

# Machine learning for the analysis of hard probes

Hannah Bossi (MIT)  
Hard Probes 2024  
Nagasaki, Japan  
September 27th, 2024



HP2024  
NAGASAKI

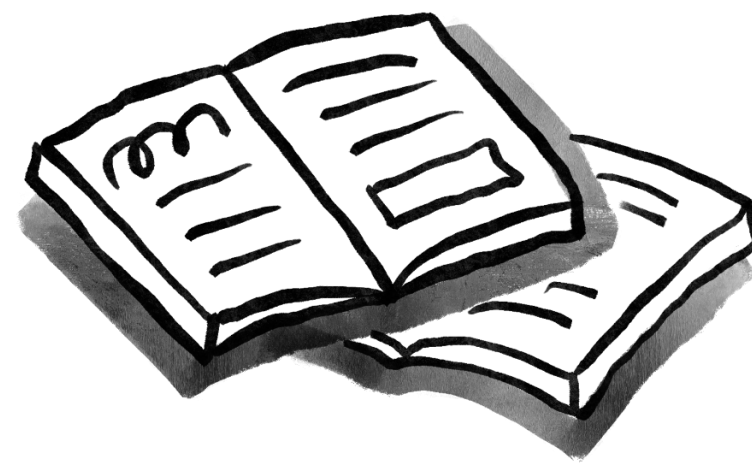


MIT HIG's work was  
supported by US DOE-NP

 Laboratory for  
Nuclear Science

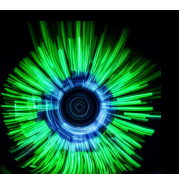
# ROADMAP

**WHAT IS AI/ML  
AND WHY IS IT  
USEFUL FOR THE  
ANALYSIS OF  
HARD PROBES?**



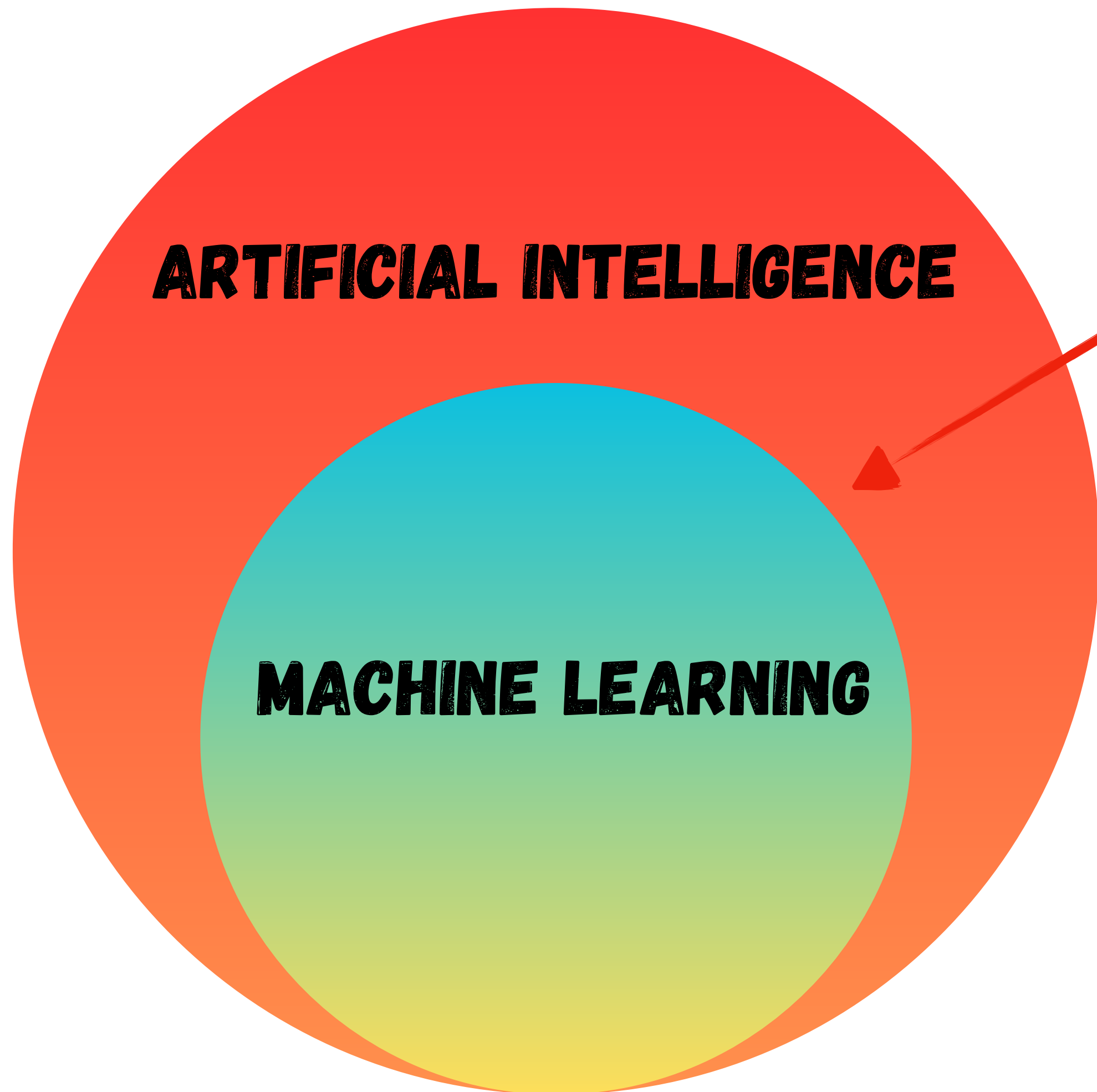
**HOW IS AI/ML  
CURRENTLY  
BEING USED FOR  
ANALYSIS?**

**WHERE ARE WE  
HEADING?**



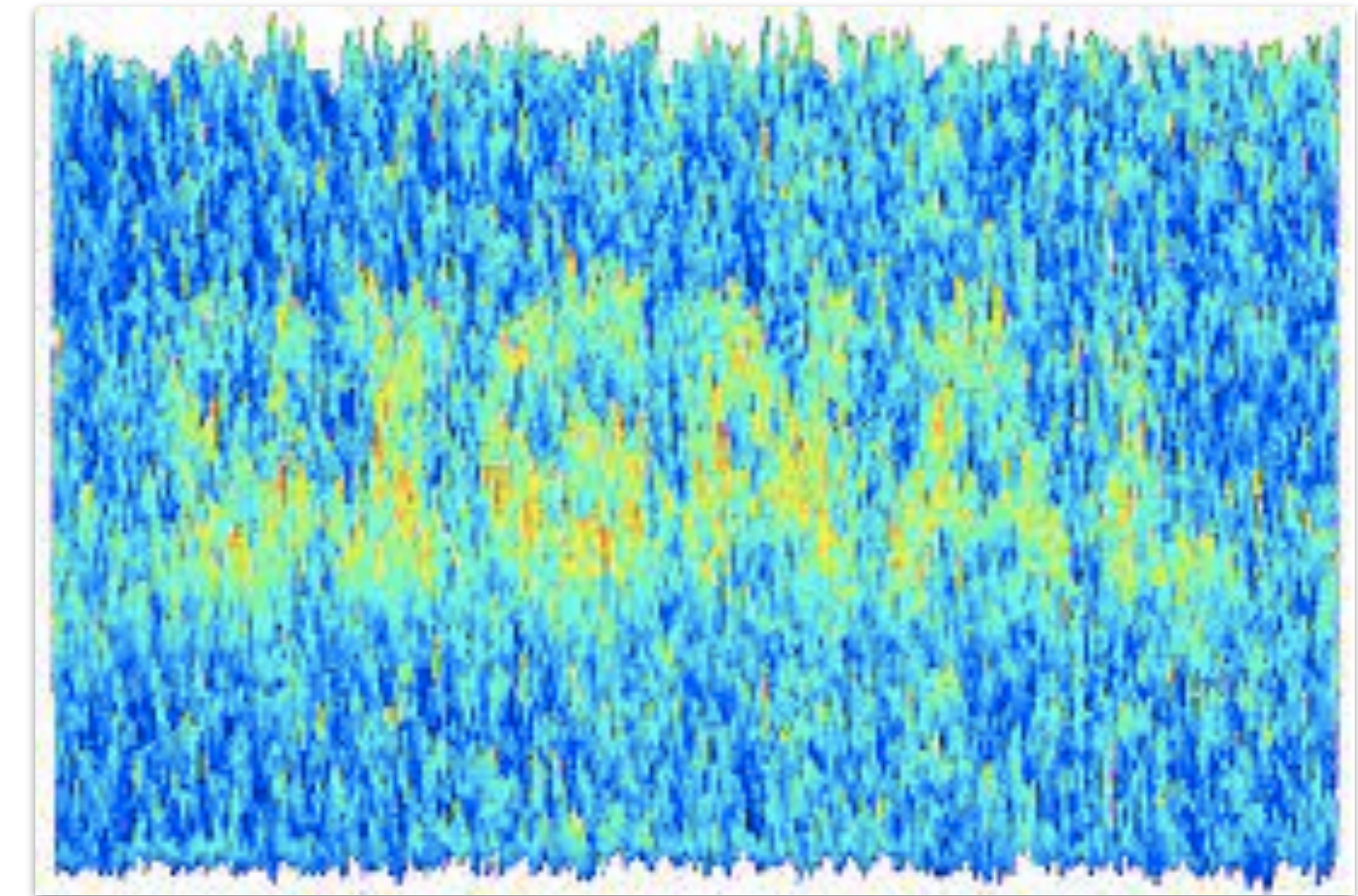
# WHAT IS AI/ML?

**Artificial Intelligence:** Programs with the ability to acquire and apply knowledge and skills.

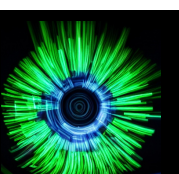


Ex: Chatbots (humans give rules)

**Machine Learning:** algorithms that imitate human learning, i.e. gradually improving accuracy over time.



At its core, pattern recognition → humans can do this by eye!



# HOW DOES THE MACHINE LEARN?

## SUPERVISED LEARNING

Algorithm learns from a labeled set of “true values”.

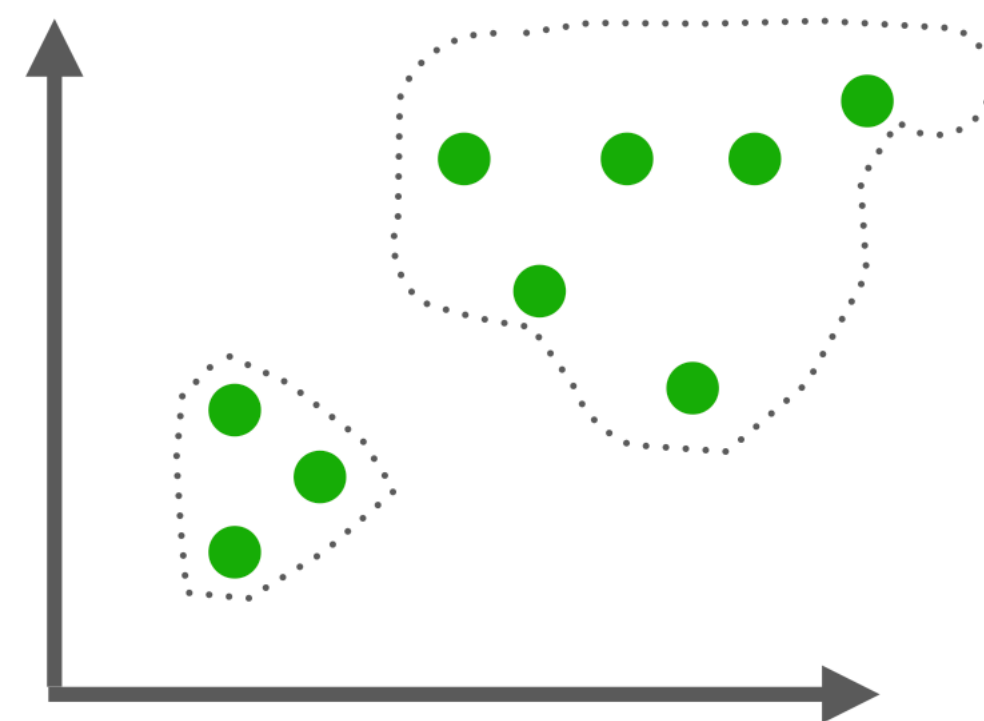


Driven by the Task

Analogy: Taking a test

## UNSUPERVISED LEARNING

Algorithm finds structure in the data without knowing the desired outcome.



Driven by the Data

Analogy: Clustering

## REINFORCEMENT LEARNING

Algorithm learns in a reward based system to determine a series of actions.



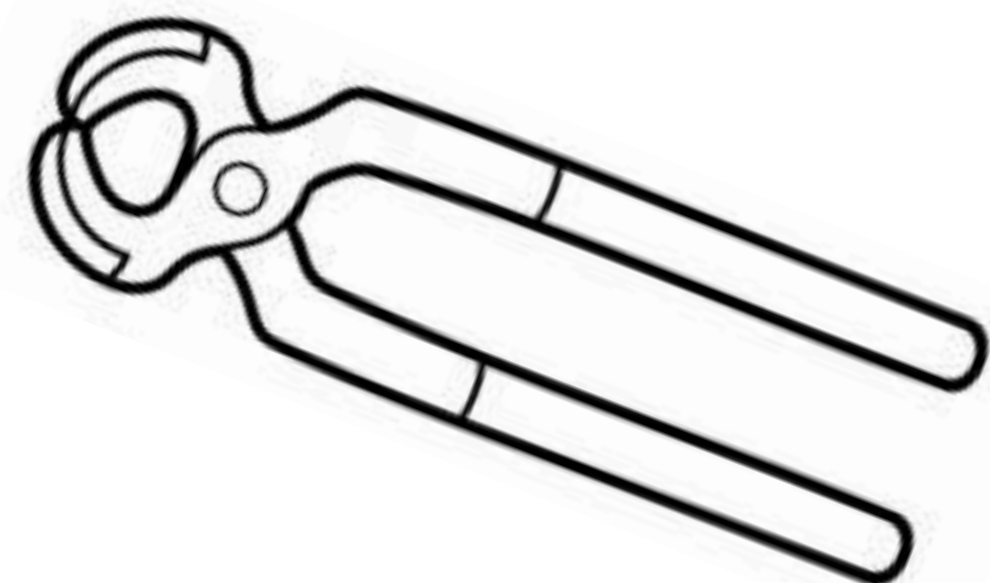
Driven by the Reward

Analogy: Dog training



# WHAT KINDS OF ML TOOLS ARE THERE?

**BOOSTED  
DECISION TREES  
(BDTS)**



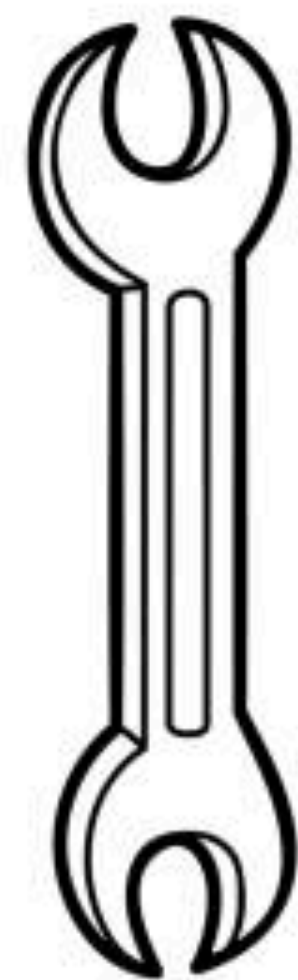
**RANDOM FORESTS**



**CONVOLUTIONAL  
NEURAL  
NETWORKS (CNN)**

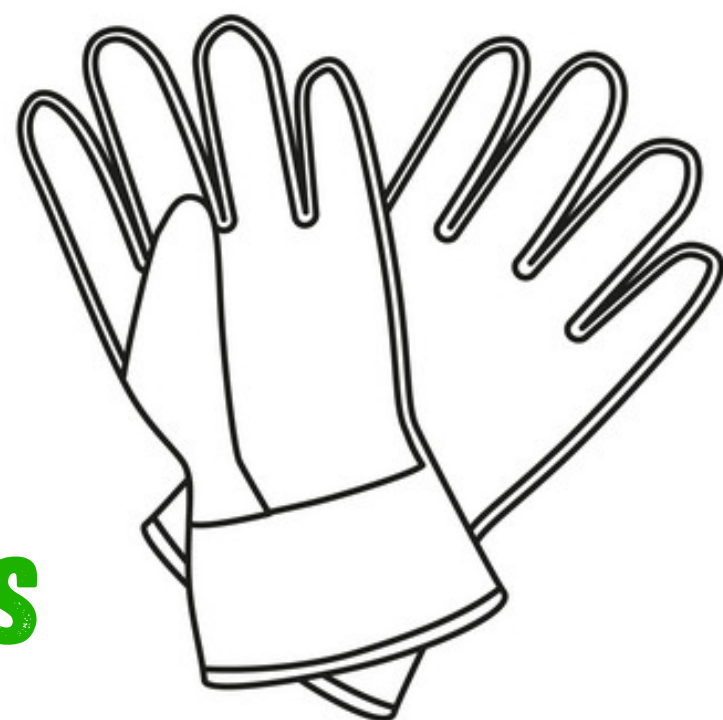


**NEURAL  
NETWORKS (NN)**

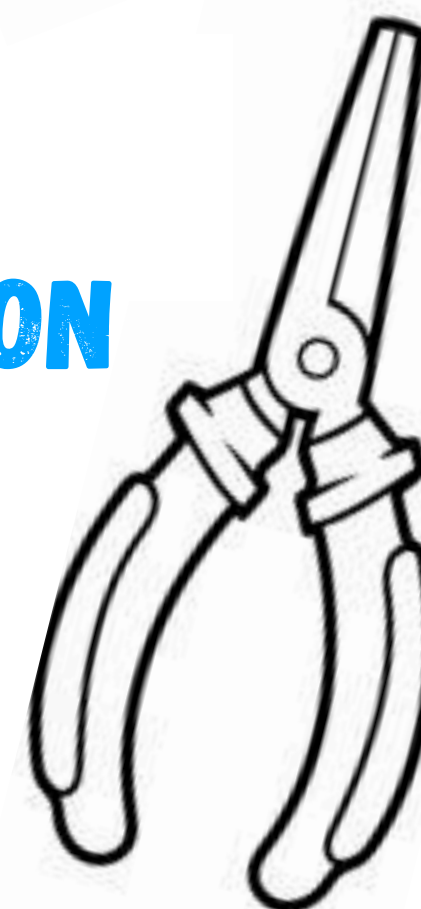


**NORMALIZING  
FLOWS**

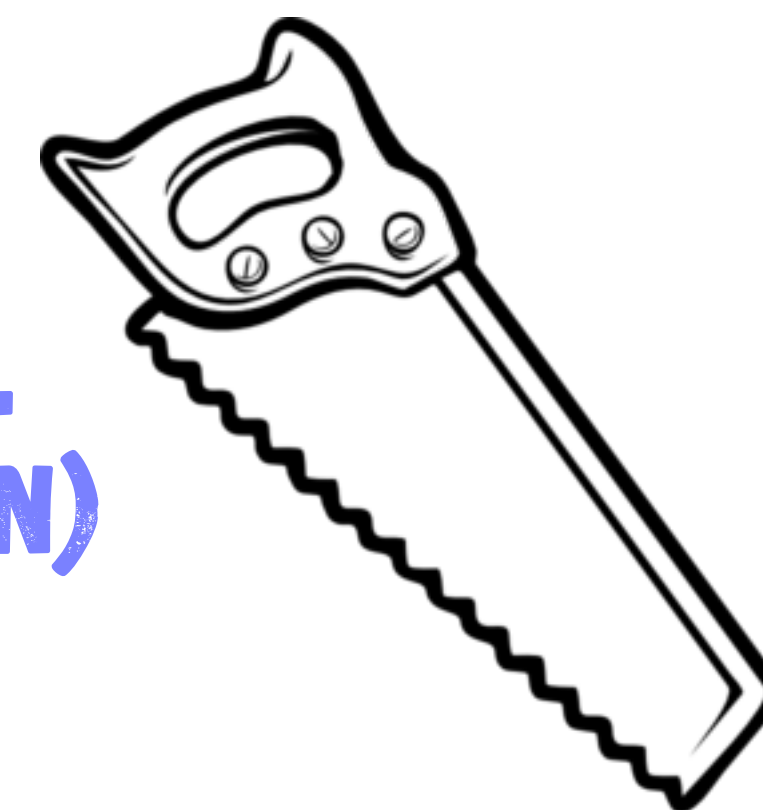
**AUTOENCODERS**



**LINEAR  
REGRESSION**



**GENERATIVE  
ADVERSARIAL  
NETWORKS (GAN)**



***Best tool depends on the problem! (Intro to these algorithms in backup)***



# WHAT CAN ML NOT DO?

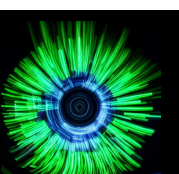


*Garbage In*

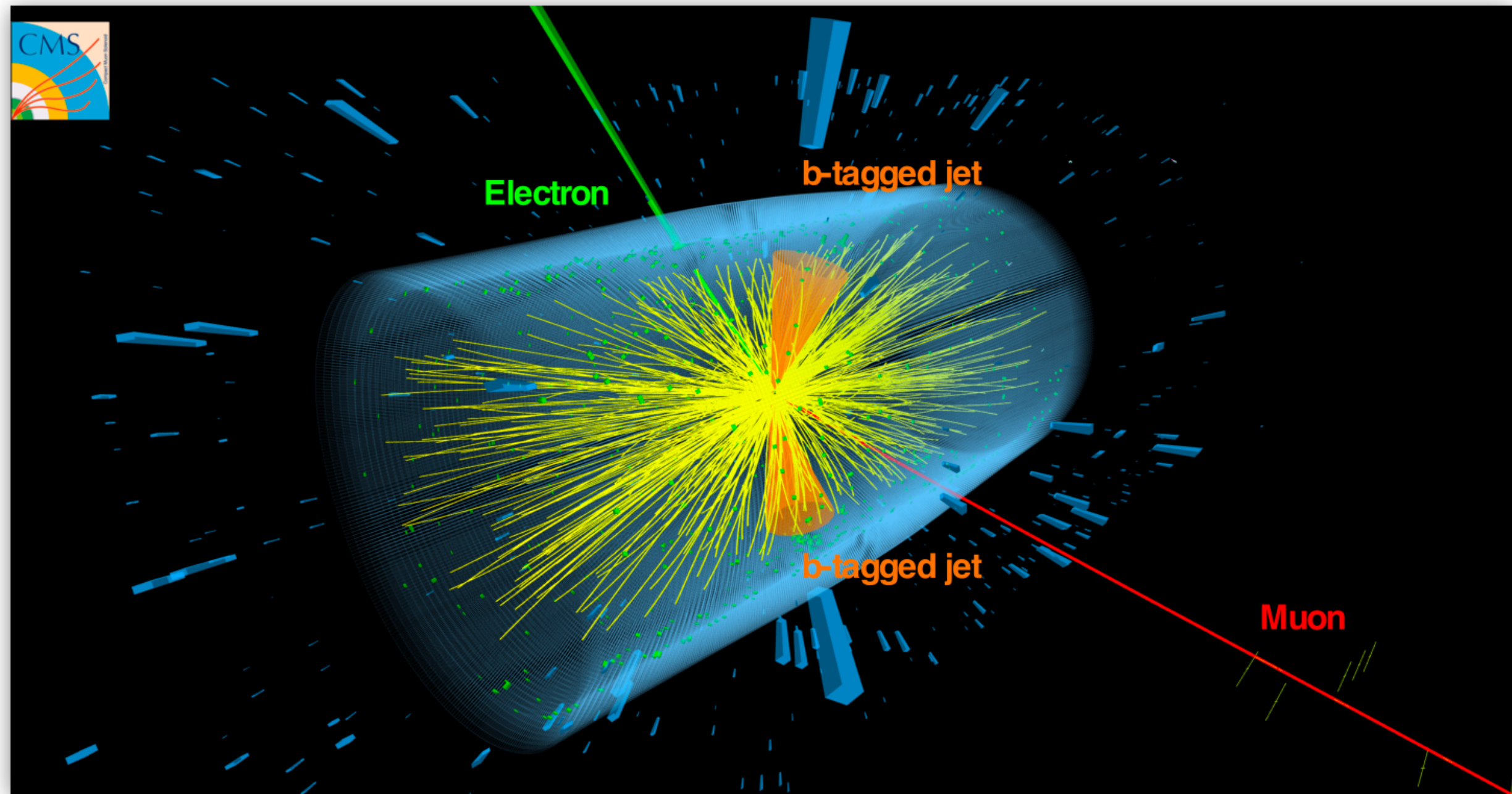


*Garbage Out*

- ❌ ML cannot replace domain knowledge.
- ❌ ML is not a causation tool.
- ❌ **ML is not a magic fix!**



# ML FOR HARD AND ELECTROMAGNETIC PROBES

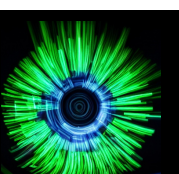


- **Hard Probes:** Products of early-stage hard scatterings that interact with the QGP medium.
- **Electromagnetic Probes:** probes that have a long mean free path relative to the size of the QGP (negligible interactions)

- **Hard and electromagnetic probes offer a clean and well-calibrated environment!**

***HI environment can be challenging for ML.***

- Higher particle multiplicities, much more complex system (even by eye)!
- Dependence on simulation used in training makes application to data difficult.

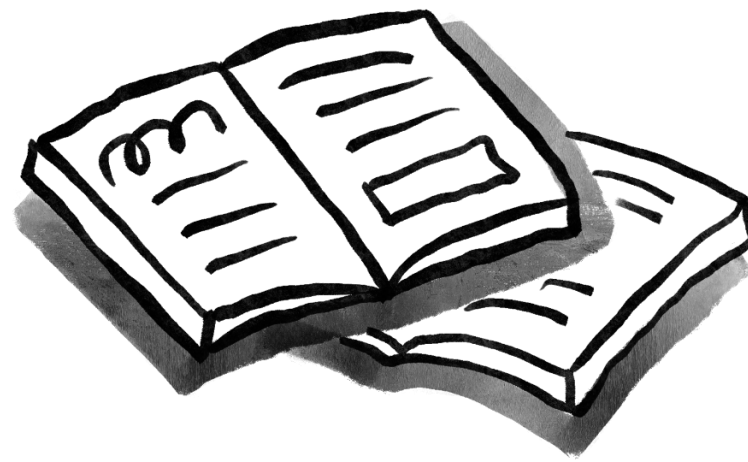


# ROADMAP

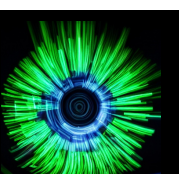
**WHAT IS AI/ML  
AND WHY IS IT  
USEFUL FOR THE  
ANALYSIS OF  
HARD PROBES?**



**HOW IS AI/ML  
CURRENTLY  
BEING USED FOR  
ANALYSIS?**



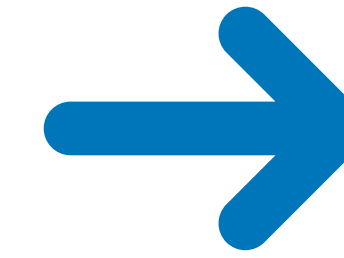
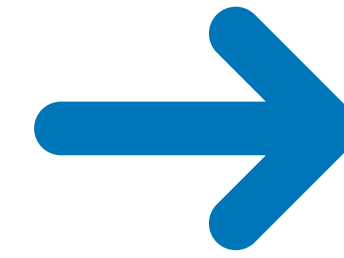
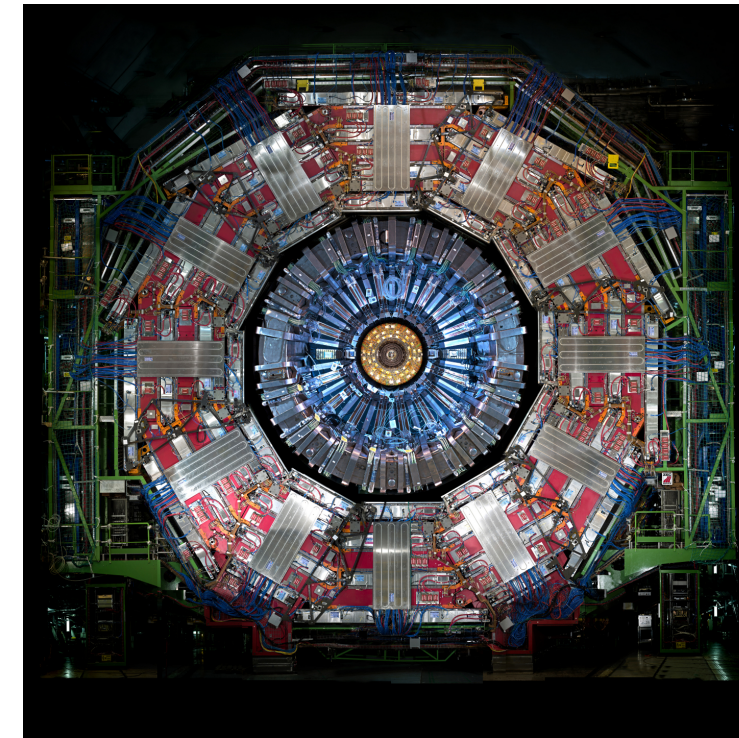
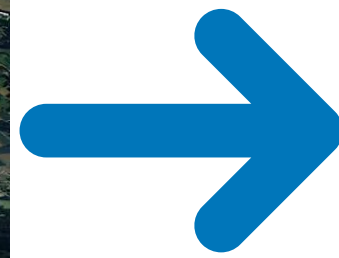
**WHERE ARE WE  
HEADING?**





# ANALYSIS PIPELINE

MACHINE LEARNING CAN BE USED THROUGHOUT THE ANALYSIS PIPELINE!

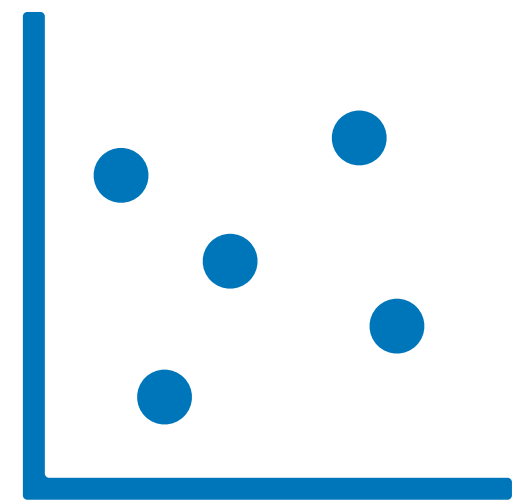


**ACCELERATOR COMPLEX**

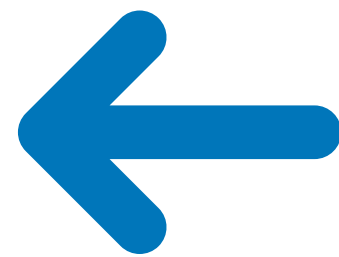
**EXPERIMENTS**

**EVENT FILTERING**

**RECONSTRUCTION**



**RESULTS**



**EVENT/ANALYSIS SELECTIONS**



**RECONSTRUCTED  
DATA**

**SIMULATION**

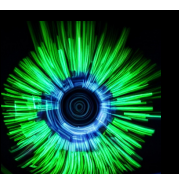
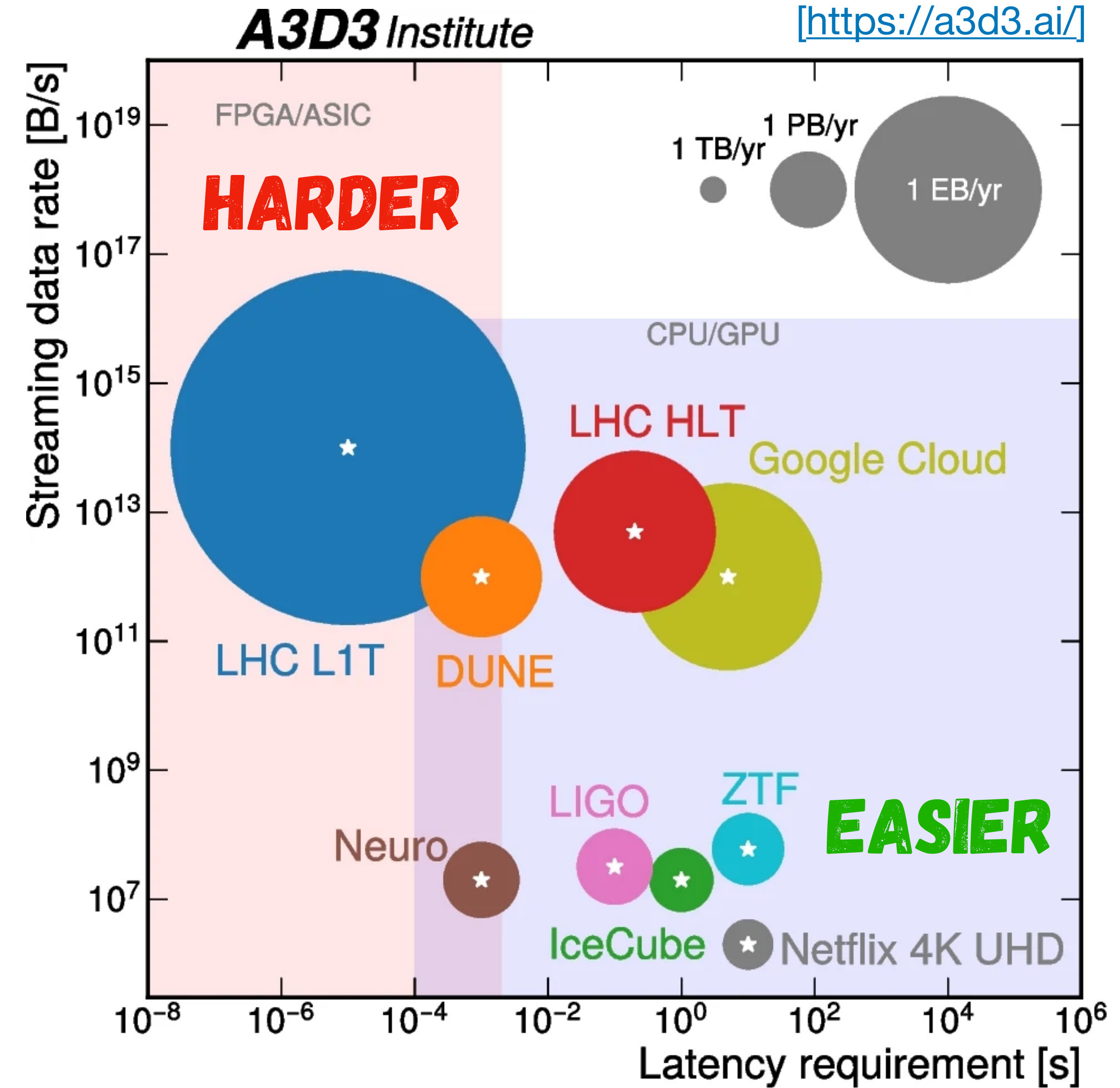


See [[Phys. Rev. C 110, 034912](#)] for full event simulation in sPHENIX using diffusion models, Published last night!



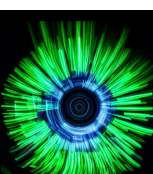
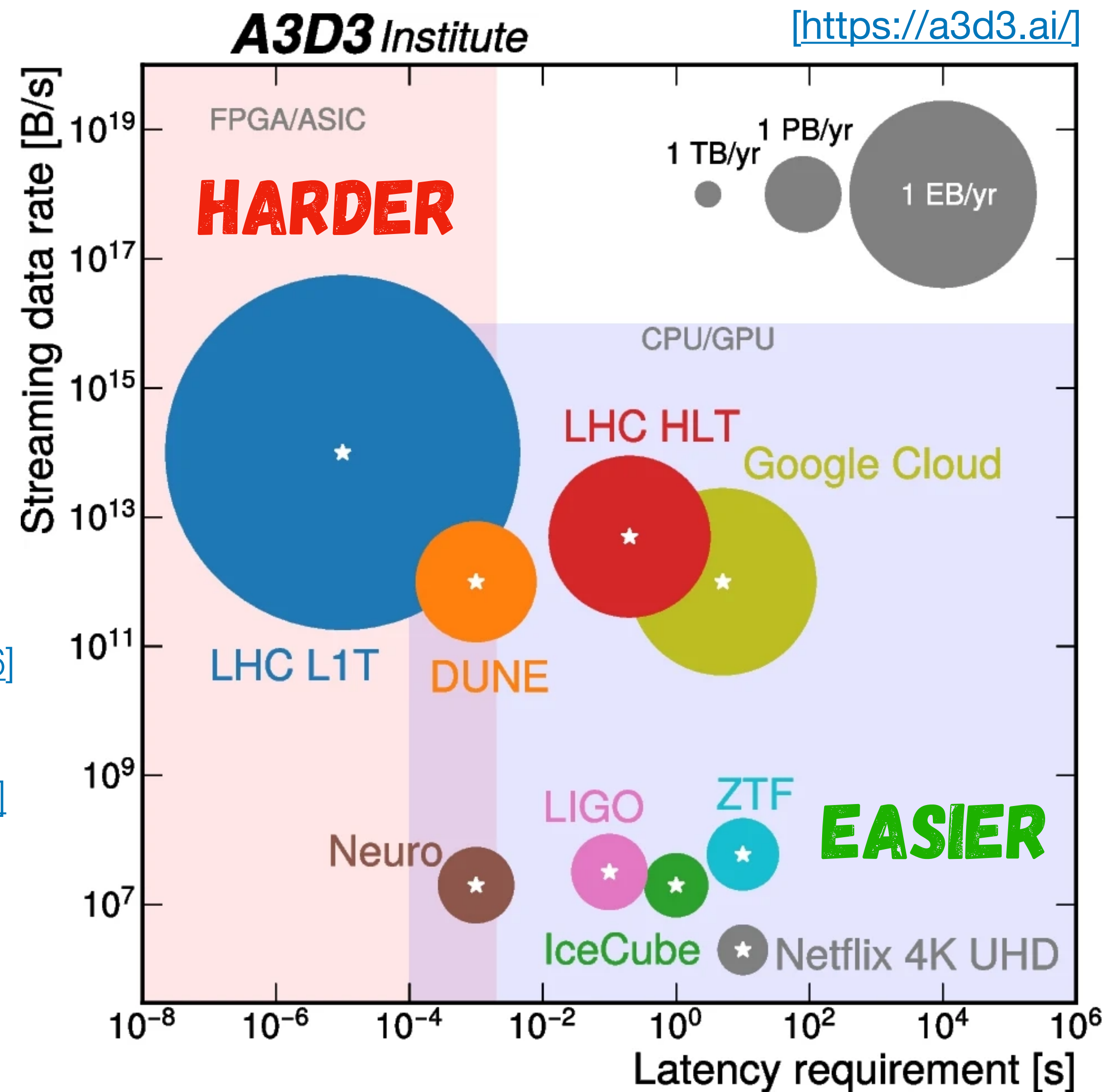
# EVENT FILTERING

- Data volume is increasing at a fast rate, need solutions for limited computing resources.
- If we took all raw data, would easily exceed storage capabilities.



