# Clustering in SPD Range System

**Georgy Golovanov**

*on behalf of SPD Muon Group*

# Outline

**Two steps of particle identification in SPD Range System**

- Clustering
- Particle tagging

**Coordinate systems for clustering**

**Reference clustering in SpdRoot**

**Clustering algorithms**

- k-Means algorithm
- DBSCAN-based algorithm

**Model performance**

- Purity and V-measure
- Muon/Pion efficiencies

**Summary**

# Particle reconstruction in Range System

**Main goal:** to develop algorithms able to reconstruct muons and hadrons based on the information from Range System *standalone* by using machine learning techniques.

This is aimed to speed up the reconstruction process compared to traditional Kalman Filter technique.

Information available from Range System:
- two coordinate electronic readout: wires and strips — provides 3-dimension hits (thanks Alexander for strip emulation in SpdRoot);
- hits in **Barrel**: *(x, y)* of wires at layers and *z* of strips (30mm pitch);
- hits in **EndCaps**: *(y, z)* of wires and *x* of strips;

**Reconstruction task is split on *two steps*:**

**Step 1: Clustering** — forms group of hits (clusters) in Range System
- K-Means;
- DBSCAN;
- etc.

**Step 2: Particle identification (muon/pion cluster tagging)** — separates muon clusters from the clusters associated with hadrons
- Boosted Decision Tree algorithms (XGBoost, CatBoost);
- Convolutional Neural Network;
- etc.

# Coordinate system for clustering

There is a choice of coordinate systems for clustering in RS:

1. **(x,y) in Barrel + (z,y) in Endcaps**          *Wire readout*

   + natural for MDT wires orientation;
   - needs cluster matching between Barrel and EndCaps.

2. **(r,φ) in Barrel + (z,θ) in EndCaps**

   + takes into account projection geometry;
   - needs cluster matching between Barrel and EndCaps.
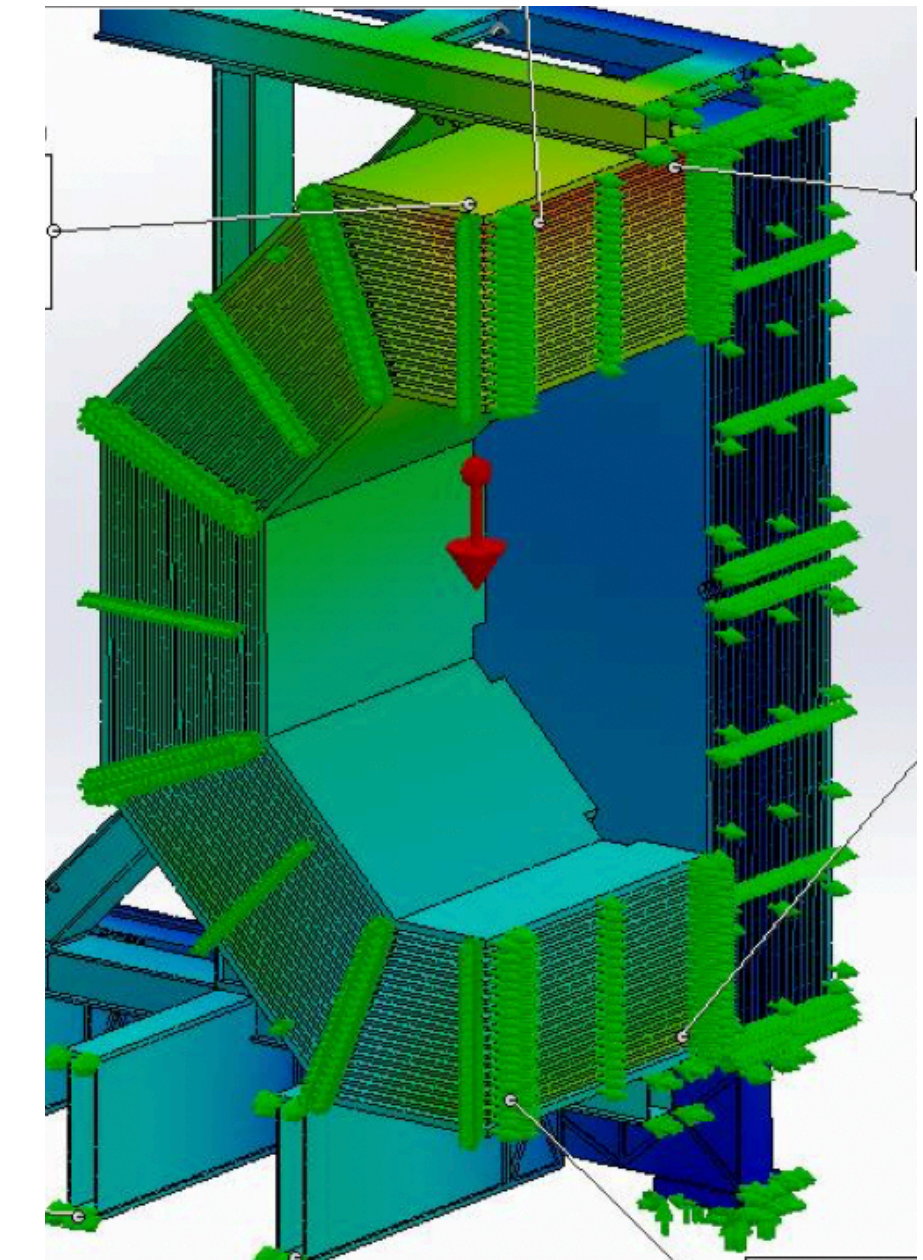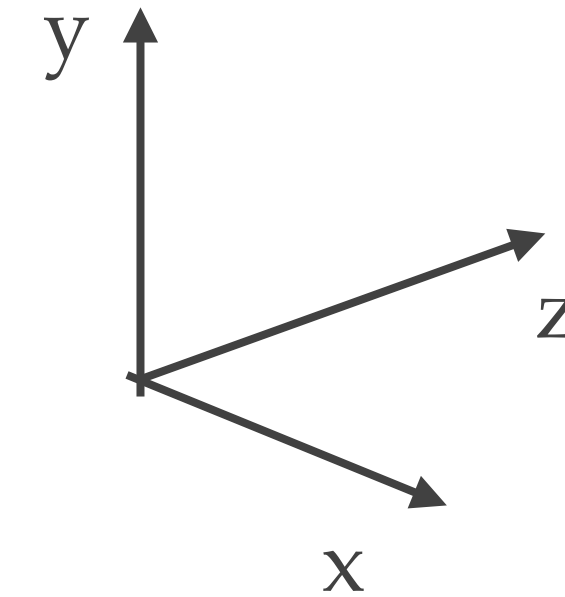
3. **(θ,φ) or (η,φ)**          {hit map}          *Wire/Strip readout*

   + no cluster split between Barrel/EndCaps;
   - clusters can be very dense since all hits are summed over 20 RS layers;
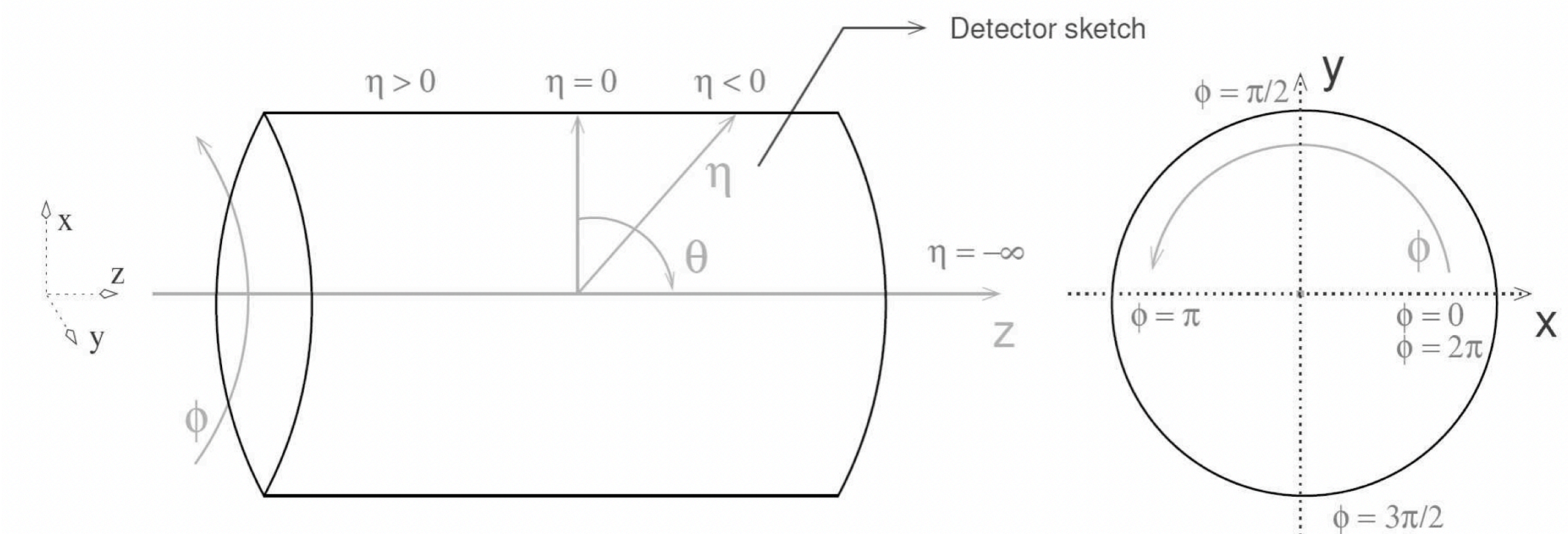
4. **(x,y,z)**          {hit space}

   + natural for MDT wires and strips orientation;
   + no cluster split between Barrel/EndCaps;
   + takes into account 3D shape of the cluster, no overlapping;
   - algorithms may be slower in 3D-space compared to 2D;

The chosen coordinate system should be in agreement with the later pion/muon separation algorithms, since they will use found clusters as an input.
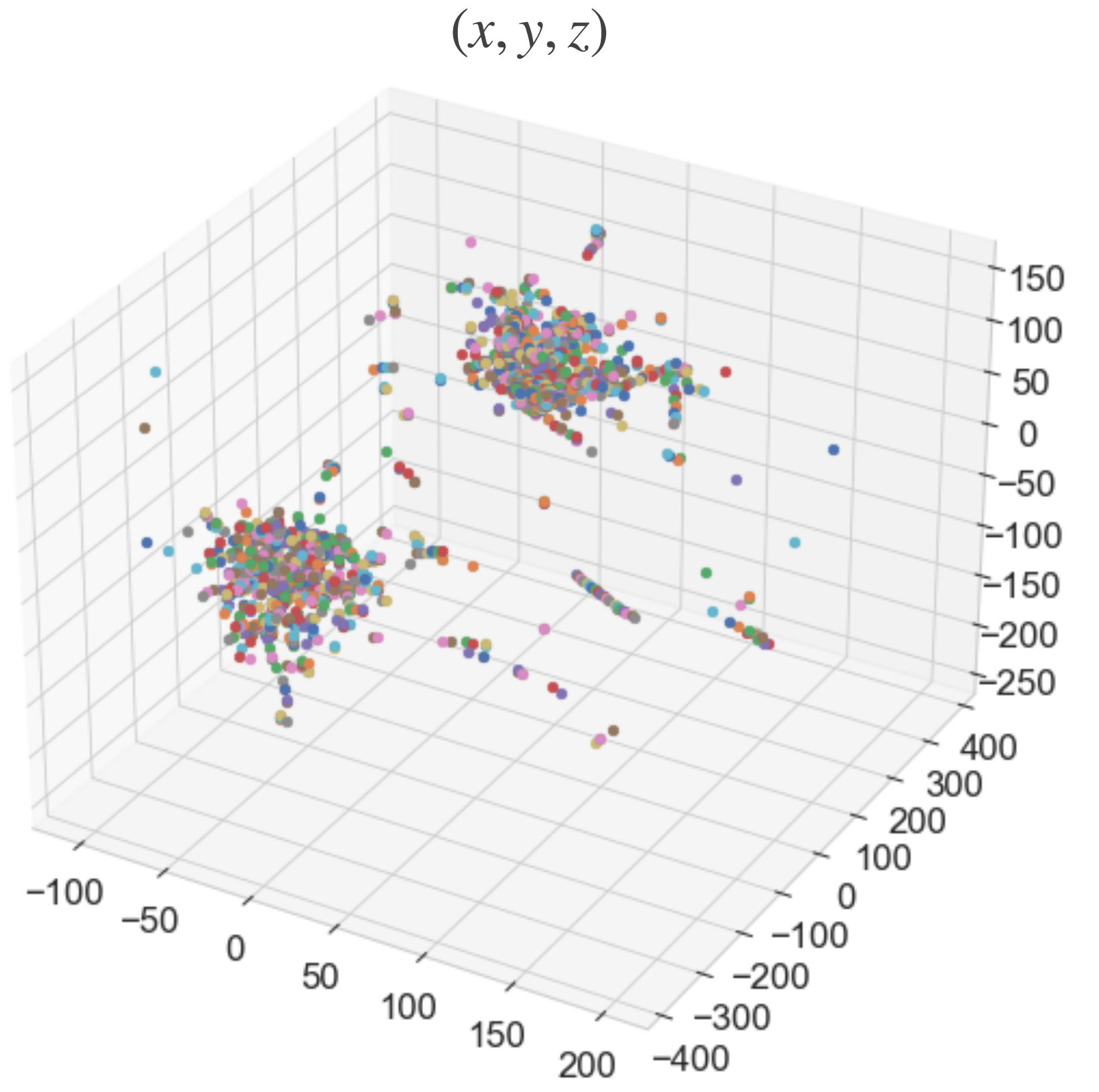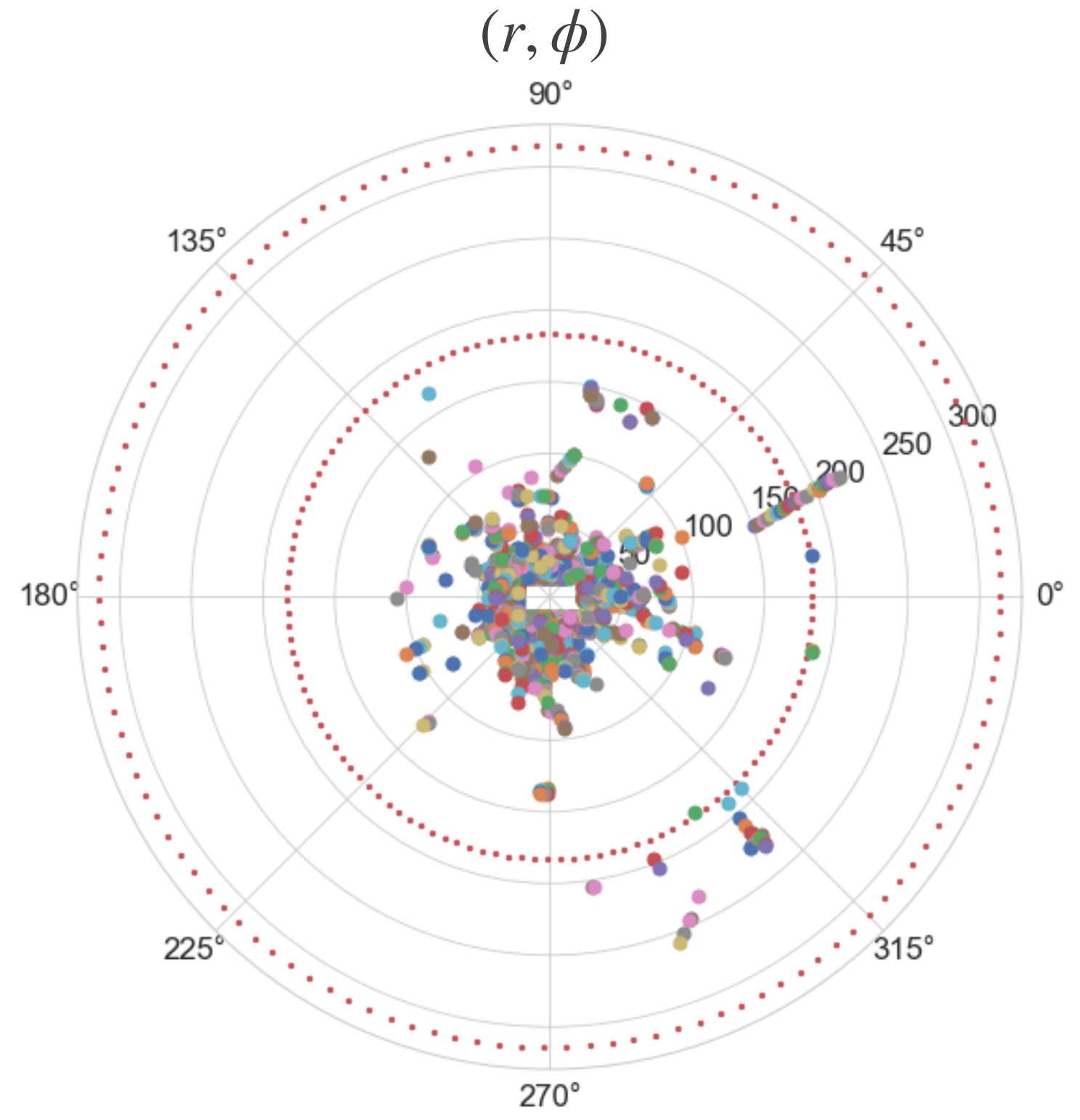
Cross section of the SPD RS

PoS (ACAT) 055

# Coordinate system for clustering

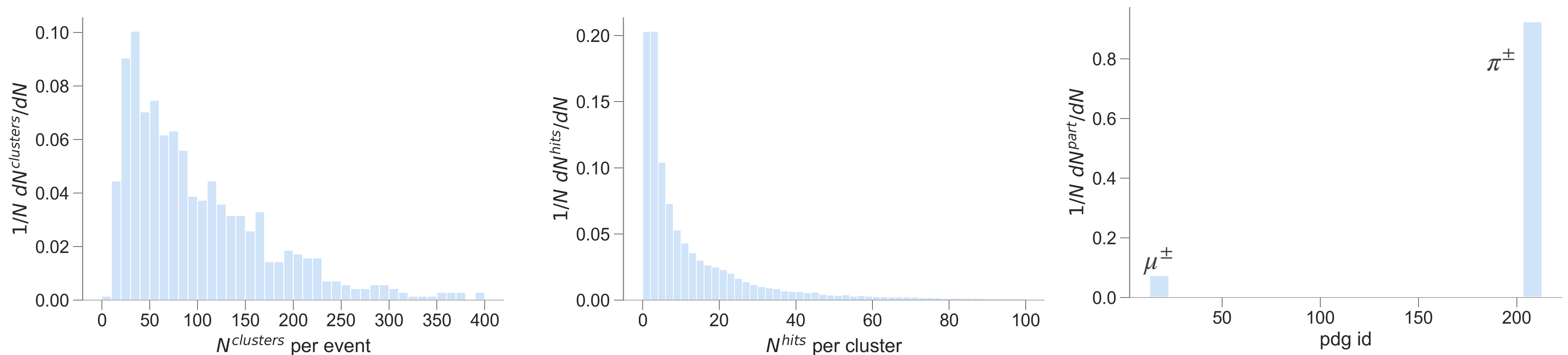A single event represented in various coordinate systems

$(x, y)$

$(r, \phi)$

$(x, y, z)$



Increased density in 2-dimenstion systems due to hit overlapping

# Reference clustering in SpdRoot

Reference clusters are defined as group of hits in Range System produced by a particle in front the RS (thanks Artur for SpdRoot realization).

