# Multi-scale Cross-Attention Transformer encoder for event classification
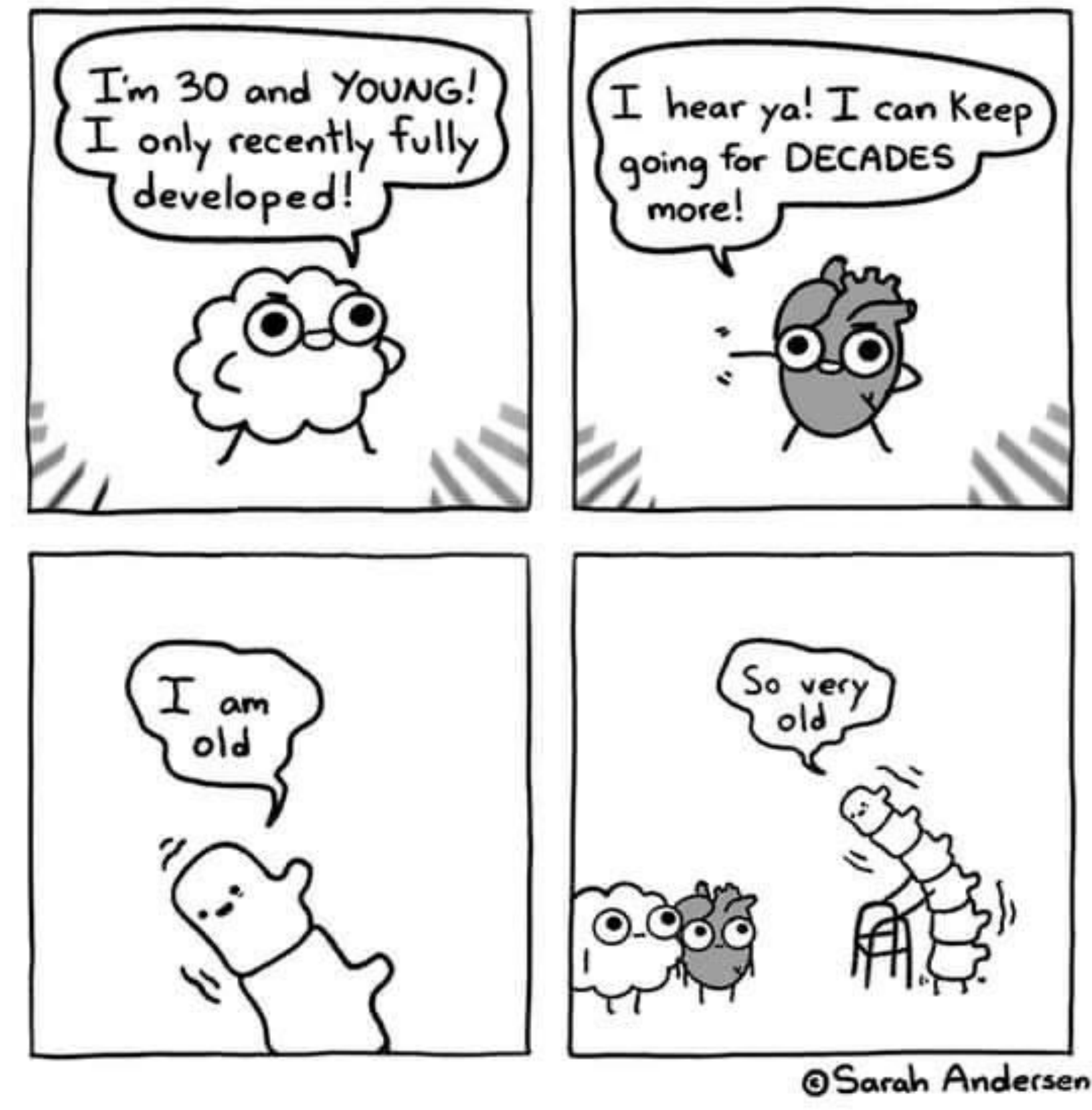
Mihoko Nojiri(IPNS, KEK), with Ahmed Hammad  Stefano Moretti
arXiv 2401.00452

# ABOUT MYSELF

- PhD Kyoto (1990) a bit old.

- PD: Supergravitiy study in heavy top era → SUSY dark matter. One of the author of first Sommerfeld effect in dark matter annihilation. (2003)

- Collider:

  - 1996: JLC study and Snowmass

  - 2002-2008 LHC BSM study in ATLAS SUSY group. BSM Convener of Les Houches TeV collider workshop twice → Jet substructure study → **Deep Learning**

- **Service**: JPS executive board member →member of Science Council of Japan(SCJ) working on Diversity Issues .

  - In KEK, we just had DEI workshop last Dec, and trying establish more DEI activities. (https://www2.kek.jp/ipns/en/news/5320/)

**"a young mind",**
**(according to Tilman Plehn)**
**but this makes me cry**

# ML(THEORY) IN JAPAN:  GRANT "MACHINE LEARNING PHYSICS "

| Overview | Organization | Events | Acheivements | Outreach |

Overview

message
Head Investigator

Koji Hashimoto

Professor
Particle Physics Theory Group
Department of physics, Kyoto University

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

Hello. My name is Koji Hashimoto, Professor of Graduate School of Science, Kyoto University. Let me explain about
the "Learning Physics Domain" that we are just now trying to create. This new transformative research area aims to
revolutionize fundamental physics by combining machine learning and physics.

B01  Akinori Tanaka (Riken  AIP)  Math and application of DL
B02 Yoshiyuki Kabashima (Tokyo) Statistical data and ML
B03 Kenji Fukushima (Tokyo) Topology and Geometry of ML
A01 Akio tomiya (IPUT Osaka)  Lattice
A02 Mihoko Nojiri  HEP
        Junichi Tanaka (ICEPP Tokyo, ATLAS)
        Masako Iawasaki (Osaka Metropolitan Belle II )
        Noriko Takemura and Hajime Nagahara (Data Science)
A03 Tomi Ohtsuki (Sophia U) Condensed Matter
A04 Koji Hashimoto   Quantum and Gravity

Ahmed  Hammad
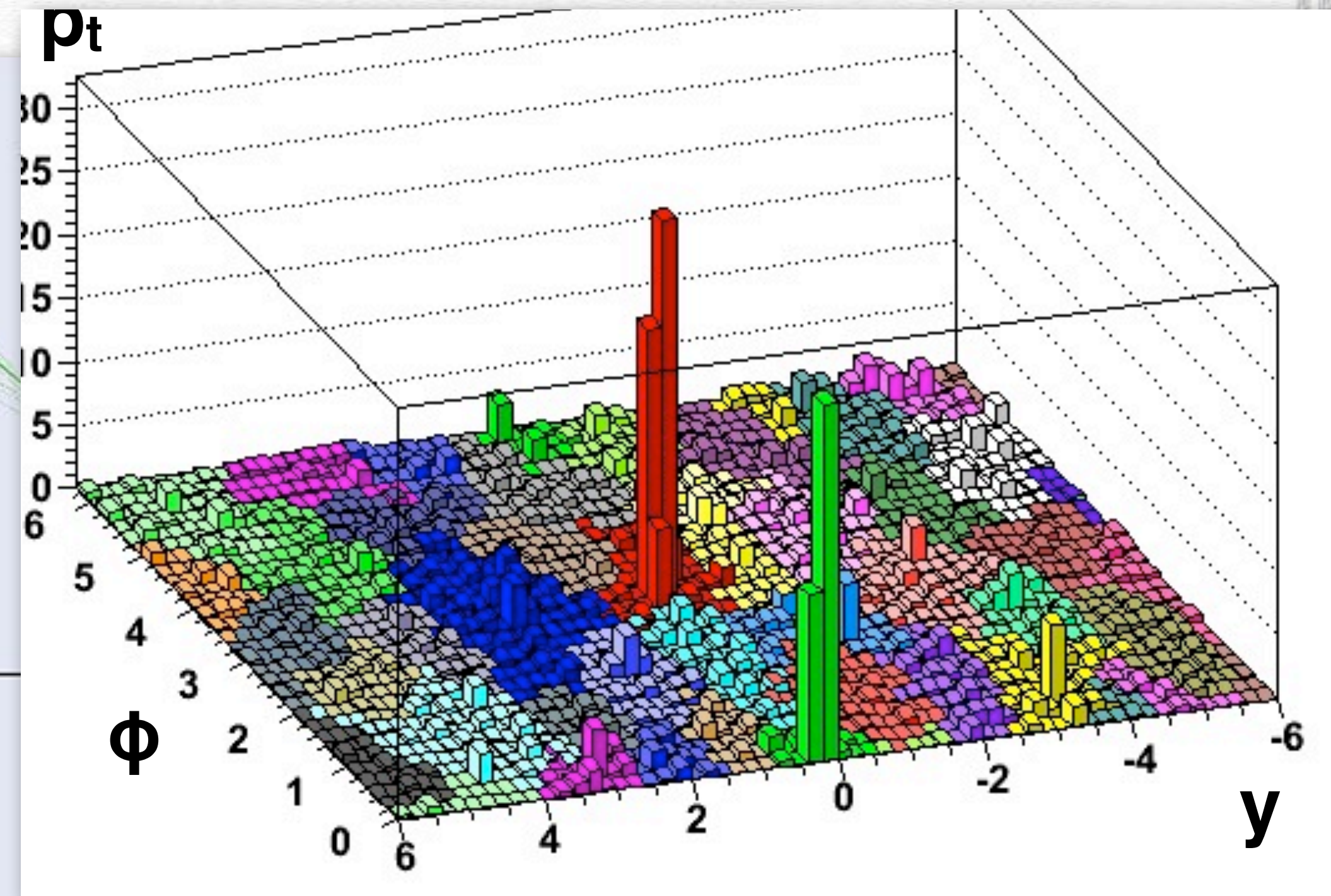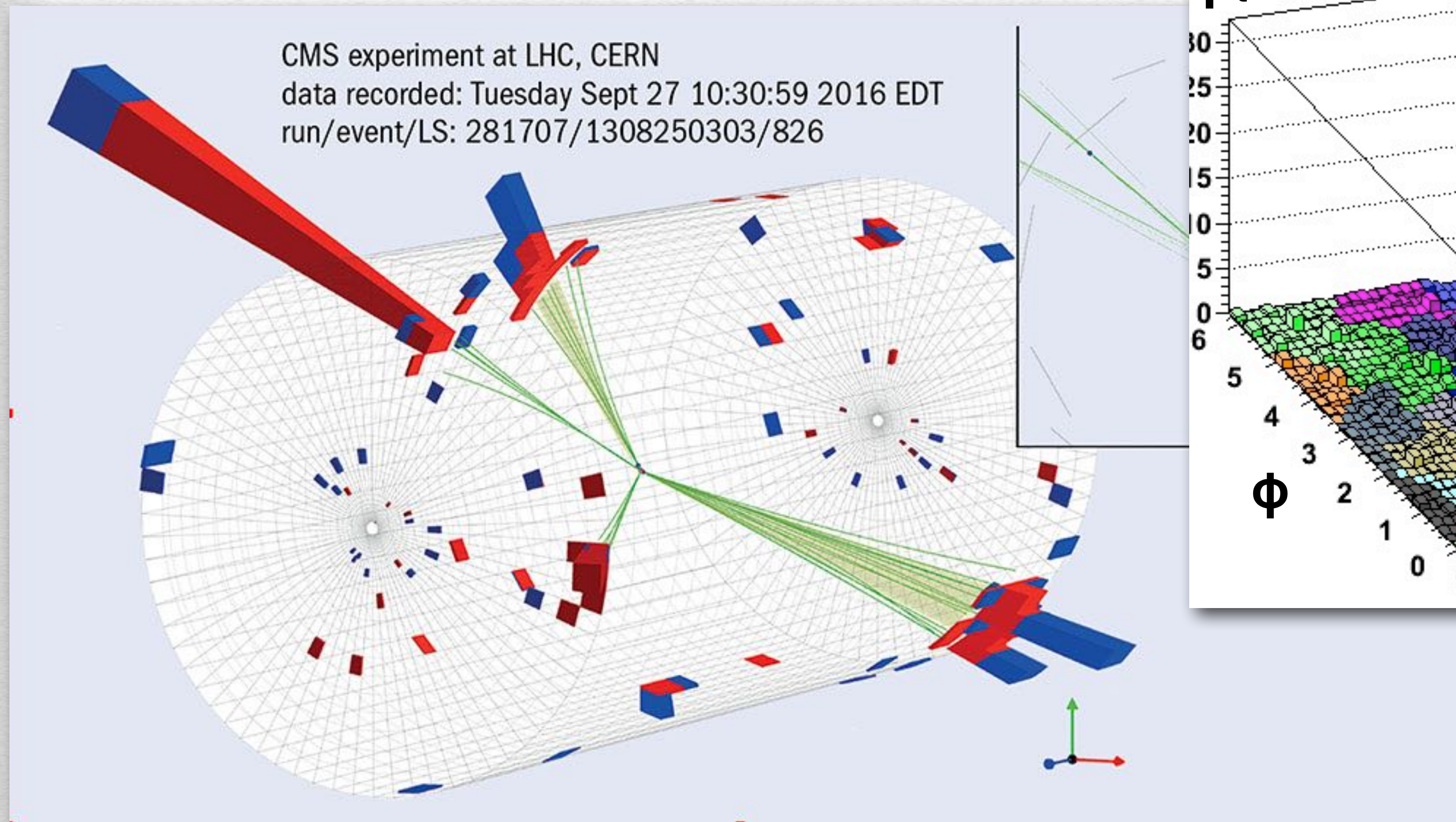2017-2020:  Ph.D Basel University,
Basel Switzerland
2020-2023: SeoulTech, Korea
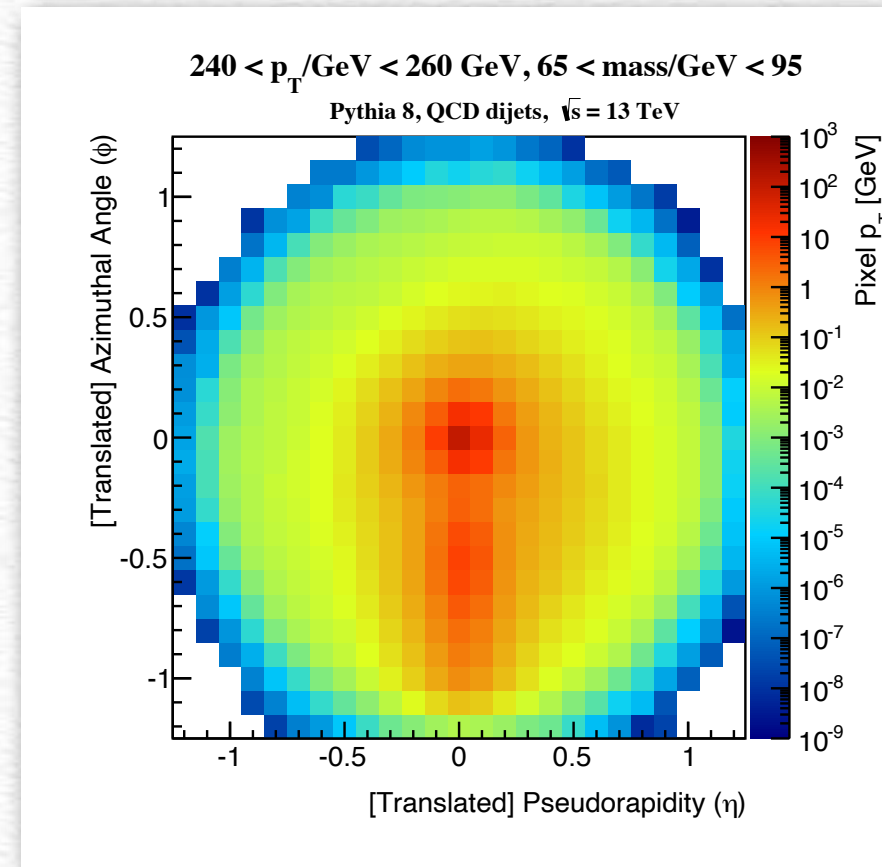2023- KEK

