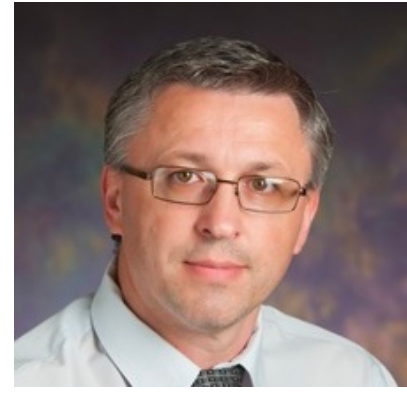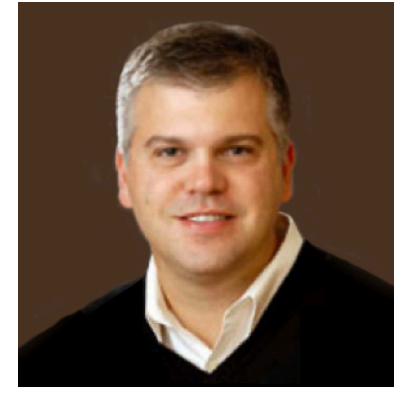Eliu Huerta
PI, ANL

Daniel S. Katz
Co-PI, NCSA

Volodymyr Kindratenko
Co-PI, NCSA

Mark Neubauer
Co-PI, UIUC

Zhizhen Zhao
Co-PI, UIUC

Priscilla Cushman
Co-PI, UMN

Andrew Furmanski
Co-PI, UMN

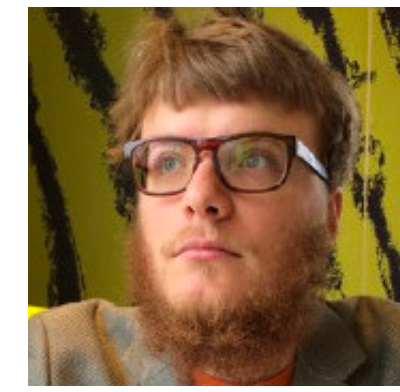Vuk Mandic
Co-PI, UMN

Roger Rusack
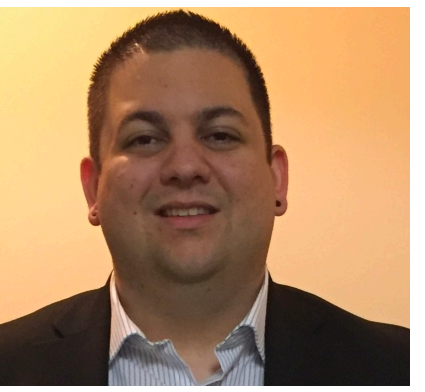Co-PI, UMN

Ju-Sun
Co-PI, UMN
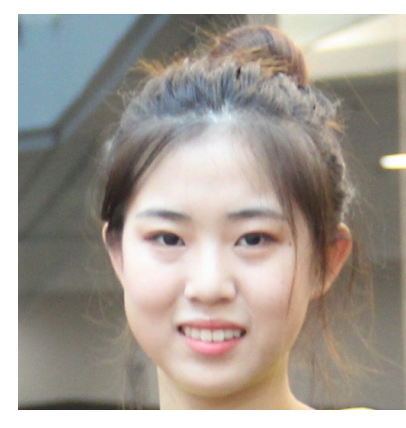
Phil Harris
Co-PI, MIT

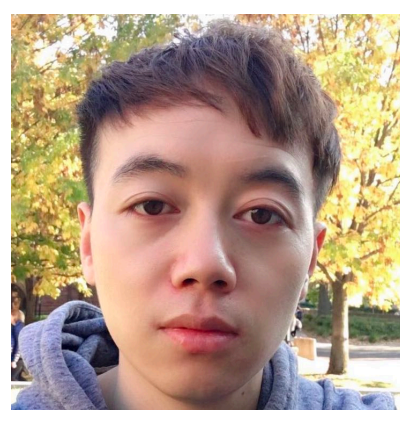Javier Duarte
Co-PI, UCSD

Andrew Evans
Postdoc, UMN

Daniel Diaz
Postdoc, UCSD

Sangeon Park
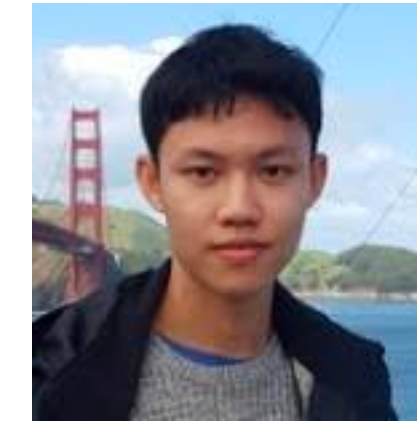PhD Student, MIT

Yifan Chen
MS Student, UIUC

Taihui Li
PhD Student, UMN

Raghav Kansal
PhD Student, UCSD

Farouk Mokhtar
PhD Student, UCSD

Steven Tsan
Undergrad, UCSD

U.S. DEPARTMENT OF ENERGY | Office of Science

‣ DOE ASCR-funded collaboration (3-year project)

  ‣ To advance our understanding of the relationship between our data and AI models by empowering scientists to explore both through the development of frameworks adhering to the principles of **f**indability, **a**ccessibility, **i**nteroperability, and **r**eusability (FAIR)

  ‣ Using HEP as the science use-case

    ‣ Investigate FAIR ways to share AI models and data

    ‣ Create an environment where novel approaches to AI can be explored and applied to new data

    ‣ Enable new insights for applying AI techniques

‣ Collaborate with partners: CERN Open Data Portal, Zenodo, DLHub

‣ Operate within larger community: Australian Research Data Commons (ARDC), Research Data Alliance (RDA)

  ‣ Note: BoF today on Fair for ML

▸ Motivation

▸ FAIR Principles

▸ FAIR4HEP Projects

    ▸ FAIR standards for data and AI in HEP

    ▸ Develop example FAIR datasets and AI models:

        ▸ H(bb) jet tagging

        ▸ Particle-flow reconstruction

        ▸ Among others

▸ Vision & Outlook

# MOTIVATION 1: REFERENCE DATA SETS FOR HEP

▸ Engage ML community for interesting, realistic tasks in experimental HEP

    ▸ As ImageNet accelerated advances in computer vision, do the same for HEP

▸ Engage ML community for interesting, realistic tasks in experimental HEP

  ▸ As ImageNet accelerated advances in computer vision, do the same for HEP

▸ Calls at many workshops for more public HEP data sets with real detector simulation for ML applications

  ▸ Example: dataset for top tagging based on Pythia+Delphes

  ▸ Example: dataset for tracking based on ACTS (kaggle TrackML challenge)

  ▸ Example: dataset for H(bb) tagging based on CMS open simulation

  ▸ Example: dataset for particle-flow based on Pythia+Delphes

  ▸ Example: dataset for ECAL crystal calibration (in preparation)

▸ Allow AI models developed for one experiment to be (re-)trained and (re-)used easily in another experiment

  ▸ Example: ATLAS recently studied GravNet developed by CMS collaborators [https://cds.cern.ch/record/2753414] for physics object localization using point cloud segmentation

▶ Easier to build upon existing work (e.g. through transfer learning)



Left: Xception pre-trained on ImageNet applied to galaxies

Right: After fine-tuning on galaxy data, two galaxy clusters can be clearly identified

▸ Easier to build upon existing work (e.g. through transfer learning)



Left: Xception pre-trained on ImageNet applied to galaxies

Right: After fine-tuning on galaxy data, two galaxy clusters can be clearly identified

▸ Share work beyond HEP

the 3D imaging clustering

▸ AI models developed for HEP-specific tasks may be useful in other domains (e.g. LiDAR point cloud data)

FINDABLE

FINDABLE

▸ F1. (meta)data have **unique** and **persistent** identifier

F2. data are described with rich metadata

F3. metadata specify the data identifier

F4. (meta)data are registered or indexed in a searchable resource

# FAIR DATA PRINCIPLES

FINDABLE

ACCESSIBLE

▸ F1. (meta)data have **unique** and **persistent** identifier

F2. data are described with rich metadata

F3. metadata specify the data identifier

F4. (meta)data are registered or indexed in a searchable resource

Image: book.fosteropenscience.eu

FINDABLE

ACCESSIBLE

▸ F1. (meta)data have **unique** and **persistent** identifier
F2. data are described with rich metadata
F3. metadata specify the data identifier
F4. (meta)data are registered or indexed in a searchable resource

▸ A1. (meta)data are retrievable using standardized protocol
A1.1 protocol is open, free, and universally implementable
A1.2 protocol allows for authentication and authorization
A2. metadata are accessible, even when the data is not

Image: book.fosteropenscience.eu

AH!

FINDABLE

HOW DO YOU OPEN A .XZQ FILE?

INTEROPERABLE

DOI 10.1007/8.7P7

ACCESSIBLE

▸ F1. (meta)data have **unique** and **persistent** identifier
F2. data are described with rich metadata
F3. metadata specify the data identifier
F4. (meta)data are registered or indexed in a searchable resource

▸ A1. (meta)data are retrievable using standardized protocol
A1.1 protocol is open, free, and universally implementable
A1.2 protocol allows for authentication and authorization
A2. metadata are accessible, even when the data is not

Image: book.fosteropenscience.eu

AH!

FINDABLE

HOW DO YOU OPEN A .XZQ FILE?

INTEROPERABLE

DOI 10.1007/8.7P7

ACCESSIBLE

▸ F1. (meta)data have **unique** and **persistent** identifier
  F2. data are described with rich metadata
  F3. metadata specify the data identifier
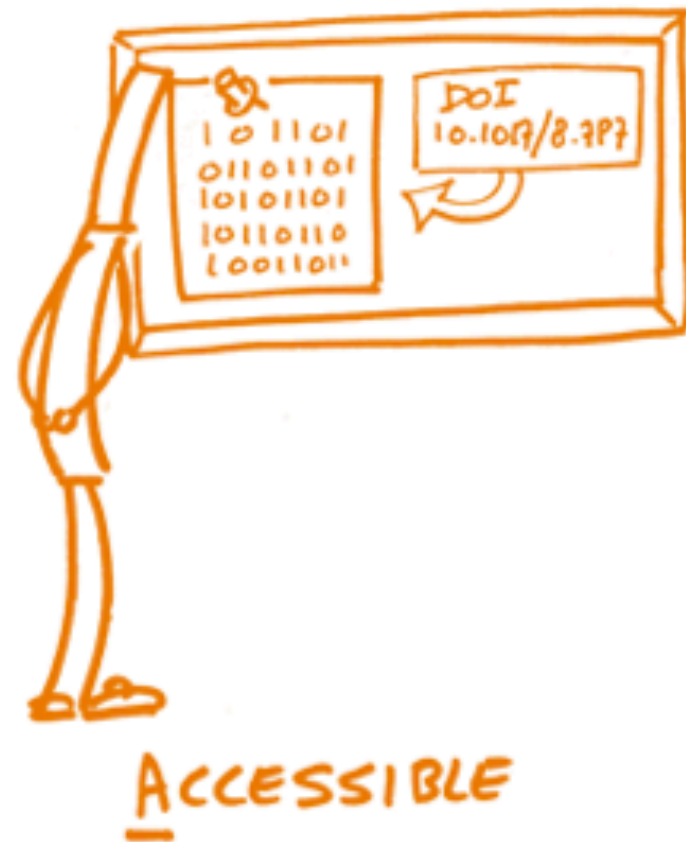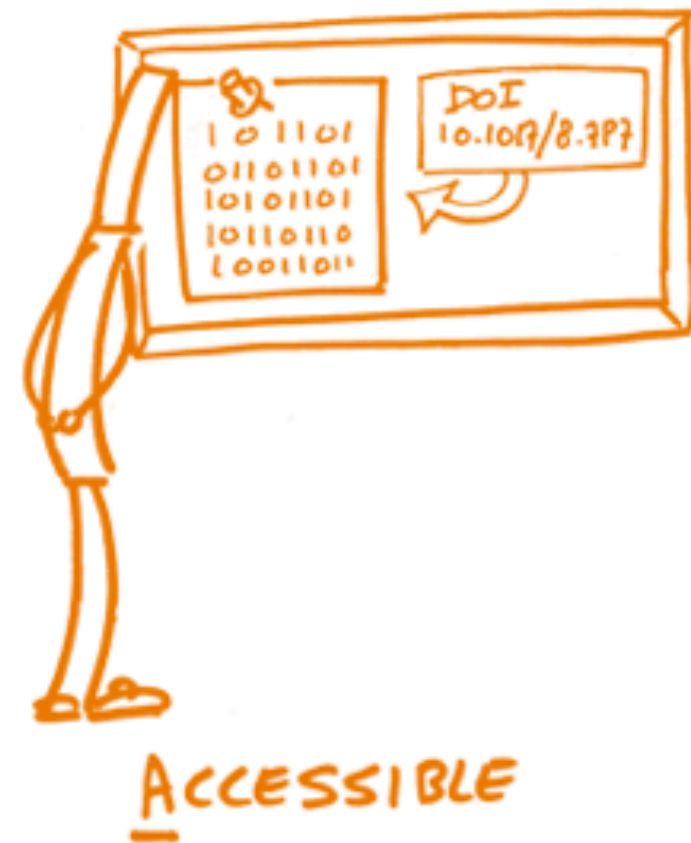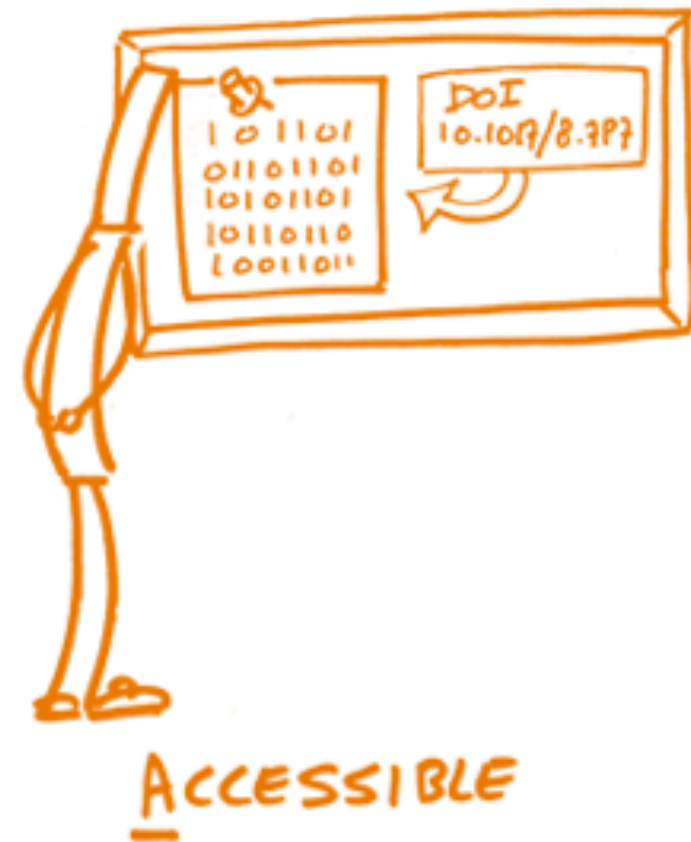  F4. (meta)data are registered or indexed in a searchable resource

