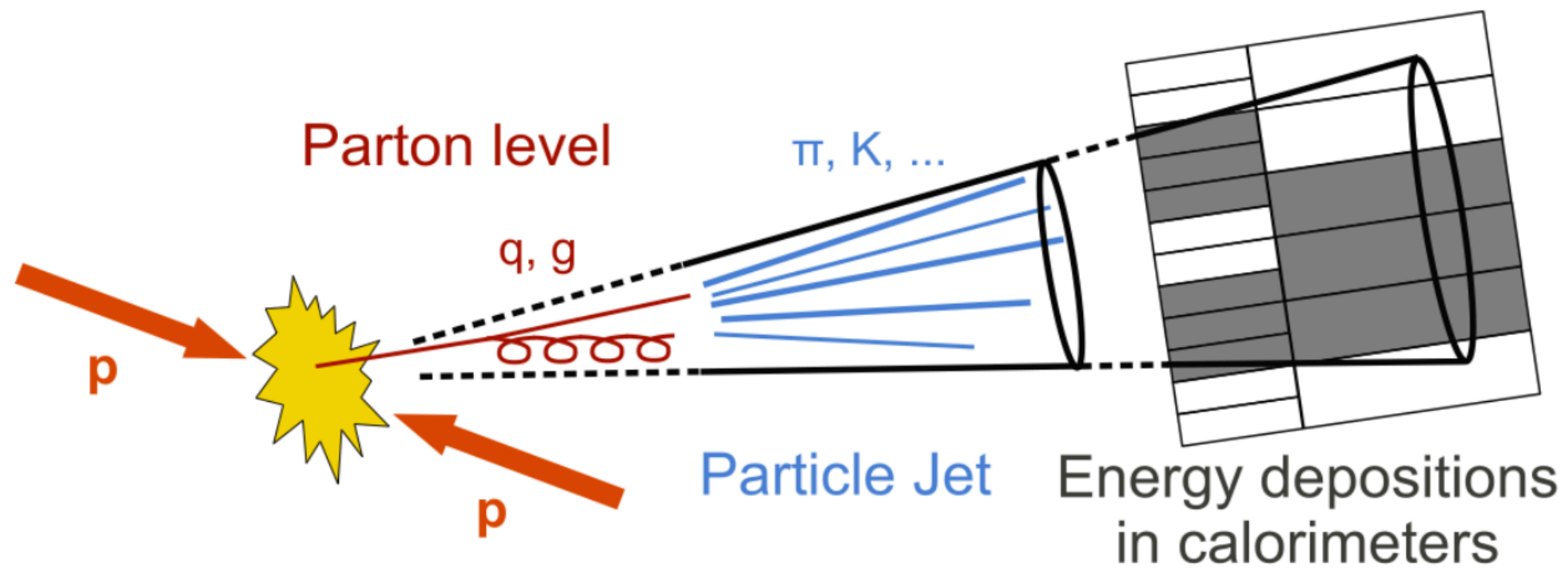# EIC jet physics and machine learning
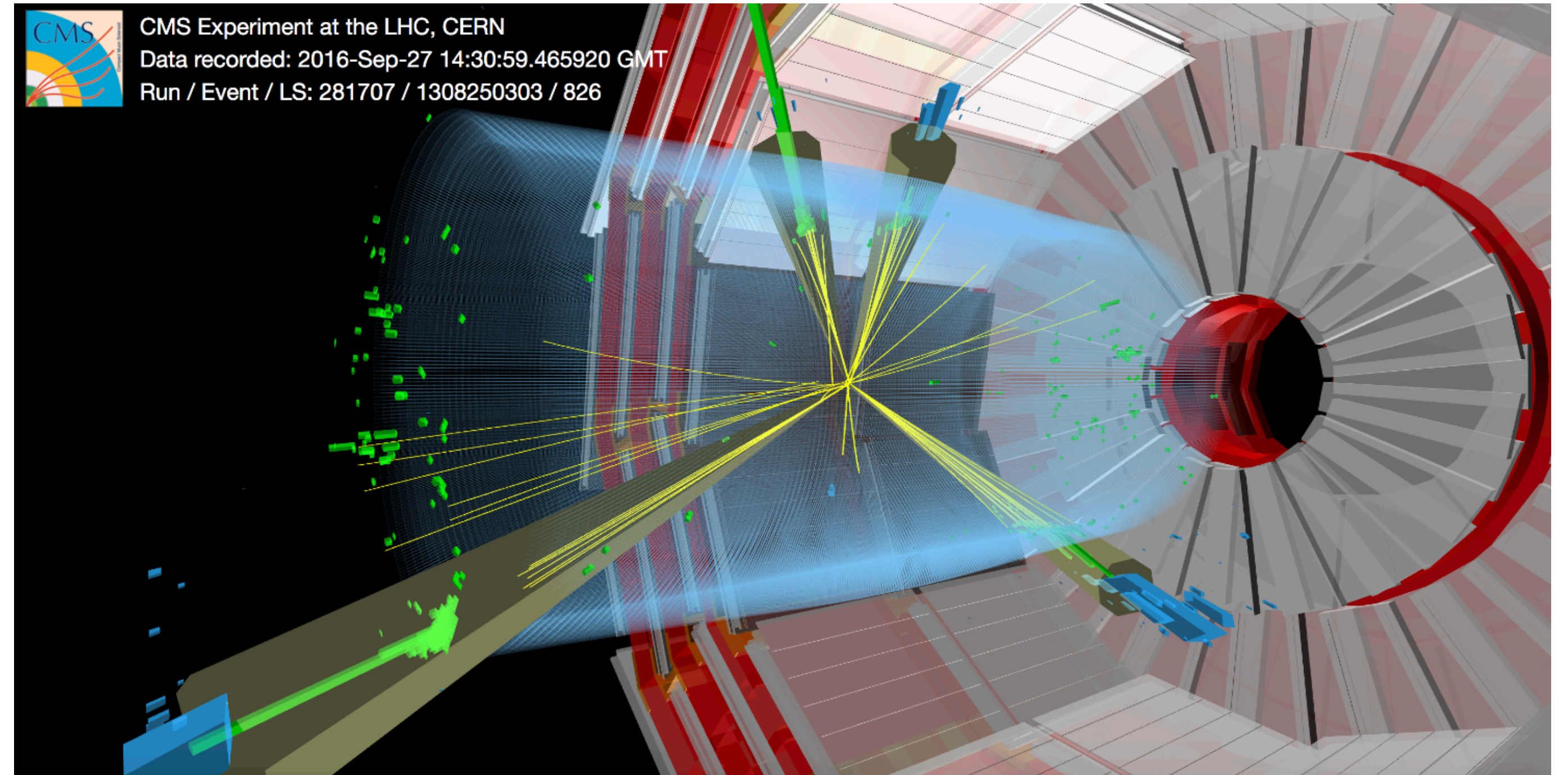
Felix Ringer

Probing Hadron Structure at the Electron-Ion
Collider, ICTS, Bangalore
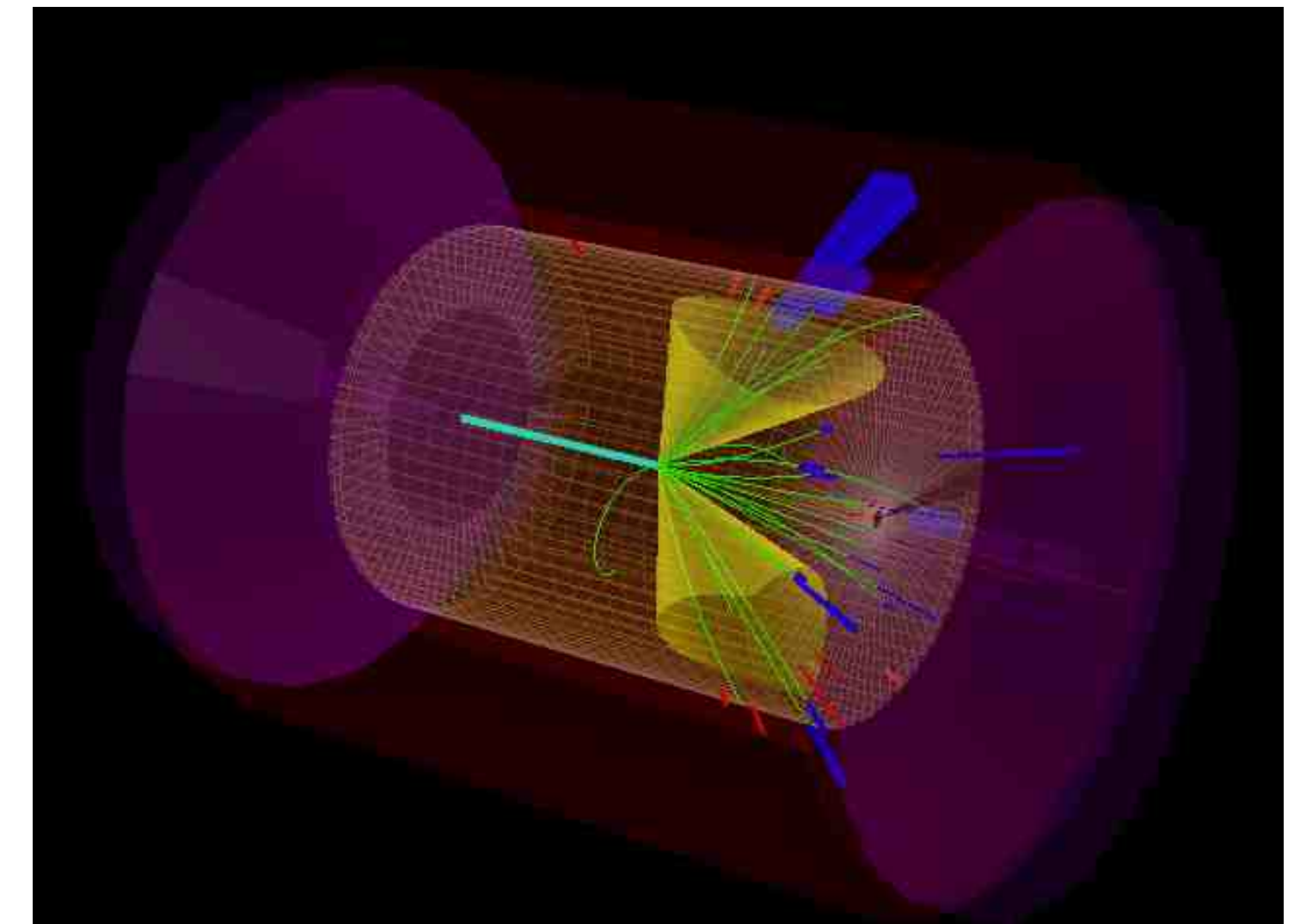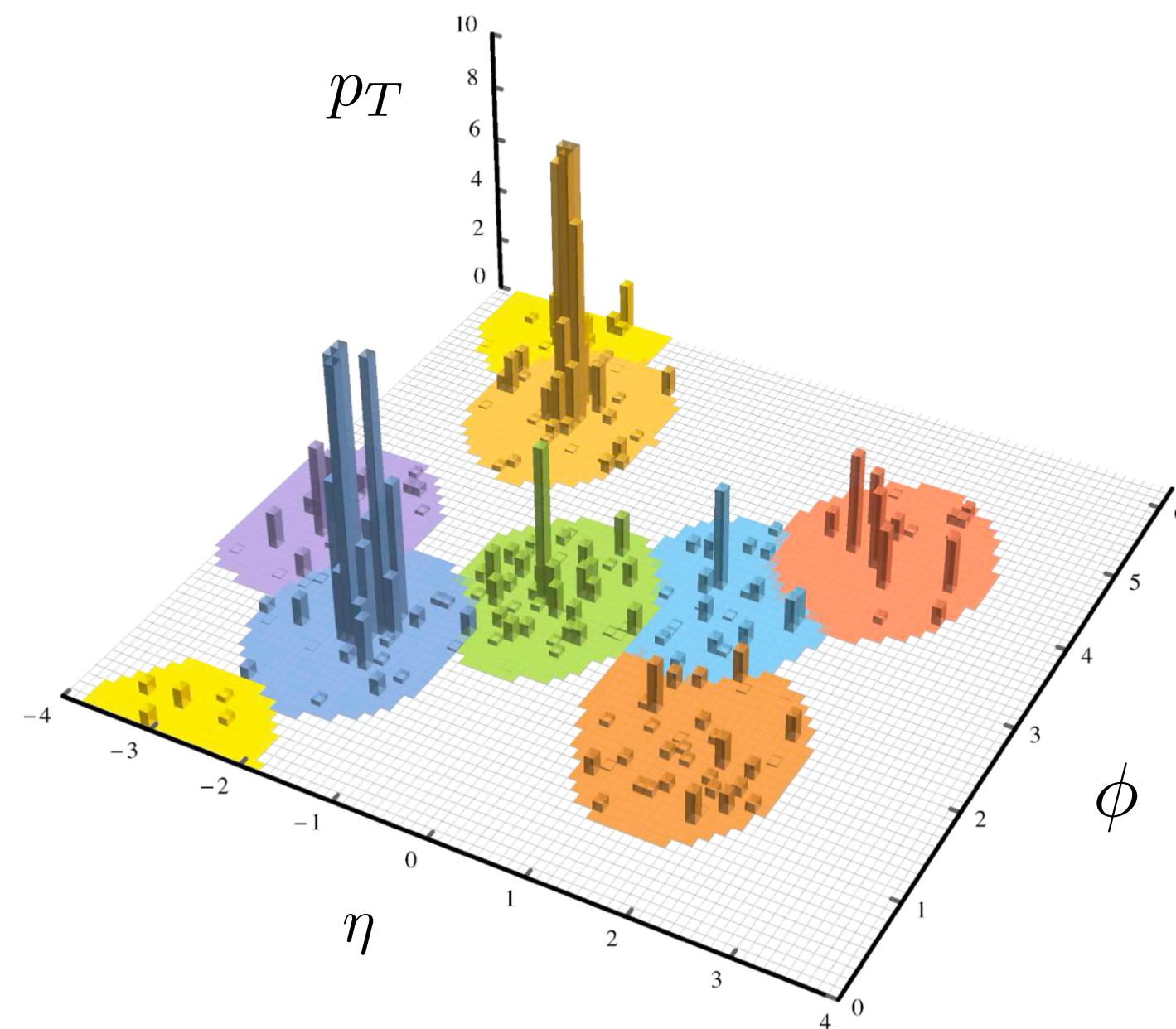
OLD DOMINION
UNIVERSITY