Should provide ground truth for clustering algorithm evaluation metrics.

- can be up to 400 clusters/event $\Longrightarrow$ high hit density events (mostly in forward region);
- single hits clusters: ~40% clusters have ≤4 hits $\Longrightarrow$ might be a problem for clustering, considered as noise;
- most pion (hadron) clusters are originated in the ECal;



Based on 5k $J/\psi \to \mu\mu$ Pythia8 MC sample from SpdRoot reference example by Igor

# RS Clustering: k-Means

**Clustering** is one of the *unsupervised machine learning* techniques.

<u>Hard clustering</u> (each hit belongs to a cluster or not) algorithms are considered.

Only geometrical informations is available in Range System
(unlike ECAL no energy is provided $\implies$ unable to use *SimpleCone, kT, anti-kT, etc.* algorithms based on the energy deposition in the tower/cell).

**Centroid-based algorithm k-Means:**
- Select a number of clusters to use and randomly initialize their respective center points;
- Each data point is classified by computing the distance between that point and each group center, and then classifying the point to be in the group whose center is closest to it;
- Recompute the group center by taking the mean of all the vectors in the group;
- Repeat these steps for a set number of iterations or until the group centers don't change much between iterations.

**Advantage**: Fast, linear complexity *O(n)*
**Downside**: Need to know number of clusters in advance

*Approximation*: a number of hits in the 1st layer of RS is taken as number of clusters;

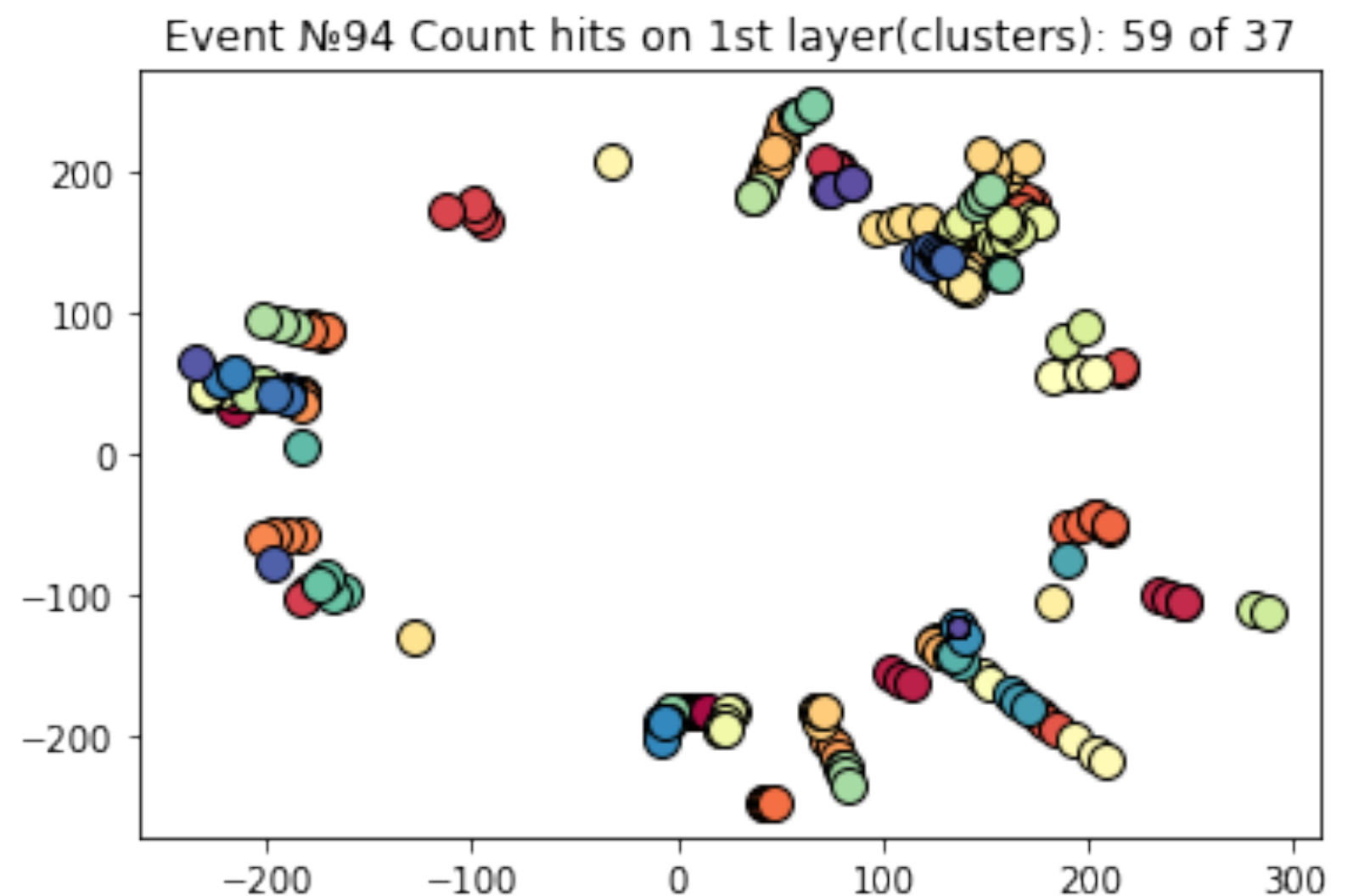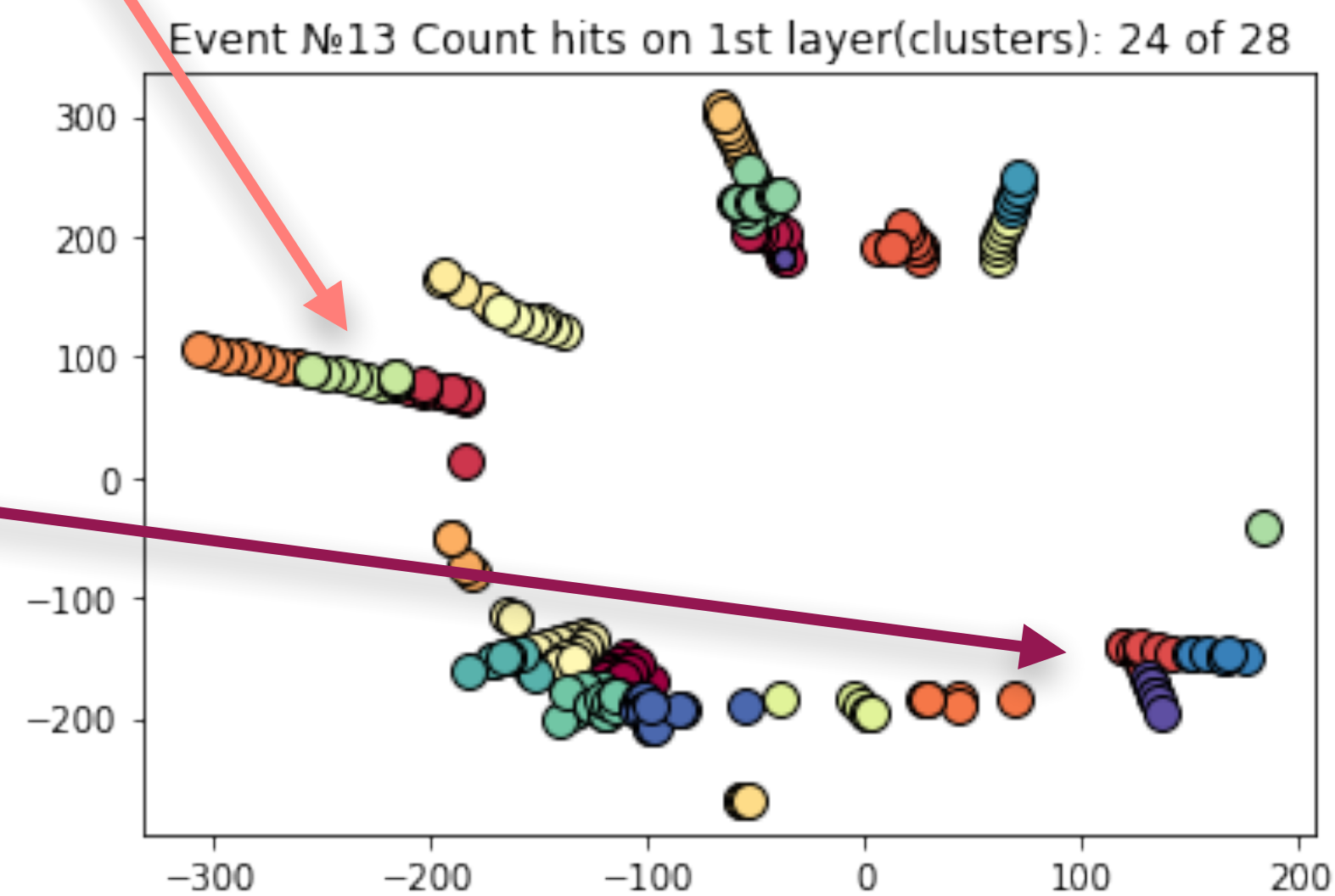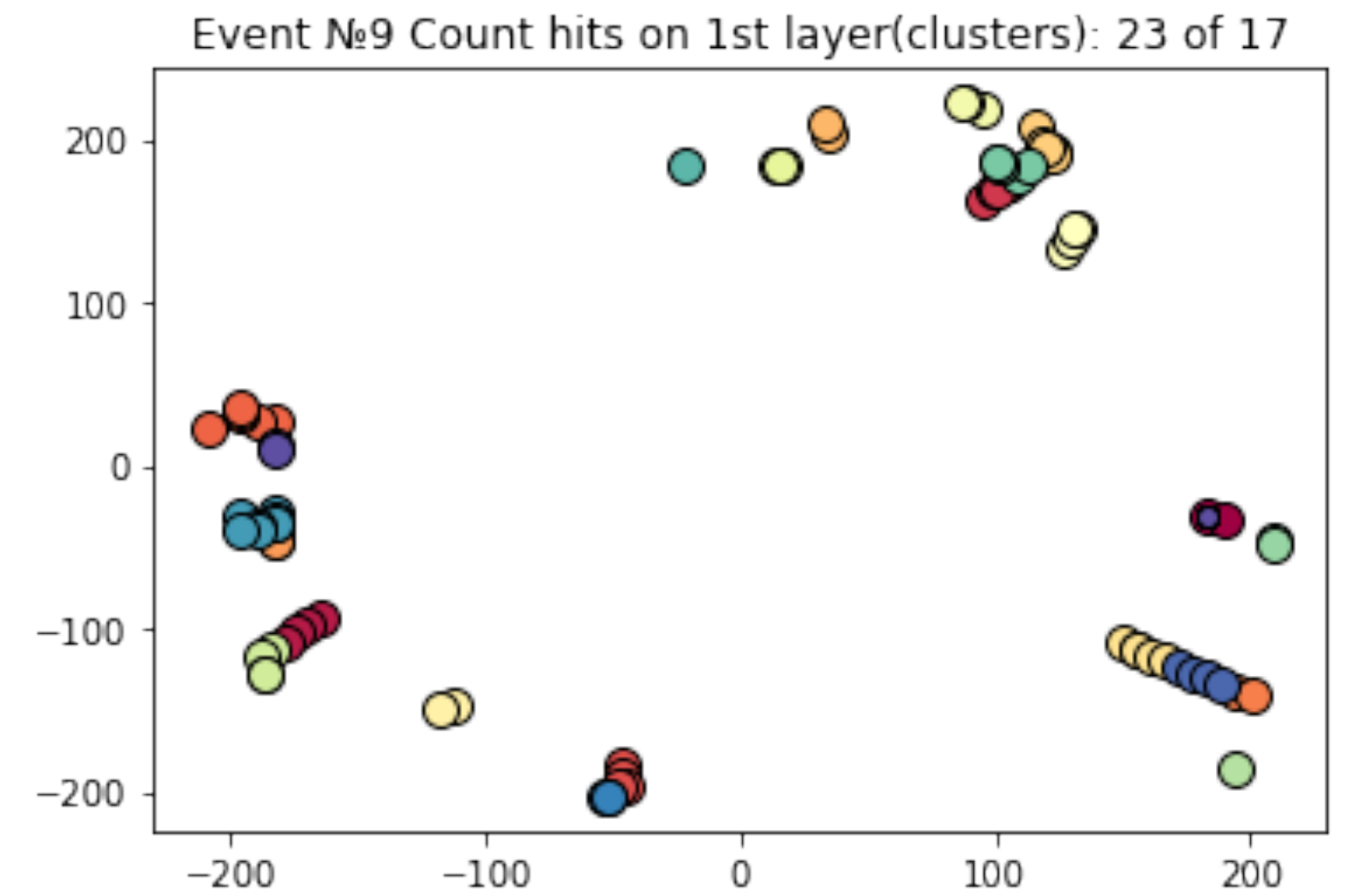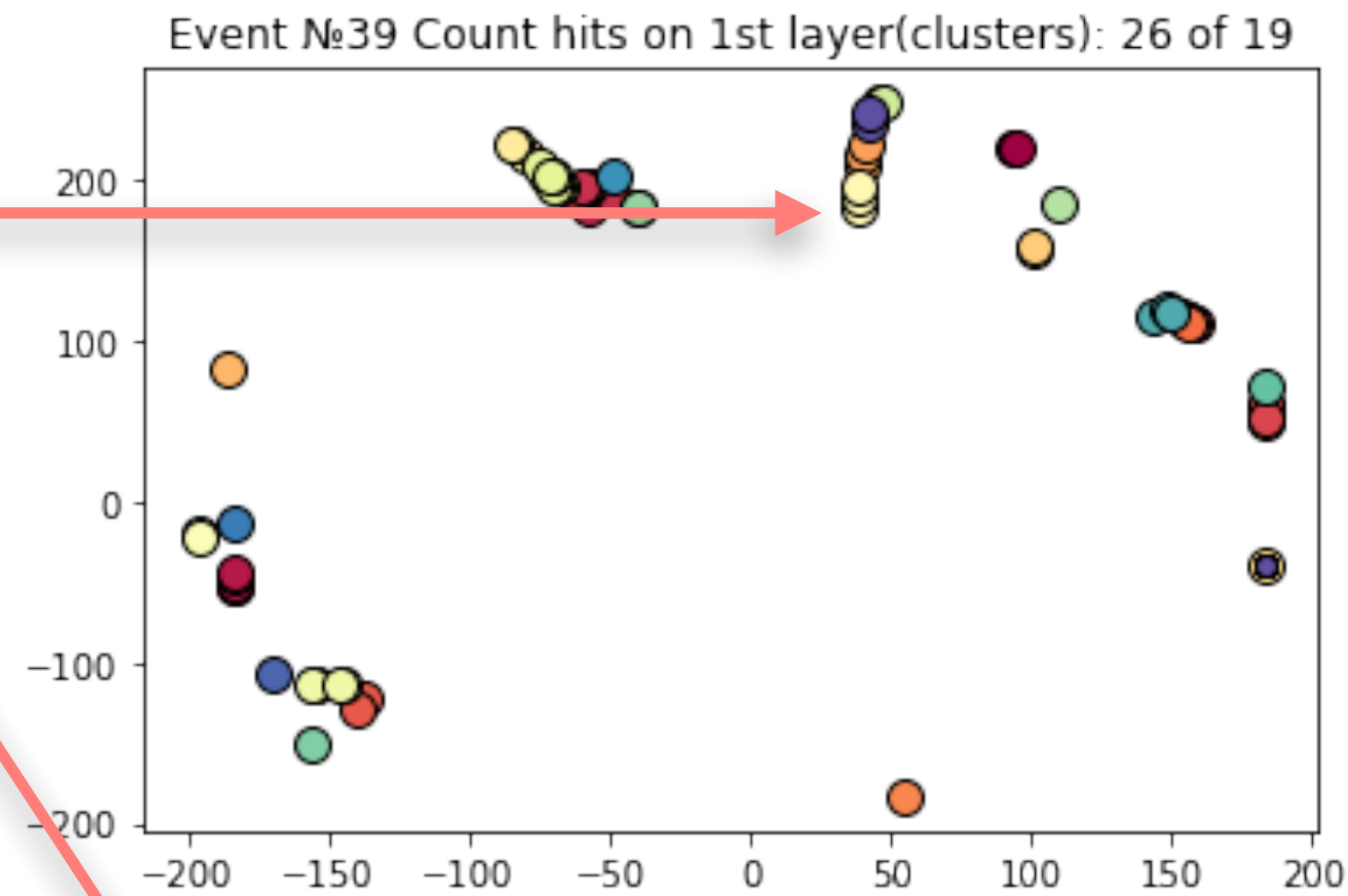This parameter should not be explicitly defined beforehand!