▸ A1. (meta)data are retrievable using standardized protocol
  A1.1 protocol is open, free, and universally implementable
  A1.2 protocol allows for authentication and authorization
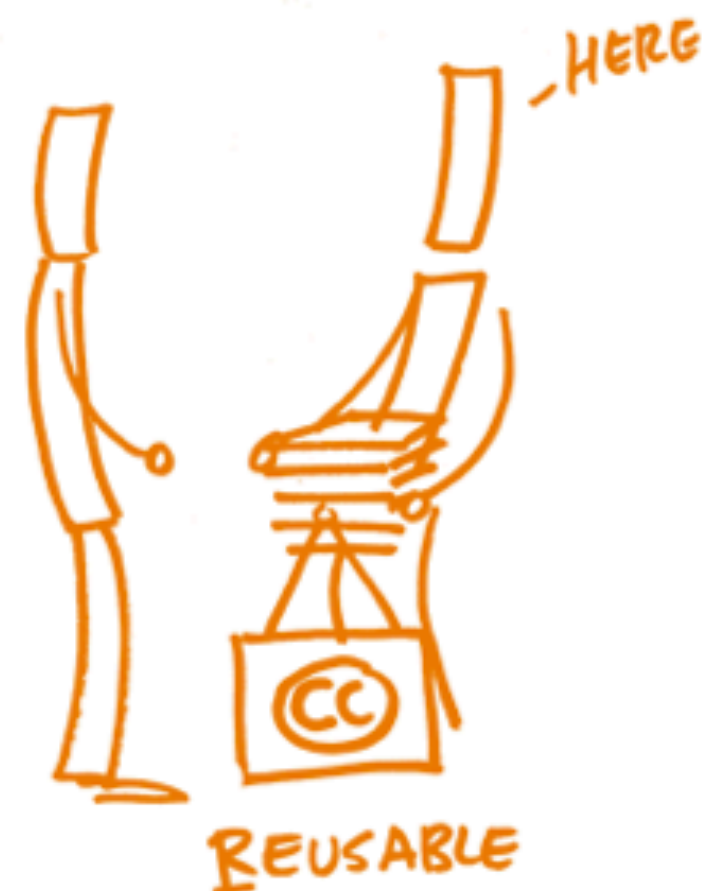  A2. metadata are accessible, even when the data is not

▸ I1. (meta)data use a formal, shared, and broadly applicable
  language for knowledge representation
  I2. (meta)data use **vocabularies** that follow FAIR principles
  I3. (meta)data include qualified references to other (meta)data

# FAIR DATA PRINCIPLES

FINDABLE

ACCESSIBLE

HOW DO YOU OPEN A .XZQ FILE?

INTEROPERABLE

HERE

REUSABLE

> F1. (meta)data have **unique** and **persistent** identifier
> F2. data are described with rich metadata
> F3. metadata specify the data identifier
> F4. (meta)data are registered or indexed in a searchable resource

> A1. (meta)data are retrievable using standardized protocol
> A1.1 protocol is open, free, and universally implementable
> A1.2 protocol allows for authentication and authorization
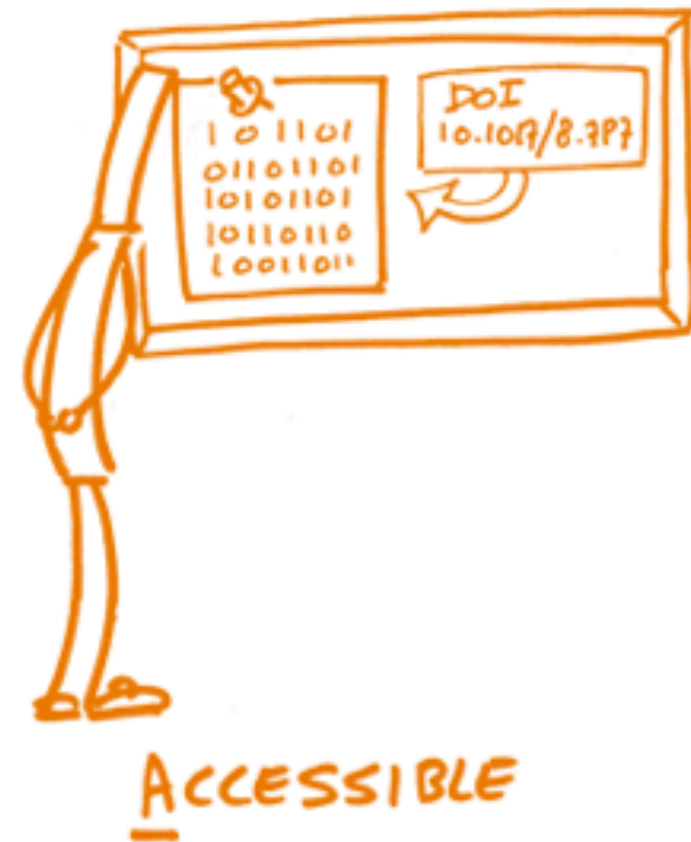> A2. metadata are accessible, even when the data is not

> I1. (meta)data use a formal, shared, and broadly applicable language for knowledge representation
> I2. (meta)data use **vocabularies** that follow FAIR principles
> I3. (meta)data include qualified references to other (meta)data

Image: book.fosteropenscience.eu

AH!

FINDABLE

DOI
10.1007/8.7P7

ACCESSIBLE

HOW DO YOU
OPEN A .XZQ FILE?

INTEROPERABLE

HERE

CC

REUSABLE

▸ F1. (meta)data have **unique** and **persistent** identifier
F2. data are described with rich metadata
F3. metadata specify the data identifier
F4. (meta)data are registered or indexed in a searchable resource

▸ A1. (meta)data are retrievable using standardized protocol
A1.1 protocol is open, free, and universally implementable
A1.2 protocol allows for authentication and authorization
A2. metadata are accessible, even when the data is not

▸ I1. (meta)data use a formal, shared, and broadly applicable
language for knowledge representation
I2. (meta)data use **vocabularies** that follow FAIR principles
I3. (meta)data include qualified references to other (meta)data

▸ R1. meta(data) have a plurality of accurate and relevant attributes
R1.1. (meta)data have clear and accessible data usage license
R1.2. (meta)data are associated with their provenance
R1.3. (meta)data meet domain-relevant community standards

Image: book.fosteropenscience.eu

# PROJECT: FAIR STANDARDS FOR DATA AND AI IN HEP

▸ Are HEP public datasets FAIR?

    ▸ Develop/refine FAIR checklist for HEP datasets

▸ What does FAIR mean for AI models (or [software generally](#))?

    ▸ Develop a standard protocol to follow to publish FAIR AI models

▸ How do users contribute their own FAIR data and AI models?

    ▸ Create mechanisms for users to contribute

▸ As users adopt these FAIR4HEP standards, it will be easier to

    ▸ Publish citable AI models (with credit, etc.)

    ▸ Retrain published AI models on new (also FAIR) datasets

    ▸ Extend published AI models for new tasks

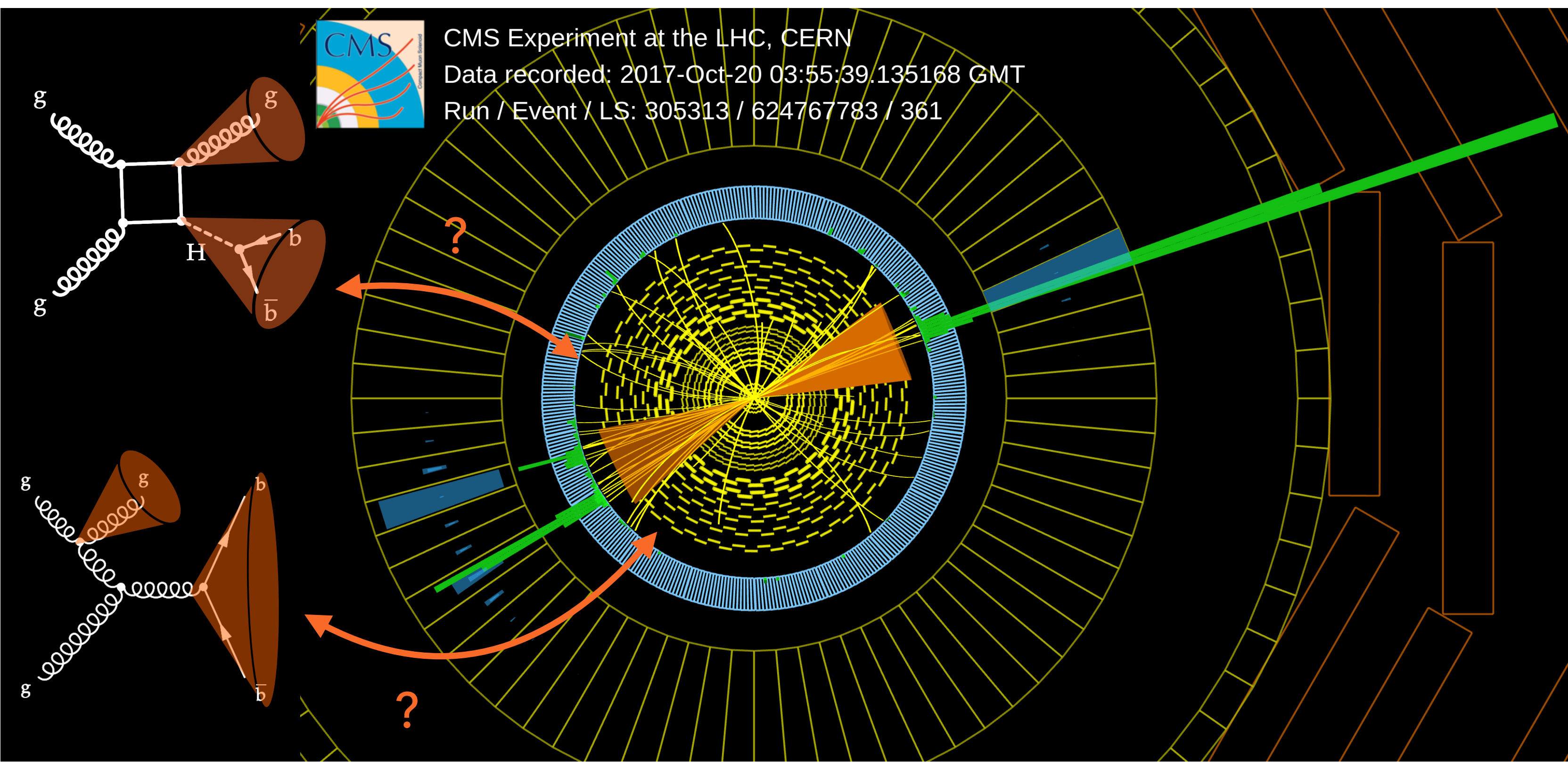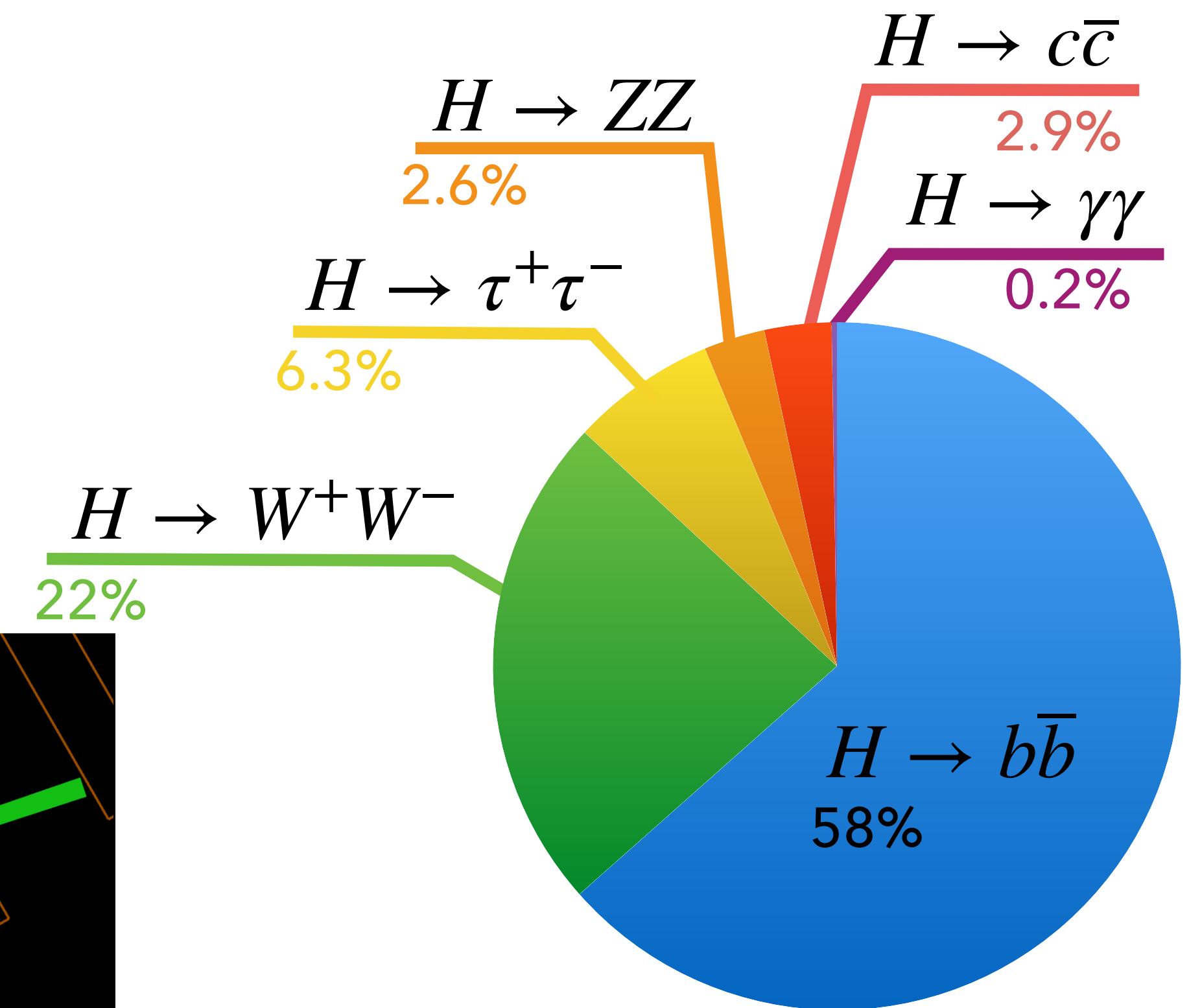    ▸ Explore the relationships between AI models and data

