

streamlined jet tagging network assisted  
by jet prong structure  
–role of cross attention–

Mihoko Nojiri(IPNS, KEK), with Ahmed Hammad and Stefano Moretti

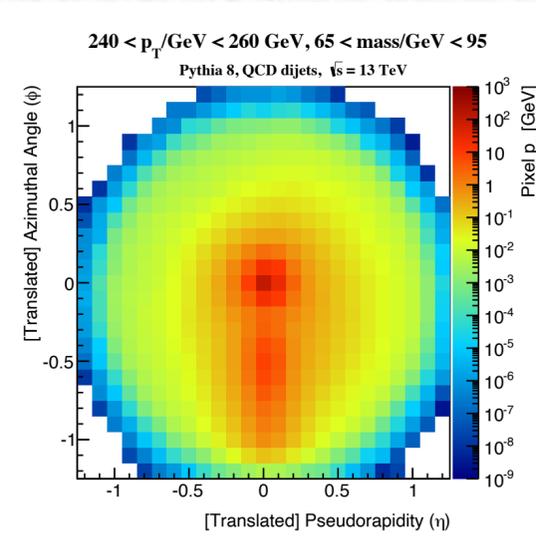
arXiv 2401.00452 JHEP 03(2024) 144

Mihoko Nojiri with Ahmed Hammad

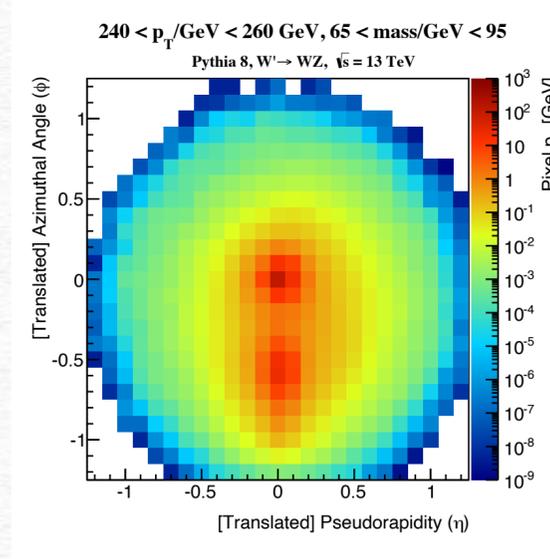
arXiv 2404.14677 JHEP 06 (2024) 176

# Jet classification using ML

QCD jet  
(in W mass region)

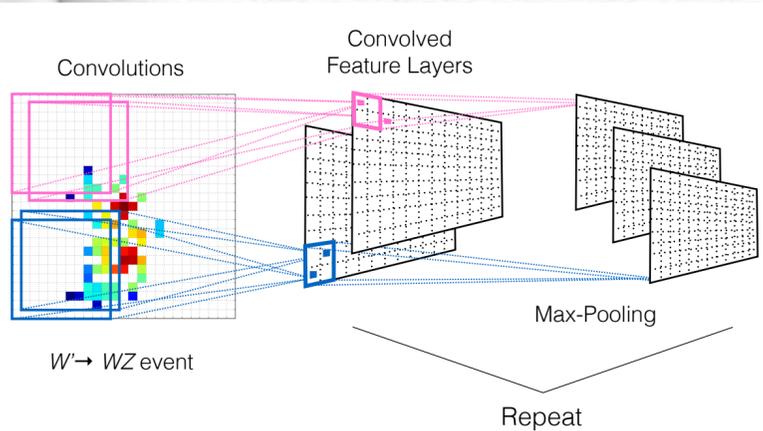


W jet



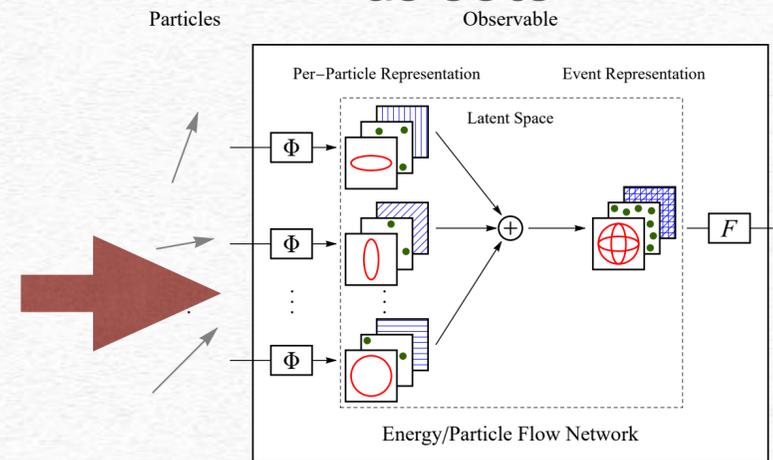
from Schwartzman et al  
<https://iopscience.iop.org/article/10.1088/1742-6596/762/1/012035>

Jet Image



CNN Oliverira et al  
(1511.05190)

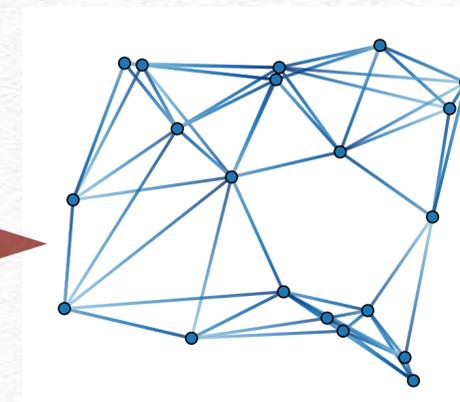
as sets



permutation invariance  
(Energy Flow Network and  
Particle Flow Network 1810.05165)