# RS Clustering: k-Means (cont'd)

Splits muon cluster into segments, otherwise includes a lot of single hits

$\implies$ as a result, predicted $N_{clust}$ is overestimated

Various colours show different clusters

Splits hadron showers



Event №39 Count hits on 1st layer(clusters): 26 of 19

Event №9 Count hits on 1st layer(clusters): 23 of 17

Event №13 Count hits on 1st layer(clusters): 24 of 28

Event №94 Count hits on 1st layer(clusters): 59 of 37

*(x,y)* view in the Barrel

# RS Clustering: DBSCAN

**DBSCAN** (Density-based Spatial Clustering Application with Noise)
- views clusters as areas of high density separated by areas of low density;
- clusters found by DBSCAN can be any shape, as opposed to k-Means which assumes that clusters are convex shaped;
- two hyper-parameters: *min_samples* and *epsilon*;
- defines a **core sample** as being a sample in the dataset such that there exist *min_samples* other samples within a distance of *epsilon*, which are defined as neighbors of the core sample.
- A cluster also has a set of **non-core samples**, which are samples that are neighbors of a core sample in the cluster but are not themselves core samples. Hits that cannot be reached are defined as **outliers** and considered as noise;

- Parameter *min_samples* primarily controls how tolerant the algorithm is towards noise (on noisy and large data sets it may be desirable to increase this parameter), usually set to *min_samples>Ndim+1*

- Parameter *epsilon* controls the local neighborhood of the points: optimal *epsilon* is about 2x÷3x distance between layers to take into account missing hits on some layers;
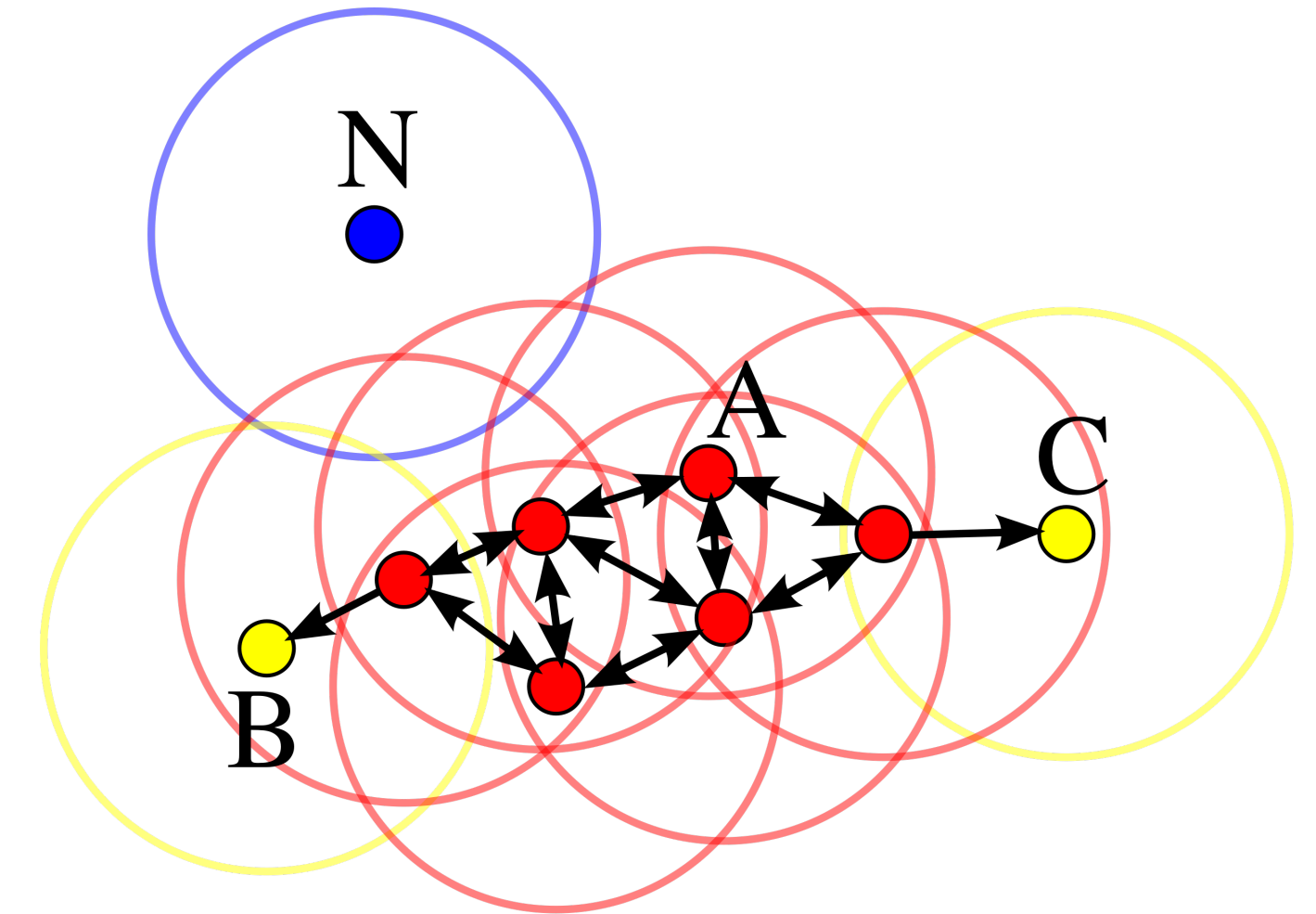
**Advantages:**
- it does not require a pre-set number of clusters at all;
- it also identifies outliers as noise;
- associates hits to clusters of arbitrary shapes;
- naturally exploits 3D coordinates of hits (wire/strips);

**Downsides:**
- it doesn't perform as well as others when the clusters are of varying density
- has two parameters to be optimized.

"A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases with Noise" Ester, M., H. P. Kriegel, J. Sander, and X. Xu, In Proceedings of the 2nd International Conference on Knowledge Discovery and Data Mining, Portland, OR, AAAI Press, pp. 226–231. 1996
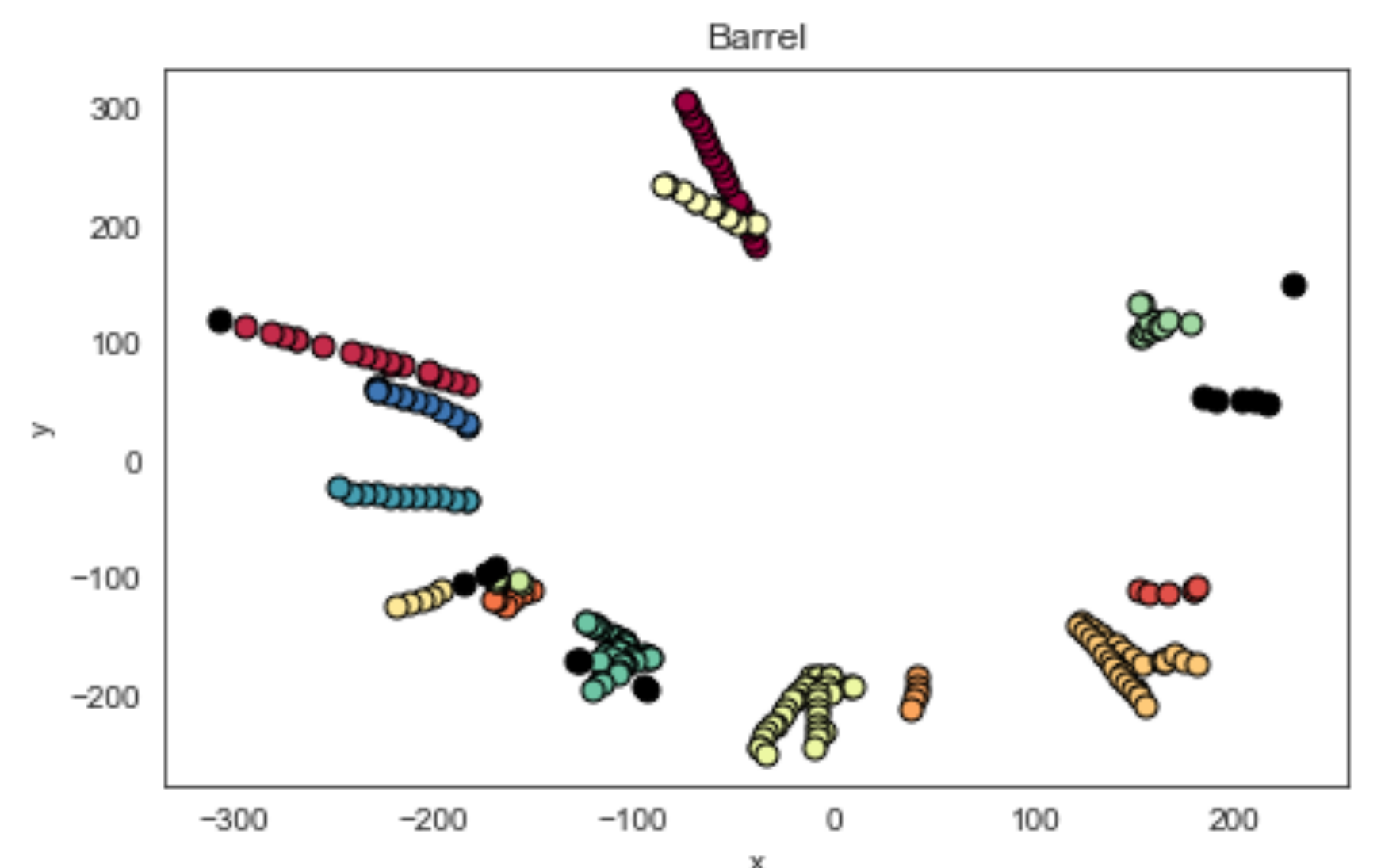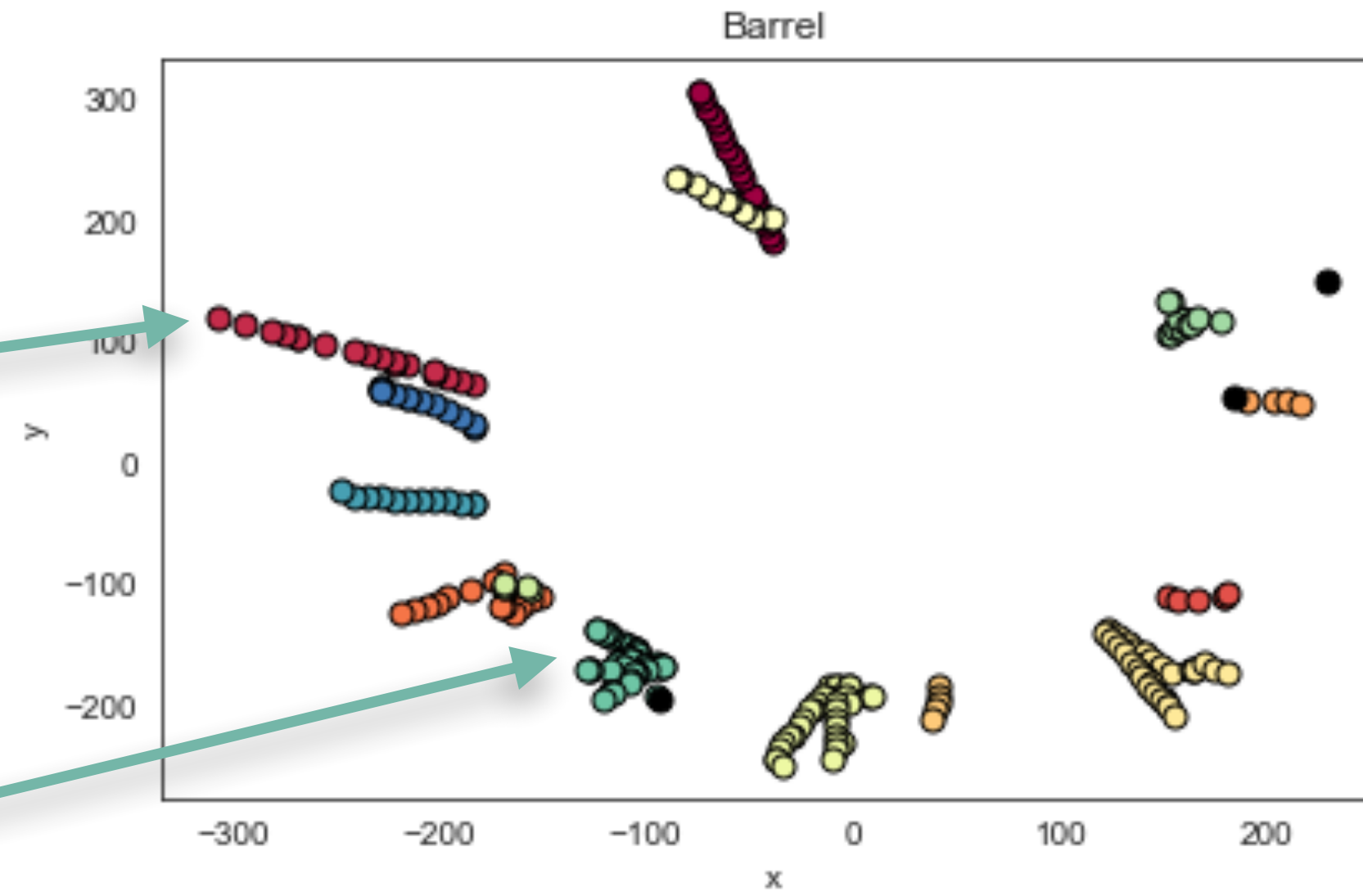


A and red points - core hits;
B,C - not core hits, but reachable from A via red hits;
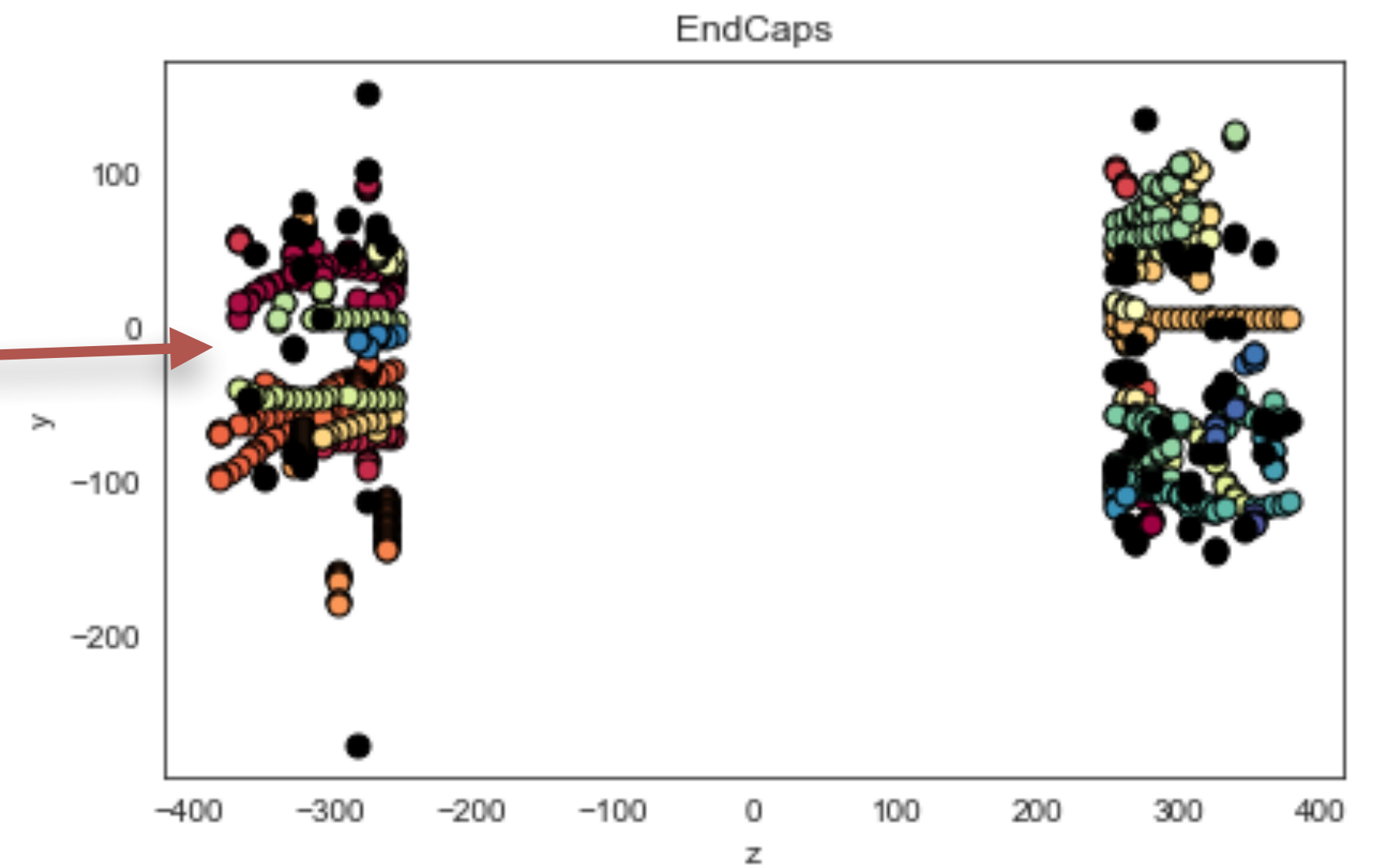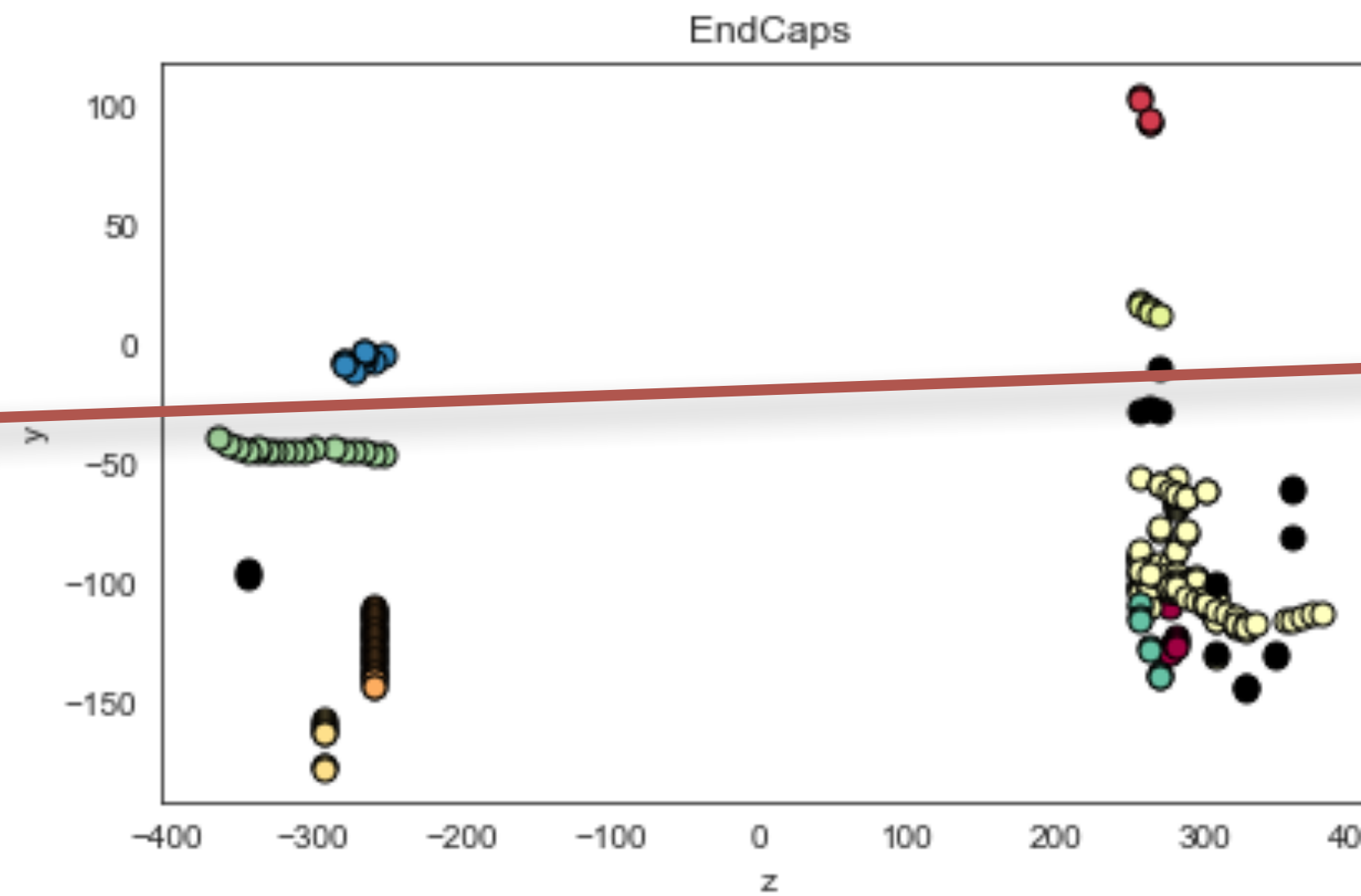N - noise.

