# Machine learning based jet and event identification at the Electron-Ion Collider

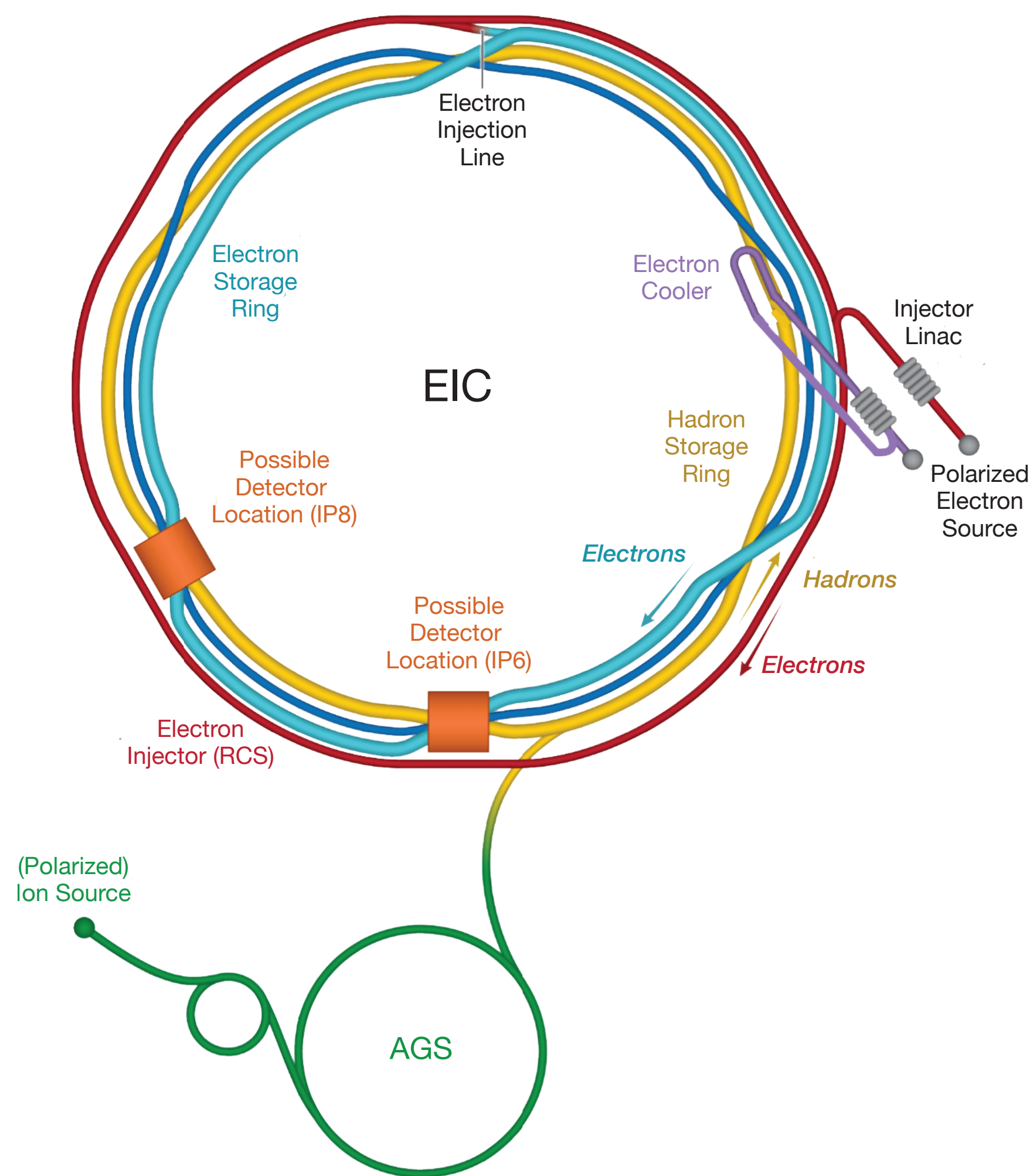## with applications to hadron structure and spin physics

K. Lee, J. Mulligan, M. Płoskoń, F. Ringer, F. Yuan

James Mulligan
UC Berkeley / LBNL

ML4Jets Workshop
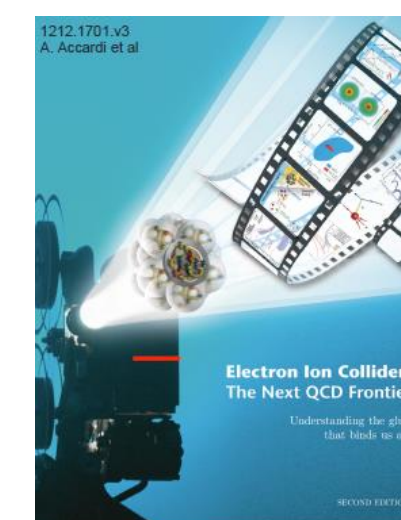Rutgers University
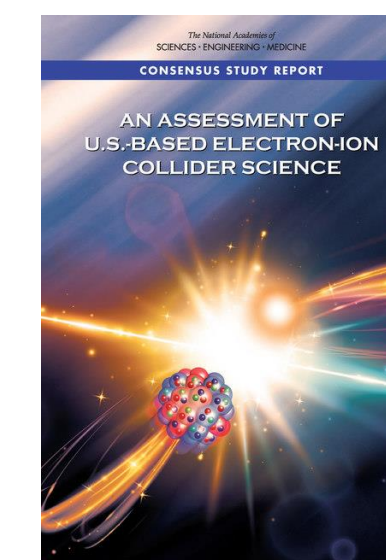Nov 3, 2022

# The Electron-Ion Collider



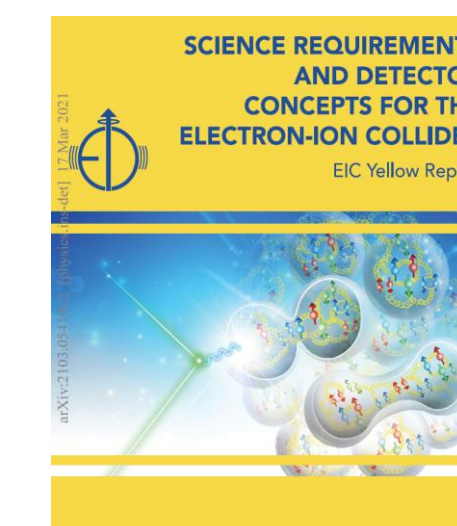**Precision QCD** with $ep$ and $eA$ collisions in the 2030s

☐ Polarized electron and proton beams
☐ Variable ion species: Au, Pb, U
☐ Variable CM energy: $20 - 140$ GeV
☐ High luminosity: $10^{33} - 10^{34}$ cm$^{-2}$s$^{-1}$

**Finding 1:** An EIC can uniquely address three profound questions about nucleons—cleons—protons—and how they are assembled to form the nuclei of atoms:

• How does the mass of the nucleon arise?
• How does the spin of the nucleon arise?
• What are the emergent properties of dense systems of gluons?