# PROJECT: DEVELOPING EXAMPLE FAIR DATASETS AND AI MODELS

‣ Advance important tasks in HEP with reference datasets and AI models to explore FAIRness criteria for both

   ‣ **H(bb) jet tagging**

   ‣ **Particle-flow reconstruction**

   ‣ ECAL crystal calibration

   ‣ Level-1 trigger jet reconstruction

   ‣ Charged particle tracking

   ‣ Among others

▸ H→bb is large, but more difficult to measure due to large backgrounds

$H \to c\bar{c}$
2.9%

$H \to ZZ$
2.6%

$H \to \gamma\gamma$
0.2%

$H \to \tau^+\tau^-$
6.3%

$H \to W^+W^-$
22%

$H \to b\bar{b}$
58%

CMS Experiment at the LHC, CERN
Data recorded: 2017-Oct-20 03:55:39.135168 GMT
Run / Event / LS: 305313 / 624767783 / 361

▸ Hosted on CERN Open Data Portal

  ▸ Collaborative effort between CERN IT-CDA and RCS-SIS groups, LHC and OPERA experiments

  ▸ Built with Invenio library management software

  ▸ Products (i.e. data, software, documentation, provenance) shared under open licenses and issued DOIs

  ▸ EOS data storage; access via XRootD, HTTP

▸ H(bb) dataset [10.7483/OPENDATA.CMS.JGJX.MS7Q]

  ▸ 182 files, 245 GB, 18 million total entries (jets)

    ▸ event features, e.g. MET, $\rho$ (average density)

    ▸ jet features, e.g. mass, $p_T$, N-subjettiness variables

    ▸ particle candidate features, e.g. $p_T$, $\eta$, $\phi$

    ▸ charged particle / track features, e.g. impact parameter
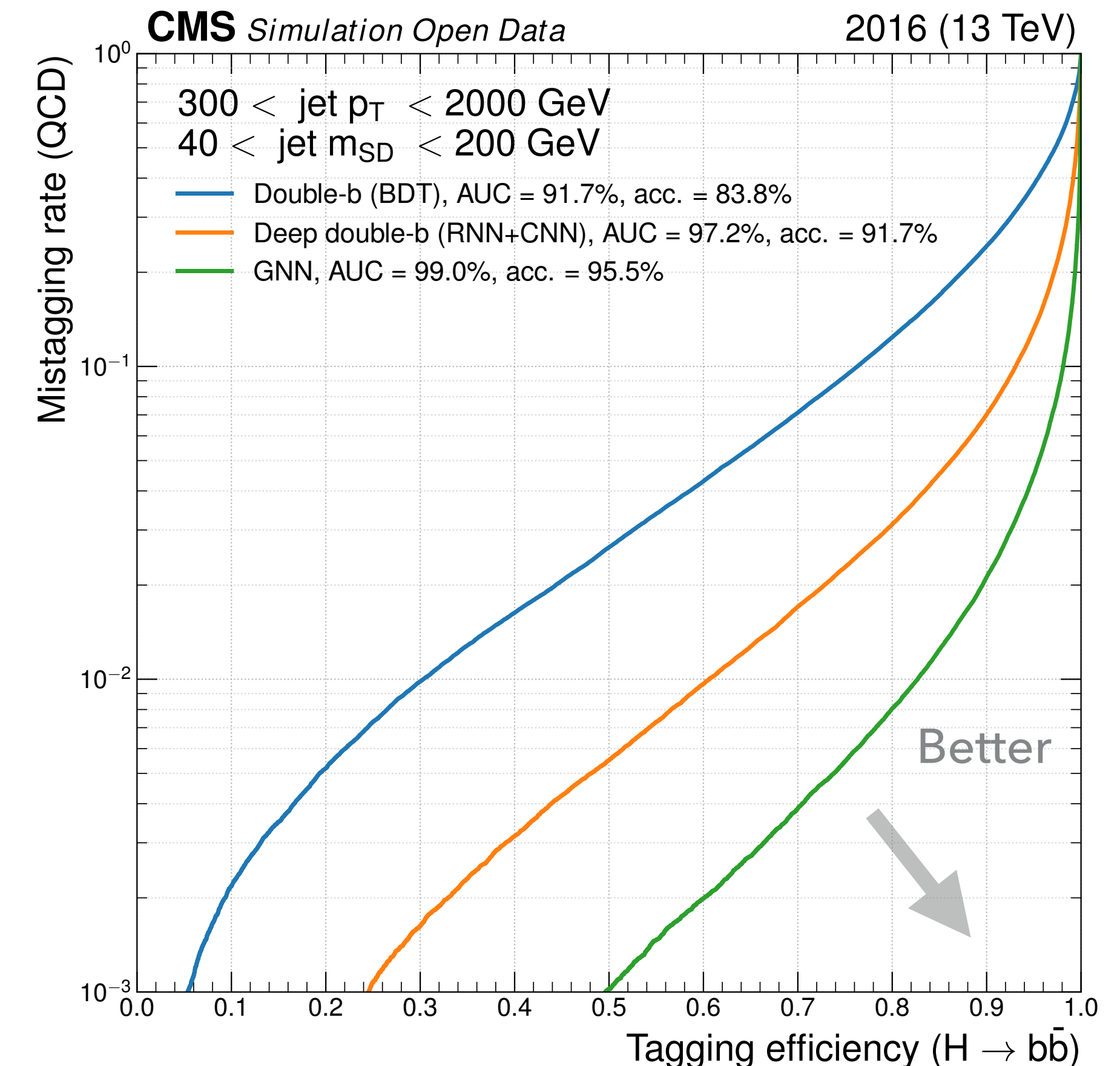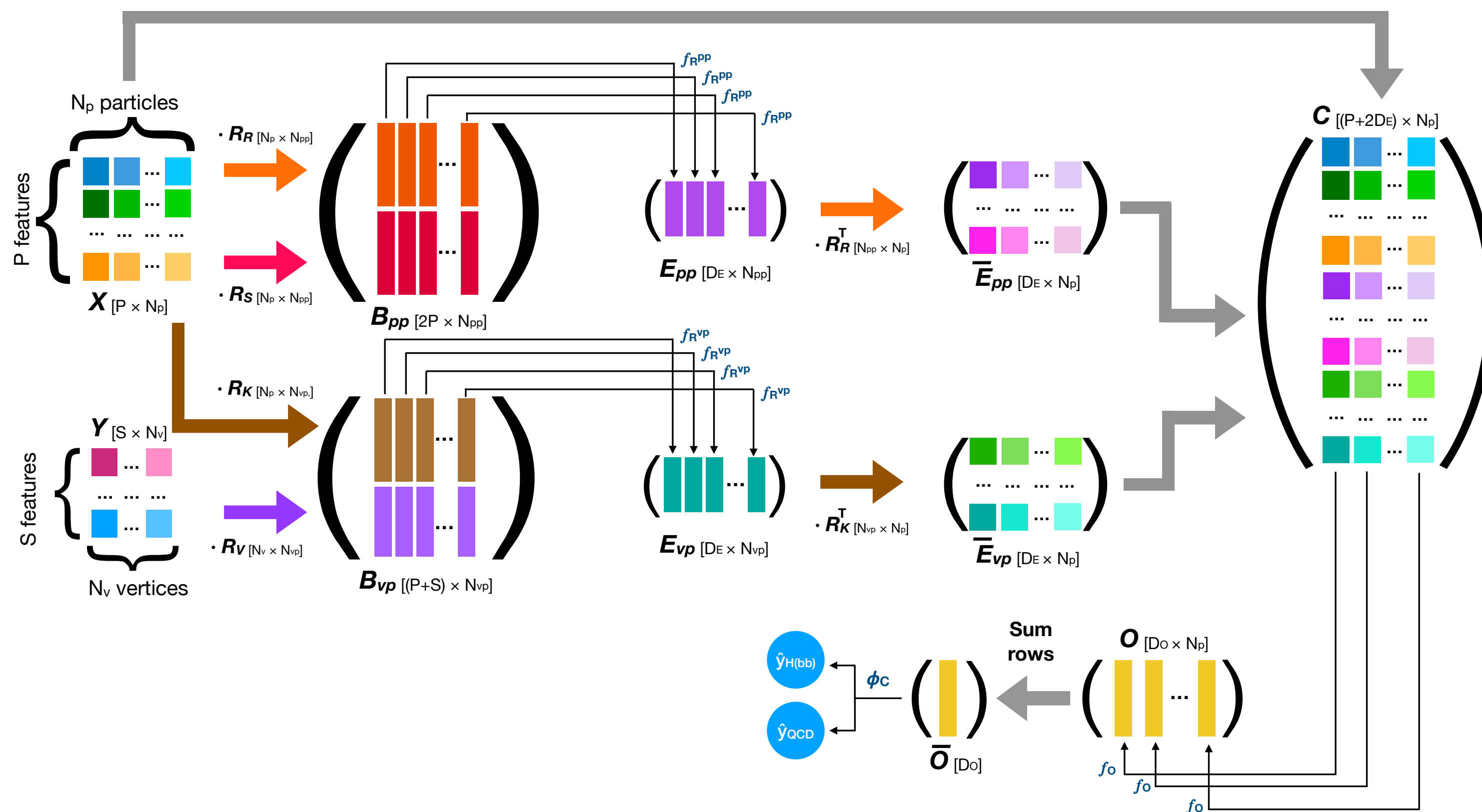
    ▸ secondary vertex features, e.g. flight distance

▸ In the process of evaluating FAIRness and contributing feedback to CERN Open Data Portal

　　▸ Similar to ARDC FAIRness self assessment tool:

　　https://ardc.edu.au/resources/working-with-data/fair-data/fair-self-assessment-tool/

▸ Take lessons learned and initiate a guide on evaluating FAIRness for the HEP community

Total across F.A.I.R

**Findable** ⓘ

Does the dataset have any identifiers assigned?
　No identifier

Is the dataset identifier included in all metadata records/files describing the data?
　No

How is the data described with metadata?
　The data is not described

What type of repository or registry is the metadata record in?
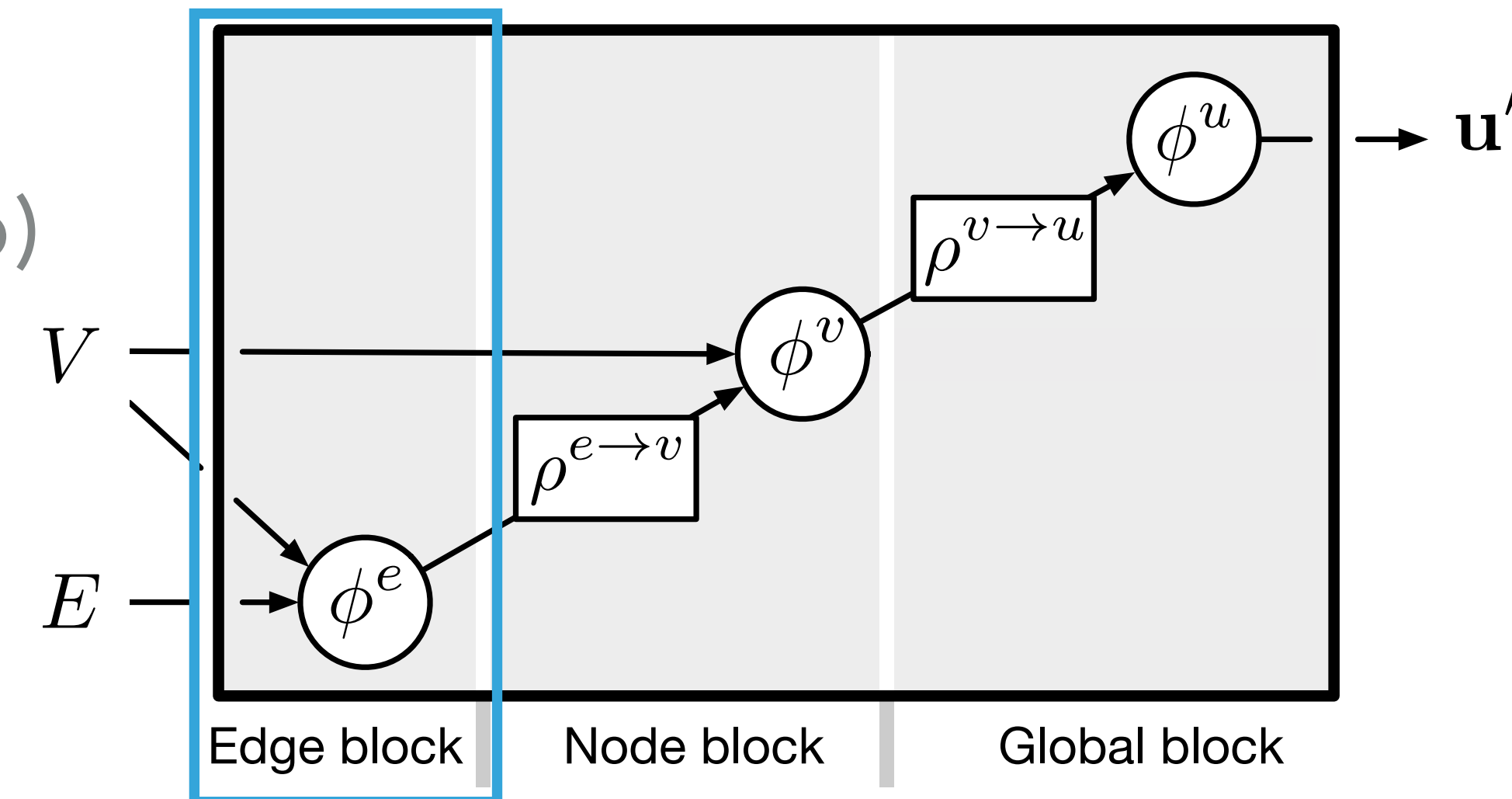　The data is not described in any repository