Detects muons even with missing hits on some layers

Detects hadron showers

Needs proper hyper-parameter tuning in high density regions

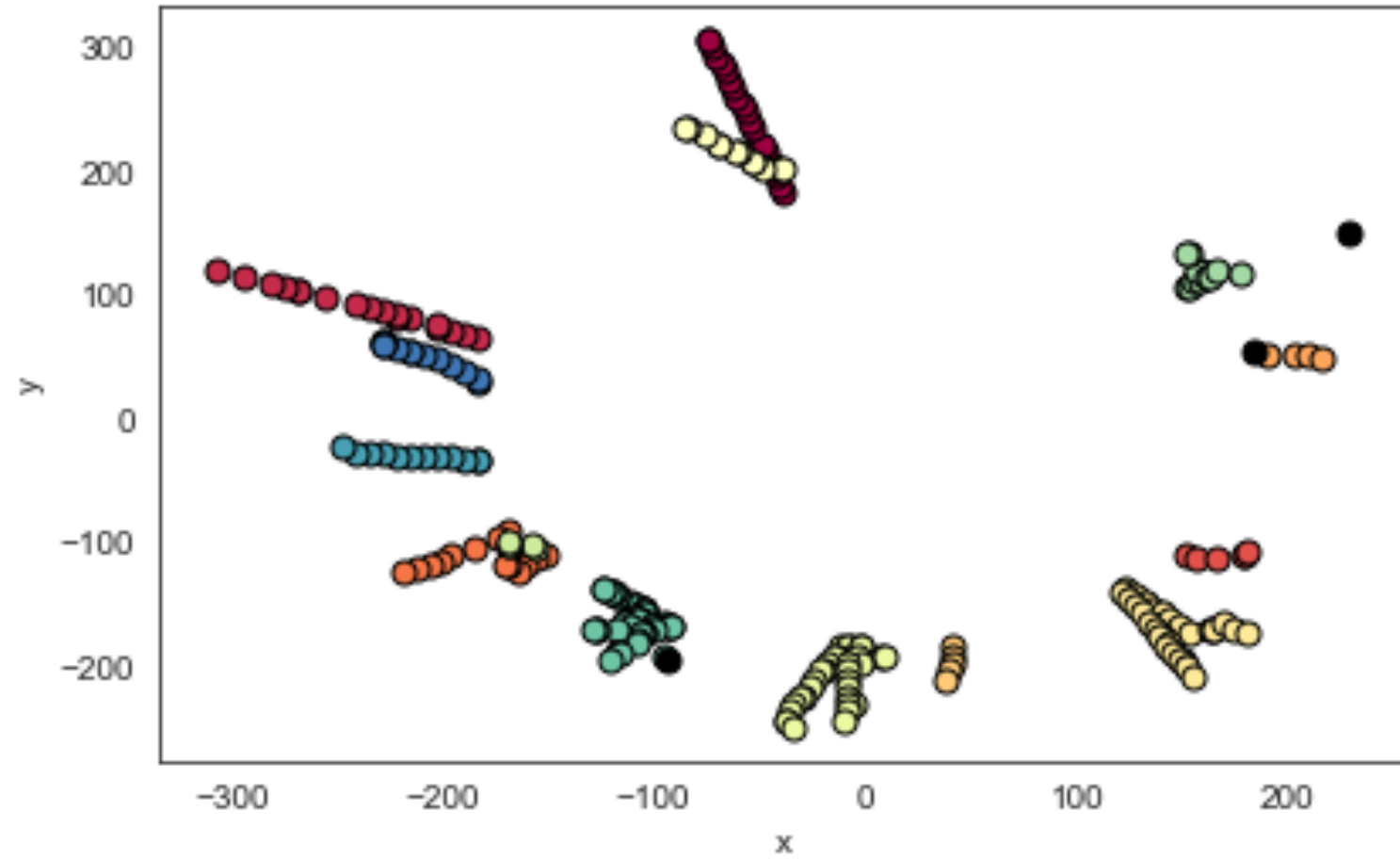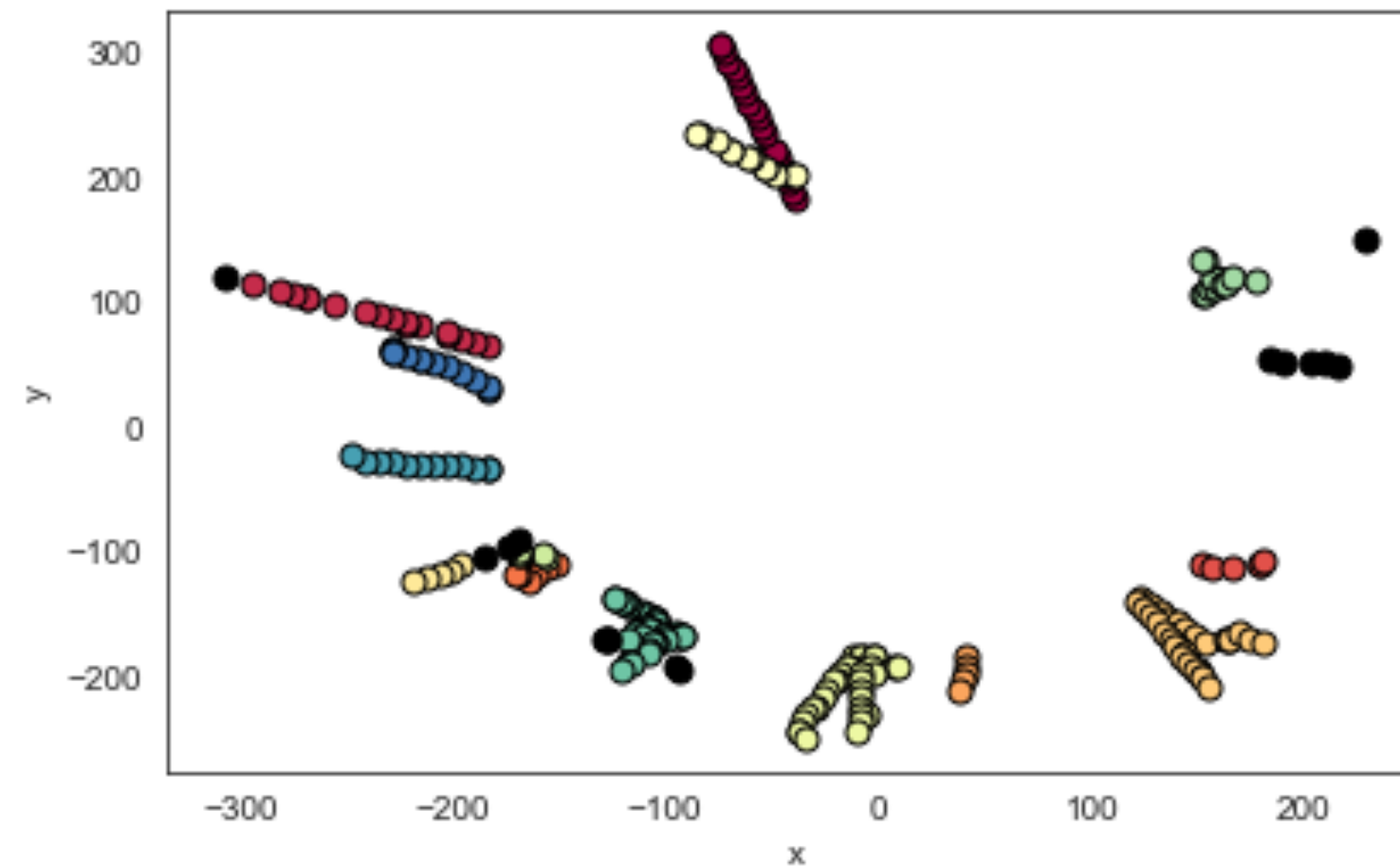# DBSCAN in Barrel and EndCaps

$\pi/8 < \theta < 7\pi/8$

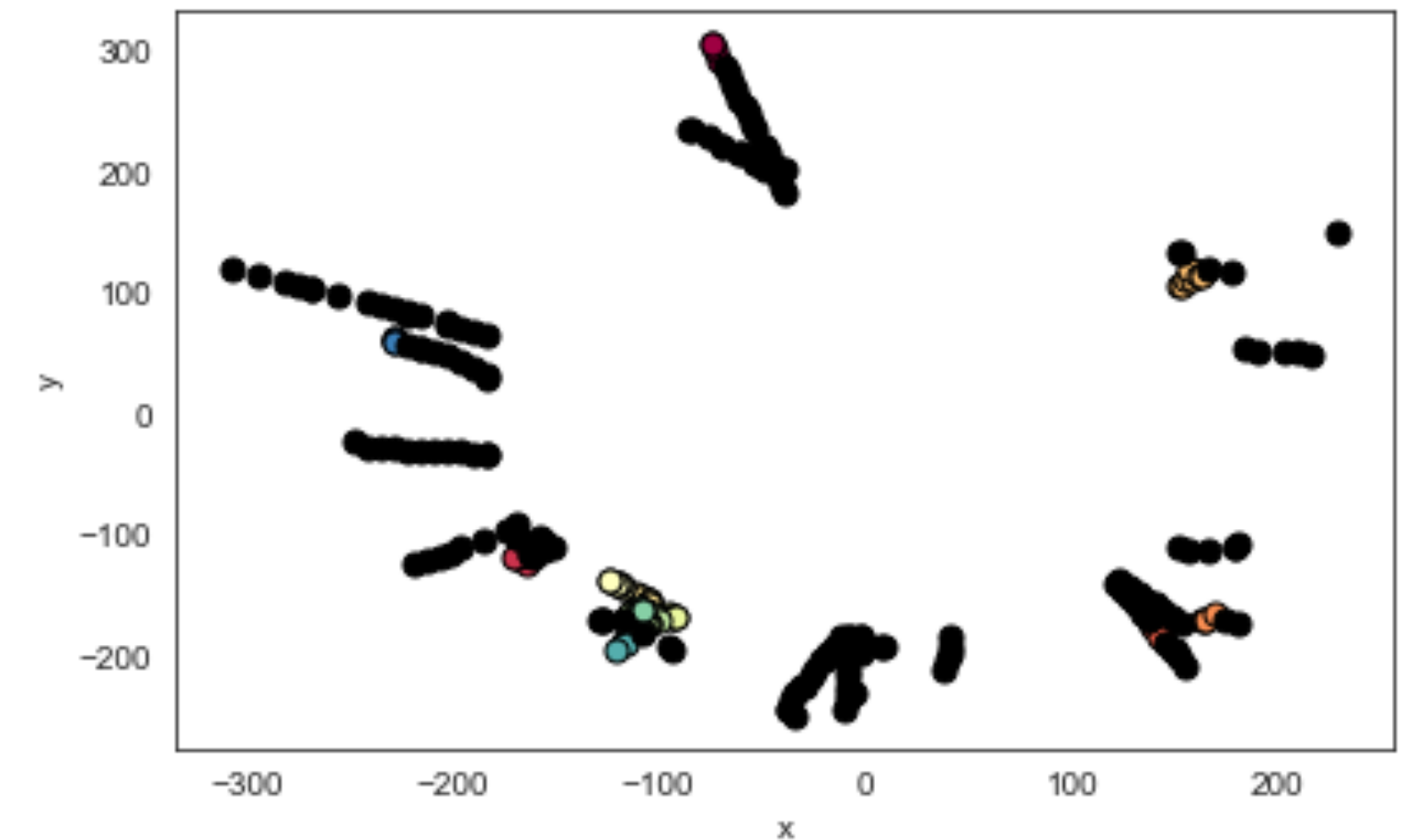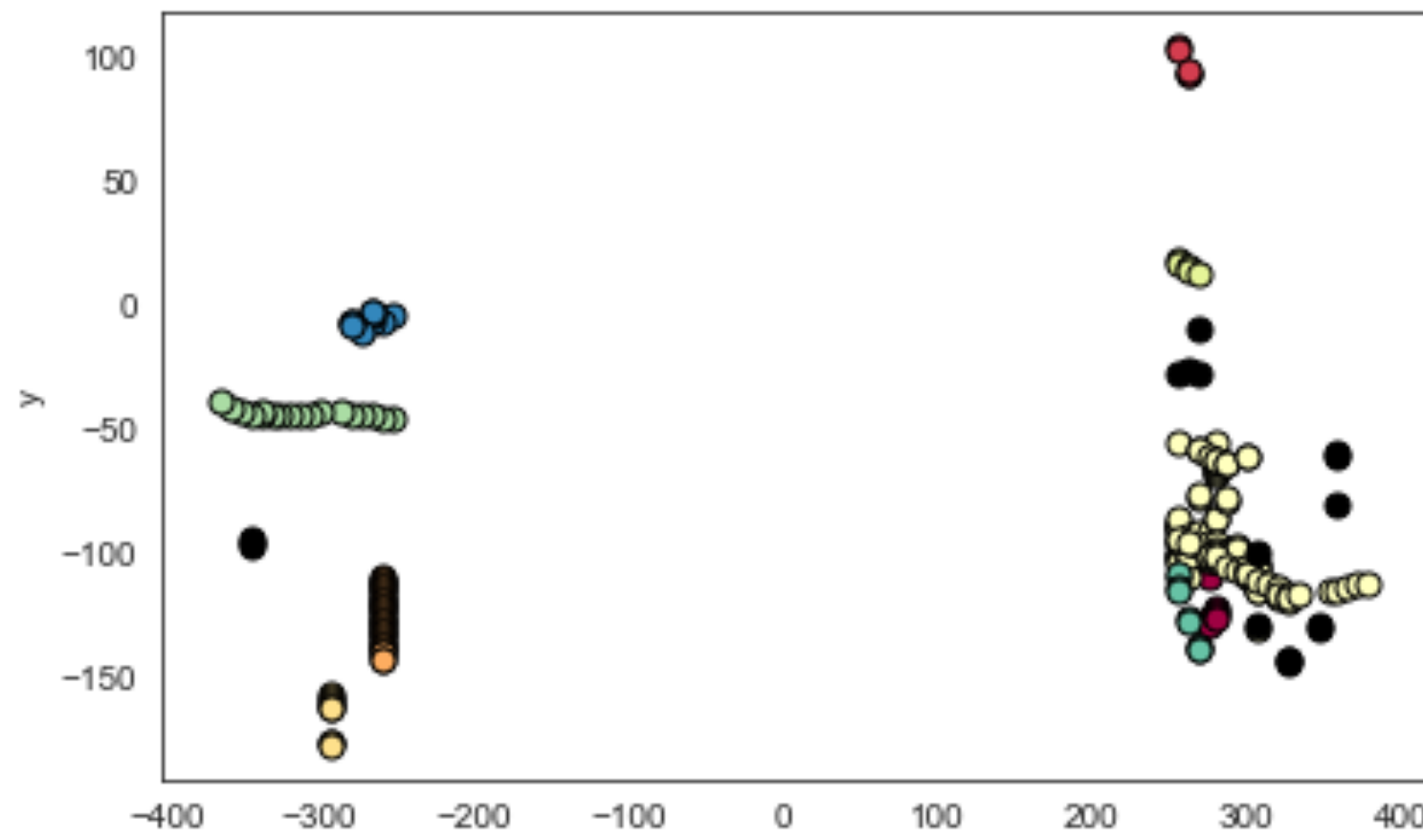$\pi/16 < \theta < 15\pi/16$