# EVENT FILTERING

- Data volume is increasing at a fast rate, need solutions for limited computing resources.
  - If we took all raw data, would easily exceed storage capabilities.
- Perform fast selection/rejection of data with ML integrated into the firmware (FPGAs)
  - Use high level synthesis packages ex: [hls4ml](#)  
*CMS L1 Trigger* [\[CMS-TDR-021\]](#) *sPHENIX HF Trigger* [\[JINST 19 C02066\]](#)  
*ATLAS Fake Track Rejection in Event Filter* [\[ATLAS-TDR-029-ADD-1\]](#)  
*LHCb track reconstruction for HLT system* [\[See website here\]](#)



# DATA REDUCTION

- One other solution is to reduce the data size, for example, with auto-encoders.

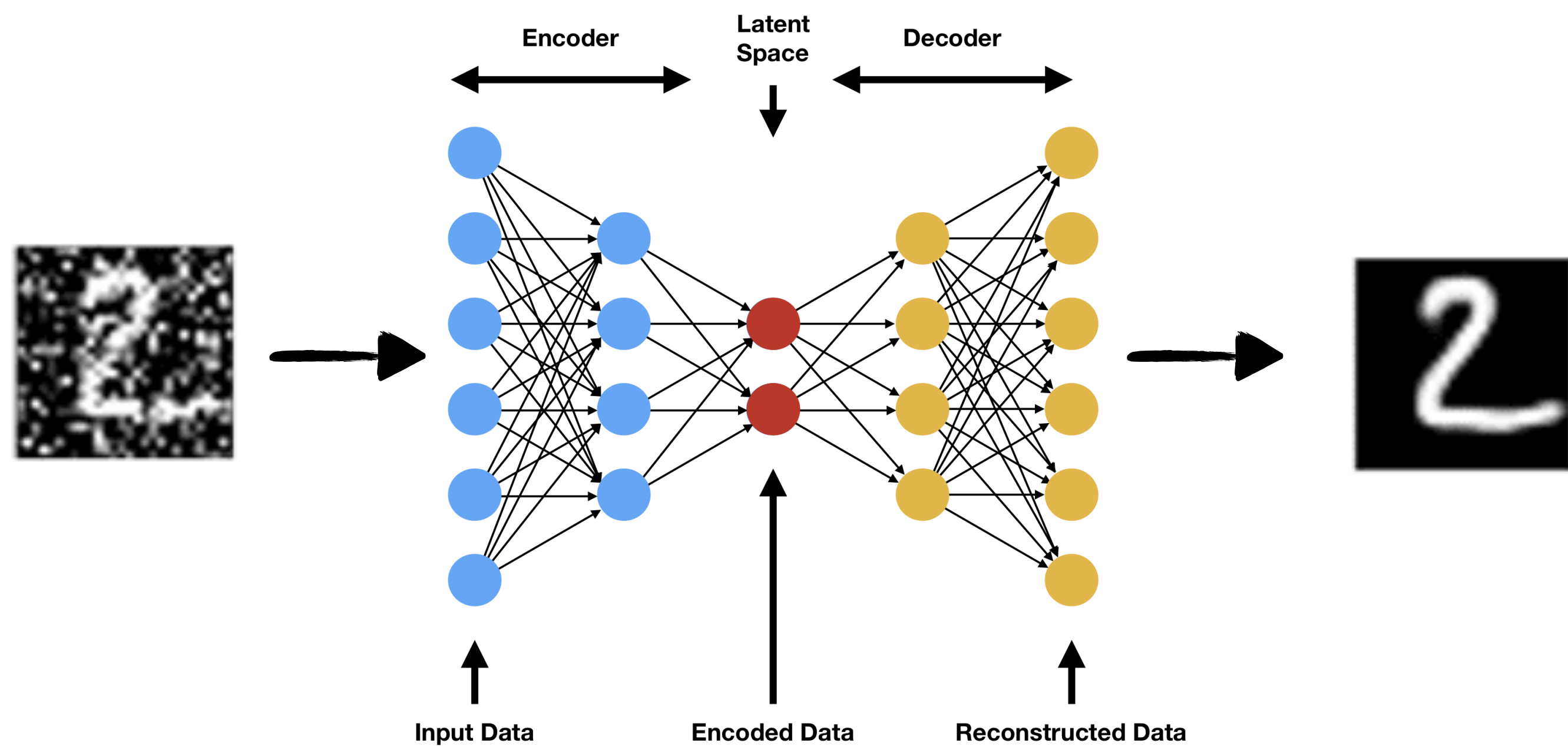
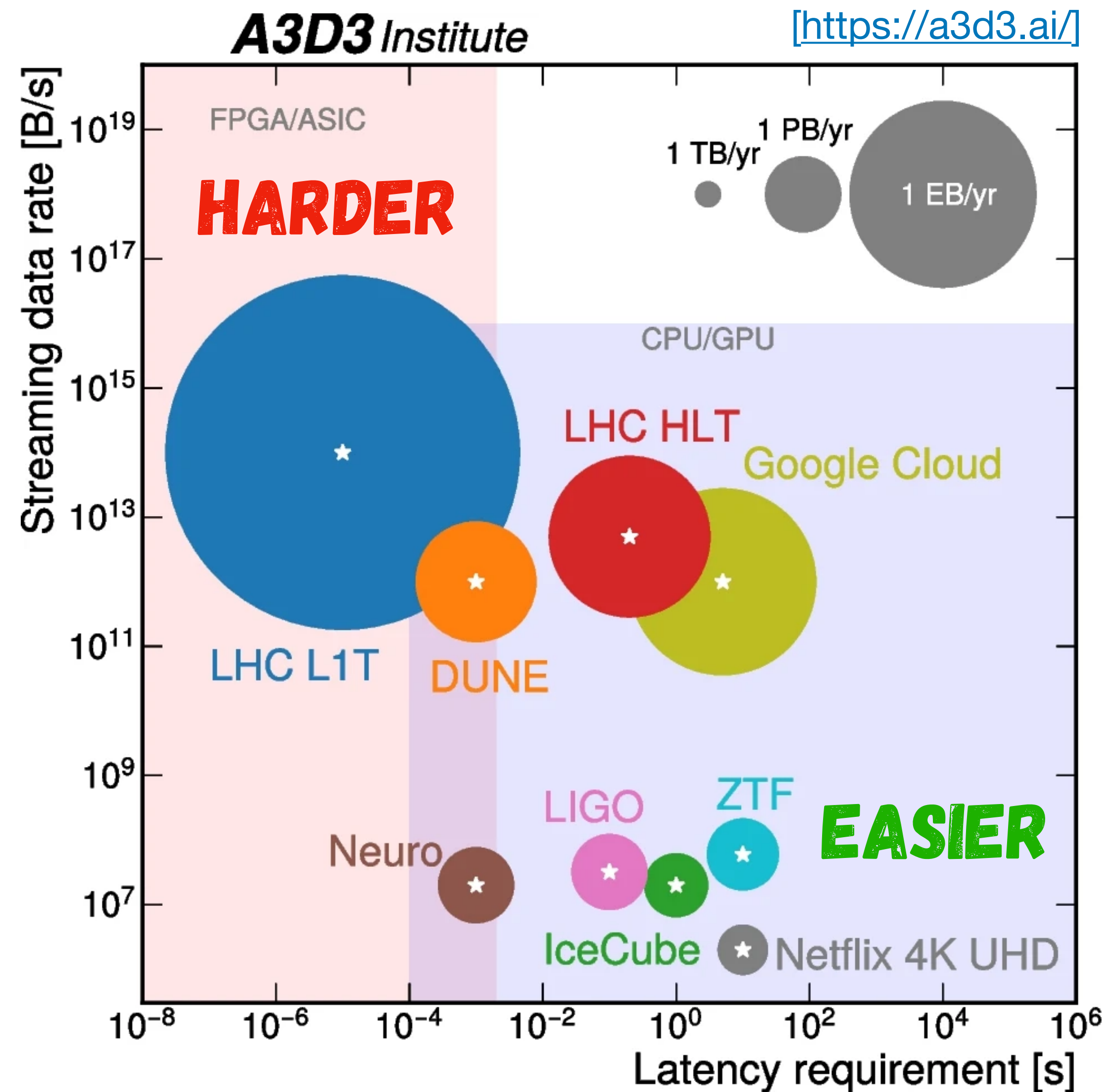


Image Credit: <https://www.compthree.com/blog/autoencoder/>

sPHENIX TPC [\[arXiv:2310.15026\]](https://arxiv.org/abs/2310.15026)

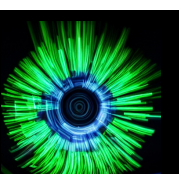
CMS HGCal upgrade [\[See talk at Fast ML 2022\]](#)

LHCb trigger system [\[See website here\]](#)



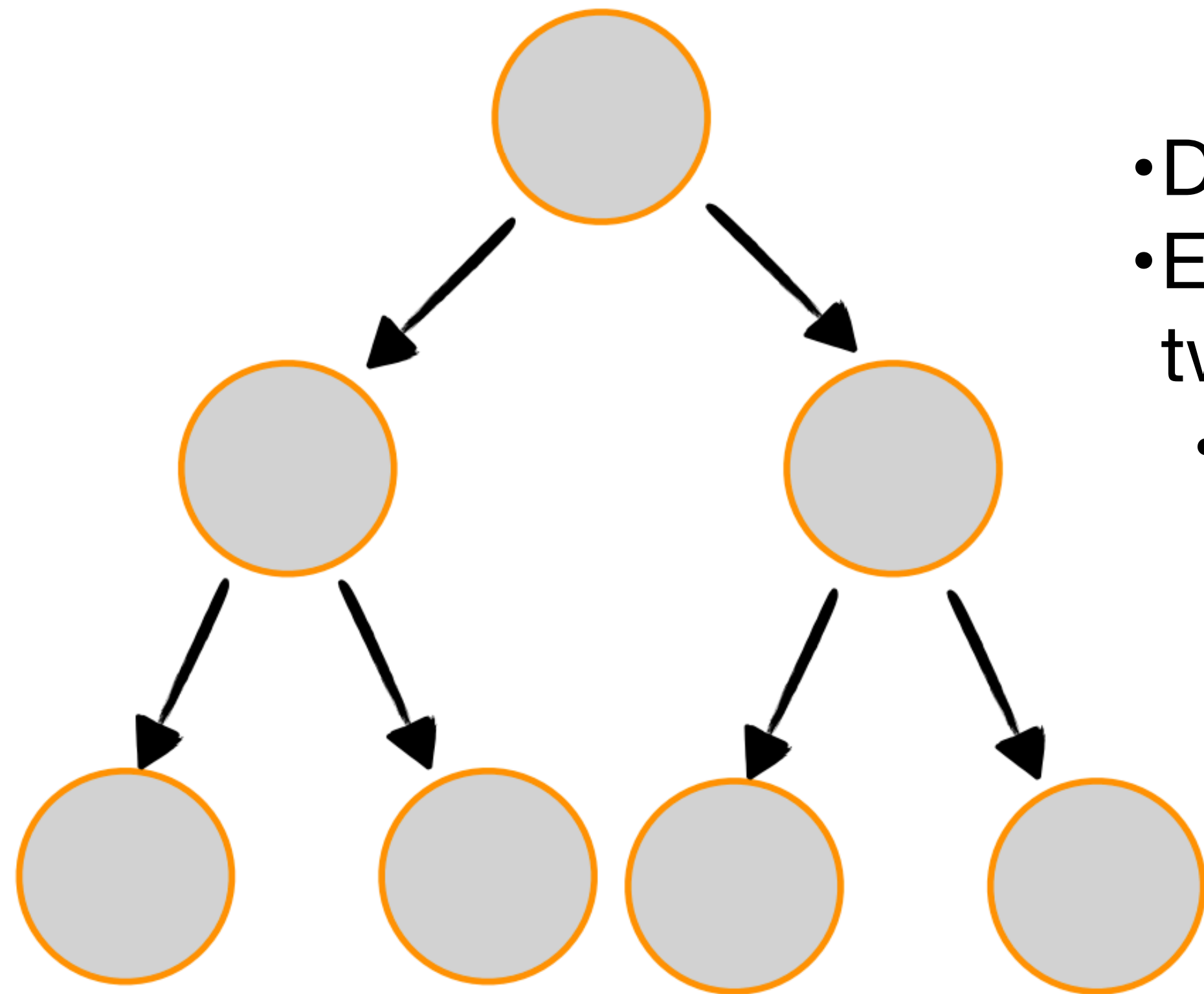
# SIGNAL/BACKGROUND DISCRIMINATION

- **Conventional approach:** Apply cuts to identify signal based on expert knowledge
  - Becomes difficult w/ complex signals or in HI environment with a large background.
- **Solution:** Employ multiple variables simultaneously - essence of ML!



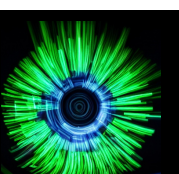
# SIGNAL/BACKGROUND DISCRIMINATION

- **Conventional approach:** Apply cuts to identify signal based on expert knowledge
  - Becomes difficult w/ complex signals or in HI environment with a large background.
- **Solution:** Employ multiple variables simultaneously - essence of ML!



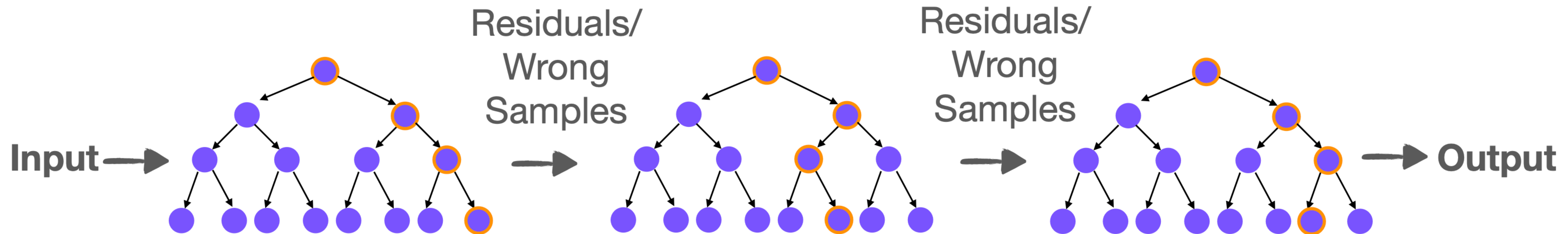
- Decision trees are commonly used for signal classification
- Each **node** is a classification rule that splits the data into two or more parts.
  - In training you determine the proper rules that maximizes the information gain and minimize entropy

$$E(x) = \sum -p(x)\log_2(p(x))$$



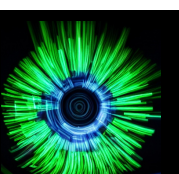
# SIGNAL/BACKGROUND DISCRIMINATION

- **Conventional approach:** Apply cuts to tag particle based on decay topology
  - Becomes difficult in heavy-ion environment with a large background.
- **Solution:** Employ multiple variables simultaneously - essence of ML!

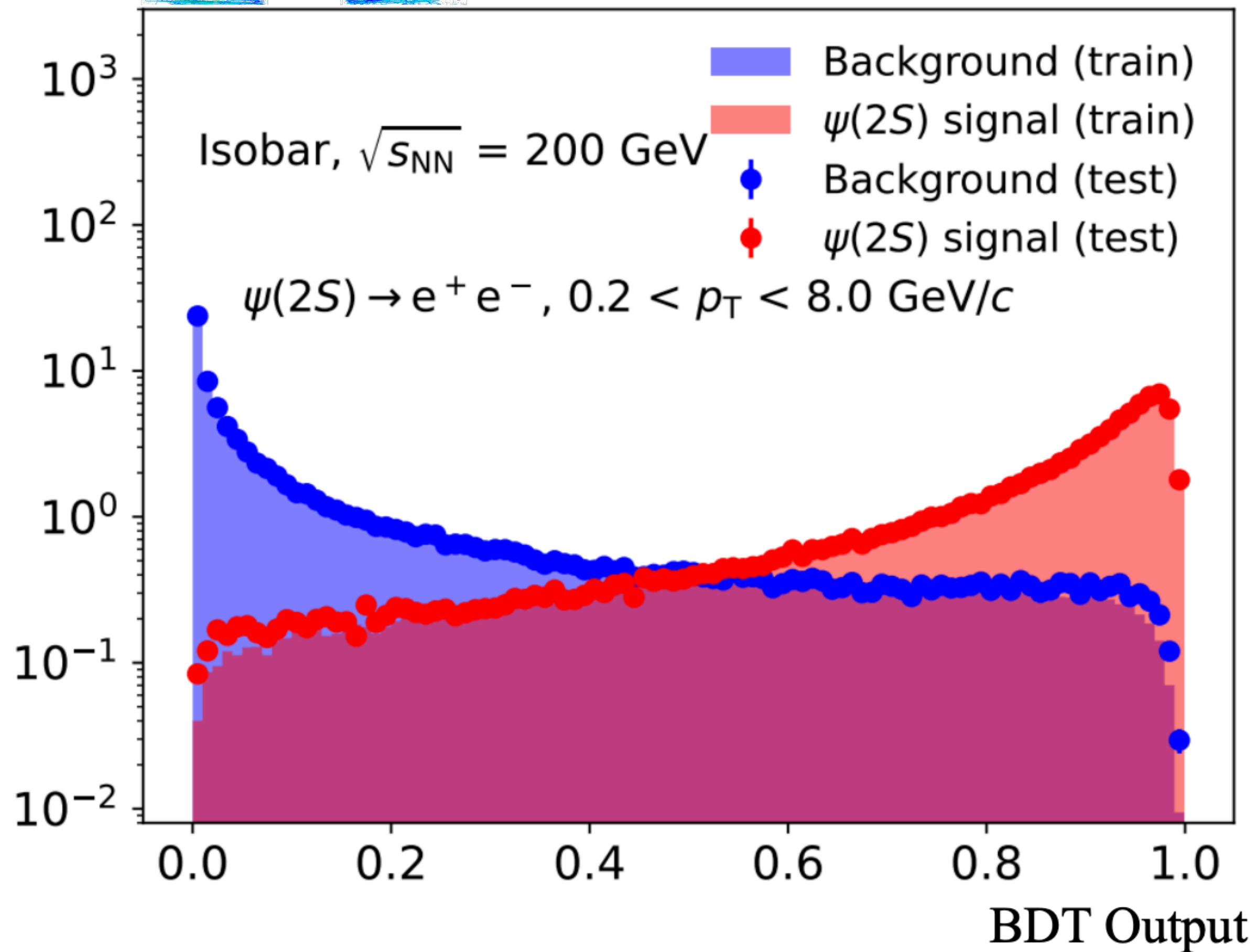
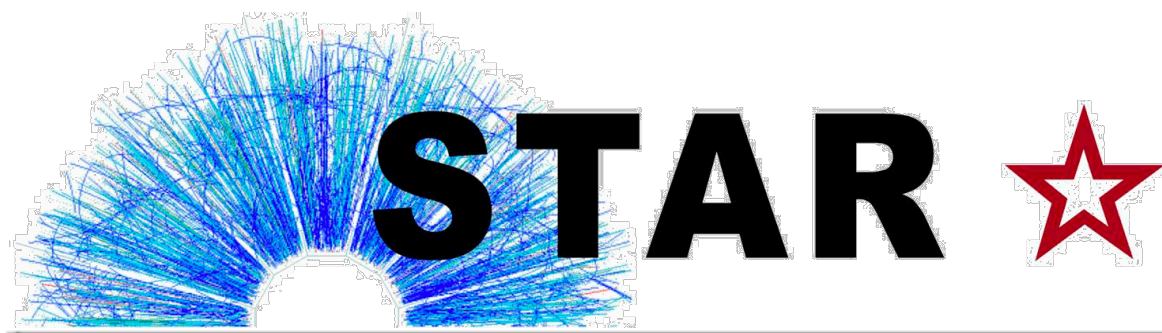


- Boosted decision trees are used when multiple weaker learners are combined in a series where each additional component seeks to minimize error of previous one.

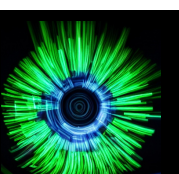
Yann Coadou, [arXiv: 2206.09645](https://arxiv.org/abs/2206.09645)



# SIGNAL/BACKGROUND DISCRIMINATION

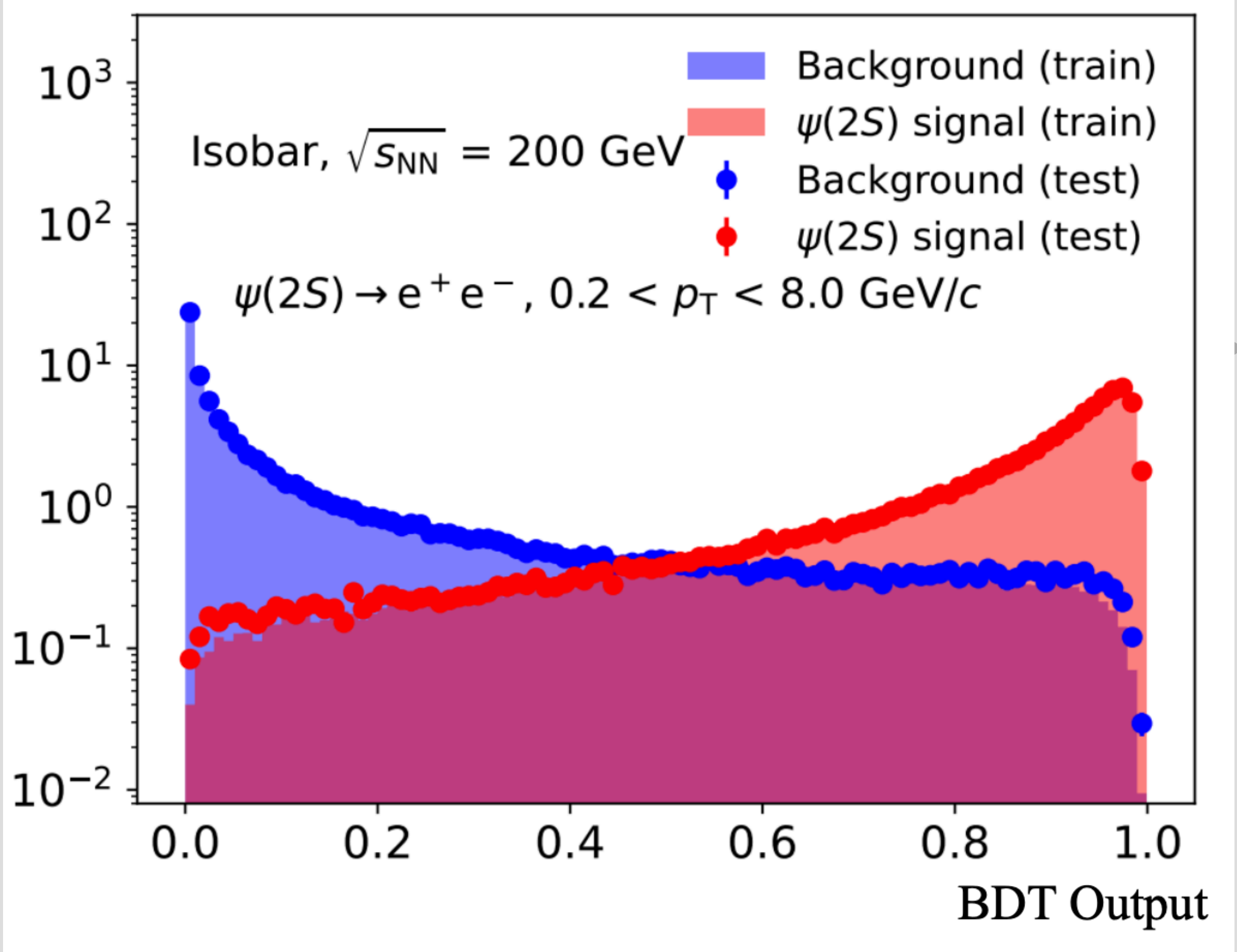


- Ex: Use BDT in order to reconstruct the  $\psi(2s)$  signal.
- XGBoost is the core of the application

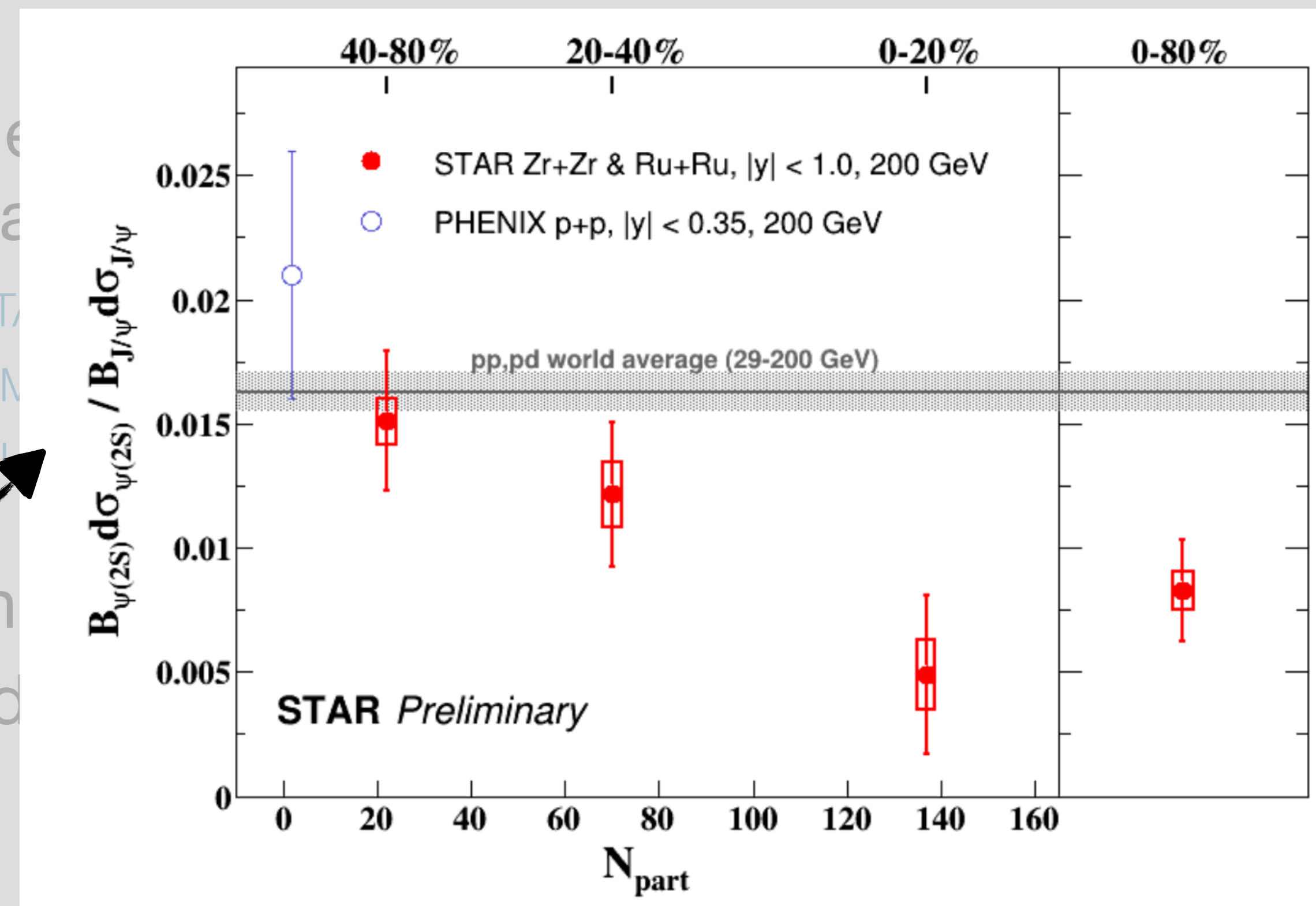




# SIGNAL/BACKGROUND DISCRIMINATION



• The  
exa  
ST  
CM  
AL  
On  
oro



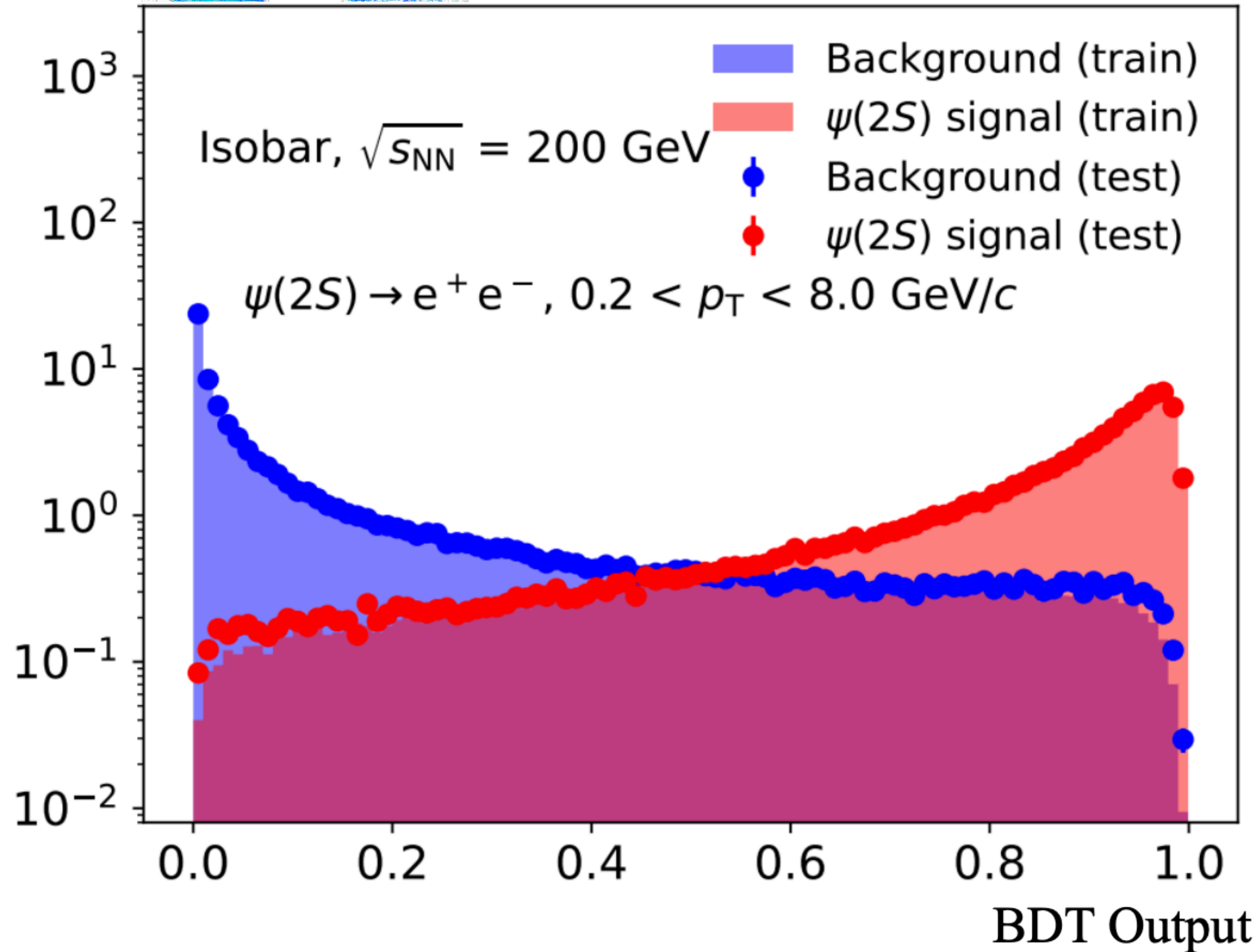
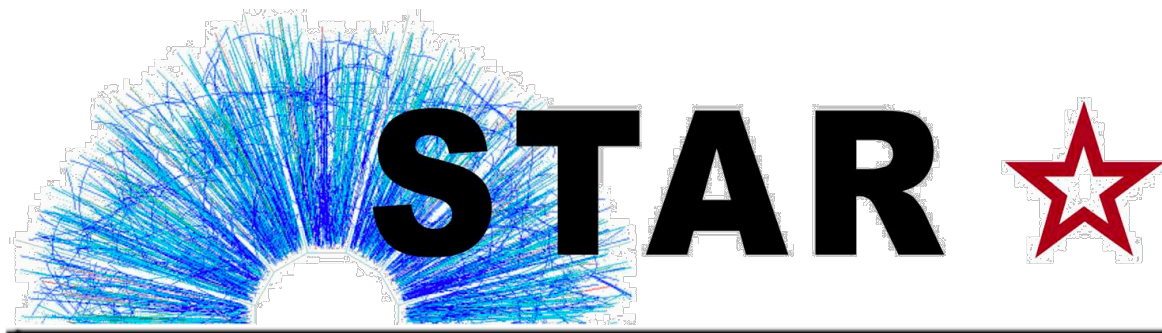
**ENABLES FIRST OBSERVATION OF  
CHARMONIUM SEQUENTIAL  
SUPPRESSION AT RHIC!**

**HARD PROBES 2024**

**W. ZHANG**

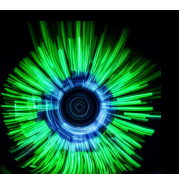
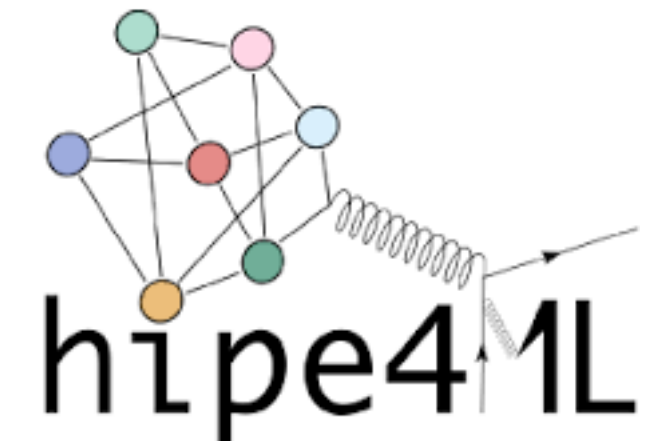
**TUES. 09:40**

# SIGNAL/BACKGROUND DISCRIMINATION



- Ex: Use BDT in order to reconstruct the  $\psi(2s)$  signal.
- XGBoost is the core of the application
- Many other examples! (Not an exhaustive list.)
  - [[Phys. Lett. B 839 \(2023\) 137796](#)]
  - [[JHEP 05 \(2021\) 220](#)]
  - [[Phys. Lett. B 782 \(2018\) 474](#)]
  - [[PRL 124, 172301 \(2020\)](#)] .....

