

Particle ID in the AHCAL using Boosted Decision Trees

CALICE Analysis meeting

Vladimir Bocharnikov, DESY
May 20, 2020

Outline

AHCAL Particle ID using BDTs

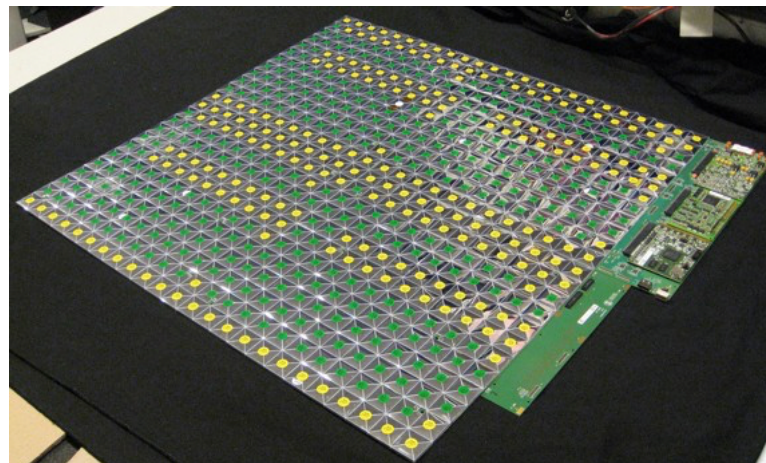
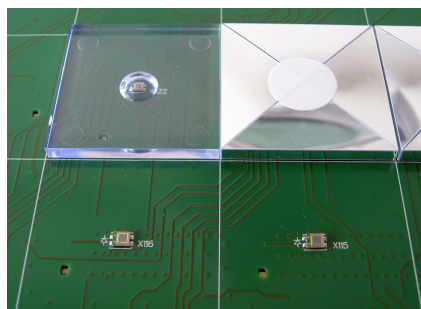
- CALICE AHCAL test beam prototype
- Particle identification
 - Motivation and method overview
 - Data preparation
 - Boosted Decision Tree method description
 - Parameters and input
 - Resulting metrics
 - Application to test beam data
 - Sources of confusion
- Summary and outlook

CALICE AHCAL

Test beam prototype.

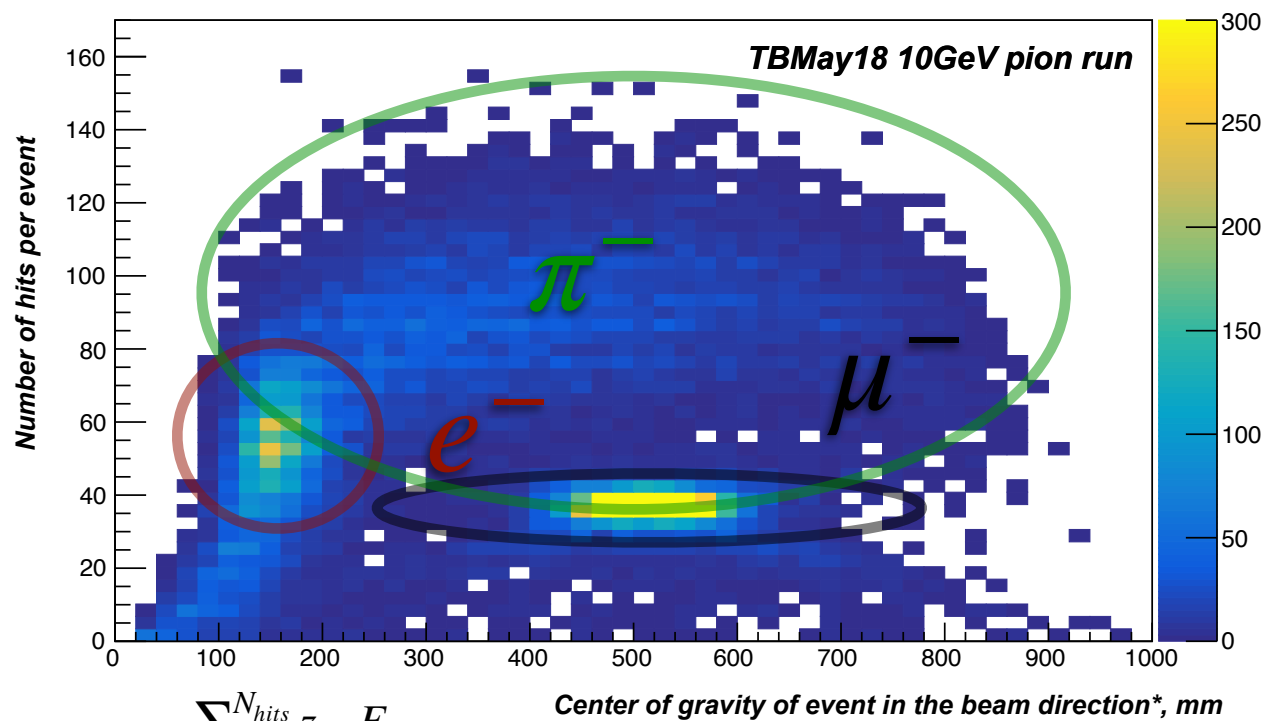
38 active layers of 24x24 scintillator tiles ($3 \times 3 \text{ cm}^2$)
alternating with 1.7 cm steel absorber + 1 “Tokyo” layer
with $6 \times 6 \text{ cm}^2$ tiles