*White paper*          *NAS report*          *Yellow report*

# This talk

1. Is machine learning based jet classification useful for the science program of the EIC?

2. How will machine learning based jet taggers perform at the relatively low EIC energies?

# This talk

1. Is machine learning based jet classification useful for the science program of the EIC?

2. How will machine learning based jet taggers perform at the relatively low EIC energies?

# Constraining TMDs with jet flavor tagging

Determining the flavor of a jet allows stronger constraints on TMDs by avoiding spin asymmetry cancellations of different flavors

$u$  $d$  $s$  $c$  $g$

$q$

**Example: Sivers function (TMD PDF)**

Burkhardt sum rule:
$$\sum_{a=q,\bar{q},g} \int_0^1 \mathrm{d}x\, f_{1T}^{\perp(1)a}(x) = 0$$

If valence quarks dominate, then $u, d$ Sivers functions have large cancellation

→ Tagging $u, d$ jets separately will allow stronger constraints on Sivers function
- Recent proposal: use jet charge
- Using ML can further boost separation

*Kang, Liu, Mantry, Shao PRL 125 242003 (2020)*
*STAR, R. Fatemi EINN 2019*

# Constraining photon PDF wi

In photoproduction, the resolved processes probe the p



**direct**

*Chu, Aschenauer, Lee, Zheng PRD 96 7, 074038 (2017)*

By classifying direct vs. resolved photoproduction processes with ML, can enhance constraints on the photon PDF relative to traditional observables

e.g. $\quad x_\gamma = \dfrac{1}{2E_e y}\left(p_{T,1}e^{-\eta_1} + p_{T,2}e^{-\eta_2}\right)$

$x_\gamma^{rec}$=0.5-0.6 (×10

$x_\gamma^{rec}$=0.4-0.5 (×10

$x_\gamma^{rec}$=0.3-0.4 (×10

$x_\gamma^{rec}$=0.2-0.3 (×10

$x_\gamma^{rec}$=0.1-0.2

**Resolved**    **Direct**

# Maximizing cold nuclear matter effects

Goal: extract transport properties of nuclear matter e.g. $\hat{q}$

*Ru, Kang, Wang, Xing, Zhang, PRD 103, L031901 (2021)*
*Li, Liu, Vitev,  PLB 816, 136261 (2021)*



## Train ML classifier to distinguish $ep$ vs. $eA$ jets

Can use *interpretable* ML:

- ☐ Gain insight about type of information responsible for differences: IRC-safe vs. IRC-unsafe, hard vs. soft

- ☐ Design maximally discriminating observables that are calculable in pQCD

$$\max_{\theta} \left| \frac{d\sigma_{eA}}{d\sigma_{ep}}(\theta) - 1 \right| \longrightarrow$$





$ep$
vs.
$eA$

## Can be applied directly on experimental data

# Maximizing spin asymmetries

Goal: Measure non-zero TSSAs associated with jets:

$$A_{UT} = \frac{\mathrm{d}\sigma^{\uparrow} - \mathrm{d}\sigma^{\downarrow}}{\mathrm{d}\sigma^{\uparrow} + \mathrm{d}\sigma^{\downarrow}}$$

*STAR PRL 99, 142003 (2007)*
*STAR arXiv 2205.11800*

**Train ML classifier to distinguish ↑ vs. ↓ jets**

Can use *interpretable* ML to design maximally discriminating observables that are calculable in pQCD

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

$|A_{UT}|$ vs. $|p_{T1} + p_{T2}|$

**Can be applied directly on experimental data**

**Can be applied at RHIC now!**

$\sigma^{\uparrow}$ vs. $\sigma^{\downarrow}$

# This talk

1. Is machine learning based jet classification useful for the science program of the EIC?

2. How will machine learning based jet taggers perform at the relatively low EIC energies?

# Setup

## Event generation

PYTHIA6
- ☐ No detector simulation
- ☐ Vary minimum particle $p_T$, PID info

For u/d/s/c tagging:
LO DIS

$p_{T,\mathrm{jet}} > 10$ GeV



For quark/gluon tagging:
low-$Q^2$ photoproduction

$p_{T,\mathrm{jet1}} > 8$ GeV
$p_{T,\mathrm{jet2}} > 5$ GeV



## Machine learning model

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

Architecture: Particle Flow Networks

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

Classifier          DNNs

*Komiske, Metodiev, Thaler JHEP 01 (2019) 121*

# Jet flavor tagging: $u$ vs. $d$

## $u$ vs. $d$ jets



Legend:
- Particle Flow Network (w/ PID)
- Particle Flow Network (w/ charge)
- Particle Flow Network (w/o PID,charge)
- Jet charge, $\kappa = 0.3$
- Jet charge, $\kappa = 0.5$
- Jet charge, $\kappa = 0.7$

True Positive Rate $= \dfrac{\text{True } d}{\text{Total } d}$

False Positive Rate $= \dfrac{\text{False } d}{\text{Total } u}$

Jet charge

$$Q_\kappa = \sum_{i \in \text{jet}} z_i^\kappa Q_i$$

☐ ML outperforms jet charge
  ☐ Charge information is crucial
  ☐ Full PID does not gain much

# Out-of-jet information

Does including out-of-jet information boost the jet flavor tagging performance?



$u$ vs. $d$ jets

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

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

| | |
|---|---|
| in-jet | ($p_{T,\text{particle}} > 0.1$ GeV) |
| in-jet + out-of-jet | ($p_{T,\text{particle}} > 0.1$ GeV) |
| in-jet | ($p_{T,\text{particle}} > 0.4$ GeV) |
| in-jet + out-of-jet | ($p_{T,\text{particle}} > 0.4$ GeV) |
| ⋯⋯ Jet charge, $\kappa = 0.3$ | |

| | |
|---|---|
| Particle Flow Network | ($p_{T,\text{particle}} > 0.1$ GeV) |
| Particle Flow Network | ($p_{T,\text{particle}} > 0.2$ GeV) |
| Particle Flow Network | ($p_{T,\text{particle}} > 0.4$ GeV) |

**Significant gain from out-of-jet information**
☐ Due to soft particles $0.1 < p_T < 0.4$ GeV

# Jet flavor tagging: $ud$ vs. $s$

$u, d$ vs. $s$ jets



Legend:
- 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

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

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

Random classifier

**For strange: ML dramatically outperforms jet charge**
- ☐ PID gives huge boost

We use precision-recall metric since there are ~40x more $ud$ than $s$
- ☐ Precision ↔ Purity
- ☐ Recall ↔ Efficiency

# Quark vs. gluon jet tagging



*q* vs. *g* jet

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

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

*Komiske, Metodiev, Thaler, JHEP 01 (2019) 121*

*Komiske, Metodiev, Thaler, JHEP 04, 013 (2018)*

ML performance not as good as at LHC, but still reasonably good

| AUC | EIC | LHC |
|---|---|---|
| **Particle Flow Network** | 0.79 | 0.91 |
| **Energy Flow Network** | 0.76 | 0.88 |
| **Energy Flow Polynomials** | 0.75 | 0.89 |