all



Igor used $|Cos\theta| < 0.9$ cut in his analysis, similar to this threshold

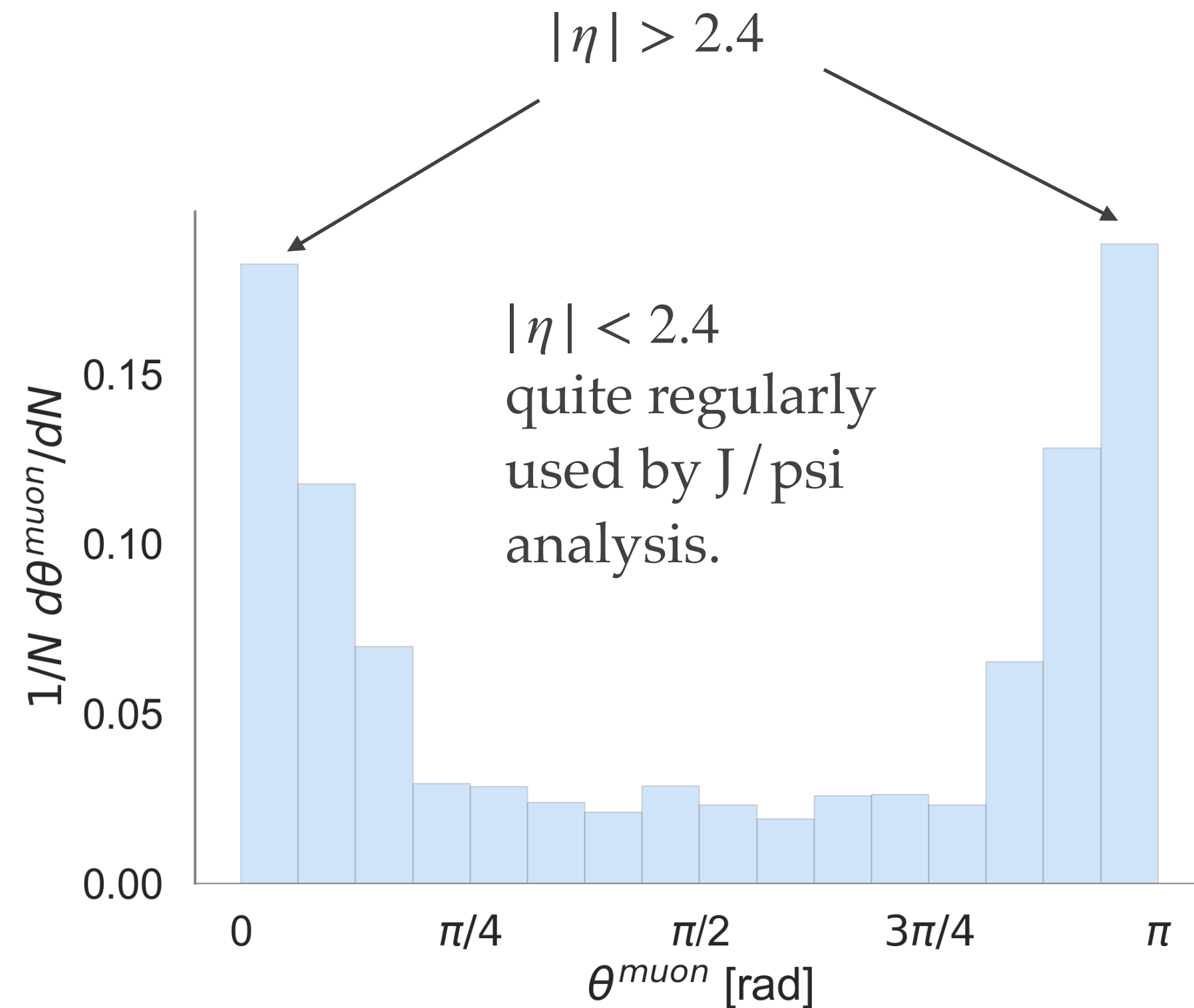Approx. corresponds to $|\eta| < 2.4$ regularly used by J/psi analysis

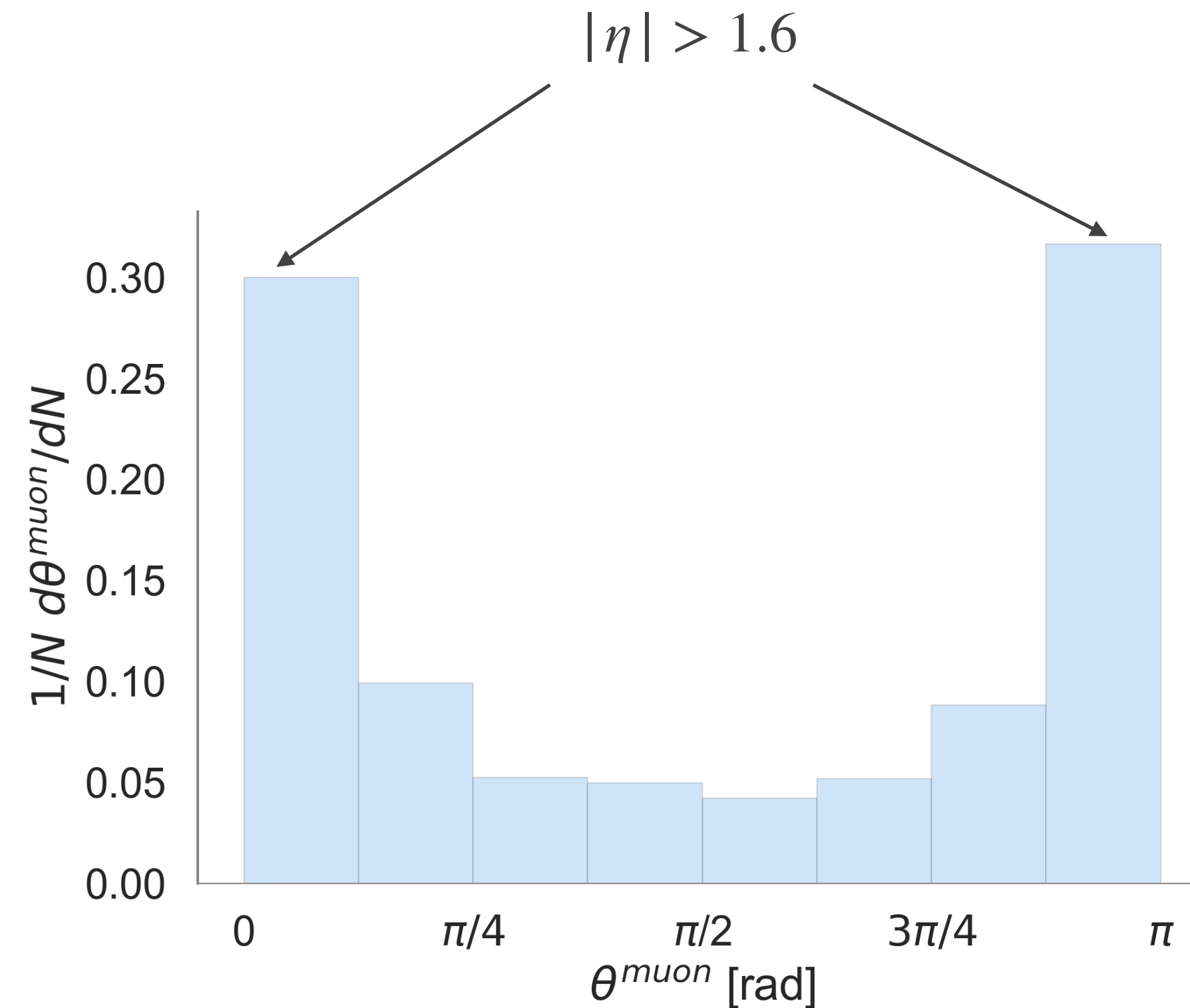Large amount of hits concentrated around the beam pipe leads the algorithm to fail

# $\theta$-Threshold

By applying the threshold on $\theta$ $(\eta)$ polar angle to get rid of dense very forward region (around beam pipe): how much muon candidates (true muon clusters) are we going to lose?

Pseudorapidity:  $\eta = -ln[tan(\theta/2)]$



$|\eta| > 2.4$

$|\eta| < 2.4$ quite regularly used by J/psi analysis.

$|\eta| > 1.6$

About 30% muons loss

About 60% muons loss

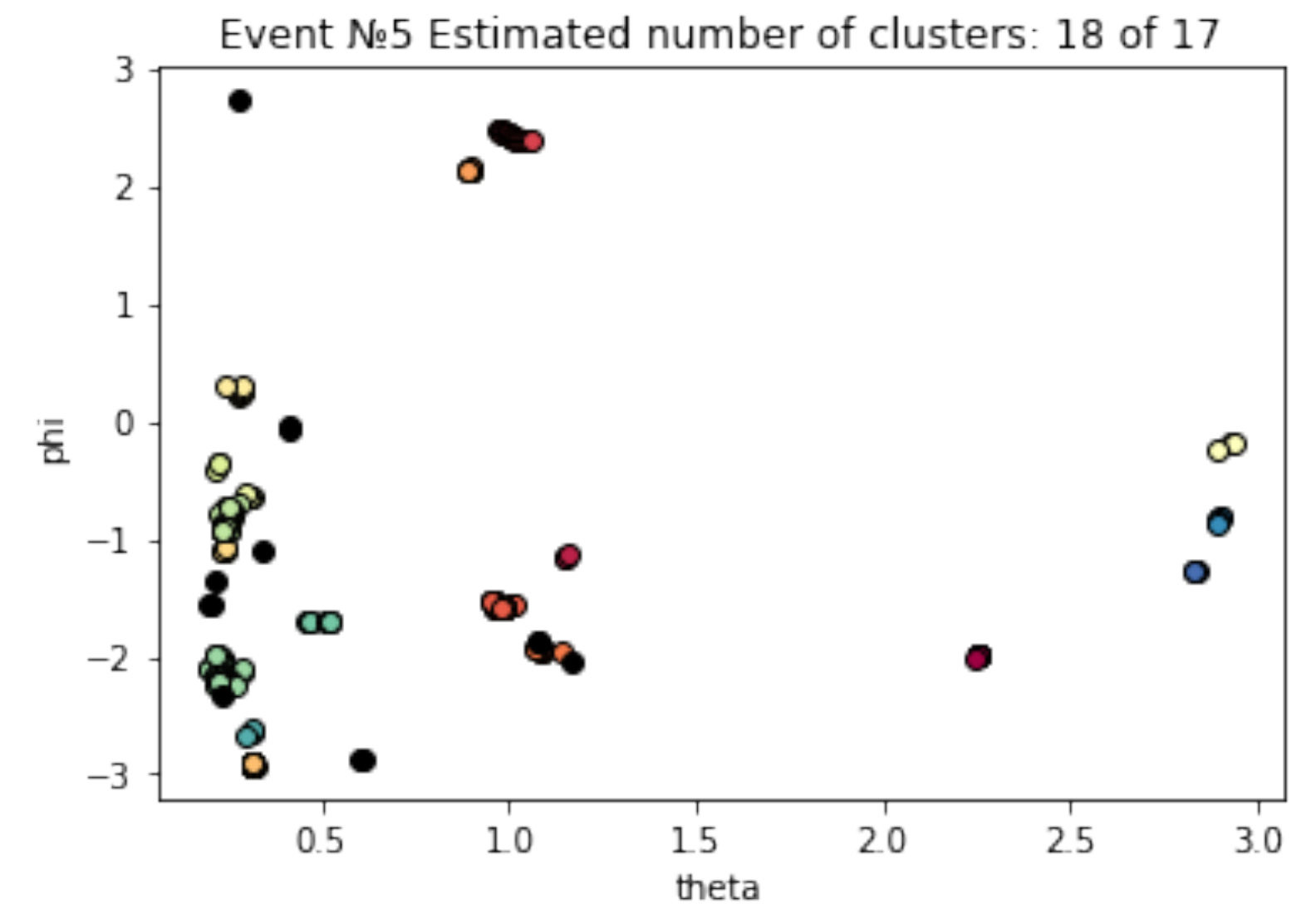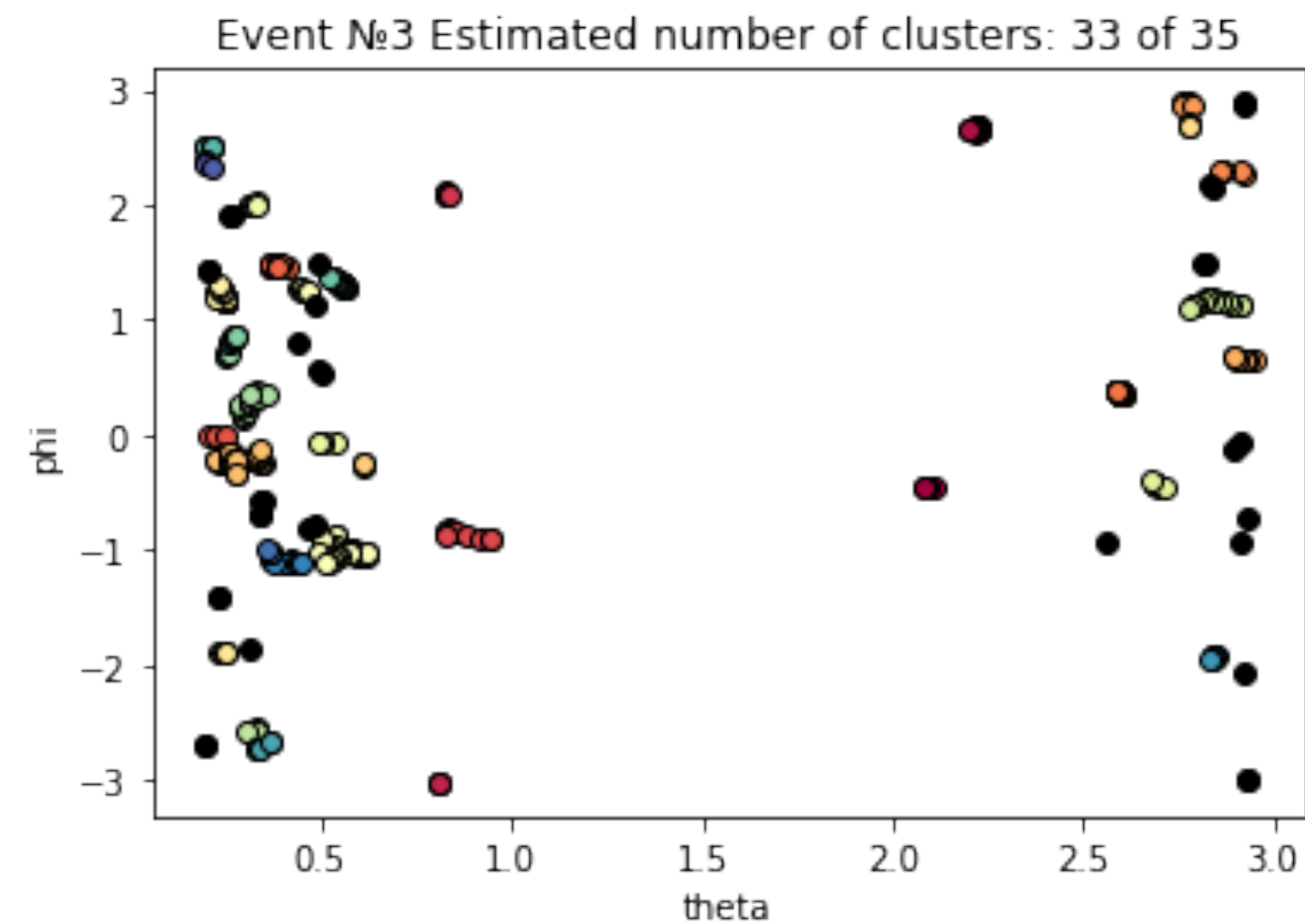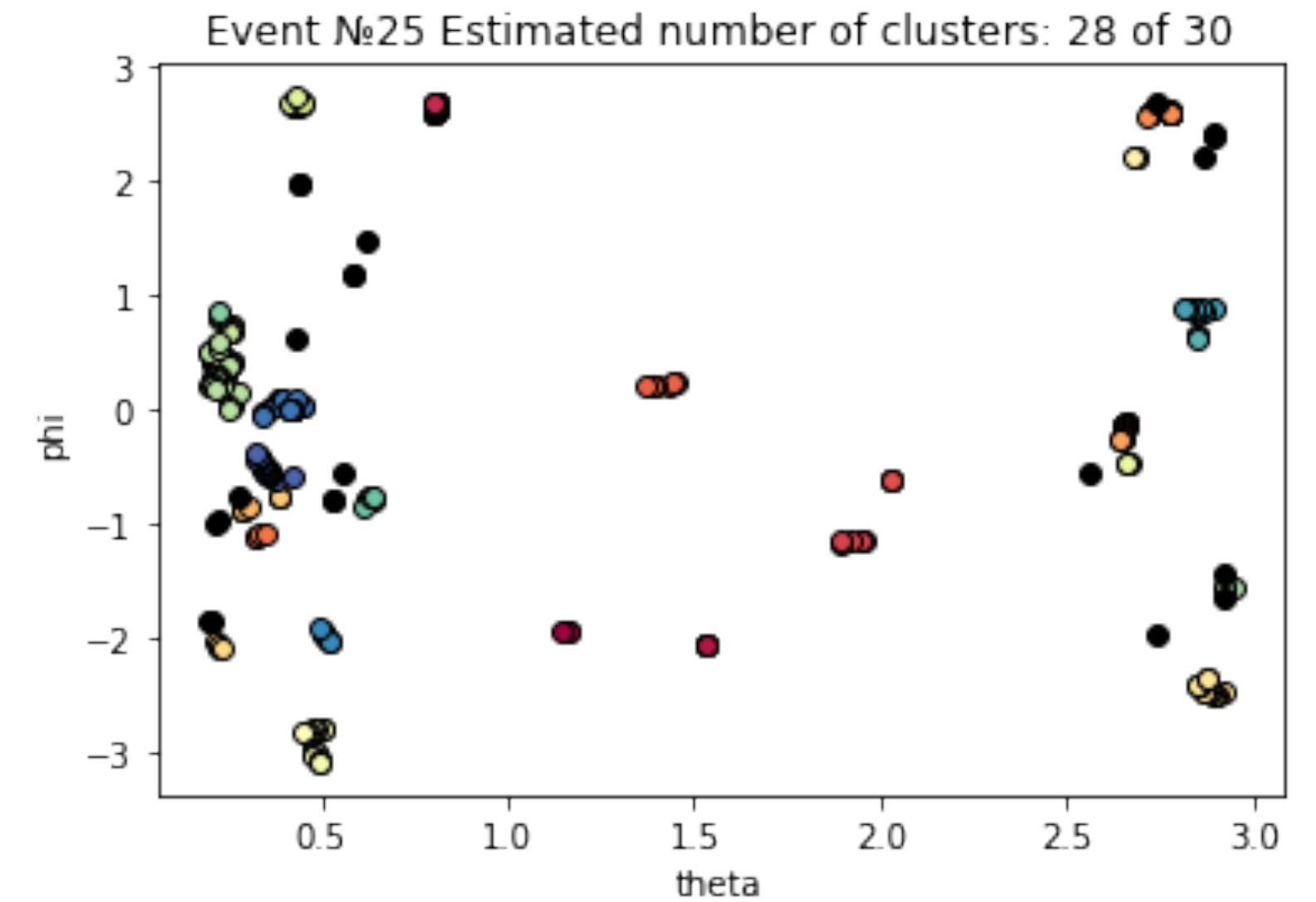Soft, low transverse momenta muons are expected around the beam pipe at the very forward region.
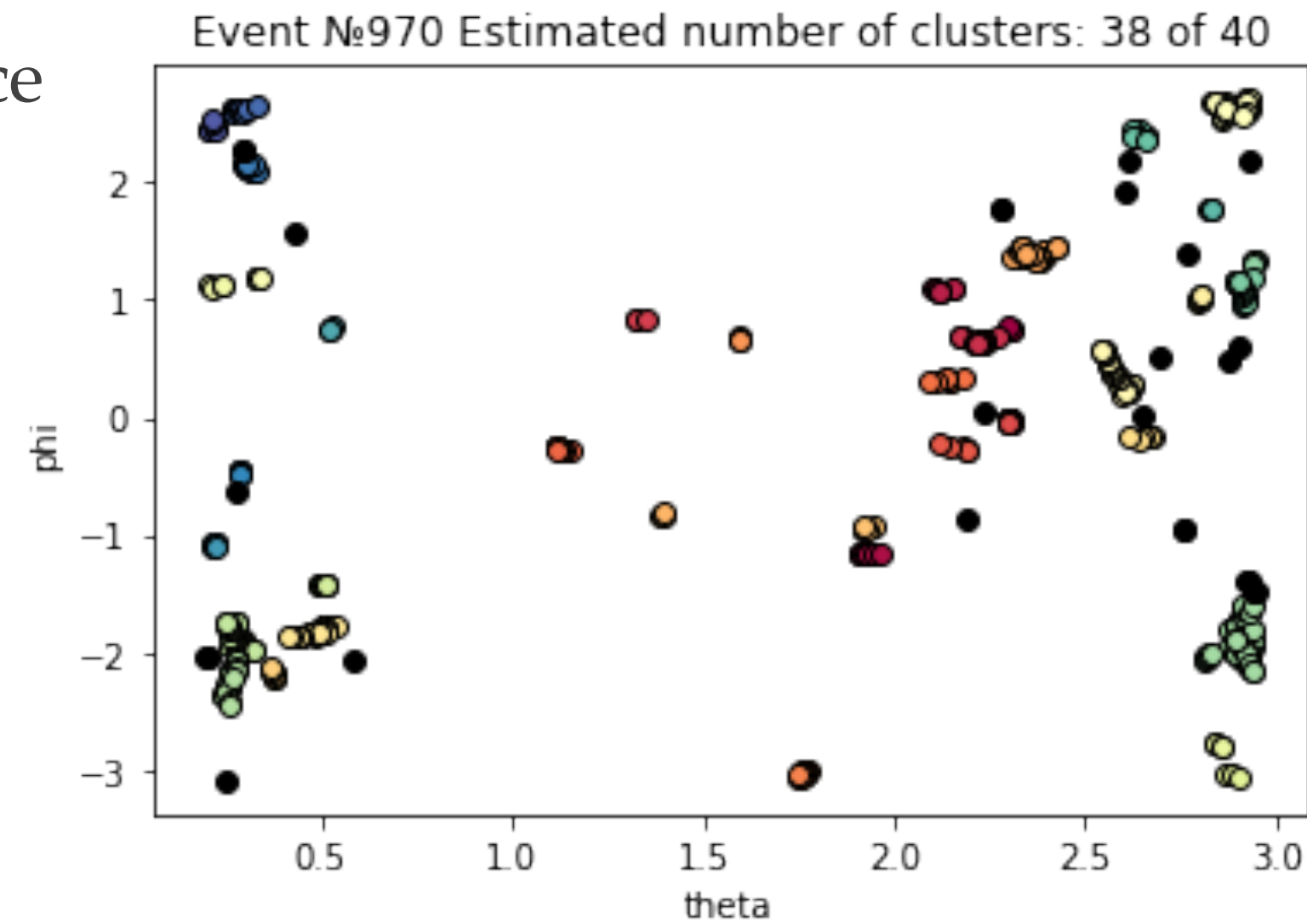
# DBSCAN-based algorithm