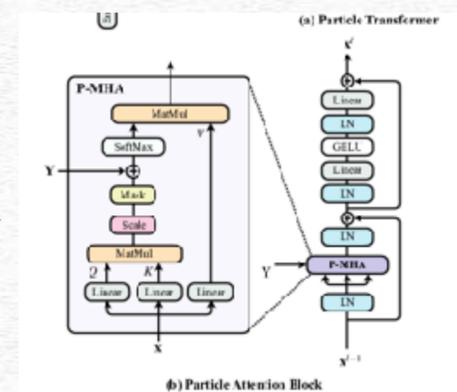
Bogatskiy et al PELICAN (2211.00454)

as graphs



sparse data  
1902.08570 Particle Net  
Dreyer et al LundNet (1807.04758)  
Gong et al LorentzNet(2201.08187)

transformer



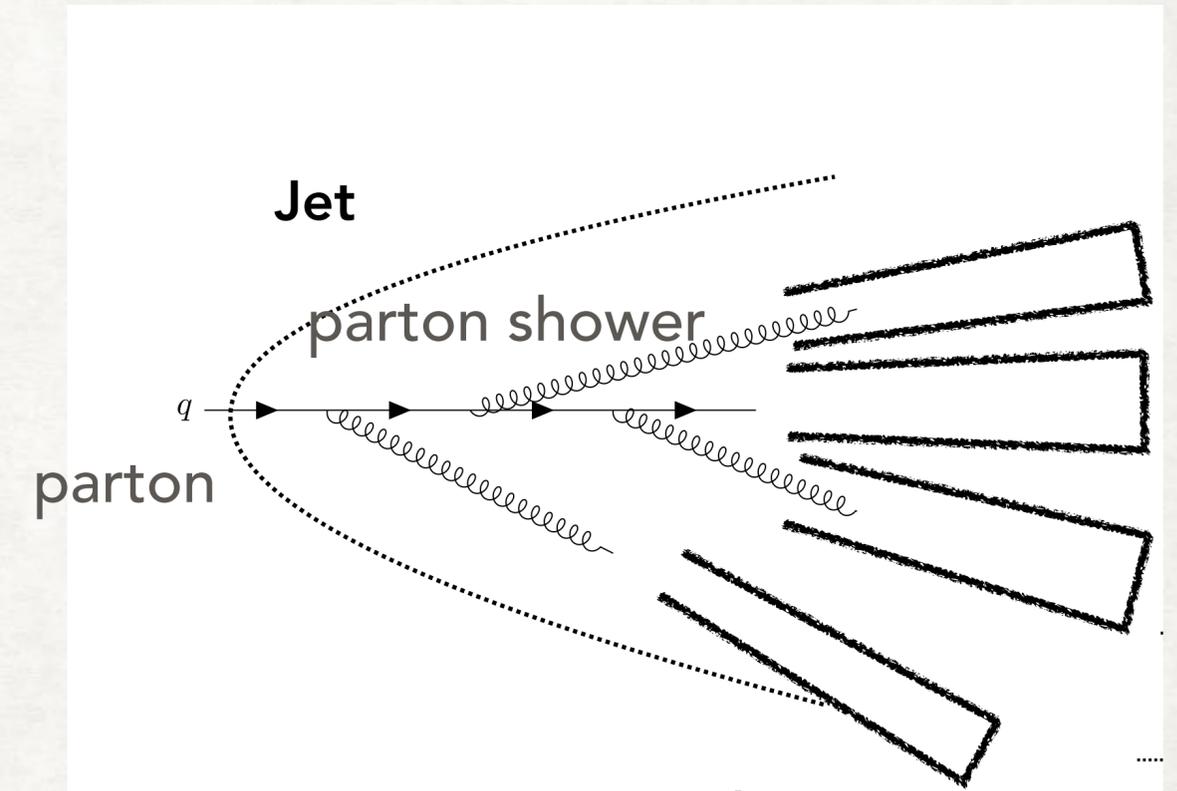
building key  
and query  
2202.03772



# BUT... PHYSICS BEFORE THE NETWORK

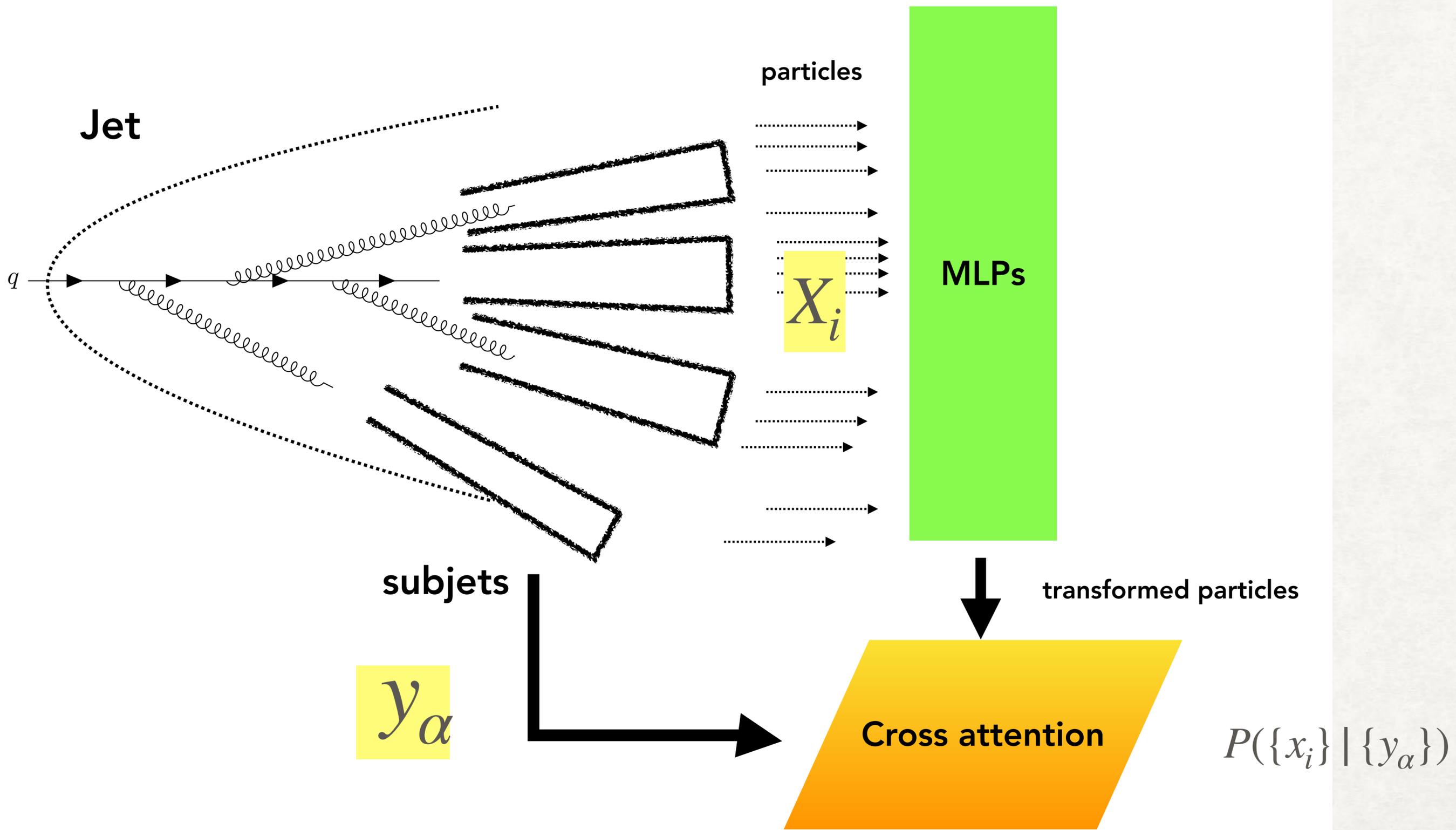
## "Physics SCALE"