# Hard process tagging



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



We classify hard processes generating $qq/q\bar{q}$ vs. $gg$ di-jets:

$$qq \rightarrow qq, q\bar{q} \rightarrow q\bar{q}, gg \rightarrow q\bar{q}, \gamma_T^* g \rightarrow q\bar{q}, \gamma_L^* g \rightarrow q\bar{q}$$

**vs.**

$$q\bar{q} \rightarrow gg, gg \rightarrow gg,$$

➤ Can use this method to tag resolved photoproduction contributions

**Significant improvement when adding subleading jet and out-of-jet particles**

# Summary

## Machine learning can improve access to hadron structure and spin physics at the EIC

- ☐ Improve jet flavor tagging performance: constrain TMDs, photon PDF, …
- ☐ Maximize the size of spin asymmetries or cold nuclear matter effects — train directly on data

## PYTHIA6 indicates that classification performance remains reasonably good at EIC

- ☐ Large performance boost from ML for strange and charm tagging when PID is included
- ☐ Large performance boost by including soft, out-of-jet particles
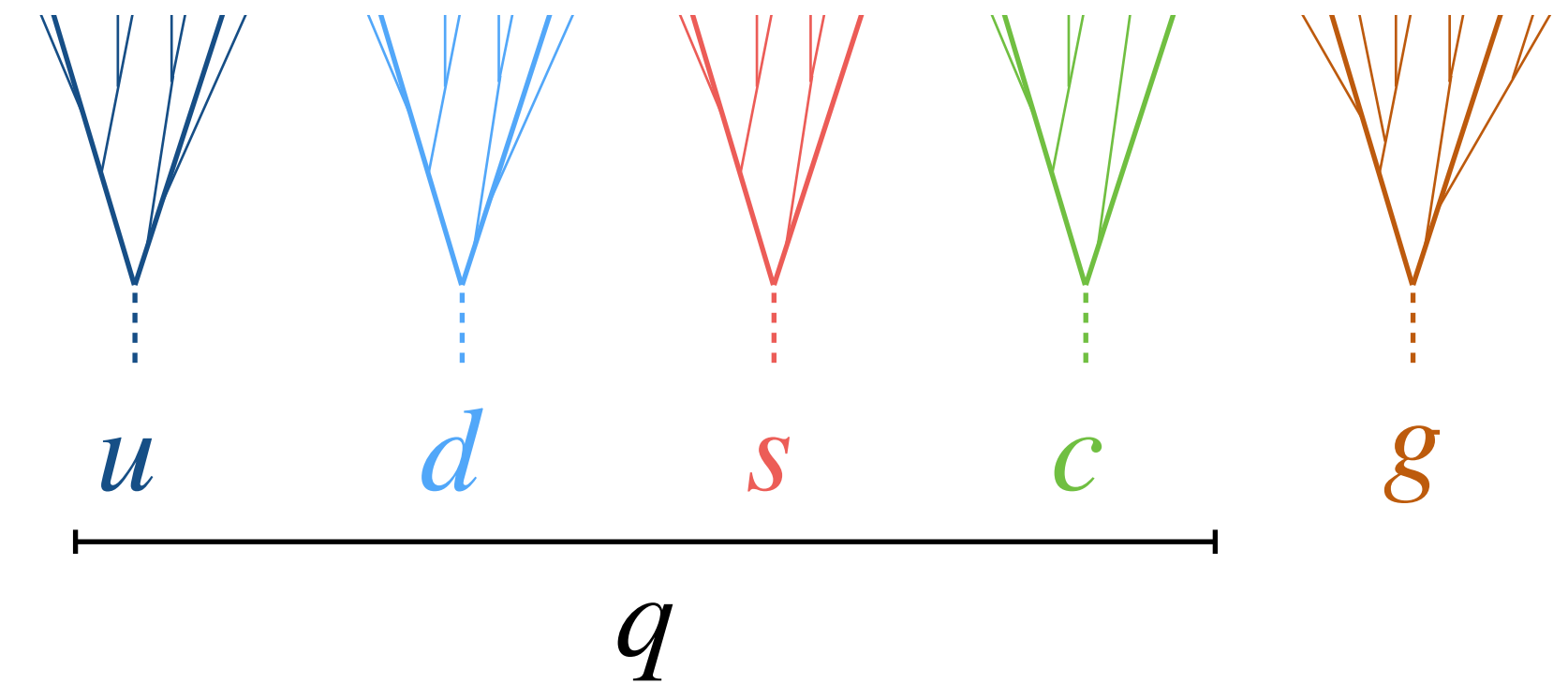
## Outlook: Study model-dependence and connect ML results to theory

- ☐ Design analytically tractable observables and/or incorporate classifiers into global fits
- ☐ Explore ML architectures — data set to be made public soon

# backup

# Constraining TMDs with jet flavor tagging

Determining the flavor of a jet allows stronger constraints on TMDs by avoiding spin asymmetry cancellations of different flavors



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

$q$

## Example: Collins fragmentation function

Schäfer-Teryaev sum rule: $$\sum_h \int_0^1 \mathrm{d}z \, H_{1,h/q}^{\perp(3)}(z) = 0$$

One usually measures identified hadrons to avoid e.g. $\pi^+$ cancellation with $\pi^-$

However the fragmentation functions still contain large parton flavor cancellations:

$$\int_0^1 \mathrm{d}z \left( H_{1,\pi^+/u}^{\perp(3)}(z) + H_{1,\pi^+/d}^{\perp(3)}(z) \right) \approx 0$$

→ Tagging jet flavor will allow stronger constraints on Collins fragmentation function

# Additional applications of jet flavor tagging

- Longitudinally polarized gluon distribution $\Delta g$ — quark flavor and quark vs. gluon

    *Zhou, Sato, Melnitchouk (JAM), PRD 105, 074022 (2022)*

- Gluon Sivers function — quark vs. gluon

    *Zheng, Aschenauer, Lee, Xiao, Yin, PRD 98, 034011 (2018)*
    *Liu, Ringer, Vogelsang, Yuan, PRL122, 192003 (2019)*

- Strange quark PDF — charm tagging

    *Arratia, Furletova, Hobbs, Olness, Sekula, PRD 103, 074023 (2021)*

- BSM searches — quark flavor

    *Li, Yan, Yuan, arXiv:2112.07747*

# Dependence on minimum particle $p_T$



$u$ vs. $d$ jets

True Positive Rate $= \dfrac{\text{True } d}{\text{Total } d}$

False Positive Rate $= \dfrac{\text{False } d}{\text{Total } u}$

Particle Flow Network    ($p_{T,\,\text{particle}} > 0.1$ GeV)

Particle Flow Network    ($p_{T,\,\text{particle}} > 0.2$ GeV)

Particle Flow Network    ($p_{T,\,\text{particle}} > 0.4$ GeV)

# Direct vs. resolved photon tagging



*direct* vs. *resolved* process

True Positive Rate $= \dfrac{\text{True } direct}{\text{Total } direct}$

False Positive Rate $= \dfrac{\text{False } direct}{\text{Total } resolved}$

Particle Flow Network (w/ PID)
Particle Flow Network (w/o PID,charge)