In total: ~ 22000 channels, $\sim 4 \lambda$



Motivation for particle ID

In test beam data



$$* z_{CoG} = \frac{\sum_{i=1}^{N_{hits}} z_i \cdot E_i}{E_{sum}}$$

We always deal with admixture of other particles.

⇒ To investigate detector response to particles of given type we need to perform particle identification

Particle ID workflow

Classification procedure



Pre-analysis

- Calculation of common observables
- Clustering and track finding*

Event filtering

- By **number of hits**:
nHits > nHits_min
- **multi-particle** and **upstream shower** event rejection



BDT multiclass model

trained on simulations (10-200GeV).

3 classifiers:

Hadron classifier

- Trained on showering pions

Electron classifier

- Trained on electrons

Muon (muon-like) classifier

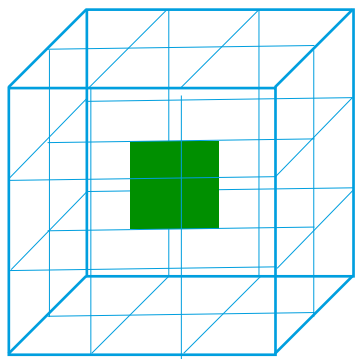
- Trained on muons

* Described during CALICE Collaboration Meeting at CERN:

https://agenda.linearcollider.org/event/8213/contributions/44343/attachments/34812/53758/VBocharnikov_CALICE_meeting_CERN.pdf

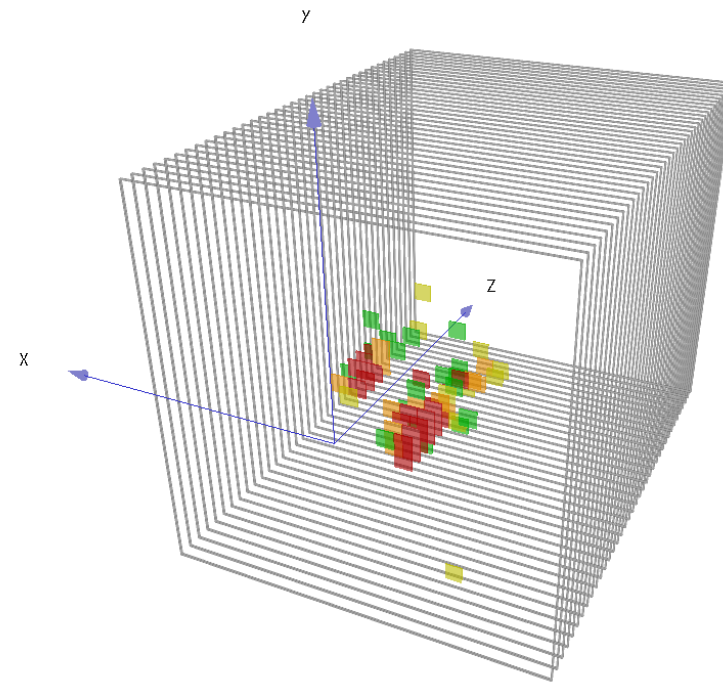
Event filtering

Simplified algorithms.



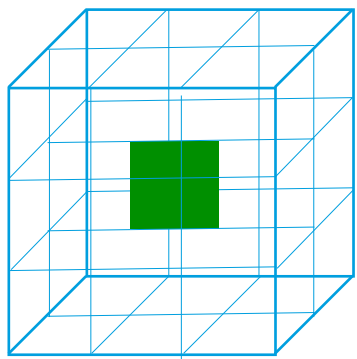
Clustering: Hits are grouped in clusters if they are neighbours in volume. First 5 layers are taken into account

If $N_{Clusters} > 1 \Rightarrow$ multi-particle event (or upstream shower)



Event filtering

Simplified algorithms.