- Hard Process = Partons  $y$
- Factorization
  - a jet:  $P(\text{hadrons in jets} \mid \text{parton } y) = P(\{x_i\} \mid y)$
  - jet with substructure  $P(\{x_i\} \mid \{y_\alpha\})$
- Maybe several fatjets in an event



$$P(\{x_i\}, \{x'_j\}, \{y_\alpha\}, \{y'_\beta\}) \sim P(\{x_i\} \mid \{y_\alpha\}) P(\{x'_j\} \mid \{y'_\beta\}) P(\{y_\alpha, y'_\beta\})$$

Why don't you construct the network focusing on QCD scale structure

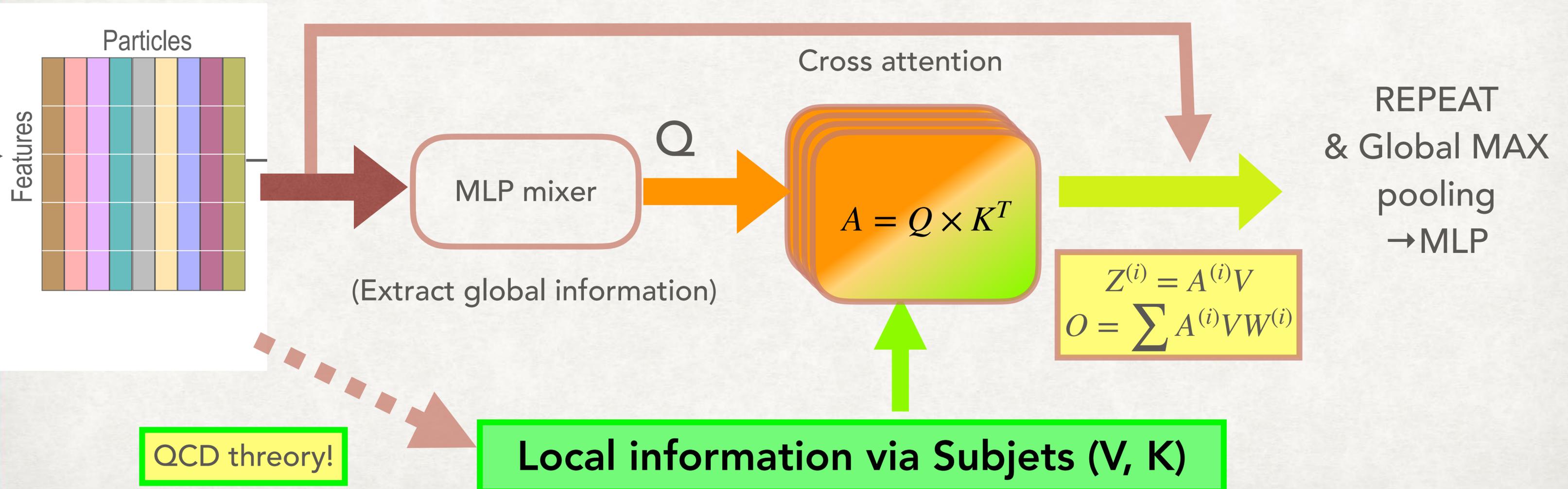


# Cross attention to focus on the P(h| (sub)jet)

ATTENTION → CROSS Attention for P(h| subjects) estimation

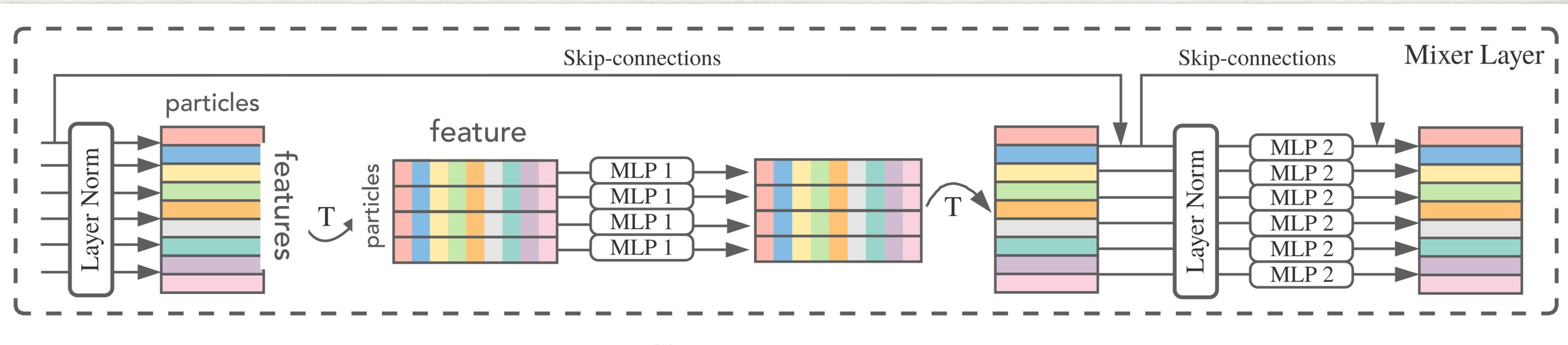
input X

skip connection  $\tilde{X} = X + O$



# MLP MIXER

The mixer layer has **only two MLP** that mix both features and Particle tokens: focus on global feature.