# Jet flavor tagging: *uds* **vs.** *c*

$u, d, s$ vs. $c$ jets



Legend:
- Particle Flow Network (w/ PID)
- Particle Flow Network (w/ charge)
- Particle Flow Network (w/o PID,charge)
- Jet charge, $\kappa = 0.3$
- Jet charge, $\kappa = 0.5$
- Jet charge, $\kappa = 0.7$

y-axis: $\text{Precision} = \dfrac{\text{True } c}{\text{True } c + \text{False } c}$

x-axis: $\text{Recall} = \dfrac{\text{True } c}{\text{Total } c}$

Random classifier
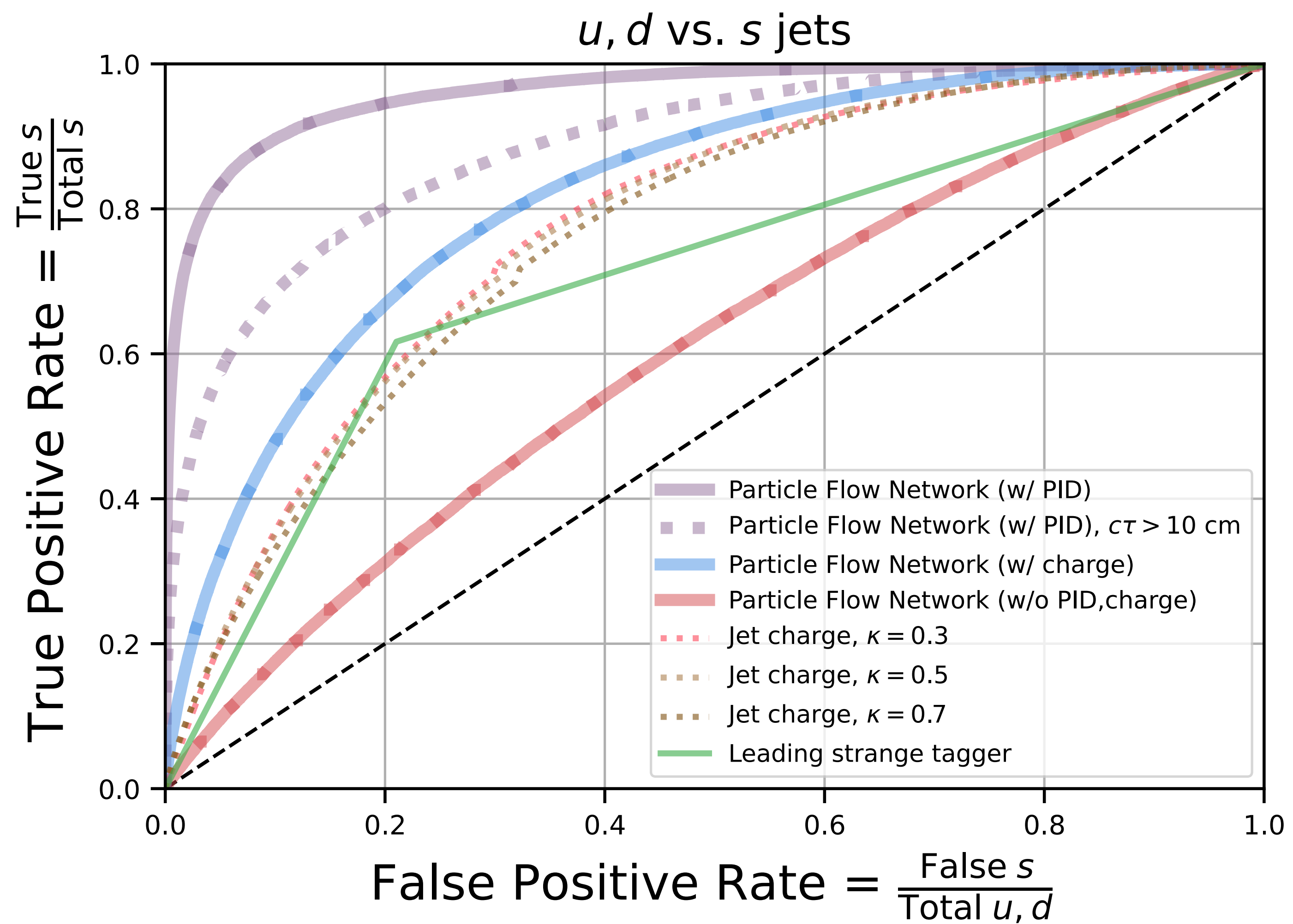
For charm: fragmentation pattern increasingly important, but PID is crucial

We use precision-recall metric since there are ~20x more *uds* than *c*
- Precision ↔ Purity
- Recall ↔ Efficiency

# ud vs. s



$u, d$ vs. $s$ jets

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

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

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

- 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

# uds vs. c

## $u, d, s$ vs. $c$ jets



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

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

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

Legend:
- Particle Flow Network (w/ PID)
- Particle Flow Network (w/ charge)
- Particle Flow Network (w/o PID,charge)
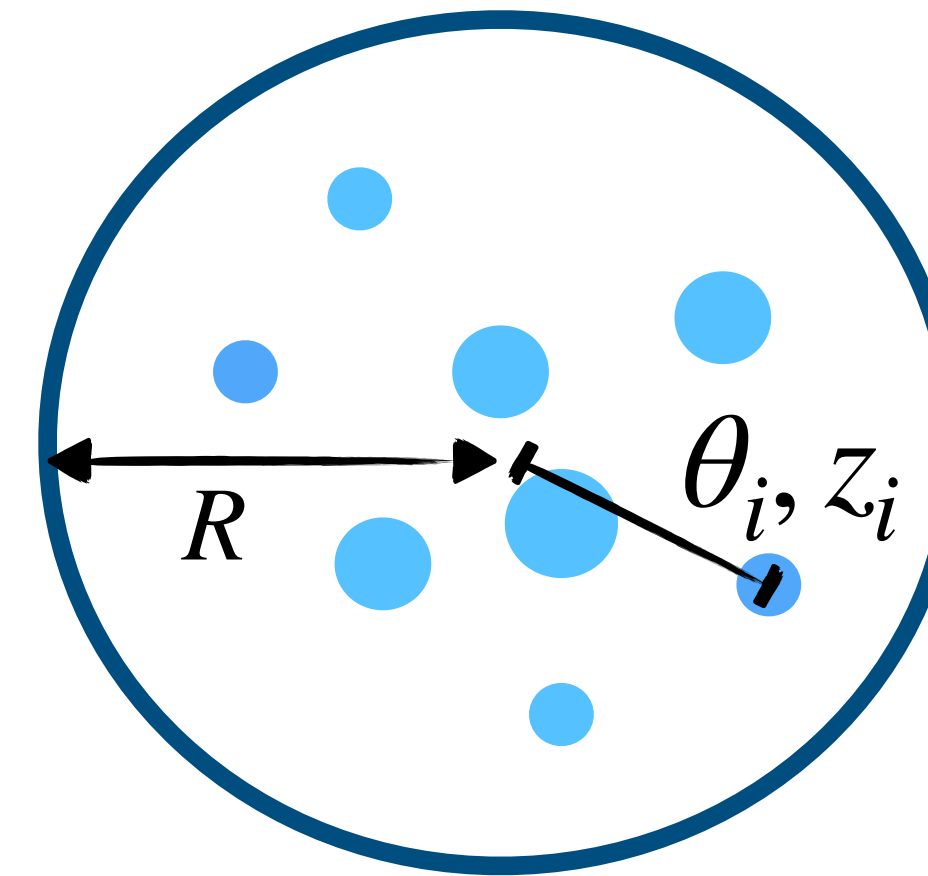- Jet charge, $\kappa = 0.3$
- Jet charge, $\kappa = 0.5$
- Jet charge, $\kappa = 0.7$