# How machine Leaning help Collider Analysis



CMS experiment at LHC, CERN
data recorded: Tuesday Sept 27 10:30:59 2016 EDT
run/event/LS: 281707/1308250303/826

Jet clustering

# Jet classification using ML

## Deep Convolutional Networks

...arning — convolutional networks in particular — cu...
...ge recognition tasks. We apply a deep convoluti...
...model selection. Below, we visualize a simple archit...

...d that architectures with large filters captured the p...
... The learned filters from the convolutional layers ex...
...that sheds light on phenomenological structures w...

QCD jet

## Physics Performance Improvements

...ur analysis shows that De...
...new physics processes co...
...enhancing the discovery po...
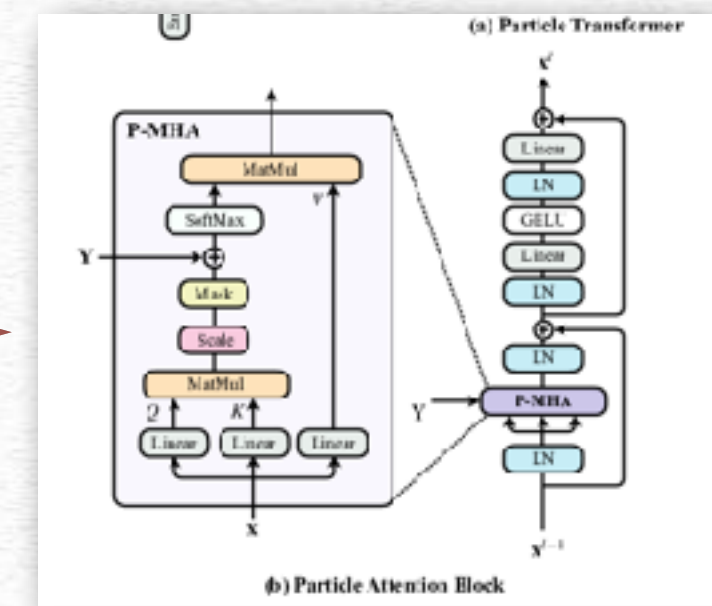...suggests that the deep co...
...physics-motivated variables...

W jet

from Schwartzman et all
https://iopscience.iop.org/article/
10.1088/1742-6596/762/1/012035

Boosted Boson Type Tagging

*Jet ETmiss*

as sets

as graphs

transformer

permutation invariance
(Energy Flow Network and
Particle Flow Network 18...5.65)

sparse data
1902.08570 Particle Net

building key
and query
2202.03772

CNN Oliveira et al
(1511.05190)

Drever et al LundNet (1807.04758)
Gong et al LorentzNet (2201.08187)

meaning that physical variables have no discrimination power. Then, we apply our learned
discriminant, and check for improvement in our figure of merit — the ROC curve.

Notice that removing out the individual effects of
the physics-related variables leads to a likelihood
performance equivalent to a random guess, but
the Deep Convolutional Network retains some
discriminative power. This indicates that the deep
network learns beyond theory-driven variables —
we hypothesize these may have to do with

Bogatskiy et al PELICAN (2211.00454)

- Non SM Higgs boson (Two Higgs doublet model)

  - pp → H (Heavy Higgs boson) → hh → 4 bjet

  - mH=600-2000 GeV, mh=125.11GeV

- Meta stable vaccum of SM → extension of Higgs sector

- why doing Deep Leaning?

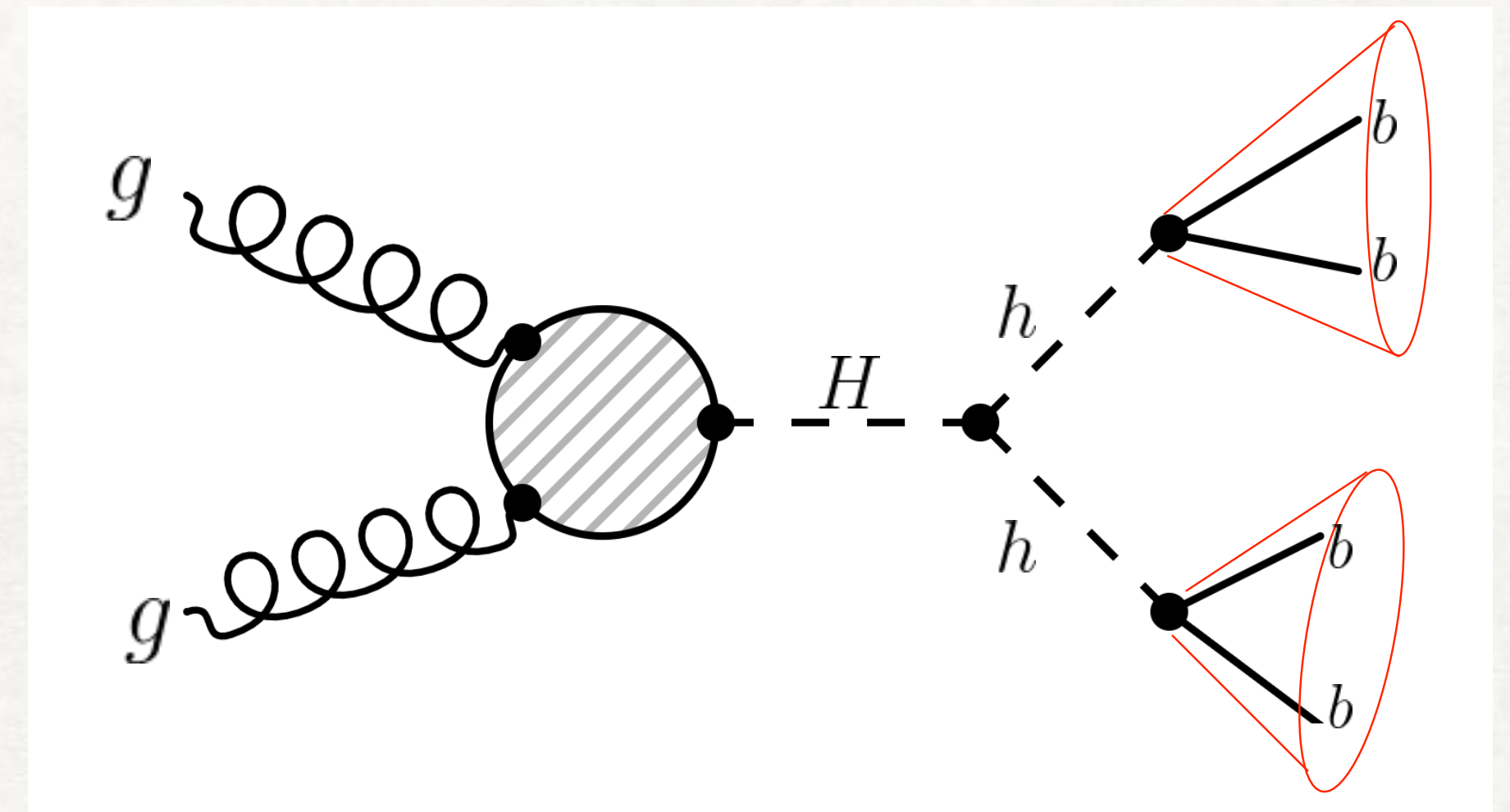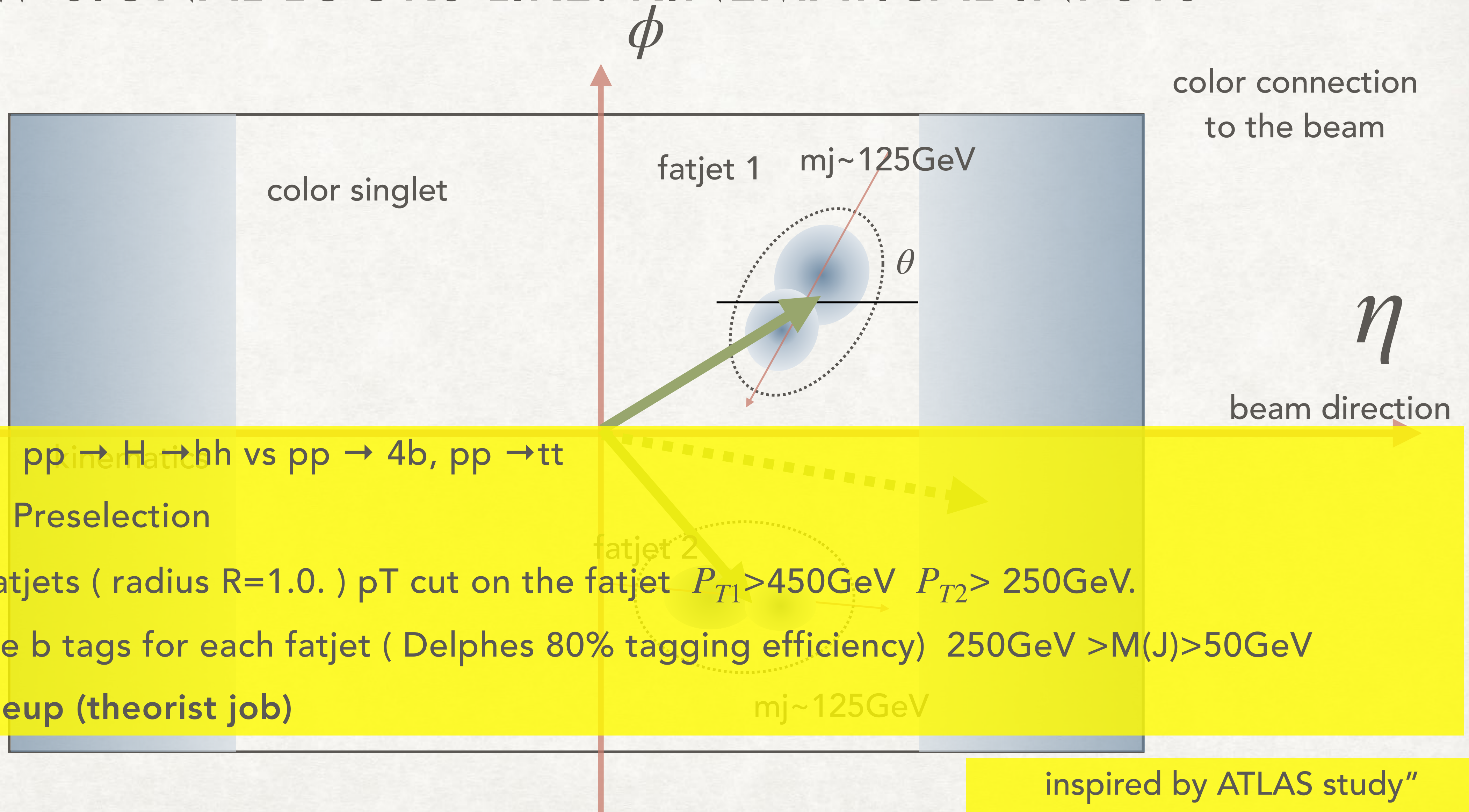  - Sensitivity under S/BG~1 scale by $1/\sqrt{N}$ with background rejection $1/N$



Figure 2: Feynman diagram for the signal process.