Jefferson Lab

# Jets at collider experiments

- Collimated sprays of particles
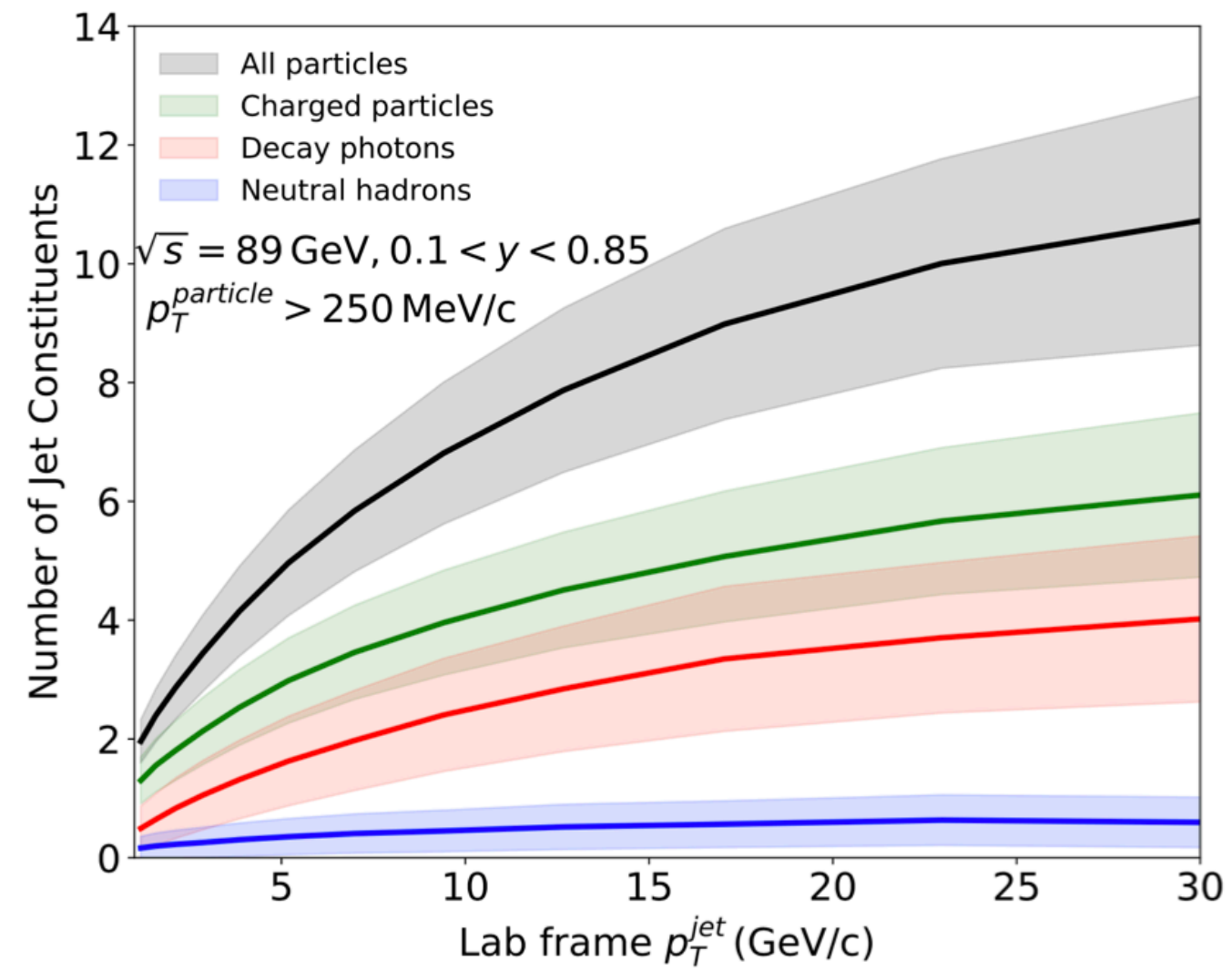- Most direct access to high-energy quarks and gluons

# EIC jet physics

- Versatile jet reconstruction algorithms & frame dependence

- Clean EIC environment

- Jet substructure & correlations

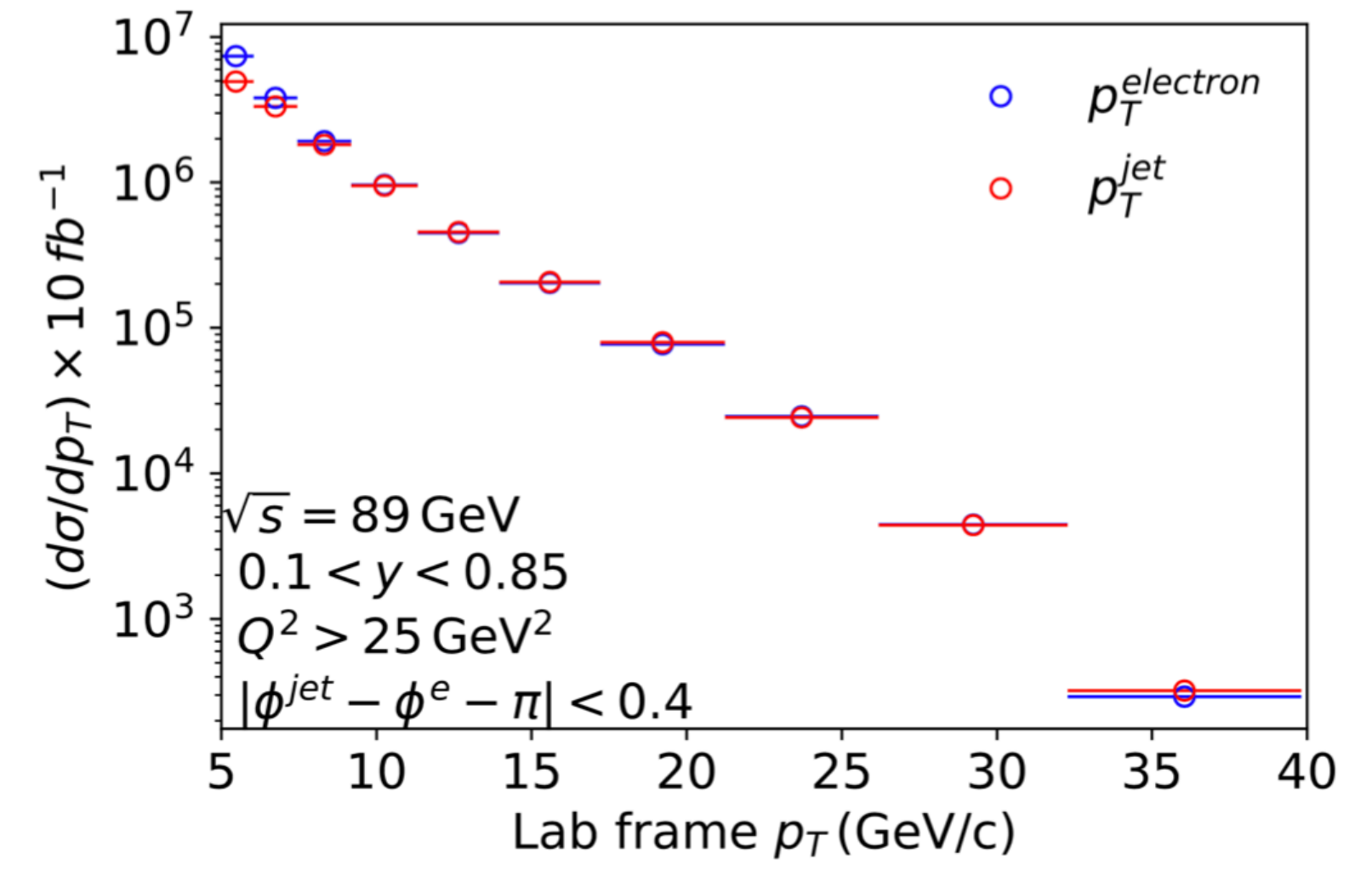- Relevant for e.g. TMDs, GPDs & hadronization
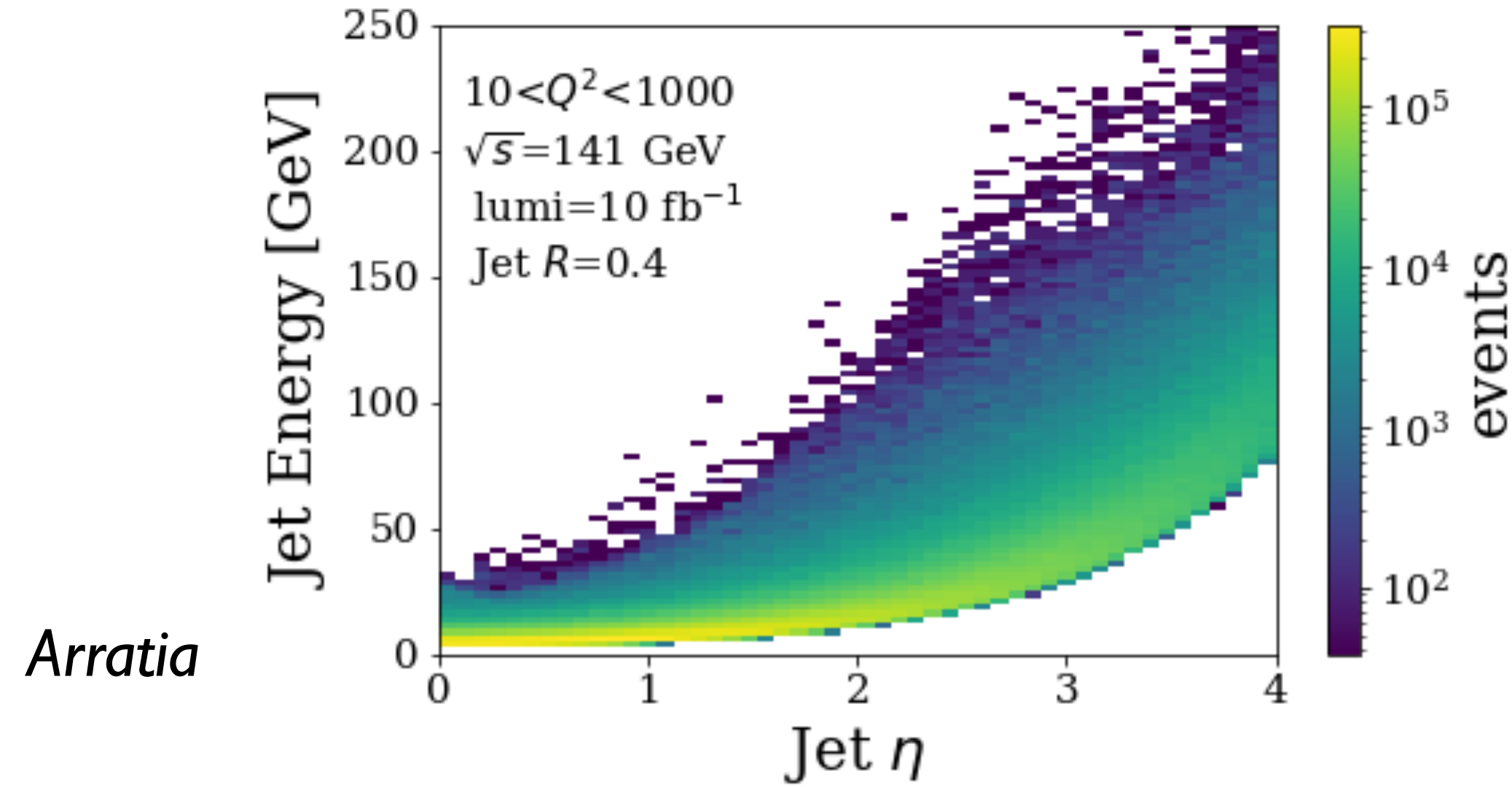
# Nature of jets at the EIC