# Jet observables and IRC safety

We are free to construct any observable
from the jet's constituents

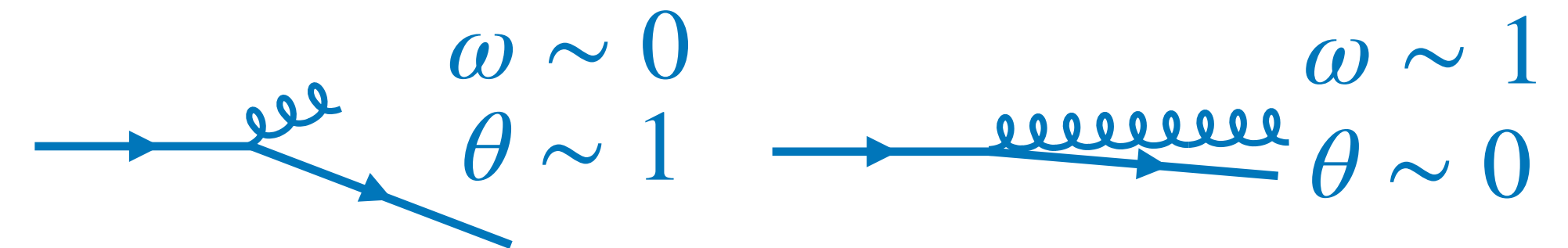e.g. $\quad \lambda_\alpha^\kappa = \sum_{i \in \text{jet}} z_i^\kappa \theta_i^\alpha$

$$\theta_i = \frac{\sqrt{\Delta y^2 + \Delta \varphi^2}}{R}$$

$$z_i = \frac{p_{\text{T},i}}{p_{\text{T,jet}}}$$

$\theta_i, z_i$

$R$

However, usually only those combinations that
obey **infrared-collinear (IRC) safety**
are calculable in perturbative QCD

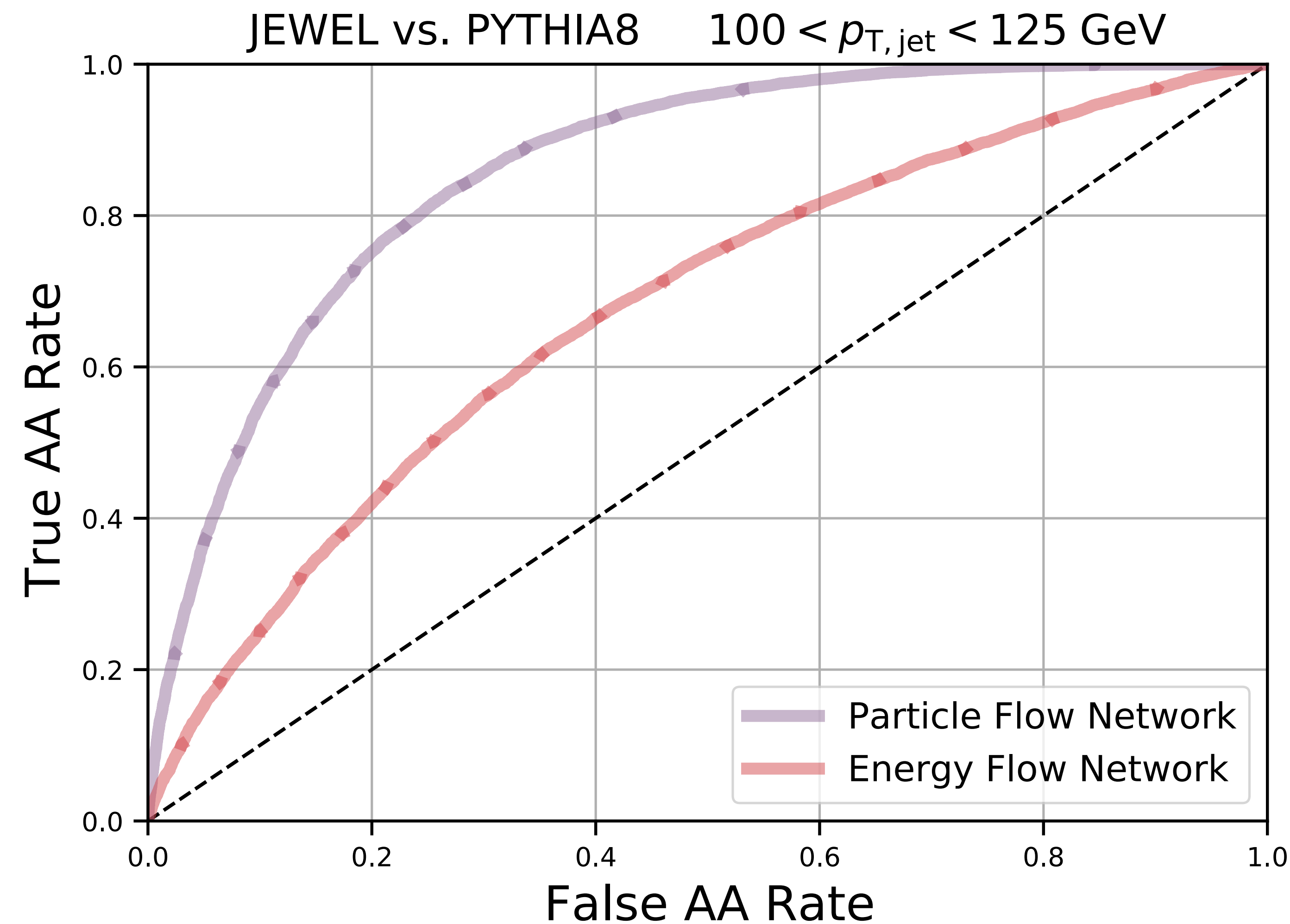e.g. $\quad \lambda_{\alpha>0}^{\kappa=1} = \sum_{i \in \text{jet}} z_i \theta_i^\alpha$

$\omega \sim 0$
$\theta \sim 1$

$\omega \sim 1$
$\theta \sim 0$

*Insensitive to soft/collinear emissions*

# IRC-safe vs. IRC-unsafe physics

*Lai, Mulligan, Płoskoń, Ringer JHEP 10 (2022) 011*

JEWEL vs. PYTHIA8    $100 < p_{\text{T, jet}} < 125$ GeV



We compare the IRC-unsafe network (PFN) to an IRC-safe network (EFN)