$\phi$

color connection
to the beam

fatjet 1    mj~125GeV

color singlet

$\theta$

$\eta$

beam direction

- Delphes  pp → H →hh vs pp → 4b, pp →tt

- Delphes Preselection

  - two fatjets ( radius R=1.0. ) pT cut on the fatjet  $P_{T1}$>450GeV  $P_{T2}$> 250GeV.

  - double b tags for each fatjet ( Delphes 80% tagging efficiency)  250GeV >M(J)>50GeV

  - **no pileup (theorist job)**

fatjet 2

mj~125GeV

# INPUT TO NETWORK : EVENT KINEMATICS

$\phi$

color connection
to the beam

color singlet

fatjet 1    mj~125GeV

$\theta$

$\eta$

beam direction

kinematics

Kinematical inputs (3, 6)
fatjet 1 = $(m_1, \eta_1, \phi_1, p_{T1}, E_1), \theta_1$
fatjet 2 = $(m_2, \eta_2, \phi_2, p_{T2}, E_2), \theta_2$
H candidate = $(m_{12}, \eta_{12}, \phi_{12}, p_{T12}, E_{12}), \theta_{12} = 0$

fatjet 2

mj~125GeV

NOTE : "5 inputs for 4 momentum" , H candidate momentum as sum of two fat jets, add θ,

$\phi$

color connection
to the beam

color singlet

jet 1

mj~125GeV

$\theta$

$\eta$

beam direction

jet 2

mj~125GeV

$\bar{\eta} = \bar{\phi} = 0$

**up to 50 constituents:**
**Regularization speed up the training and reduce**
**the required events.**
1. shift coordinate to (0,0)
2. rotate jet based on covariant matrix
3. flip $\eta$ so that E($\bar{\eta} > 0$) > E($\bar{\eta} < 0$)
4. particles are ordered by pT and we take up to 50
$p_i = (\bar{\eta}_i, \bar{\phi}_i, p_{Ti}, \log p_{Ti}) \rightarrow$ (50, 4) data

Naive approach "simple concatenation"

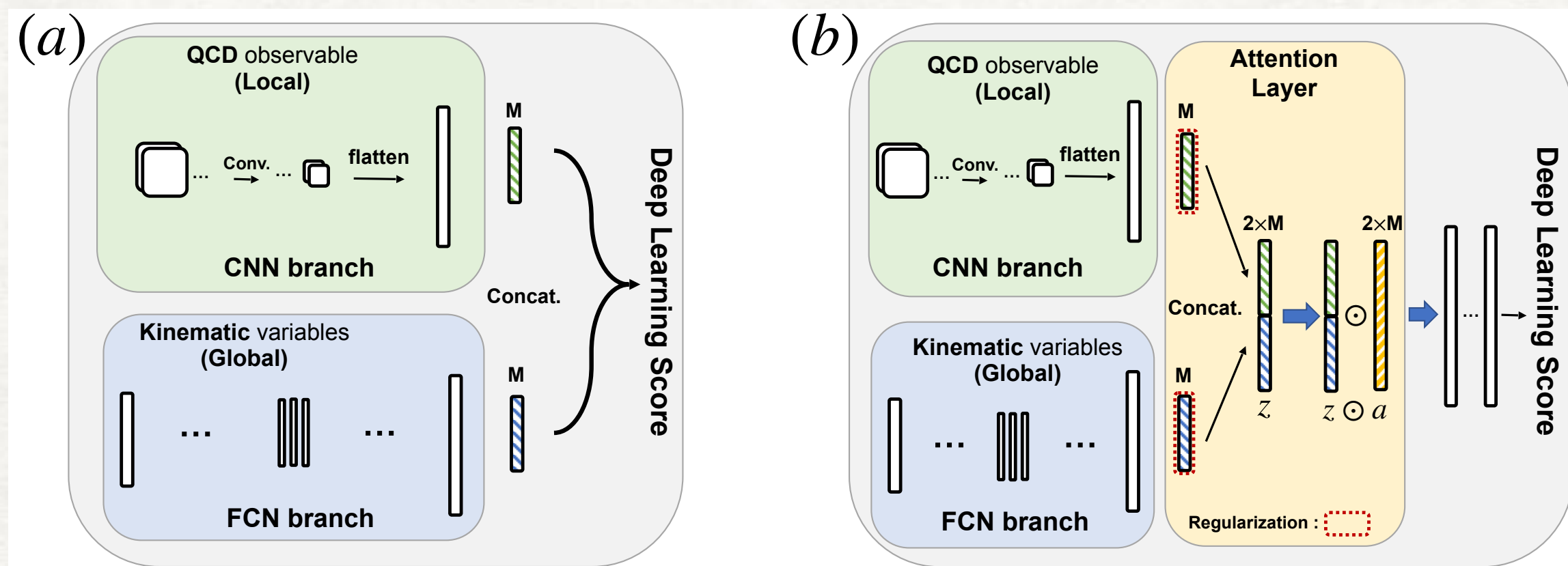2311.16674[hep-ph] K. Ban, KC Kong, M Park, S.C. Park



FIG. 2. The schematic plots for neural network structures: (a) conventionally used one in previous studies only with concatenation and (b) our proposed one with a regularized attention mechanism.

a) [Jet momentum (parton momentum) ]+[jet concatenation] does not work.
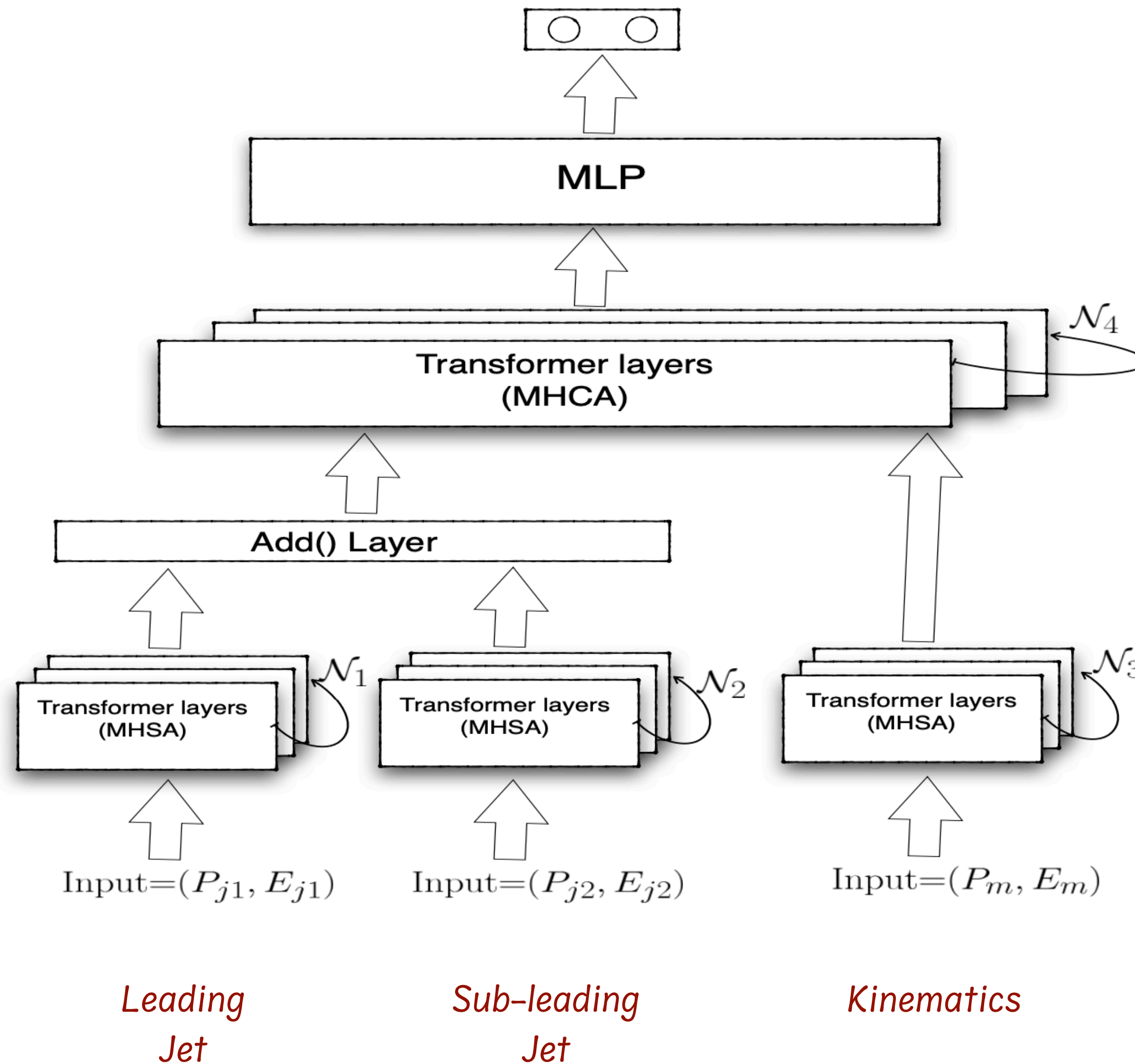because of imbalance of "importance" of two information → the minor one can be ignored in the training.
Pre-training and freeze substructure analysis? We would loose the correlation to global kinematics.

# OUR CROSS ATTENSION MODEL

# TAKEAWAYS

- use "cross attention" when you combine the "high scale information" to the "low energy scale", because cross attention layer gives extra emphasis to the information linked to the high energy kinematics.

- skip connection and Interpretation : Skip connection helps to maintain some connection to the inputs

- More Physics: Heavy particles decay into colored particles (discovery, spin, color structure? ) Cross attention network probably more useful to resolve correlation of jet structures.
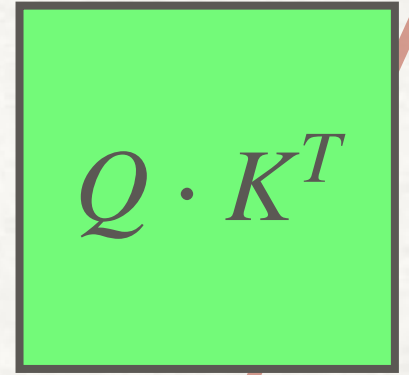
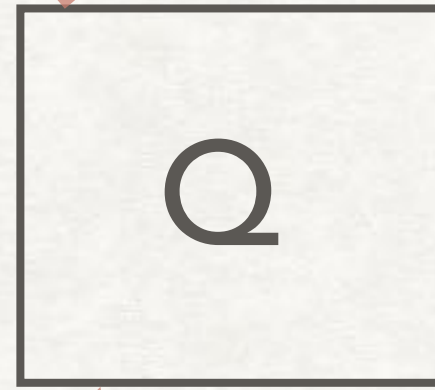# Self Attention