**Particle #**



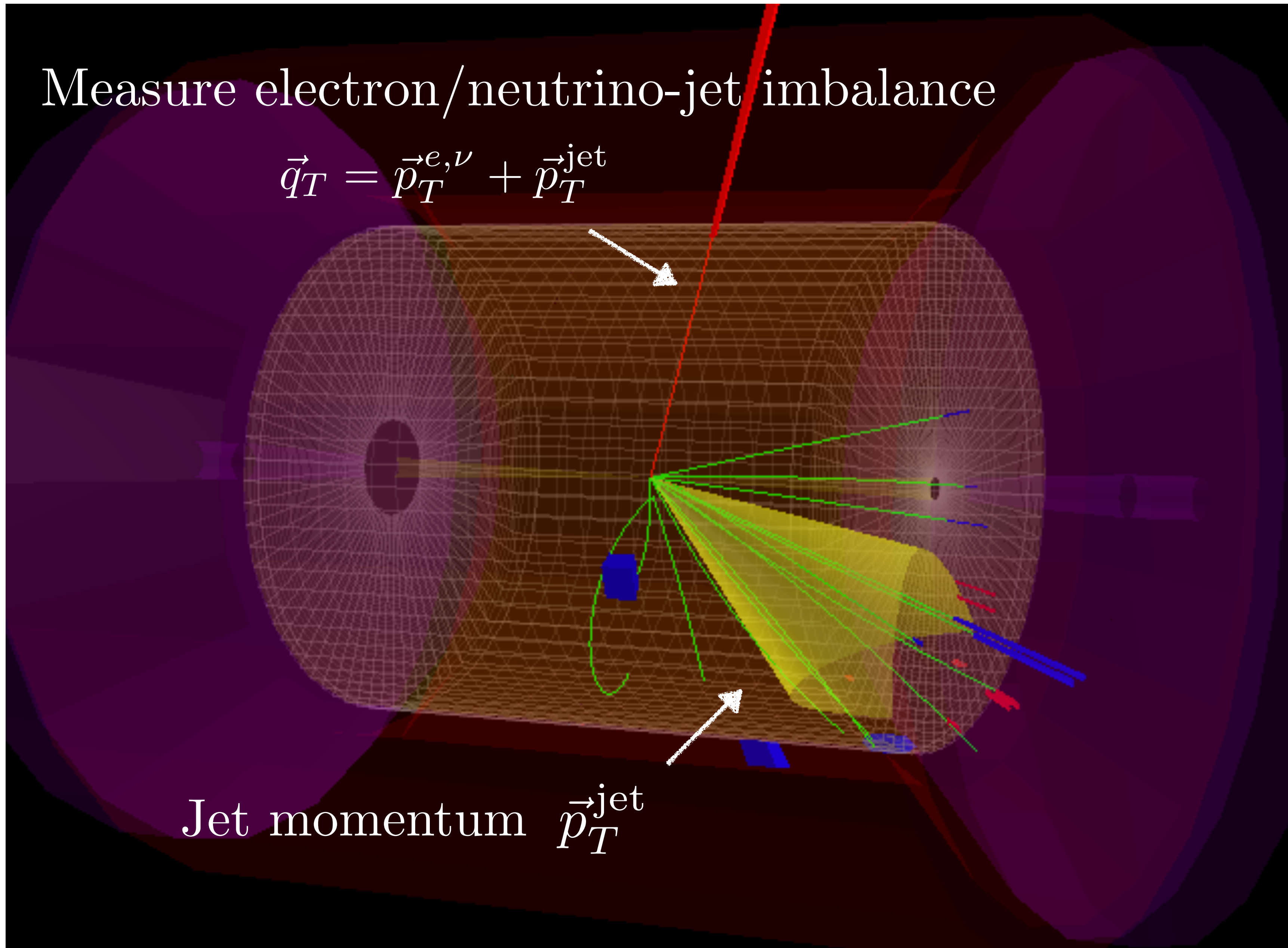**Transverse momentum**



**Jet energy**



*Arratia*

Hard scale $p_T$
and/or $Q^2$

*Arratia, Jacak, FR, Song `19*
*see also Aschenauer et al.*

Laboratory frame

Measure electron/neutrino-jet imbalance

$$\vec{q}_T = \vec{p}_T^{\,e,\nu} + \vec{p}_T^{\,\text{jet}}$$

Jet momentum $\vec{p}_T^{\,\text{jet}}$

# Electron-jet correlations

- Electron-jet imbalance at the EIC

$$\vec{q}_T = \vec{p}_T^{\,e} + \vec{p}_T^{\,\text{jet}}$$

- Sensitivity to TMD PDFs but no TMD FF

- TMD factorization

$$F_{UU} = \sigma_0 \, H_q(Q, \mu) \sum_q e_q^2 \, J_q(p_T^{\text{jet}} R, \mu)$$

$$\times \int \frac{\mathrm{d}^2 \vec{b}_T}{(2\pi)^2} \, e^{i \vec{q}_T \cdot \vec{b}_T} \, f_q^{\text{TMD}}(x, \vec{b}_T, \mu) \, S_q(\vec{b}_T, y_{\text{jet}}, R, \mu)$$
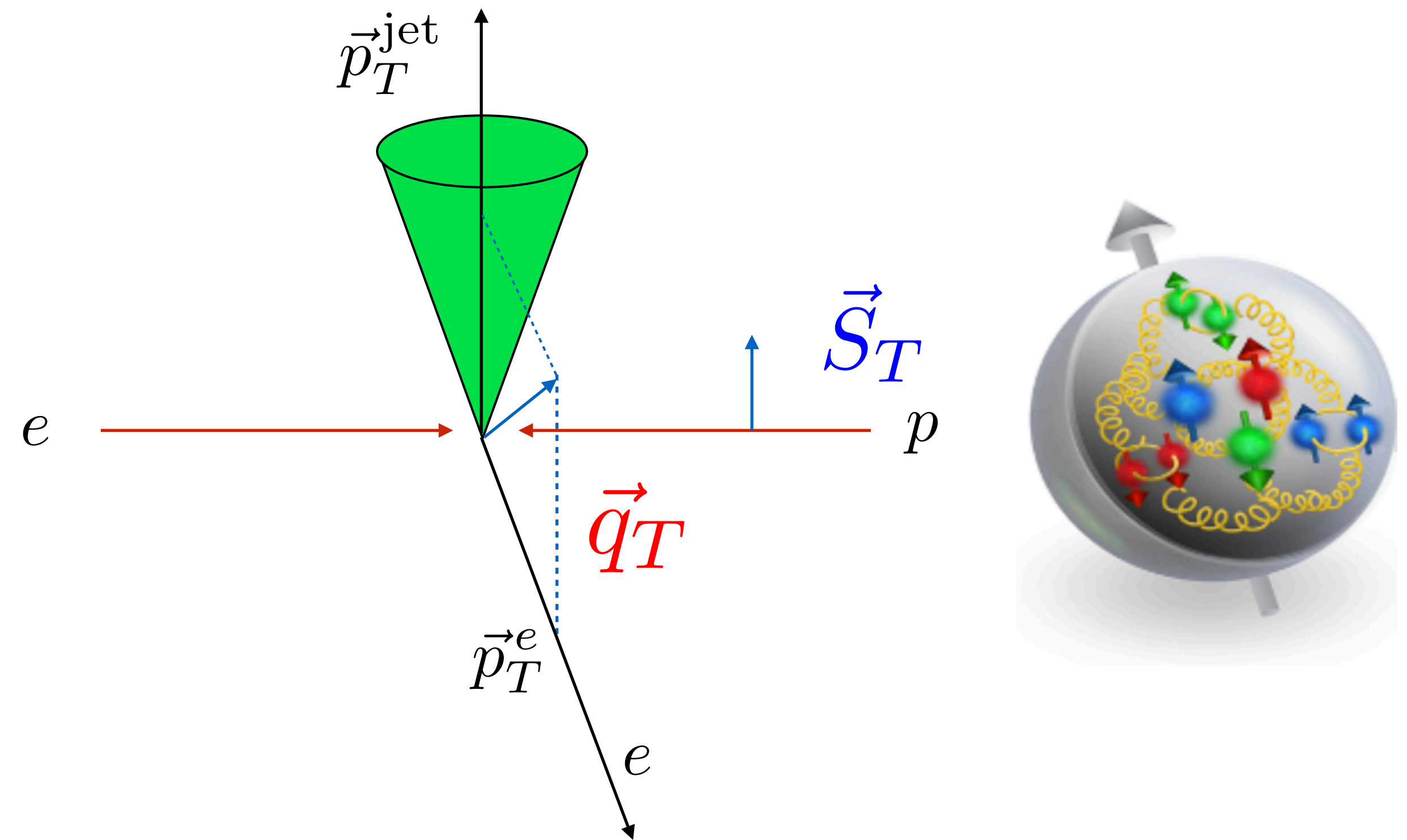
# Electron-jet correlations

- Electron-jet imbalance at the EIC

$$\vec{q}_T = \vec{p}_T^{\,e} + \vec{p}_T^{\,\text{jet}}$$

- Sensitivity to TMD PDFs but no TMD FF

- TMD factorization

$$F_{UU} = \sigma_0 \, H_q(Q, \mu) \sum_q e_q^2 \, J_q(p_T^{\text{jet}} R, \mu)$$

$$\times \int \frac{\mathrm{d}^2 \vec{b}_T}{(2\pi)^2} \, e^{i\vec{q}_T \cdot \vec{b}_T} \, f_q^{\text{TMD}}(x, \vec{b}_T, \mu) \, S_q(\vec{b}_T, y_{\text{jet}}, R, \mu)$$



*Liu, FR, Vogelsang, Yuan `18, `20*
*Arratia, Kang, Prokudin, FR`20*
*H1, PRL 128 (2022) 13, 132002*

# Jets & spin asymmetries

- E.g. Sivers asymmetries can be small due to large flavor cancellations

Burkardt sum rule `04

$$\sum_{a=q,\bar{q},g} \int_0^1 dx\, f_{1T}^{\perp(1)a}(x) = 0$$
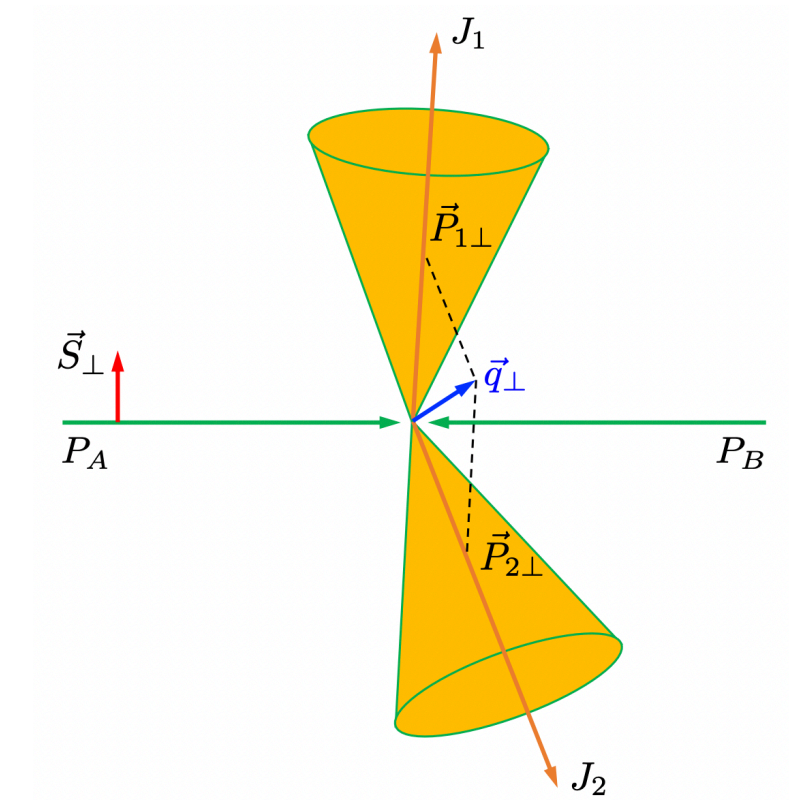
*Fatemi EINN `19, Liu DNP `19*
*see also Kang et al., Yuan et al.*





Can we obtain better constraints with ML-based jet classification?



Fundamental QCD parameters

# Jet physics & Machine learning

- Various jet classifiers have been developed

- Typically ML significantly outperformed traditional observables

- Use full event-by-event information instead of low-dimensional projections (observables)



*Fig. Komiske, Metodiev, Schwartz*

# Jet physics & Machine learning

- Various jet classifiers have been developed

  - Example: Quark vs. gluon jet classification

- Quantify using a ROC curve



$u \quad d \quad s \quad c \quad g$

V.S.