MLP 1 :mix feature only acts for all particles

MLP 2: mix particles acts for all features

transformer-like (add any information, apply repeatedly )

"subject information" take care cluster information

# Performace comparable to Particle Transformer but much faster and lighter

Models	AUC	R50%	#Parameter	Time (GPU%)
ParT	0.9858	413+-16	2.14M	612
Mixer+subjet (CA)	0.9856	392+-5	86.03K	33
(AK)	0.9854	375+-5	86.03K	33
(HDBSCAN)	0.9859	416+-5	86.03K	33
LorentzNet	0.9868	498+-18	224K	
PELICAN (Lorents Invariance)	0.9869	-	45K	-

\*Subjet cone size  $R=0.3$

\*HDBSCAN is algorithm without distance measure

Performace comparable to Particle Transformer but much faster and lighter

Models	AUC	R50%	#Parameter	Time (GPU%)
ParT	0.9858	413+-16	2.14M	612
Mixer+subjet (CA)	0.9856	392+-5	86.03K	33
(AK)	0.9854	375+-5	86.03K	33
(HDBSCAN)	0.9859	416+-5	86.03K	33
LorentzNet	0.9868	498+-18	224K	-
PELICAN (Lorents Invariance)	0.9869	-	45K	-

SMALL SIZE

612

FAST

\*Subjet cone size  $R=0.3$

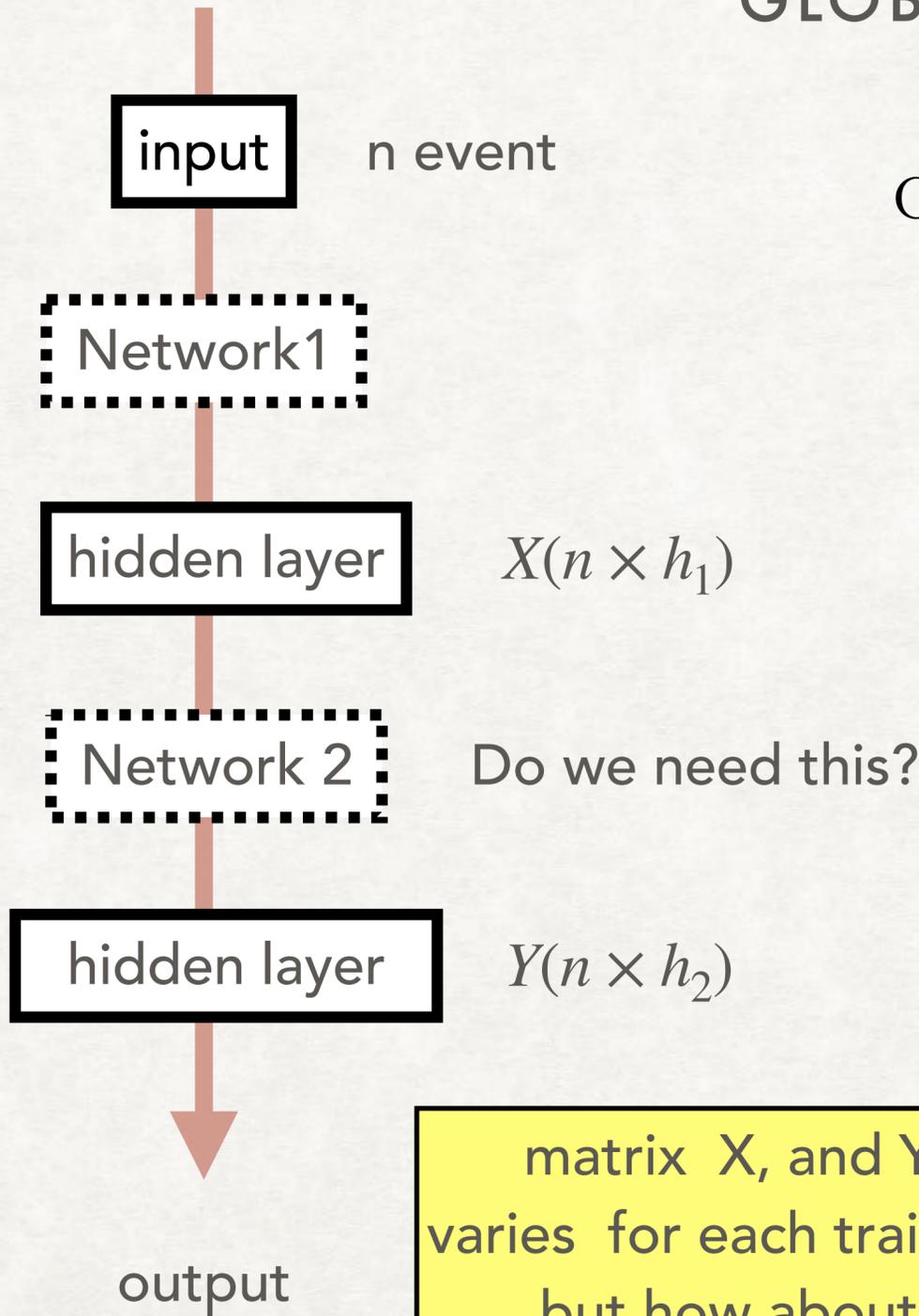
\*HDBSCAN is algorithm

HIGH PERFORMANCE WITHOUT USING LORENTS INVARIANCE

without distance measure

# INTERPRETATION USING CKA SIMILARITY

GLOBAL(MIXER) AND LOCAL(SUBJET)