output size = input size

○

d

n



MATMUL

n

$Q \cdot K^T$

Attention
matrix

n

MATMUL

d

n | K

n | Q

n | V

d

$W_Q$

$W_K$

d

$W_V$

n=50
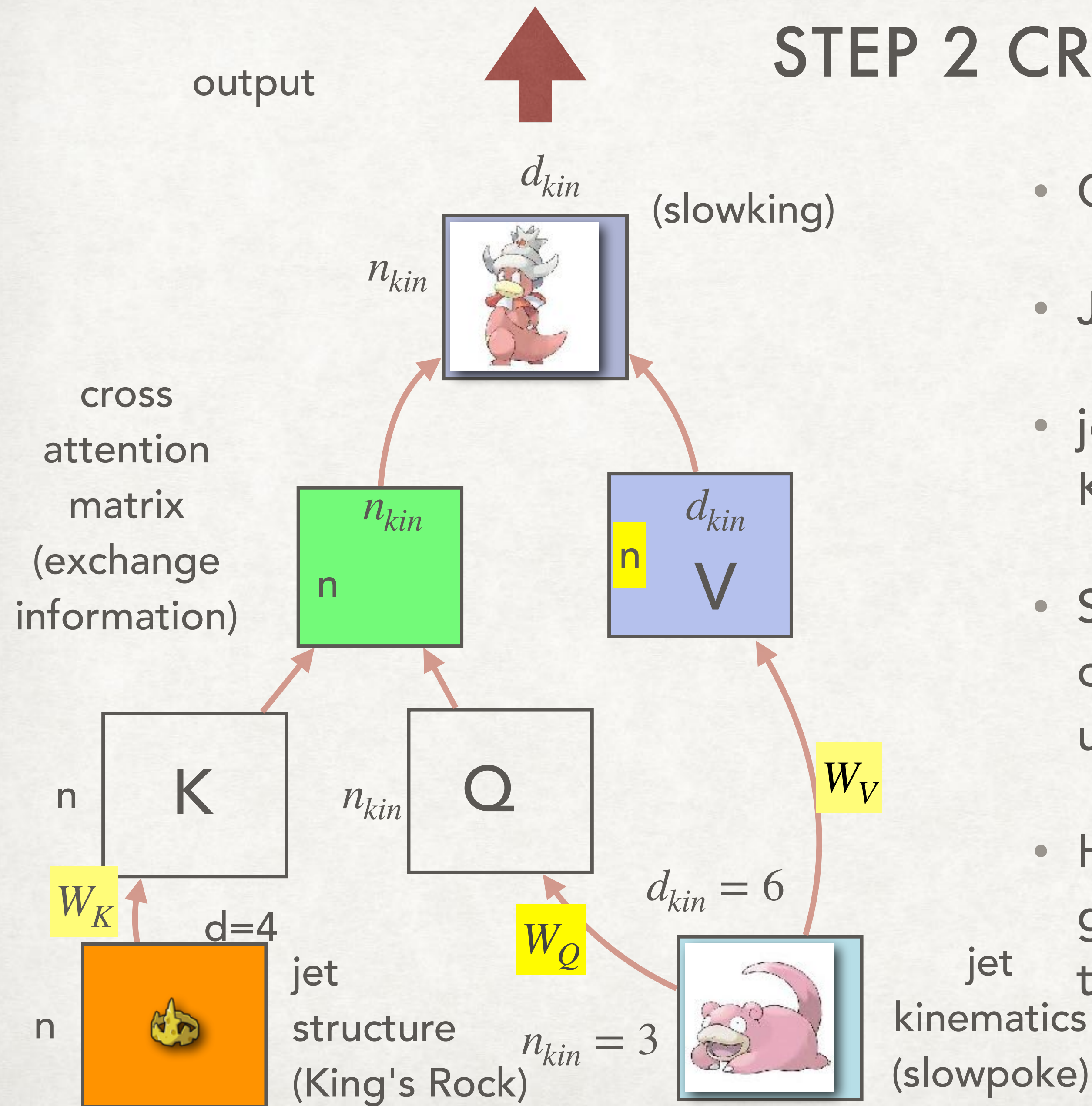
input data
(constituent momentums )

$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

- ATTENTION Matrix mix all features. Higher attention elements indicates important correlations

- transformation  V → V'  does not change the dimension. Structure of V retained for the next transformation.

- We adopt 50x50 self attention for jet and 3x3 self attention for kinematics, with $n_{head} = 5$

$W_v$

# STEP 2 CROSS ATTENTION LAYERS

- Choose cross attention (jet kin) x (jet str. )

- Jet momentum : hard physics of partons Q , V

- jet substructure: parton shower, hadronization K

- Substructure output K and Jet kinematics output Q make attention matrix. The pairs update V (jet Kin)

- High scale feature relevant for classification gives extra weight to the corresponding jets though backward propagation

## Naive approach "simple concatenation"

2311.16674[hep-ph] K. Ban, KC Kong, M Park, S.C. Park



$(a)$

QCD observable
(Local)

Conv. ... flatten

M

CNN branch

Concat.

Kinematic variables
(Global)

M

FCN branch

Deep Learning Score

$(b)$

QCD observable
(Local)

Conv. ... flatten

M

CNN branch

Attention
Layer

2×M      2×M

Concat.

Kinematic variables
(Global)

M

FCN branch

$z$      $z \odot a$

Regularization :

Deep Learning Score



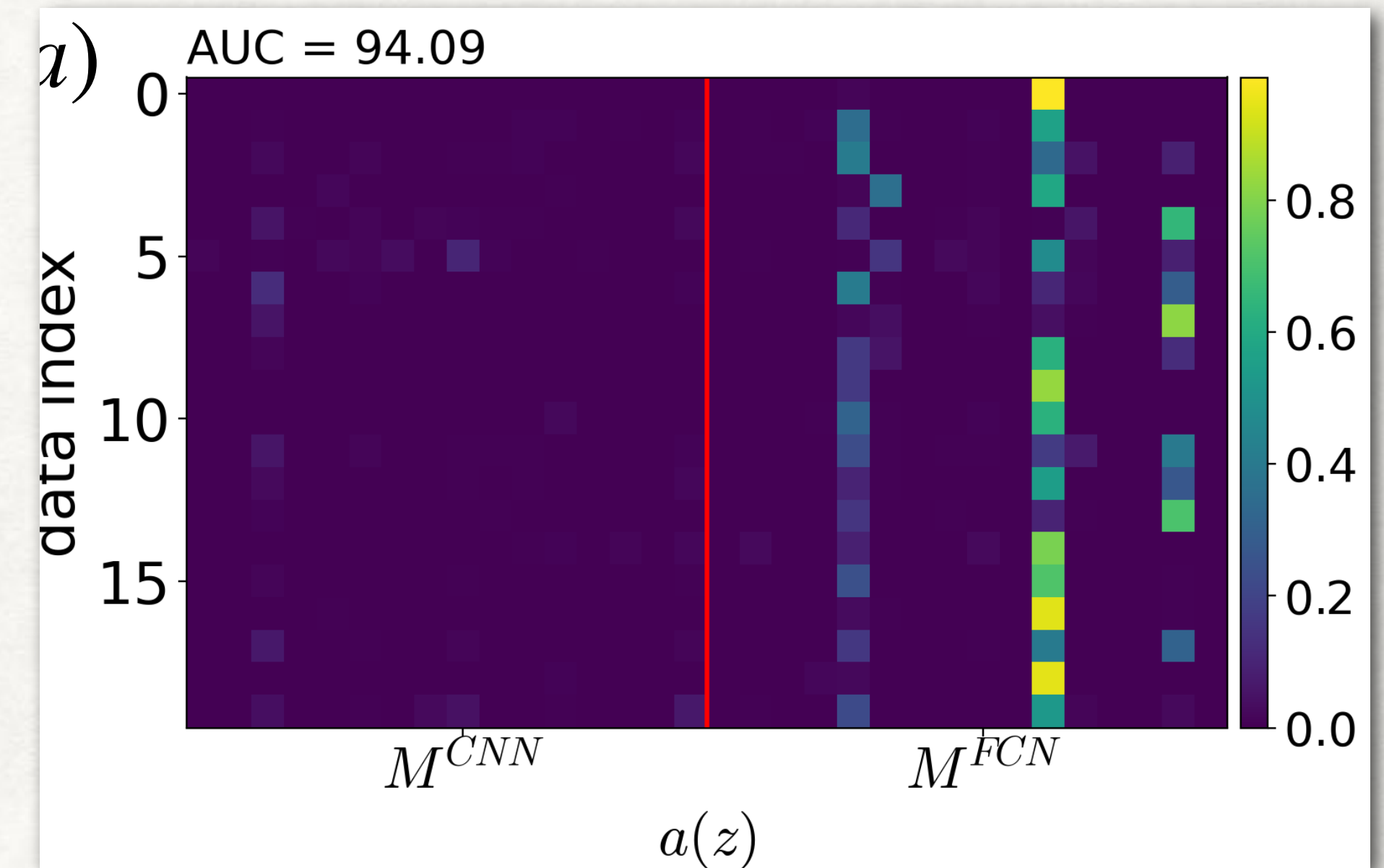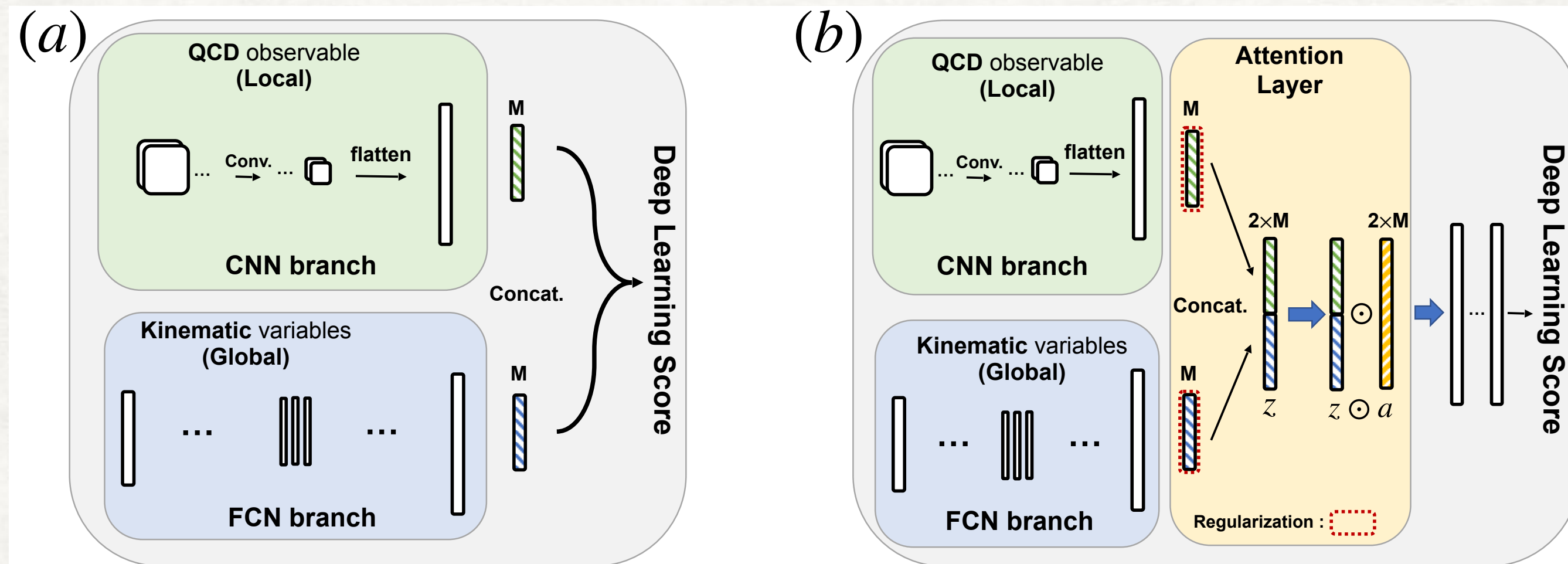$a)$

AUC = 94.09

data index

$M^{CNN}$      $M^{FCN}$

$a(z)$

G. 2. The schematic plots for neural network structures: (a) conventionally used one in previous studies only with concatenation and (b) our
posed one with a regularized attention mechanism.

## (b) self attention matrix of combined information

our network kill this term
and keep off diagonal part only

$$A\,V = \begin{pmatrix} Q(Sub) \times K(Sub) & Q(Kin \times K(Sub) \\ Q(Sub) \times K(Kin) & Q(Kin)\,K(Kin) \end{pmatrix} V = Q(kin)\,K(kin)\,V(kin) + \ldots$$