| FAIR Principle | F1: (Meta)data are assigned globally unique and persistent identifiers. |
|---|---|
| *METRIC* | *Pass/Fail with Comments* |
| **Identifier Uniqueness** : this metric measures whether there is a scheme to uniquely identify the digital resource. | **Pass**. There is a unique URL to a registered identifier scheme. The DOI: 10.7483/OPENDATA.CMS.JGJX.MS7Q of the data that resolves to this URL is also available. |
| **Identifier Persistence**: it measures whether there is a policy that describes what the provider will do in the event an identifier scheme becomes deprecated. | **Pass**. What the provider will do in the event an identifier scheme becomes deprecated? DOI provide persistent interoperable identifier. |
| FAIR Principle | F2: Data are described with rich metadata |
| *METRIC* | *Pass/Fail with Comments* |
| **Machine-readability of Metadata**: a URL to a document containing machine-readable metadata for the digital resource must be provided. | **Pass.** URL for JSON format metadata with REST API: http://opendata.cern.ch/api/records/12102. Also, running the url through https://search.google.com/test/rich-results shows the data page is eligible for rich results and the fields of metadata are machine readable. |
| **Richness of Metadata** | **(Newly added) Partially Pass.** The metadata can be improved with richer fields. Reviewing the datacite metadata for the DOI shows a pretty sparse record. |
| FAIR Principle | F3: Metadata clearly and explicitly include the identifier of the data they describe. |
| *METRIC* | *Pass/Fail with Comments* |
| **Resource Identifier in Metadata**: it measures whether the metadata document contains the identifier for the digital resource that meets F1 principle. | **Pass** The association between a metadata and the dataset is made explicit that the dataset's globally unique and persistent identifier can be found in the metadata. Specifically, the DOI is a top-level and a mandatory field in the metadata record. |
| FAIR Principle | F4: (Meta)data are registered or indexed in a searchable resource |
| *METRIC* | *Pass/Fail with Comments* |
| **Index in a searchable resource**: it measures the degree to which the digital resource can be found using web-based search engines | **Pass.** The dataset is indexed by Google Dataset Search engine. |

▸ Edge convolutions for particle-particle and particle-vertex connections update particle features; summed particle features used to predict H(bb) or QCD prob.

▸ GNN improves on previous methods

▸ Model and code: https://github.com/eric-moreno/IN [10.5281/zenodo.3891869]



CMS *Simulation Open Data*    2016 (13 TeV)

$300 <$ jet $p_T < 2000$ GeV
$40 <$ jet $m_{SD} < 200$ GeV

Double-b (BDT), AUC = 91.7%, acc. = 83.8%
Deep double-b (RNN+CNN), AUC = 97.2%, acc. = 91.7%
GNN, AUC = 99.0%, acc. = 95.5%

Mistagging rate (QCD)

Tagging efficiency (H → b$\bar{\text{b}}$)

Better

- ▸ GNN tutorial with PyTorch Geometric: <u>UCSD Data Science Capstone</u>
- ▸ Environment specified with docker and conda
- ▸ CI deployed in GitHub Actions
- ▸ Expand this example into a FAIR AI model (via e.g. DLHub)
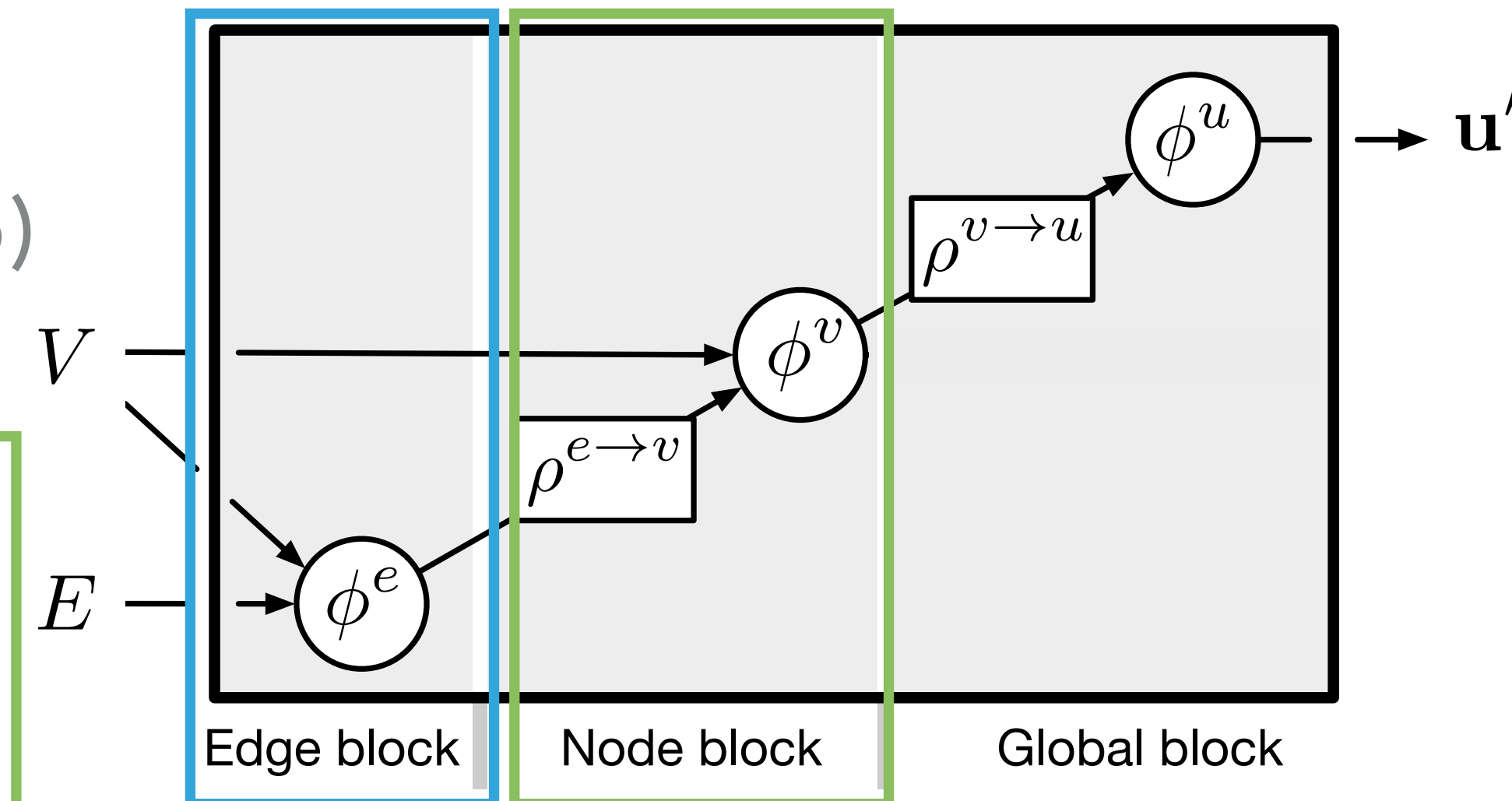
```python
class EdgeBlock(torch.nn.Module):
    def __init__(self):
        super(EdgeBlock, self).__init__()
        self.edge_mlp = Seq(Lin(inputs*2, hidden),
                            BatchNorm1d(hidden),
                            ReLU(),
                            Lin(hidden, hidden))

    def forward(self, src, dest, edge_attr, u, batch):
        out = torch.cat([src, dest], 1)
        return self.edge_mlp(out)
```



Edge block    Node block    Global block

▸ GNN tutorial with PyTorch Geometric: [UCSD Data Science Capstone](UCSD Data Science Capstone)

▸ Environment specified with docker and conda

▸ CI deployed in GitHub Actions

▸ Expand this example into a FAIR AI model (via e.g. DLHub)

```python
class EdgeBlock(torch.nn.Module):
    def __init__(self):
        super(EdgeBlock, self).__init__()
        self.edge_mlp = Seq(Lin(inputs*2, hidden),
                            BatchNorm1d(hidden),
                            ReLU(),
                            Lin(hidden, hidden))

    def forward(self, src, dest, edge_attr, u, batch):
        out = torch.cat([src, dest], 1)
        return self.edge_mlp(out)
```

```python
class NodeBlock(torch.nn.Module):
    def __init__(self):
        super(NodeBlock, self).__init__()
        self.node_mlp_1 = Seq(Lin(inputs+hidden, hidden),
                              BatchNorm1d(hidden),
                              ReLU(),
                              Lin(hidden, hidden))
        self.node_mlp_2 = Seq(Lin(inputs+hidden, hidden),
                              BatchNorm1d(hidden),
                              ReLU(),
                              Lin(hidden, hidden))

    def forward(self, x, edge_index, edge_attr, u, batch):
        row, col = edge_index
        out = torch.cat([x[row], edge_attr], dim=1)
        out = self.node_mlp_1(out)
        out = scatter_mean(out, col, dim=0, dim_size=x.size(0))
        out = torch.cat([x, out], dim=1)
        return self.node_mlp_2(out)
```



$\phi^u \rightarrow \mathbf{u'}$

$\rho^{v \rightarrow u}$

$V$ — $\phi^v$

$\rho^{e \rightarrow v}$

$E$ — $\phi^e$

Edge block  Node block  Global block

▸ GNN tutorial with PyTorch Geometric: [UCSD Data Science Capstone](UCSD Data Science Capstone)