See also heavy ion physics environment for machine learning ([hipe4ml](#))



# HEAVY FLAVOR JET TAGGING

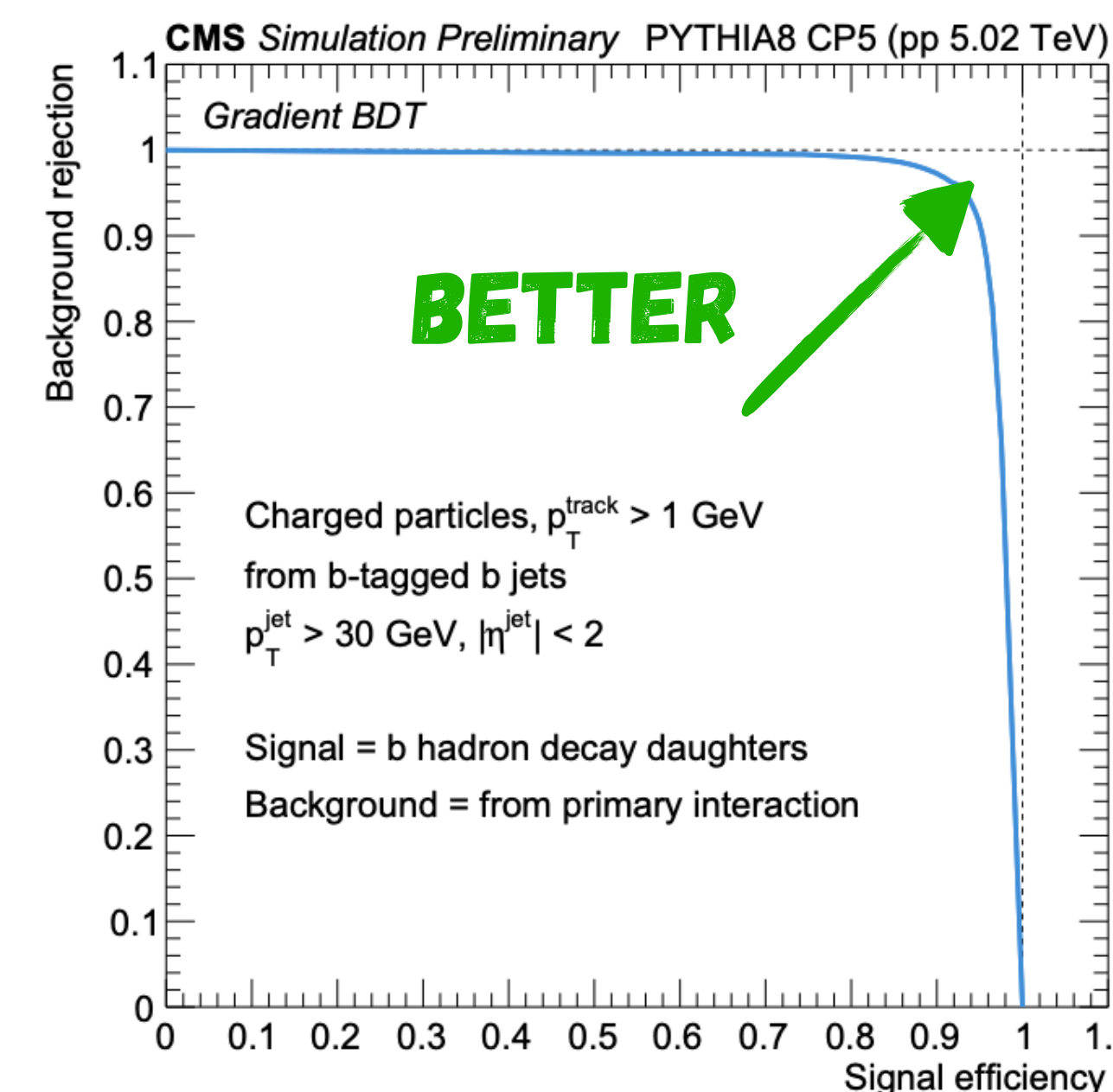
**Goal:** identify jets initiated by a heavy-quark

**Conventional approach:** Apply cuts to select jets with displaced decay vertices and large impact parameter tracks.

**ML approach:** Learn from low-level features in a supervised approach using BDT or a GNN



[CMS-PAS-HIN-24-005](#)



[JINST 13 (2018) 05, P05011]



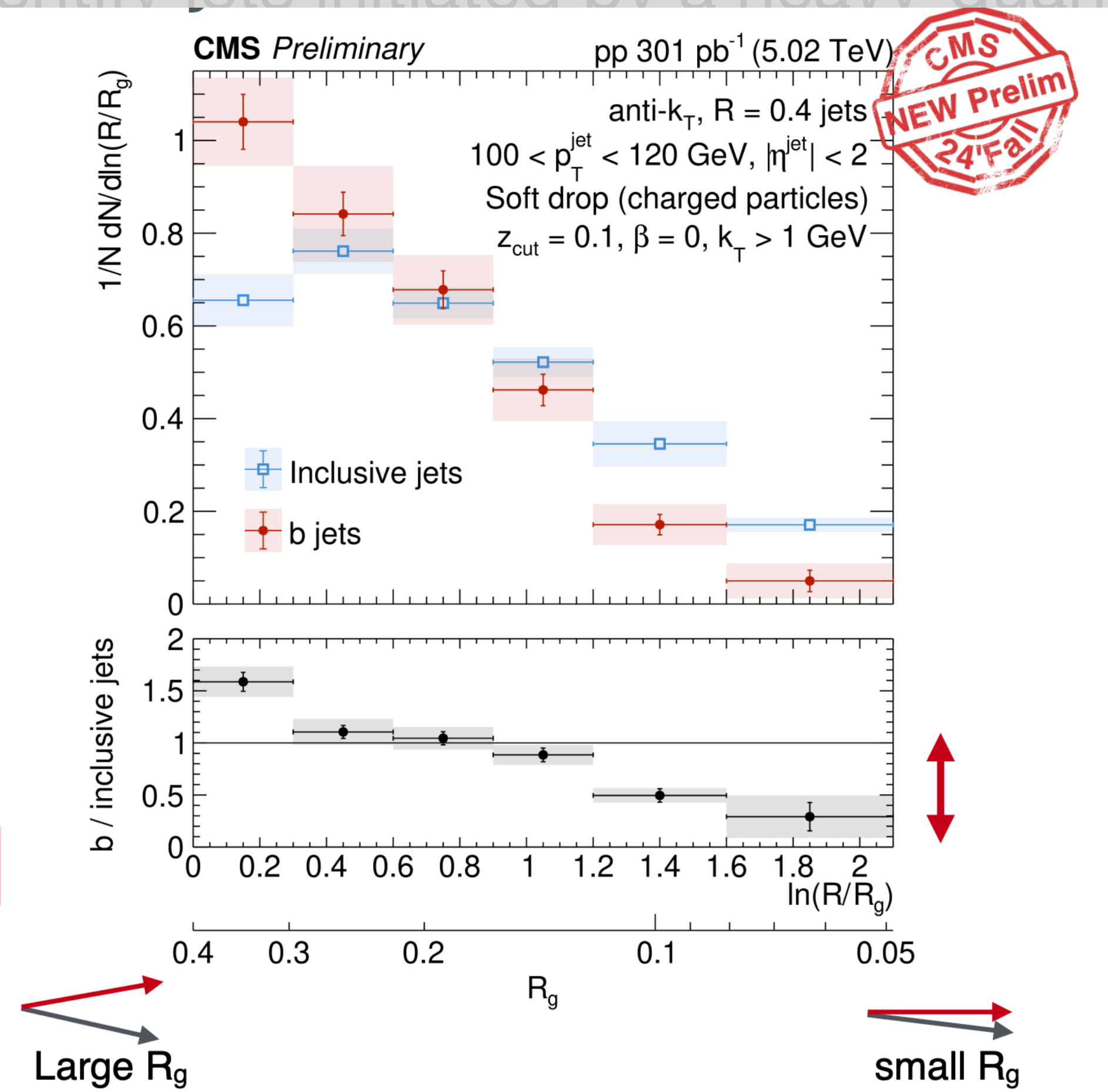
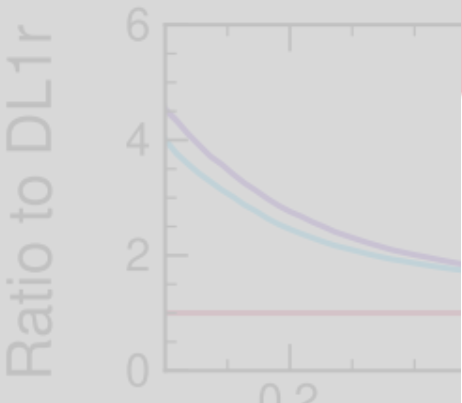
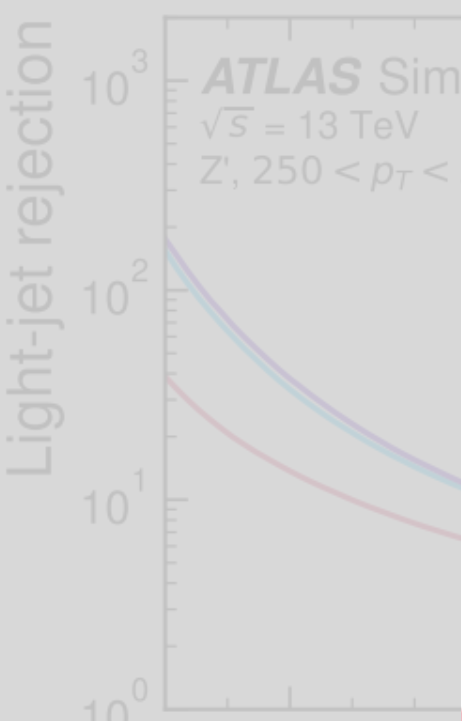
# HEAVY FLAVOR JET TAGGING

Goal: identify jets initiated by a heavy-quark (HF jet)

Conventional  
decay veto

ML approach

[ATLAS]

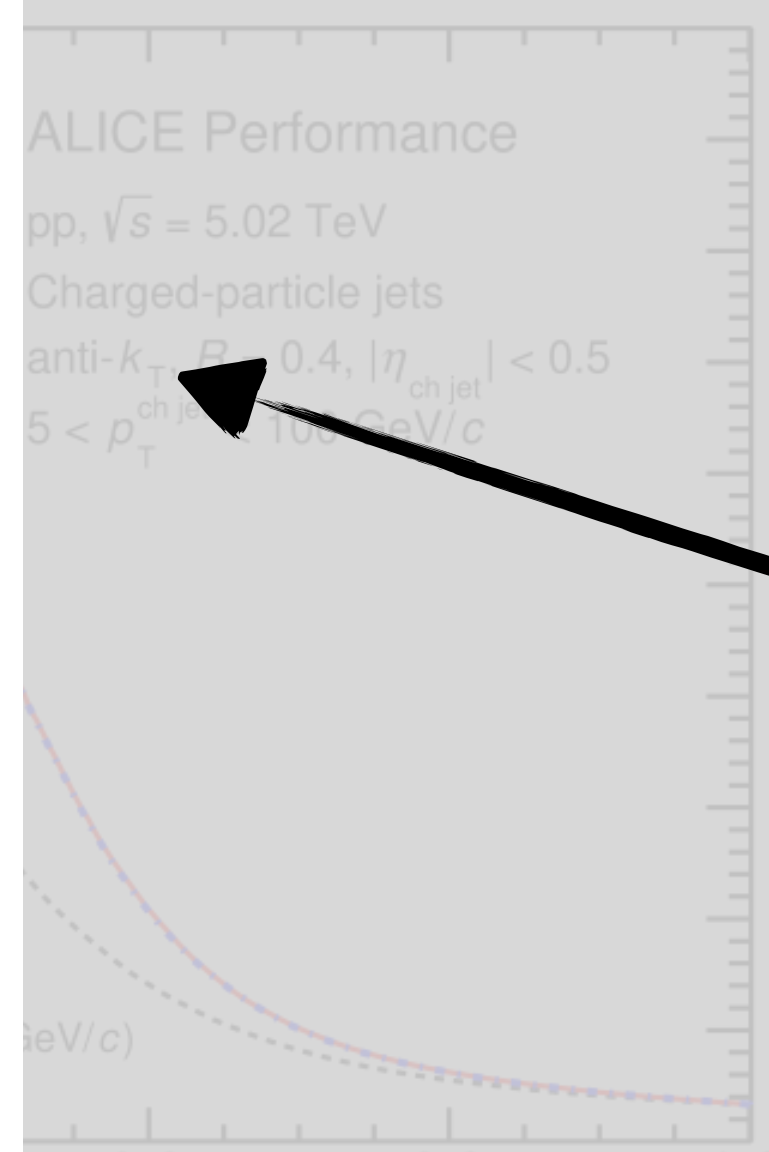


c-jet tagging efficiency

Identify jets with displaced tracks.

achieve using BDT or a GNN

ALICE Performance



b-jet tagging efficiency

**HARD PROBES 2024**

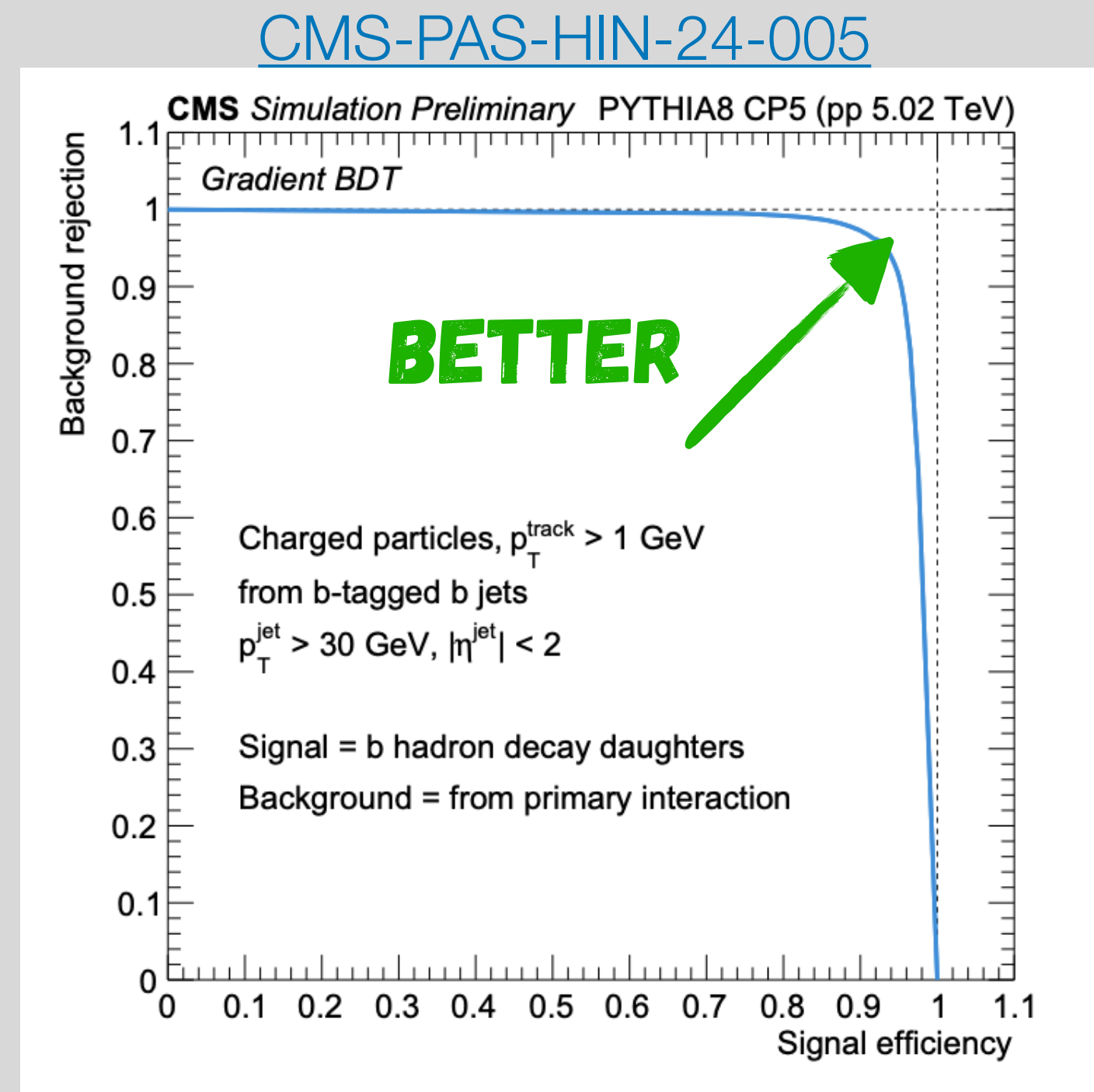
**L. KALIPOLITI**

**WED. 9:40**

**HARD PROBES 2024**

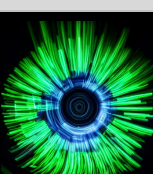
**C. CHOI**

**POSTER**



**ALLOWS FOR THE FIRST OBSERVATION OF B-QUARK DEAD CONE!**

[JINST 13 (2018) 05, P05011]



# HEAVY FLAVOR JET TAGGING



**Goal:** identify jets initiated by a heavy-quark (HF jet)

**Conventional approach:** Apply cuts to select jets with displaced decay vertices and large impact parameter tracks.

**ML approach:** Learn in a supervised approach using BDT or a GNN

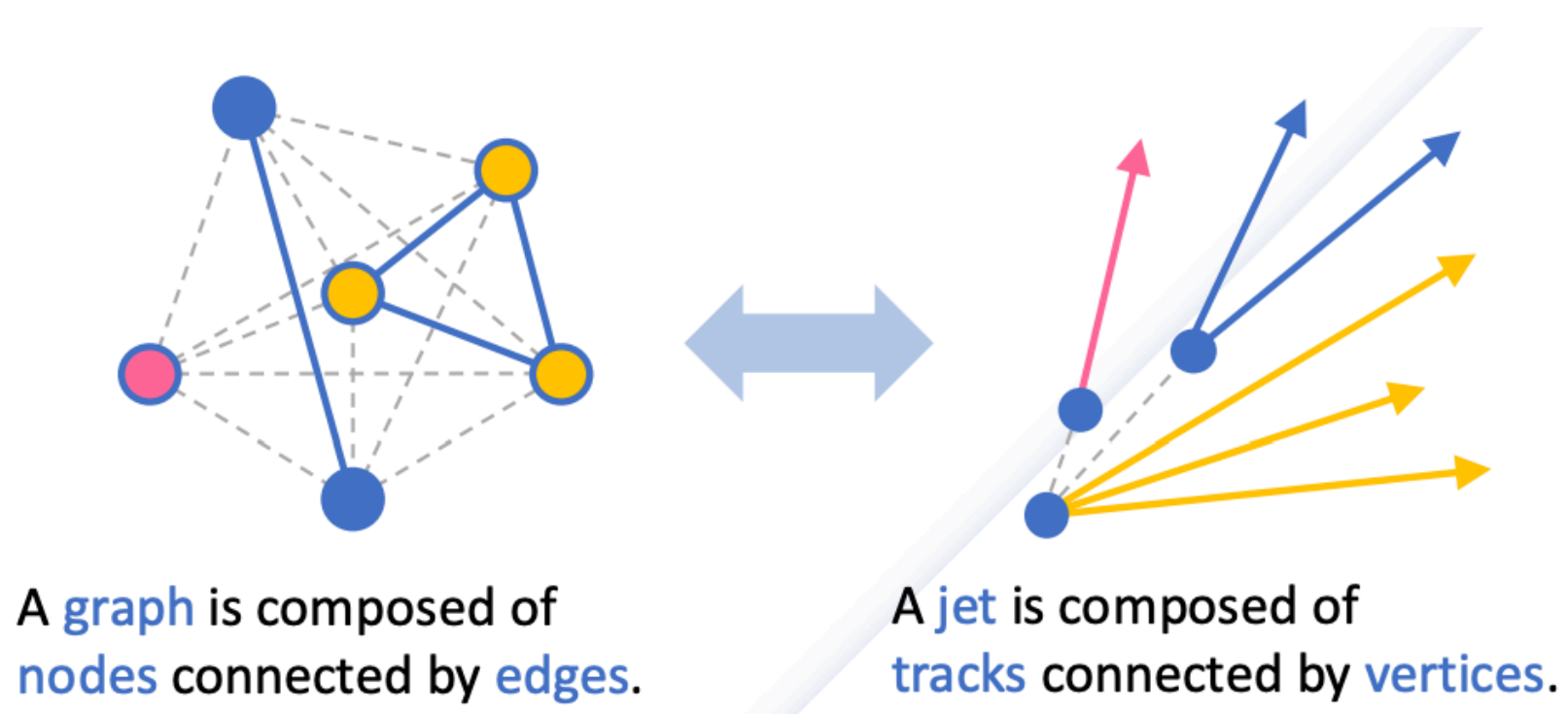
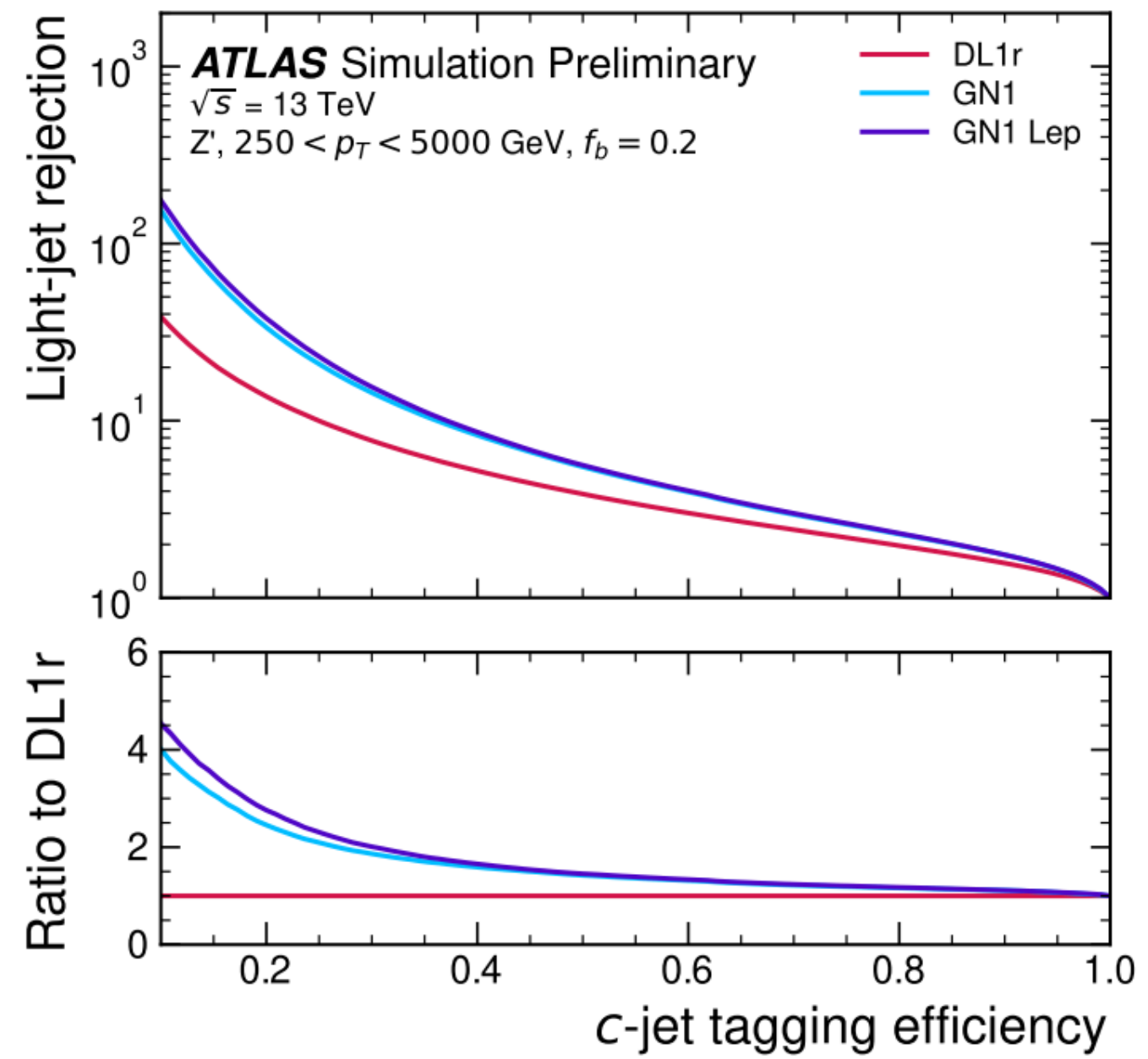
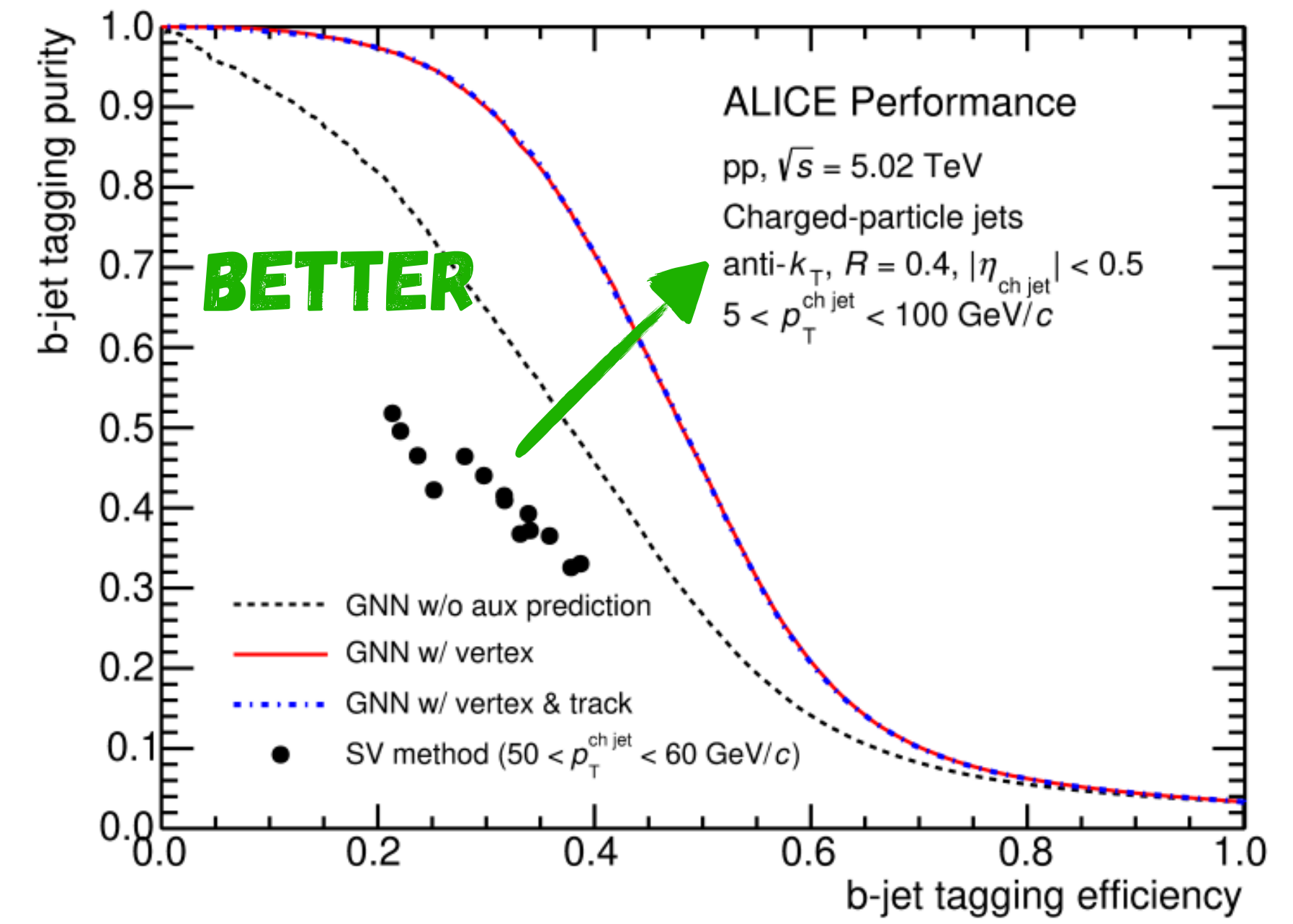


Image from Changwhan Choi, Poster

[ATL-PHYS-PUB-2022-027]

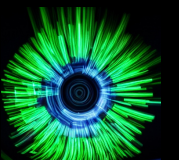


ALICE Collaboration, Preliminary



ALI-PERF-579779

[JINST 13 (2018) 05, P05011]

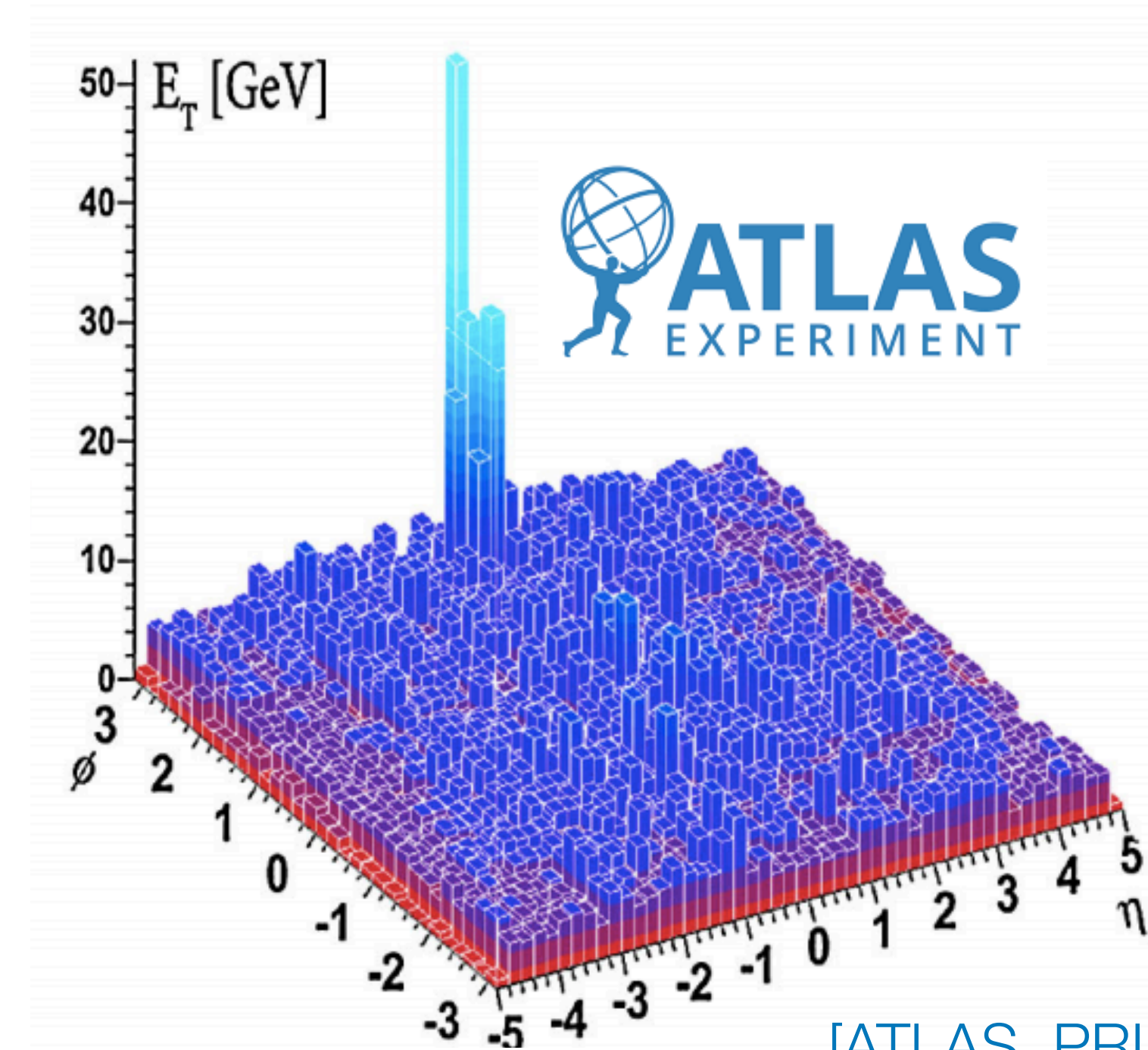


# JET BACKGROUND CORRECTION

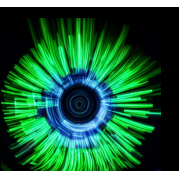
- **Differential measurements of jets are key to understanding jet quenching effects!**
- These often involve pushing to large  $R$  and/or low  $p_T$ , where background contribution is difficult to subtract.

By now many methods in which ML can be used to solve this problem! We will discuss two.

See also [[Phys. Rev C. 108.L021901 \(2023\) 6](#)]



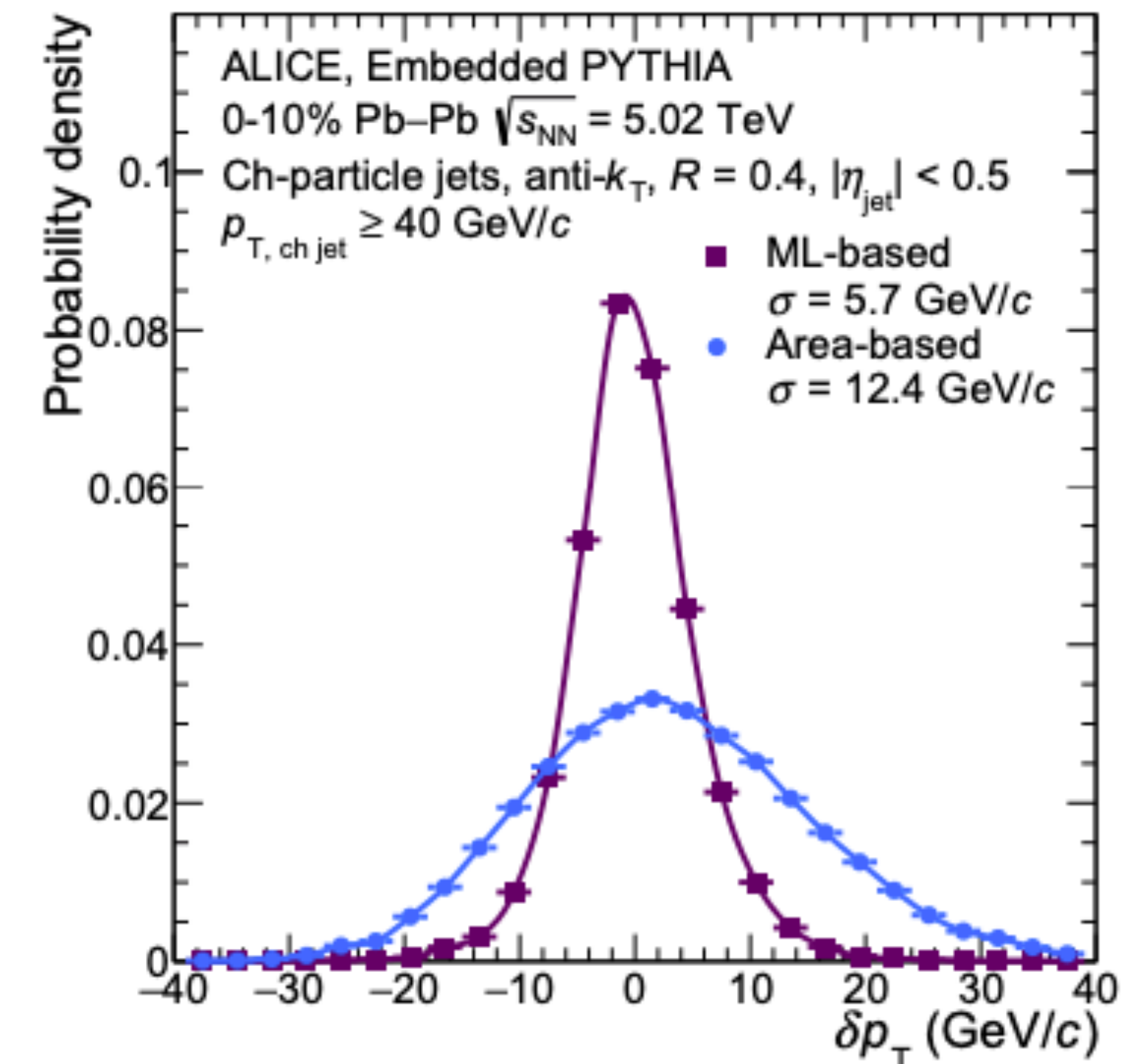
[[ATLAS, PRL 105, 252303 \(2010\)](#)]



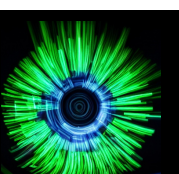
# JET BACKGROUND CORRECTION

- **Method 1:** Shallow NN in [scikit-learn](#) (simple tools) trained on PYTHIA embedded into HI background [\[PRC 99, 064904 \(2019\)\]](#)

ALICE, [\[PLB 849 \(2024\) 138412\]](#)