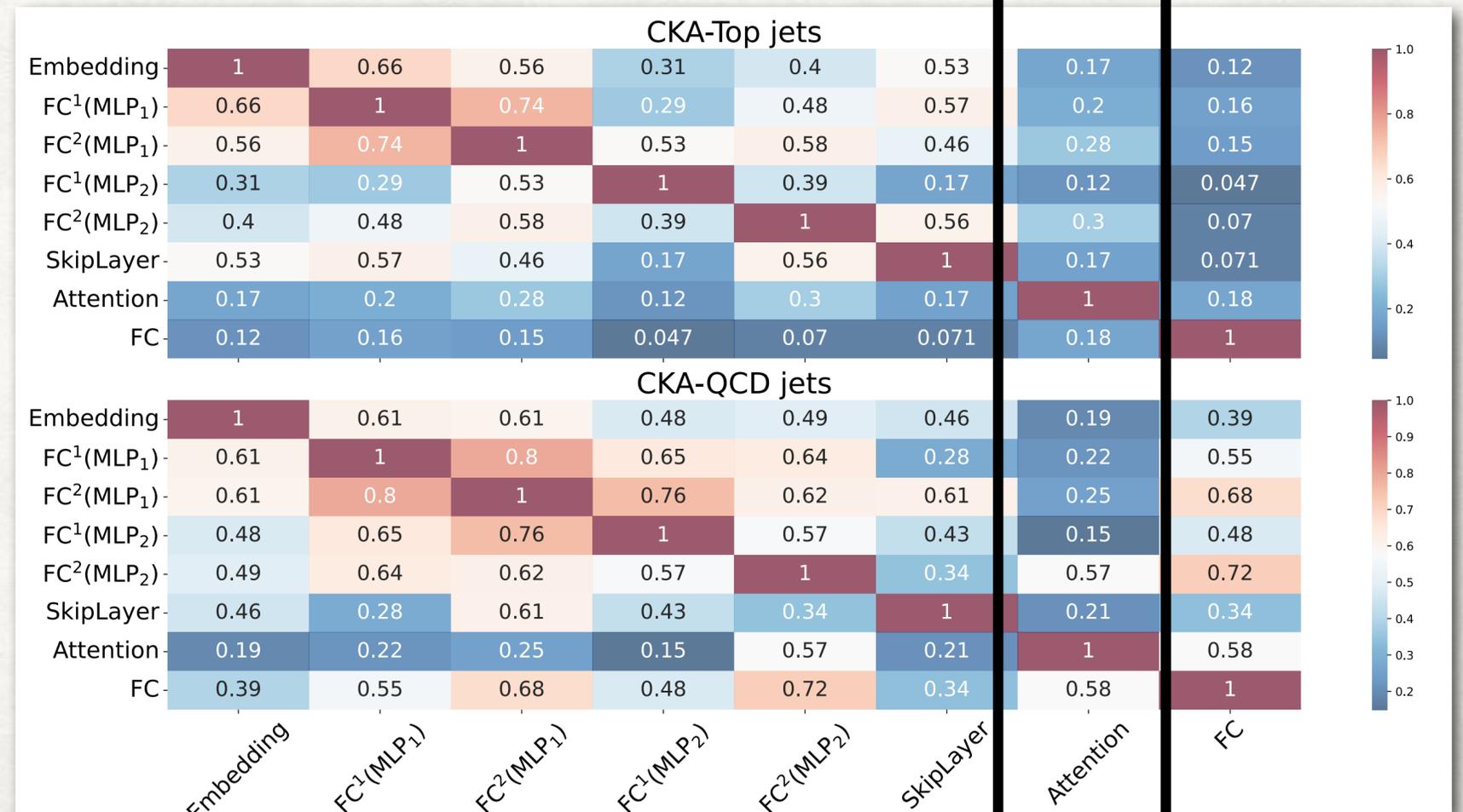
$$CKA(M,N) = \frac{HSIC(M,N)}{\sqrt{HSIC(M,M)HSIC(N,N)}},$$

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

$$H = \delta_{ij} - \frac{1}{d}$$

1 if they are same (no improvement)

Efficient!