AI/ML

Traditional observable

*Gallicchio, Schwartz*
*Komiske, Metodiev, Thaler `19*

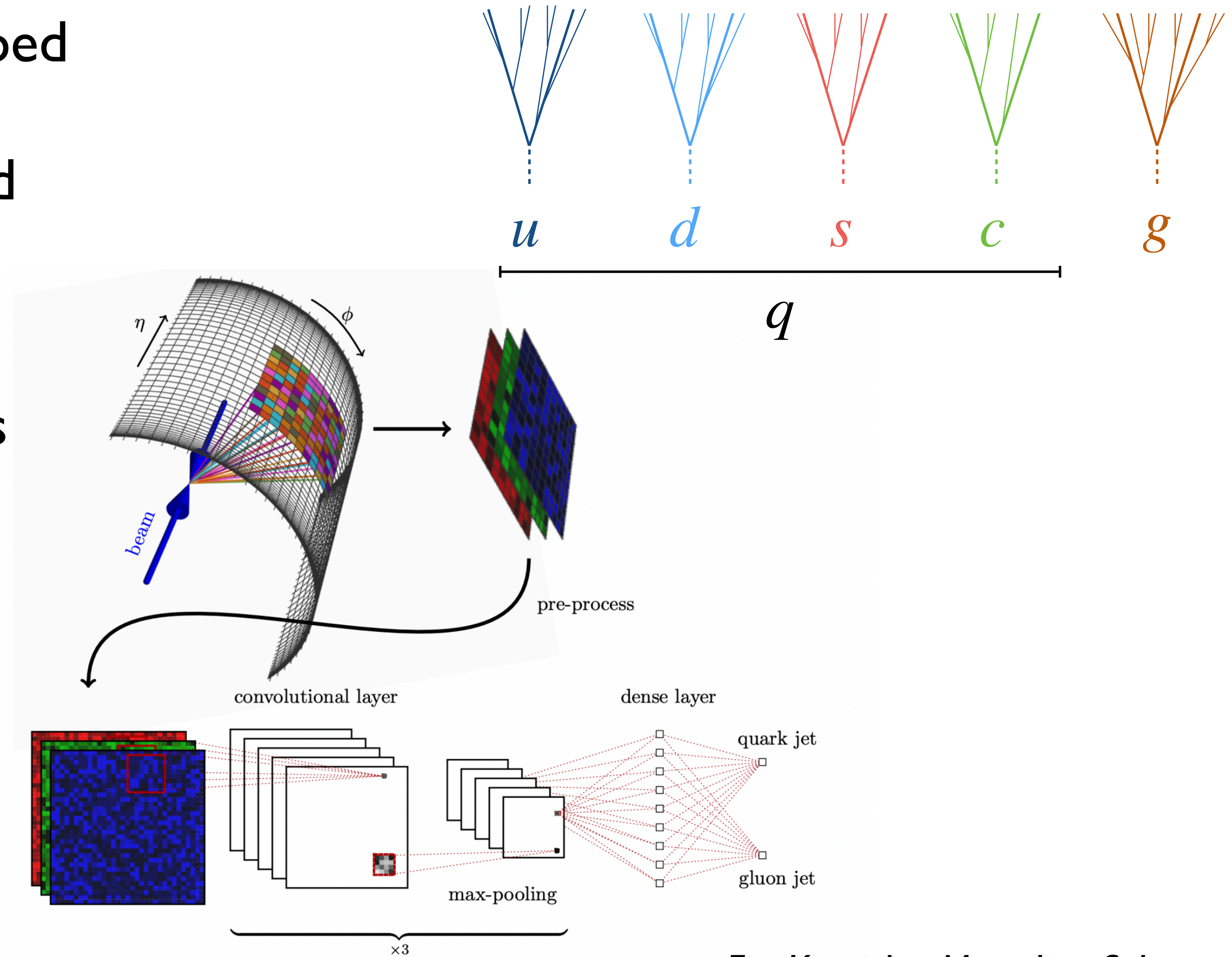# Data sets using Pythia

*Lee, Mulligan, Ploskon, FR, Yuan `22*

- Relatively low particle multiplicities at the EIC

- Pythia 6
  - ▫ No detector simulation
  - ▫ Laboratory frame jets
  - ▫ Particle $(p_{Ti}, \eta_i, \phi_i, \mathrm{PID}_i)$

$$Q^2 > 25 \ \mathrm{GeV}^2,$$
$$p_T > 10 \ \mathrm{GeV}$$

$e'$

$e$ ———————— $p$

jet

$ud, s$ jet classification

Photoproduction, low $Q^2$

jet₂

$\gamma$

$e$ ——— $p$

$e'$

jet₁

$q, g$ jet classification

# Machine learning architecture

- Binary classification: $u$ vs. $d$, $ud$ vs. $s$, …

- Deep sets or Particle Flow Networks

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

Classifier

Neural networks



Particles    Observable

Per−Particle Representation    Event Representation

Latent Space

$\Phi$

$\Phi$

$\Phi$

$F$

Energy/Particle Flow Network

*Komiske, Metodiev, Thaler JHEP 01 (2019) 121*
*Permutation invariant Deep Sets*
*See also GNNs, transformers*

# Quark vs. gluon jet tagging

*Lee, Mulligan, Ploskon, FR, Yuan `22*



$q$ vs. $g$ jet

better

$$\text{True Positive Rate} = \frac{\text{True } q}{\text{Total } q}$$

$$\text{False Positive Rate} = \frac{\text{False } q}{\text{Total } g}$$

Leading jet
— Particle Flow Network (w/ PID)
— Energy Flow Network
— Energy Flow Polynomials (DNN), $d = 7$
⋯ Jet mass

☐ Some improvement with ML
☐ Relatively few particles per jet, less information

## Data & code available

*https://zenodo.org/record/7538810#.Y8RcaS-B2gQ*

# Quark vs. gluon event tagging

*Lee, Mulligan, Ploskon, FR, Yuan `22*

## $qq, q\bar{q}$ vs. $gg$ process



Axis labels:
- $\text{True Positive Rate} = \frac{\text{True } qq, q\bar{q}}{\text{Total } qq, q\bar{q}}$
- $\text{False Positive Rate} = \frac{\text{False } qq, q\bar{q}}{\text{Total } gg}$

Legend:
- Leading jet
- Leading jet + subleading jet
- All event particles

☐ Significant gain with ML!

☐ Use full event information

☐ Quantifies total information content

☐ Motivates further theory efforts

## Data & code available

*https://zenodo.org/record/7538810#.Y8RcaS-B2gQ*

# Strange jet tagging

$u, d$ vs. $s$ jets

$u, d$ vs. $s$ jets

better

$$\text{True Positive Rate} = \frac{\text{True } s}{\text{Total } s}$$

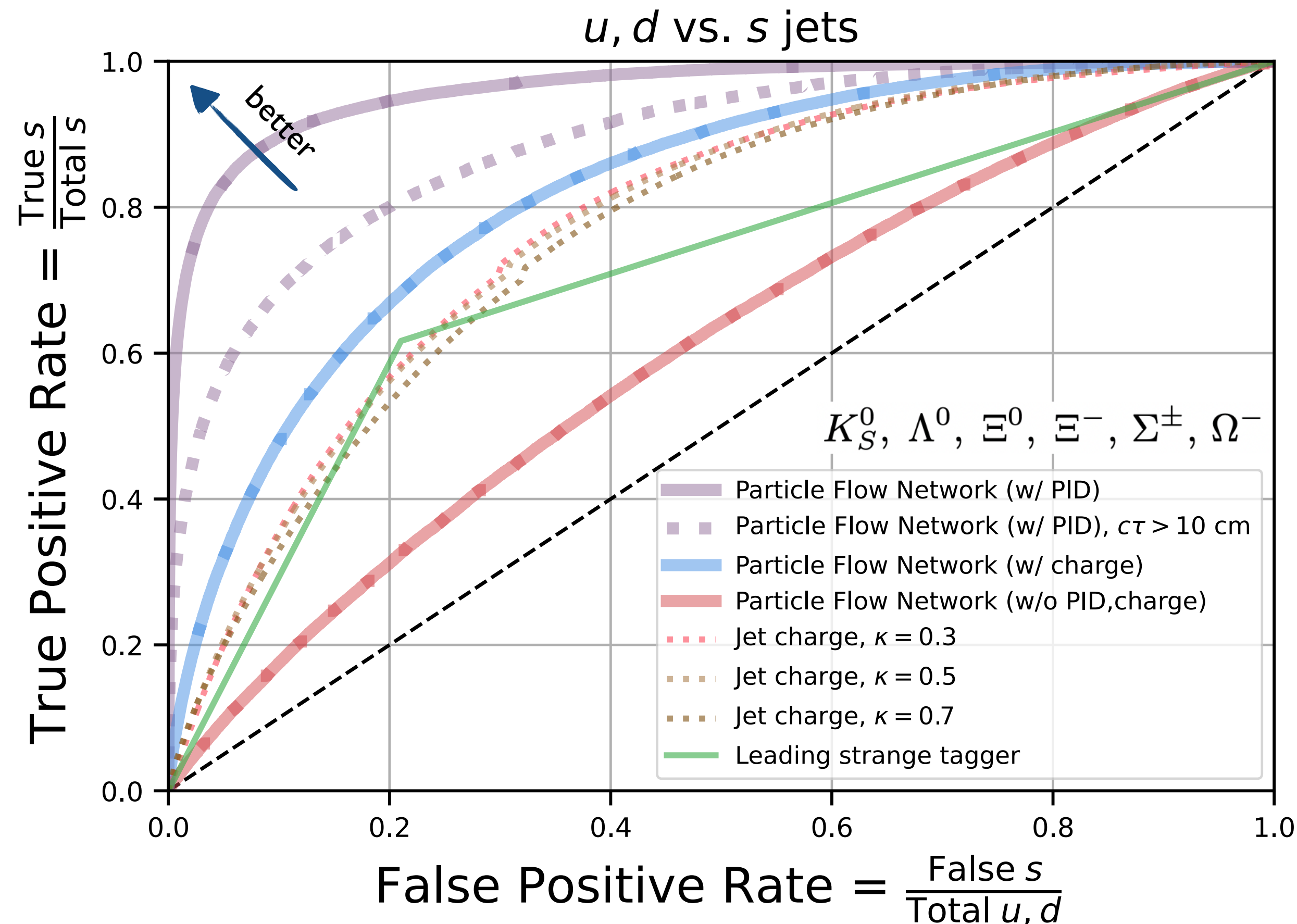$K_S^0,\ \Lambda^0,\ \Xi^0,\ \Xi^-,\ \Sigma^\pm,\ \Omega^-$

Particle Flow Network (w/ PID)
Particle Flow Network (w/ PID), $c\tau > 10$ cm
Particle Flow Network (w/ charge)
Particle Flow Network (w/o PID,charge)
Jet charge, $\kappa = 0.3$
Jet charge, $\kappa = 0.5$
Jet charge, $\kappa = 0.7$
Leading strange tagger

$$\text{False Positive Rate} = \frac{\text{False } s}{\text{Total } u, d}$$

Particle Flow Network (w/ PID)
Particle Flow Network (w/ PID), $c\tau > 10$ c
Particle Flow Network (w/ charge)
Particle Flow Network (w/o PID,charge)
Jet charge, $\kappa = 0.3$
Jet charge, $\kappa = 0.5$
Jet charge, $\kappa = 0.7$
Leading strange tagger

$$\text{Precision} = \frac{\text{True } s}{\text{True } s + \text{False } s}$$

Data & code available

$$\text{Recall} = \frac{\text{True } s}{\text{Total } s}$$

# Information content of jets & events

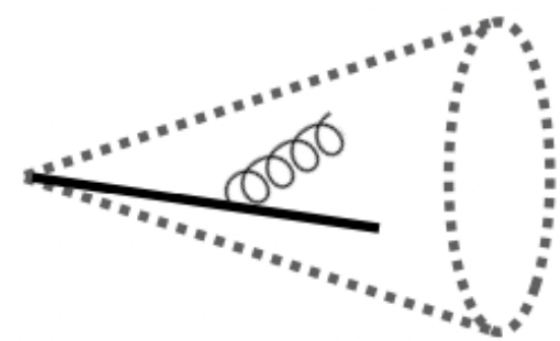How can we make use of all this additional information?

# Information content of jets & events

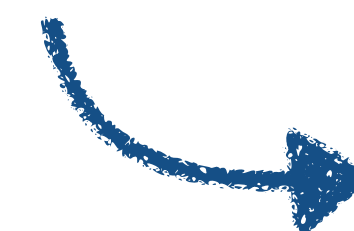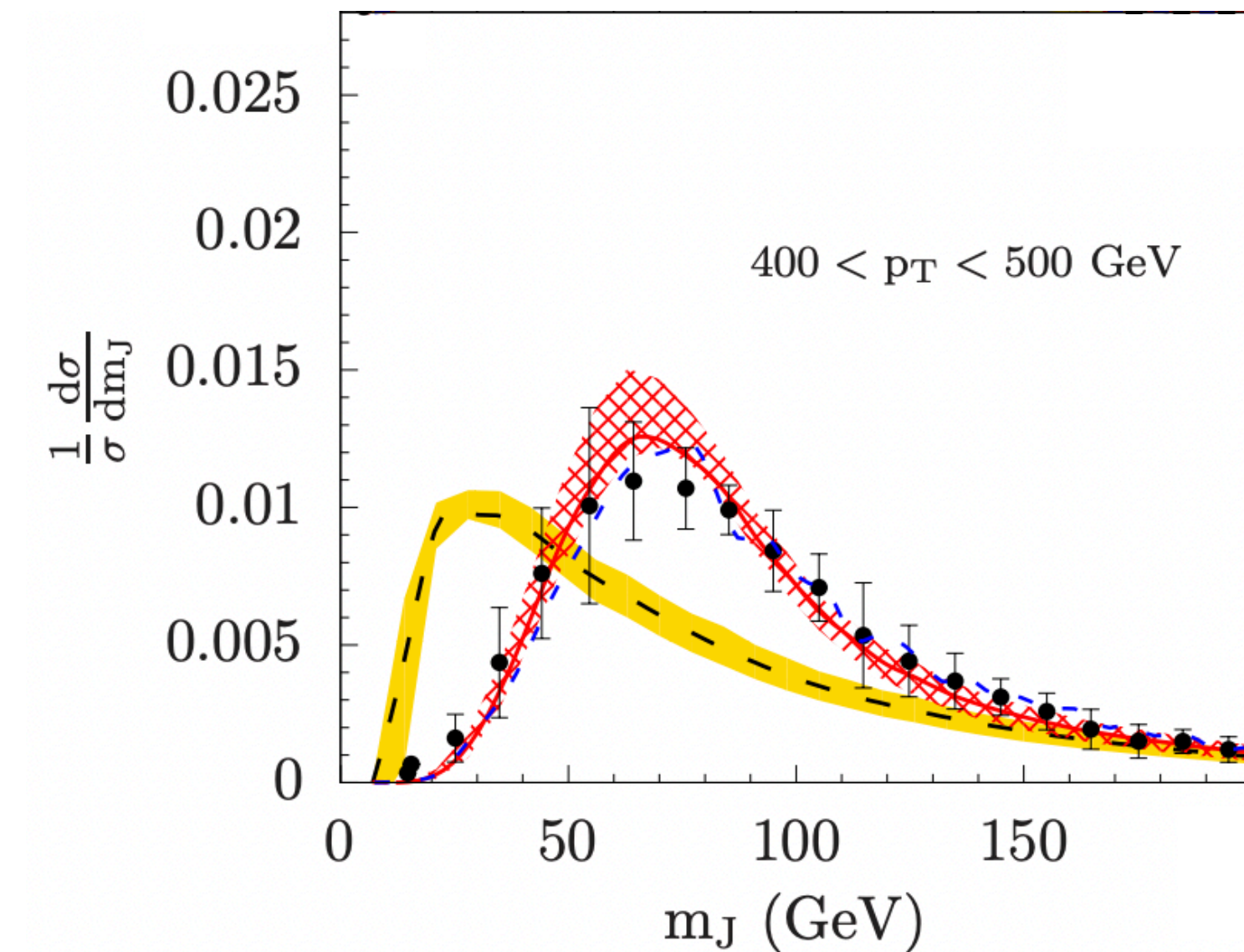How can we make use of all this additional information?