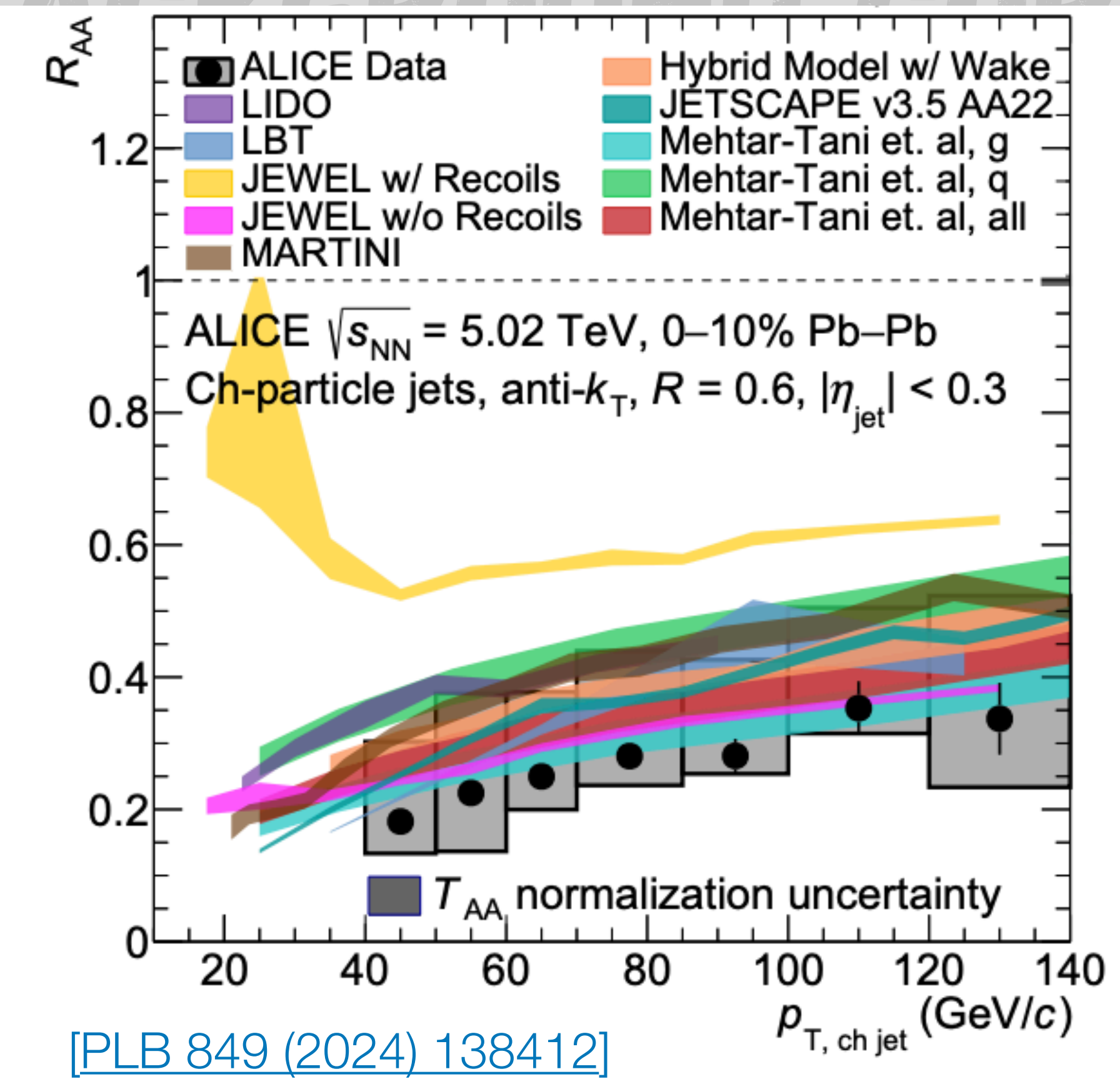


$$\delta p_T = p_{T,rec} - p_{T,true}$$

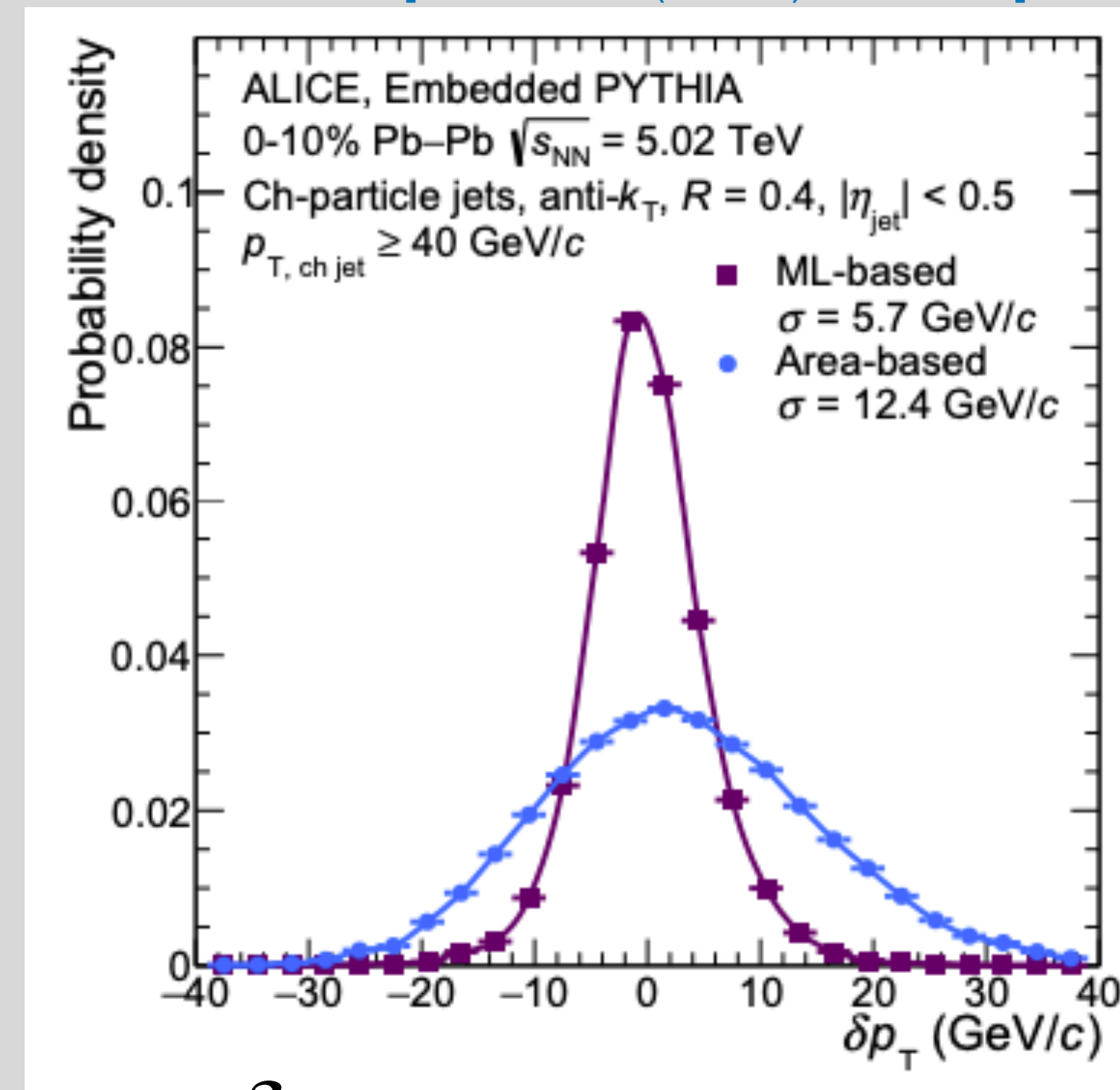


# JET BACKGROUND CORRECTION

- Method 1: PYTHIA embedded
- Fragmentation with realistic
- This is a



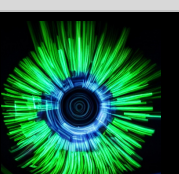
ALICE, [PLB 849 (2024) 138412]



$$\delta p_T = p_{T, \text{rec}} - p_{T, \text{true}}$$

## MAKES LARGE RADIUS JET MEASUREMENTS ( $R = 0.6$ ) POSSIBLE W/ ALICE

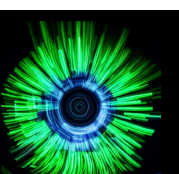
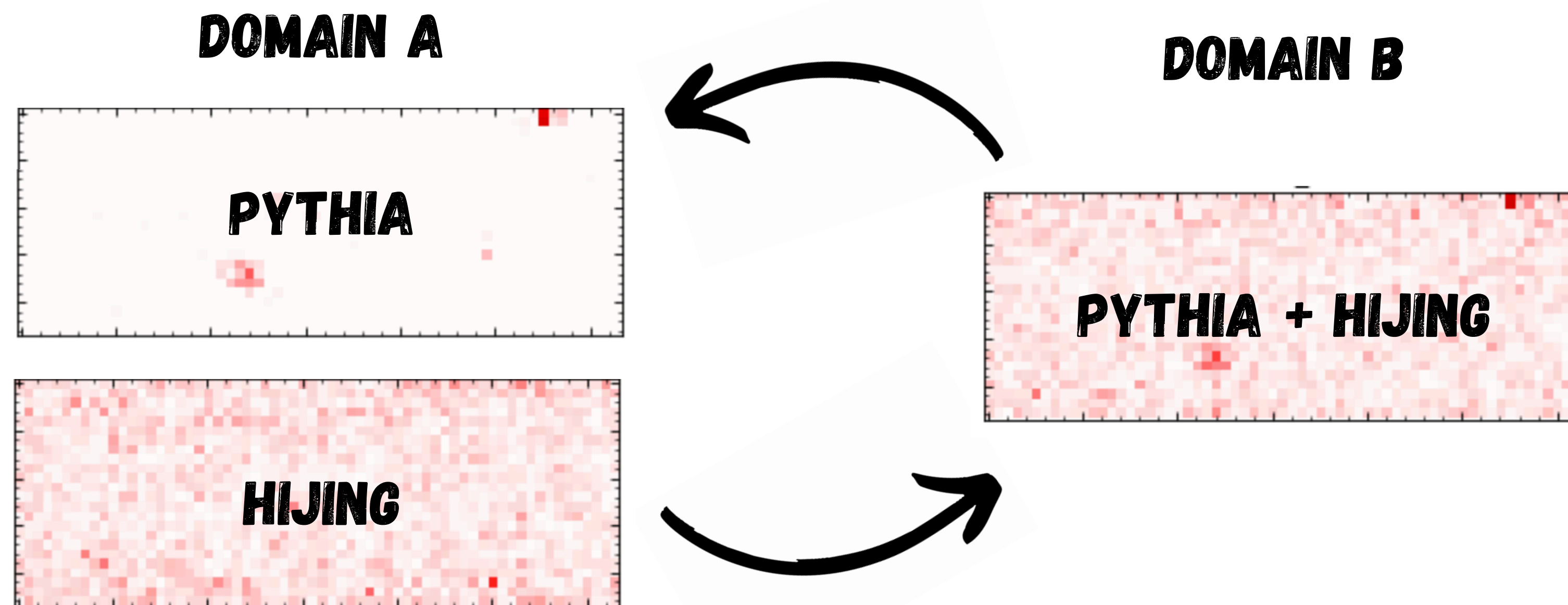
- Simulation bias quantified with a PYTHIA-based toy model with realistic fragmentation variations
- **This is an important source of uncertainty!**





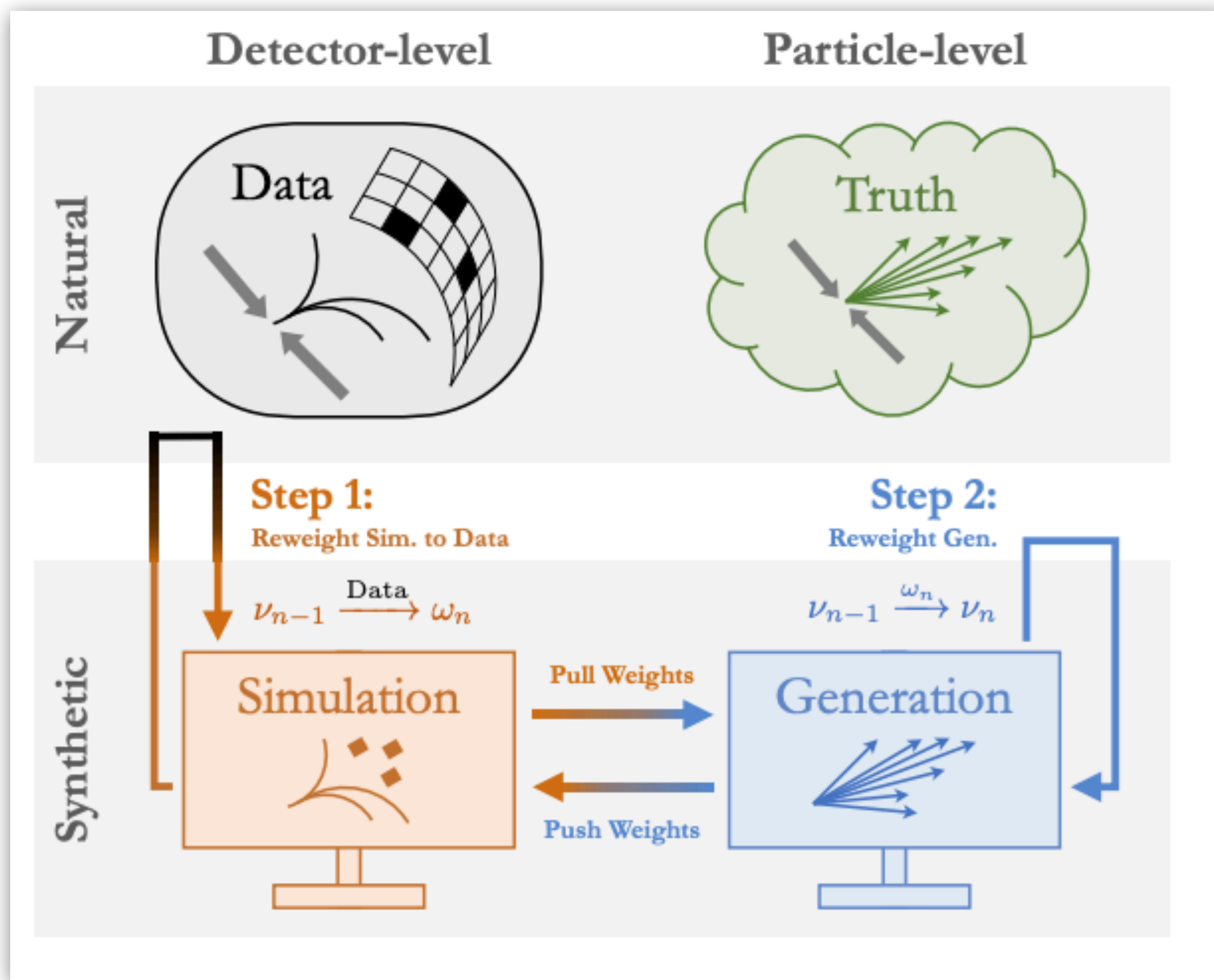
# JET BACKGROUND CORRECTION

- **Method 2:** Use generative AI (unpaired image-to-image translation, cycleGANs) to subtract jet background in an *unsupervised* way.
  - Composed of two generator-discriminator pairs w/ cyclic closure (i.e.  $A \rightarrow B \rightarrow A \sim A$ )
    - One to translate from domain  $A \rightarrow B$
    - One to translate from domain  $B \rightarrow A$



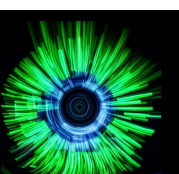
# UNFOLDING WITH ML

[PRL 124, 182001 (2020)]



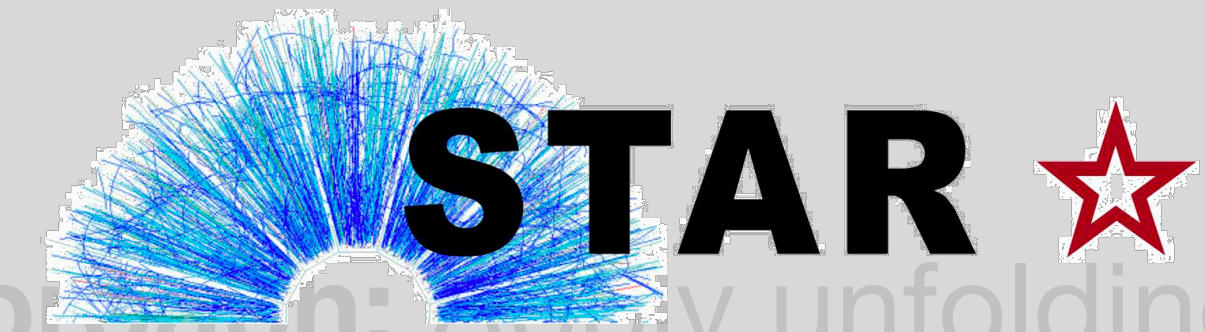
**Conventional Approach:** Apply unfolding procedure on a binned distribution and repeat for each observable.

**ML-based Approach:** Use ML to calculate weighting factors and unfold the phase space all at once, before the choice of binning or observable!



# UNFOLDING WITH ML

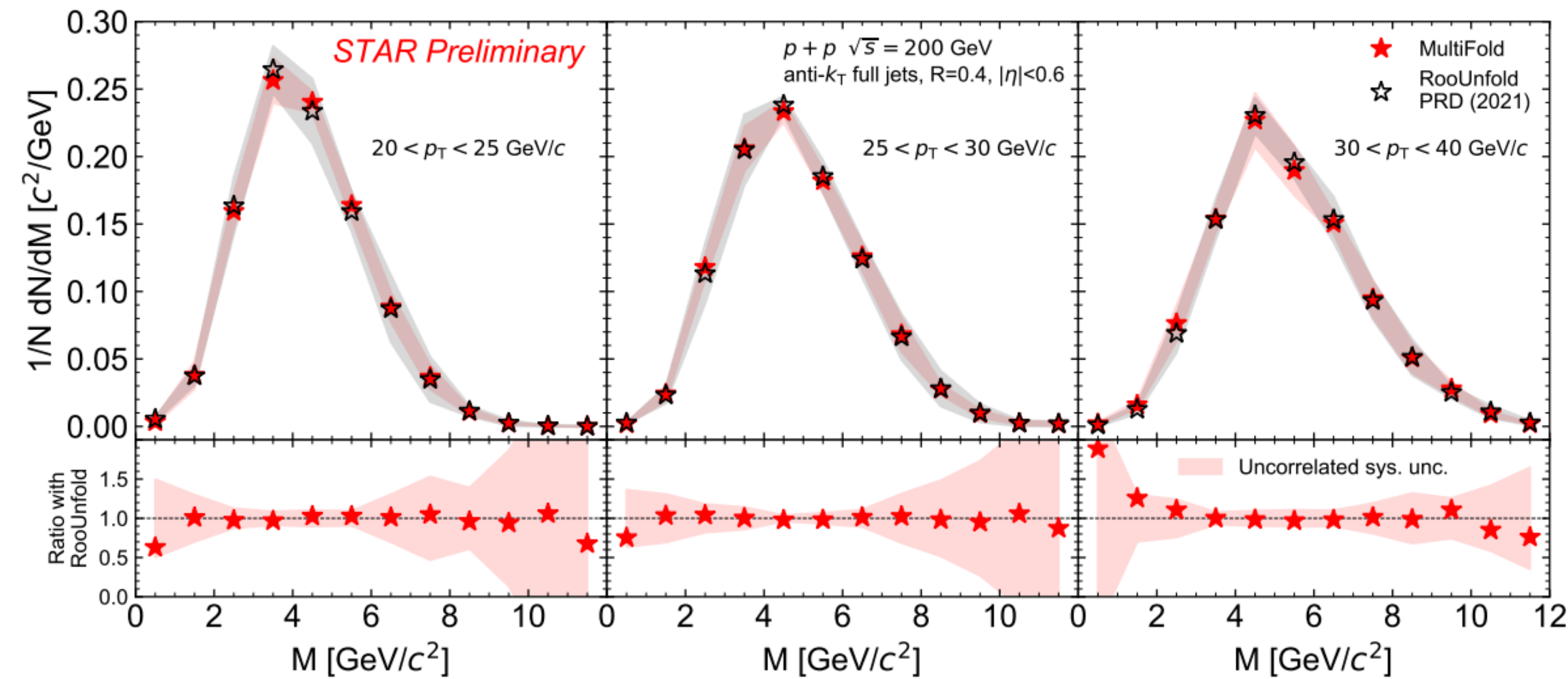
[PRL 124, 182001 (2020)]



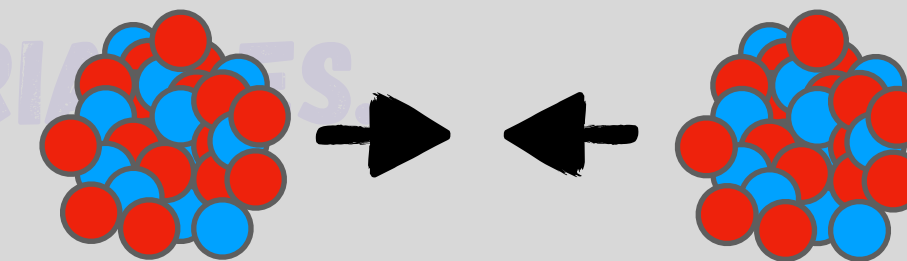
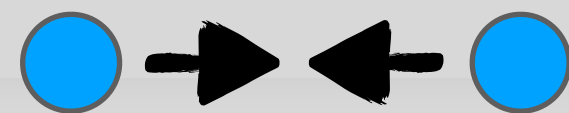
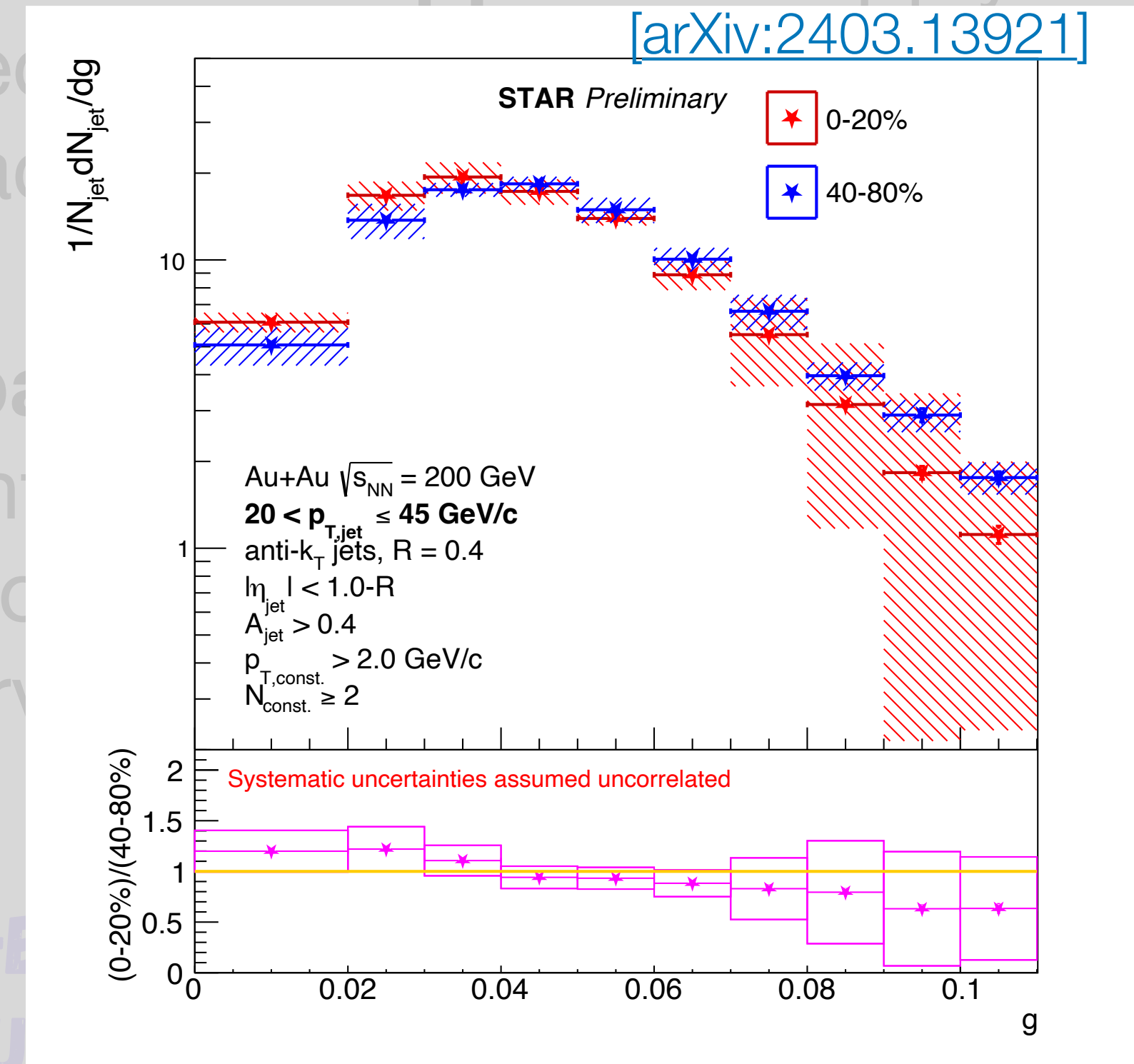
Conventional Approach: Apply unfolding procedure and repeat

Detector-level

Particle-level



Youqi Song, DIS 2023, [arXiv:2307.07718]



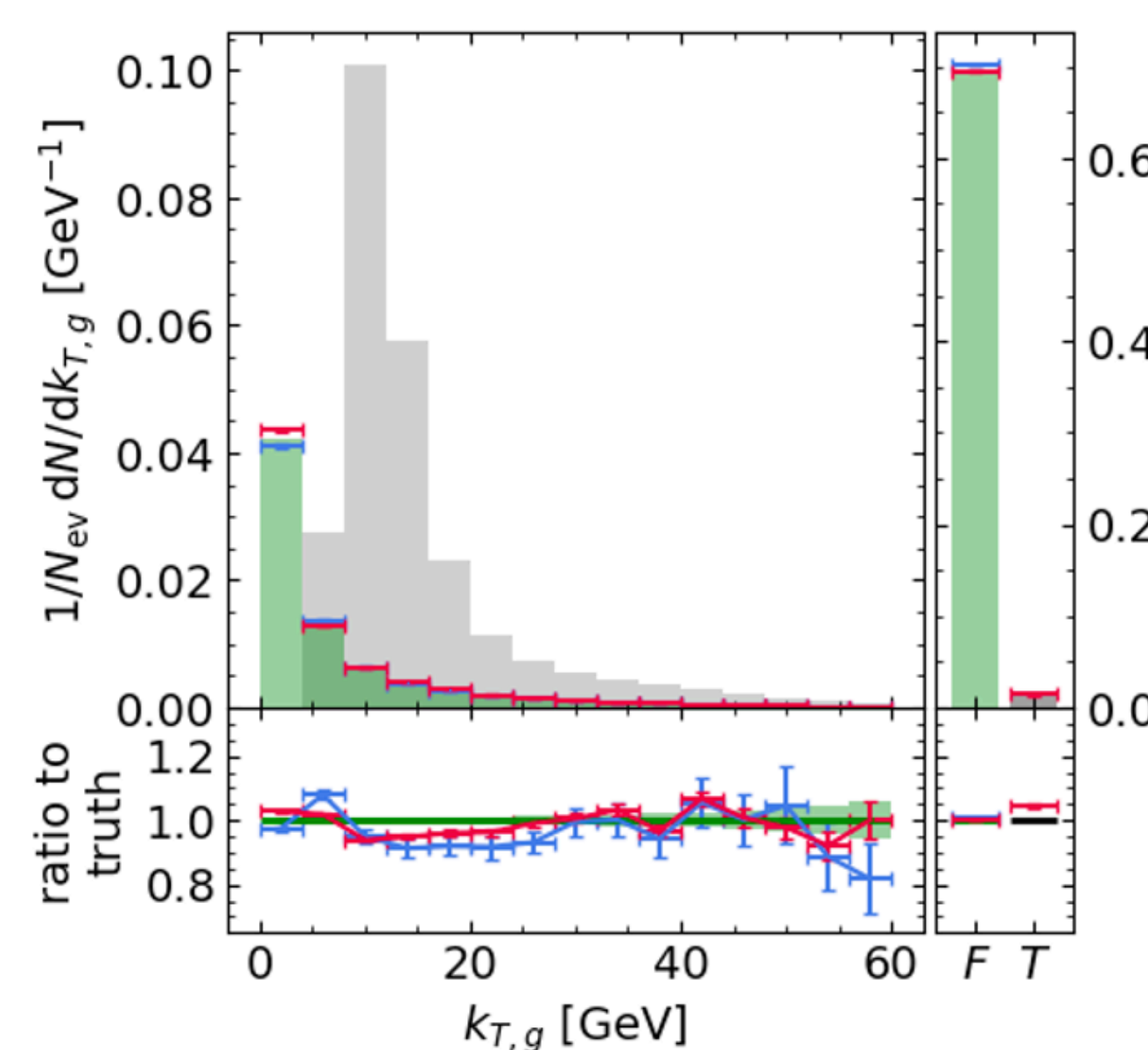
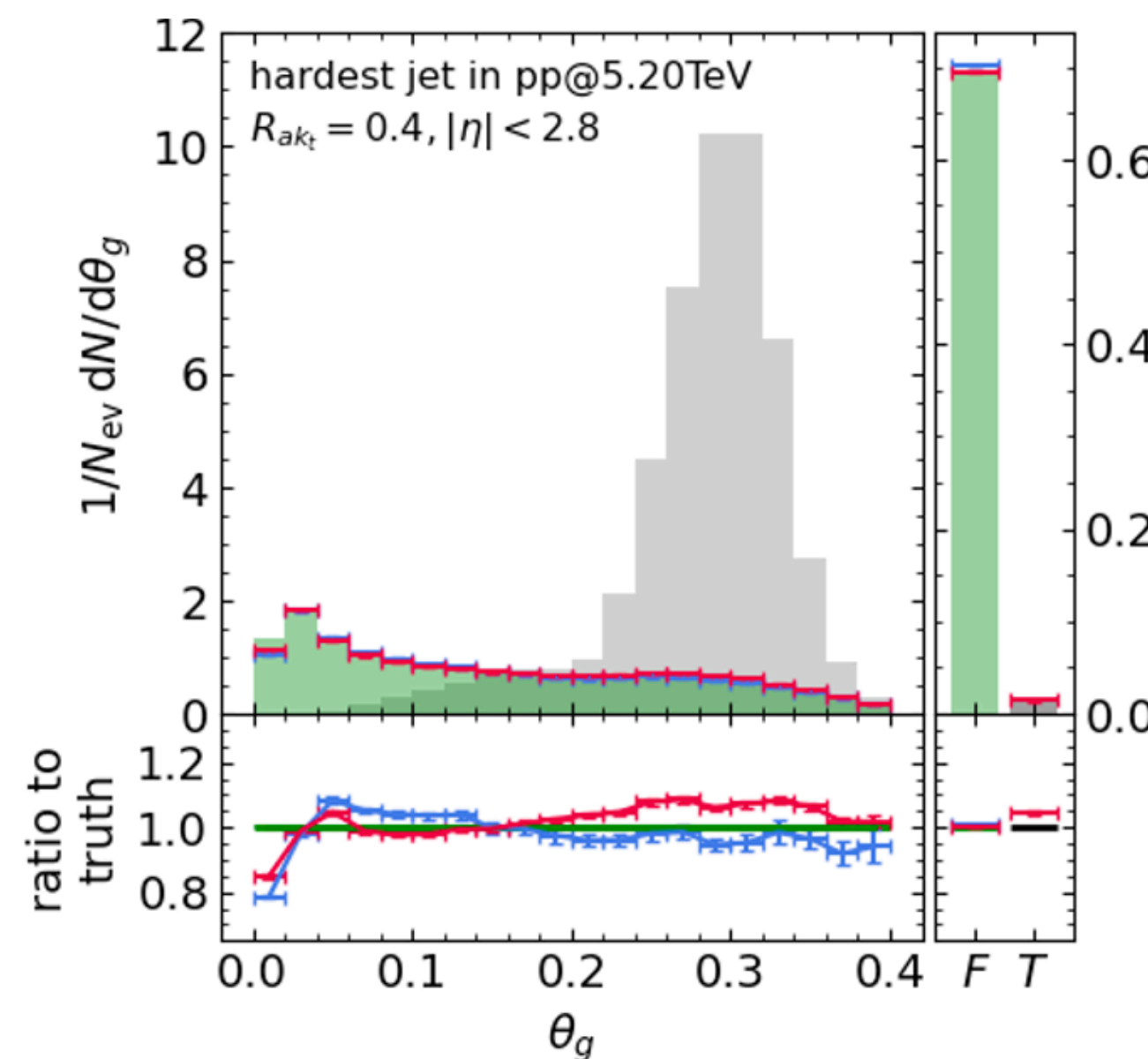
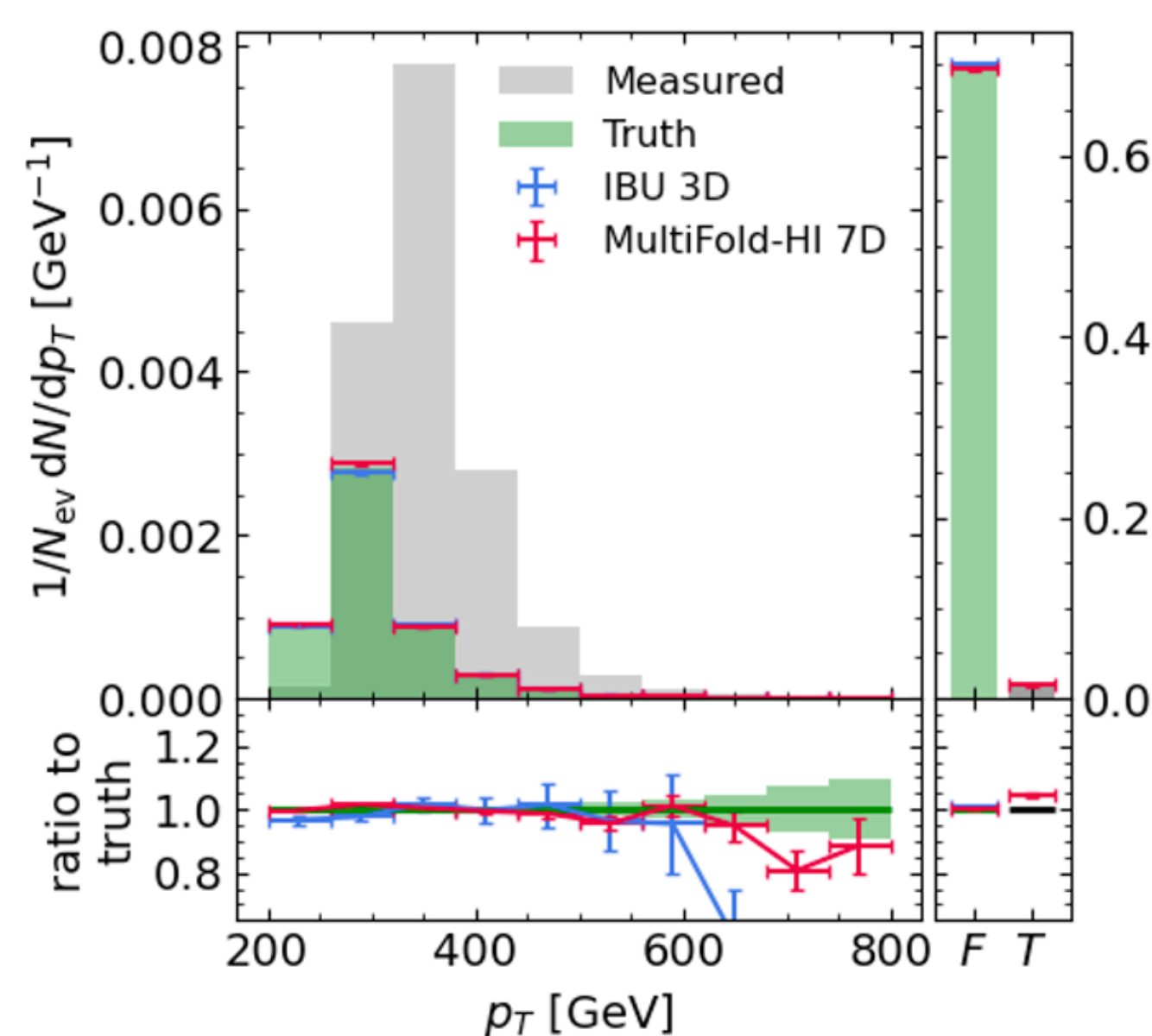
**ALLOWS FOR MULTI-DIFFERENTIAL MEASUREMENTS OF JET SUBSTRUCTURE IN PP AND AU+AU\***

\* model uncertainty not yet evaluated in Au+Au

# UNFOLDING WITH ML



- Tested for the first time on HI environment (PYTHIA/HERWIG + thermal background), similar or better performance to Bayesian unfolding in 3D.
- Modify the approach in [\[PRL 124 182001 \(2020\)\]](#) in order to also treat the case of ...
  - Measured events without true match (fakes, F)
  - True events that are not measured (trash, T)
- **No explicit background subtraction, built into MultiFold-HI!**

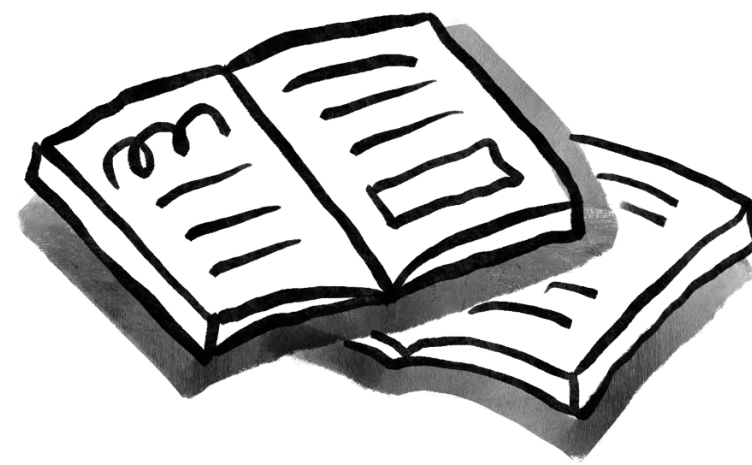


Truth/Measured: Herwig7 + thermal bkg. + DelpheS  
 Smeared/Generated: Pythia8 + thermal bkg. + DelpheS