matrix X, and Y varies for each training but how about  $M = XX^T, N = YY^T$

# TOWARD GLOBAL EVENT ANALYSIS

A HAMMAD S. MORETTI MN JHEP 03 (2024) 144

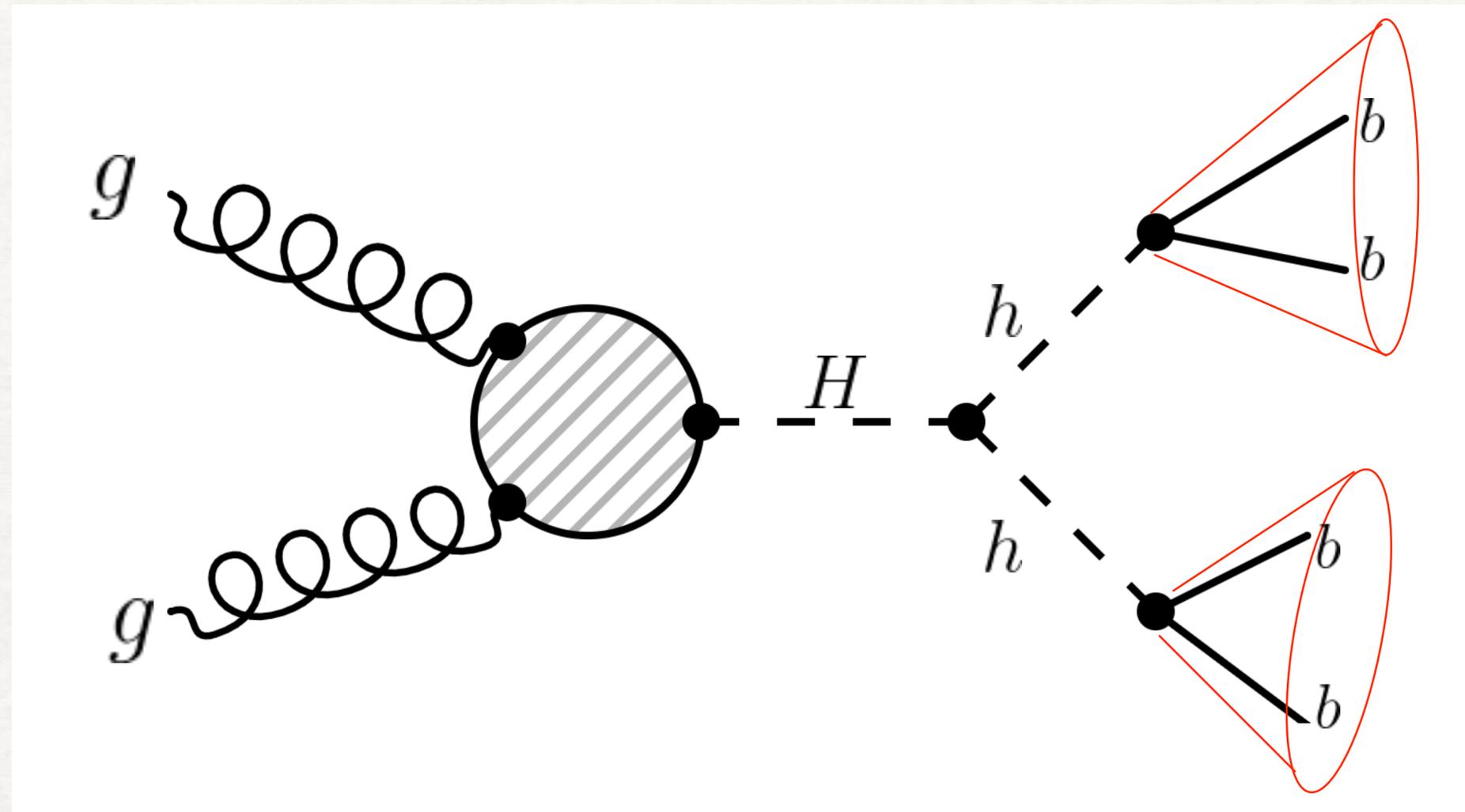
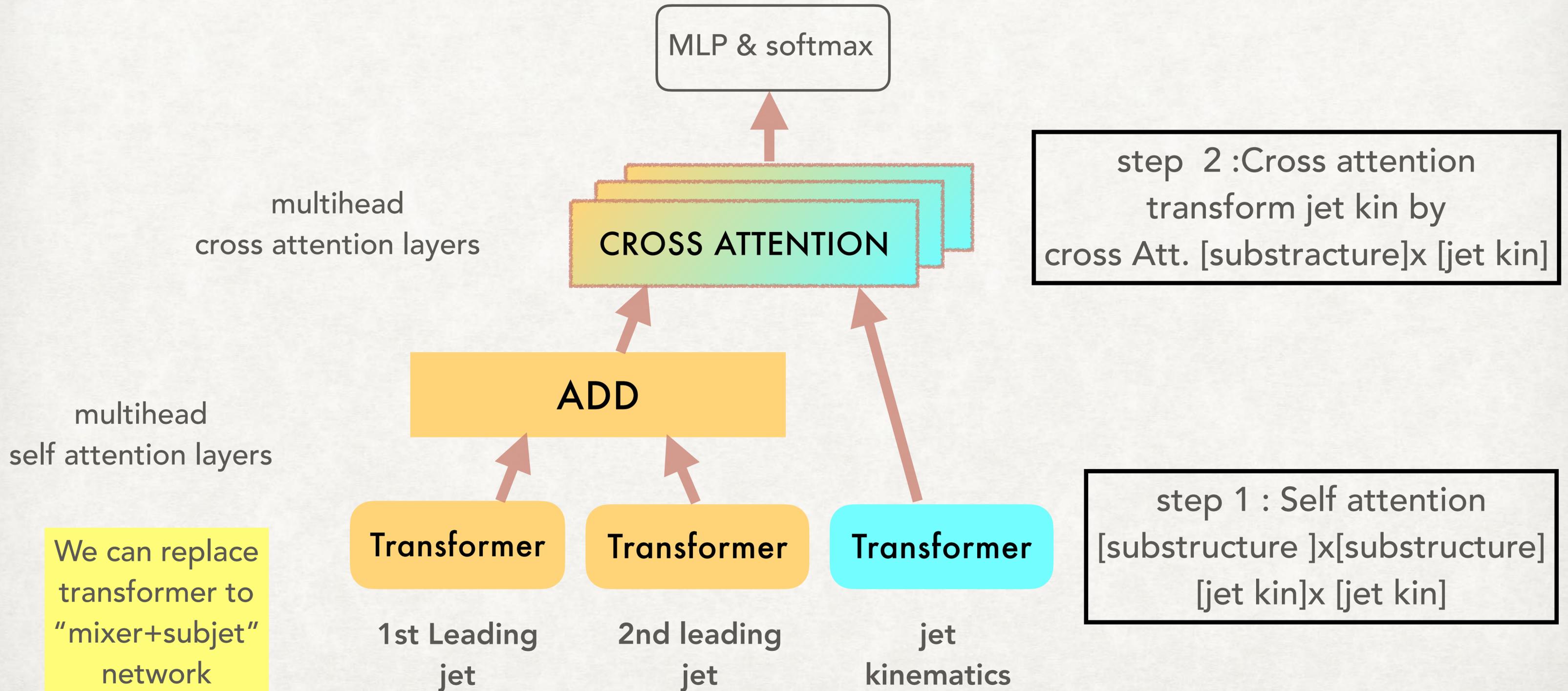


Figure 2: Feynman diagram for the signal process.

# cross attention motivation for 2 fatjet events



# INPUT TO NETWORK : EVENT KINEMATICS

Kinematical inputs (3, 6)

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

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

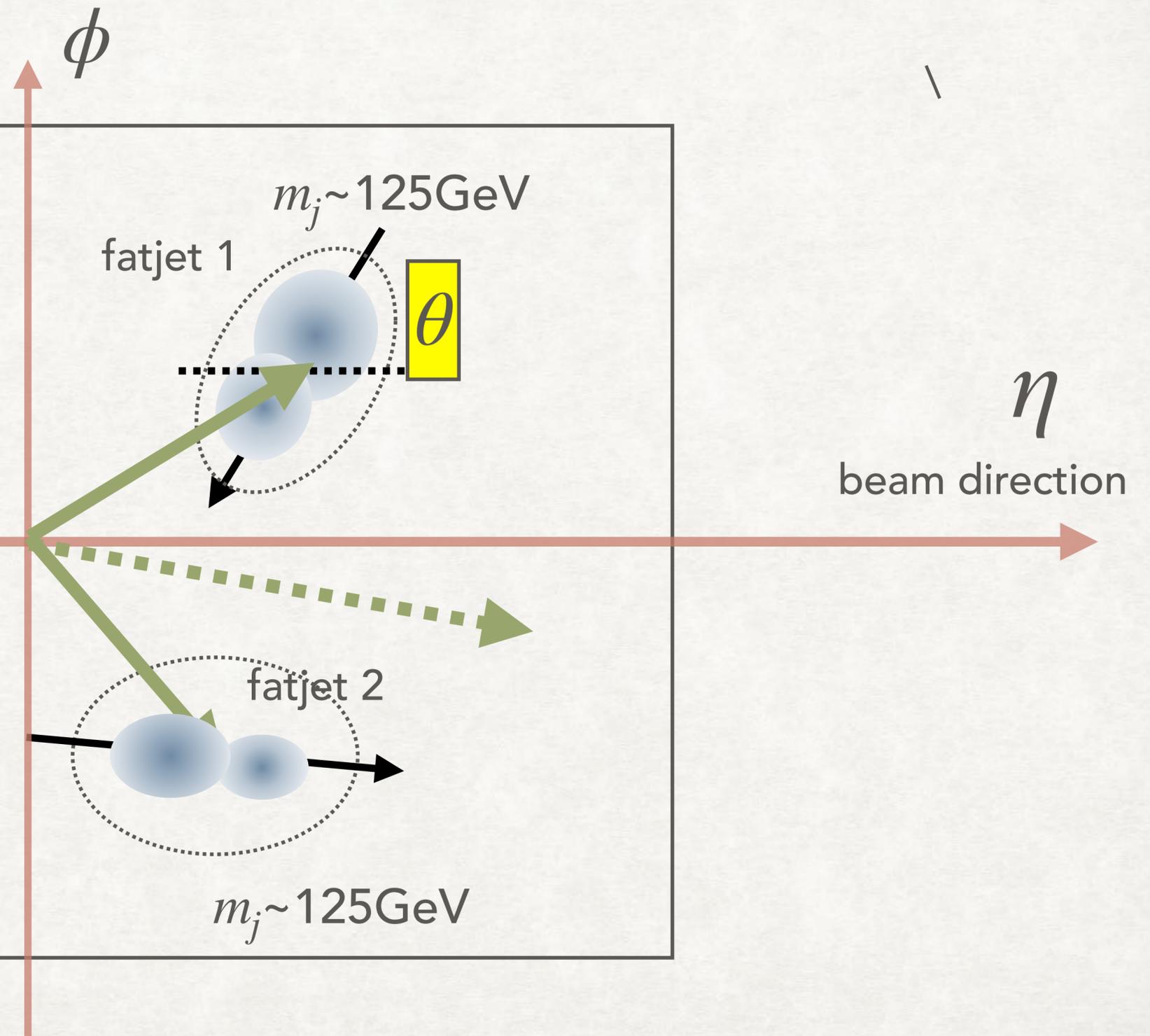
H candidate =  $(m_{12}, \eta_{12}, \phi_{12}, p_{T12}, E_{12}), \theta_{12} = 0$

