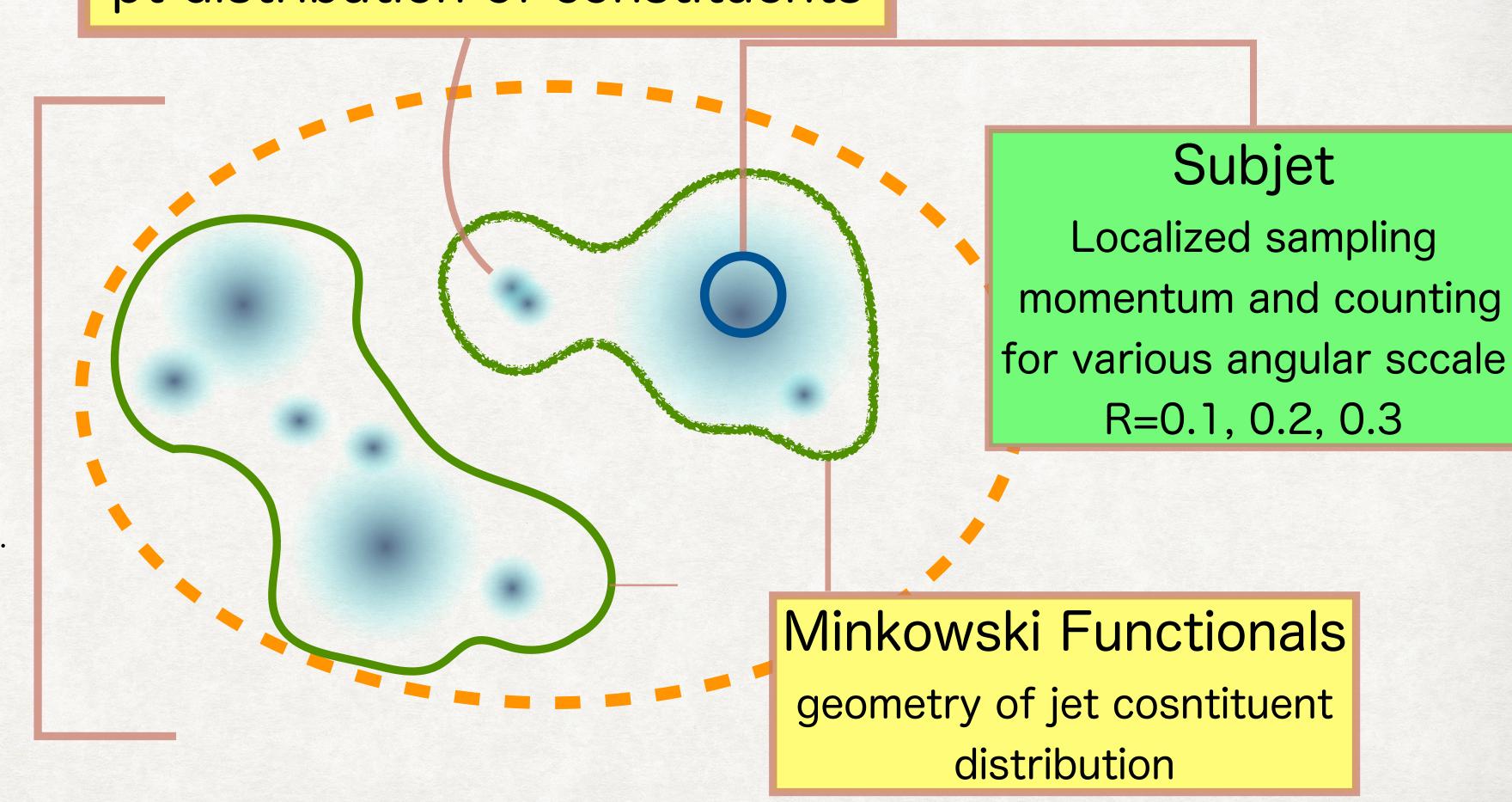
two point Energy

correlation

(unlocalized sampling )