- Need complete sets of observables

- Observables at the level of events vs. ensemble



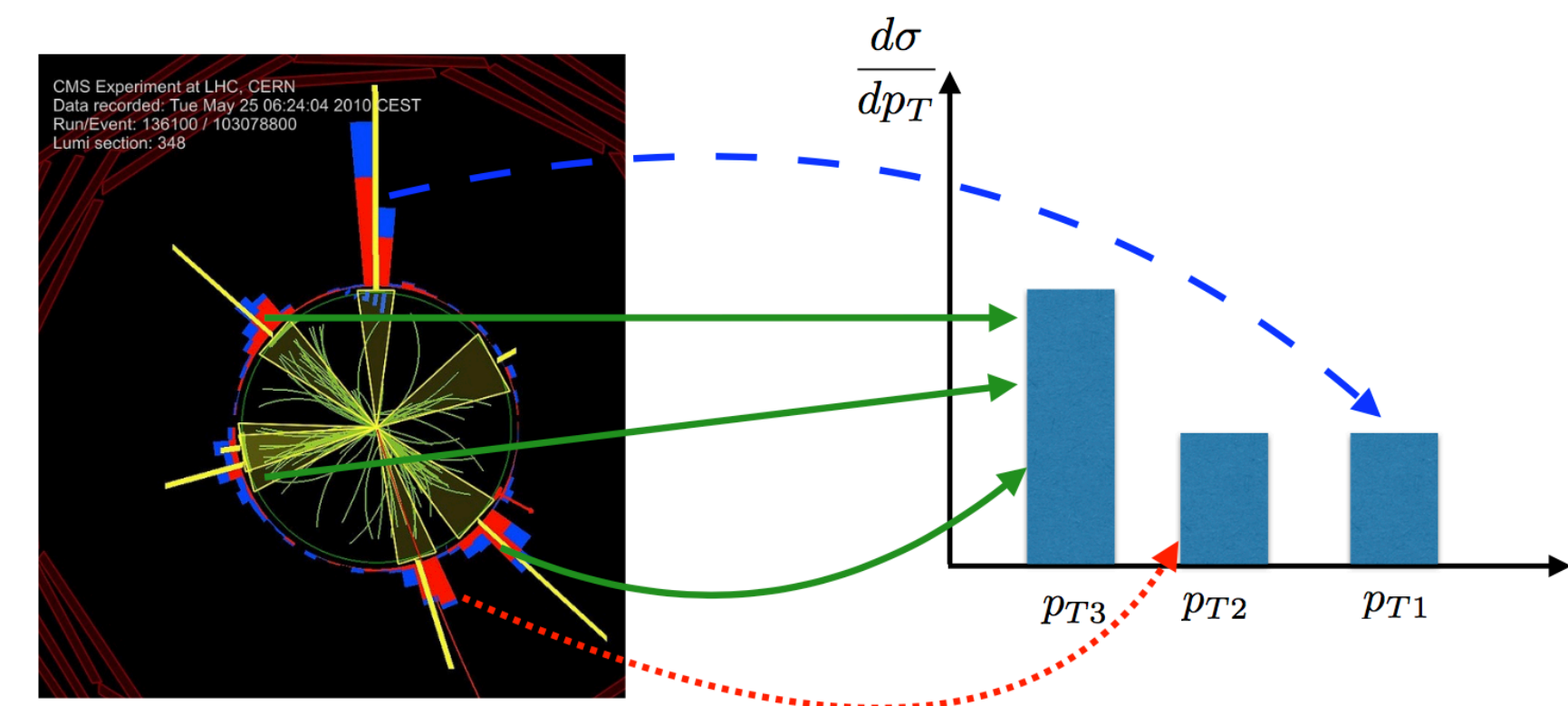Measure $m_J$ per jet
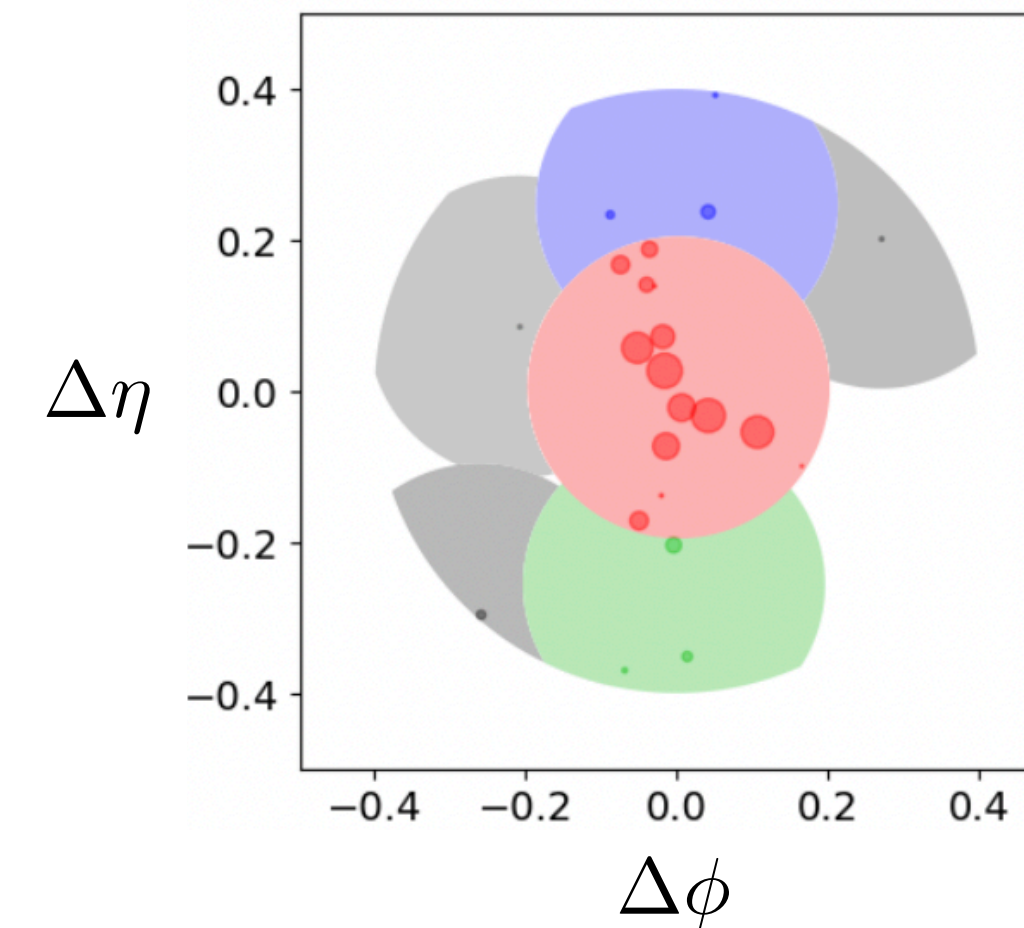
Histogram event samples

Compare to QCD calculation

# Information content of jets & events

How can we make use of all this additional information?

- Need complete sets of observables

- Observables at the level of events vs. ensemble

  - <u>Event only</u>: Position information $(\eta_i, \phi_i)$

  - <u>Ensemble only</u>: Inclusive jets and correlators

# Information content of jets & events

How can we make use of all this additional information?



- Need complete set of observables
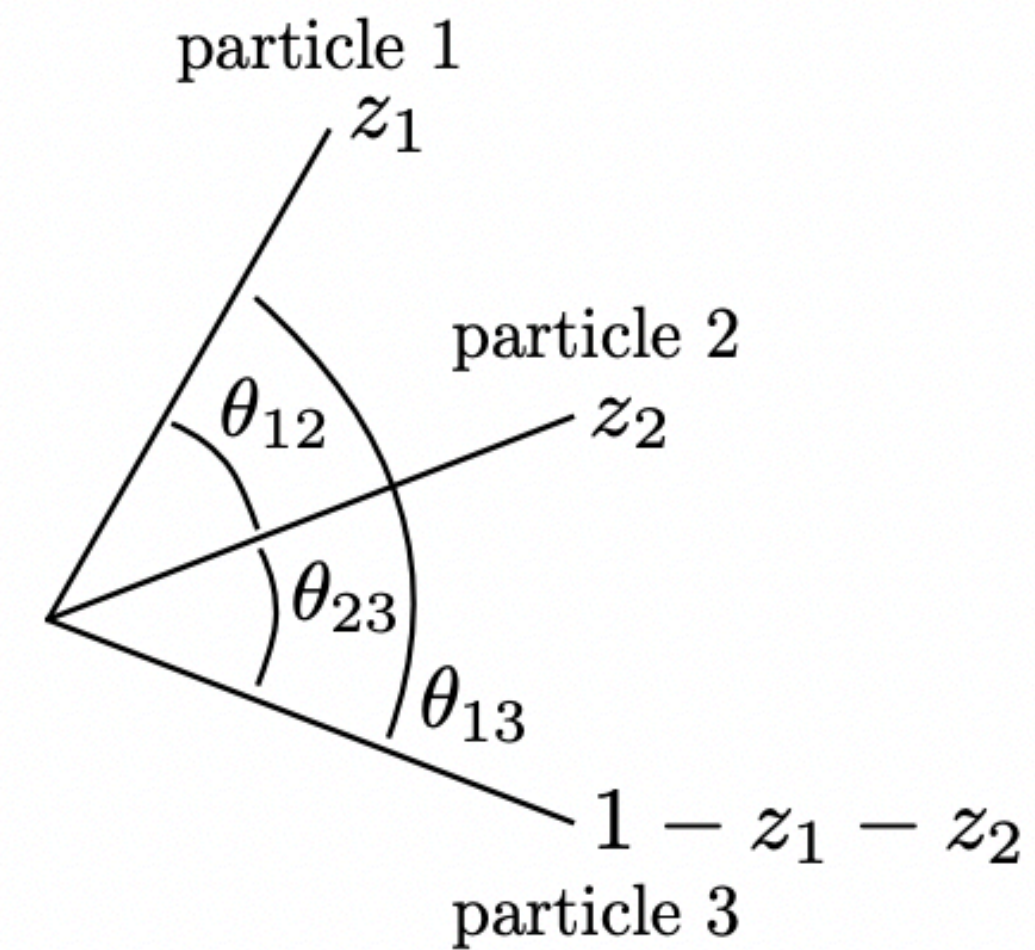
- N-jettiness basis & Energy Flow Polynomials

*Datta, Larkoski `17*                    *Metodiev, Komiske, Thaler `18*

Both are IRC safe & defined at the event and ensemble level

→        Can use AI to identify the most useful observables!