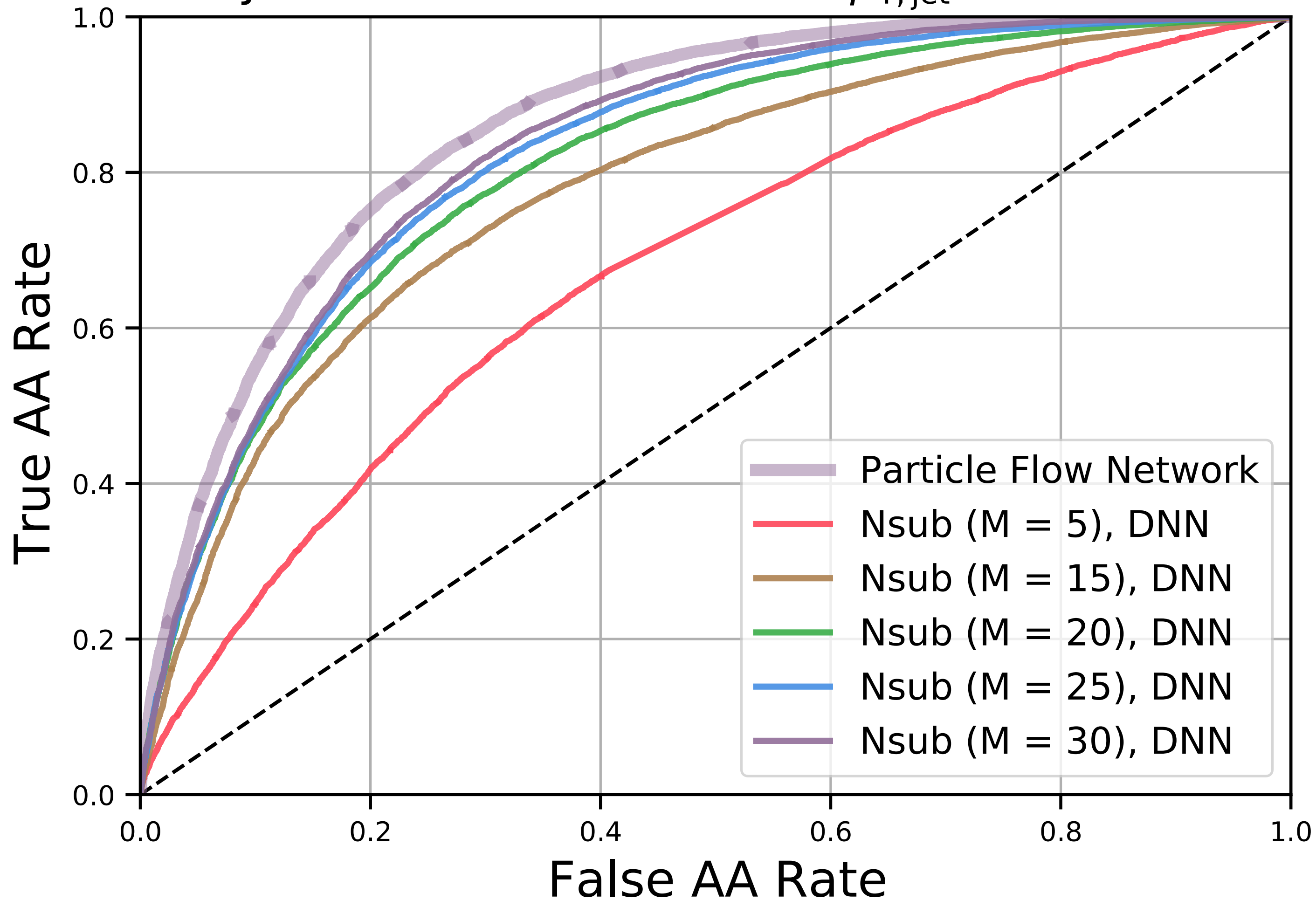$$f(p_1, \ldots, p_M) = F\left( \sum_{i=1}^{M} z_i \Phi\left(\hat{p}_i\right) \right)$$

Classifier          DNNs

IRC-unsafe information contains significant discriminating power

# Hard vs. soft physics

JEWEL vs. PYTHIA8    $100 < p_{\mathrm{T, jet}} < 125$ GeV

Legend:
- Particle Flow Network
- Nsub ($M = 5$), DNN
- Nsub ($M = 15$), DNN
- Nsub ($M = 20$), DNN
- Nsub ($M = 25$), DNN
- Nsub ($M = 30$), DNN

X-axis: False AA Rate
Y-axis: True AA Rate

How many observables does one need to measure to saturate information?

DNN with $3M - 4$ *N*-subjettiness basis observables as input:

$$\left\{ \tau_1^{(0.5)}, \tau_1^{(1)}, \tau_1^{(2)}, \tau_2^{(0.5)}, \tau_2^{(1)}, \tau_2^{(2)}, \ldots, \tau_{M-2}^{(0.5)}, \tau_{M-2}^{(1)}, \tau_{M-2}^{(2)}, \tau_{M-1}^{(0.5)}, \tau_{M-1}^{(1)} \right\}$$
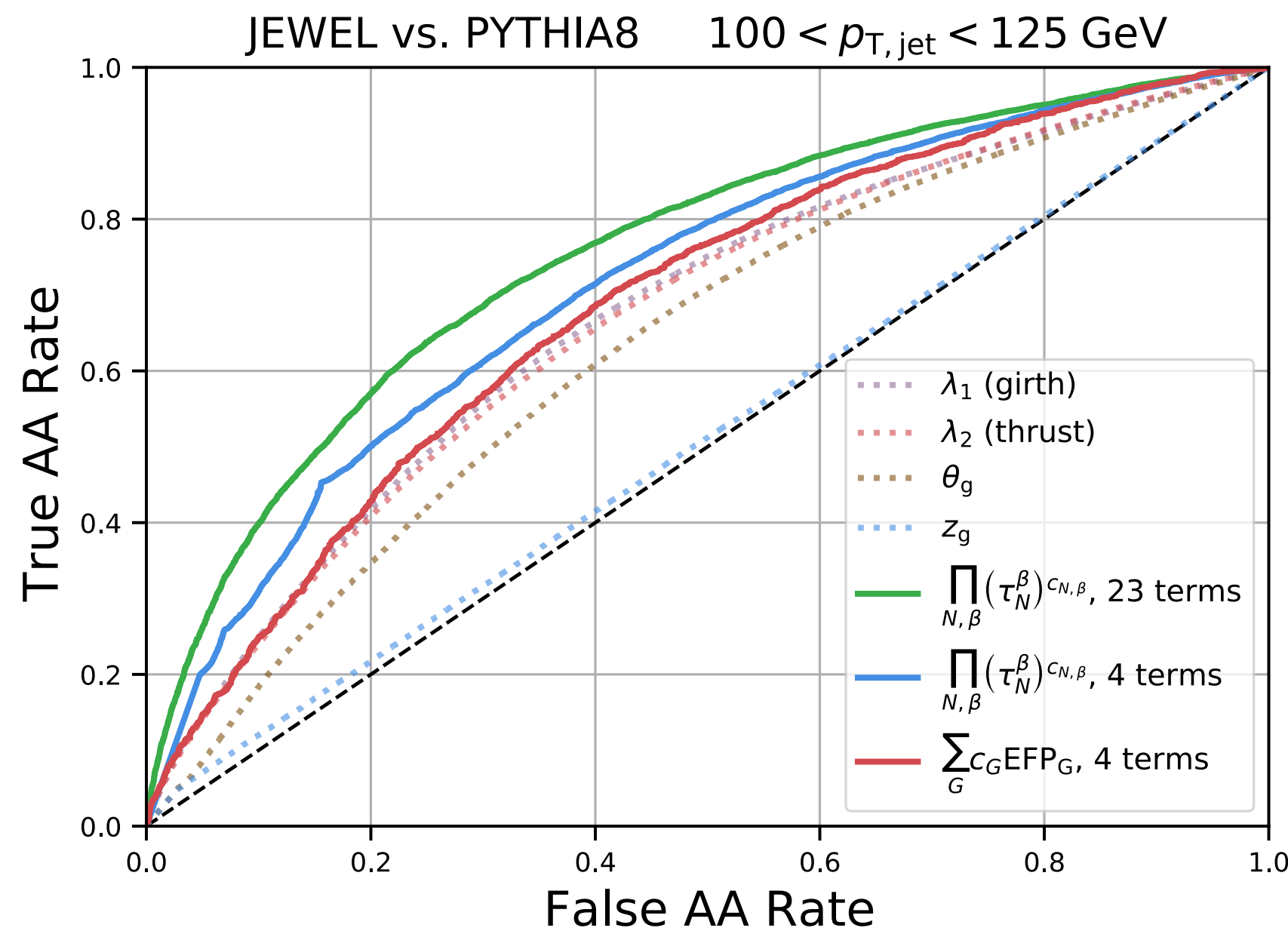
Significant information in quenched jets up to $M \approx 25$