= EFP with N=2

$$S_{2,ab}(R) \stackrel{\text{def}}{=} \sum_{i \in a} \sum_{j \in b} p_{T,i} p_{T,j} \delta(R - R_{ij}).$$



#### PYTHIA VS HERWIG UNDER ML

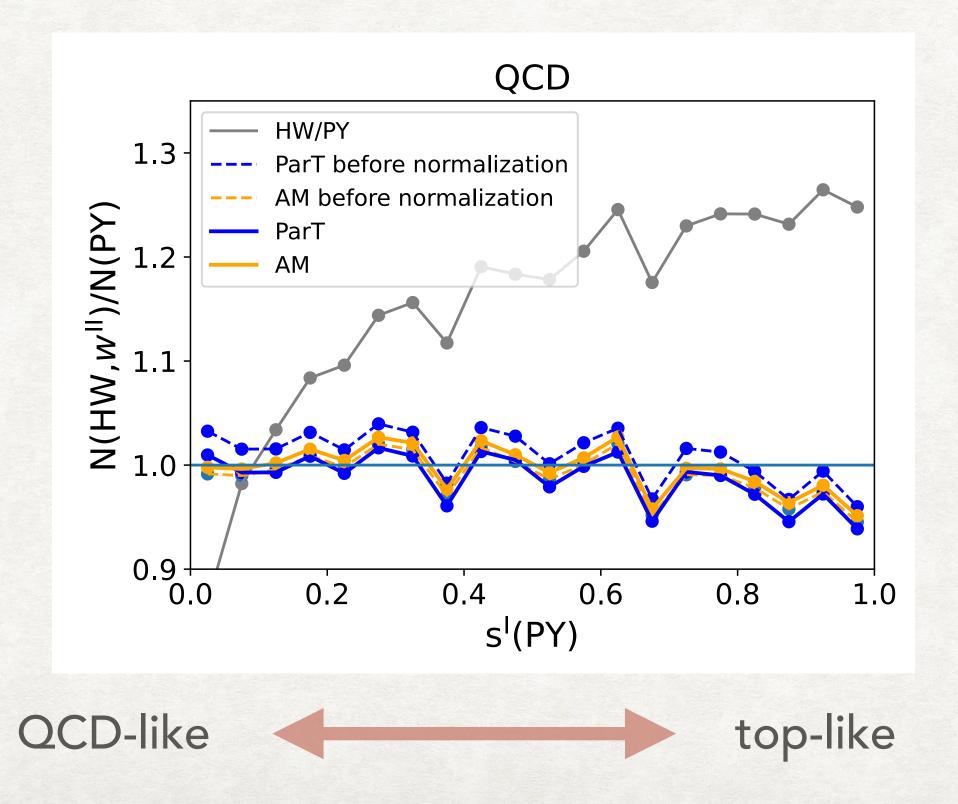
 $s_{gen}(x)$ : generator classifier output