▸ Environment specified with docker and conda

▸ CI deployed in GitHub Actions

▸ Expand this example into a FAIR AI model (via e.g. DLHub)

```python
class EdgeBlock(torch.nn.Module):
    def __init__(self):
        super(EdgeBlock, self).__init__()
        self.edge_mlp = Seq(Lin(inputs*2, hidden),
                            BatchNorm1d(hidden),
                            ReLU(),
                            Lin(hidden, hidden))

    def forward(self, src, dest, edge_attr, u, batch):
        out = torch.cat([src, dest], 1)
        return self.edge_mlp(out)
```
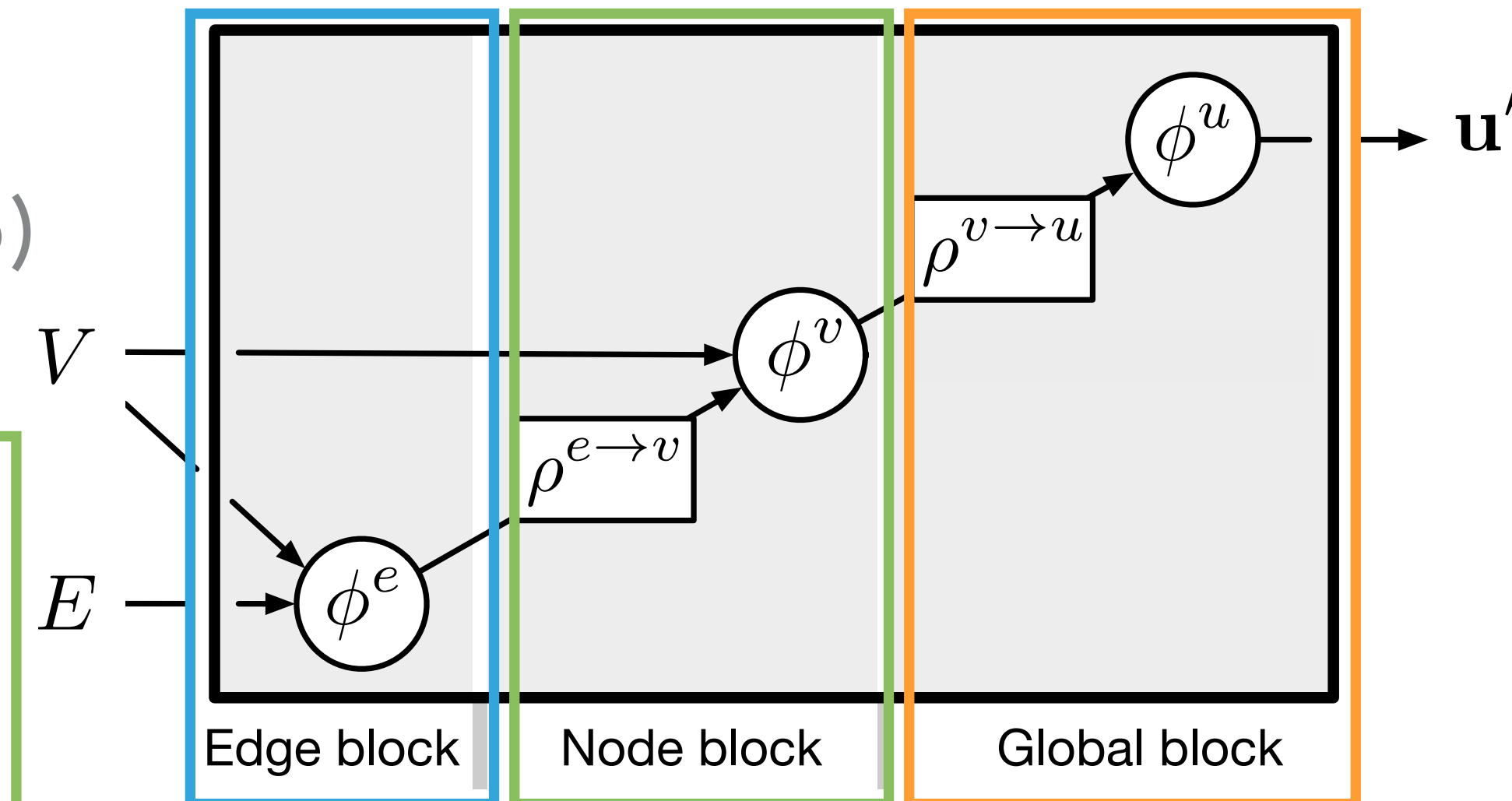
```python
class GlobalBlock(torch.nn.Module):
    def __init__(self):
        super(GlobalBlock, self).__init__()
        self.global_mlp = Seq(Lin(hidden, hidden),
                              BatchNorm1d(hidden),
                              ReLU(),
                              Lin(hidden, outputs))

    def forward(self, x, edge_index, edge_attr, u, batch):
        out = scatter_mean(x, batch, dim=0)
        return self.global_mlp(out)
```

```python
class NodeBlock(torch.nn.Module):
    def __init__(self):
        super(NodeBlock, self).__init__()
        self.node_mlp_1 = Seq(Lin(inputs+hidden, hidden),
                              BatchNorm1d(hidden),
                              ReLU(),
                              Lin(hidden, hidden))
        self.node_mlp_2 = Seq(Lin(inputs+hidden, hidden),
                              BatchNorm1d(hidden),
                              ReLU(),
                              Lin(hidden, hidden))

    def forward(self, x, edge_index, edge_attr, u, batch):
        row, col = edge_index
        out = torch.cat([x[row], edge_attr], dim=1)
        out = self.node_mlp_1(out)
        out = scatter_mean(out, col, dim=0, dim_size=x.size(0))
        out = torch.cat([x, out], dim=1)
        return self.node_mlp_2(out)
```
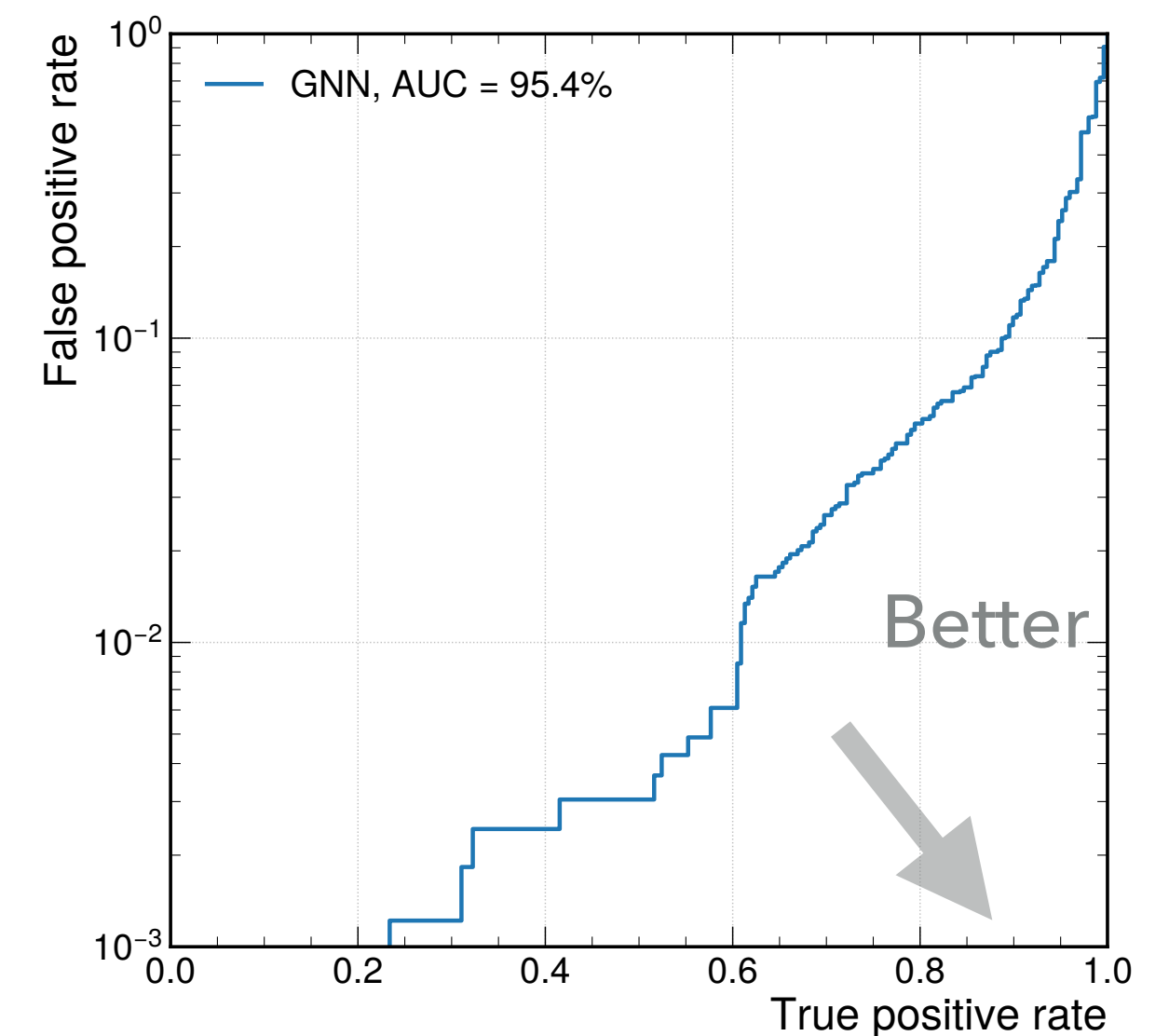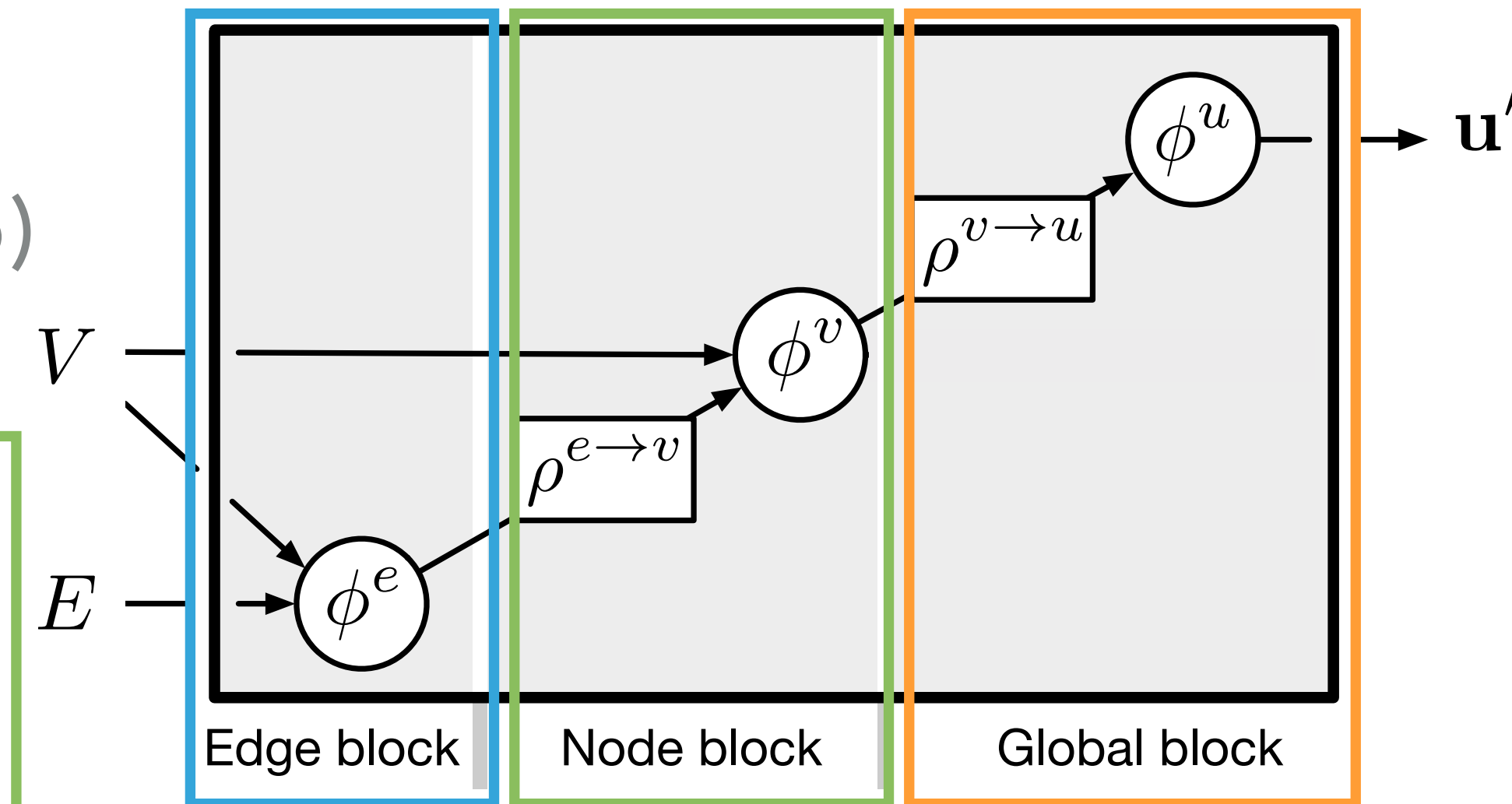
▸ GNN tutorial with PyTorch Geometric: <u>UCSD Data Science Capstone</u>

▸ Environment specified with docker and conda

▸ CI deployed in GitHub Actions

▸ Expand this example into a FAIR AI model (via e.g. DLHub)

```python
class EdgeBlock(torch.nn.Module):
    def __init__(self):
        super(EdgeBlock, self).__init__()
        self.edge_mlp = Seq(Lin(inputs*2, hidden),
                            BatchNorm1d(hidden),
                            ReLU(),
                            Lin(hidden, hidden))

    def forward(self, src, dest, edge_attr, u, batch):
        out = torch.cat([src, dest], 1)
        return self.edge_mlp(out)
```

```python
class NodeBlock(torch.nn.Module):
    def __init__(self):
        super(NodeBlock, self).__init__()
        self.node_mlp_1 = Seq(Lin(inputs+hidden, hidden),
                              BatchNorm1d(hidden),
                              ReLU(),
                              Lin(hidden, hidden))
        self.node_mlp_2 = Seq(Lin(inputs+hidden, hidden),
                              BatchNorm1d(hidden),
                              ReLU(),
                              Lin(hidden, hidden))

    def forward(self, x, edge_index, edge_attr, u, batch):
        row, col = edge_index
        out = torch.cat([x[row], edge_attr], dim=1)
        out = self.node_mlp_1(out)
        out = scatter_mean(out, col, dim=0, dim_size=x.size(0))
        out = torch.cat([x, out], dim=1)
        return self.node_mlp_2(out)
```

```python
class GlobalBlock(torch.nn.Module):
    def __init__(self):
        super(GlobalBlock, self).__init__()
        self.global_mlp = Seq(Lin(hidden, hidden),
                              BatchNorm1d(hidden),
                              ReLU(),
                              Lin(hidden, outputs))

    def forward(self, x, edge_index, edge_attr, u, batch):
        out = scatter_mean(x, batch, dim=0)
        return self.global_mlp(out)
```
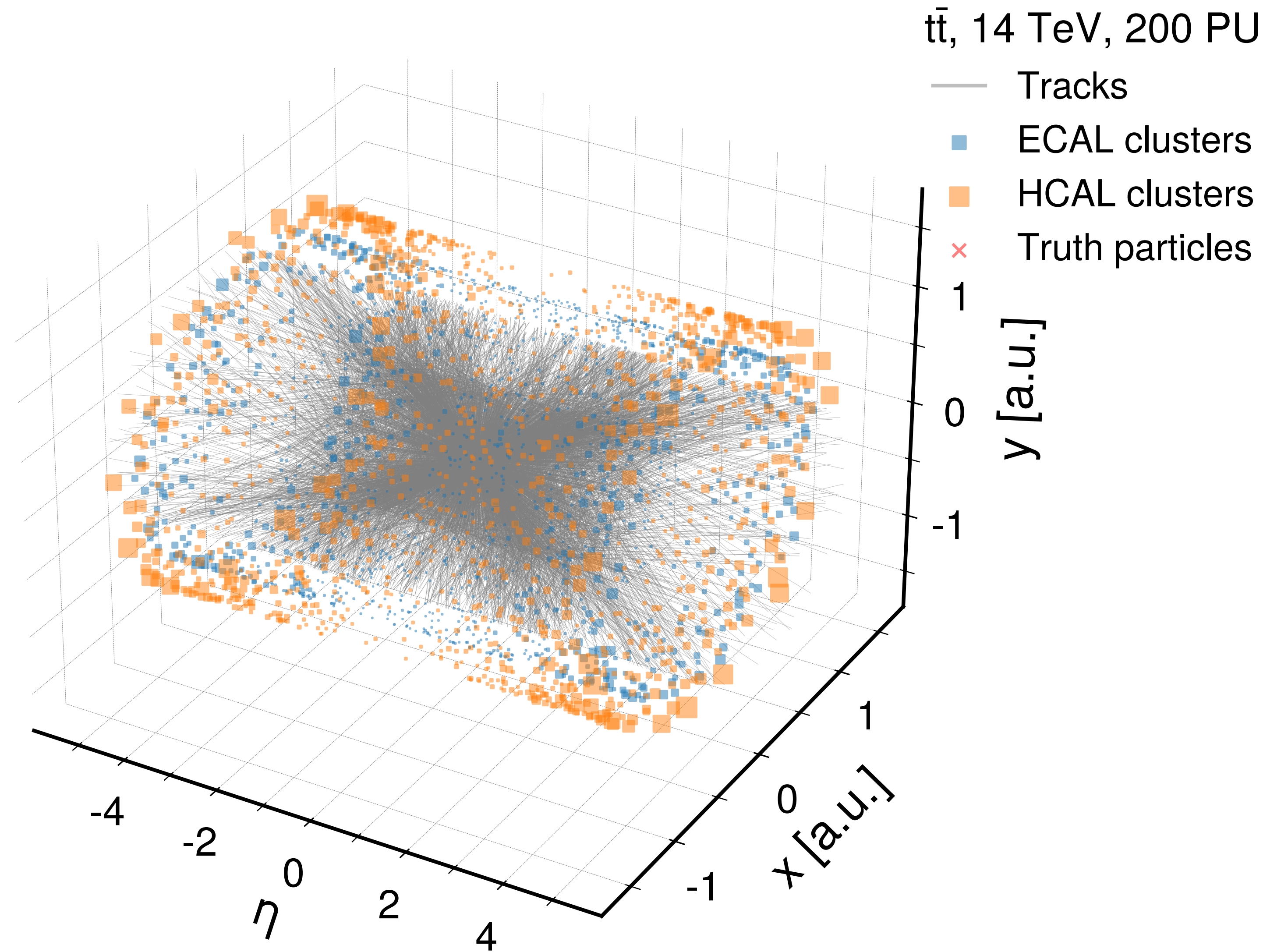
```python
class InteractionNetwork(torch.nn.Module):
    def __init__(self):
        super(InteractionNetwork, self).__init__()
        self.interactionnetwork = MetaLayer(EdgeBlock(), NodeBlock(), GlobalBlock())
        self.bn = BatchNorm1d(inputs)

    def forward(self, x, edge_index, batch):
        x = self.bn(x)
        x, edge_attr, u = self.interactionnetwork(x, edge_index, None, None, batch)
        return u
```
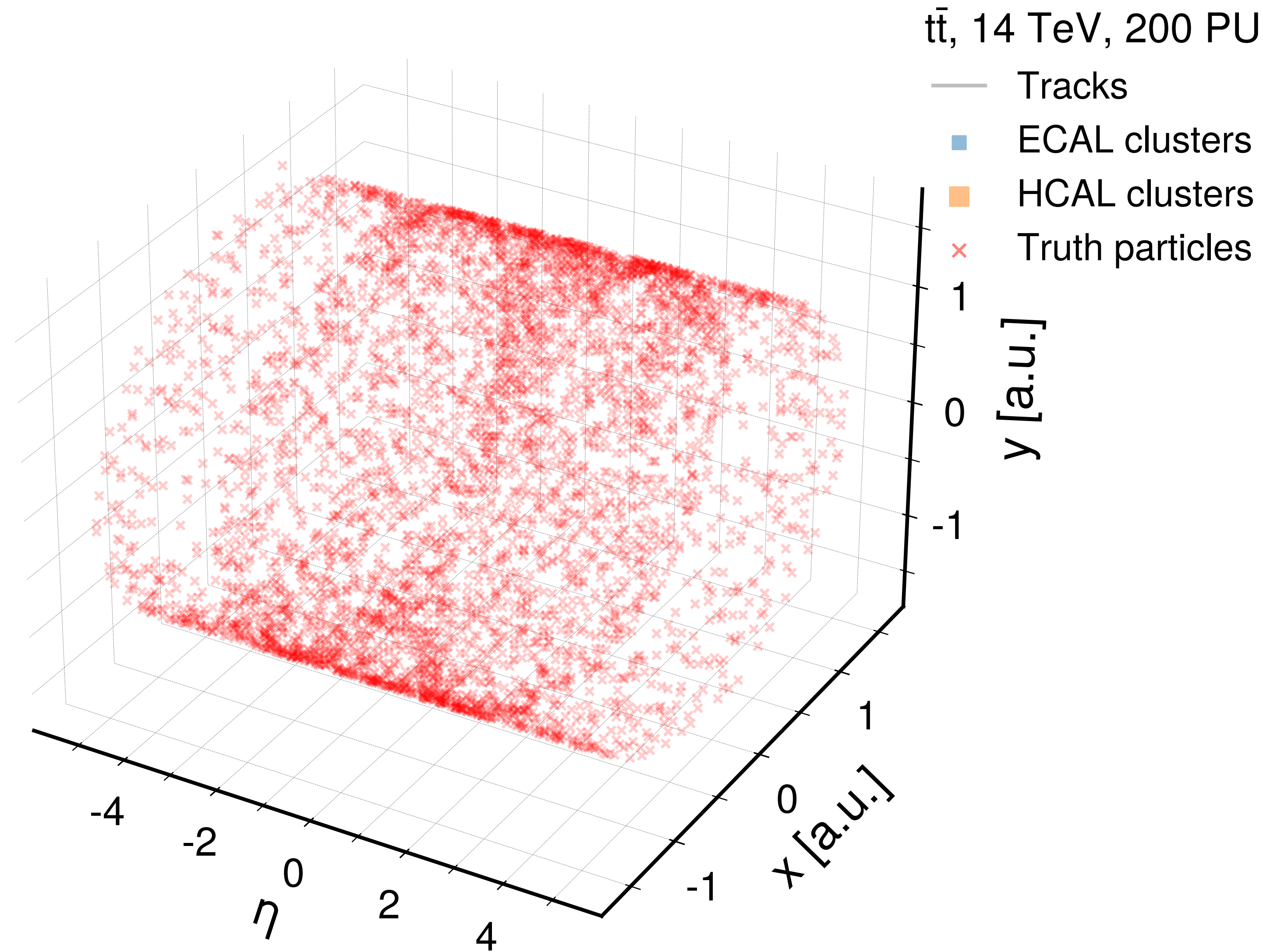