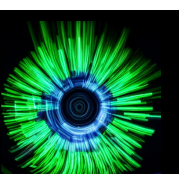
# ROADMAP

**WHAT IS AI/ML  
AND WHY IS IT  
USEFUL FOR THE  
ANALYSIS OF  
HARD PROBES?**



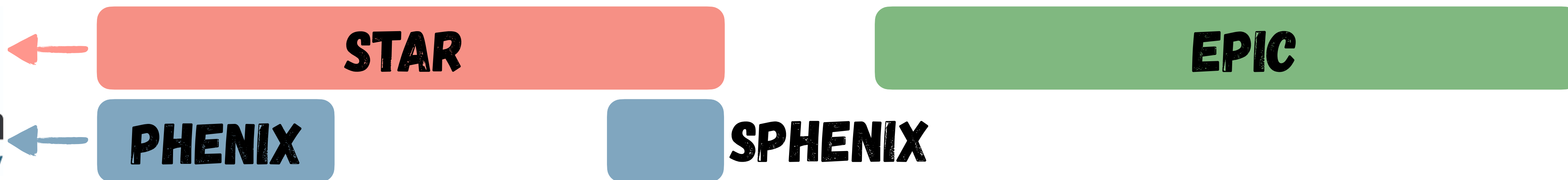
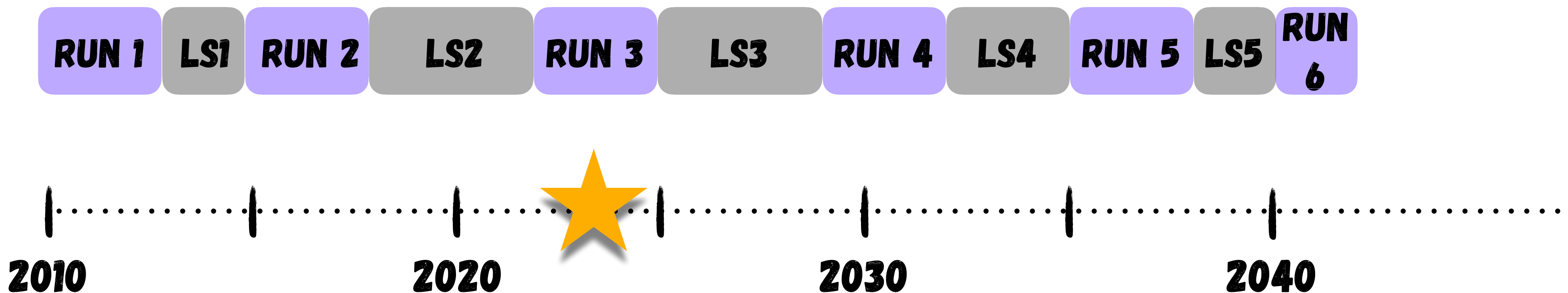
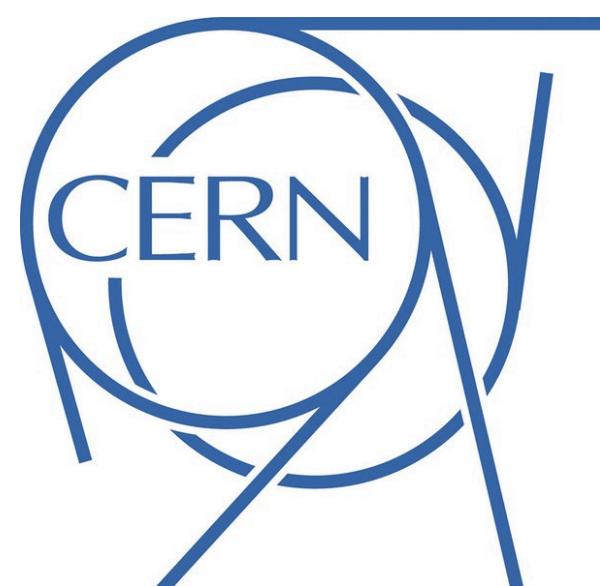
**HOW IS AI/ML  
CURRENTLY  
BEING USED FOR  
ANALYSIS?**

**WHERE ARE WE  
HEADING?**



# WHERE ARE WE GOING?

## LARGE HADRON COLLIDER

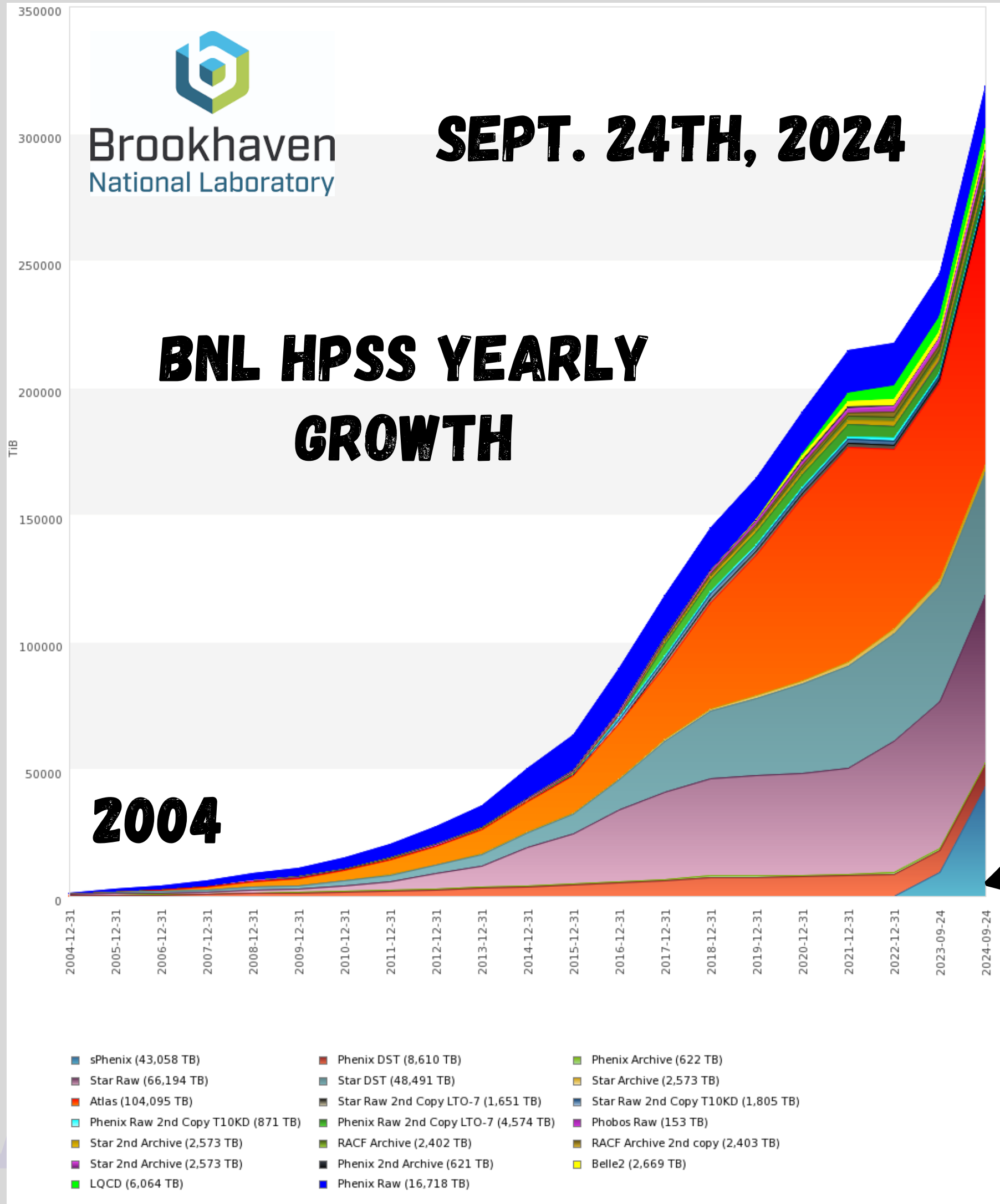


## RELATIVISTIC HEAVY ION COLLIDER

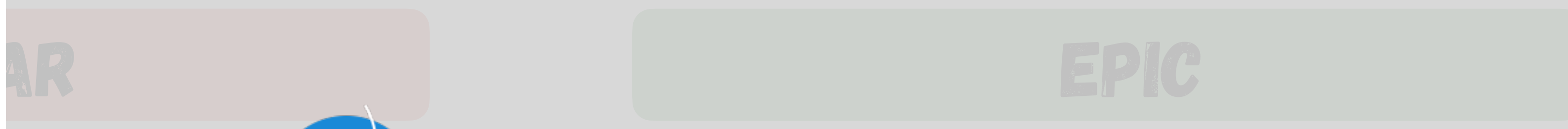
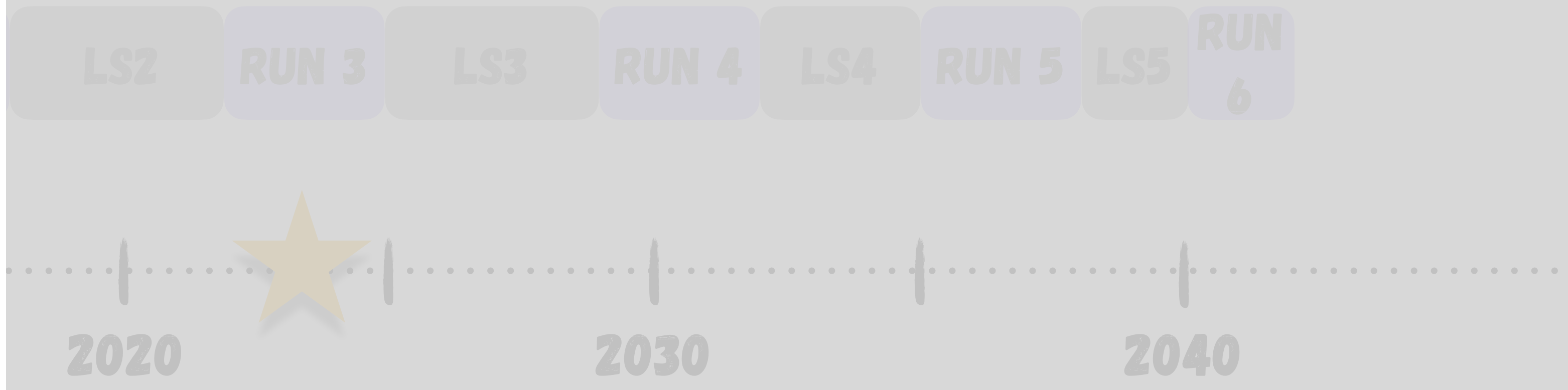
## ELECTRON ION COLLIDER

Very large volumes of will be taken and analyzed in the decades to come - new tools will be increasingly important!

# WHERE ARE WE GOING?



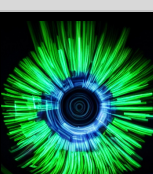
## LARGE HADRON COLLIDER



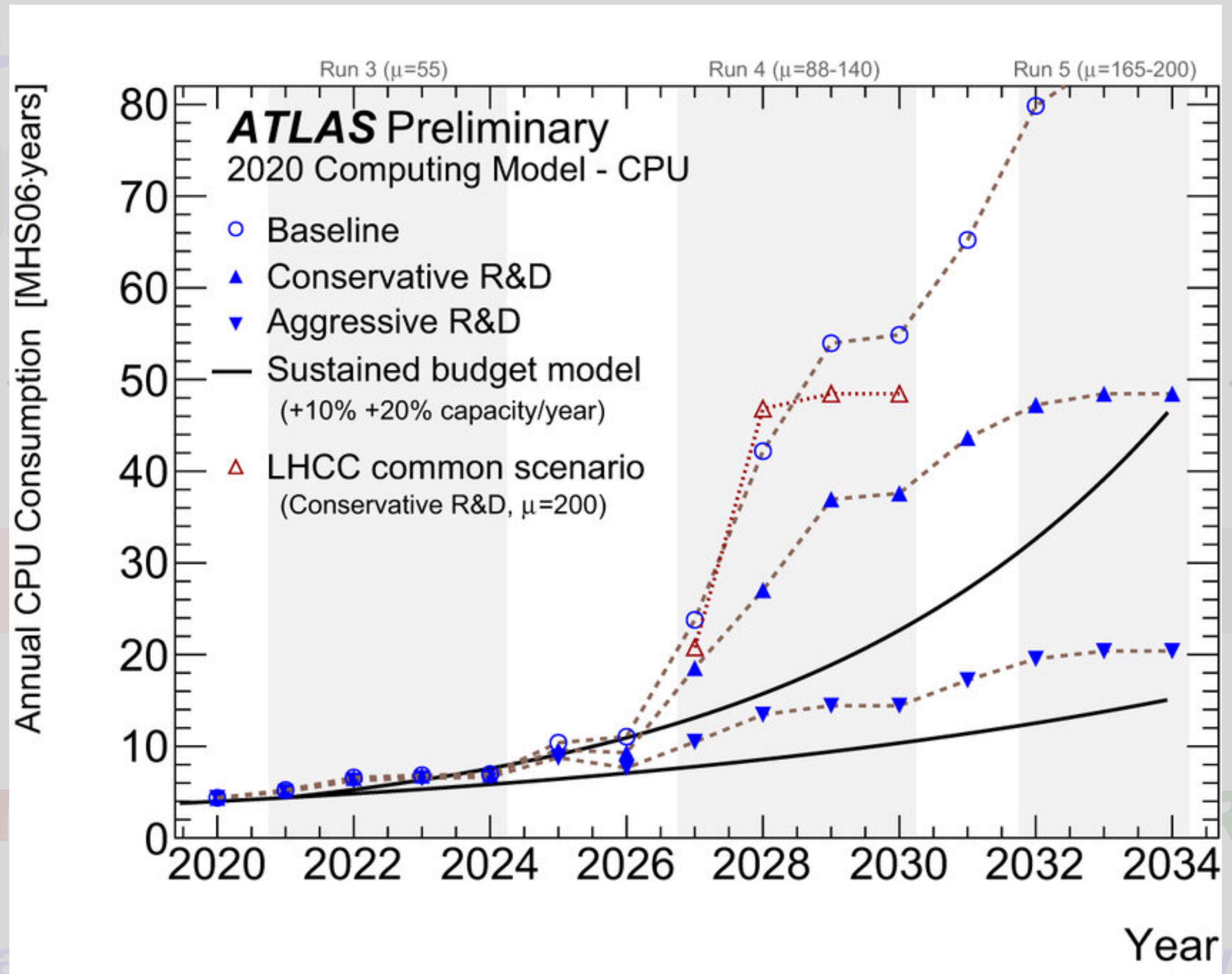
## HEAVY ION COLLIDER

## ELECTRON ION COLLIDER

**SEEING EXPONENTIAL GROWTH IN AMOUNT OF DATA STORED!**  
**SPHENIX QUICKLY BECOMING SIZABLE FRACTION OF THE TOTAL**



# WHERE ARE WE GOING?



**SAME TRENDS TRUE AT THE LHC!**



# OPEN QUESTIONS FOR NEXT ~5 YEARS

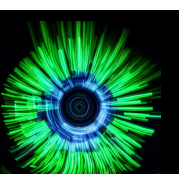
**HOW DO WE ASSIGN A SYSTEMATIC  
UNCERTAINTY FOR THE ML?**

**HOW DO WE CONSTRUCT  
MORE INTERPRETABLE  
MODELS?**



**HOW CAN WE MAKE ML-BASED  
APPLICATIONS REPRODUCIBLE?**

**DO WE NEED TO  
STANDARDIZE ML  
APPLICATIONS ACROSS  
EXPERIMENTS?**



# ML FOR UNDERLYING PHYSICS

“Data”-based learning complements simulation-based inference.

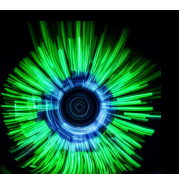
~ Given an answer  
~ “White Box” ML  
~ Underlying physics

~ Domain knowledge  
~ “Black Box” ML  
~ Answer

- Learning from data is difficult due to systematic experimental biases.
- Helpful in understanding uncertainties or shortcomings of models!

Proof of concept identifying the AP splitting function exists [\[PLB 829 \(2022\) 137055\]](#)

**THIS IS A LONG TERM EFFORT!**

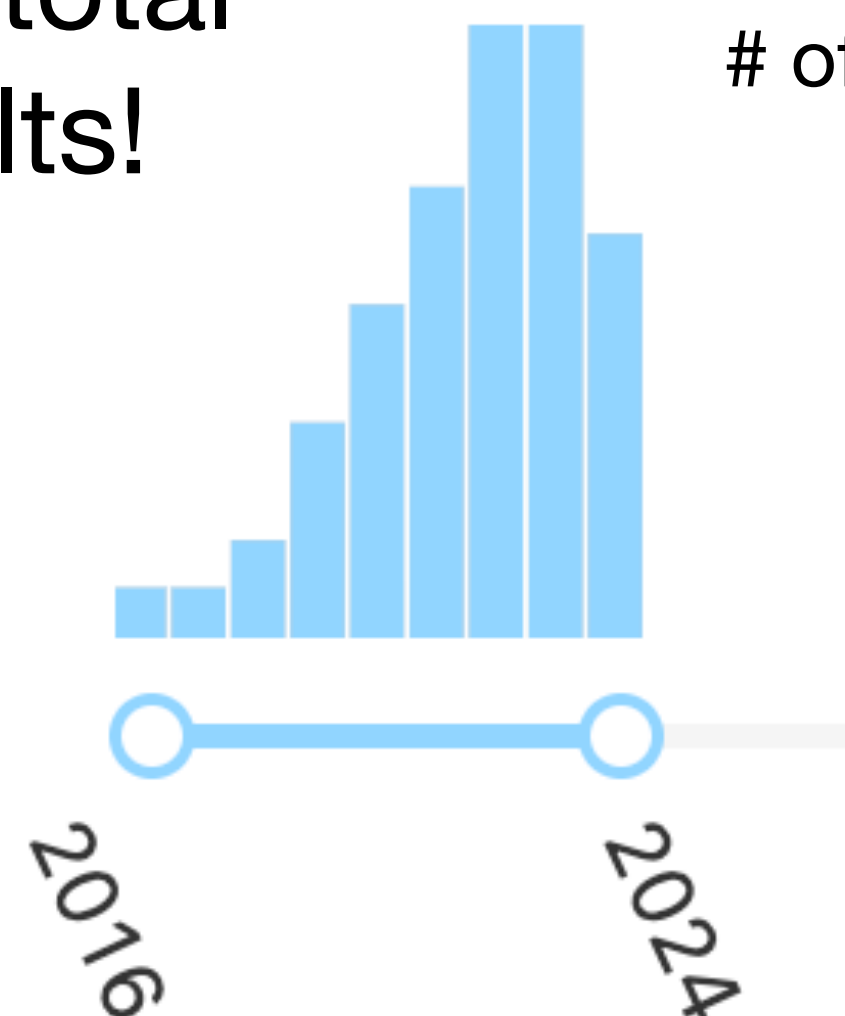


# CONCLUSIONS

- We are taking more data and making more complex measurements than ever before!
- Machine learning has led to new physics insights and can be used throughout the whole analysis pipeline!
  - Many great examples at this conference!
- Will be crucial at future facilities such as the HL-LHC and the EIC!

161 total results!

# of papers



Inspire HEP search results for “machine learning heavy ion”

**FUTURE IS VERY BRIGHT!**

**HP 2025?**





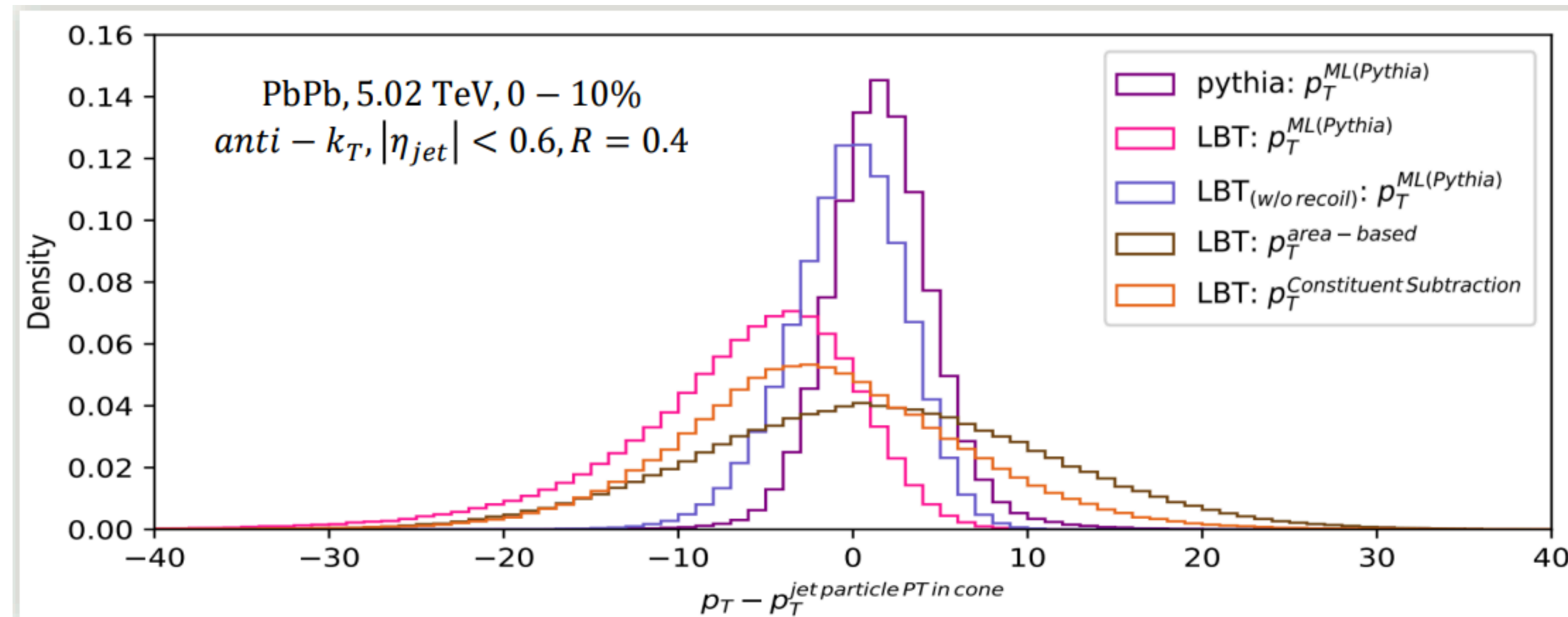
# Thank you!!

**Special thanks to Fabio Catalano, Changwhan Choi, Raymond Ehlers, Alexandre Falcão, Yeonju Go, Laura Havener, Maja Karwowska, Diptanil Roy, Youqi Song, Adam Tackas and MIT Heavy Ion group for useful discussions and feedback!**

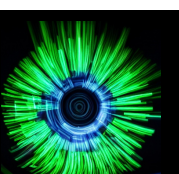
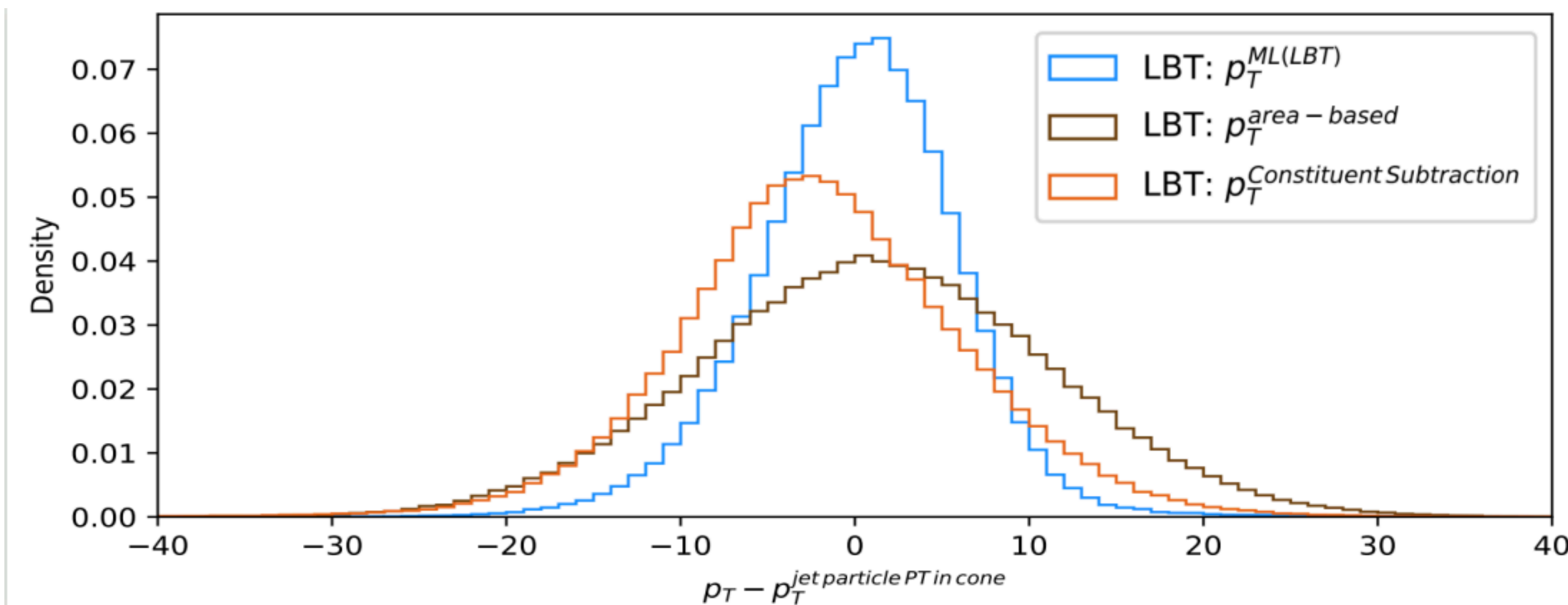


**Backup**

# STUDIES WITH NN JET PT RECONSTRUCTION



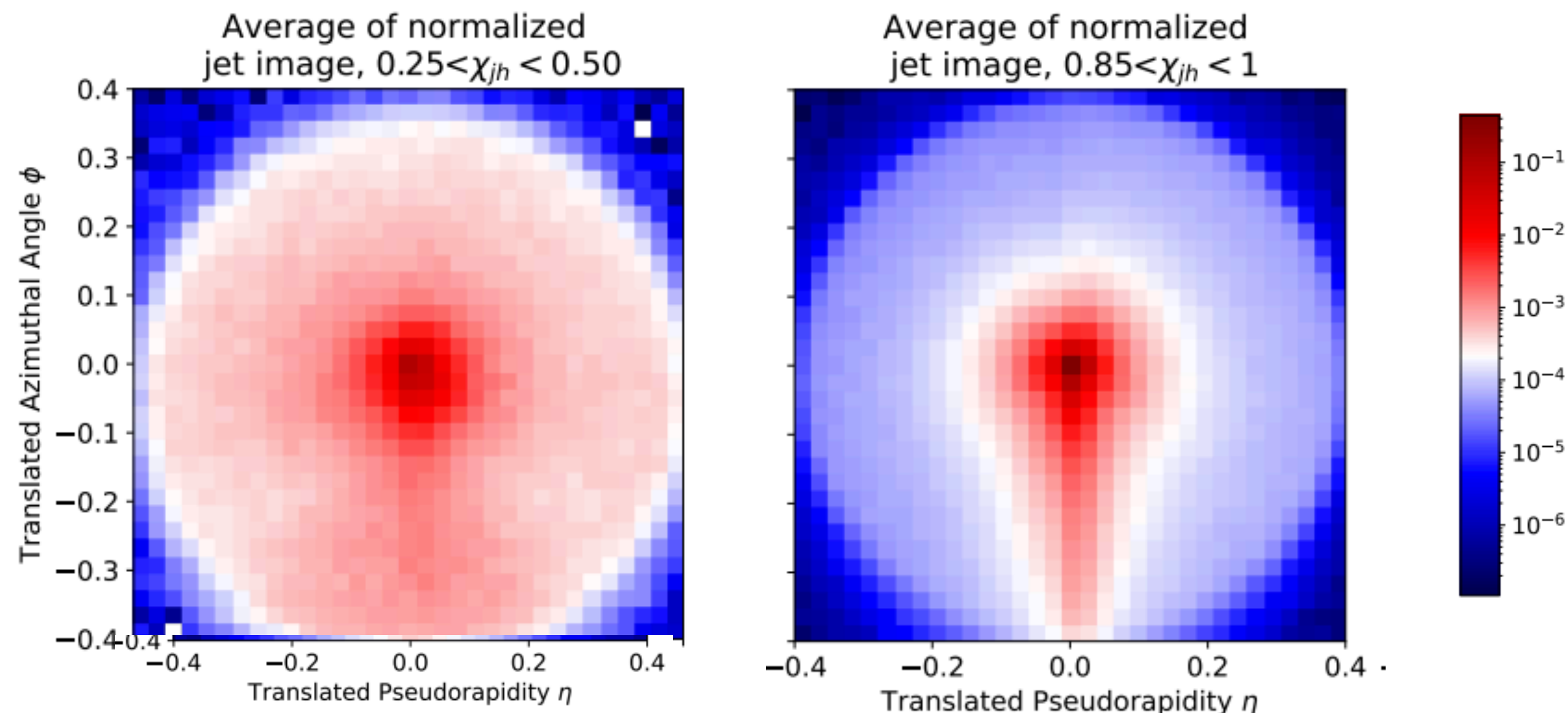
- See offset (bias) in  $\delta p_T$  when ML is trained in PYTHIA vs. LBT.
- Crucial for applications in data to correct for this bias in an unfolding procedure.
  - Apply same model on your data and the response matrix.



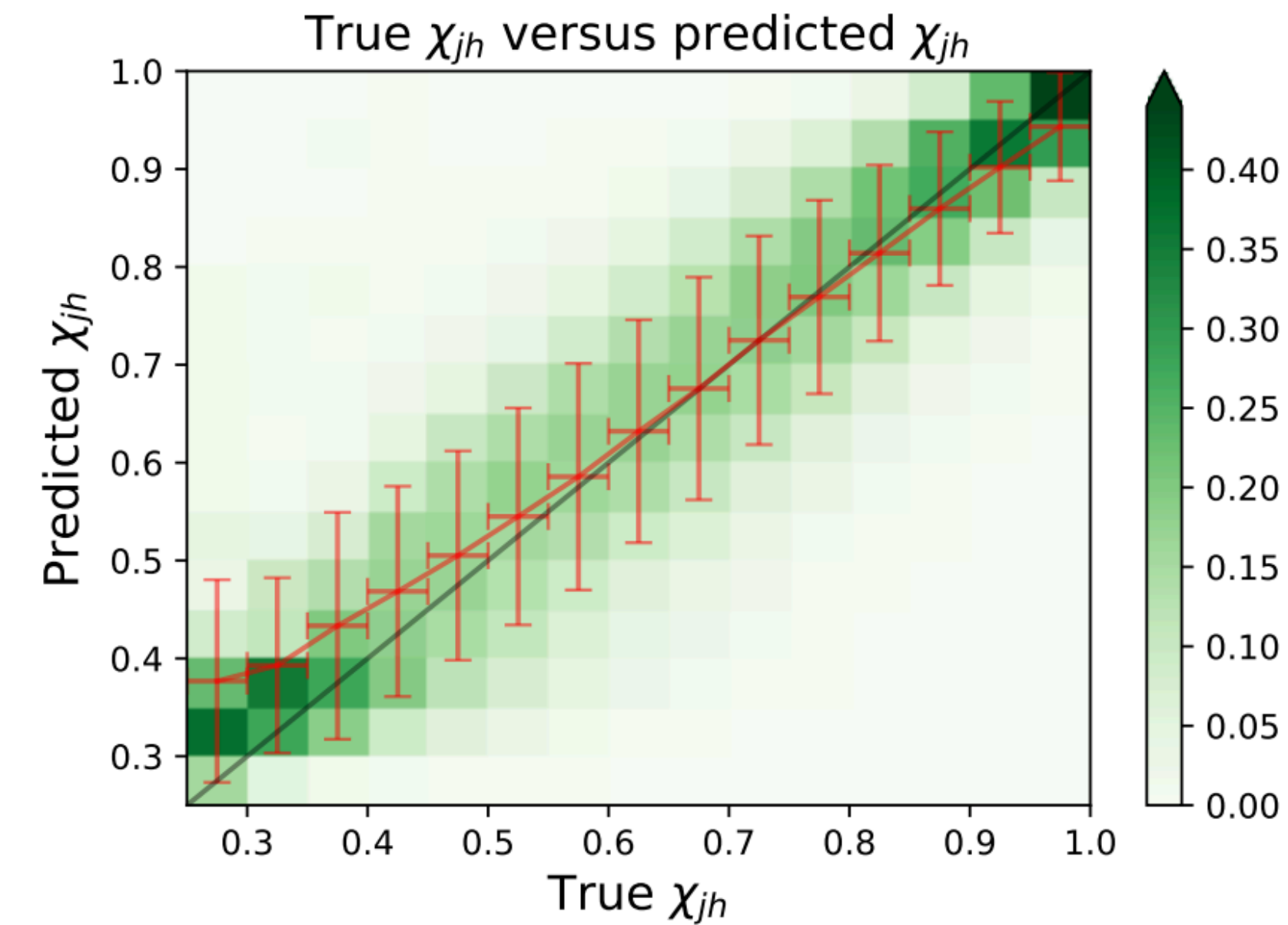
# DEEP LEARNING JET MODIFICATIONS

Use supervised learning on jet images with a CNN to perform the regression task of predicting the energy loss ratio in HI collisions (hybrid model).

[JHEP 2021, 206 (2021)]

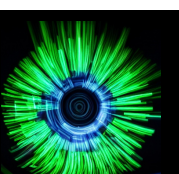


$$\chi_{jh} = \frac{E_f^h}{E_i^h}$$

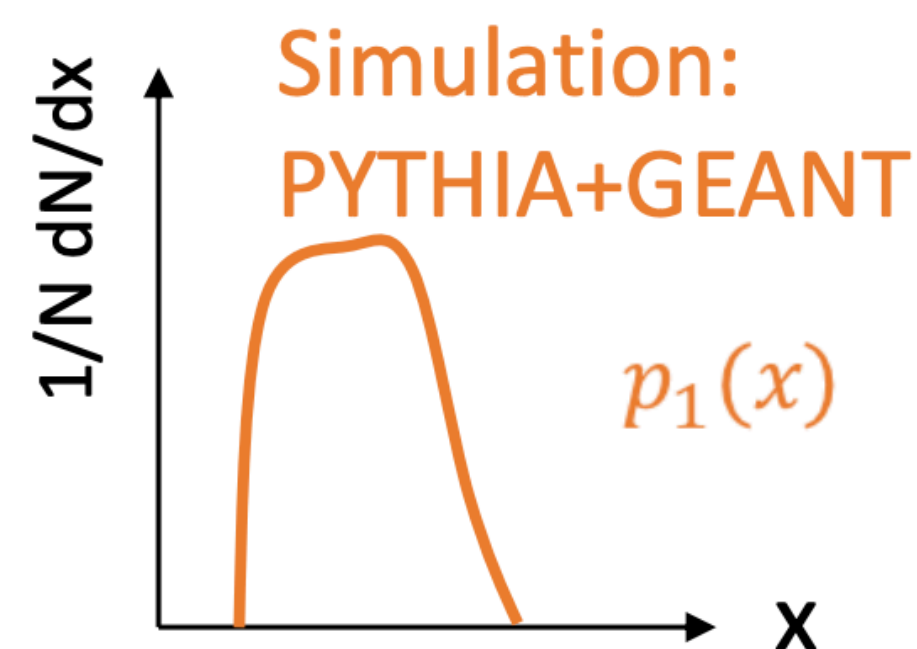
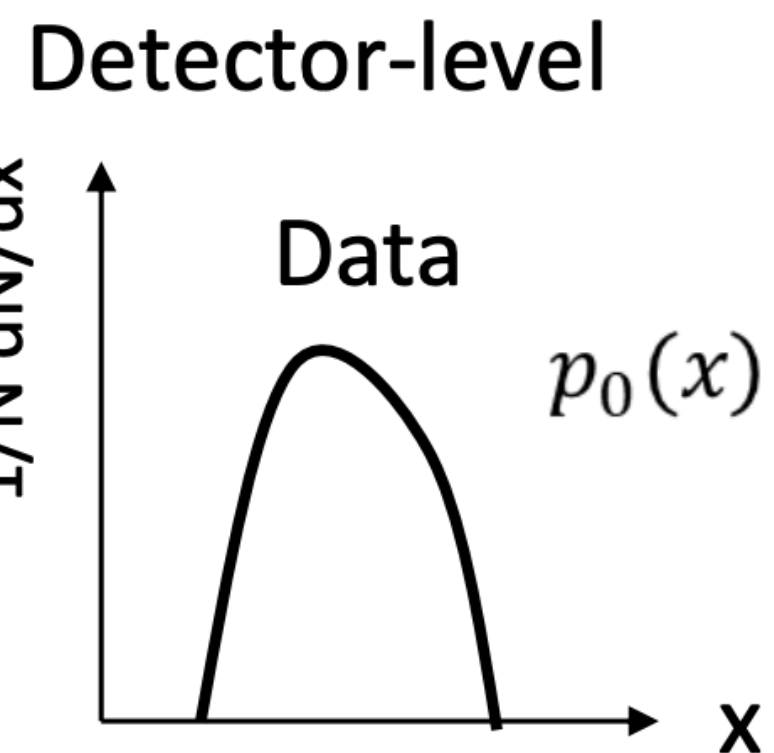
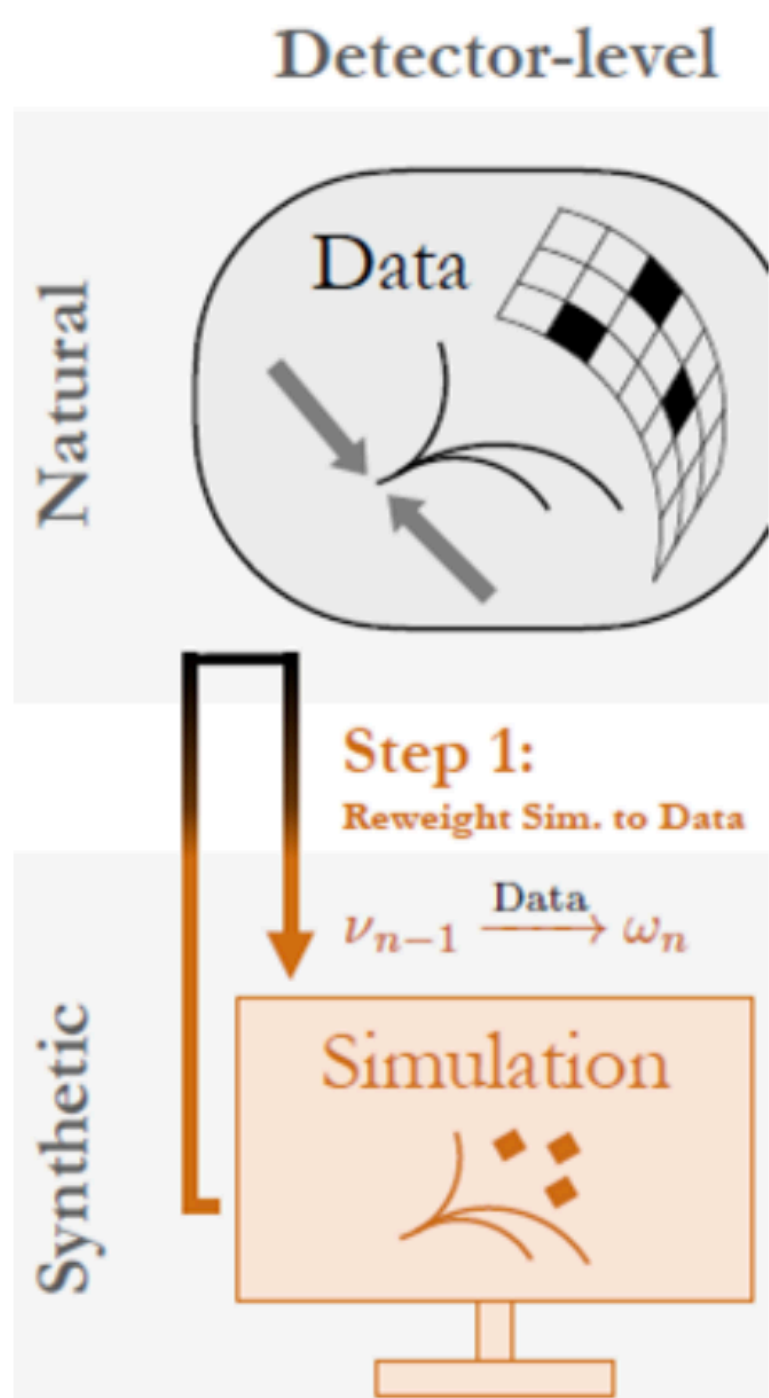


Shows good performance!