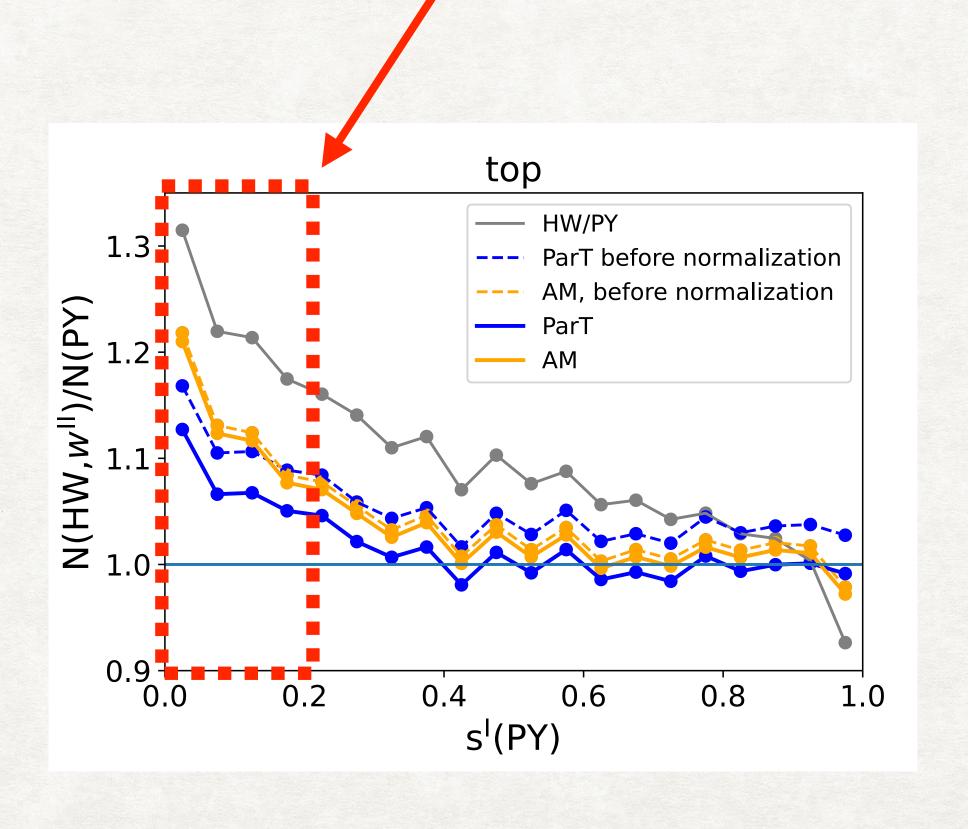
$$v = \frac{s_{gen}(x)}{1 - s_{gen}(x)}$$
 estimated probability ratio

poor overlap between PY vs HW in QCD-like top jet

top-like

agreement of reweighted HW distribution to PY quantify DL infering





QCD-like

# HL feature for generator classification

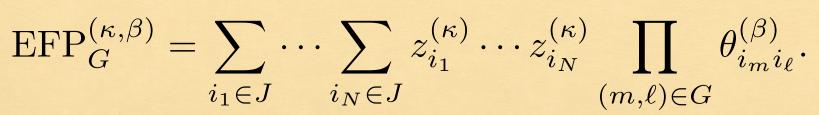
Amon Furuichi, Sung Hak Lim, Mihoko M. Nojiri JHEP 07 (2024) 146 2312.11760