# PHYSICS

- a jet:

$$\text{P(hadrons in jets | parton or jet )} = P(\{x_i\} \,|\, y)$$

- a fatjet  or a jet with substructure

$$P(\{x_i\} \,|\, \{y_\alpha\})$$

- two fatjets in an event

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

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

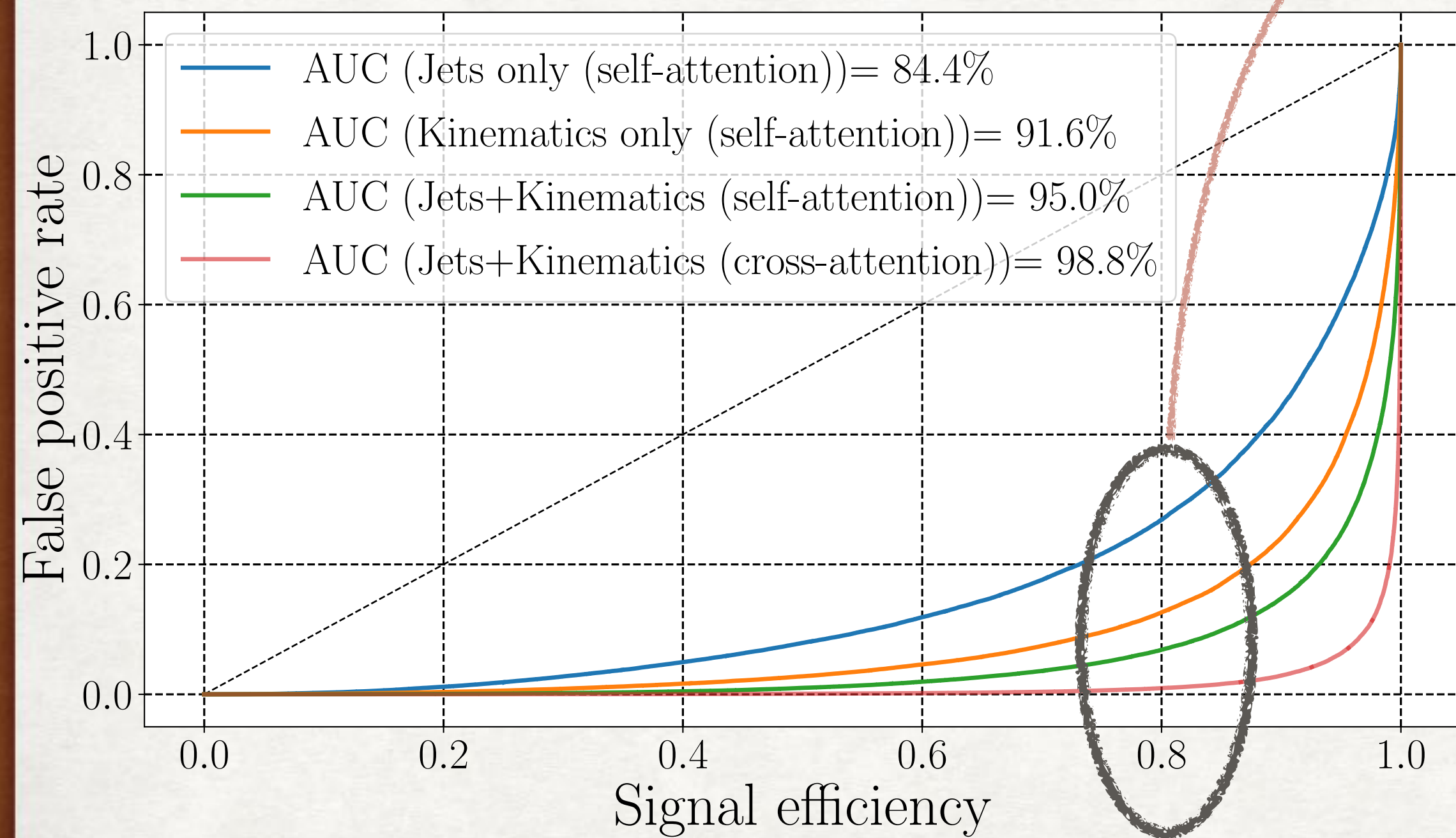cross attention                                        jet kinematics
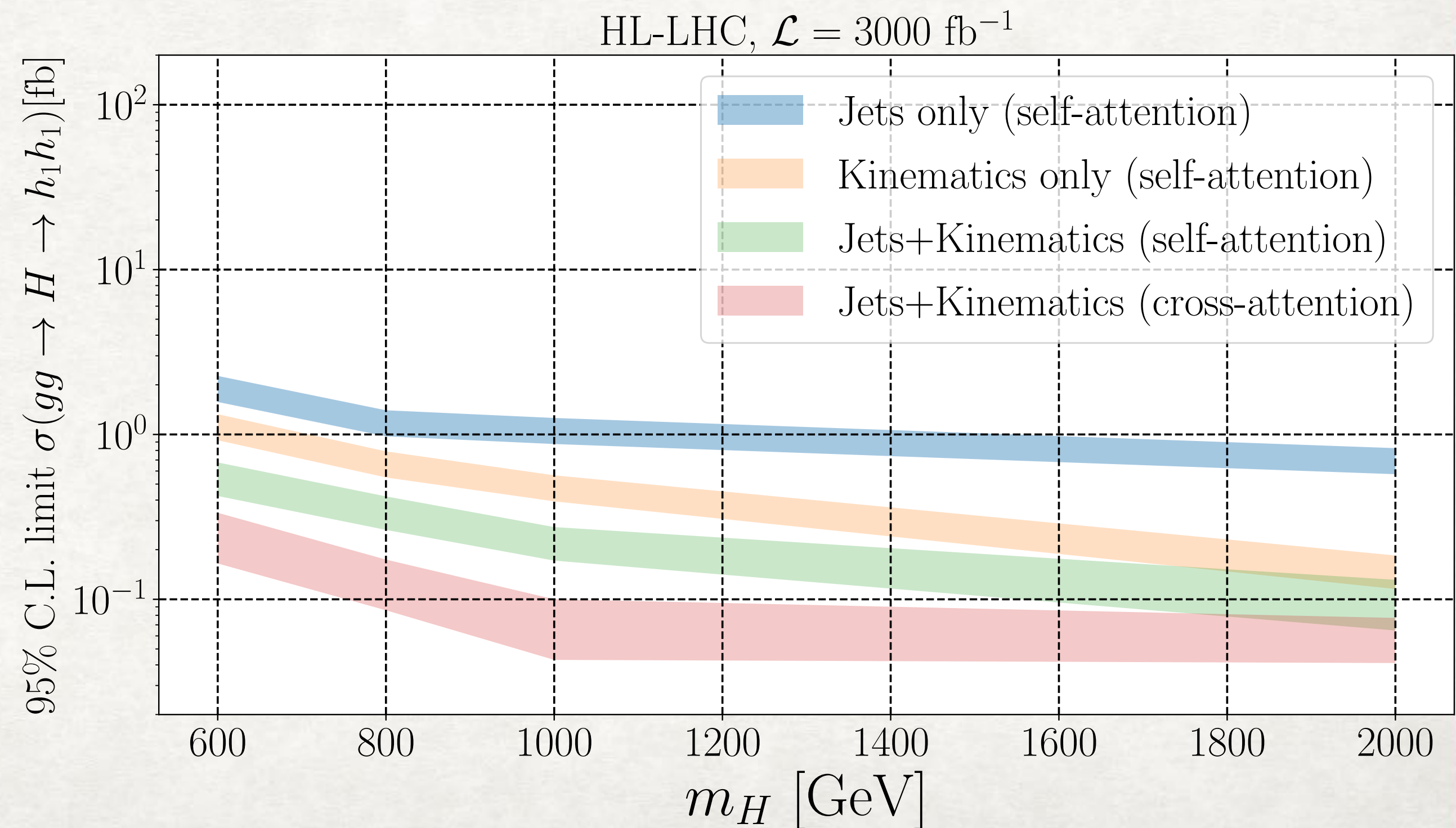
# IMPROVEMENT USING CROSS ATTENTION

green: self attention of Jet str. and Kin
→ concatenate and MLP
red line : cross attention

factor 5 improvement at the same acceptance.
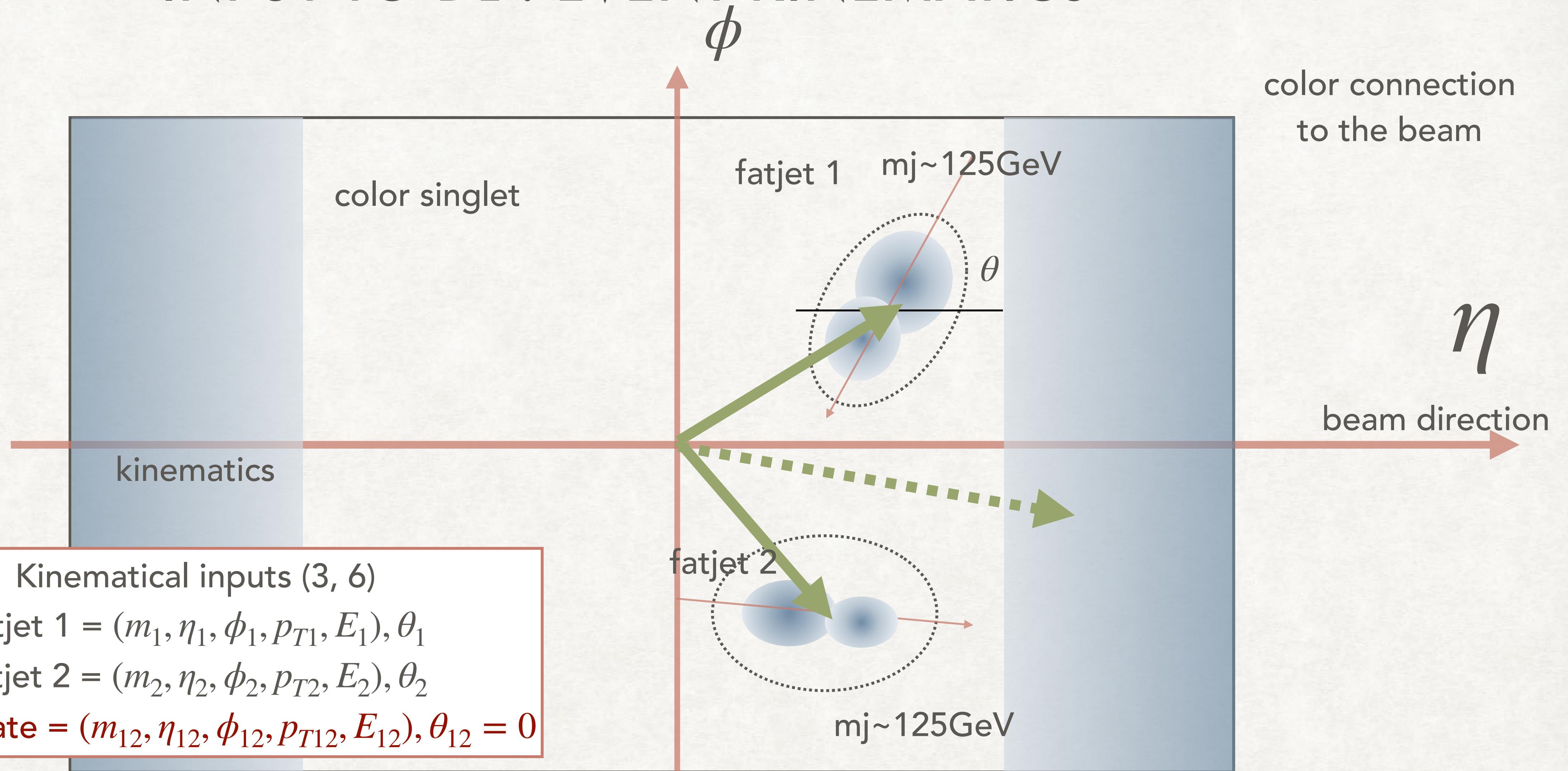
Simple estimation of the upper limits



AUC (Jets only (self-attention))= 84.4%
AUC (Kinematics only (self-attention))= 91.6%
AUC (Jets+Kinematics (self-attention))= 95.0%
AUC (Jets+Kinematics (cross-attention))= 98.8%

Cross attention improve the rejection
efficiently

# INPUT TO DL : EVENT KINEMATICS

color connection
to the beam

color singlet

fatjet 1    mj~125GeV

$\theta$

$\eta$

beam direction

kinematics

fatjet 2

Kinematical inputs (3, 6)

fatjet 1 = $(m_1, \eta_1, \phi_1, p_{T1}, E_1), \theta_1$

fatjet 2 = $(m_2, \eta_2, \phi_2, p_{T2}, E_2), \theta_2$

H candidate = $(m_{12}, \eta_{12}, \phi_{12}, p_{T12}, E_{12}), \theta_{12} = 0$

mj~125GeV

NOTE : "5 inputs for 4 momentum" , H candidate momentum as sum of two fat jets, add θ,