- Very useful to separate and study **quenched** vs. **unquenched** jets as well as extracting the initial energy of the jet. (Ideal probe of selection bias!)



# OMNIFOLD/MULTIFOLD



E.g., Iteration 1, step 1:

Weights:  $w(x) = p_0(x)/p_1(x)$

$\approx f(x)/(1 - f(x))$

Ok for 1D

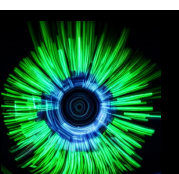
[\(Andreassen and Nachman PRD 101, 091901 \(2020\)\)](#)

where  $f(x)$  is a neural network and trained with the binary cross-entropy loss function

to distinguish jets coming from data vs from simulation

Unfolding → Reweighting histograms  
→ Classification → Neural network

Where does the machine learning part come in?

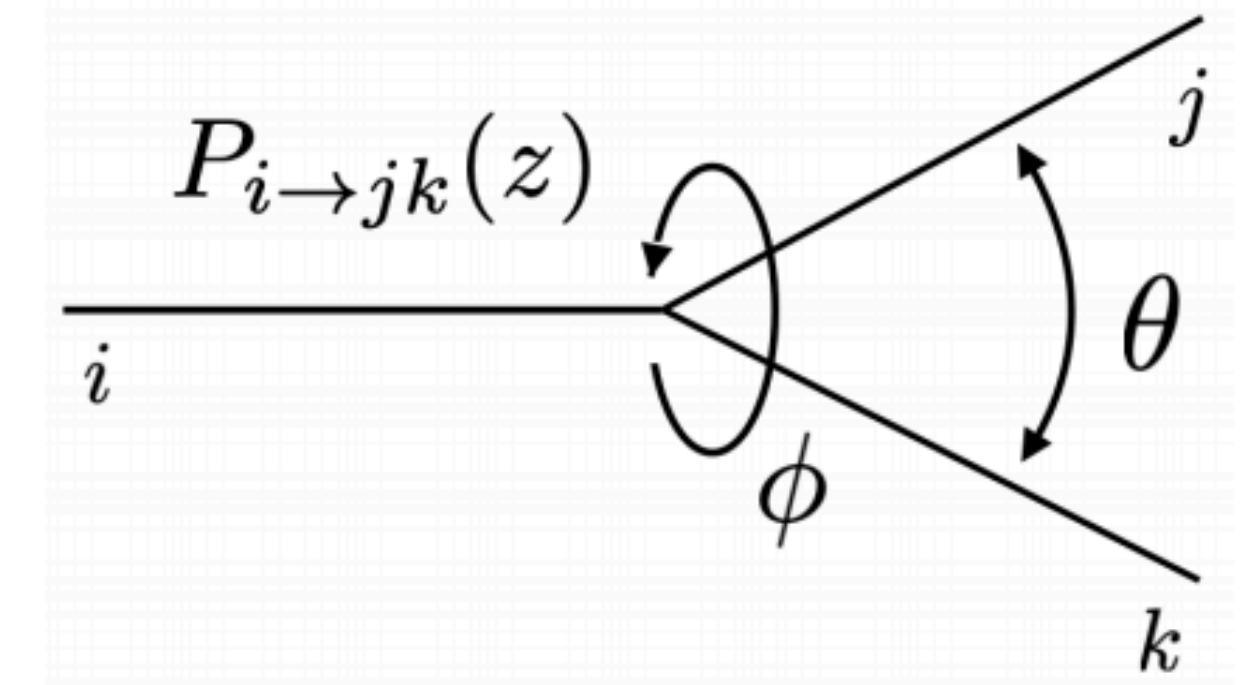




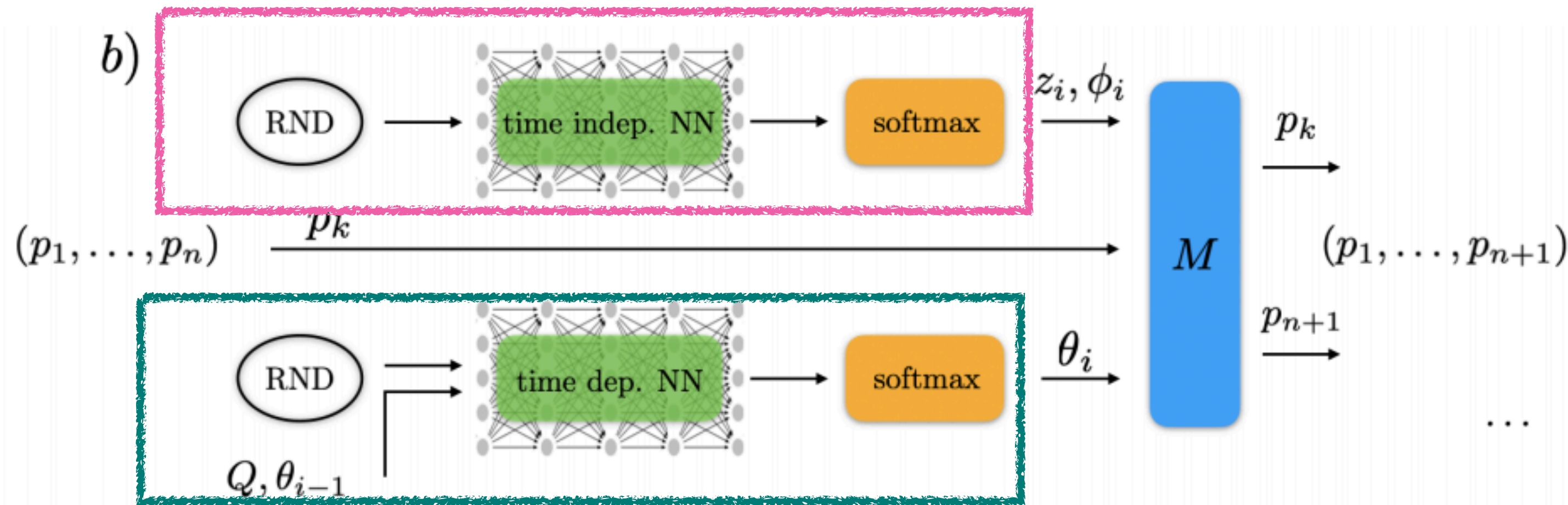
# PROOF OF CONCEPT

- Extract splitting function from the network in white-box ML.

Done with a GAN split into two components.

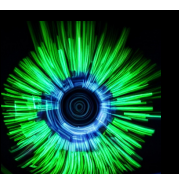
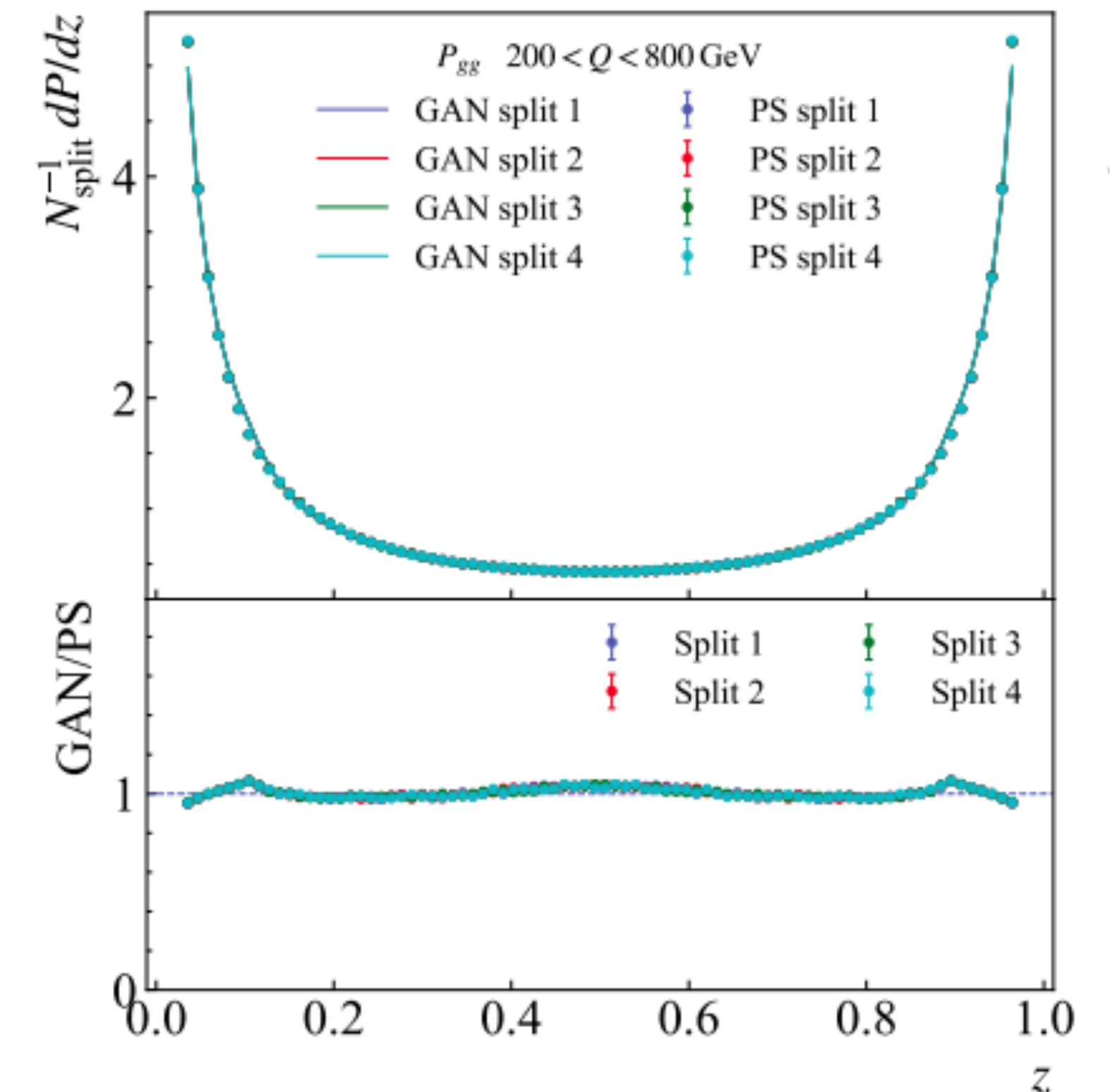


## 1. Time independent learns the $z, \phi$



## 2. Time dependent learns the $\theta$

Was able to reproduce AP splitting function.



# EVENT CLASSIFICATION AT THE EIC

[JHEP 03 (2023) 085]

- Study the effectiveness of ML-based classifiers to
  - Identify the flavor of the jet
  - Identify the underlying hard process of the collision

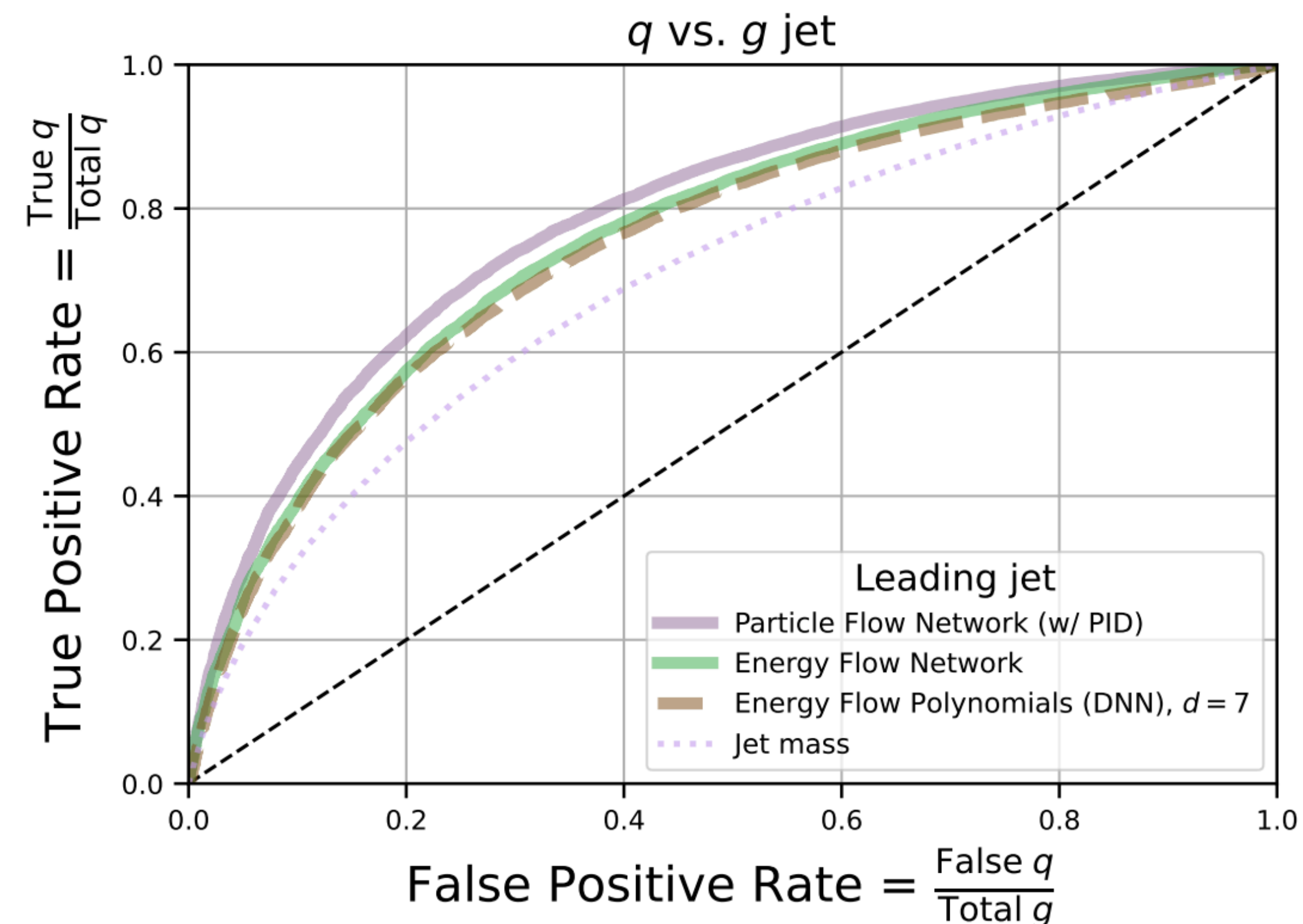
- Additionally study the effectiveness of different ways of representing information

- Particle Flow Networks [JHEP 01 (2019) 121]

$$f(p_1, \dots, p_N) = F\left(\sum_{i=1}^N \Phi(p_i)\right) \quad p_i = (z_i, \eta_i, \phi_i, \text{PID}_i)$$

- Energy Flow Polynomials [JHEP 04 (2018) 013]

$$\text{EFP}_G = \sum_{i_1} \cdots \sum_{i_V} z_{i_1} \cdots z_{i_V} \prod_{(k,l) \in E} \theta_{i_k i_l}$$



Indications that ML-based methods will have an improved performance over traditional techniques!

See also event classification with large language models, [arXiv:2404.05752]



# ML FOR PARTICLE IDENTIFICATION

TRACK DATA



NN TRAINED FOR  
A SPECIFIC  
PARTICLE SPECIES



CERTAINTY VALUE

- Use NN trained on a specific particle type to predict a certainty value that is then compared to a pre-set threshold.
- Decide threshold based on efficiency/purity tradeoff.
- Takes into account particles from different sub-detectors (here TPC, TOF, TRD of ALICE) , robust against missing data.

Model	Precision	Recall
Standard	99.99 ± 0.01	78.37 ± 0.01
Ensemble	97.47 ± 0.25	99.46 ± 0.21
Mean	97.31 ± 0.07	99.52 ± 0.07
Proposed	97.49 ± 0.06	99.54 ± 0.05
Regression	97.33 ± 0.06	99.49 ± 0.07

Pion

Model	Precision	Recall
Standard	99.40 ± 0.01	59.72 ± 0.03
Ensemble	97.16 ± 0.46	93.74 ± 0.30
Mean	97.85 ± 0.41	93.34 ± 0.32
Proposed	97.80 ± 0.44	93.86 ± 0.27
Regression	97.38 ± 0.40	93.67 ± 0.38

Proton

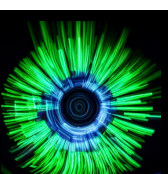
Model	Precision	Recall
Standard	92.87 ± 0.01	60.37 ± 0.05
Ensemble	91.18 ± 02.00	82.72 ± 01.42
Mean	90.83 ± 01.71	82.32 ± 0.96
Proposed	91.55 ± 0.71	83.68 ± 0.82
Regression	91.17 ± 01.00	81.78 ± 0.21

Kaon

- When comparing the **standard method** to the **proposed method**, proposed method has better balance of precision (purity) and recall (efficiency)!

See also LHCb NN to identify calo hits [[Int. J. Mod. Phys. A 30, 1530022 \(2015\)](#)], ATLAS Electron PID w/ CNN [[ATL-PHYS-PUB-2023-001](#)]

CMS Deep NN to identify hadronic  $\tau$ -lepton decays [[JINST 17 \(2022\) P07023](#)]



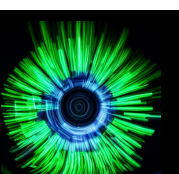
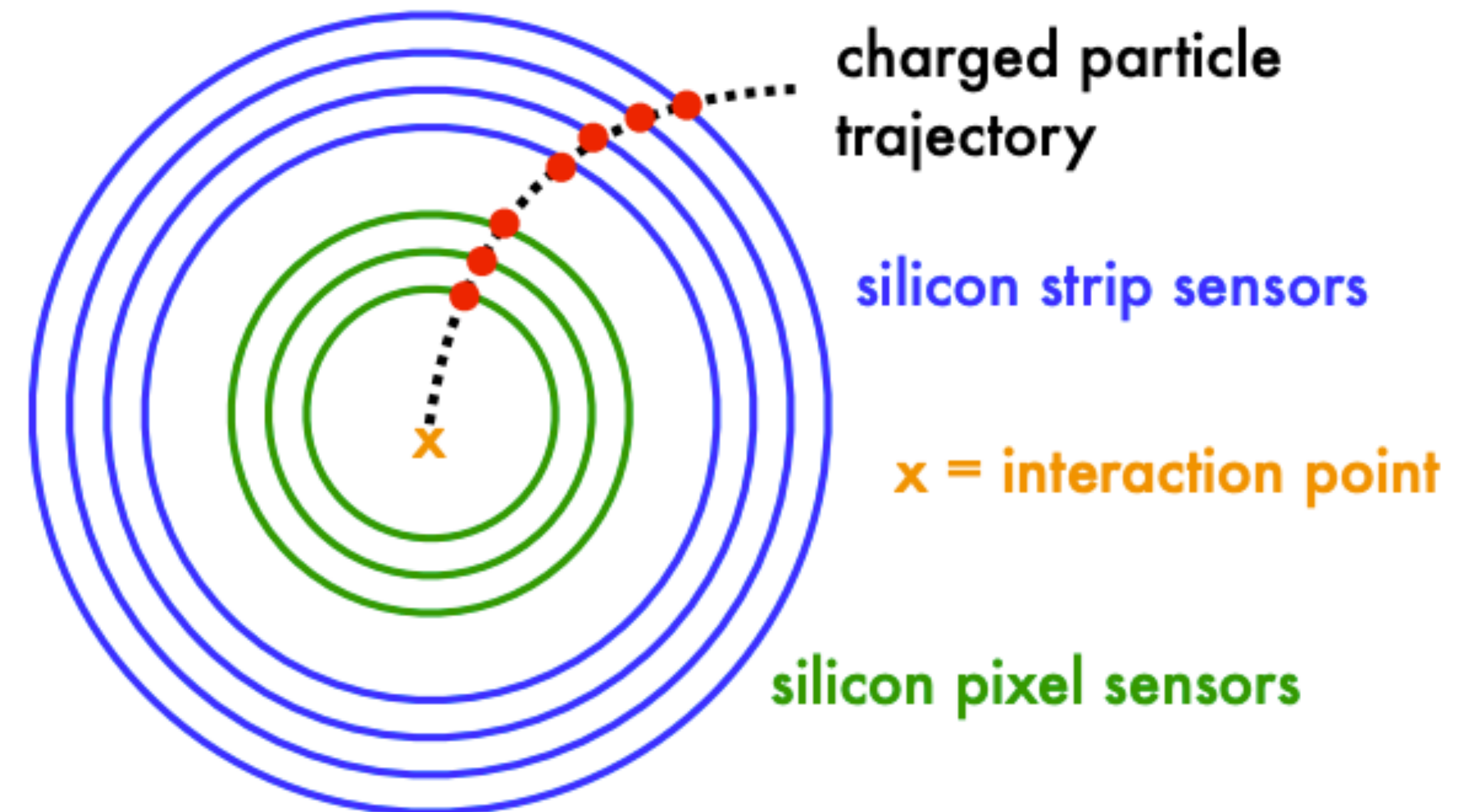
# TRACK RECONSTRUCTION AT THE HL-LHC

- Data volume and reconstruction will also be a problem for the HL-LHC
  - Reconstructing charged particle trajectory is computationally expensive - increases with the power of the multiplicity.

**Standard approach:** Kalman Filter used to locate hits in charged particle trajectory

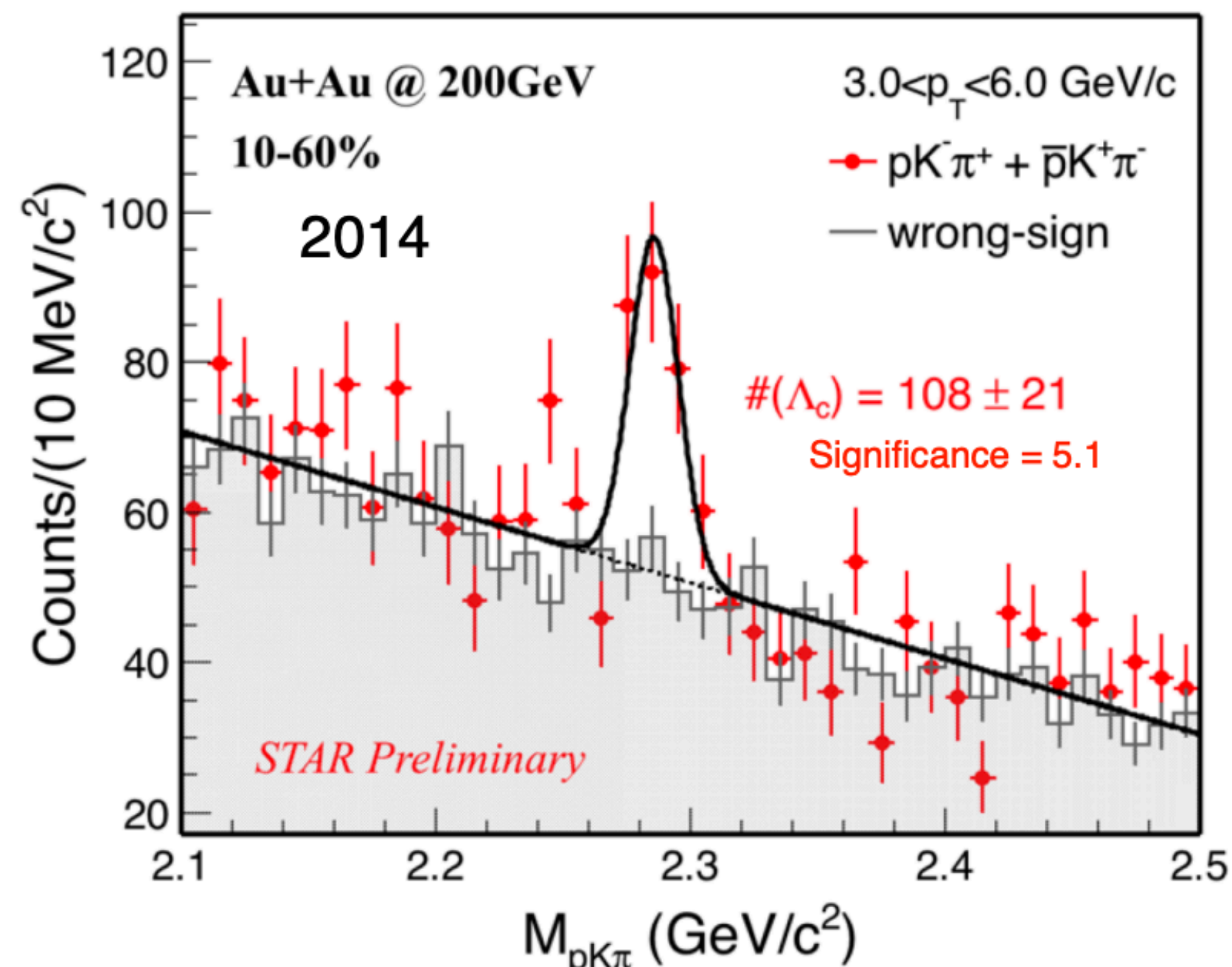
**ML-based approach:** Use ML tools to speed this up such as...

- Recurrent Neural Network [\[arXiv:2212.02348\]](https://arxiv.org/abs/2212.02348)
- Convolutional neural network [\[See Here\]](#)

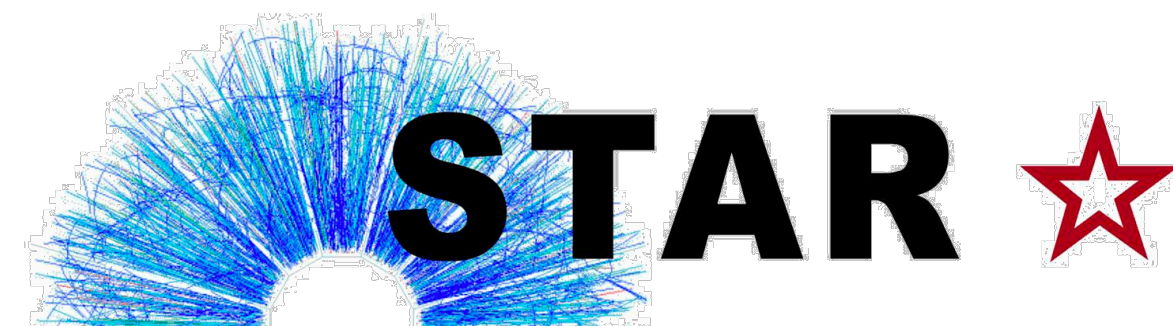
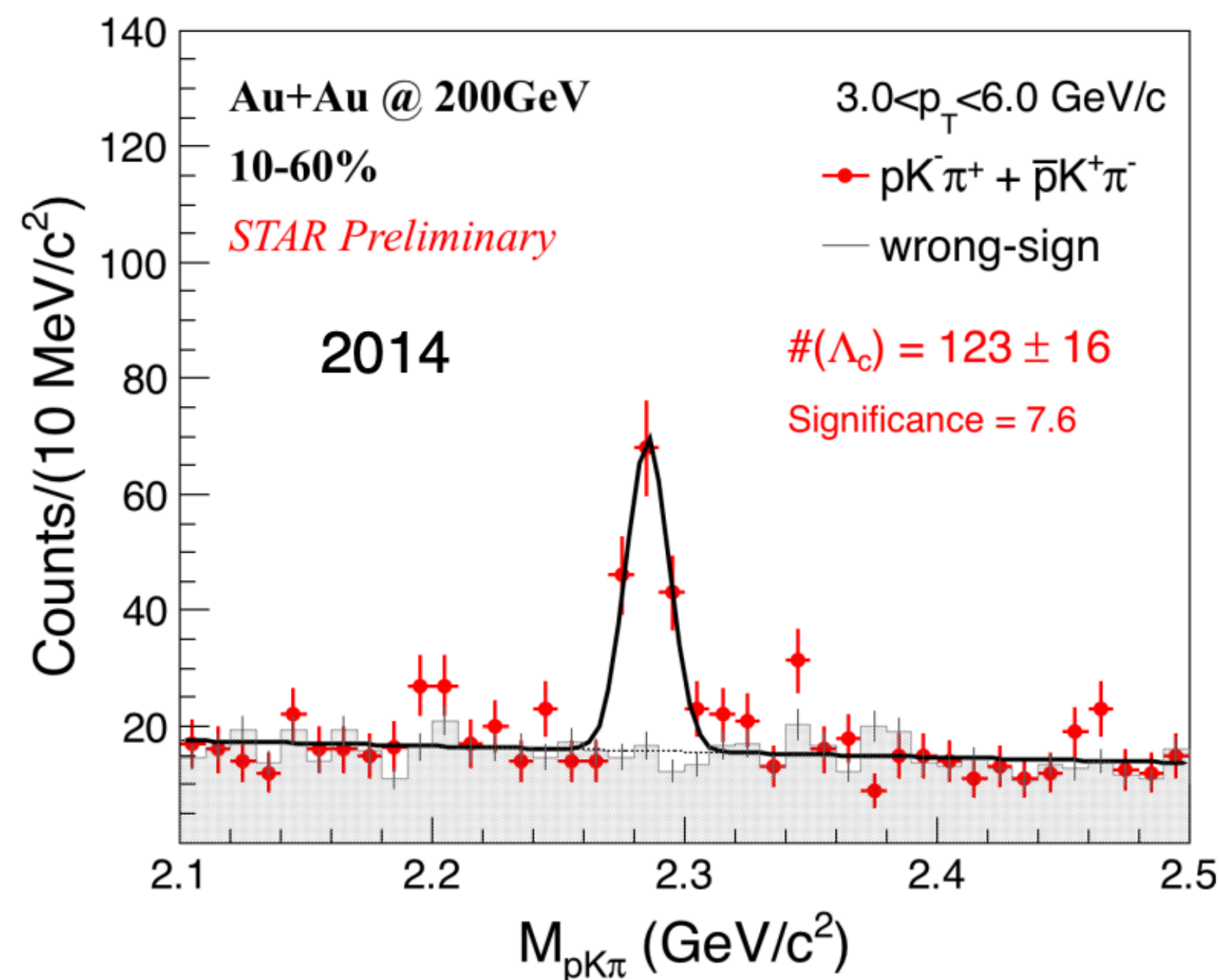


# SIGNAL/BACKGROUND DISCRIMINATION

Traditional Techniques

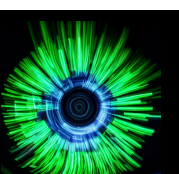


With BDT



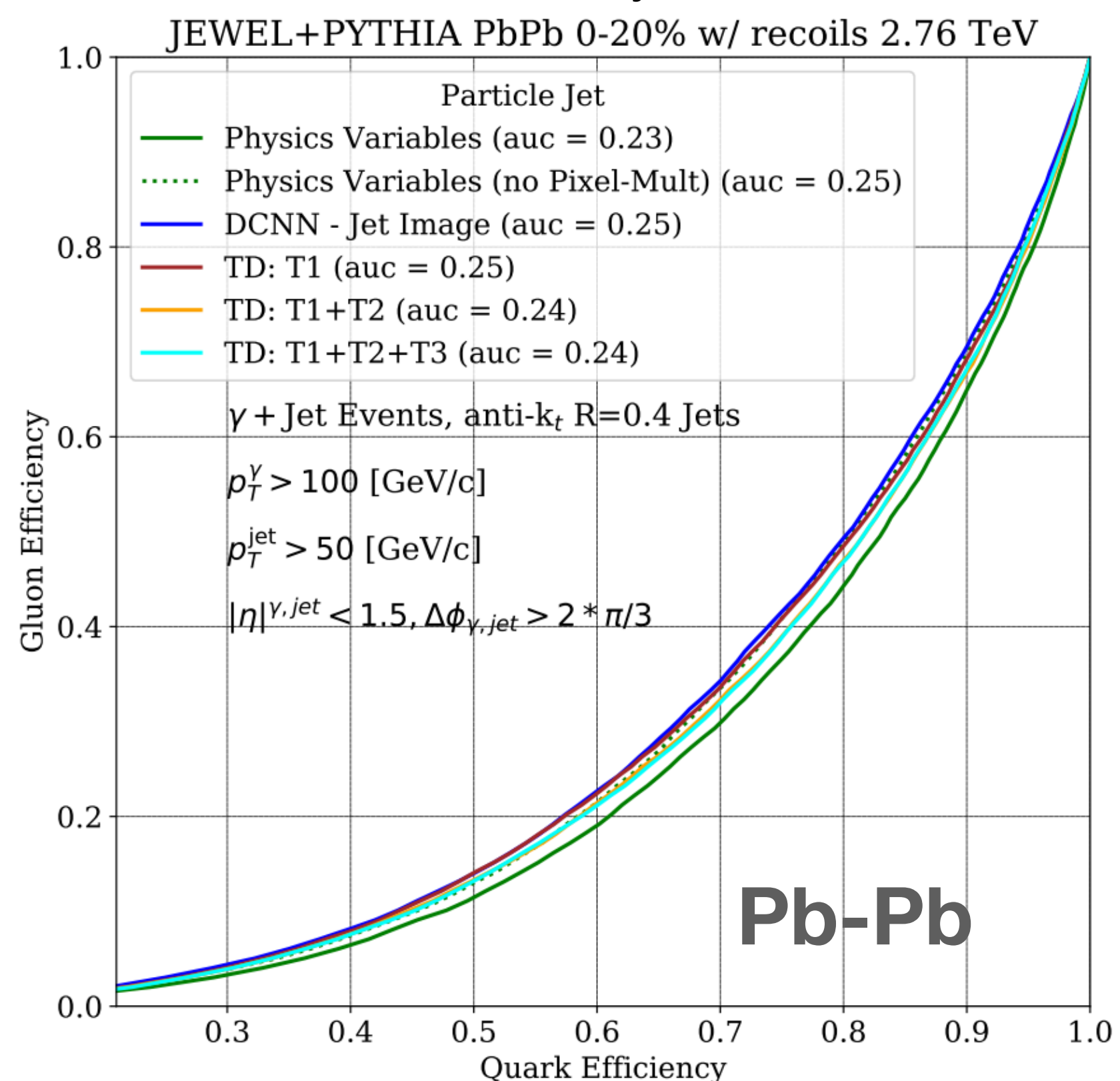
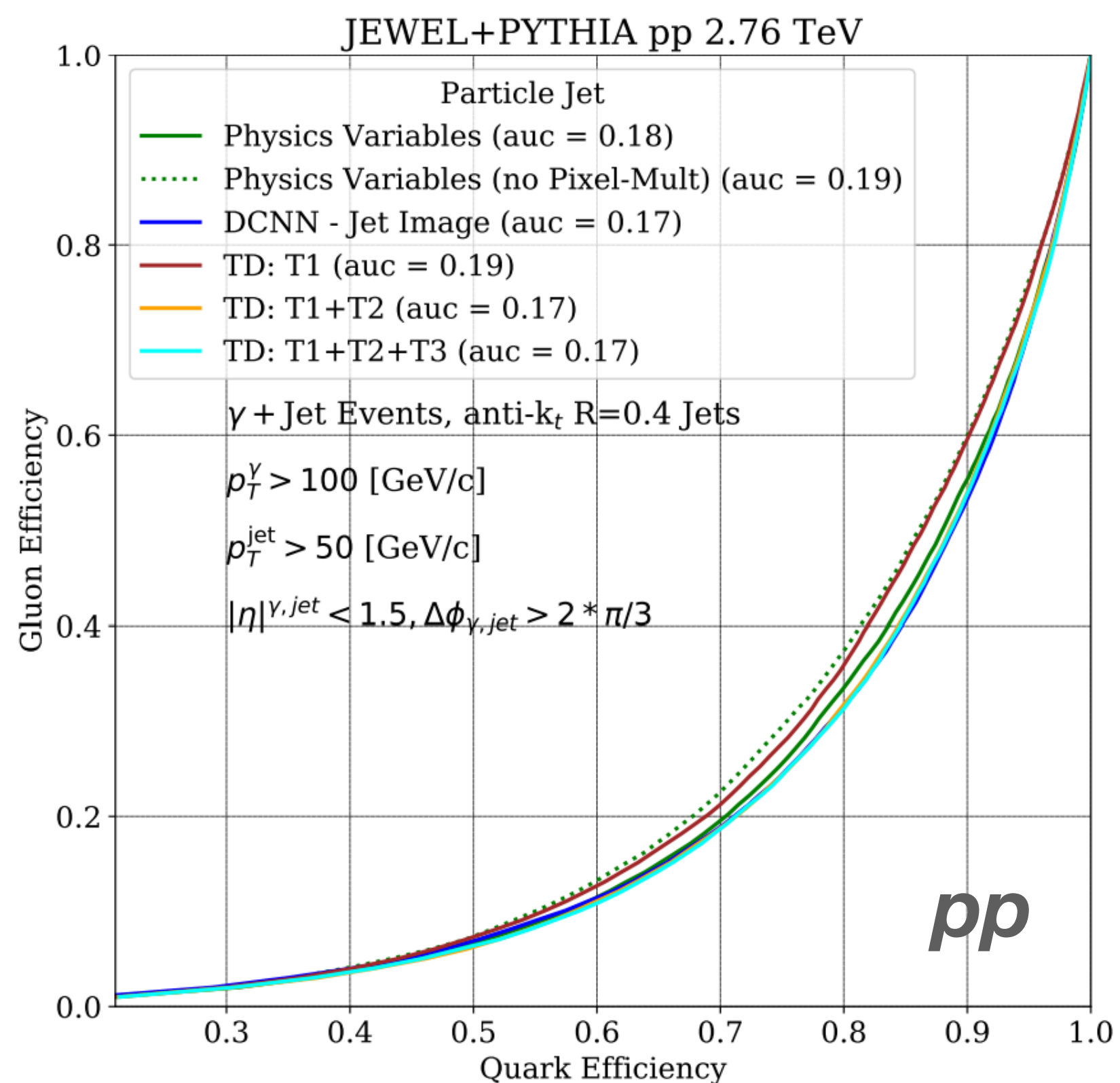
[[PRL 124, 172301 \(2020\)](#)]

- **Boosted Decision Tree** implemented in **ROOT TMVA** to optimize signal for  $\Lambda_c$  baryon production.
- Trained in a supervised manner with **EvtGen**
- 50% increase in signal significance with ML!