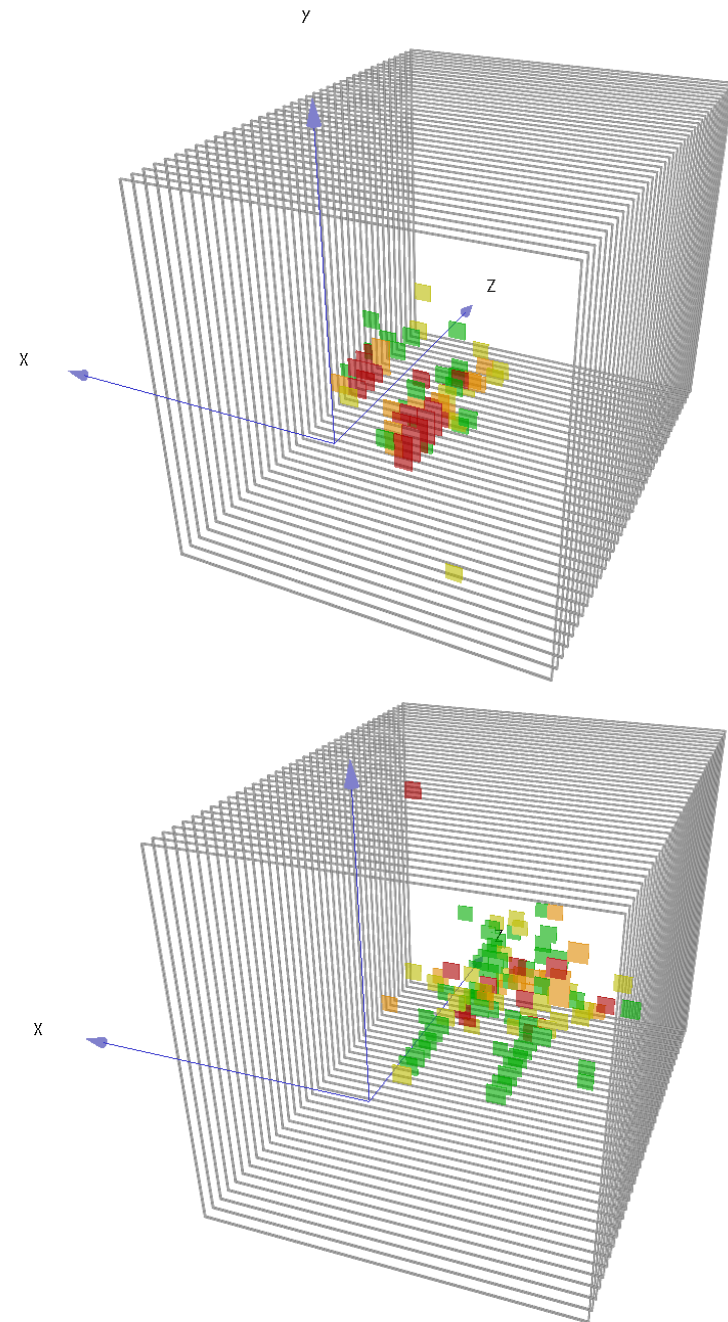


Clustering: Hits are grouped in clusters if they are neighbours in volume. First 5 layers are taken into account

If $N_{Clusters} > 1 \Rightarrow$ multi-particle event (or upstream shower)

MIP tracking: Construct towers with same x and y coordinates. First 5 layers are taken into account.

If $N_{MIPTracks} > 1 \Rightarrow$ multi-particle event



BDT classification

Model and input. TBJune18.


Software and model:

- **LightGBM** package
- Multi-class **Gradient Boosted Decision Tree**
- **Multi log** loss function

BDT classification

Model and input. TBJune18.

Software and model:

- **LightGBM** package
 - Multi-class **Gradient Boosted Decision Tree**
 - **Multi log** loss function
- 

Gradient Boosting:

Method combines many sequential decision trees with weights. Weights are optimised during training by calculating the gradient of loss function

BDT classification

Model and input. TBJune18.

Software and model:

- **LightGBM** package
- Multi-class **Gradient Boosted Decision Tree**
- **Multi log** loss function

Gradient Boosting:

Method combines many sequential decision trees with weights. Weights are optimised during training by calculating the gradient of loss function

Multi log loss:

$$L = -\frac{1}{N} \sum_i^N \sum_j^3 Y_{ij} \ln(p_{ij})$$

Where N - number of events in the test sample, 3 - number of classes, Y_{ij} is binary variable with the expected labels and p_{ij} is the classification probability output by the classifier for the i -instance and the j -label.