# ROLE OF $\theta$

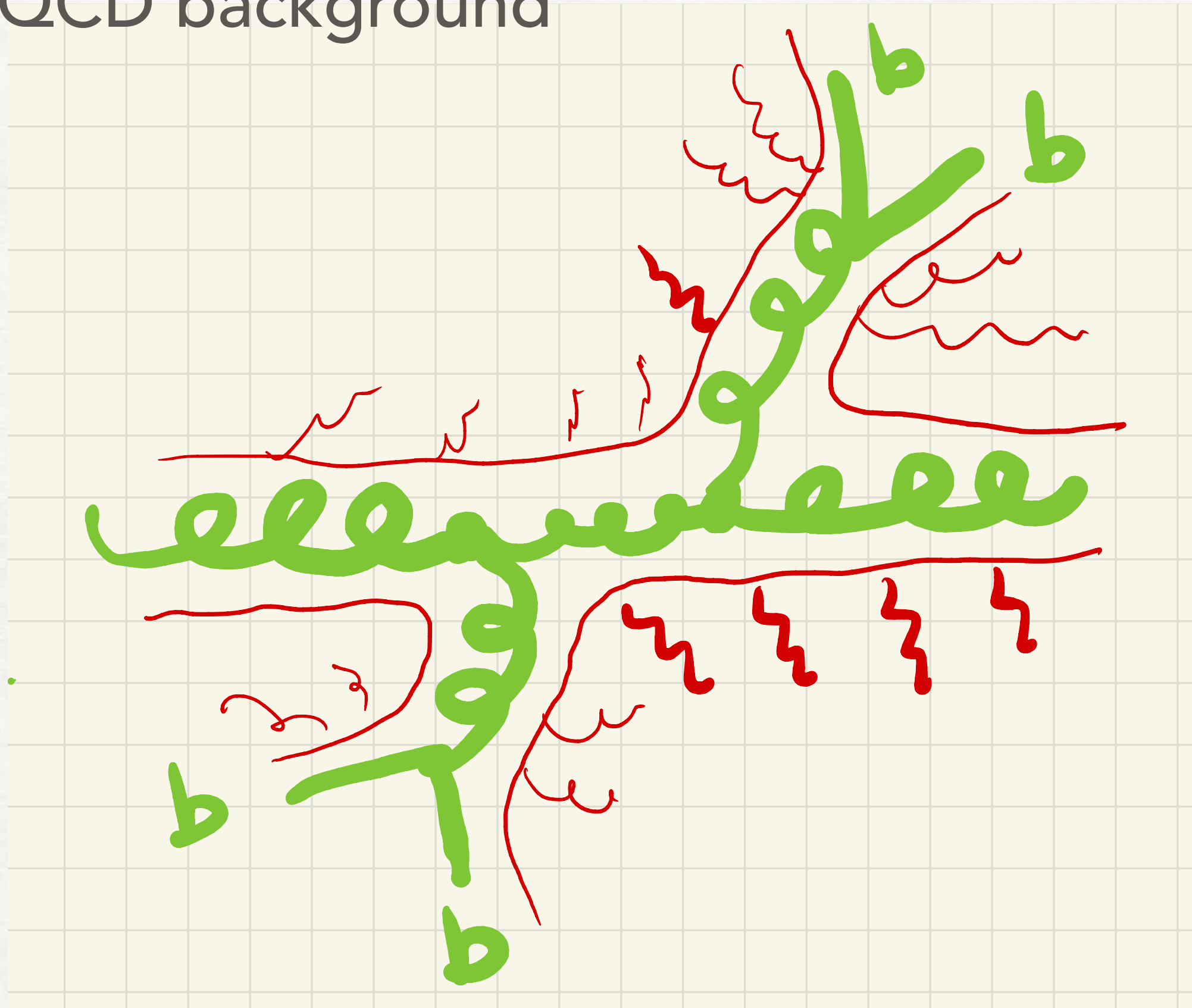|  | Kinematics | Kin +θ | jet str.+kin | jet str +Kin + $\theta$ |
|---|---|---|---|---|
| ROC | **91.01%** | **91.6** | 97.23-98.16 | 98.68-99.28 |



adding rotation angle $\theta$ improve classification when both jet str and kinematical information available.

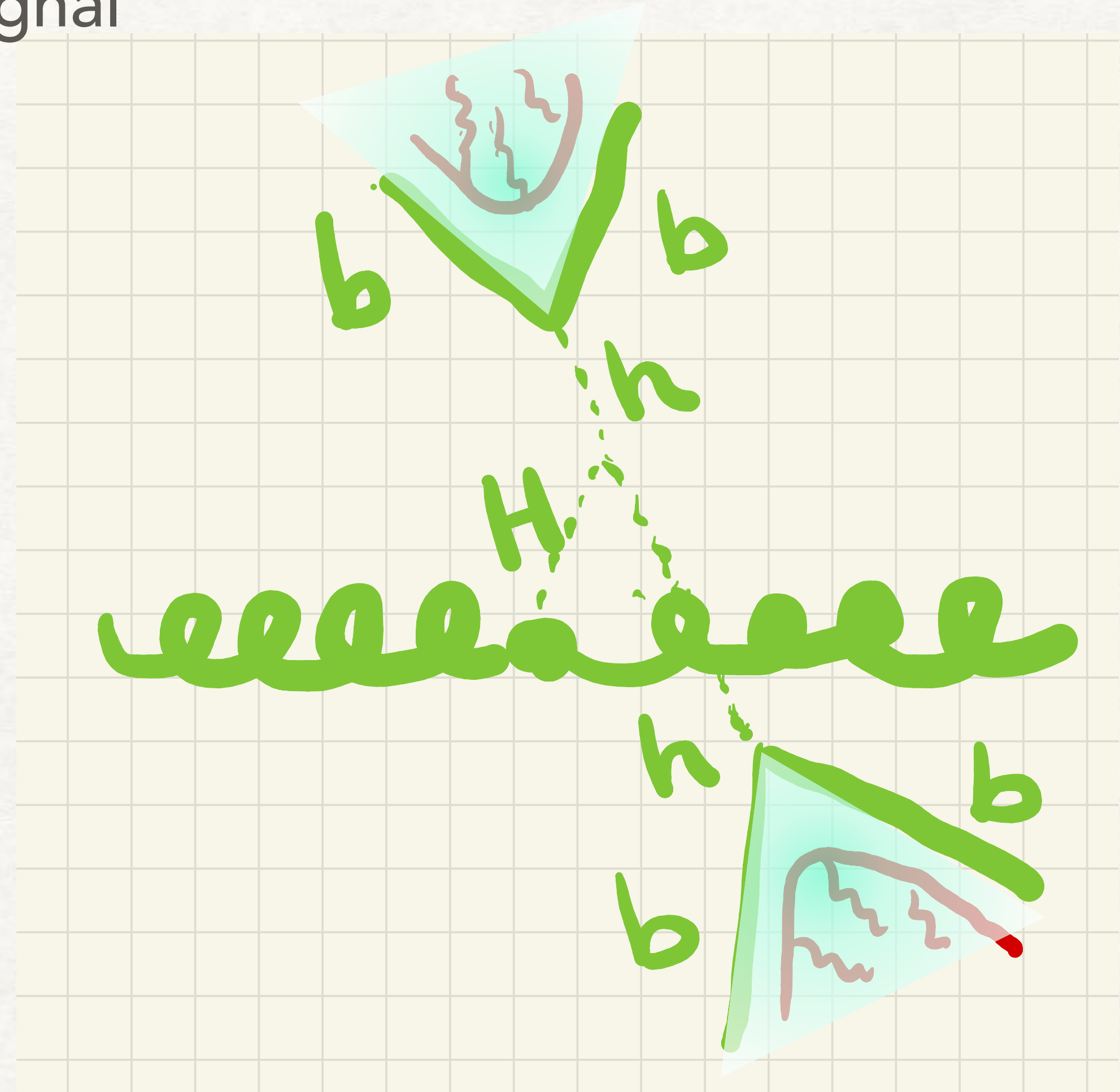We are working on to idententify the origin. (color connection? momentum resolution? )

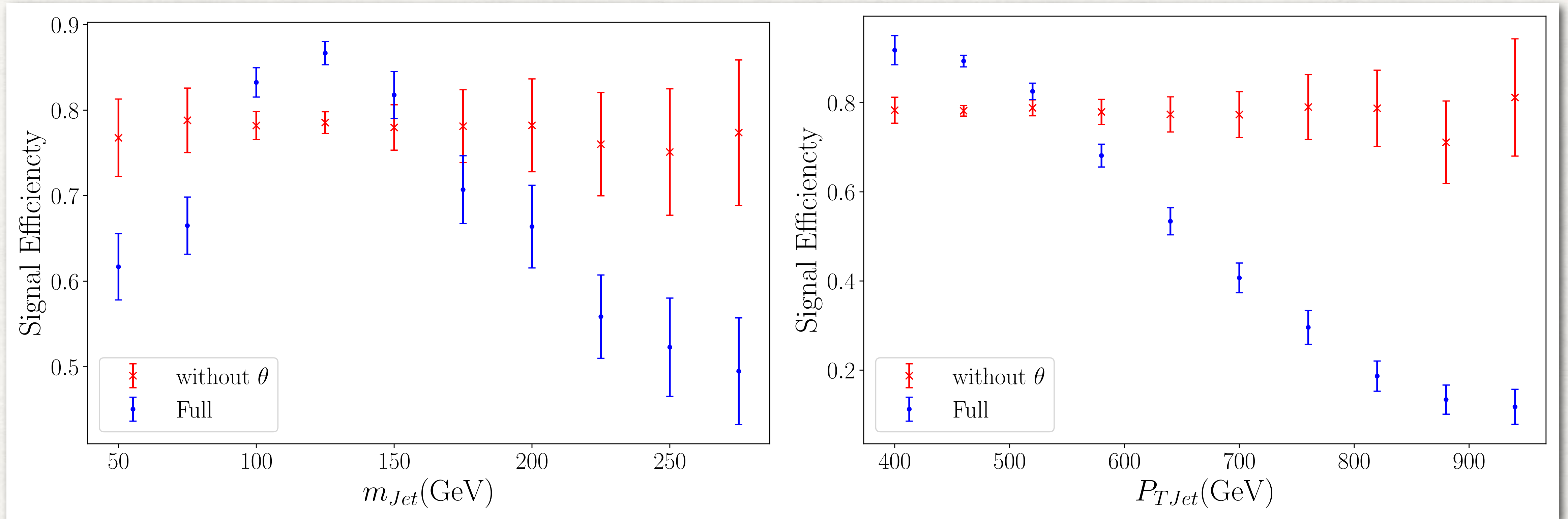# event color structure



QCD background

signal

For QCD and top event, fatjets are likely
color connected to
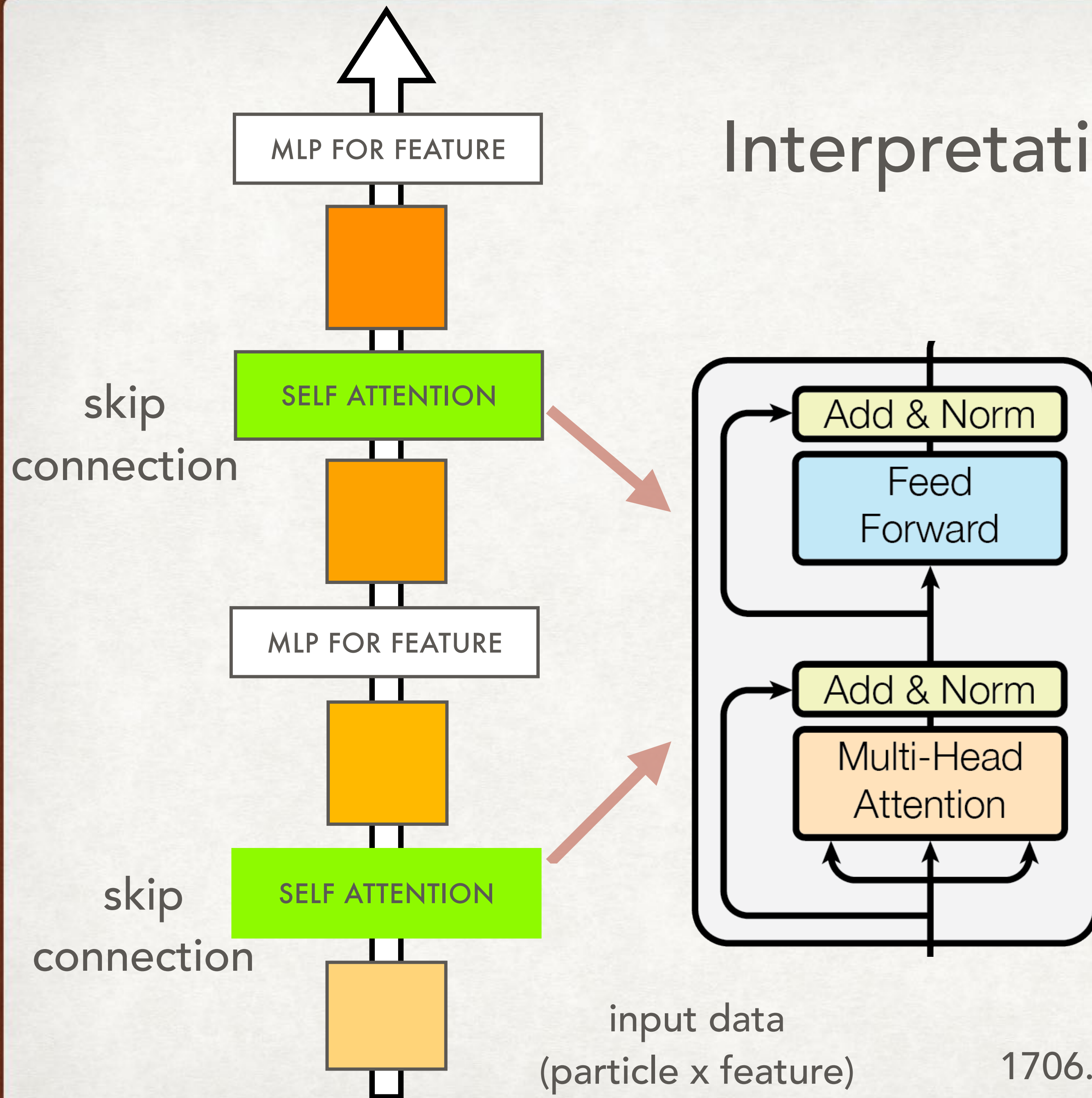the other activities of the event

Higgs bosons are color isolated.

better selection of Higgs mass

rejecting high PT events

# Interpretation and Skip Connection



- Deep Learning suffers low interpretability and it is always annoying.