DBSCAN application in the $(\theta, \varphi)$ space

$\pi/16 < \theta < 15\pi/16$

Muon and pion clusters look more similar, especially in Barrel

Muons can have some length in the EndCaps

Event №970 Estimated number of clusters: 38 of 40

Event №25 Estimated number of clusters: 28 of 30

Event №3 Estimated number of clusters: 33 of 35

Event №5 Estimated number of clusters: 18 of 17

# Clustering evaluation metrics

**1. Purity:**

$$P = \frac{\sum_i N_{i,\ hits}^{correct}}{N_{hits}^{total}} = \frac{ClusterA + ClusterB + ClusterC}{Total}$$

**Advantages**: straightforward and transparent metrics, easy to calculate;
**Downsides**: increases as the number of cluster increases.

**2. V-measure:**

- **homogeneity**: each cluster contains only members of a single class.
- **completeness**: all members of a given class are assigned to the same cluster.
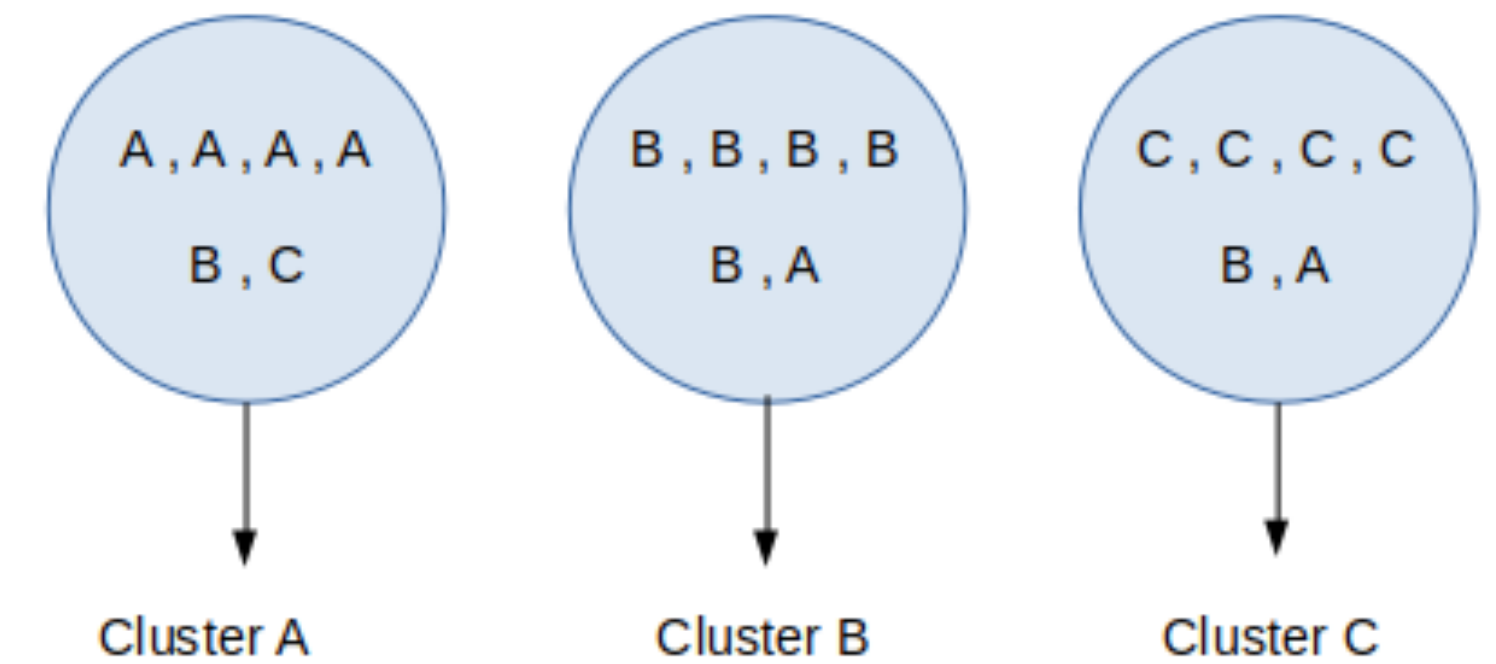
Harmonic mean between the homogeneity and completeness:

$$v = \frac{(1 + \beta) * homogeneity * completeness}{(\beta * homogeneity + completeness)},$$

where by default $\beta = 1$.

**Advantages**: normalized [0, 1]; can be used to compare different clustering models that have different number of clusters;
**Downsides**: random labelling won't yield zero scores especially when the number of clusters is large

A , A , A , A
B , C
Cluster A

B , B , B , B
B , A
Cluster B

C , C , C , C
B , A
Cluster C

There are a lot of other, more sophisticated metrics.
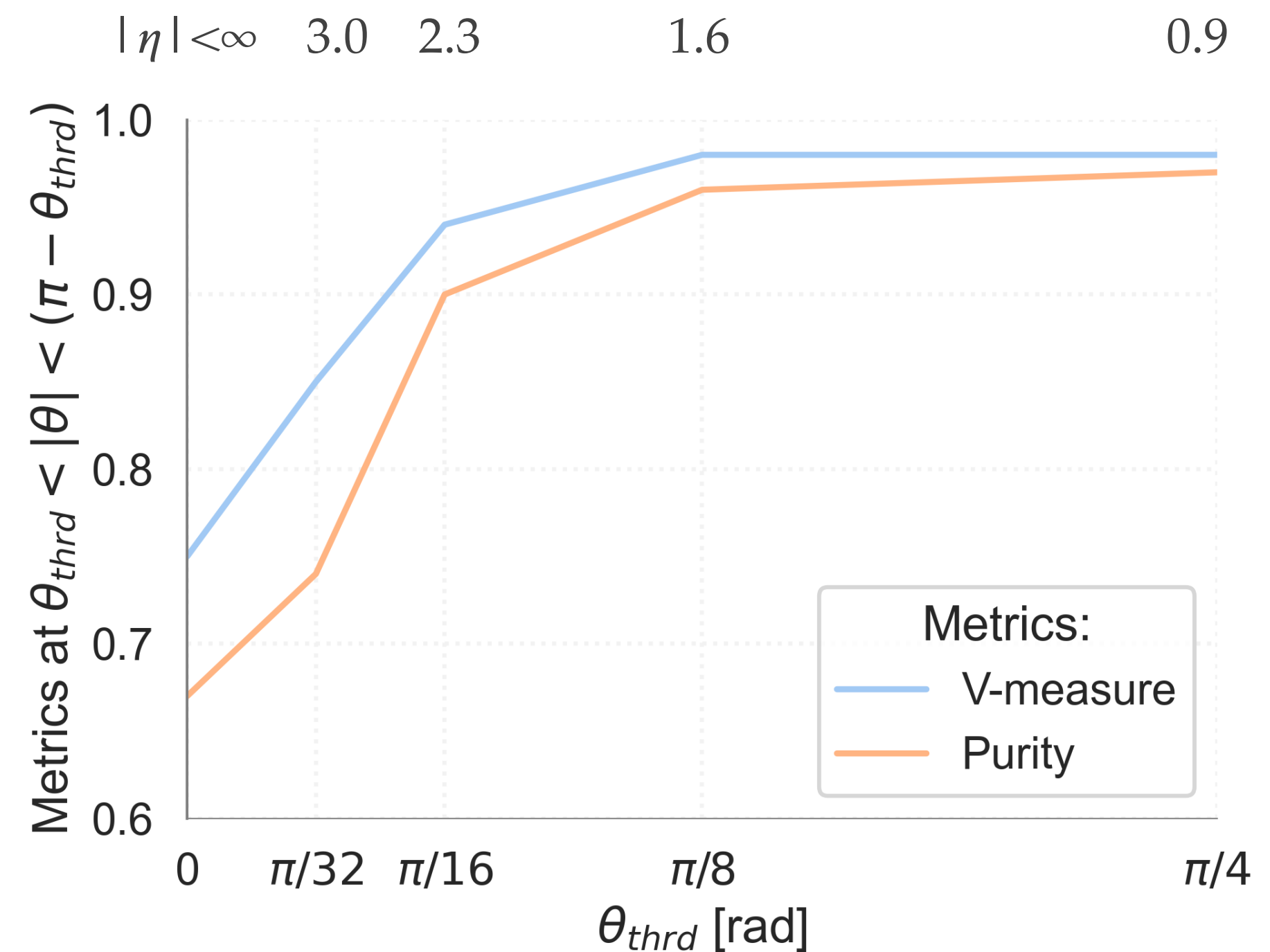Will be considered later.

# DBSCAN evaluation metrics

Using the ground truth from reference clustering one can evaluate the performance of the algorithm

Purity of the clusters is at the level of 90% for $\theta_{thrd} = \pi/16$ (approximately corresponds to $|\eta| < 2.4$ requirement) and up to 96% for $|\eta| < 1.6$;

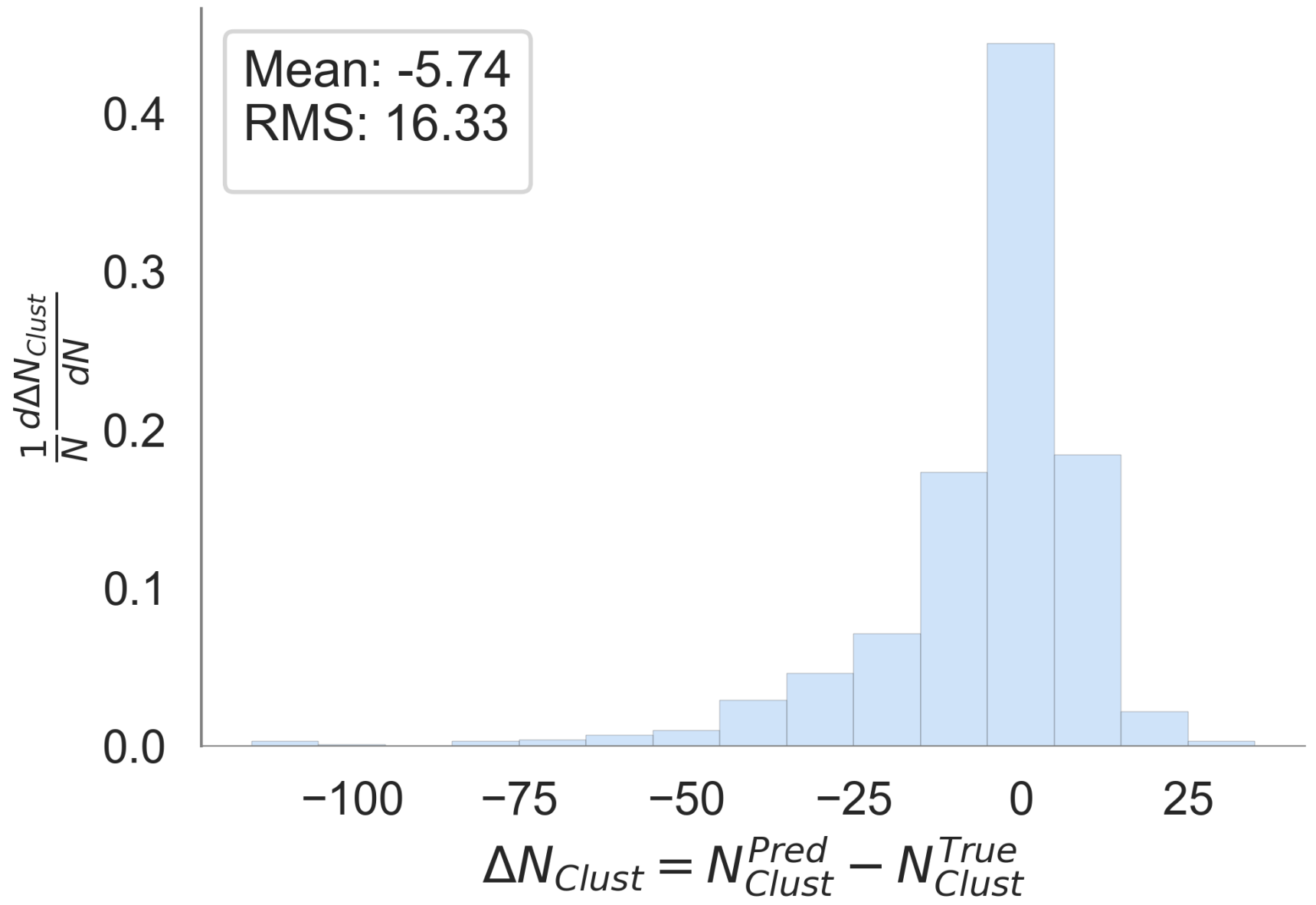Quite ordinary performance with no threshold requirement due to high hit density in very forward region;

Similar values for V-measure metrics;