Edge block    Node block    Global block

▸ Particles interact with detector, leaving energy deposits (ECAL clusters, HCAL clusters) and tracks

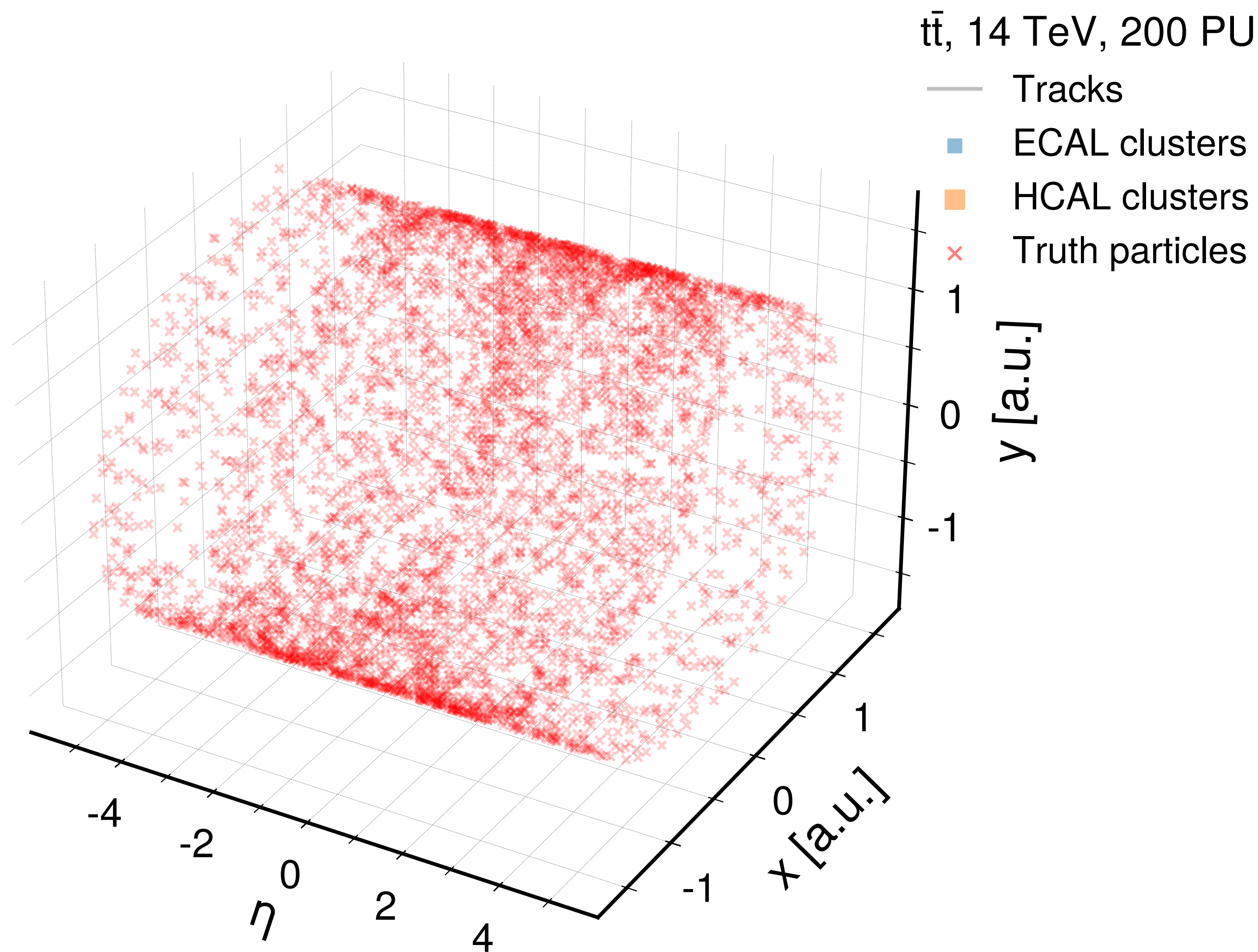▸ Particles interact with detector, leaving energy deposits (ECAL clusters, HCAL clusters) and tracks

▸ PF: combine info from complementary detector subsystems to produce a holistic, particle interpretation of the event (truth particles)



t$\bar{\text{t}}$, 14 TeV, 200 PU

— Tracks
■ ECAL clusters
■ HCAL clusters
× Truth particles

▸ Particles interact with detector, leaving energy deposits (ECAL clusters, HCAL clusters) and tracks

▸ PF: combine info from complementary detector subsystems to produce a holistic, particle interpretation of the event (truth particles)

▸ Goal: construct a mapping that minimizes some distance between truth particles and reconstructed particles



t̄t, 14 TeV, 200 PU
— Tracks
■ ECAL clusters
■ HCAL clusters
× Truth particles

▸ Particles interact with detector, leaving energy deposits (ECAL clusters, HCAL clusters) and tracks

▸ PF: combine info from complementary detector subsystems to produce a holistic, particle interpretation of the event (truth particles)

▸ Goal: construct a mapping that minimizes some distance between truth particles and reconstructed particles

▸ Make public dataset: [10.5281/zenodo.4559324] and AI model [10.5281/zenodo.4559587] **FAIR**



$t\bar{t}$, 14 TeV, 200 PU
— Tracks
▪ ECAL clusters
▪ HCAL clusters
× Truth particles
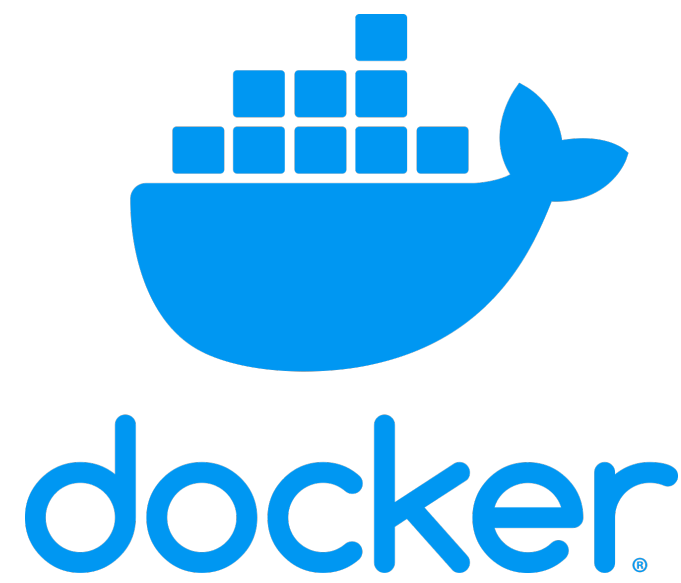
# VISION AND OUTLOOK

Data repositories

Data repositories



Deployable AI models

Data repositories

Deployable AI models

Papers / indexing / search / discovery

Data repositories

Deployable AI models



Papers With Code

CONNECTED PAPERS

Data repositories

Papers / indexing /
search / discovery

Competitions?

Deployable AI models

Papers With Code

CONNECTED
PAPERS

Papers / indexing / search / discovery

Data repositories

Competitions?

Deployable AI models

Automated ML (retraining) flows
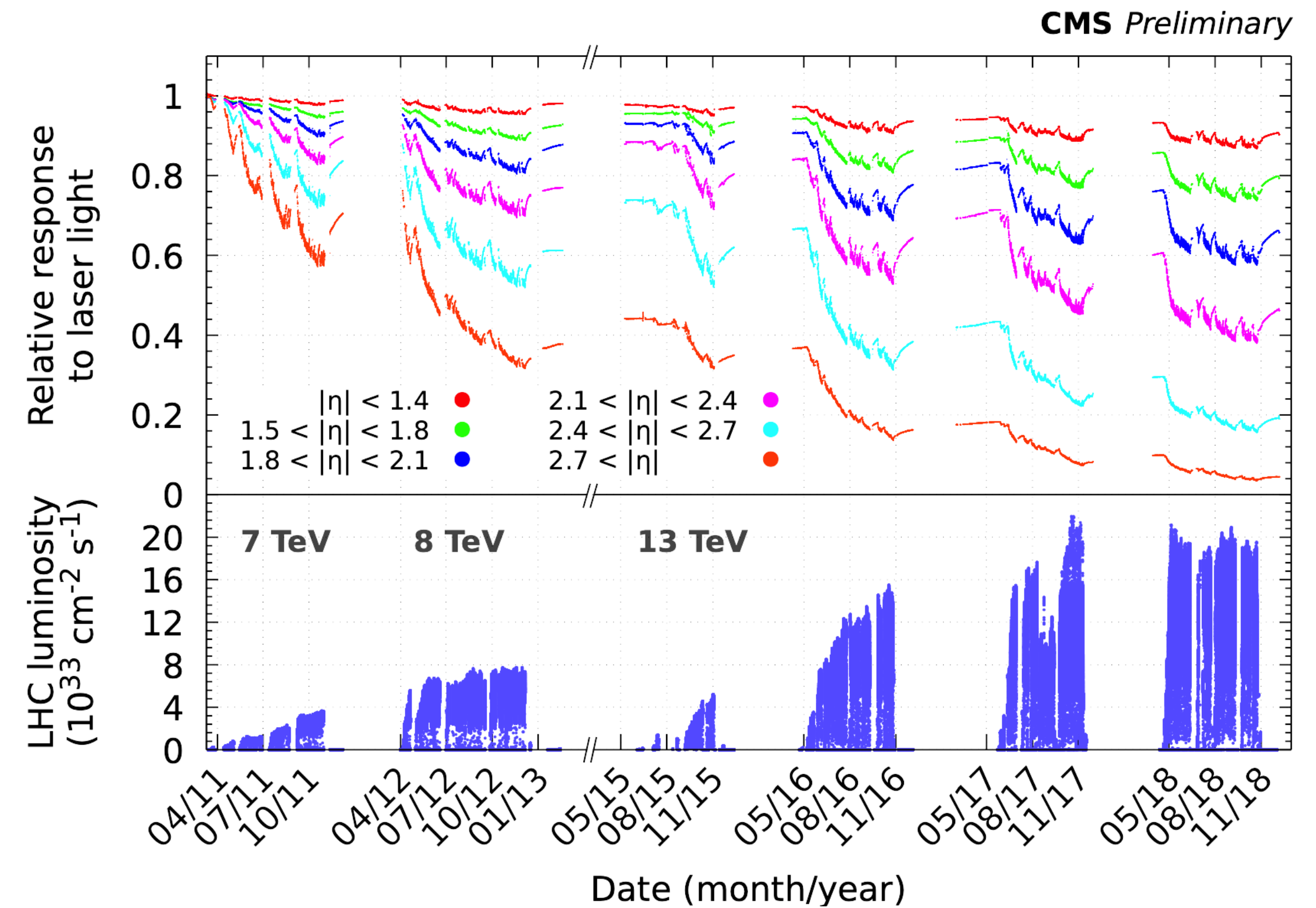
Papers With Code