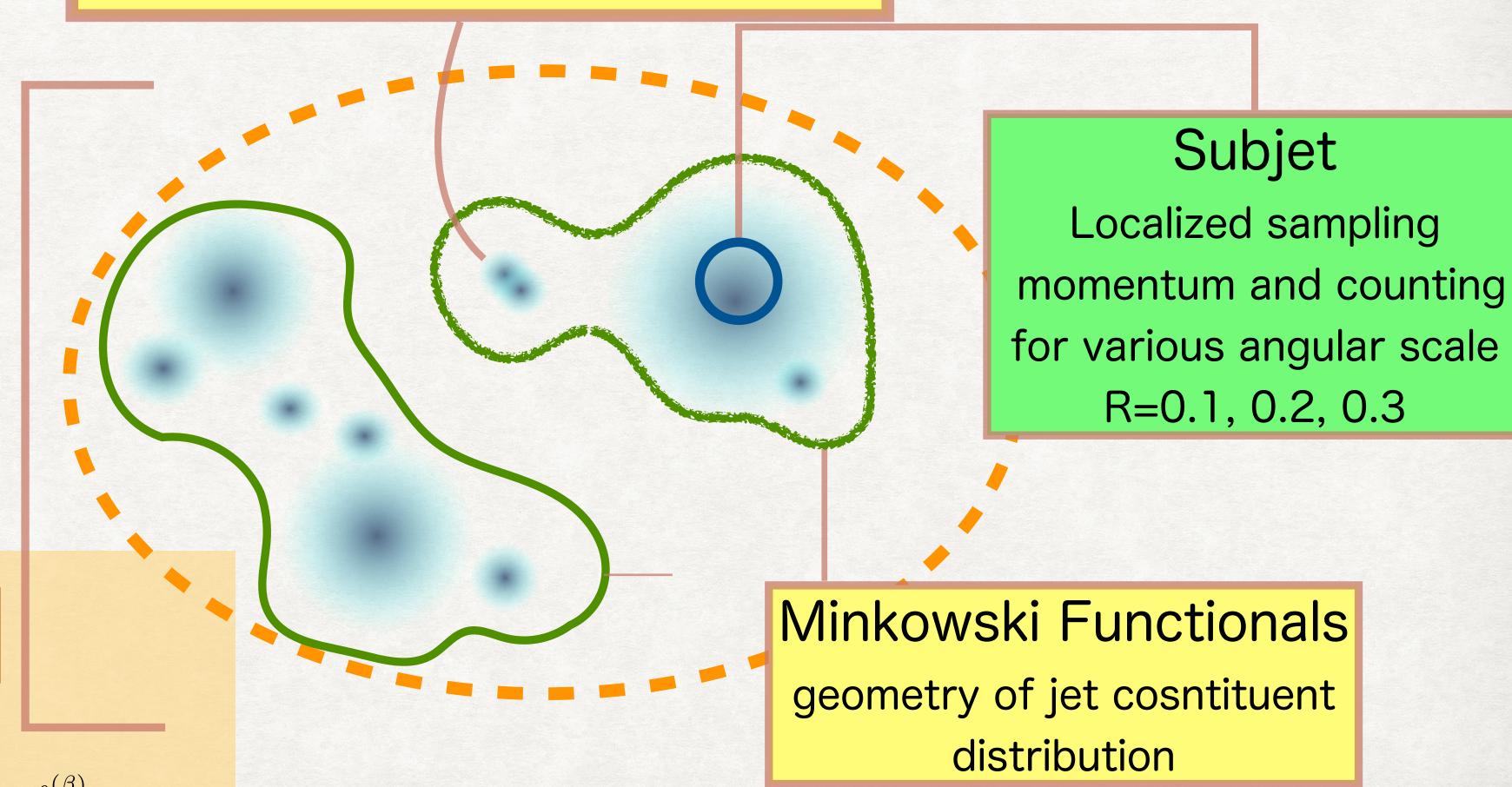
pt distribution of constituents

Jet spectrum
two point Energy
correlation
(unlocalized sampling)

$$S_{2,ab}(R) \stackrel{\text{def}}{=} \sum_{i \in a} \sum_{j \in b} p_{T,i} p_{T,j} \delta(R - R_{ij}).$$