# Information content of jets & events

- N-jettiness basis  *Datta, Larkoski `17*

Systematically map 3M-4 phase space
variables to a set of observables e.g.

$$z(1-z) = \frac{\left(\tau_1^{(1)}\right)^2}{4\tau_1^{(2)}}, \qquad \theta = \frac{2\tau_1^{(2)}}{\tau_1^{(1)}}$$

emission 1
$z$
$\theta$
$1 - z$
emission 2

emission 1
$z_1$
$\theta_{12}$
emission 2
$z_2$
$\theta_{23}$
$\theta_{13}$
$1 - z_1 - z_2$
emission 3

$$\tau_N^{(\beta)} = \frac{1}{p_T} \sum_{i \in \text{jet}} p_{Ti} \min\left\{ R_{1i}^\beta, R_{2i}^\beta, \ldots, R_{Ni}^\beta \right\}$$

$$\left\{ \tau_1^{(0.5)}, \tau_1^{(1)}, \tau_1^{(2)}, \tau_2^{(0.5)}, \tau_2^{(1)}, \tau_2^{(2)}, \ldots \right\} \longrightarrow$$

Use as input to a neural network for
classification and feature selection
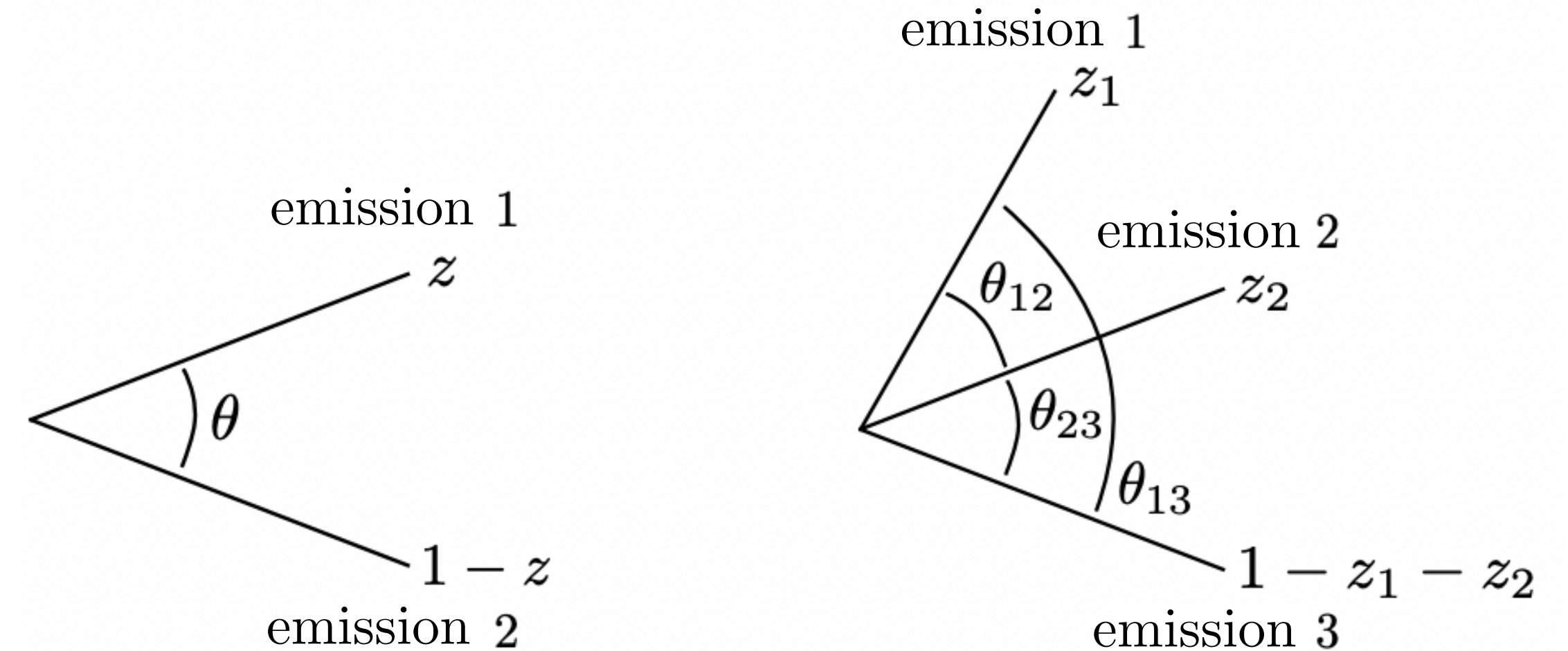
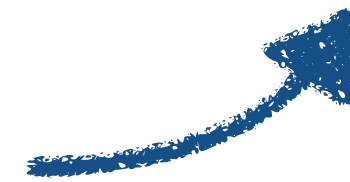*see e.g. Lai, Mulligan, Ploskon, FR `21*

# Information content of jets & events

- N-jettiness basis  *Datta, Larkoski `17*

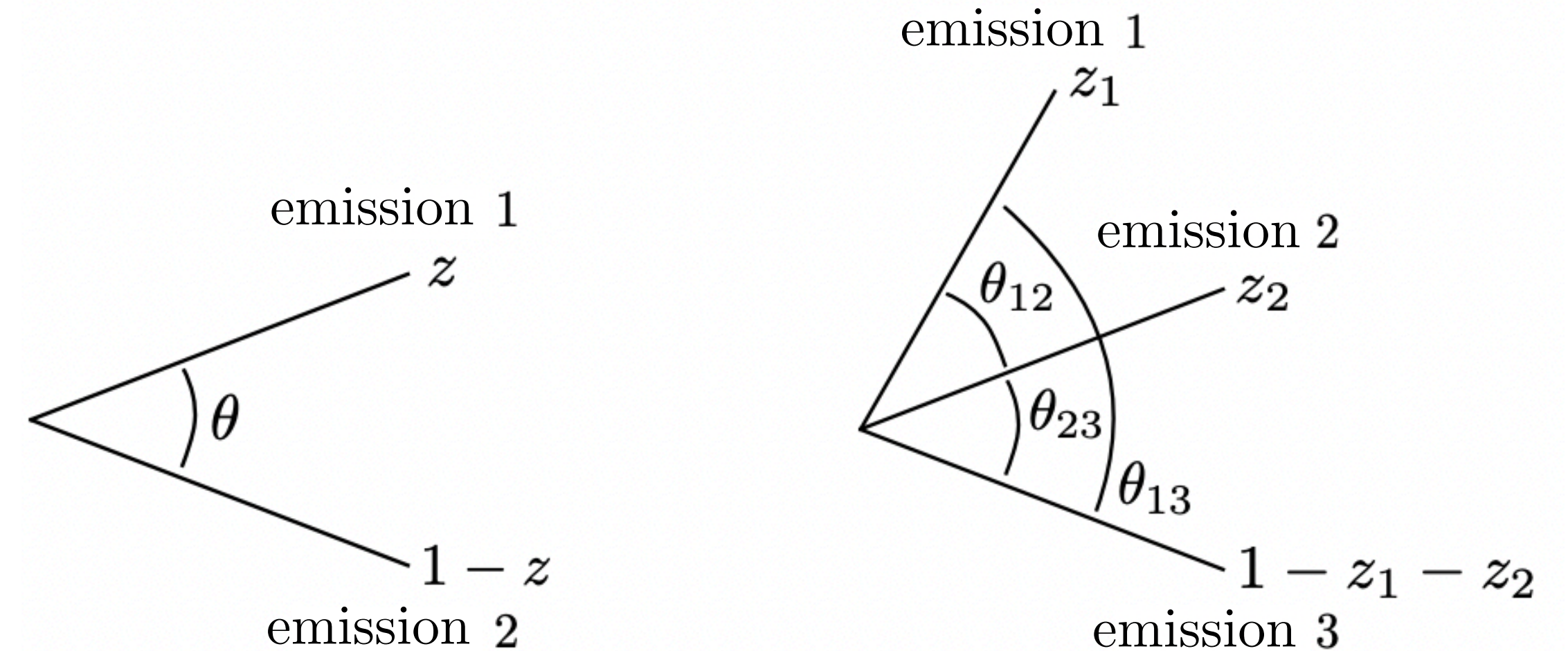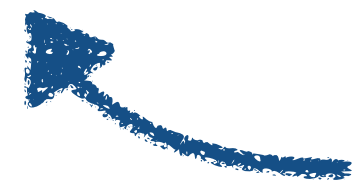Systematically map 3M-4 phase space variables to a set of observables e.g.



$$z(1-z) = \frac{\left(\tau_1^{(1)}\right)^2}{4\tau_1^{(2)}}, \qquad \theta = \frac{2\tau_1^{(2)}}{\tau_1^{(1)}}$$

Sudakov safe

$$\left\{\tau_1^{(0.5)}, \tau_1^{(1)}, \tau_1^{(2)}, \tau_2^{(0.5)}, \tau_2^{(1)}, \tau_2^{(2)}, \ldots\right\}$$

IRC safe

$$\tau_N^{(\beta)} = \frac{1}{p_T} \sum_{i \in \text{jet}} p_{Ti} \min\left\{R_{1i}^\beta, R_{2i}^\beta, \ldots, R_{Ni}^\beta\right\}$$

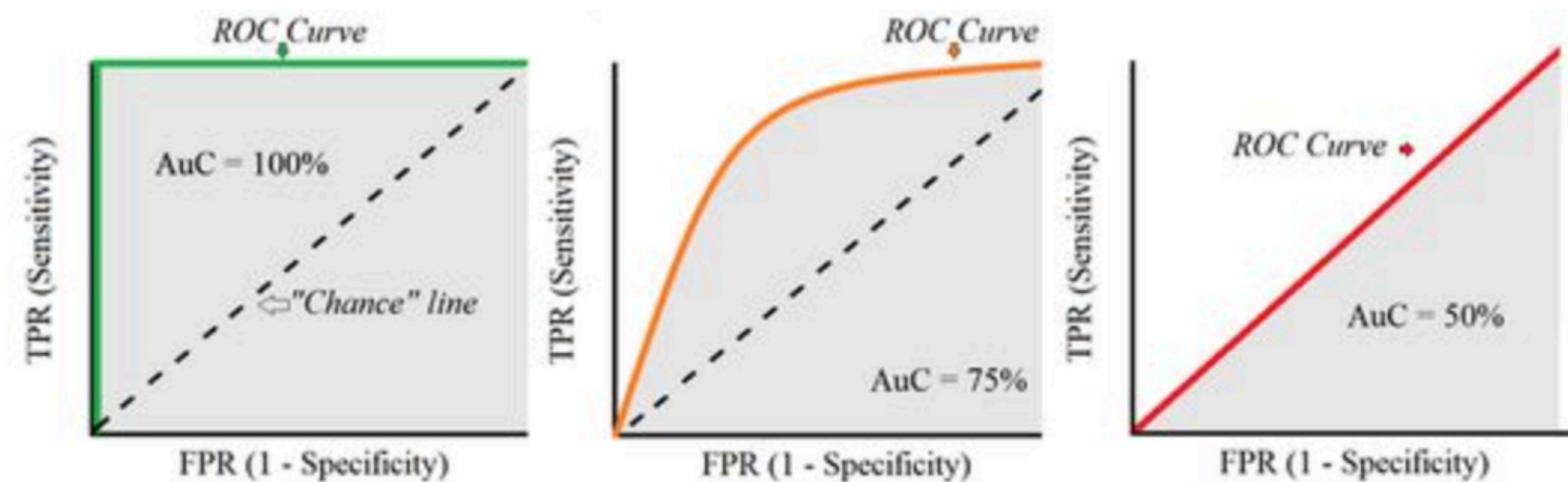Use as input to a neural network for classification and feature selection

*see e.g. Lai, Mulligan, Ploskon, FR `21*

# Information content of jets & events

- N-jettiness basis  *Datta, Larkoski `17*

… but there appears to be a
performance gap  *e.g. Metodiev, Komiske, Thaler `18*

Quantify using the area
under the ROC curve

| Model | AUC |
|---|---|
| PFN-ID | **0.9052** $\pm$ 0.0007 |
| PFN-Ex | 0.9005 $\pm$ 0.0003 |
| PFN-Ch | 0.8924 $\pm$ 0.0001 |
| PFN | 0.8911 $\pm$ 0.0008 |
| EFN | 0.8824 $\pm$ 0.0005 |
| RNN-ID | 0.9010 |
| RNN | 0.8899 |
| EFP | 0.8919 |
| DNN | 0.8849 |
| CNN | 0.8781 |
| $M$ | 0.8401 |
| $n_{SD}$ | 0.8297 |
| $m$ | 0.7401 |

IRC unsafe
classifier

⋮

N-jettiness
observables

# Information content of jets & events

- N-jettiness basis  *Datta, Larkoski `17*

… but there appears to be a
performance gap  *e.g. Metodiev, Komiske, Thaler `18*

The gap could be due to…

- IRC safety?
- the type of input?
- the network architecture?

| Model | AUC |
|-------|-----|
| PFN-ID | **0.9052** $\pm$ 0.0007 |
| PFN-Ex | 0.9005 $\pm$ 0.0003 |
| PFN-Ch | 0.8924 $\pm$ 0.0001 |
| PFN | 0.8911 $\pm$ 0.0008 |
| EFN | 0.8824 $\pm$ 0.0005 |
| RNN-ID | 0.9010 |
| RNN | 0.8899 |
| EFP | 0.8919 |
| DNN | 0.8849 |
| CNN | 0.8781 |
| $M$ | 0.8401 |
| $n_{SD}$ | 0.8297 |
| $m$ | 0.7401 |

← IRC unsafe classifier

:

← N-jettiness observables

# Is IRC-safe information all you need for jet classification?

- Use the same ML algorithm as the best classifier
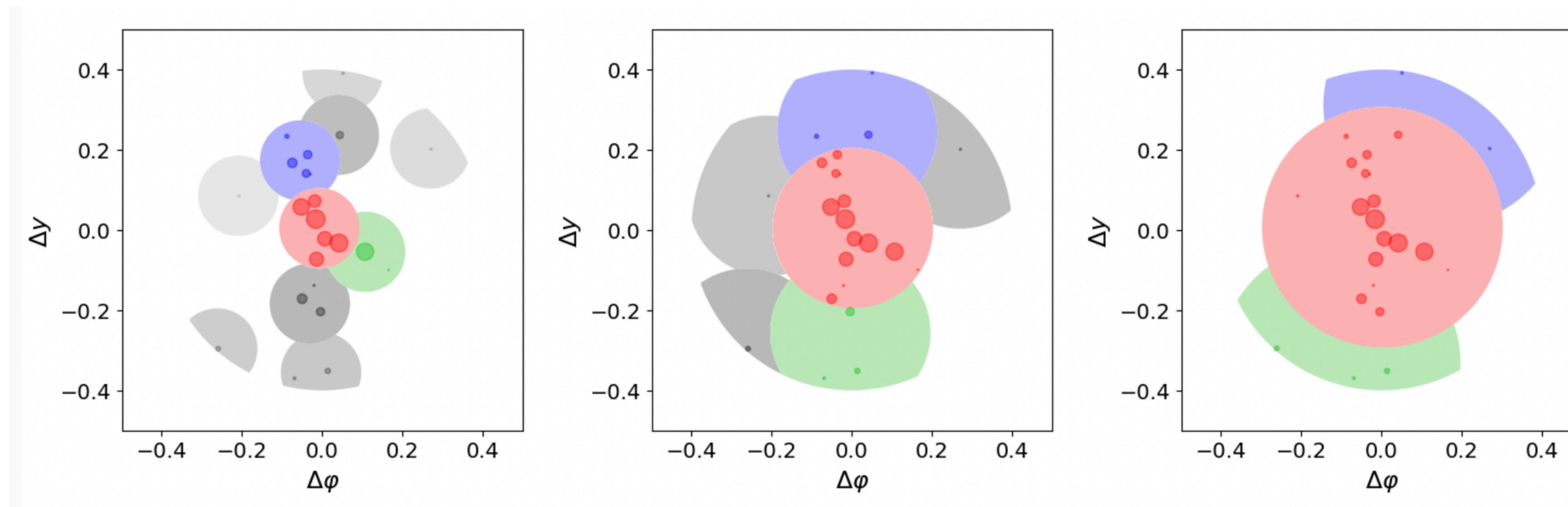
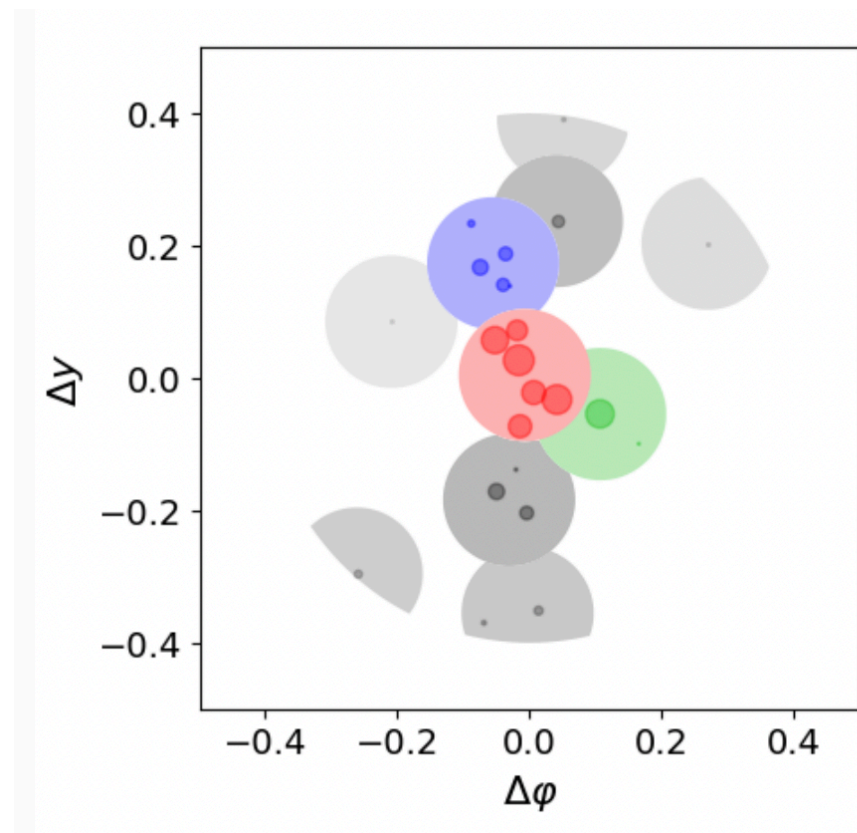- … but cluster jet constituents into subjets first
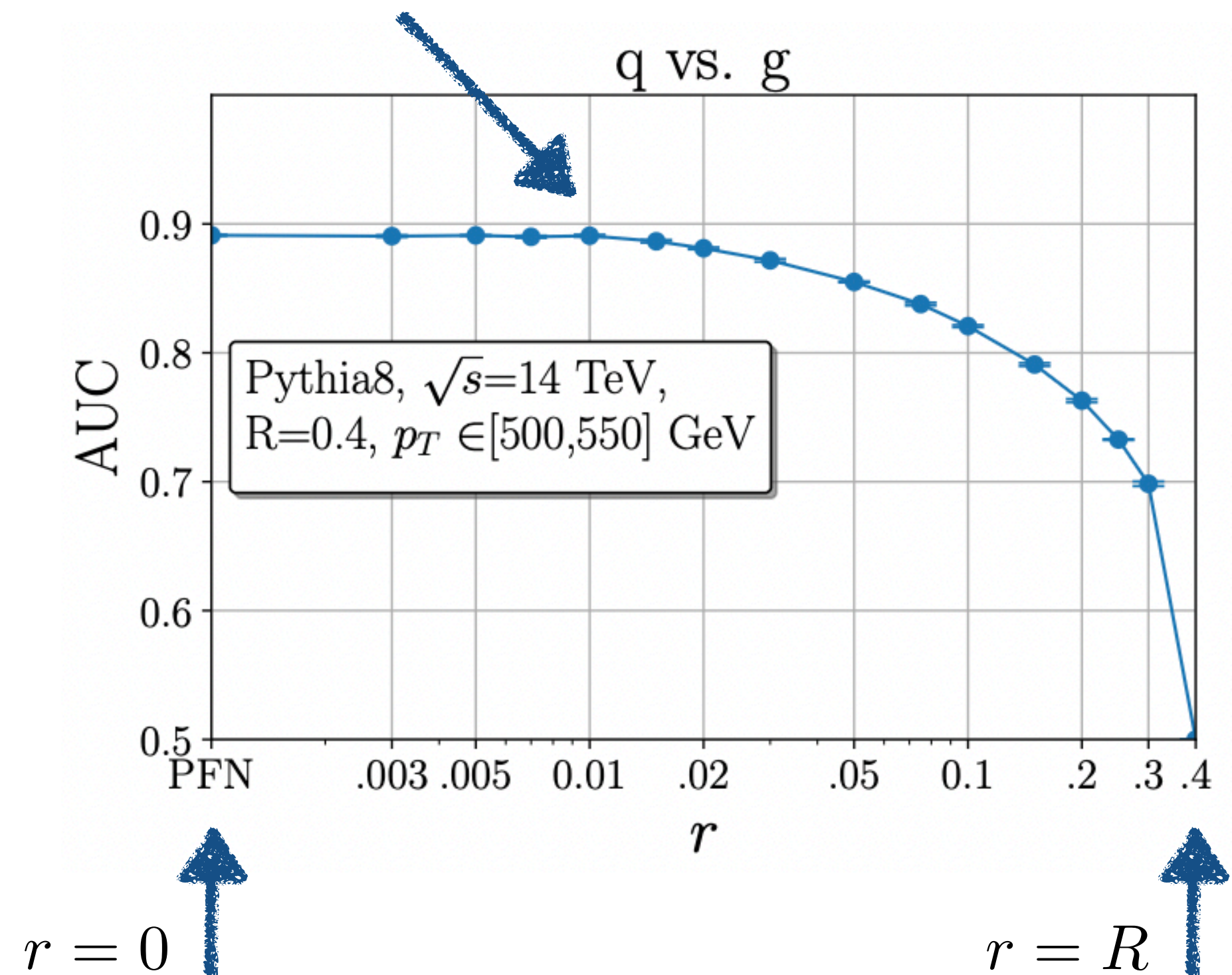


Decrease subjet radius $r$

# Is IRC-safe information all you need for jet classification?

*Athanasakos, Larkoski, Mulligan, Ploskon, FR `23*

- Cluster jet constituents into subjets

- Train deep sets on $(\eta_i, \phi_i)$ of subjets with different radii

- Max performance, IRC-unsafe limit obtained for $r \to 0$

Performance plateaus for a finite subjet radius!



$r = 0$          $r = R$

# Is IRC-safe information all you need for jet classification?

- Cluster jet constituents into subjets

- Train deep sets on $(\eta_i, \phi_i)$ of subjets with different radii

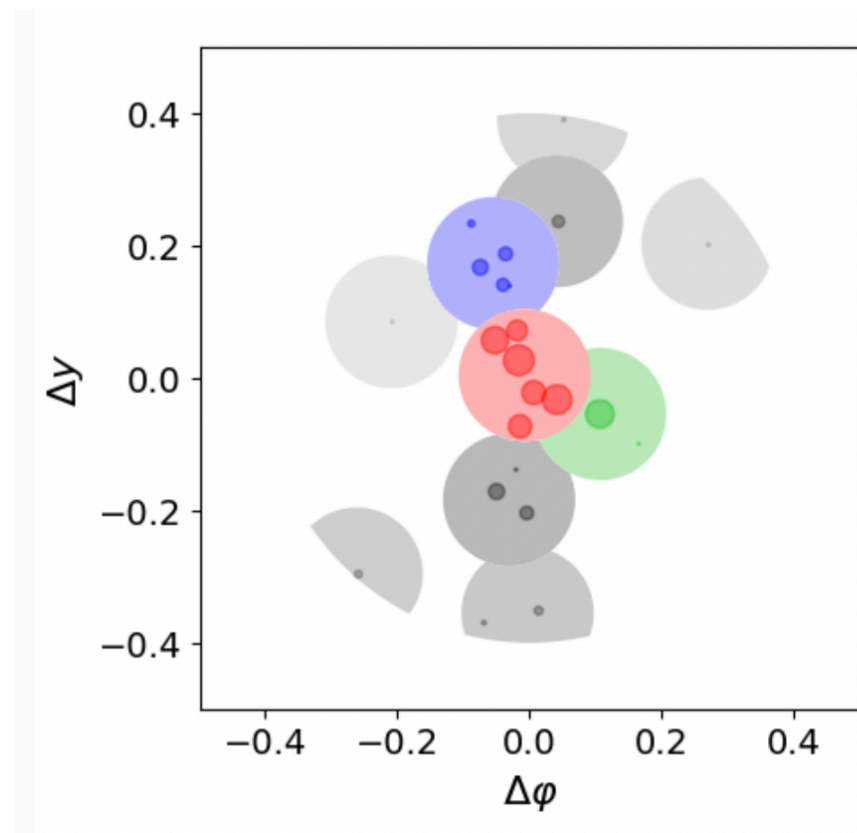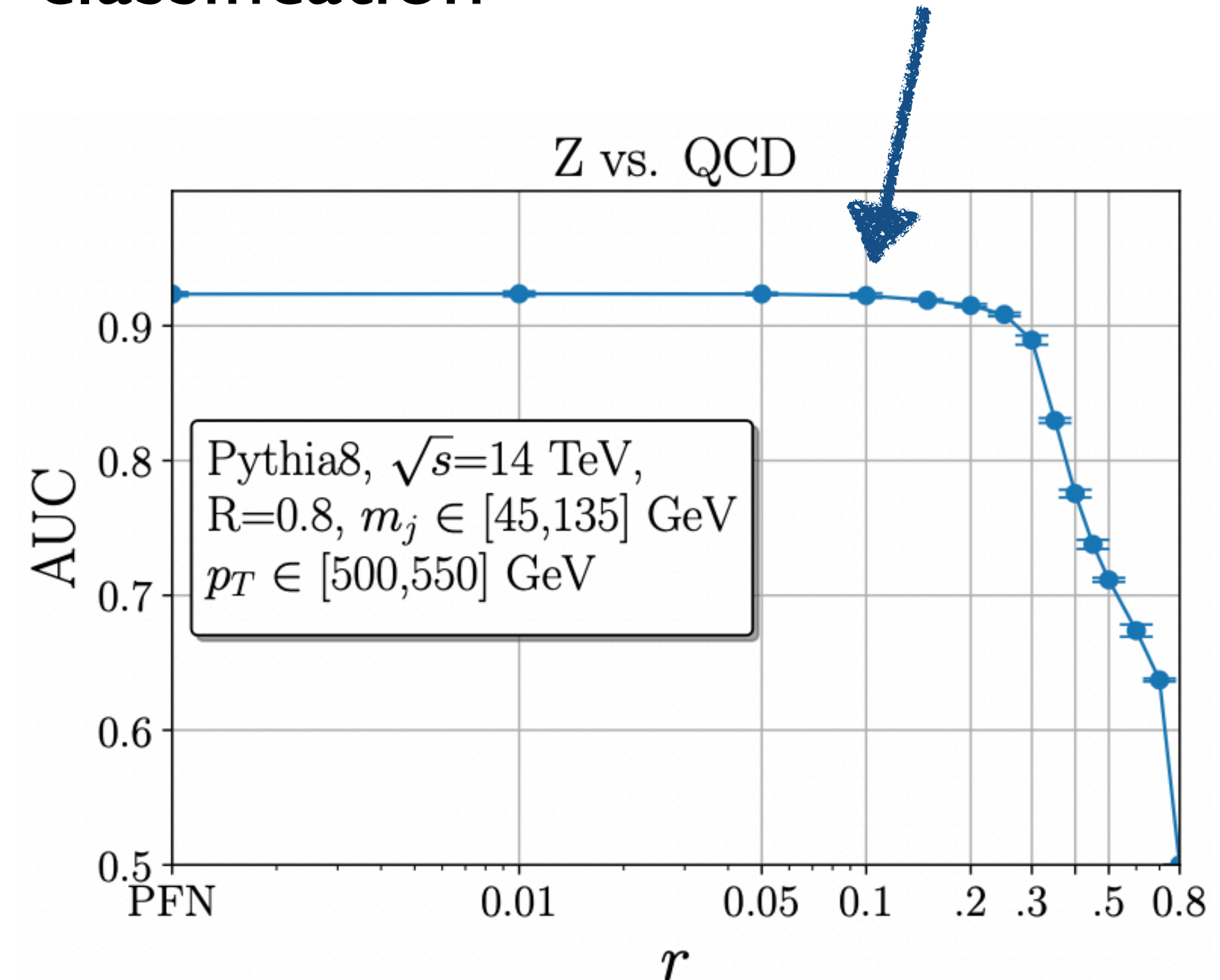- Max performance, IRC-unsafe limit obtained for $r \rightarrow 0$