CONNECTED PAPERS

▸ Goal of FAIR4HEP is to interpret and refine what FAIR means for HEP data/models

  ▸ Enable "plug and play" datasets: allow for combinations of different computing resources

▸ Vision: connected services linking datasets, benchmark models (code), deployment servers, and publications to make everything more FAIR

  ▸ Simpler discovery of new datasets and models

▸ Projects

  ▸ Evaluate FAIRness of existing public datasets

  ▸ Standardize FAIR publication of AI models in HEP

  ▸ Create example FAIR datasets and AI models

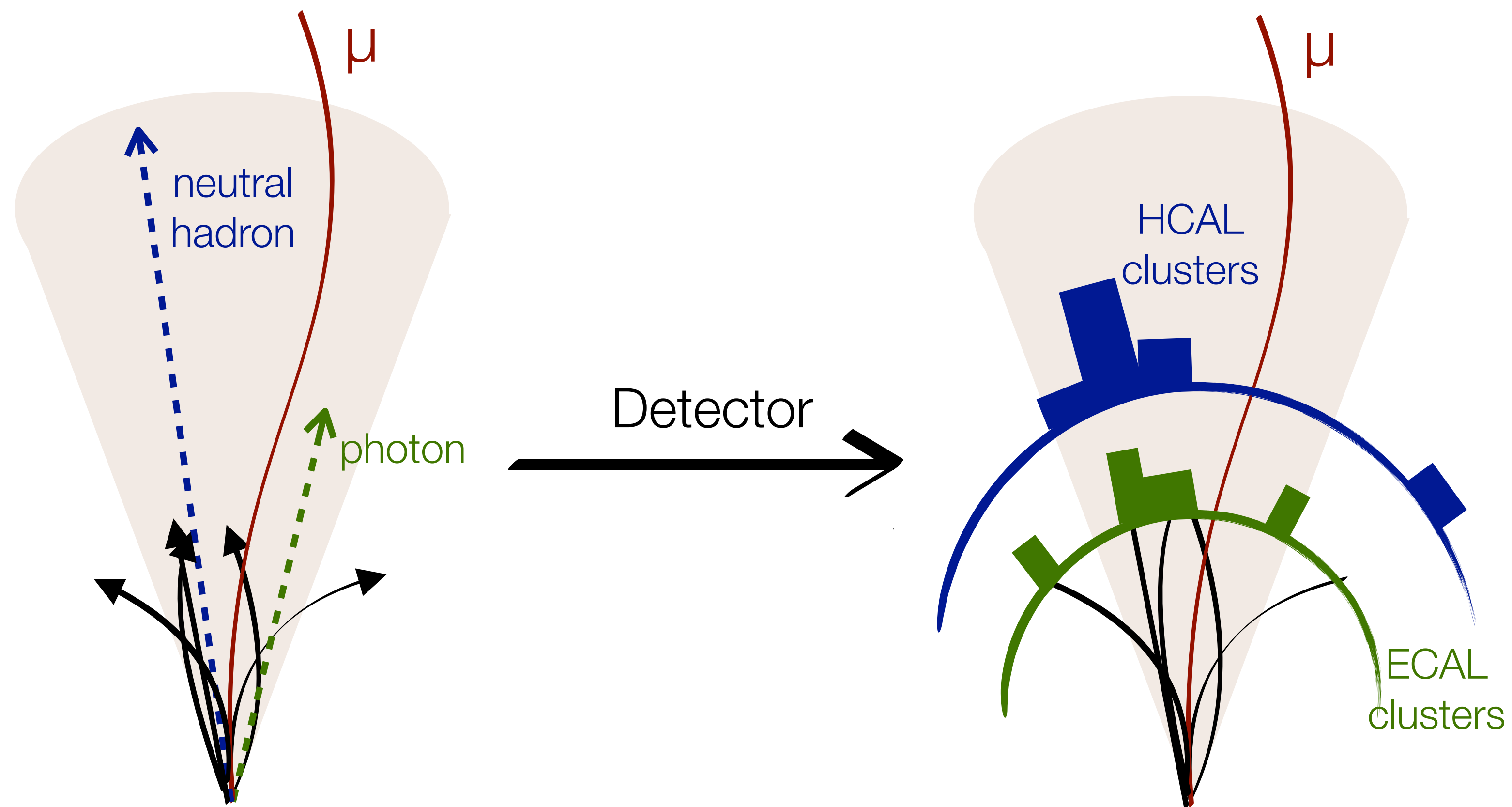  ▸ Enhance existing services to make them more FAIR

▸ Welcome feedback!

# JAVIER DUARTE
# IRIS-HEP TOPICAL MEETING
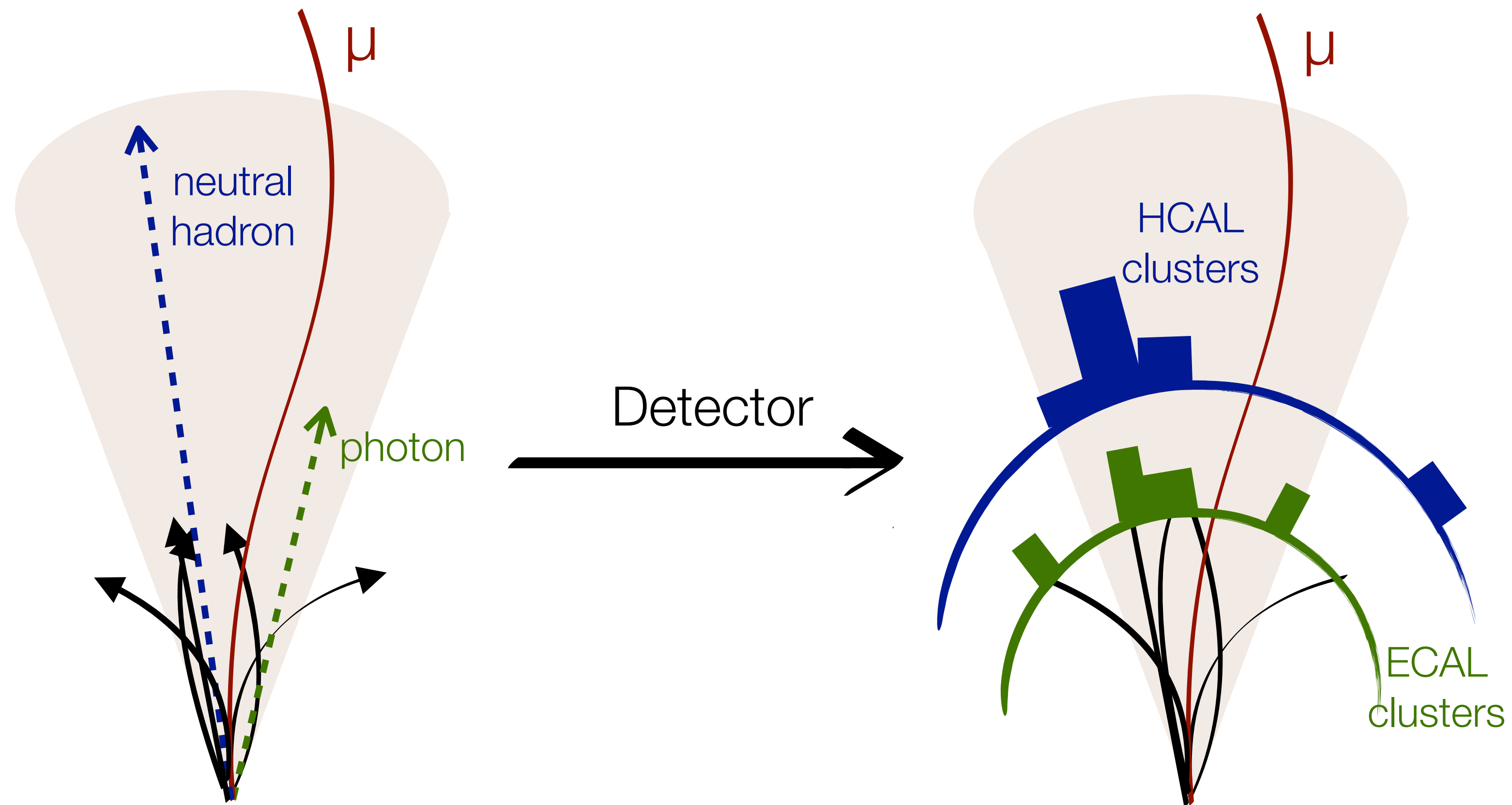# APRIL 21, 2021

# BACKUP

▸ High radiation dose causes
transparency loss in ECAL $\eta$ crystals

    ▸ Crystals recover over time

    ▸ ~70,000 crystals ×
~10,000 calibrations per year =
~700,000,000 learnable parameters

▸ Developing a public FAIR dataset (and
AI model) to study (and predict) time
dependence of transparency loss



CMS *Preliminary*

Legend:
$|\eta| < 1.4$ ●
$1.5 < |\eta| < 1.8$ ●
$1.8 < |\eta| < 2.1$ ●
$2.1 < |\eta| < 2.4$ ●
$2.4 < |\eta| < 2.7$ ●
$2.7 < |\eta|$ ●

7 TeV    8 TeV    13 TeV

Relative response to laser light

LHC luminosity ($10^{33}$ cm$^{-2}$ s$^{-1}$)