Energy flow polynomials





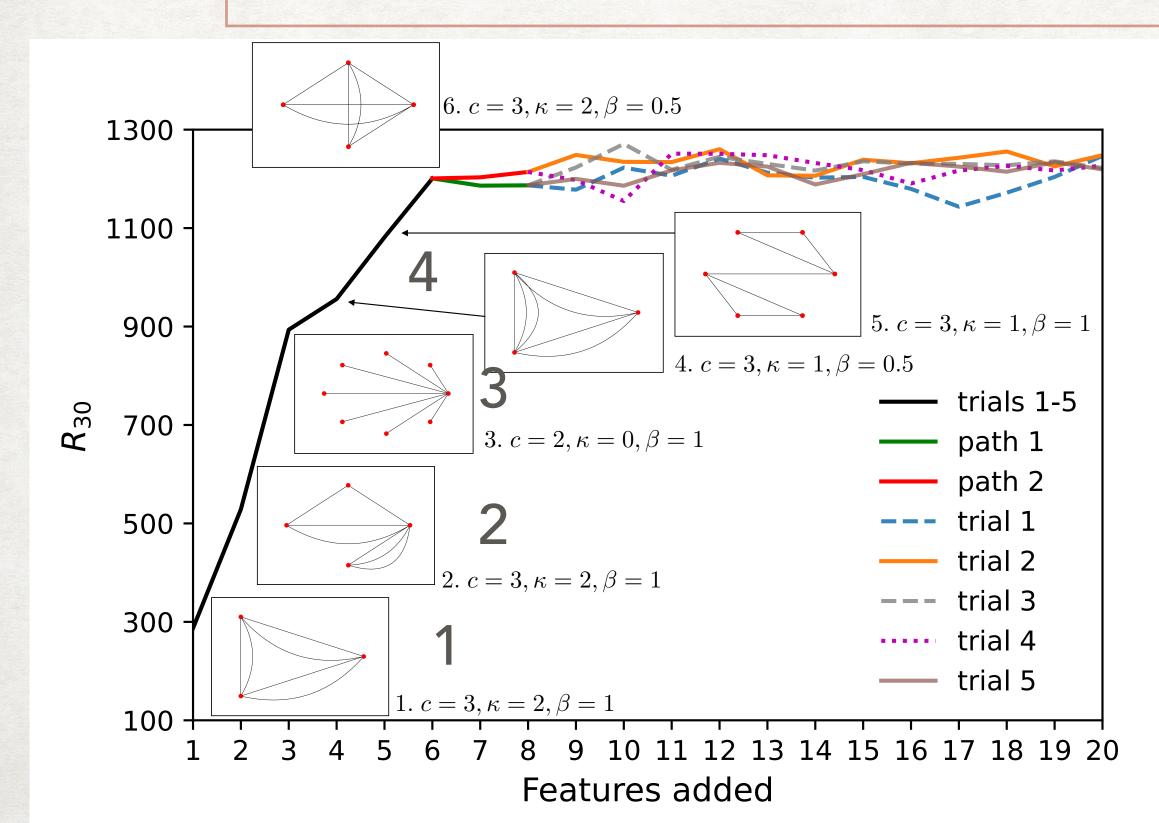
#### ENERGY FLOW POLYNOMIAL AND TOP JET CLASSIFICATION

"Feature selection with distance correlation", Ranit Das, Gregor Kasieczka, David Shih *Phys.Rev.D* 109 (2024) 5, 054009

#### Energy flow polynomial

$$EFP_G^{(\kappa,\beta)} = \sum_{i_1 \in J} \cdots \sum_{i_N \in J} z_{i_1}^{(\kappa)} \cdots z_{i_N}^{(\kappa)} \prod_{(m,\ell) \in G} \theta_{i_m i_\ell}^{(\beta)}.$$

take all possible combination of N constituent and weight by common power k for energy of the particle and common power  $\beta$  for angle according to the graph