DBSCAN evaluation metrics applied in *(x,y,z)* coordinates, as a function of polar angle threshold.
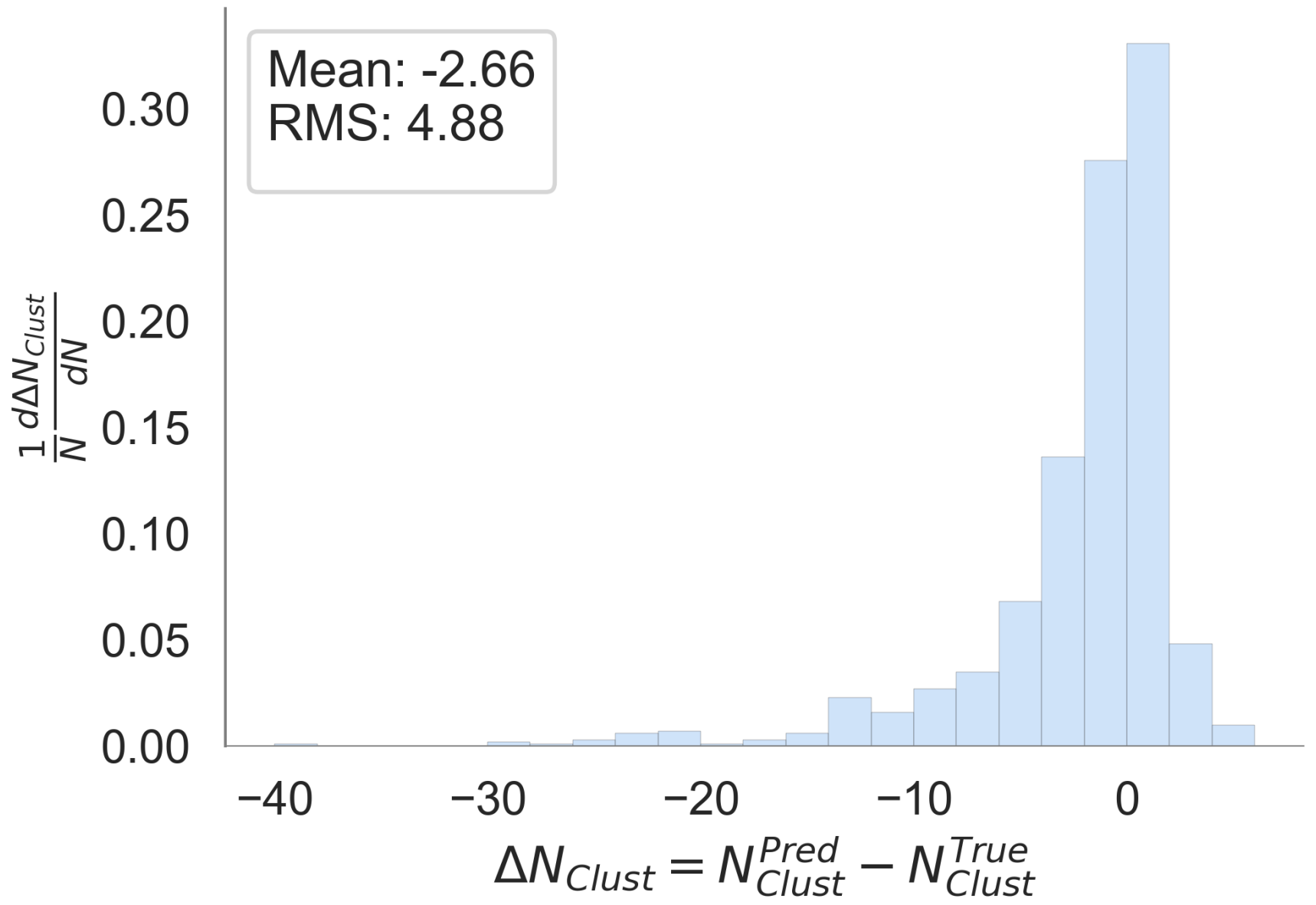
# Clustering efficiency

all

$$\pi/16 < \theta < 15\pi/16$$

Mean: -5.74
RMS: 16.33

$\frac{1}{N}\frac{d\Delta N_{Clust}}{dN}$

$$\Delta N_{Clust} = N_{Clust}^{Pred} - N_{Clust}^{True}$$

Mean: -2.66
RMS: 4.88

$\frac{1}{N}\frac{d\Delta N_{Clust}}{dN}$

$$\Delta N_{Clust} = N_{Clust}^{Pred} - N_{Clust}^{True}$$

Relative error: $\dfrac{\Delta N_{Clust}}{N_{Clust}} \sim 18\,\%$

Relative error: $\dfrac{\Delta N_{Clust}}{N_{Clust}} \sim 15\,\%$

# Clustering algorithms performance

| Algorithm | Threshold | Purity | V-measure | $N_{clust}$ relative error (%) |
|:---:|:---:|:---:|:---:|:---:|
| K-Means | — | 0.76 | 0.76 | 0.41 |
| DBSCAN $(\theta,\varphi)$ | — | 0.58 | 0.71 | 0.37 |
| DBSCAN $(\theta,\varphi)$ | $|\eta| < 2.4$ | 0.89 | 0.93 | 0.13 |
| DBSCAN $(\theta,\varphi)$ | $|\eta| < 1.6$ | 0.96 | 0.97 | 0.08 |
| DBSCAN $(x,y,z)$ | — | 0.67 | 0.75 | 0.18 |
| DBSCAN $(x,y,z)$ | $|\eta| < 2.4$ | 0.94 | 0.94 | 0.15 |
| DBSCAN $(x,y,z)$ | $|\eta| < 1.6$ | 0.96 | 0.98 | 0.09 |

# Conclusion

First look made at the centroid- and density-based clustering algorithms application to the SPD RS:

- k-Means clustering algorithm was tested and showed relatively poor performance for RS hit clustering.
- Density-based algorithm (DBSCAN) is found to be more promising to that end;
- One shows relatively good performance in $(\eta, \phi)$ and 3-dimensional $(x,y,z)$ coordinate frames;
- Additional hyper-parameter optimization and fine tuning is needed to improve performance.

The DBSCAN-based algorithm might be useful for Online Filter tasks;

As a product we need to organize a pipeline {clustering + $\pi/\mu$ separation algos} and evaluate combined metrics focused on muons.

DBSCAN implementation in C++ is ~200 lines of code; no complex external libraries dependences are needed; can be quite easily incorporated within SpdRoot.

There is a faster version DBSCAN++ (additional research is needed);
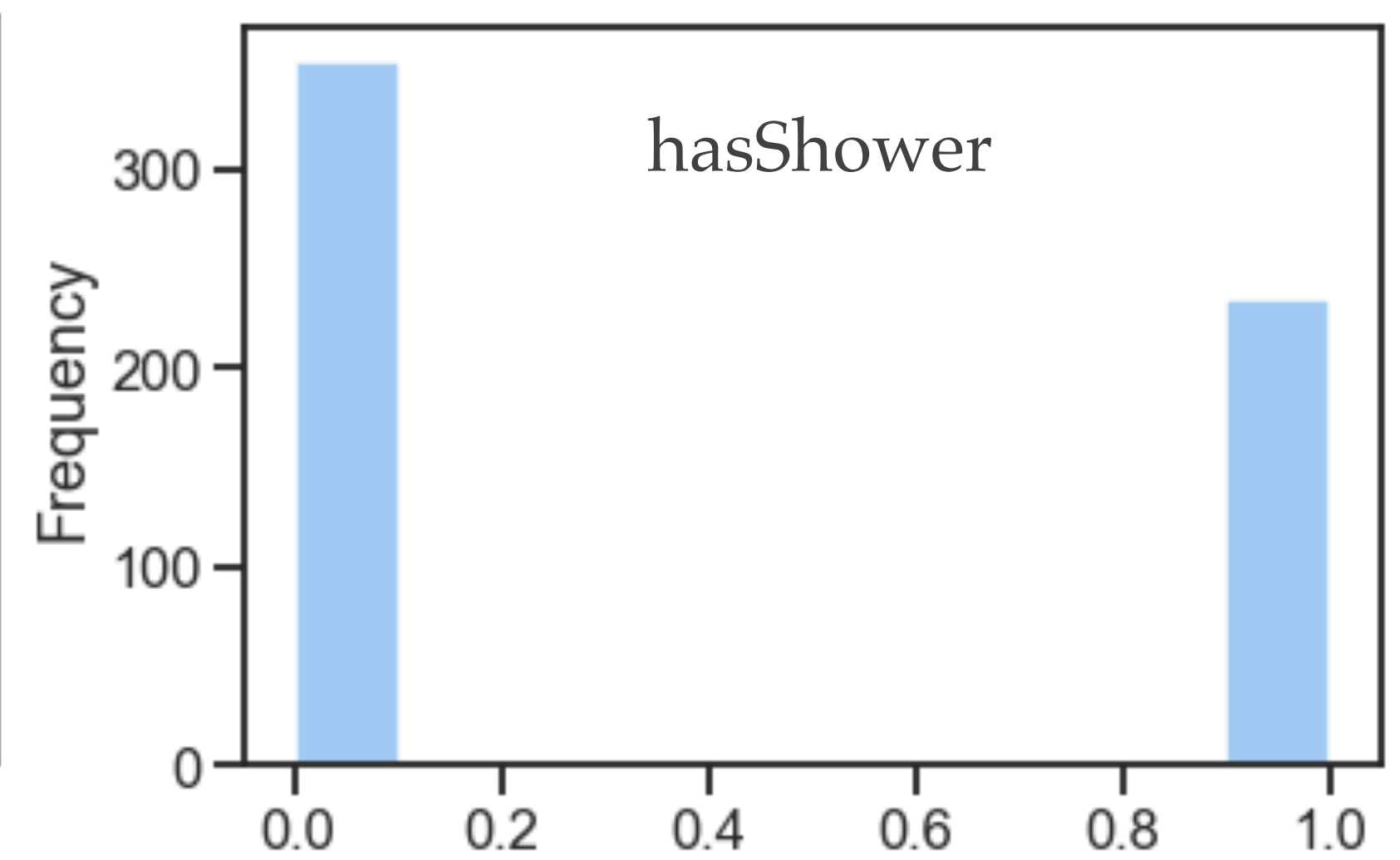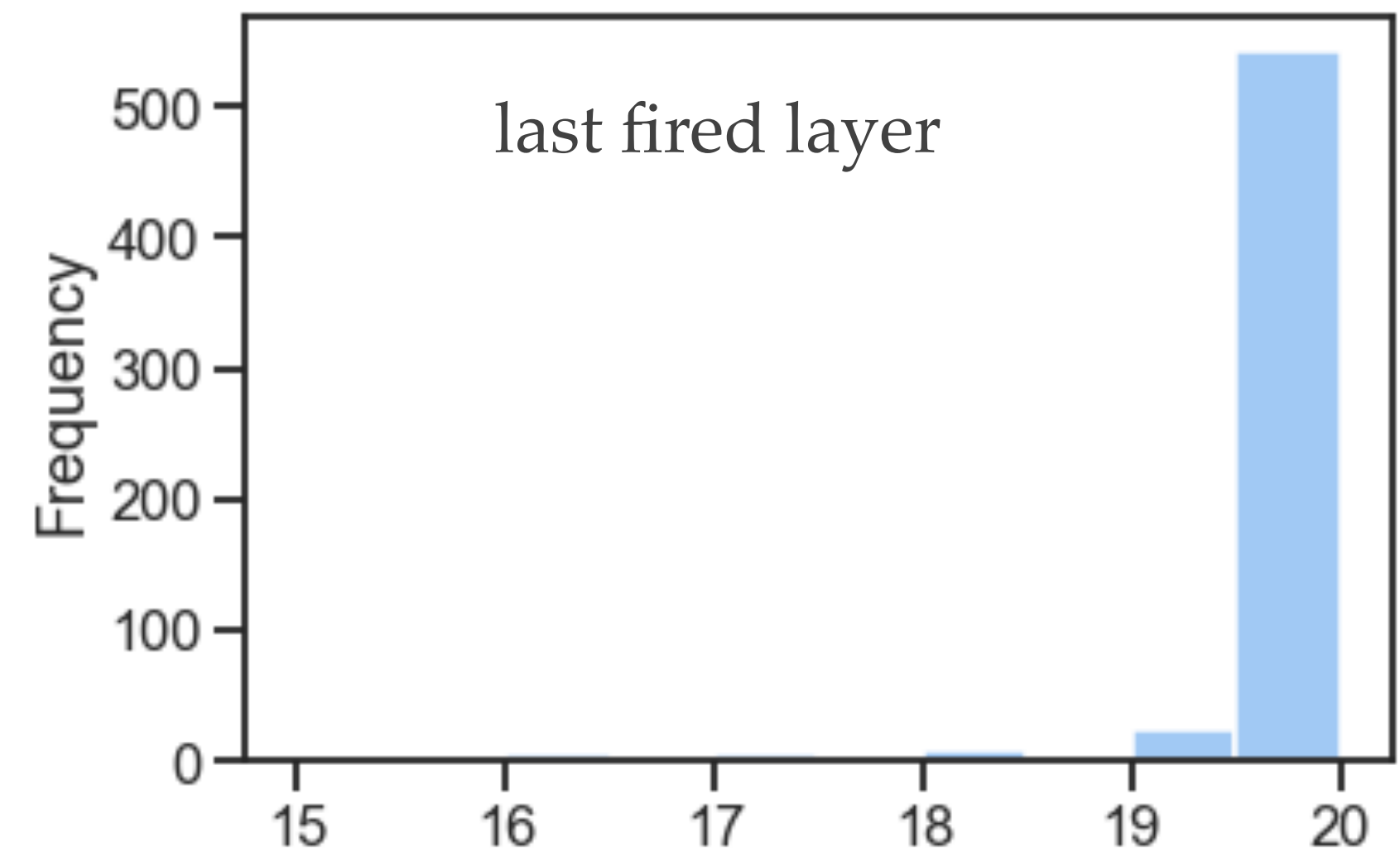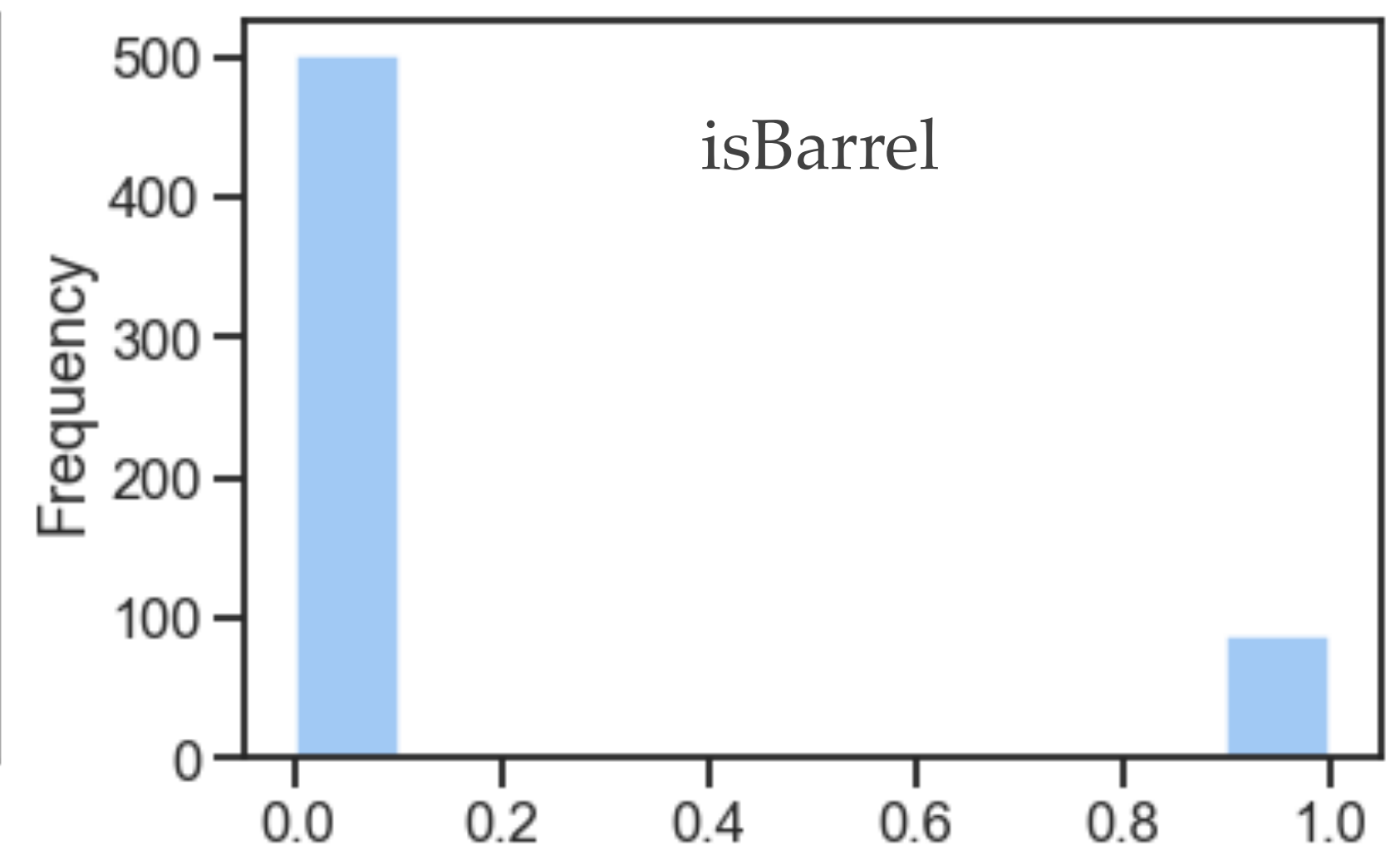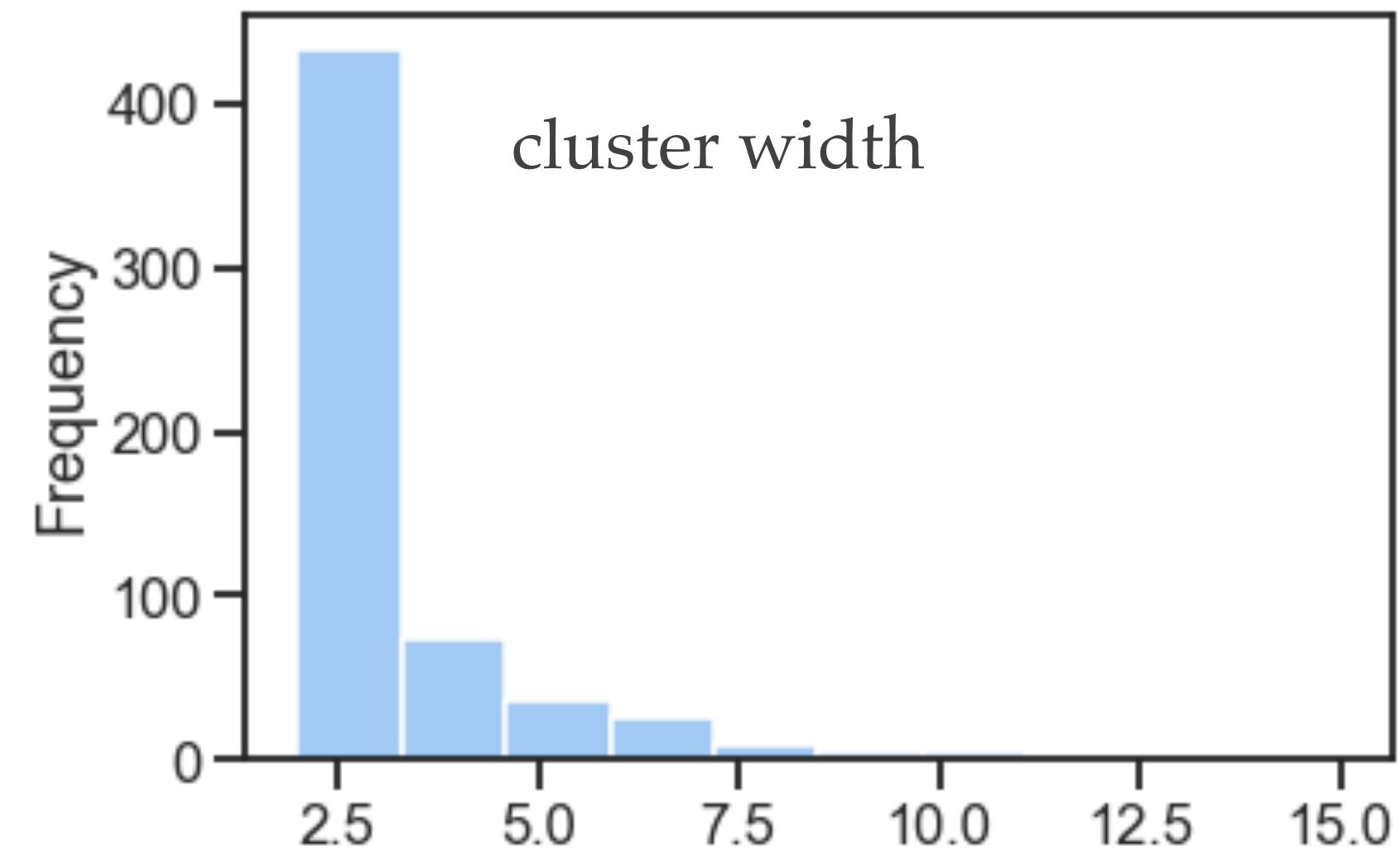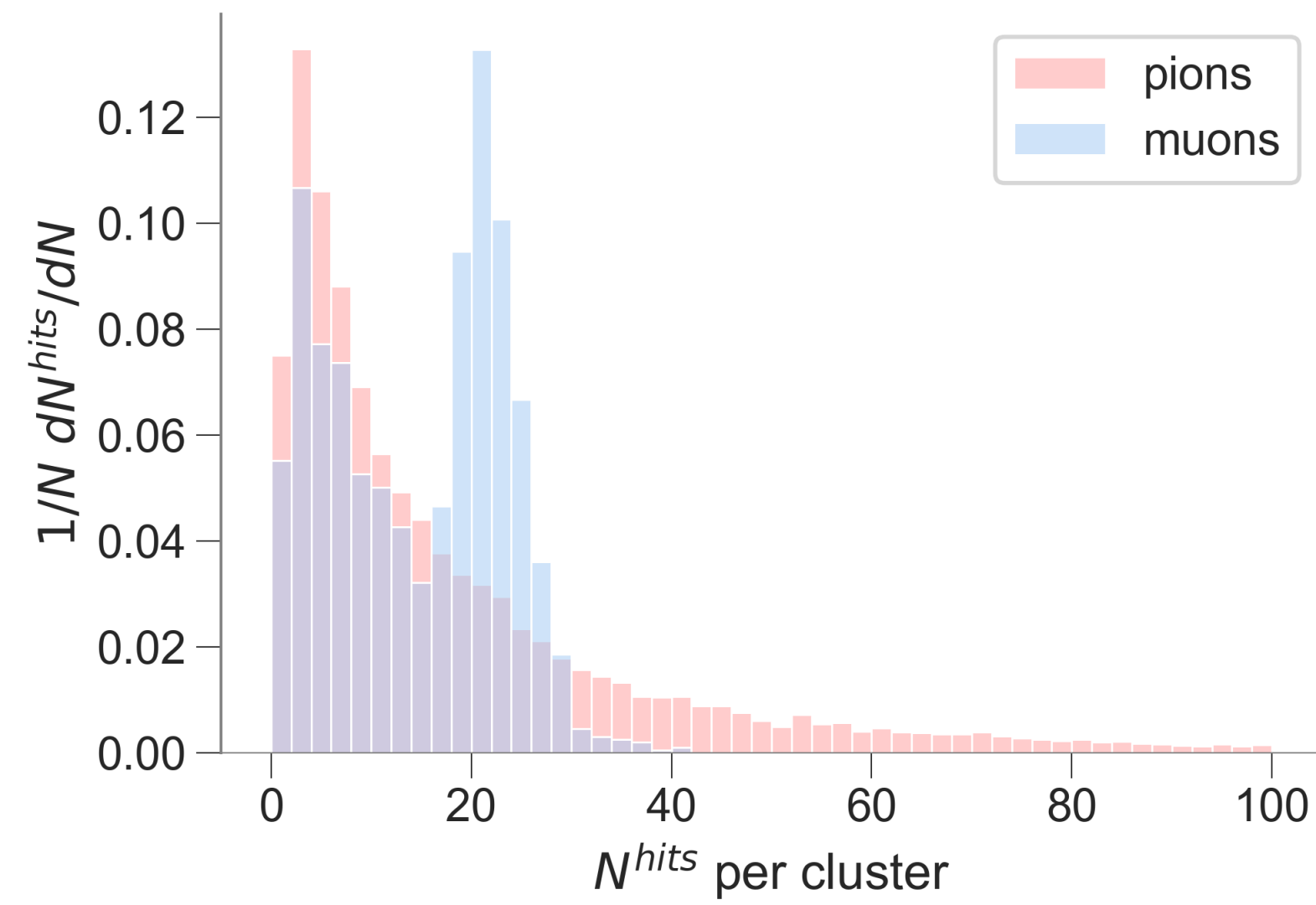
Viktoriya Zel (Balandina) from Muon Group has joined this task. This might be her Master thesis topic at Dubna Univ.