BDT classification

Model and input. TBJune18.

Software and model:

- **LightGBM** package
- Multi-class **Gradient Boosted Decision Tree**
- **Multi log** loss function

Training and test set:

- **MC** particles **10-200GeV** QGSP_BERT_HP physics list simulated and reconstructed using June 2018 setup:
 - **pions** ($st \leq 40$)
 - **electrons**
 - **muons**
- Simulated data is split **50/50 - test/train**

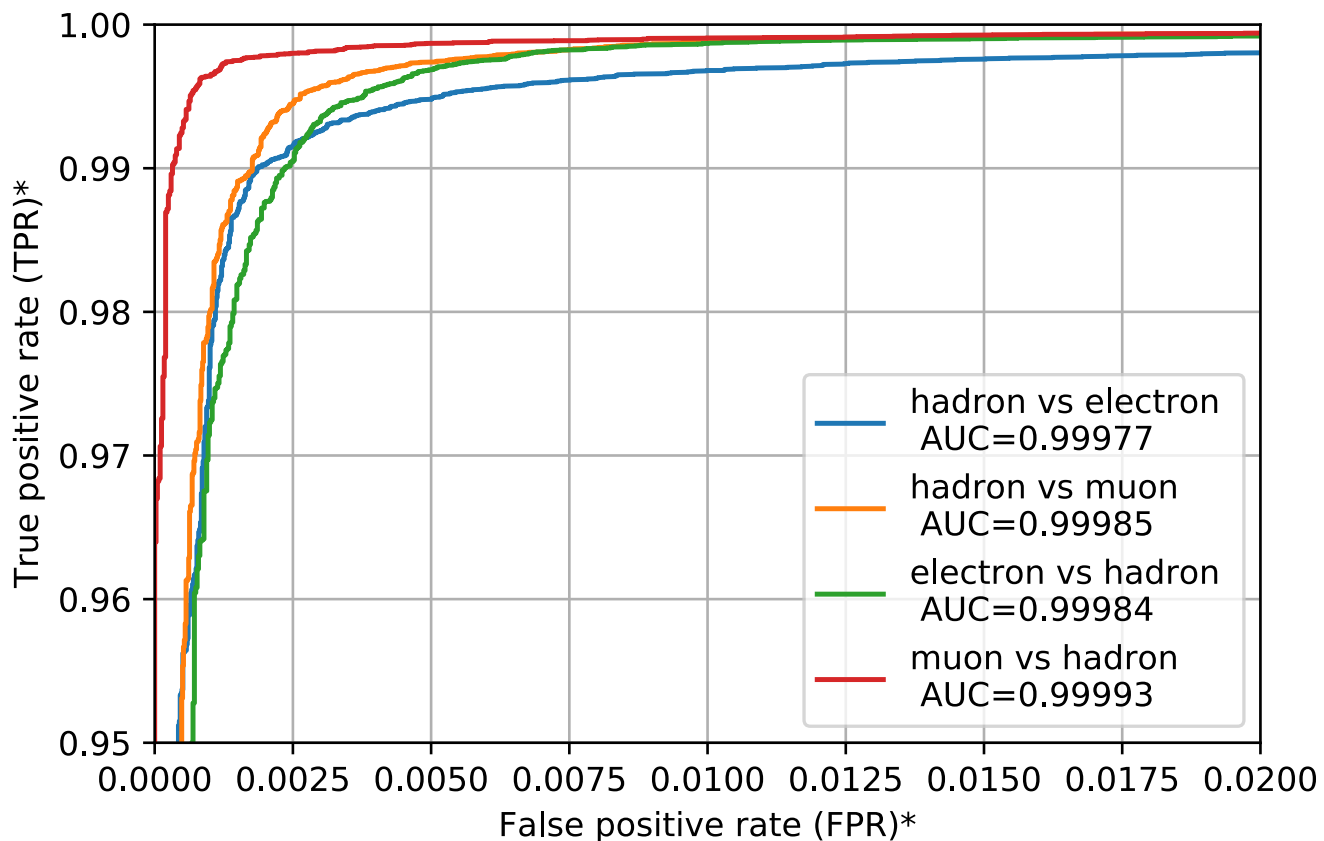
Observables:

- Number of hits
- Shower start
- Event radius
- Center of gravity in z
- Energy fraction in first 22 layers
- Energy fraction in shower center
- Energy fraction in shower core
- Fraction of track hits
- Number of track hits
- Number of layers with hits from last 5
- Mean hit energy after shower start