Similar for QCD vs. Z jet classification

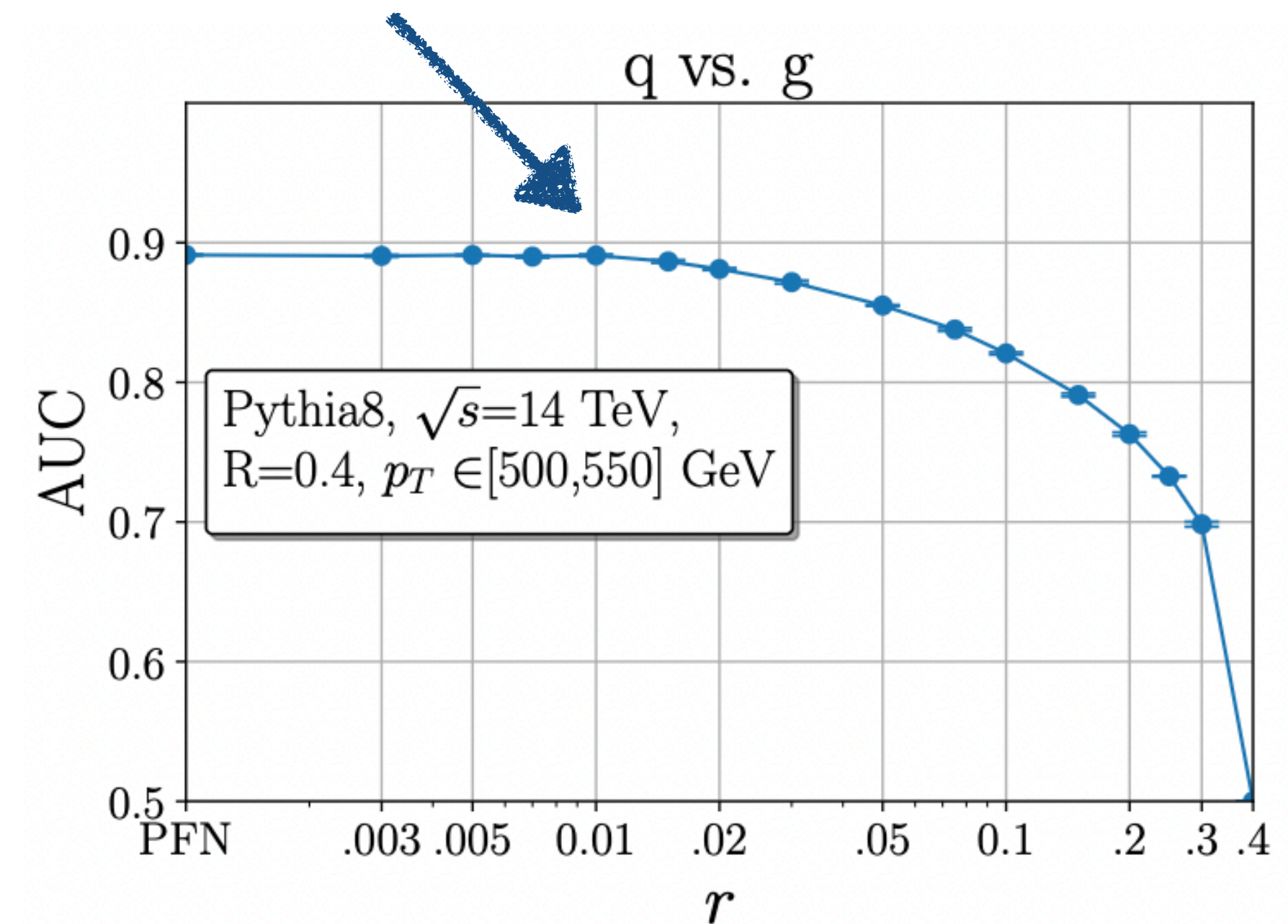# Is IRC-safe information all you need for jet classification?

*Athanasakos, Larkoski, Mulligan, Ploskon, FR `23*

- Tentatively, the answer is - Yes!

  *Theoretical perspective, see Metodiev, Larkoski `19*

- Emissions below some angular scale do not contain relevant information

- Jet Flow Networks are "gapless"

- Can identify the scale of the onset of the plateau

$$p_T \cdot r \sim 5 \text{ GeV}$$

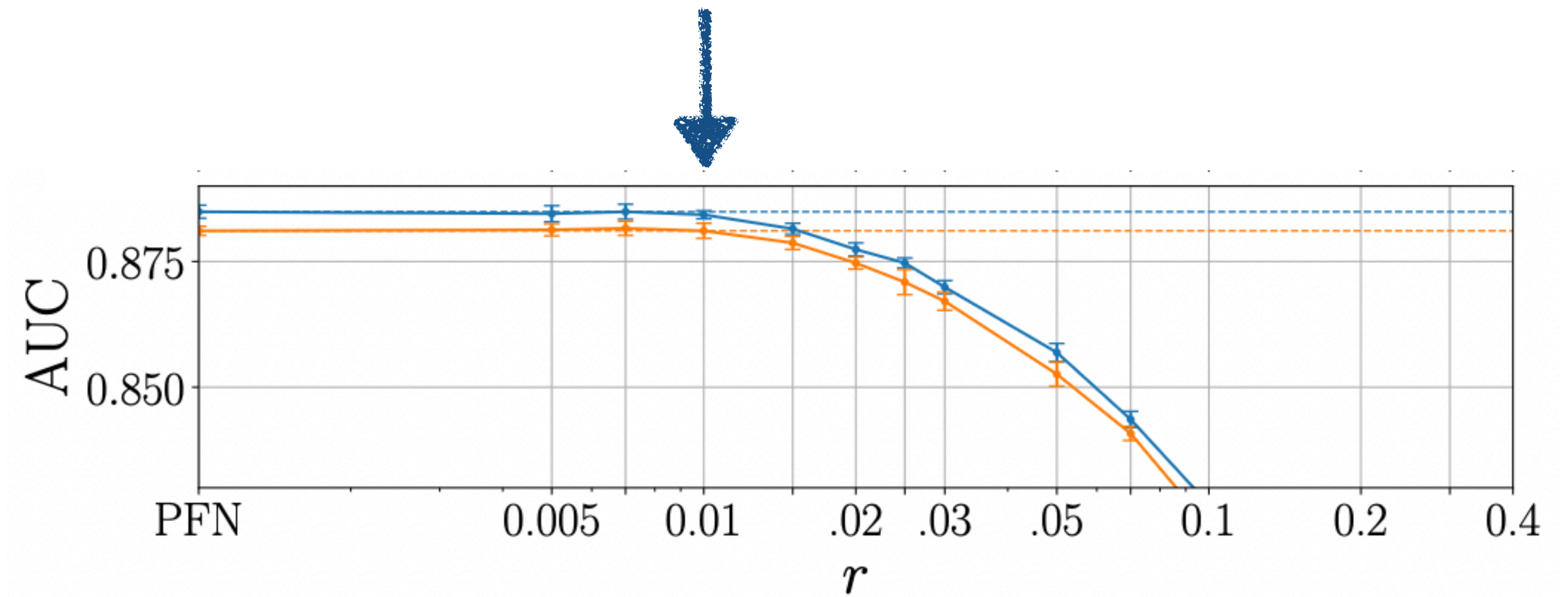# Is IRC-safe information all you need for jet classification?

*Athanasakos, Larkoski, Mulligan, Ploskon, FR `23*

- Tentatively, the answer is - Yes!

*Theoretical perspective, see Metodiev, Larkoski `19*

- Emissions below some angular scale do not contain relevant information

- Jet Flow Networks are "gapless"

- Can identify the scale of the onset of the plateau

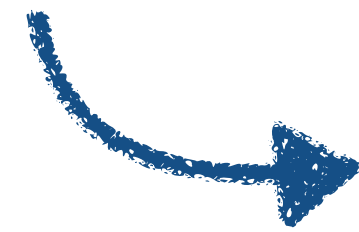$$p_T \cdot r \sim 5 \text{ GeV}$$



Scale is independent of the shower cutoff in Pythia $p_T^{\text{cut}}$

# Information content of jets & events

- N-jettiness basis  *Datta, Larkoski `17*

- The performance gap could be due to…

- ~~IRC safety?~~ ✅

  - the type of input? ❓

  - the network architecture? ❓

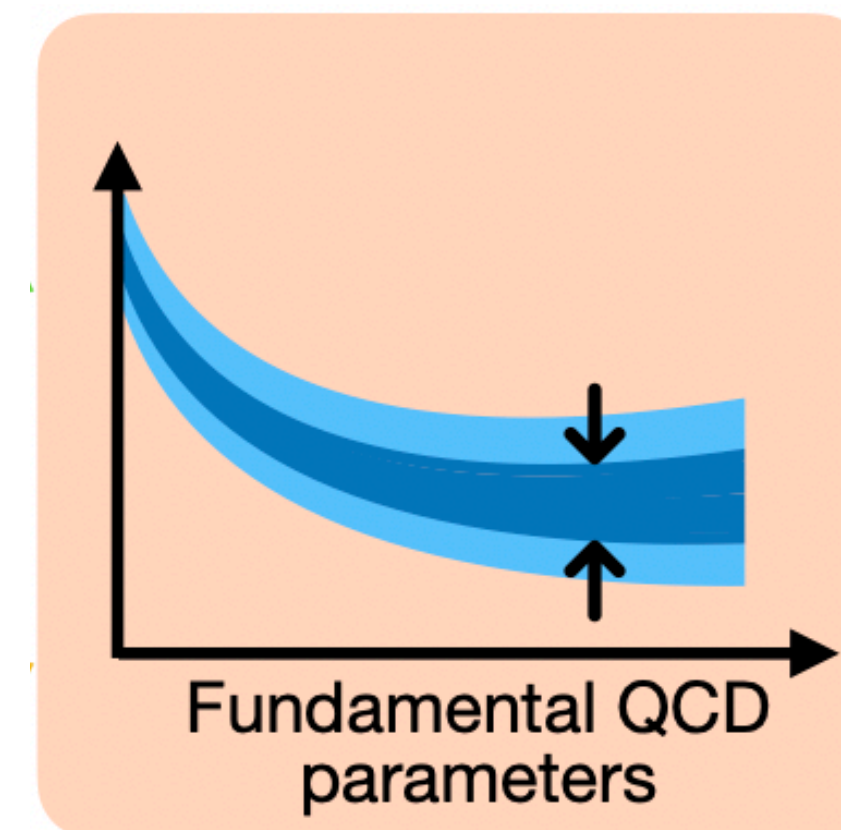Answers will provide guidance for making use of the full information

| Model | AUC |
|-------|-----|
| PFN-ID | **0.9052** $\pm$ 0.0007 |
| PFN-Ex | 0.9005 $\pm$ 0.0003 |
| PFN-Ch | 0.8924 $\pm$ 0.0001 |
| PFN | 0.8911 $\pm$ 0.0008 |
| EFN | 0.8824 $\pm$ 0.0005 |
| RNN-ID | 0.9010 |
| RNN | 0.8899 |
| EFP | 0.8919 |
| DNN | 0.8849 |
| CNN | 0.8781 |
| $M$ | 0.8401 |
| $n_{\text{SD}}$ | 0.8297 |
| $m$ | 0.7401 |

IRC unsafe classifier

⋮

N-jettiness observables

*Metodiev, Komiske, Thaler `18*



Fundamental QCD parameters

# ML for spin spin physics

- How can we apply these techniques to spin-dependent observables?

1. Supervised machine learning

2. Train on data   e.g.   $A_{UT} = \dfrac{\mathrm{d}\sigma^{\uparrow} - \mathrm{d}\sigma^{\downarrow}}{\mathrm{d}\sigma^{\uparrow} + \mathrm{d}\sigma^{\downarrow}}$

- Reformulate regression task as classification problem

$$\max_{\theta} \; |A_{UT}(\theta)|$$

→  Upper limit on what can possibly be achieved

→  Identify new observables

# Summary

- Jets will be versatile tools at the EIC

- Can take advantage of the EIC's clean environment, high luminosity etc.

- AI/ML can complement hadron structure & spin physics program

- Requires coordination with experiment

- …and can inform detector design?