The OPTICS algorithm - a generalization of the DBSCAN that relaxes *epsilon* from a single value to a range - can be another option for improvement; HDBSCAN and hierarchal clustering algorithms are other candidates.
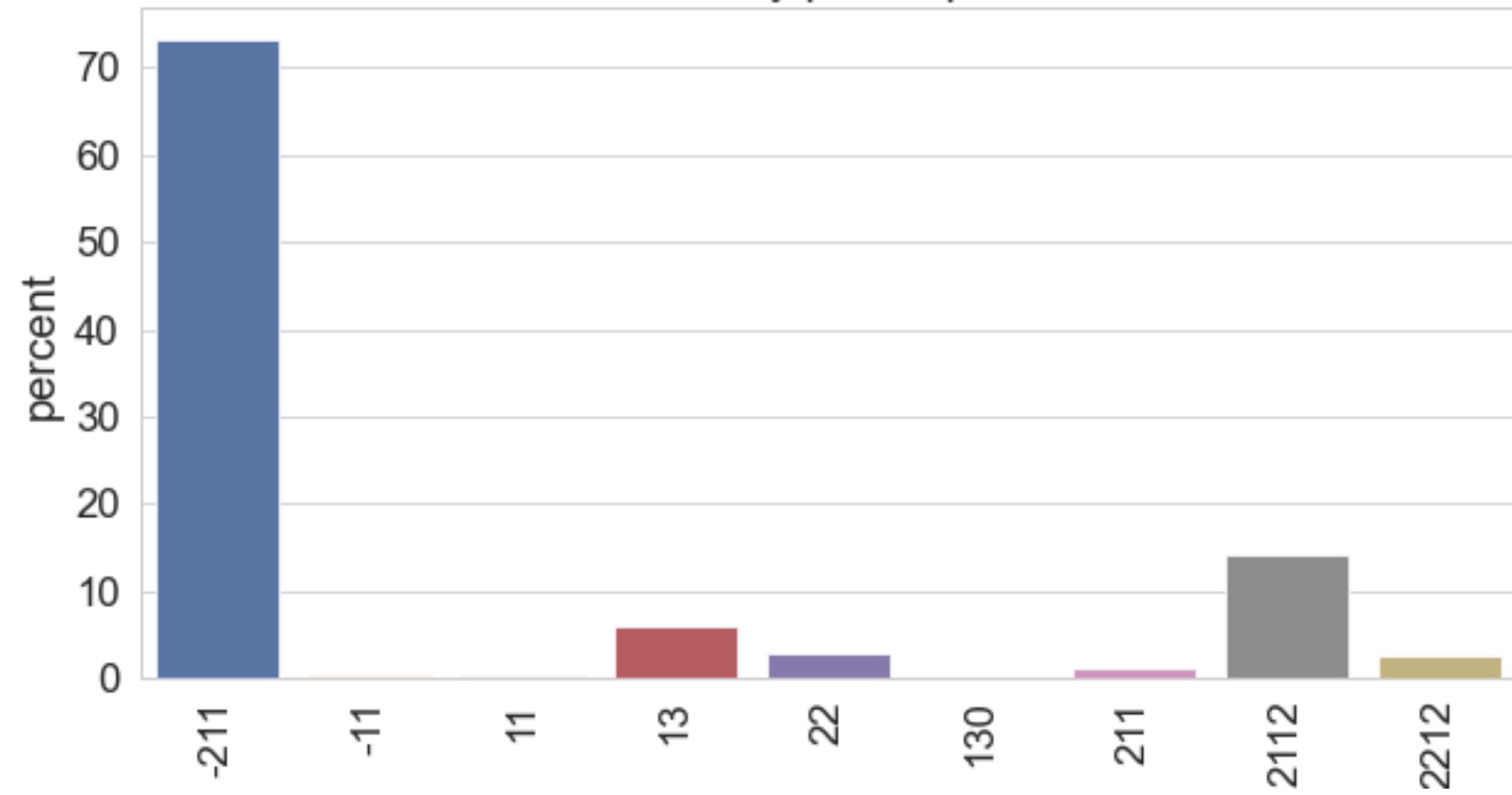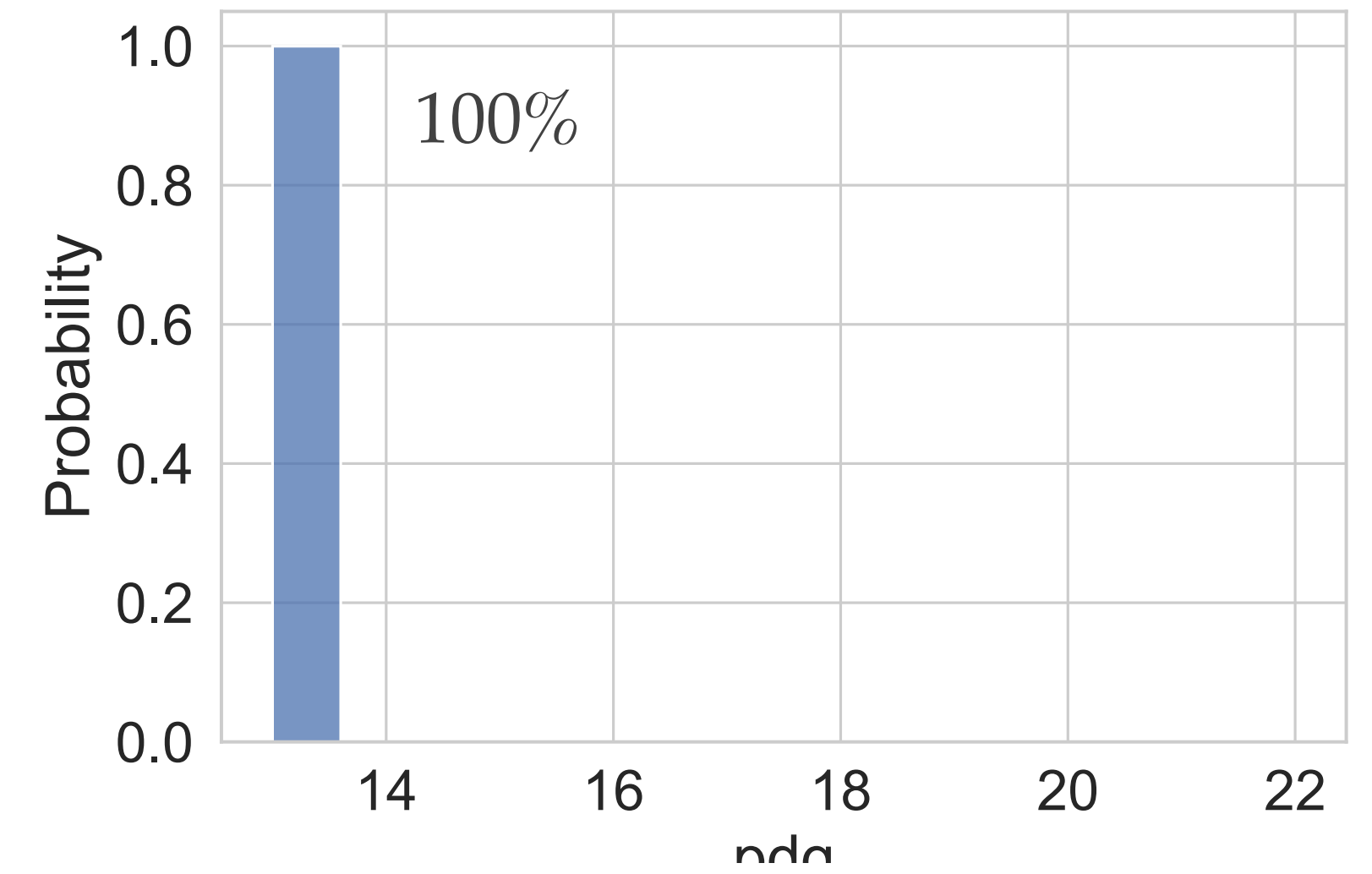
Thank you for attention!

# Backup

# Backup

## Single Pions

## Single Muons



1 GeV
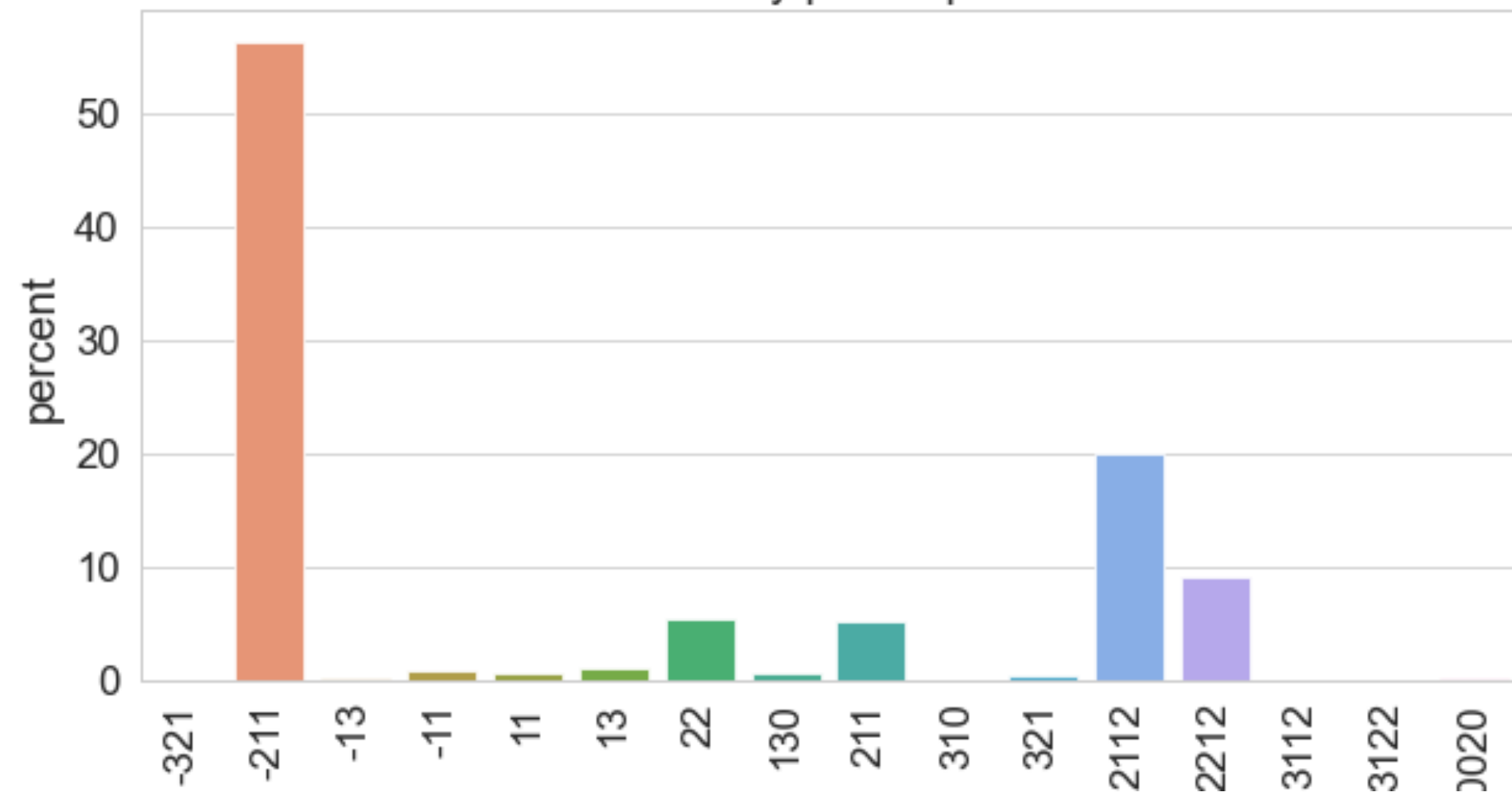


3 GeV