- skip connection of attention blocks helps connecting input data to extracted feature(transformed quantity) in **some level.**

1706.03762 Vaswani et all "Attention is all you need

self attention map

axis: odering  of modified particles

First few "particle" token express
Higgs nature efficiently

Cross attention map:Particle in the jet (50)
and parent particle (3)

sum of fatjet momenta capture signal

Signal

Cross Attention maps

BG



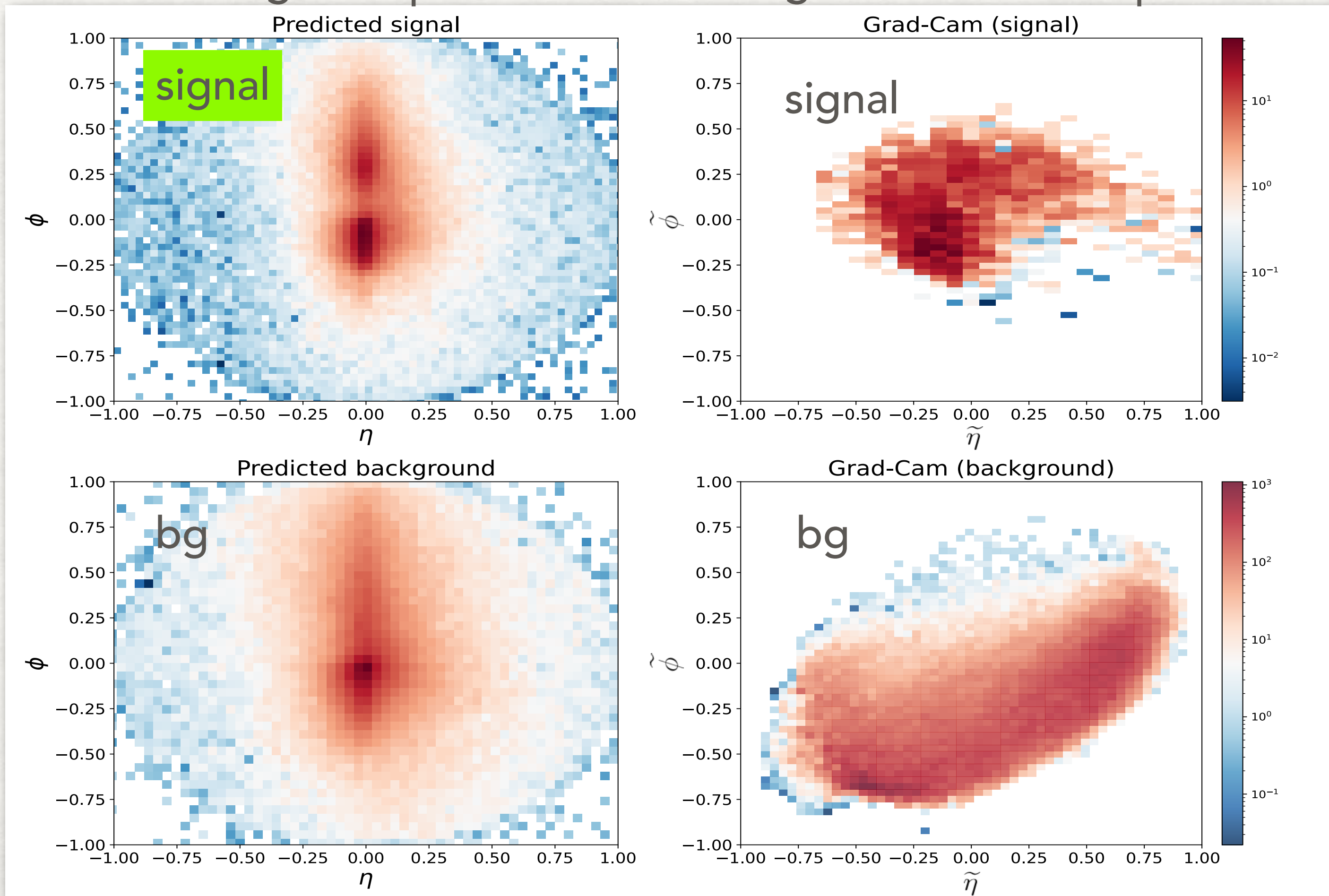maybe number of averaged particles are different

# GRAD-CAM (1610.02391)

- Output of last attension layers (some correlation with original inputs)

5000 signal events

Original inputs      grad-cam heatmap



Y : class scaore

F:output from last attention layer

$\widetilde{\eta}, \widetilde{\phi}$    transformed corrdinate

$$\alpha_k(\widetilde{\eta}, \widetilde{\phi}) = \frac{1}{Z} \sum \frac{\partial Y_c}{\partial F_k(\widetilde{\eta}, \widetilde{\phi}, \widetilde{p_T})}$$

$$\text{Grad-CAM}(\widetilde{\eta}, \widetilde{\phi}) = \frac{1}{k} \sum_k \alpha_k(\widetilde{\eta}, \widetilde{\phi}) F_k(\widetilde{\eta}, \widetilde{\phi}, \widetilde{p_T})$$

Still see some connection

# TAKEAWAYS

- **use "cross attention"** when you combine the "high scale information" to the "low energy scale", because cross attention layer gives extra emphasis to the information linked to the high energy kinematics.

- **skip connection and Interpretation** : Skip connection helps to maintain some connection to the inputs

- **More Physics:** Heavy particles decay into colored particles (discovery, spin, color structure? ) Cross attention network probably more useful to resolve correlation of jet structures.

- Result looks very good to me and I am still worrying about bugs…

# NEED TO BE IMPROVED

- Current GPU requirement: **2 x NVIDIA RTX A6000 (48GB) with 80% and 30% utilization in tensor flow mirror strategy.** 96% consumption /card 20min/ training.

- We definitely have to change "jet substructure part" to simpler one, keeping cross attention structure(this part is generic)

- Ex: "Modulated Network of HL variables"

  - QCD vs top, Amon Furuichi(Nagoya), Sung Hak Lim(Rutgers) and M. Nojiri arXiv 2312.11760[hep-ph] work as good as Particle Transformer.

- ......  but are they robust for color connection?
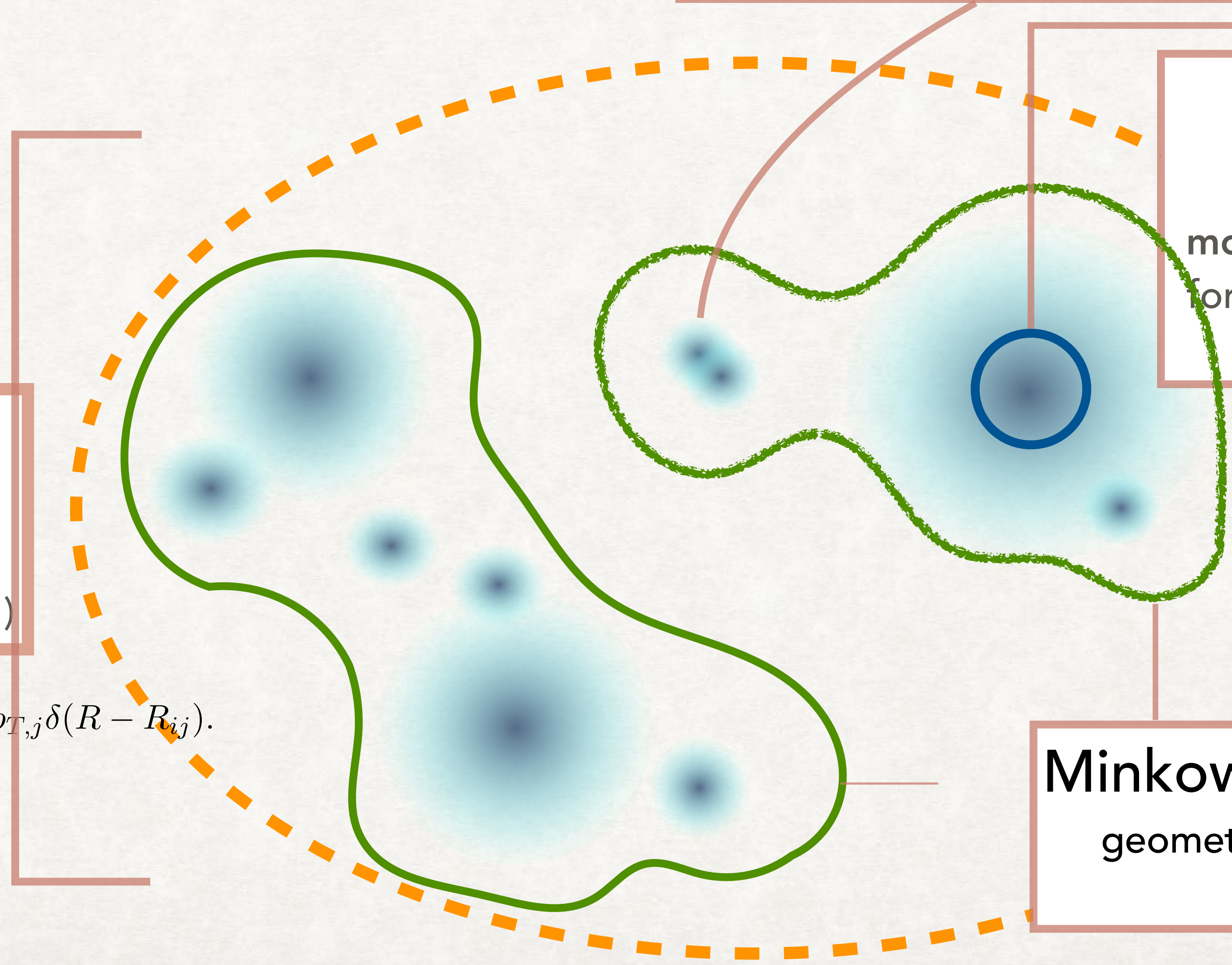
# BACK UP SLIDES

JET High Level variables

pt distribution of constituents

Subjet
Localized sampling
**momentum and counting**
for various angular sccale
R=0.1, 0.2, 0.3