# Observable design

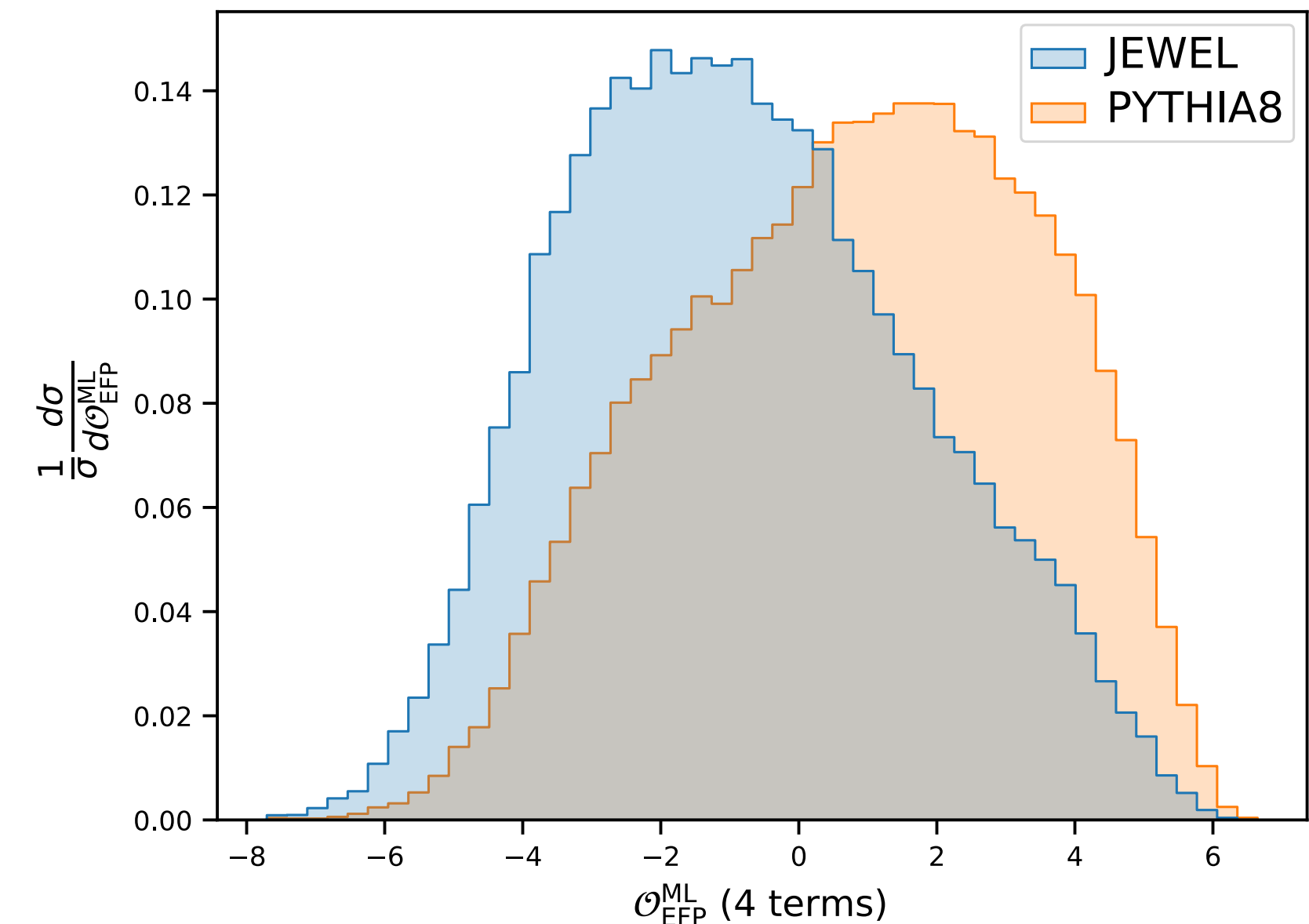*Lai, Mulligan, Płoskoń, Ringer JHEP 10 (2022) 011*

## By balancing the tradeoff of discriminating power and complexity, we can design the *most strongly modified* calculable observable