# QUARK VS. GLUON JETS WITH ML

Y. Chien, R. Elayavalli: 1803.3589



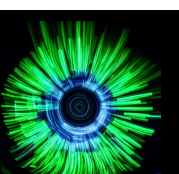
→ Lower the curve, the better the performance.

→ All methods explored tend to perform consistently, indicating that they may be picking up on similar features.

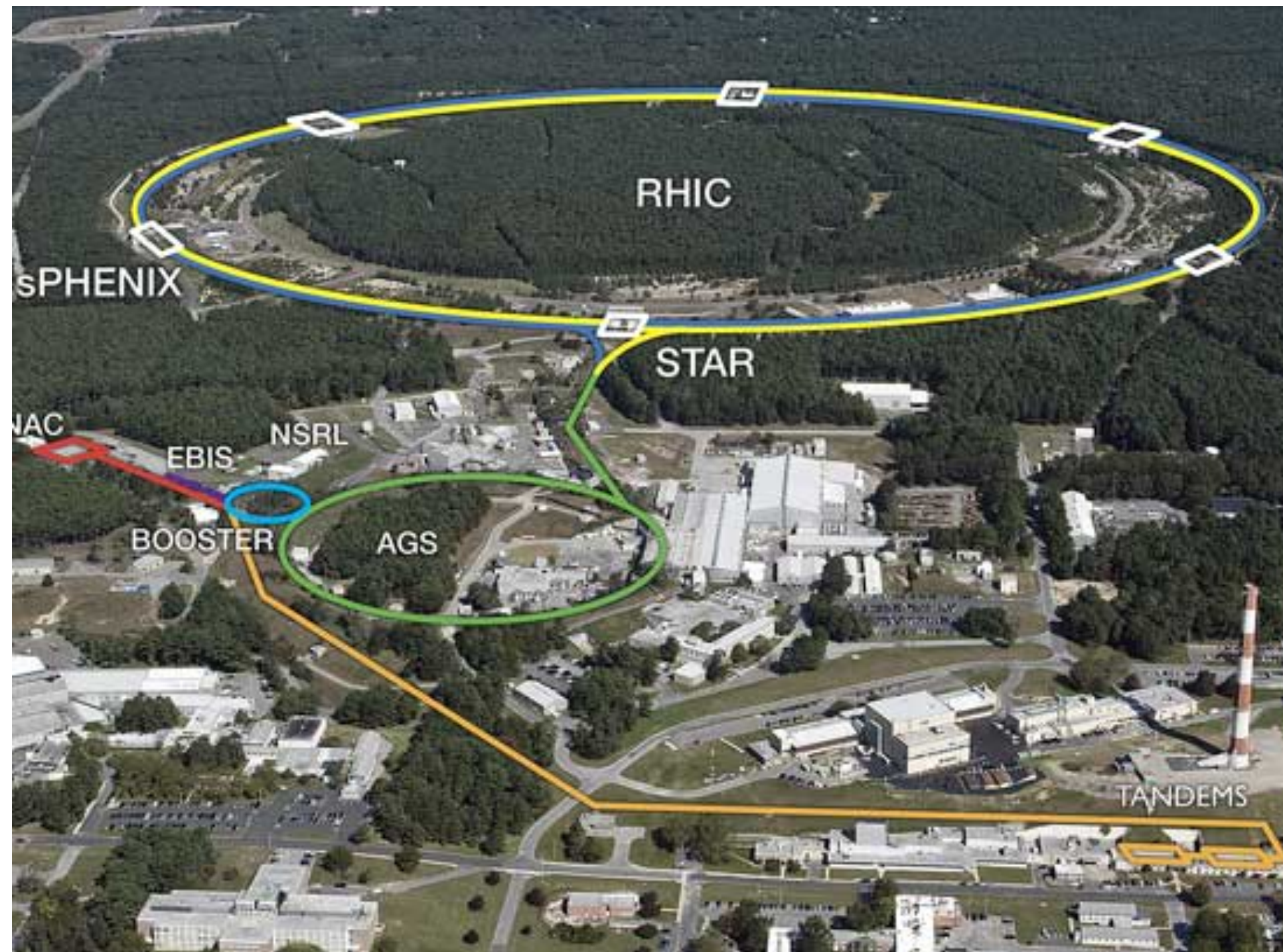
→ The performance worsens for Pb—Pb, due to the large UE.

*Quark and gluon discrimination is a difficult and ongoing effort in HIs!*

Future: Apply these methods to data in pp and Pb—Pb!



# ML @ THE RHIC ACCELERATOR COMPLEX



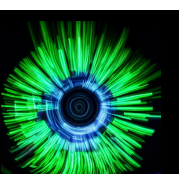
- **Boosted Decision Trees** to identify and predict magnet quenches from historical data.
- Combined with **Autoencoders** used to identify signs indicative of future quenches.

[JACoW IPAC2023 (2023) WEPA10]

- **Autoencoders** and PCA used for dimensionality reductions to see which parameters are useful for beam cooling. [JACoW NAPAC2022 (2022) 260-262]

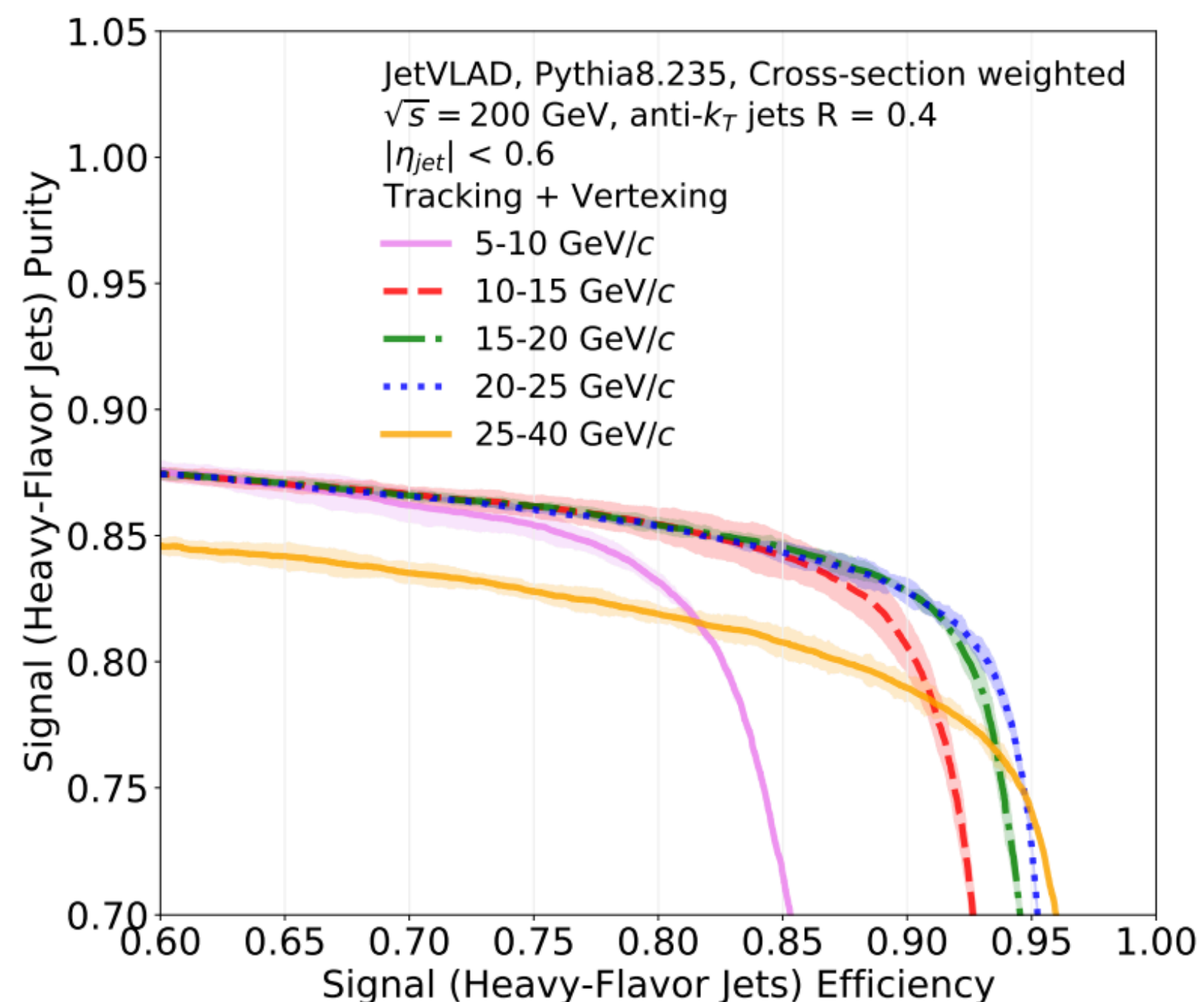
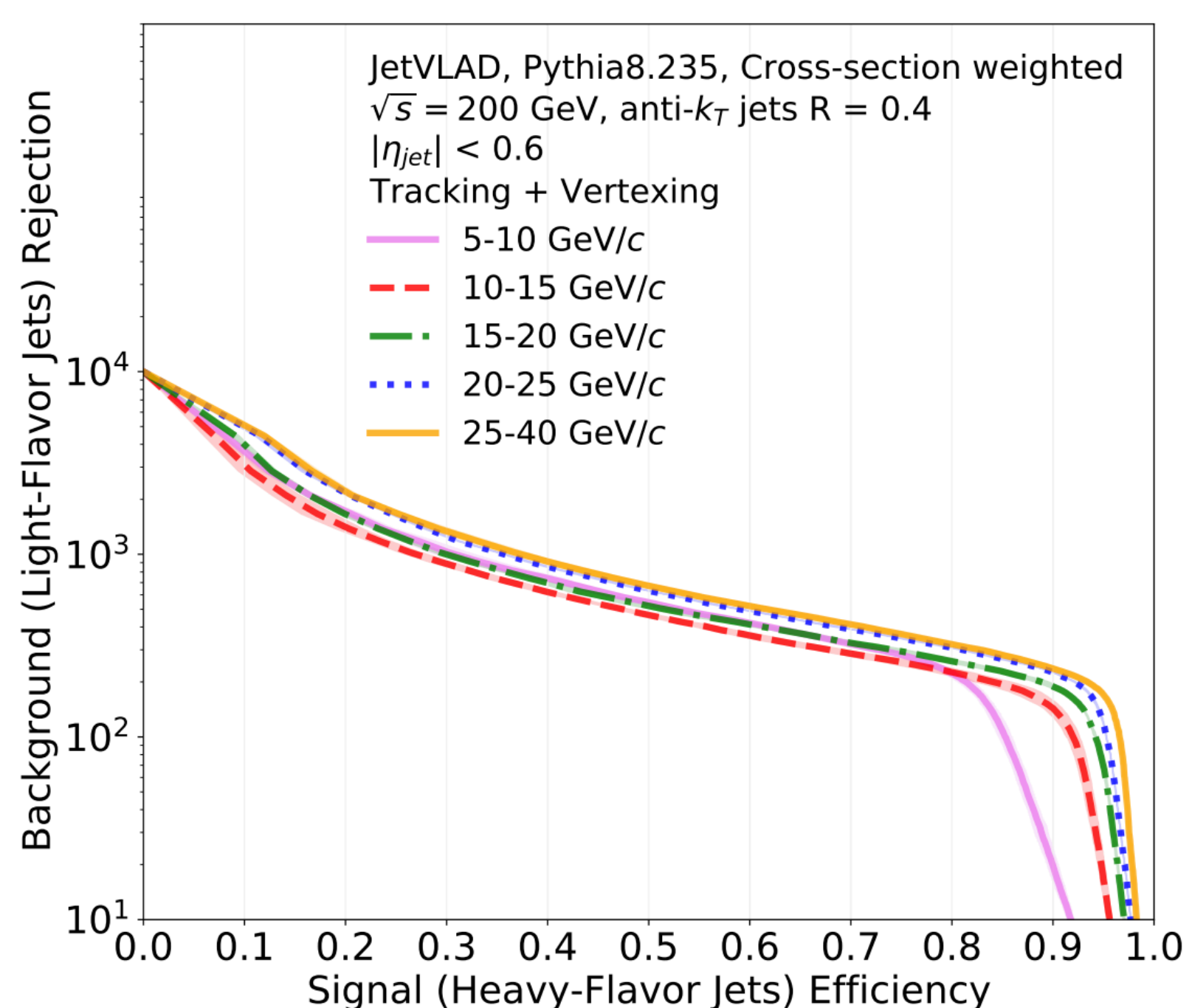
- Algorithms need to be robust to machine parameters.
  - Reinforcement or unsupervised learning useful.
- Need machine development time, can use simulations.

[JACoW ICALEPCS2023 (2023) FR2AO04]



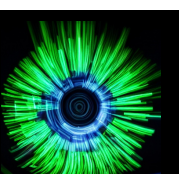
For input to the model treat the jet as a set of particles  $\mathcal{J} = \{(p_{T,i}, \eta_i, \phi_i, \dots)\}_{i=1}^n$

Model includes pooling layer that takes set of feature descriptors as an input and returns a fixed-length feature vector that characterizes each set.



- For higher  $p_T$  HF jets, background rejection increases, but purity decreases
- Fragmentation changes as function of  $p_T$  leads to an overlap of feature space

**This is a challenging problem! Especially in Au+Au!**





# ML AT THE EIC

Electron Ion Collider is a future facility being designed with future techniques in mind!

## Ongoing Activities w/ AI

- Detector design
- Simulation
- Reconstruction
- Particle Identification
- Analysis



See [\[AI4EIC\]](#) for a comprehensive overview

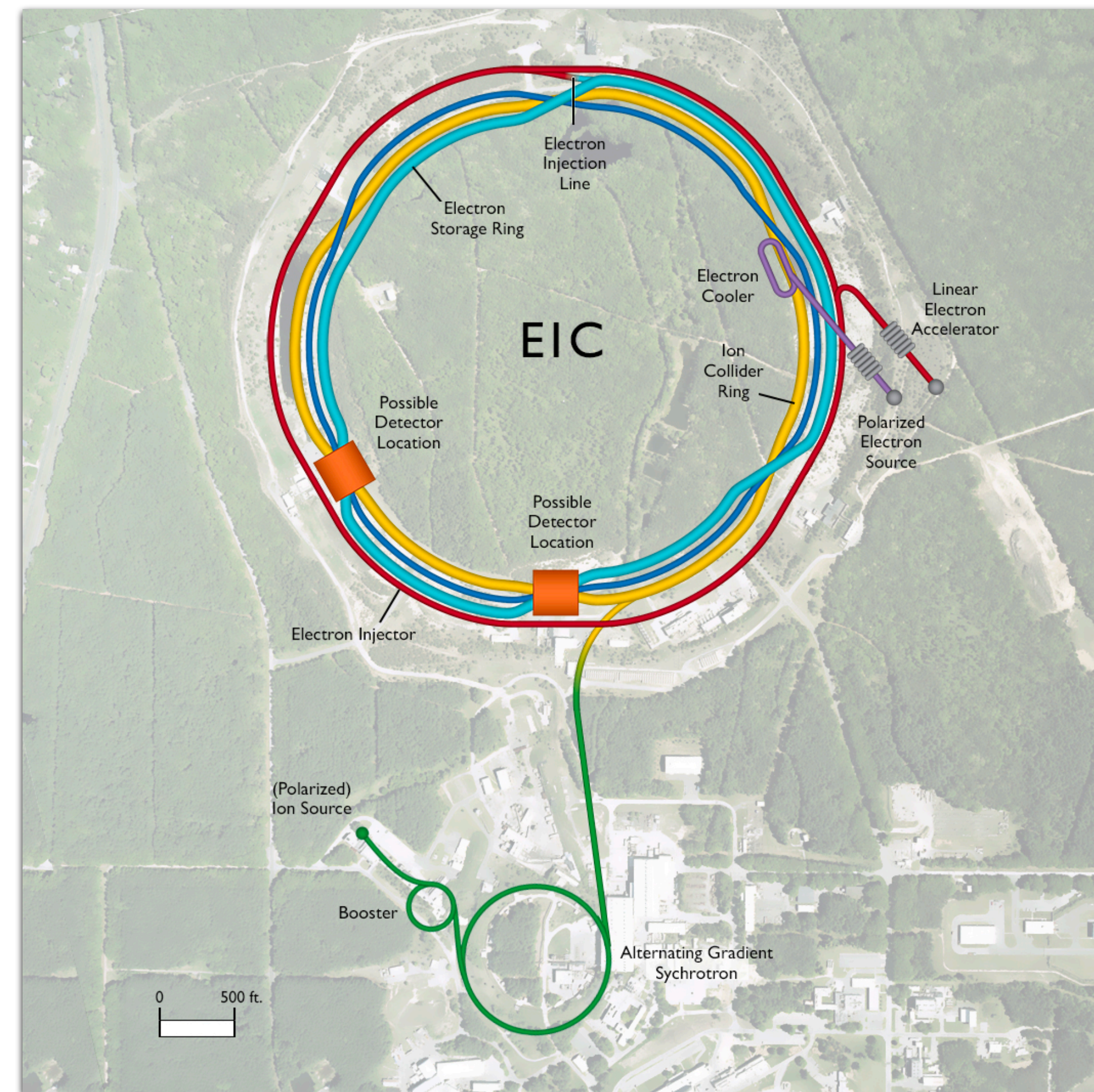
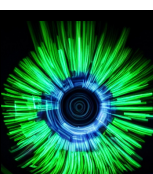


Image Credit: [\[Brookhaven National Lab\]](#)



# BIG DATA

Industry

Academia

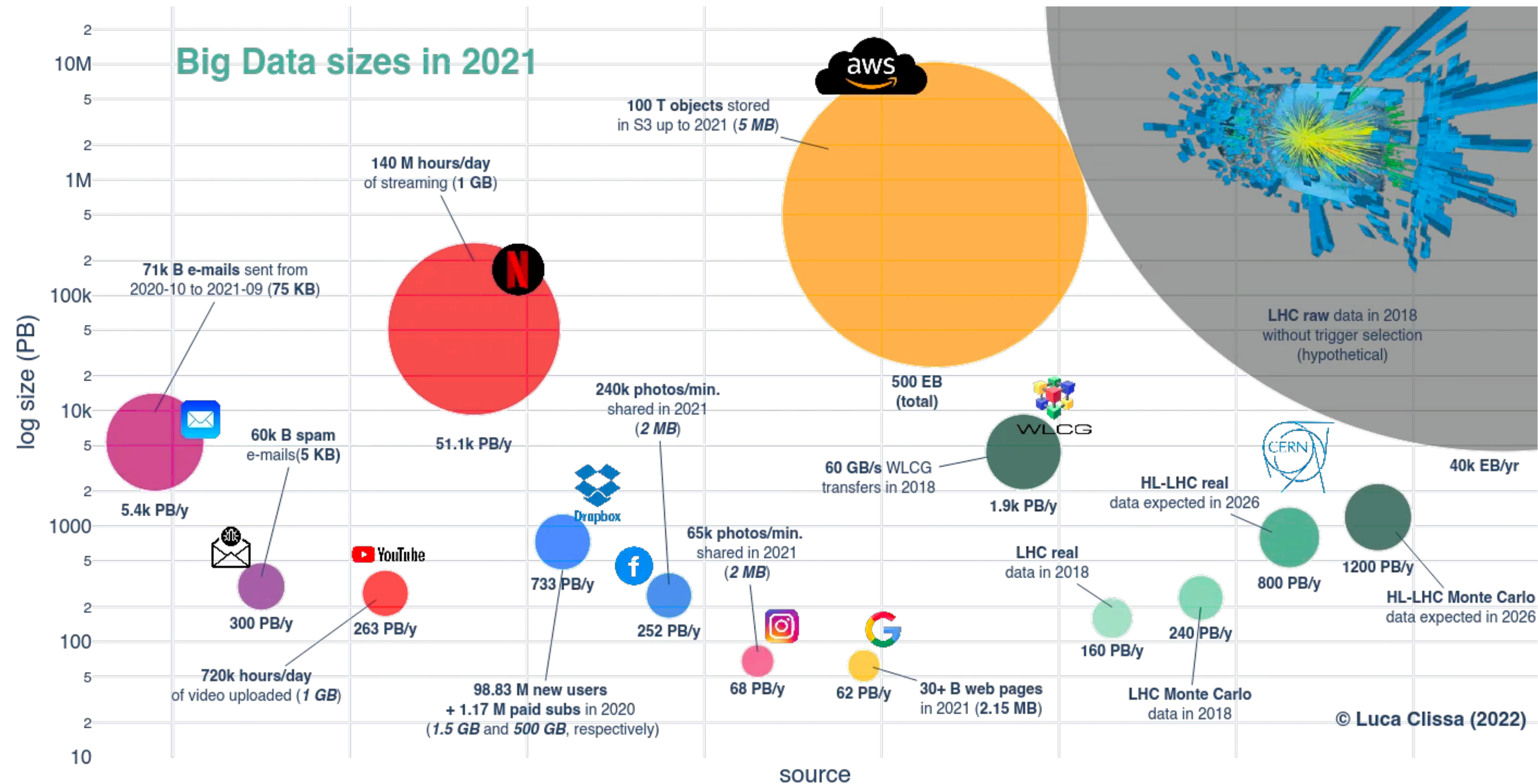
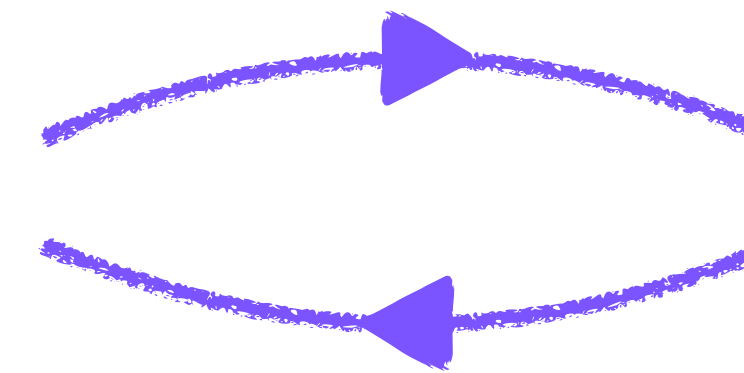
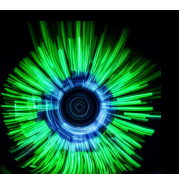


Image Credit: [Towards Data Science]

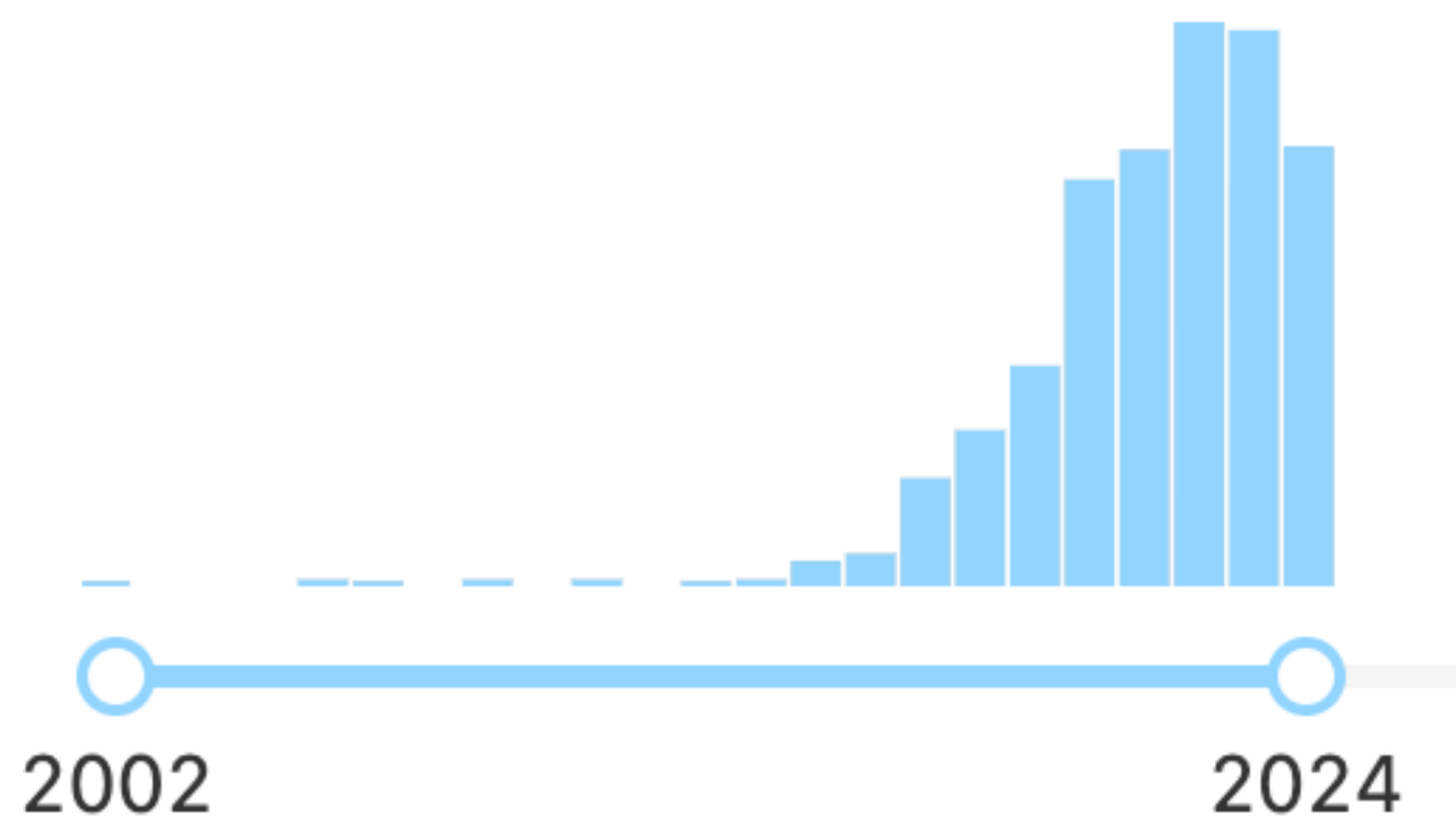
Data volumes comparable to medium-sized industry applications.



# ML ON THE RISE

1,594 total

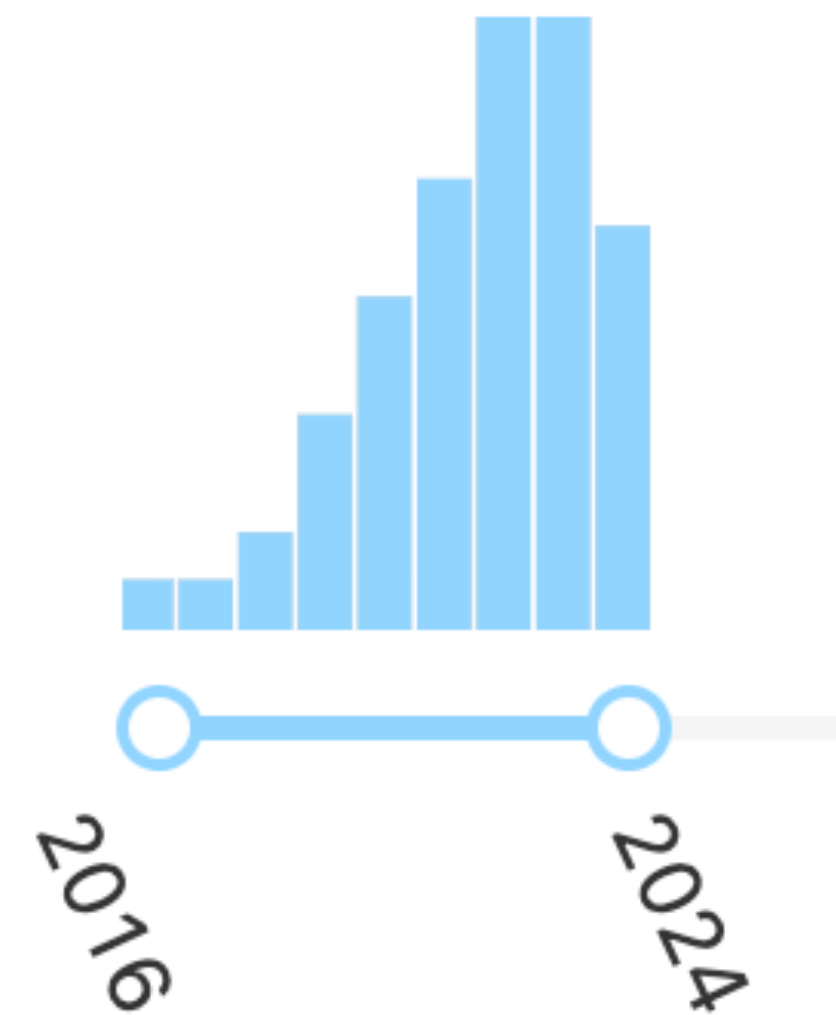
# of papers



161 total results!

# of papers

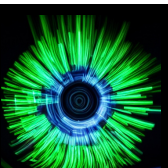
&



Inspire HEP search results for  
“machine learning HEP”

Inspire HEP search results for  
“machine learning heavy ion”

***ML is a rapidly growing field for HEP and heavy-ions!***



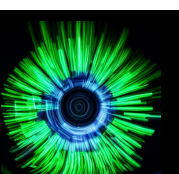
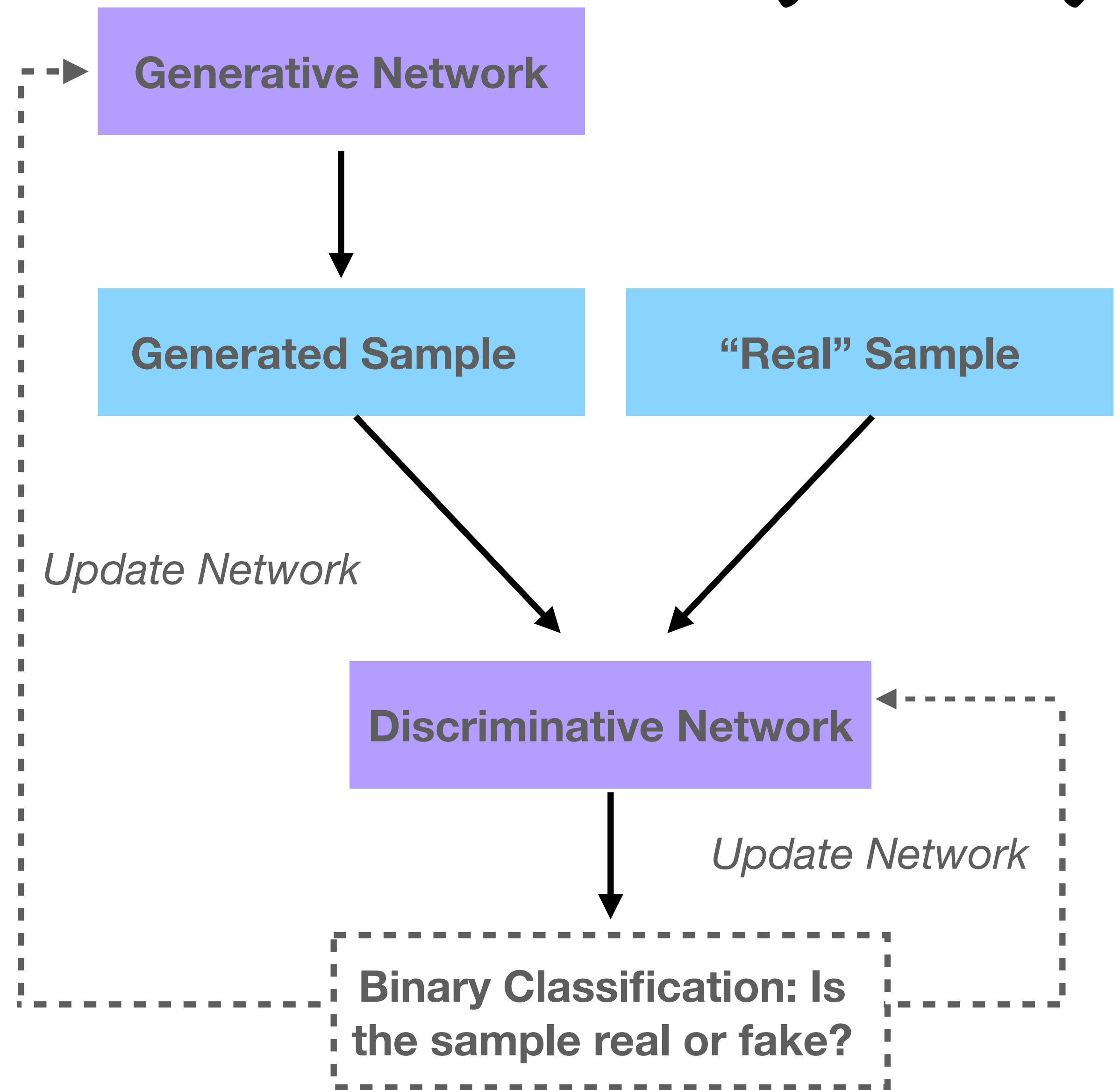
# GENERATIVE ADVERSARIAL NETWORKS (GANs)

Two networks compete with one another in a game.

The generative network seeks to fool the discriminative network.

The discriminative network seeks to find the real sample from the generated samples.

Indirect training → generative network never sees the true distribution!