Jet spectrum
two point **Energy correlation**
(unlocalized sampling )
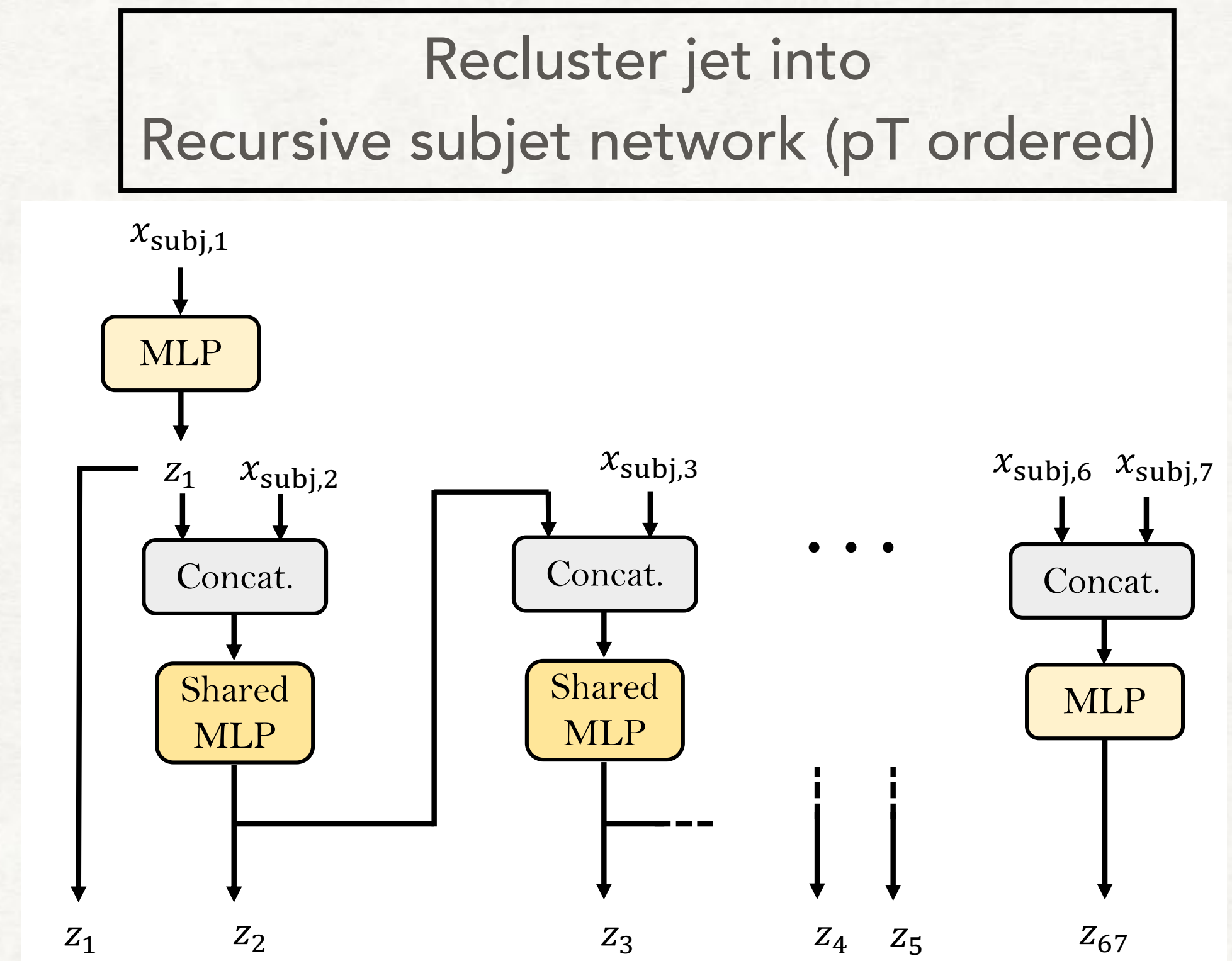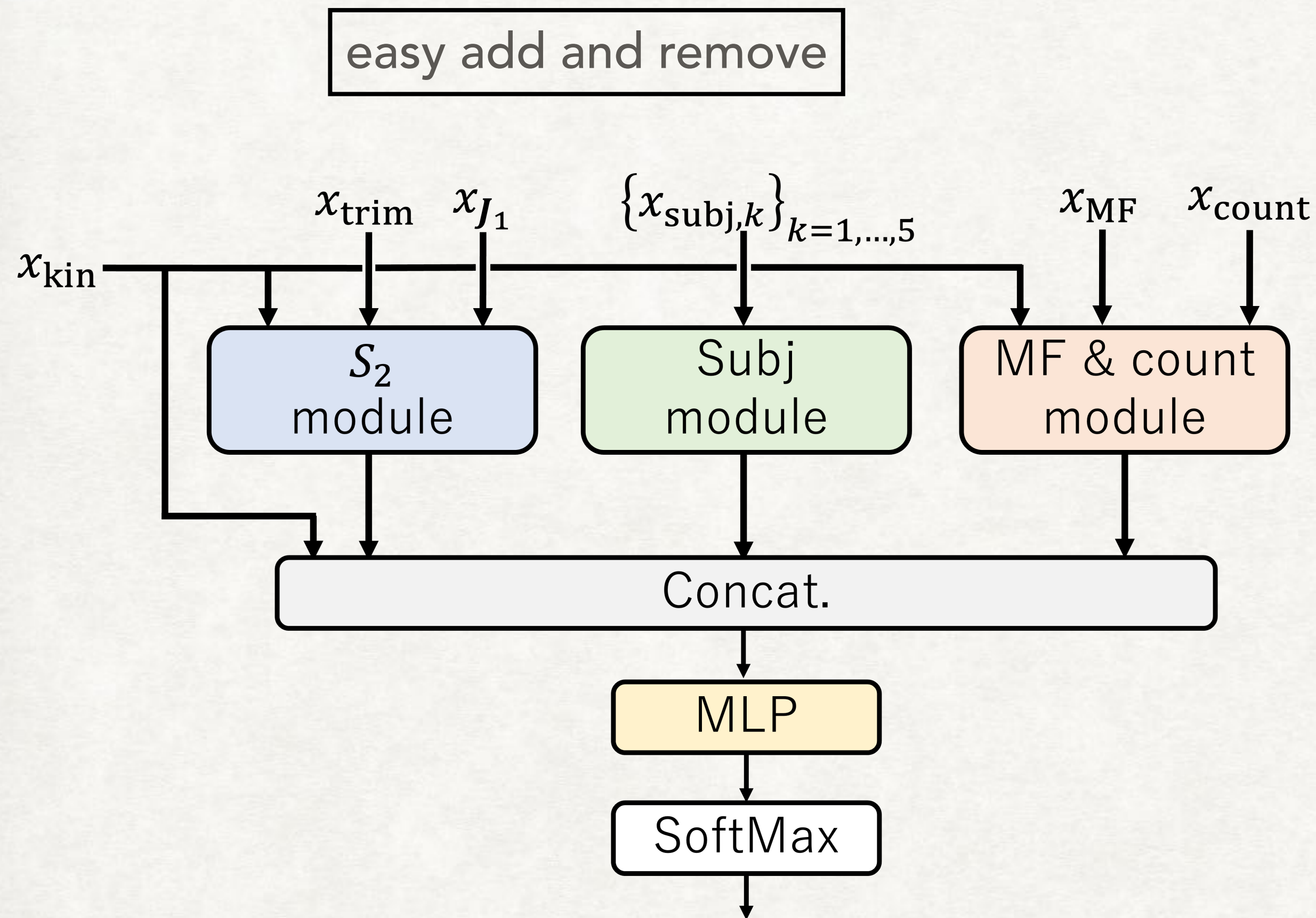
$$S_{2,ab}(R) \overset{\text{def}}{=} \sum_{i \in a} \sum_{j \in b} p_{T,i} p_{T,j} \delta(R - R_{ij}).$$

Minkowski Functionals
geometry of jet cosntituent distribution
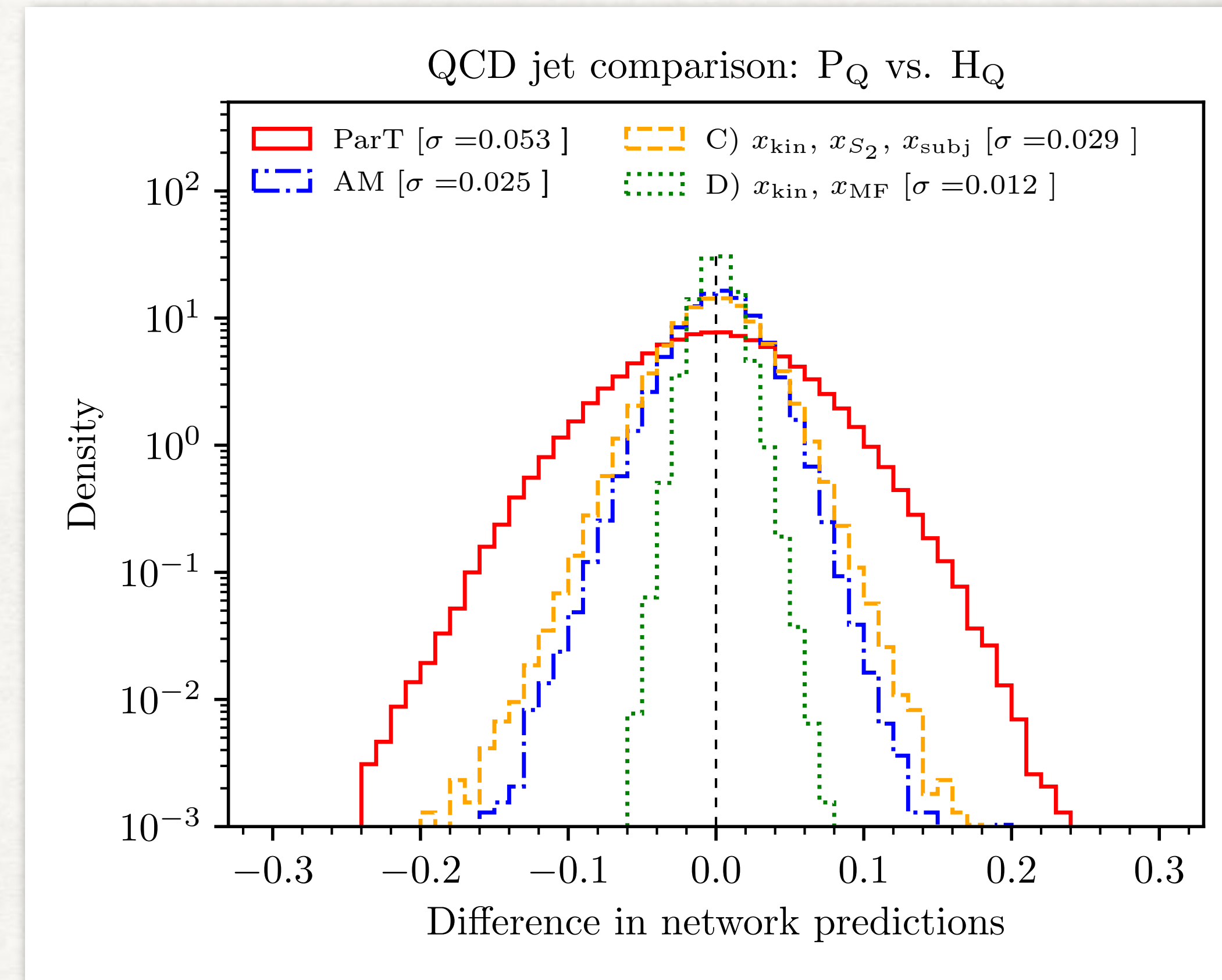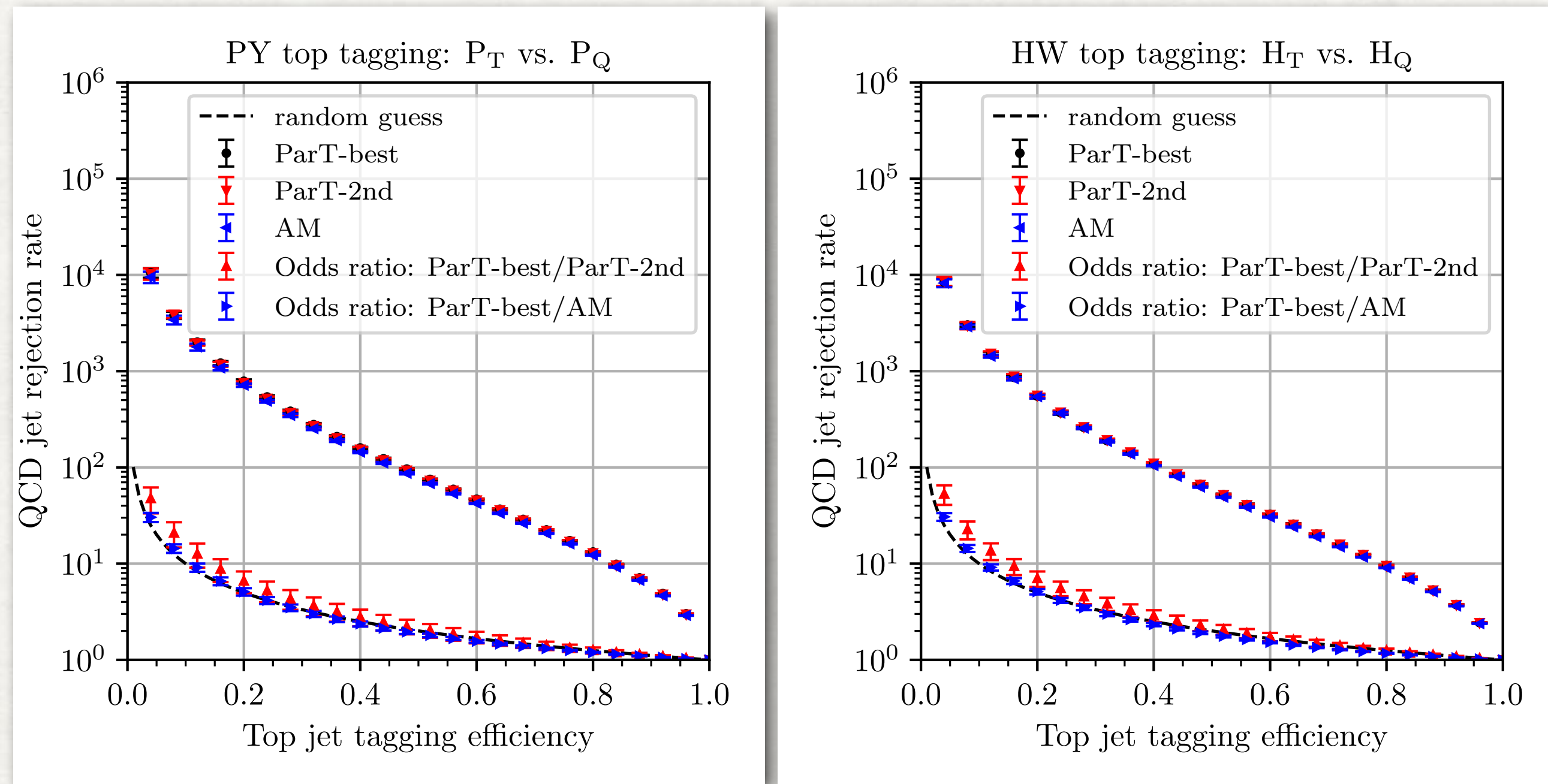
# NETWORK USING HL INPUTS (ANALYSIS MODEL=AM)



(b) A schematic diagram of subjet recursive module.

⭐input: subjet with multiple cone size
(R=0.1, 0.2, 0.3)  =information of clustering
⭐Shared MLP for 2nd to 5th subjet to reduce paramters

PY top tagging: $P_T$ vs. $P_Q$

HW top tagging: $H_T$ vs. $H_Q$

QCD jet comparison: $P_Q$ vs. $H_Q$

- ParT [$\sigma = 0.053$]
- AM [$\sigma = 0.025$]
- C) $x_{kin}, x_{S_2}, x_{subj}$ [$\sigma = 0.029$]
- D) $x_{kin}, x_{MF}$ [$\sigma = 0.012$]

QCD jets: $P_Q$ vs. $H_Q$

Top jets: $P_T$ vs. $H_T$

AM model : **1GB GPU memory** on GeForce 1080Ti GPU(11.3TFLOPS) with 35% GPU utilization. need lots of preprocessing

ParT: **14GB GPU memory** RTX A6000( 38.7TFLOPS) GPU utilization 95%