Approximate classifier with small number of features

$\longrightarrow$

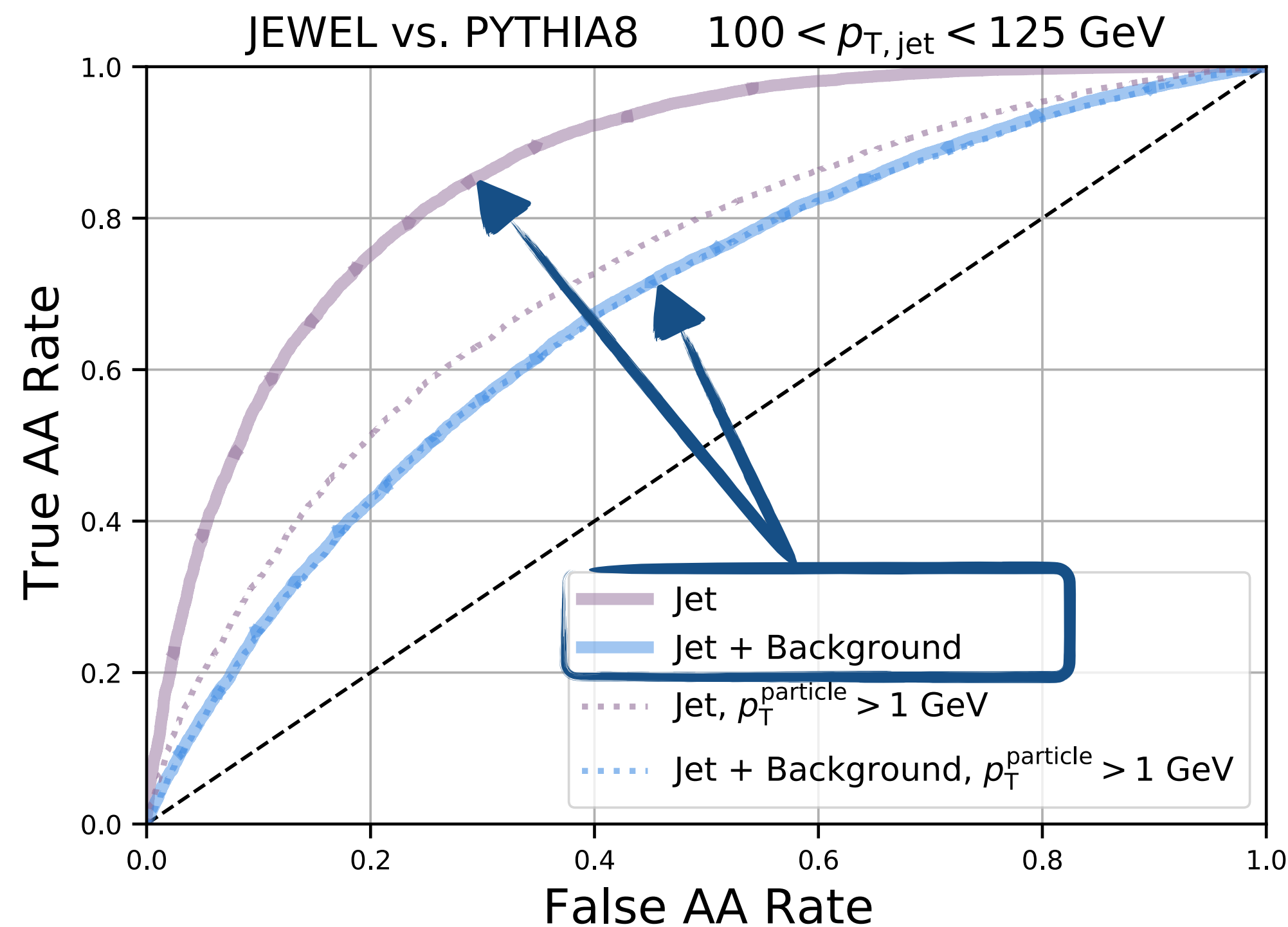"Symbolic regression" using Lasso

**ML-assisted observable design provides guidance to experiments and theory — can then measure and calculate designed observables using traditional methods**
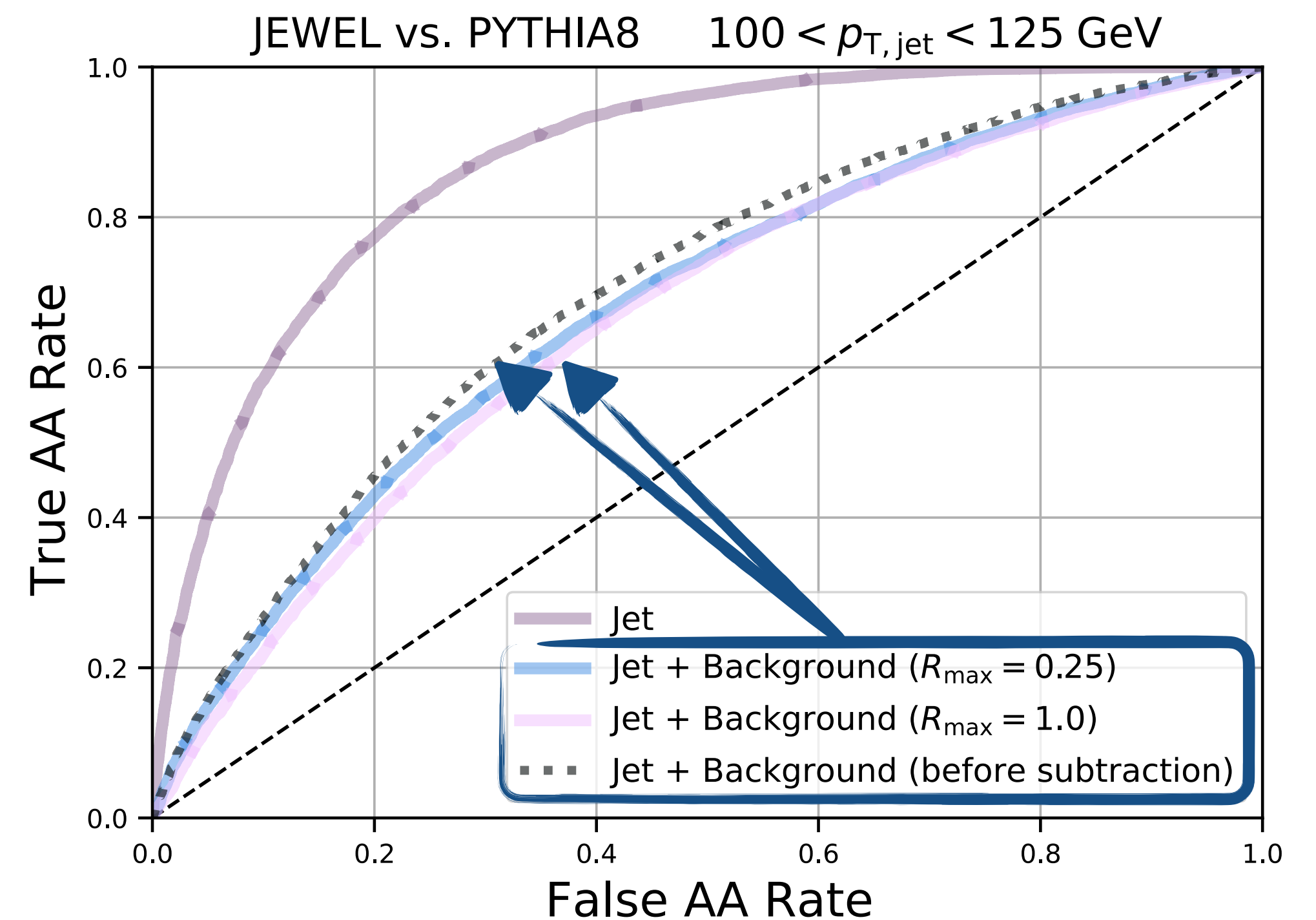
# Information loss due to background

*Lai, Mulligan, Płoskoń, Ringer JHEP 10 (2022) 011*

**Discriminating power is highly reduced by the fluctuating underlying event**



**Background subtraction algorithms remove small but significant information**



Delicate challenge: soft information crucial, yet background prevents from being accessed

New metric to assess background subtraction algorithms