Date (month/year)
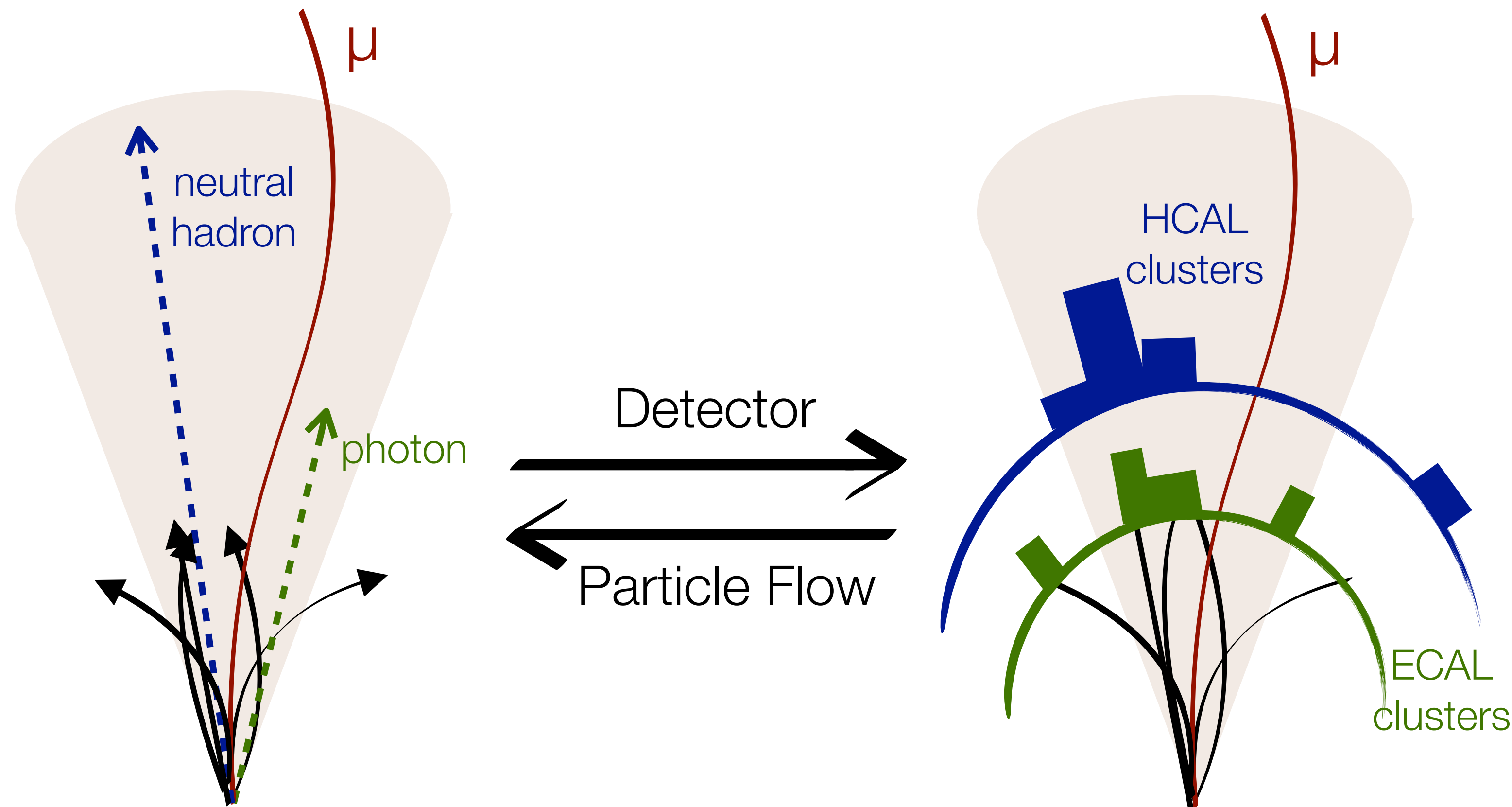
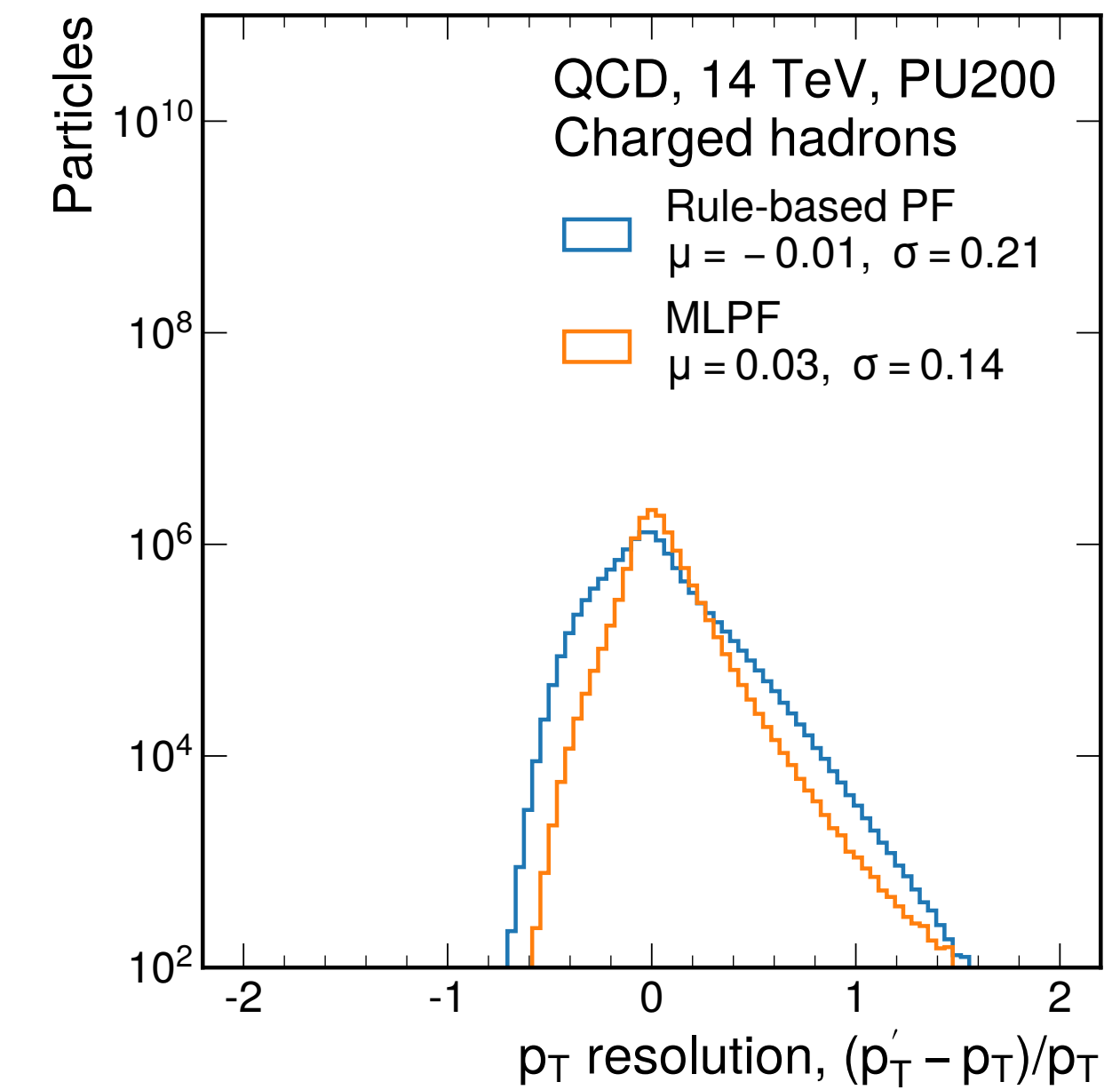▸ Particles interact with detector, leaving energy deposits and tracks
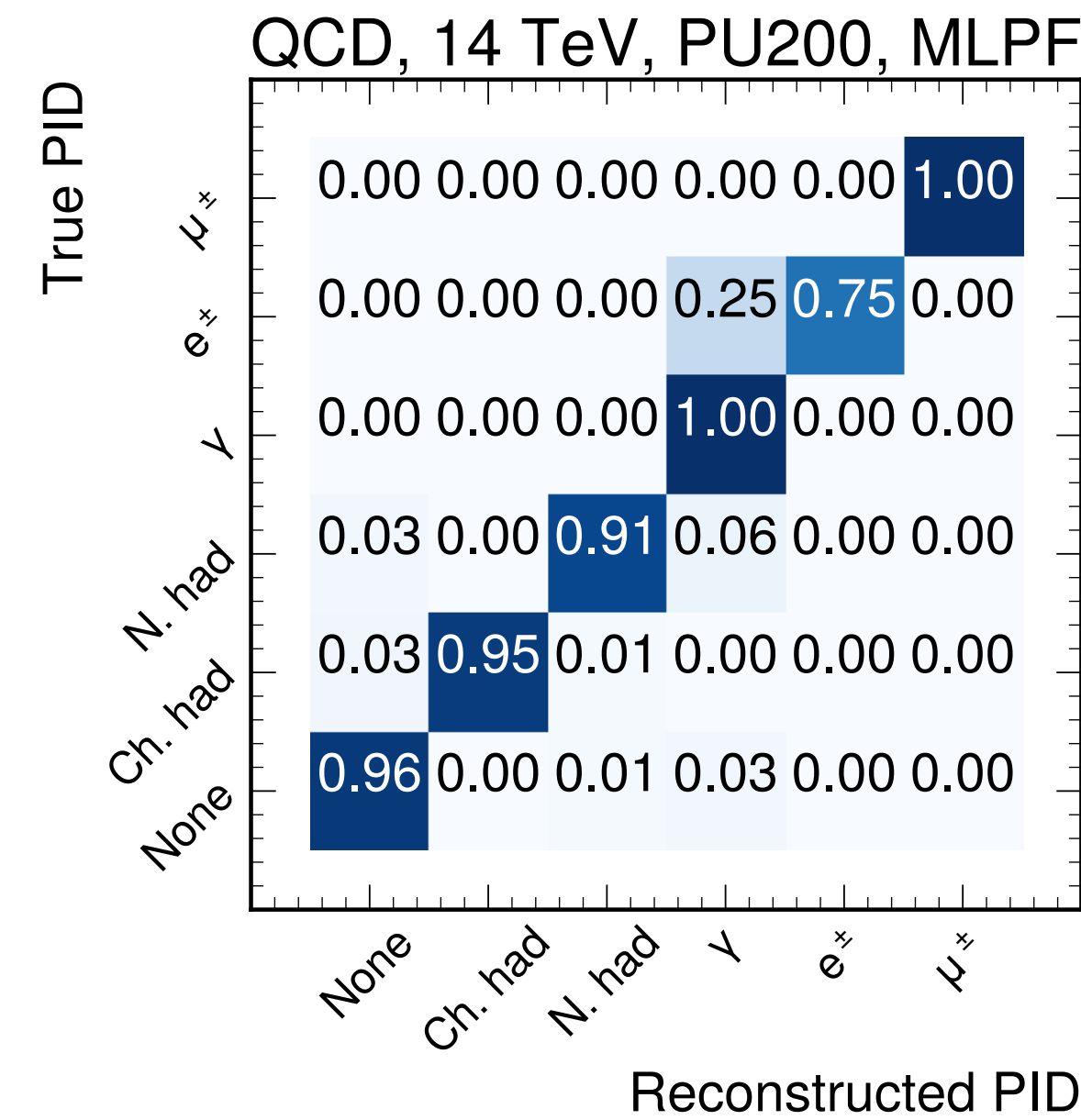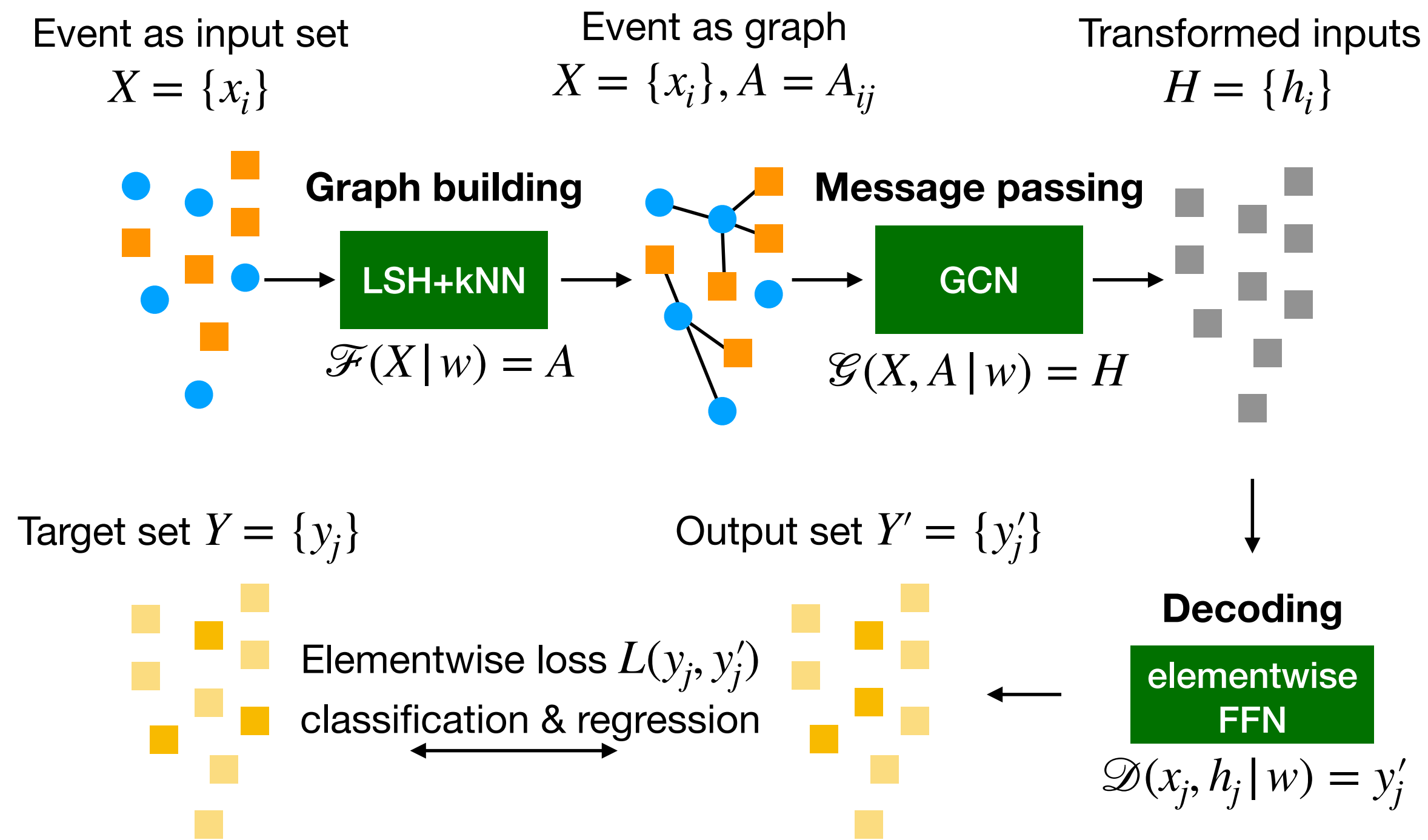
▸ Particles interact with detector, leaving energy deposits and tracks

▸ Efficient combination of info. from complementary detector subsystems to produce a holistic, particle interpretation of the event (that improves on any individual subsystem)

- Pythia8+Delphes3 50,000 $t\bar{t}$+jets and 5,000 QCD events produced in pp collisions at 14 TeV with 200 pileup

- Partial provenance (Pythia+Delphes configuration, versions etc.)

- Detector calorimeter towers and tracks as input and generator particles as ground truth

  - bzip2-compressed python pickle including

$X = \{x_i\}$

$X = \{x_i\}, A = A_{ij}$

$\mathcal{F}(X \mid w) = A$

$\mathcal{G}(X, A \mid w) :$

$Y = \{y_j\}$

$Y' = \{y_j'\}$

$L(y_j, y_j')$

$\mathcal{D}(x_j, h_j \mid w) = y_j'$

$x_i = [\text{type}, p_{\text{T}}, E_{\text{ECAL}}, E_{\text{HCAL}}, \eta, \phi, \eta_{\text{outer}}, \phi_{\text{outer}}, q, \ldots], \quad \text{type} \in \{\text{track, cluster}\}$

$y_i = [\text{PID}, p_{\text{T}}, E, \eta, \phi, q, \ldots], \quad \text{PID} \in \{\text{none, charged hadron, neutral hadron}, \gamma, e^{\pm}, \mu^{\pm}\}$

$h_i \in \mathbb{R}^{256}$

$\mathcal{F}, \mathcal{G}, \mathcal{D}$

Event as input set
$X = \{x_i\}$

Event as graph
$X = \{x_i\}, A = A_{ij}$

Transformed inputs
$H = \{h_i\}$

**Graph building**

LSH+kNN

$\mathcal{F}(X \mid w) = A$

**Message passing**

GCN

$\mathcal{G}(X, A \mid w) = H$

**Decoding**

elementwise
FFN

$\mathcal{D}(x_j, h_j \mid w) = y'_j$

Target set $Y = \{y_j\}$

Output set $Y' = \{y'_j\}$

Elementwise loss $L(y_j, y'_j)$
classification & regression

$x_i = [\text{type}, p_{\text{T}}, E_{\text{ECAL}}, E_{\text{HCAL}}, \eta, \phi, \eta_{\text{outer}}, \phi_{\text{outer}}, q, \ldots], \quad \text{type} \in \{\text{track}, \text{cluster}\}$

$y_j = [\text{PID}, p_{\text{T}}, E, \eta, \phi, q, \ldots], \quad \text{PID} \in \{\text{none}, \text{charged hadron}, \text{neutral hadron}, \gamma, e^{\pm}, \mu^{\pm}\}$

$h_i \in \mathbb{R}^{256}$

$\mathcal{F}, \mathcal{G}, \mathcal{D}$

QCD, 14 TeV, PU200, MLPF

True PID

|  | None | Ch. had | N. had | γ | e± | μ± |
|---|---|---|---|---|---|---|
| μ± | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 1.00 |
| e± | 0.00 | 0.00 | 0.00 | 0.25 | 0.75 | 0.00 |
| γ | 0.00 | 0.00 | 0.00 | 1.00 | 0.00 | 0.00 |
| N. had | 0.03 | 0.00 | 0.91 | 0.06 | 0.00 | 0.00 |
| Ch. had | 0.03 | 0.95 | 0.01 | 0.00 | 0.00 | 0.00 |
| None | 0.96 | 0.00 | 0.01 | 0.03 | 0.00 | 0.00 |

Reconstructed PID

Particles

QCD, 14 TeV, PU200
Charged hadrons

Rule-based PF
$\mu = -0.01, \; \sigma = 0.21$

MLPF
$\mu = 0.03, \; \sigma = 0.14$

$p_{\text{T}}$ resolution, $(p'_{\text{T}} - p_{\text{T}})/p_{\text{T}}$

Classification and regression
performance: MLPF same or better
than rule-based PF algorithm

▸ Models and training code available:

▸ https://github.com/jpata/particleflow
  [10.5281/zenodo.4559587]

▸ Environment specified with singularity

▸ CI deployed through GitHub Actions