NOTE :

1. "5 inputs for 4 momentum" ,

2. H candidate momentum as sum of the fat jet momentum.

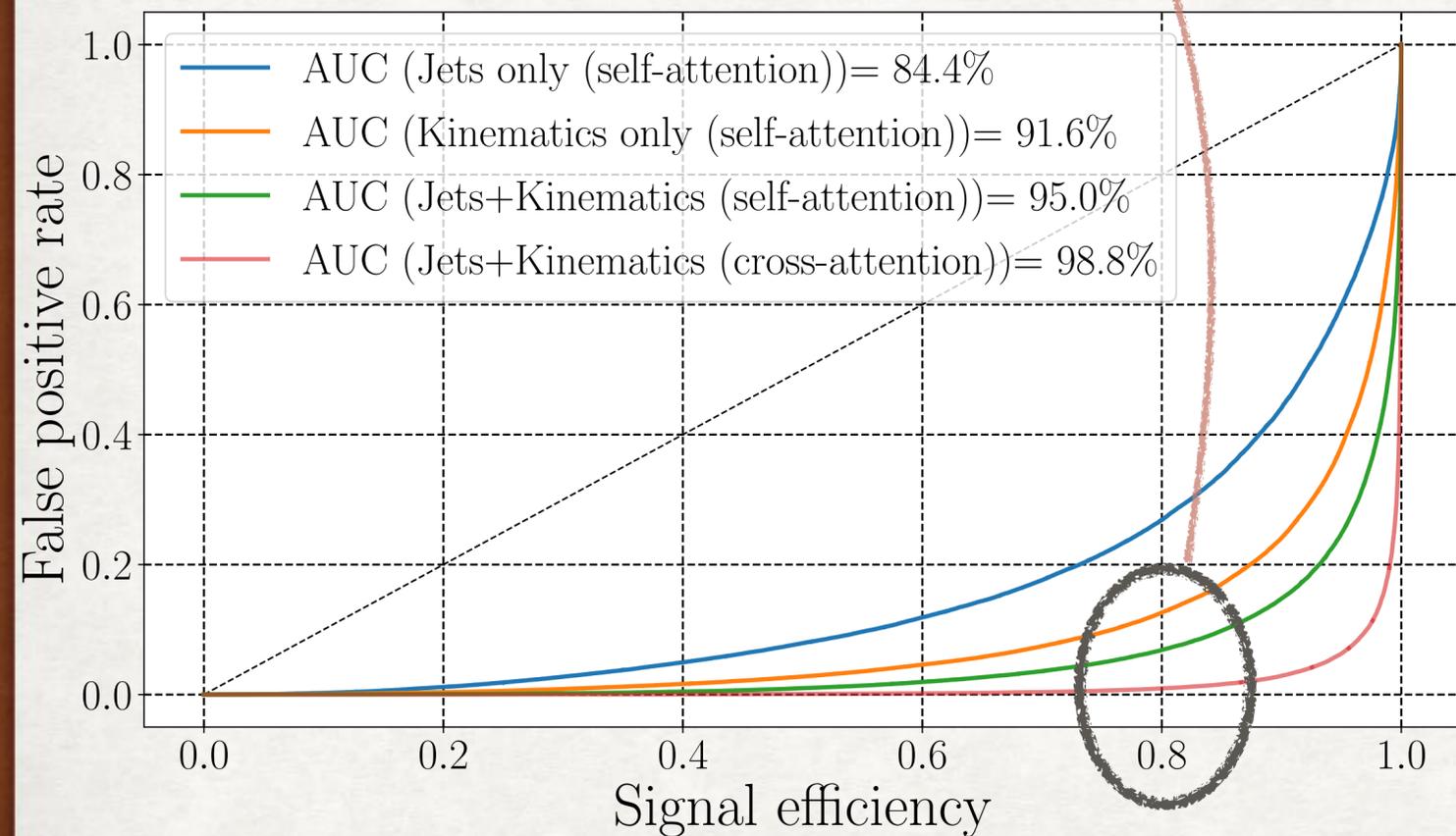
3. add " $\theta$ " :the correlation beyond a subjet