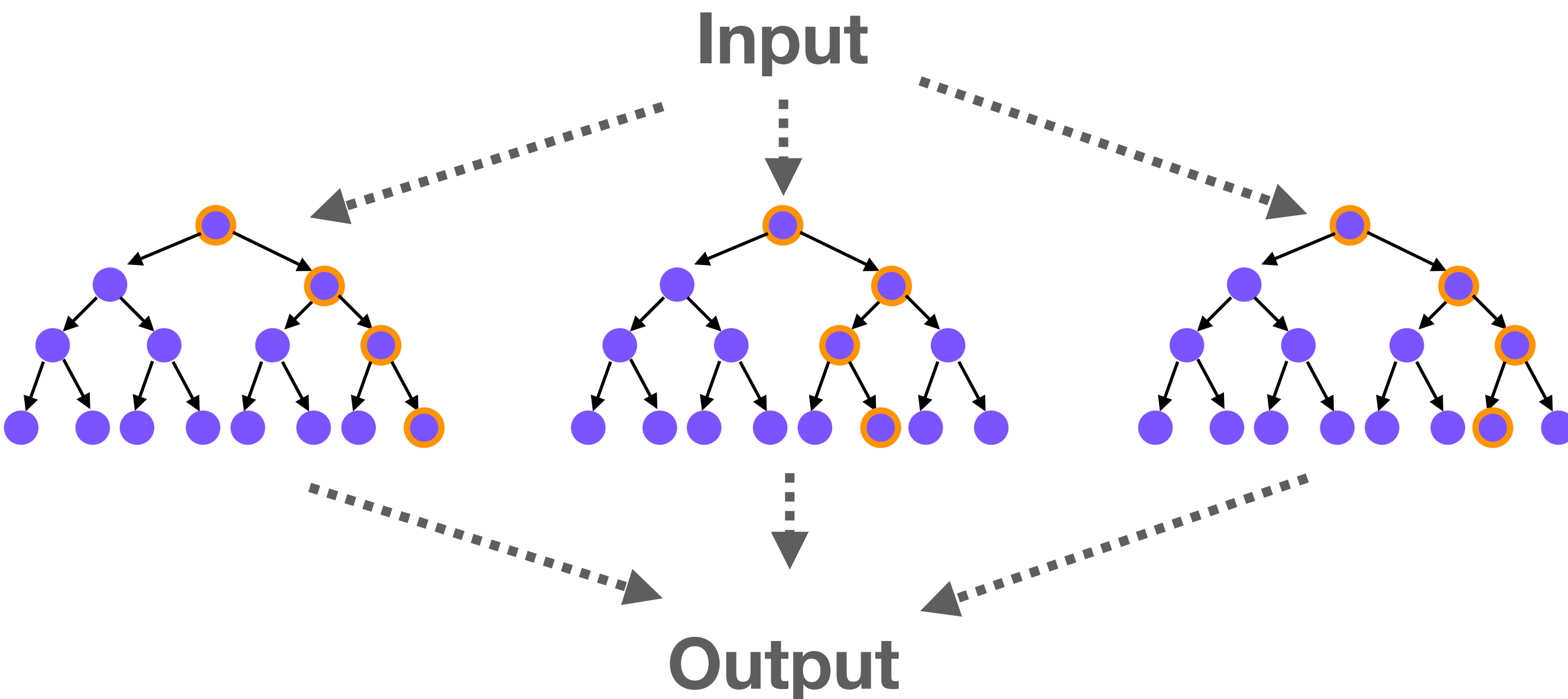
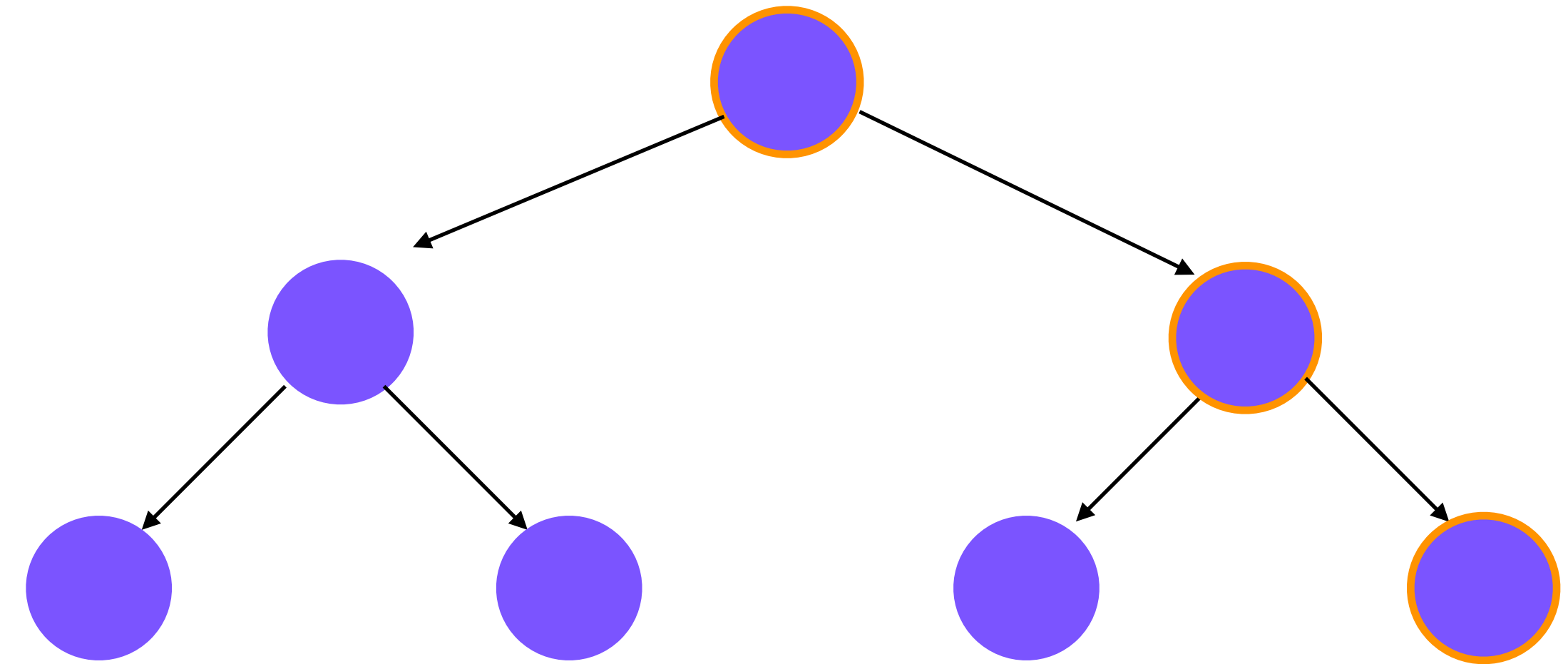


# INTRO TO RANDOM FOREST

Random forests are composed of decision trees.

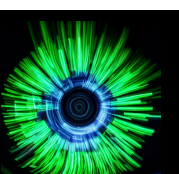
Decision trees are a set of rules organized in a tree structure.

Each node is a rule which subdivides the dataset into two or more parts (think 20 questions).



Output of the random forest is a combination of the output of each of the decision trees.

In training, the algorithm sets up the rules of each decision tree.



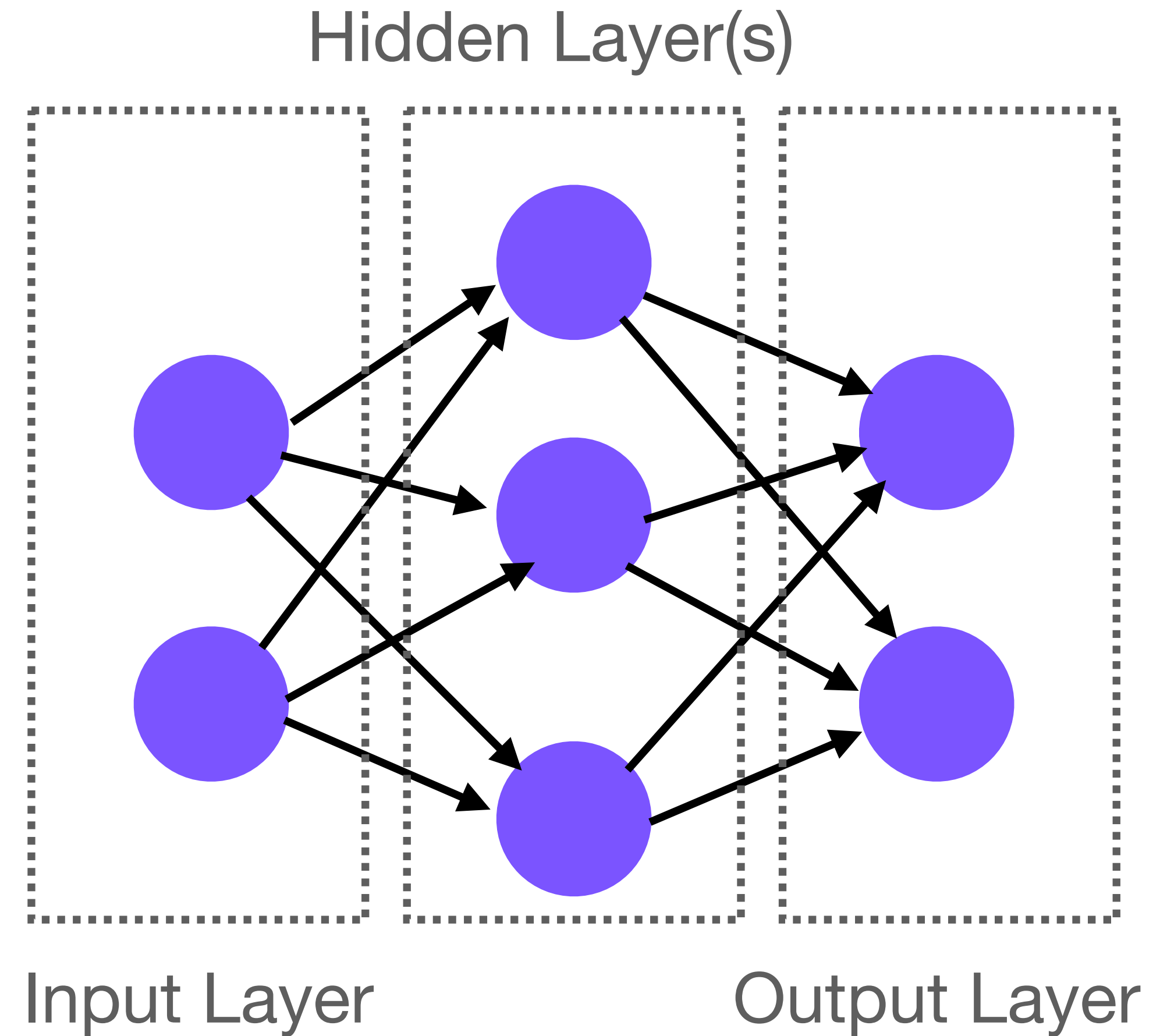
# NEURAL NETWORKS

Flow of information happens between **nodes**.

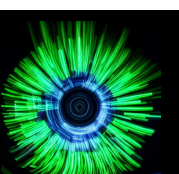
A weight is associated with each input to a given node.

The output of each node is a function of the weighted inputs. The output of a node  $j$ , is generally written something like

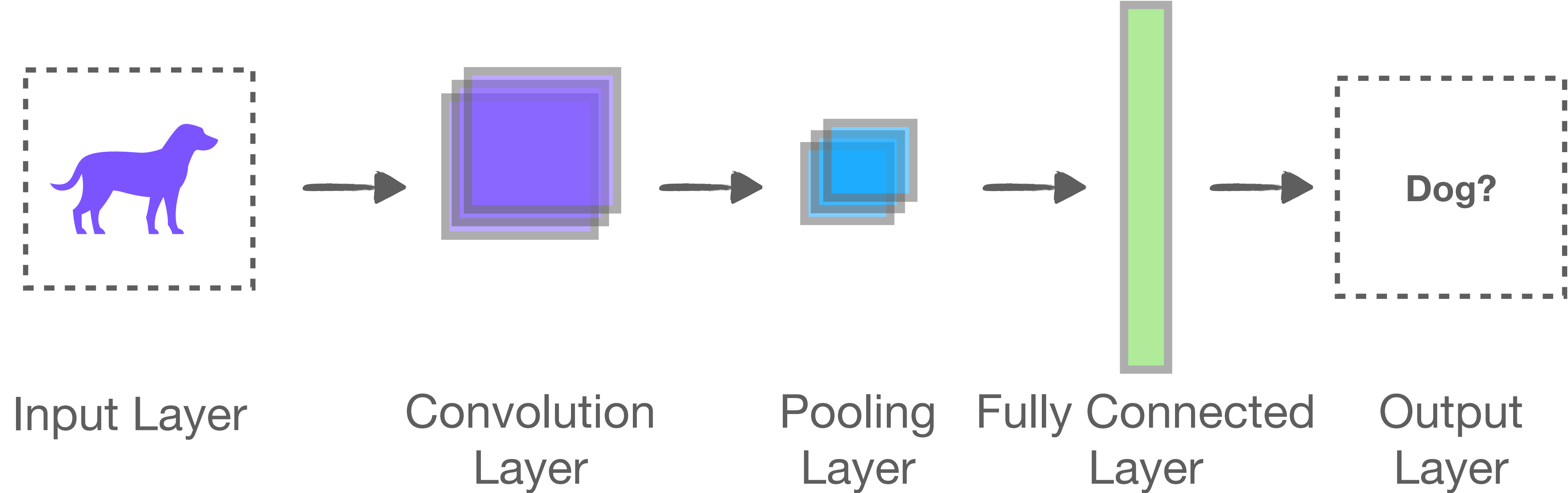
$$O_j = \sum_{i=0}^{N-1} w_{ij} O_i$$



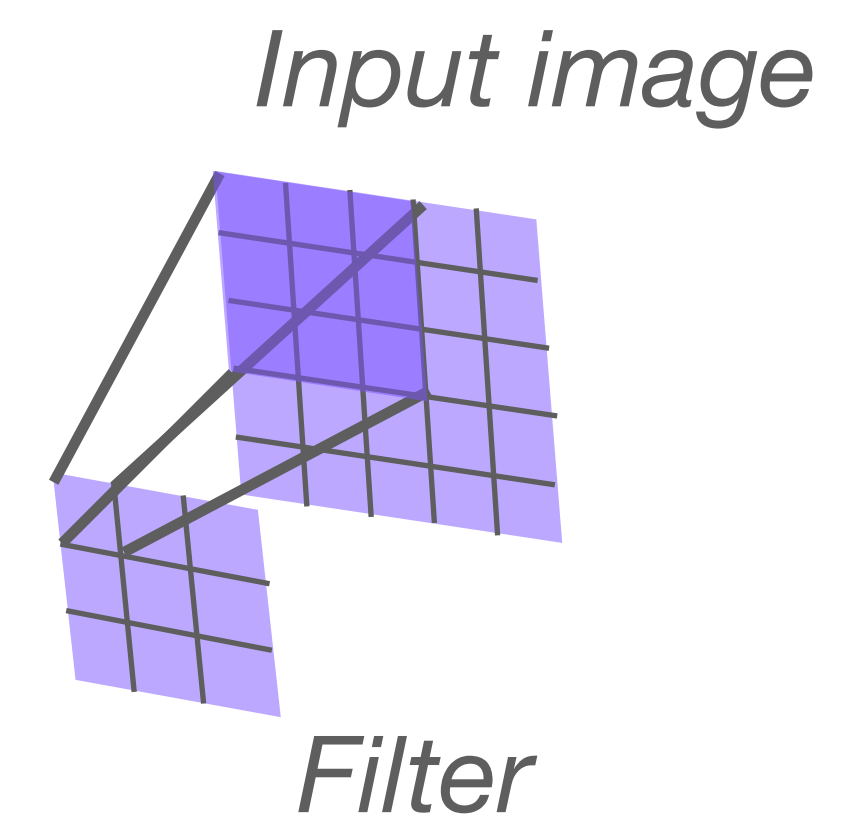
In training we seek to learn the set of weights which minimize the total error of the network.



# CONVOLUTIONAL NEURAL NETWORKS (CNNs)



→ Key component of a CNN is the **convolution layer**, which (with the help of a filter) will determine if a feature/pattern is present.



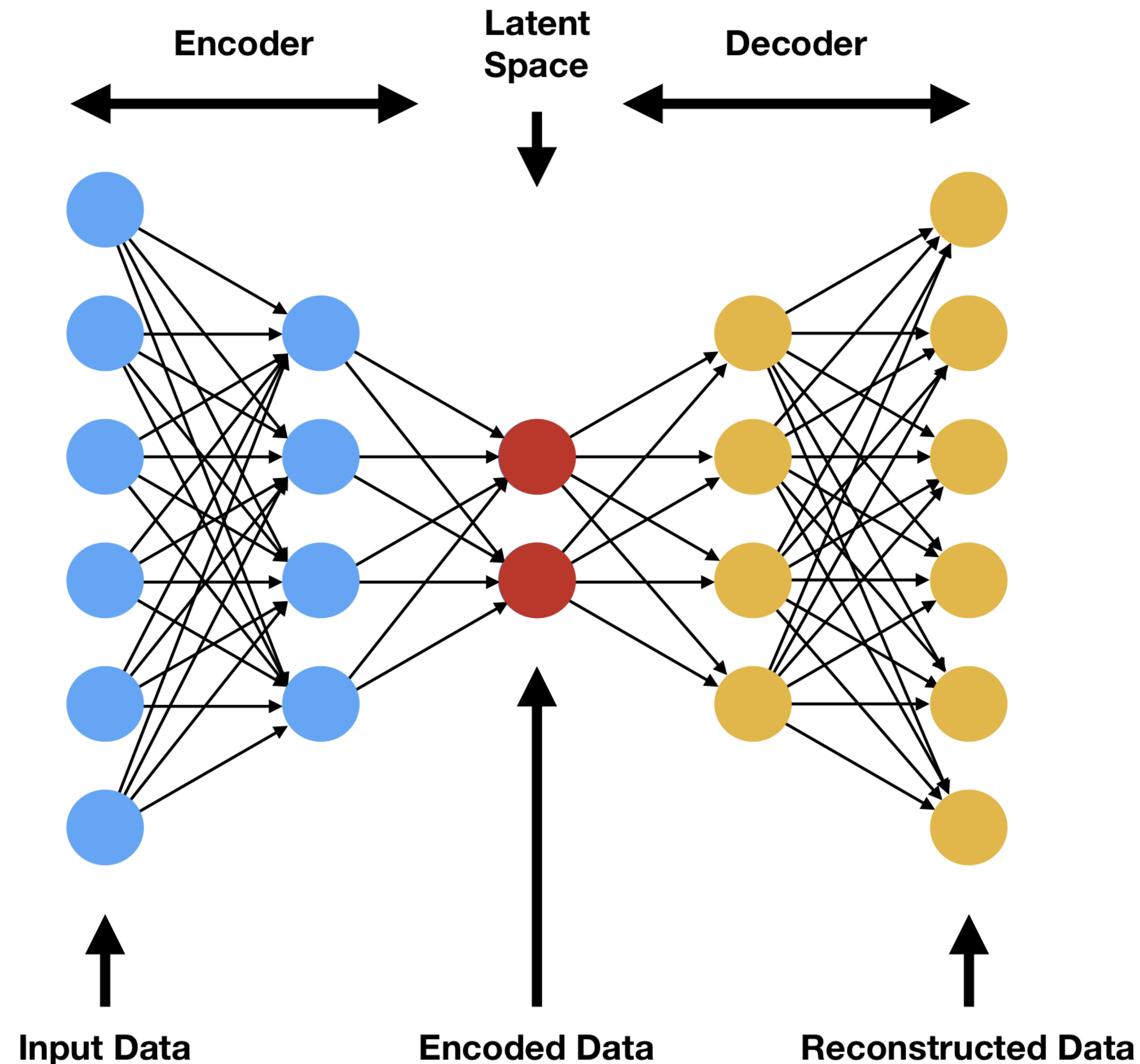
# AUTO-ENCODERS

**Simple task:** NN architecture trained to copy inputs to outputs!

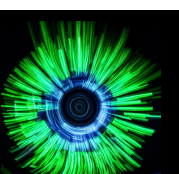
Encoder takes the input and dramatically reduces its complexity via a NN.

Decoder takes the encoded data and reconstructs outputs like the data.

*Does not require labeled data as input!*



<https://www.compthree.com/blog/autoencoder/>

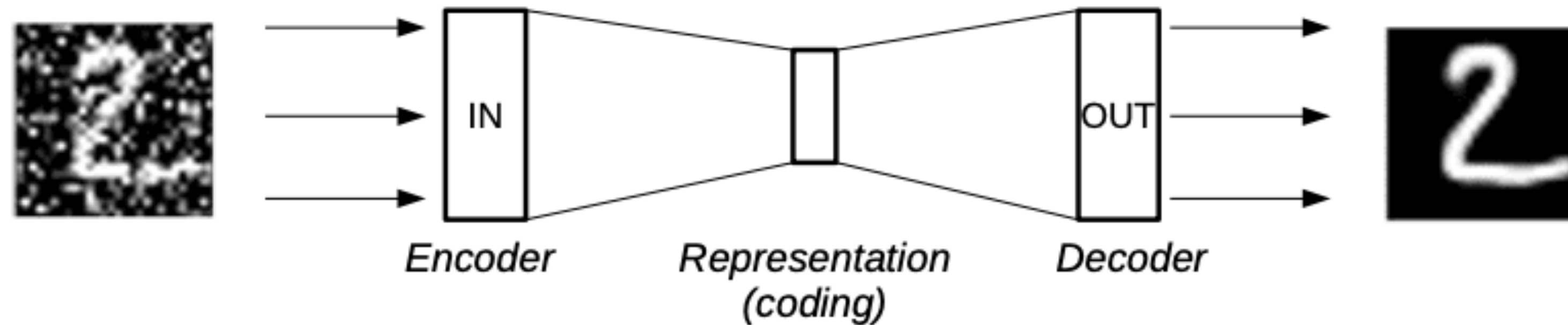




# USES OF AUTO-ENCODERS

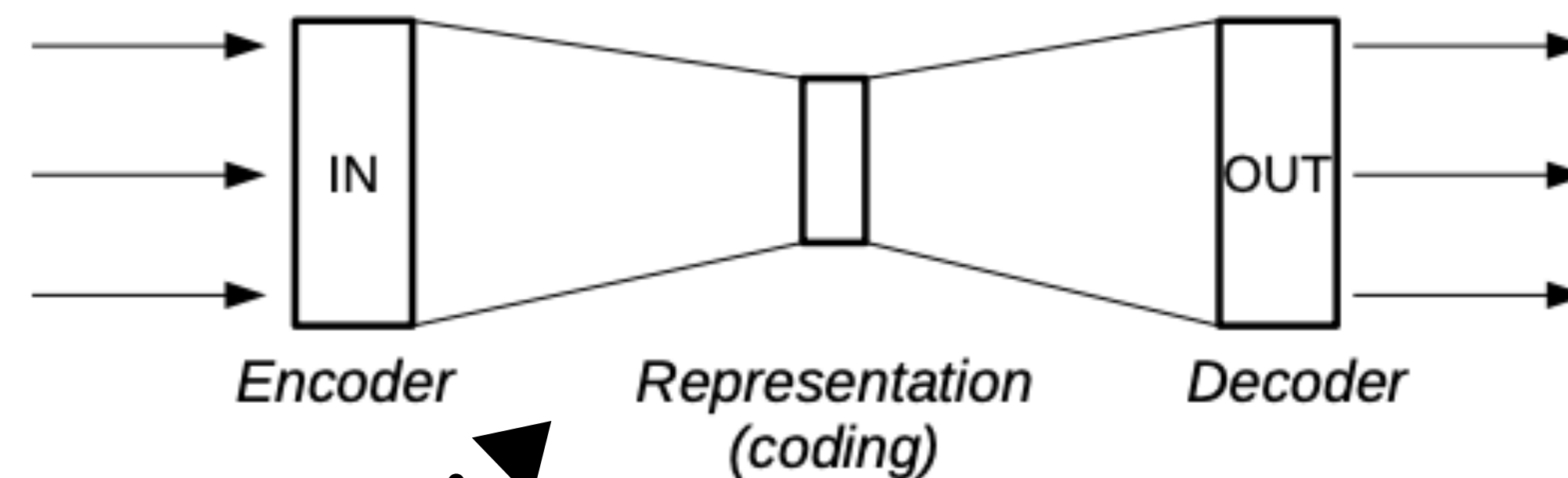
Used to learn efficient representations of some input data.

① De-noising inputs

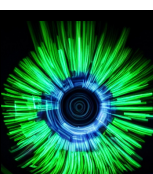


② Unsupervised learning

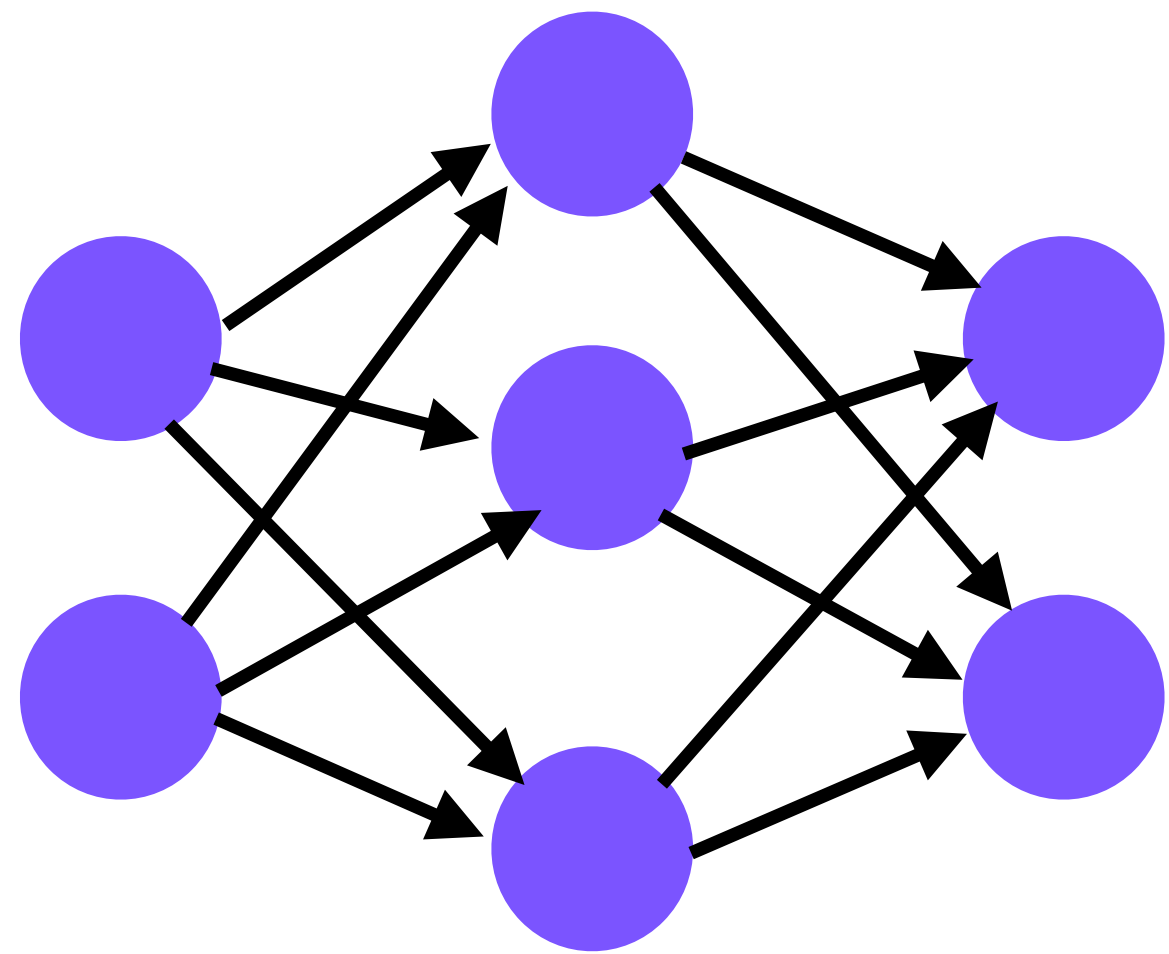
*Sort items into classes here*



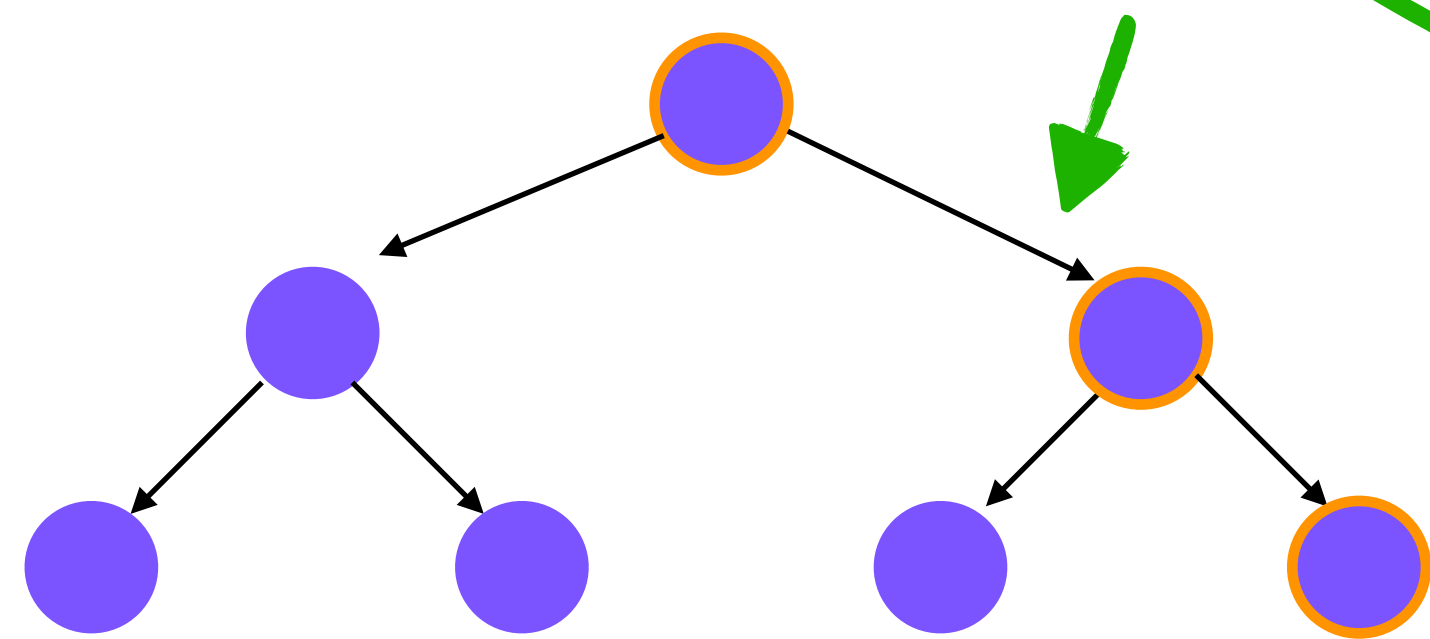
③ Anomaly detections: If you fail to reconstruct data in the decoding step you have an anomaly!



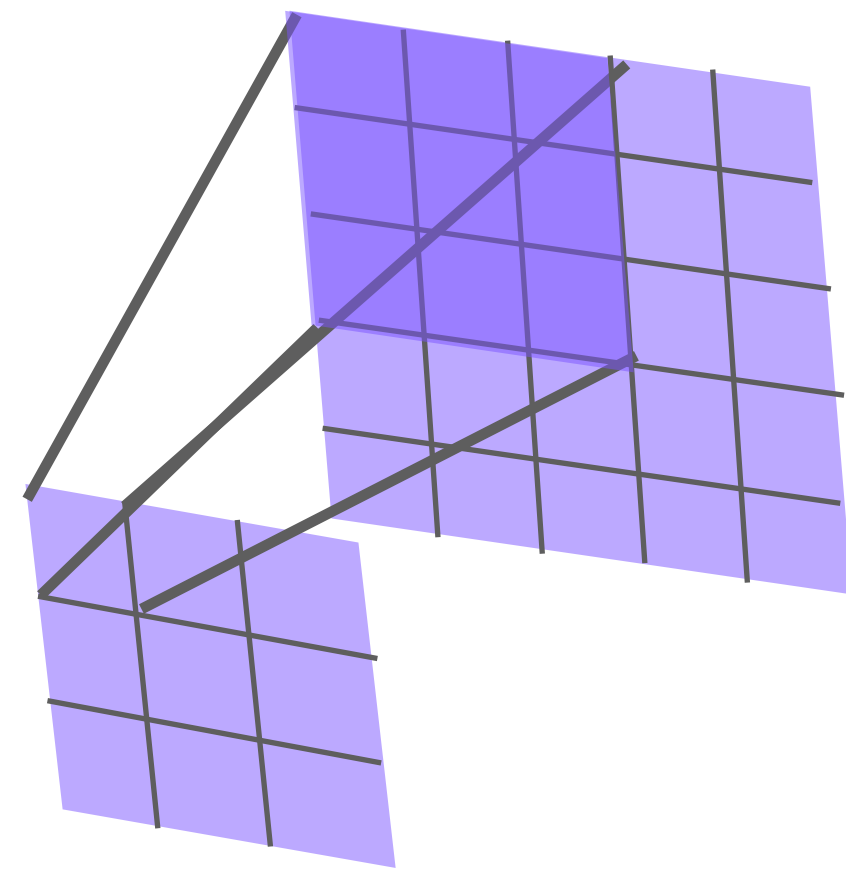
# DIFFERENT ALGORITHMS FOR DIFFERENT PROBLEMS!



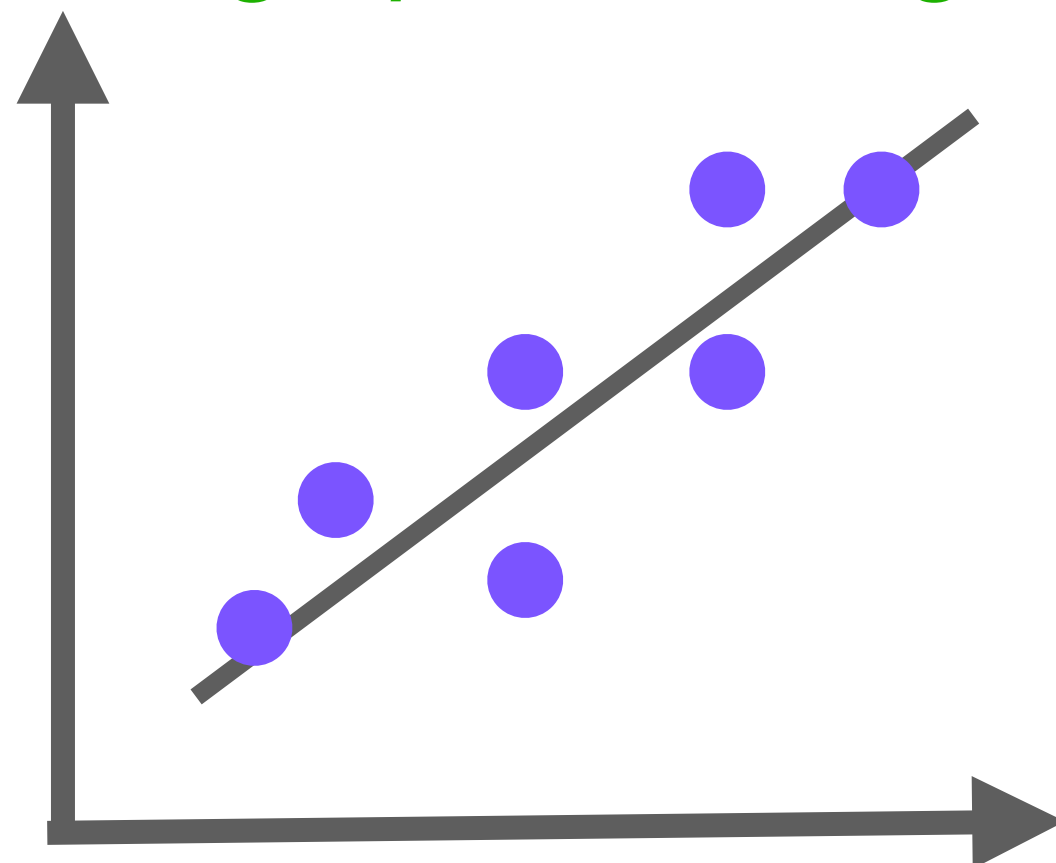
(Shallow or Deep) Neural Networks → *Great for making predictions!*



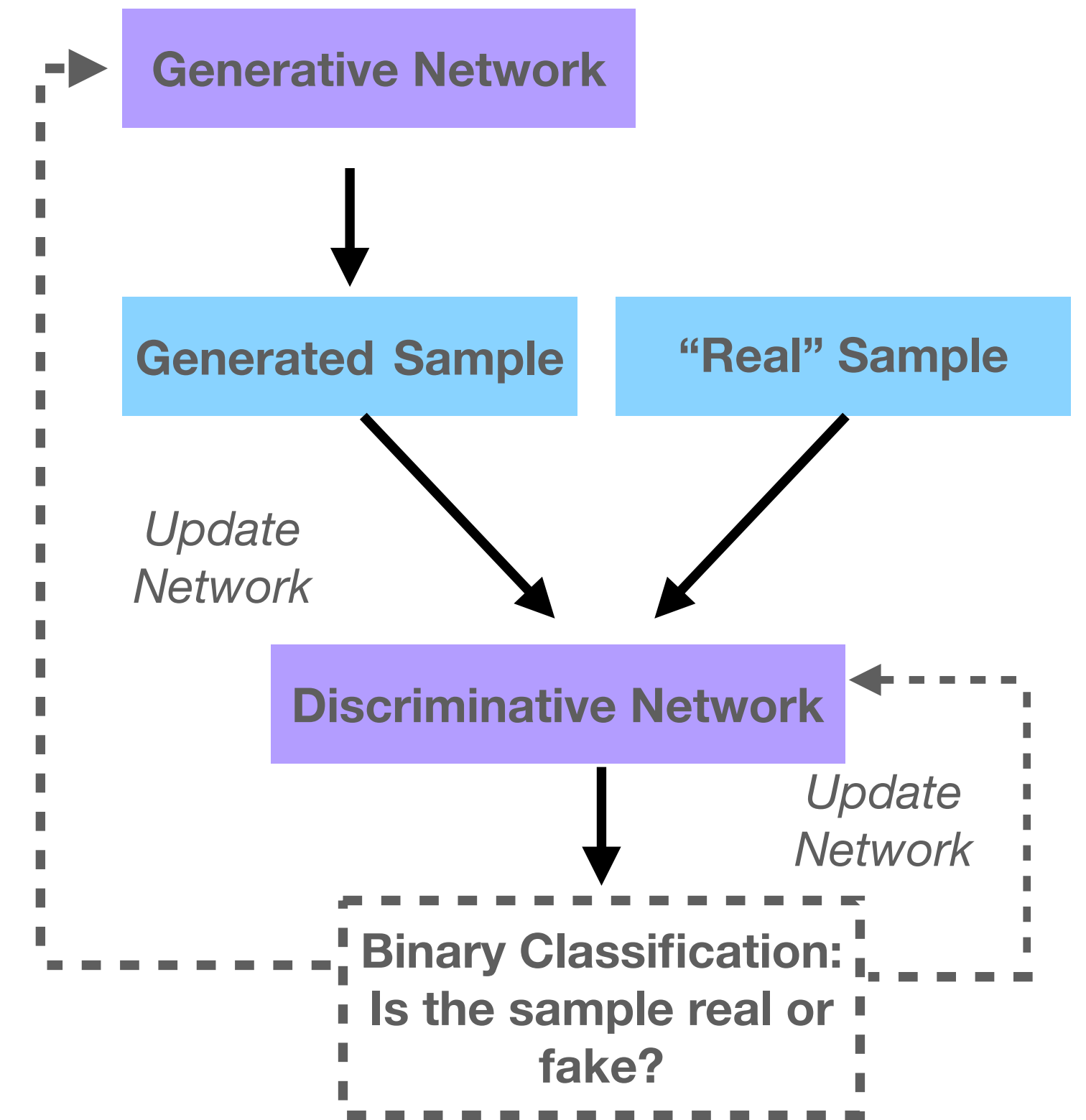
Random Forest (Decision Trees)



Convolutional Neural Networks (CNNs) → *Great for image processing!*



Linear Regression



Generative Adversarial Networks (GANs) → *Powerful tool for generating samples!*