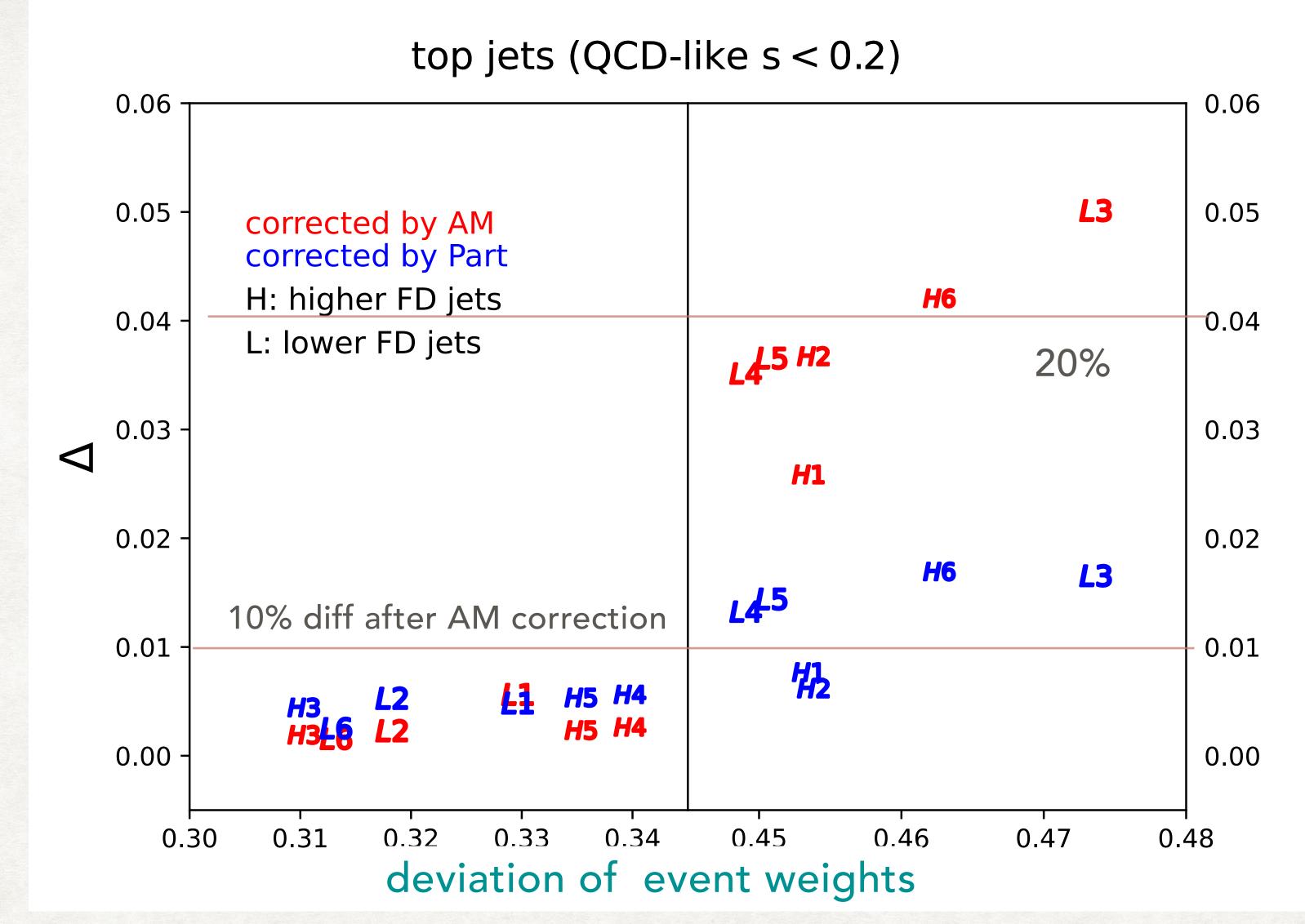


After reweigting by AM Pythia and HW difference remains for "QCD like top event" with

Higher FD1, FD2, FD6 ( $\kappa = 2$ , 3 or 4 points) Low FD3, FD4, FD5 ( $\kappa = 0.1 \ \theta$  power ~3.5,7, 8)

corresponding to narrow jet where a few jet constituents take most of jet energy

reweighting deviation after



$$\Delta_{\mathrm{L}}(\mathrm{FD}_n) \coloneqq \left(\frac{N_{\mathrm{L}}(\mathrm{FD}_n|\mathrm{HW};w)}{N_{\mathrm{L}}(\mathrm{FD}_n|\mathrm{PY})} - 1\right)^2$$



smaller overlap between PY & HW evenets

#### TAKE AWAY MESSAGE

- 1. fast, lightweight, while keeping performance
- 2. Incorporate physics picture
- 3. Jet analysis → event analysis.(H→hh)

RESPCETING QCD

Cross attention is important

- 4. Respect symmetry Replacing "attention from generic features"

  → "pairwise boost invariant information " (IAFormer) Symmetry
- 5. Reduce valiance in training

Improved stability within DL

6. Identify the key parameters for classifications

Identify Important variables in DL era Improving MC simulation