Resulting metrics

After training

ROC curves for the test data



$$*TPR = \frac{TP}{TP + FN}, \quad FPR = \frac{FP}{FP + FN}$$

Multi log loss:

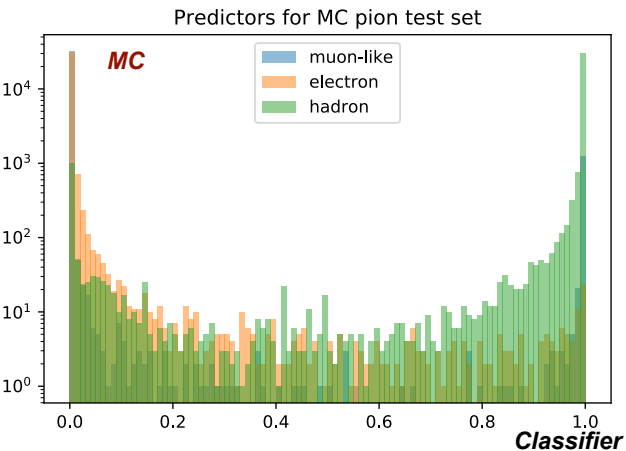
$$L = -\frac{1}{N} \sum_i \sum_j^3 Y_{ij} \ln(p_{ij}) = 0.0086$$

Where N - number of events in the test sample, 3 - number of classes, Y_{ij} is binary variable with the expected labels and p_{ij} is the classification probability output by the classifier for the i -instance and the j -label.

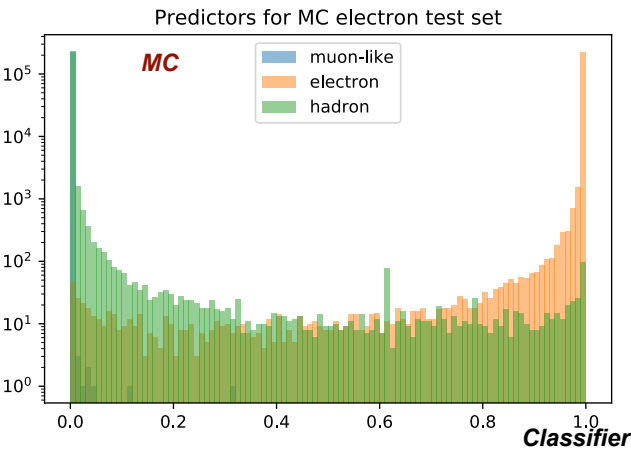
BDT classification

Output. Comparison with data.