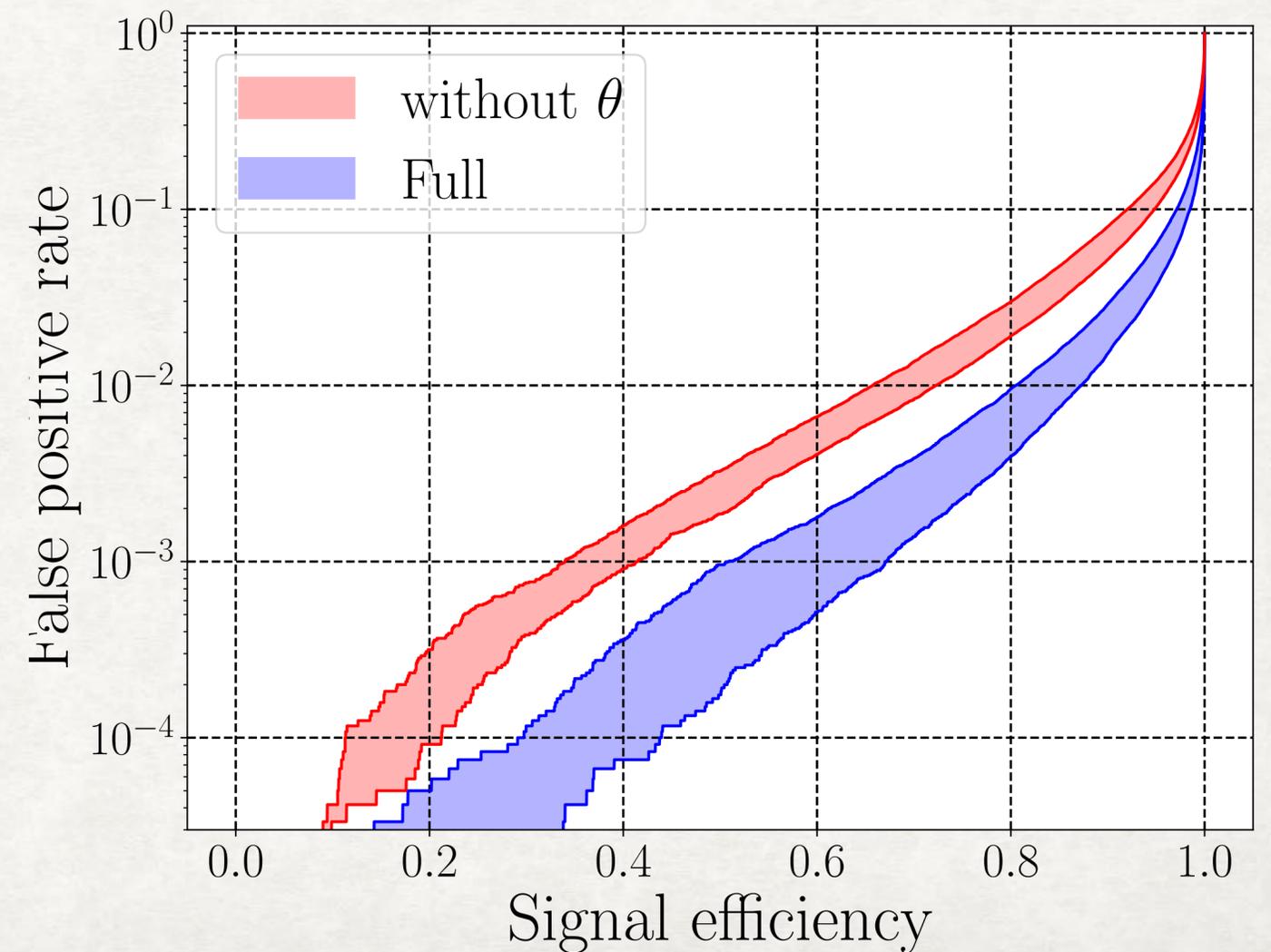
# IMPROVEMENT USING CROSS ATTENTION

factor 5 improvement at the same acceptance.

Decay correlation is important  
(because QCD background is correlated)



Cross attention improves the rejection efficiency significantly

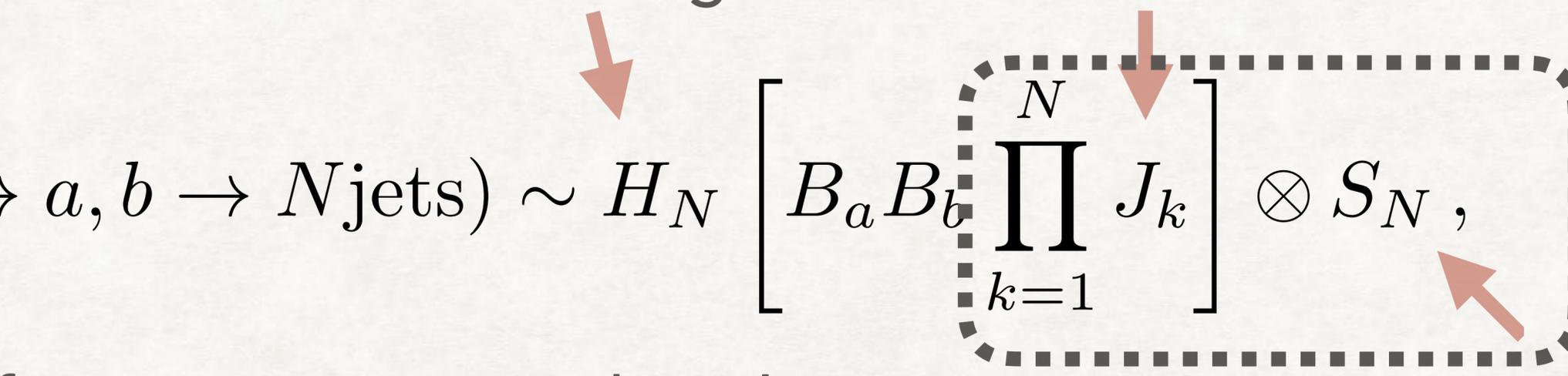


# SUMMARY

LHC process

Hard scattering

Jet function Parton shower

$$\sigma(pp \rightarrow a, b \rightarrow N\text{jets}) \sim H_N \left[ B_a B_b \prod_{k=1}^N J_k \right] \otimes S_N,$$


Cross attention

for P( constituents | (sub)jets ~ partons)

Something soft

constituent  
information

$$A = QK^T$$

Local information via (sub)jets K

# SUMMERY

## Mixer+ Subjet network

- Small, first, and high perfomance (you can test it on your computer!)
- Can apply repeatedly without losing information.
- you can stack all information (vertex, track, etc )