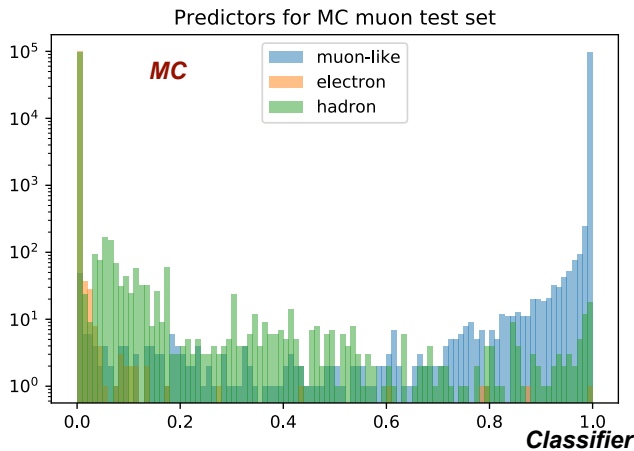
Hadrons



Electrons

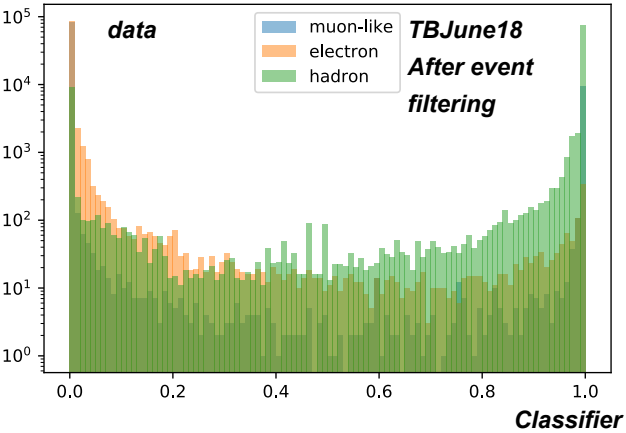


Muons

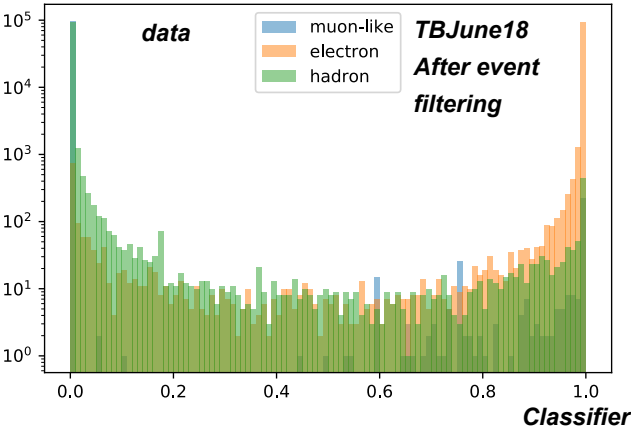


- Similar response on data and simulations

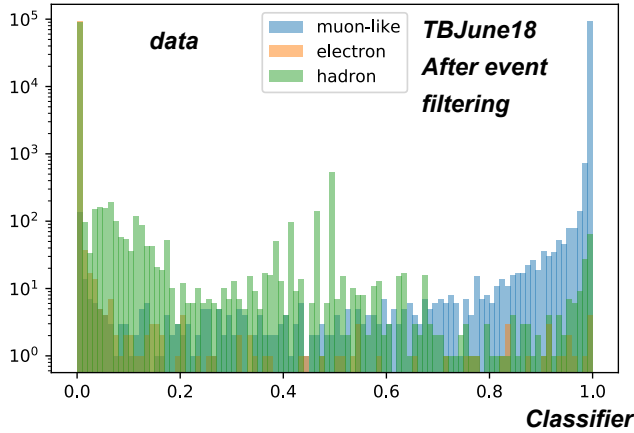
Hadrons



Electrons



Muons

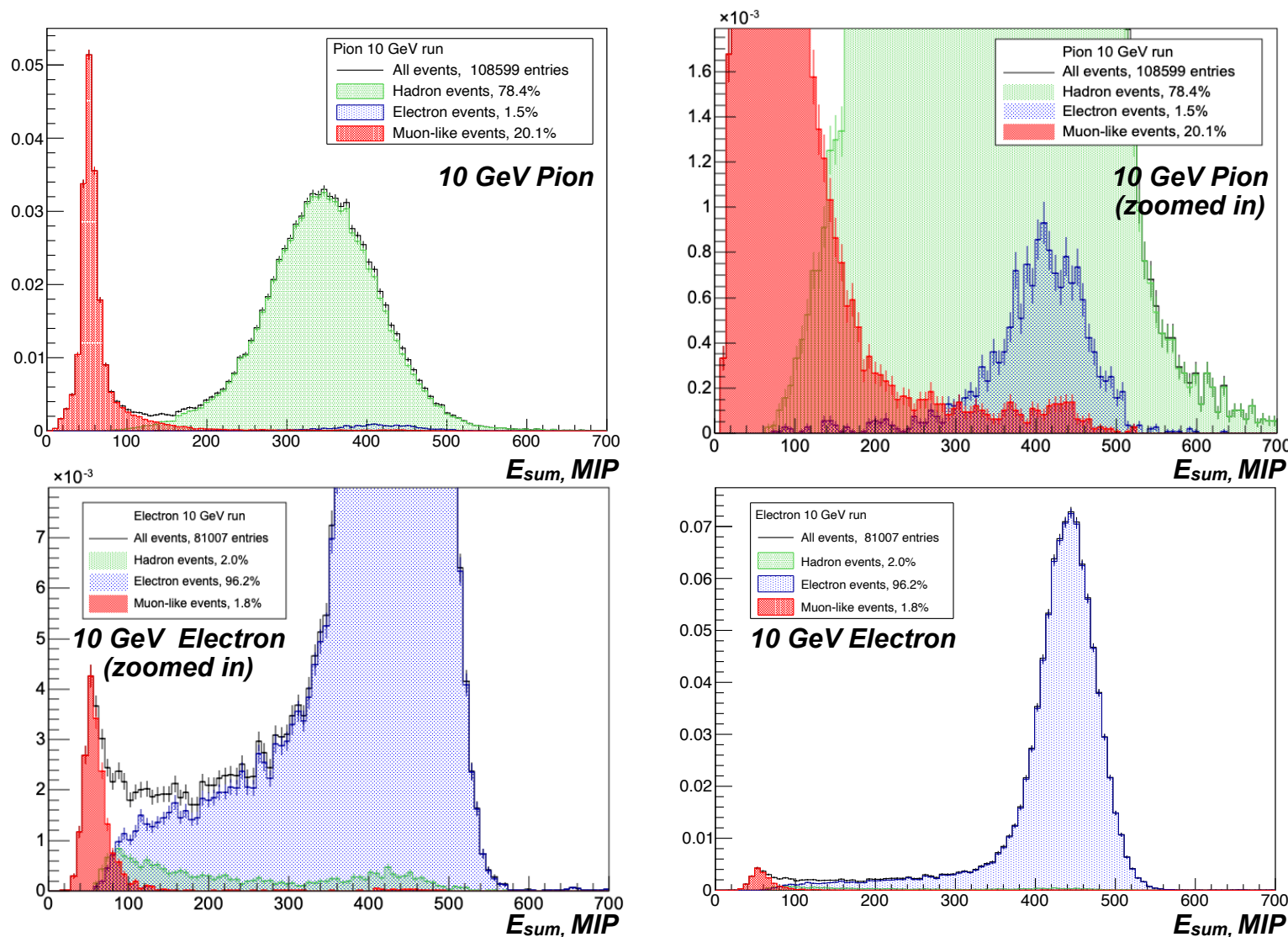


Classifiers:

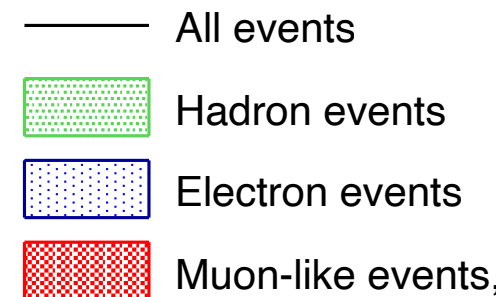
- muon-like
- electron
- hadron

Results on test beam data taken in June 2018

Energy sum distributions for 10GeV runs

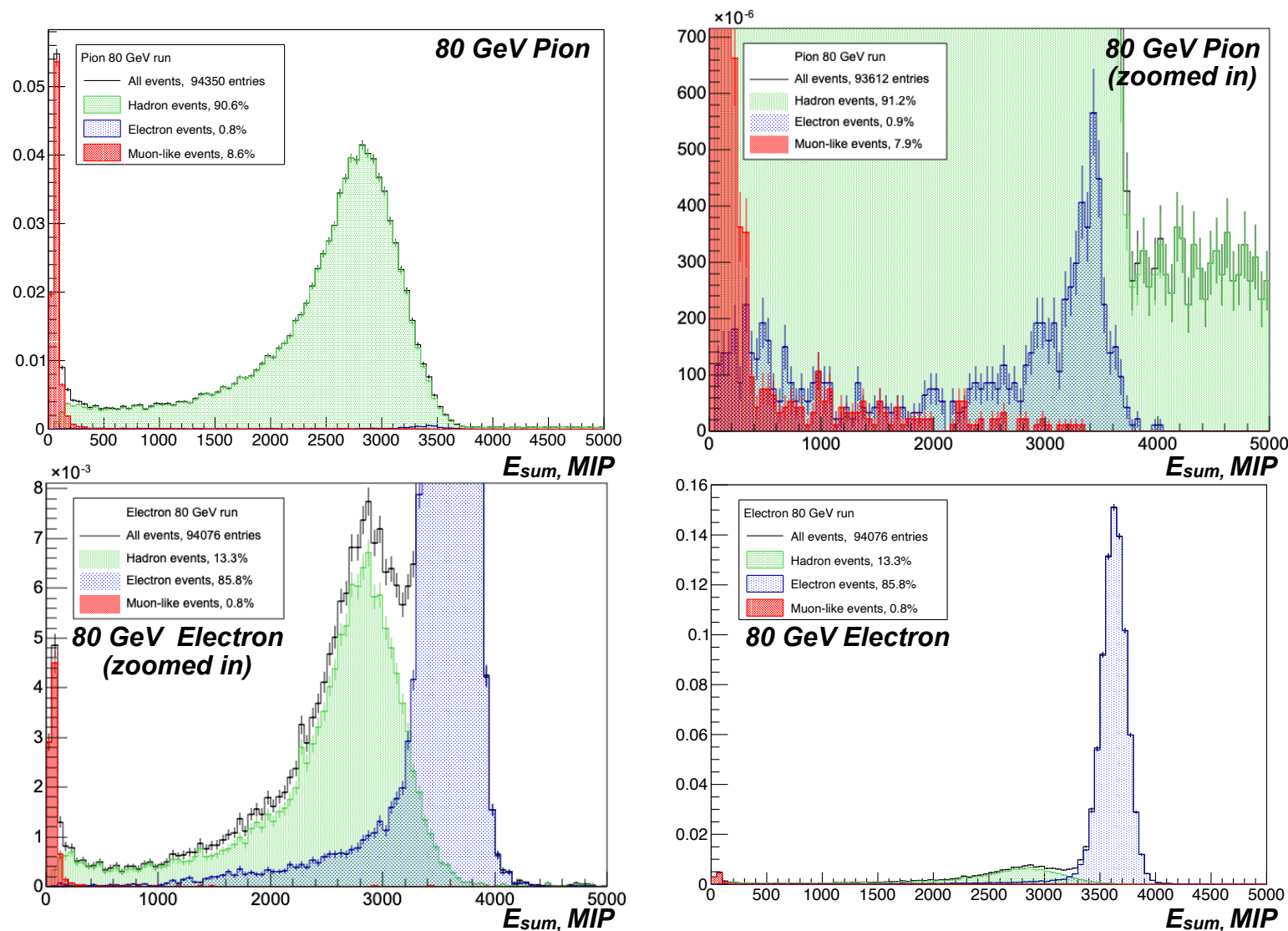


- Energy expectation for electron events in pion run is close to real electron run
- Long high energy tail of muon-like events
- Low energy tail for electrons
- Most of hadron events in electron run are at low energy



Results on test beam data taken in June 2018

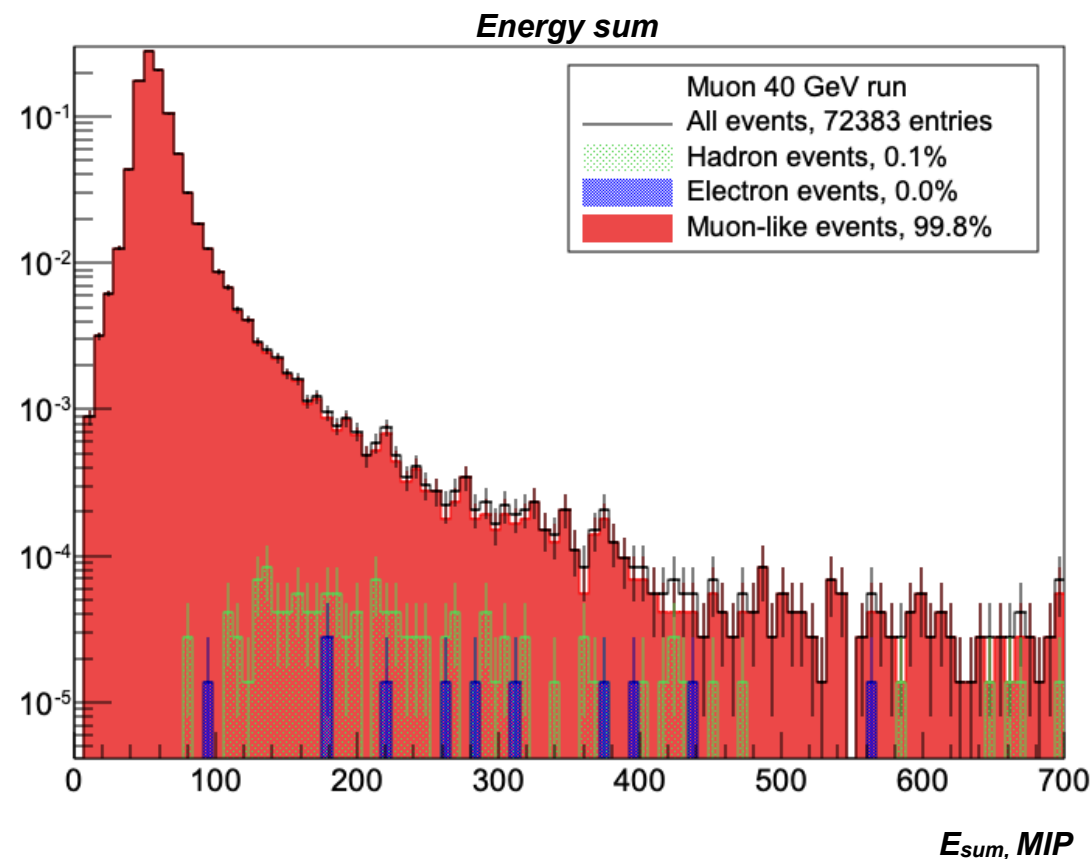
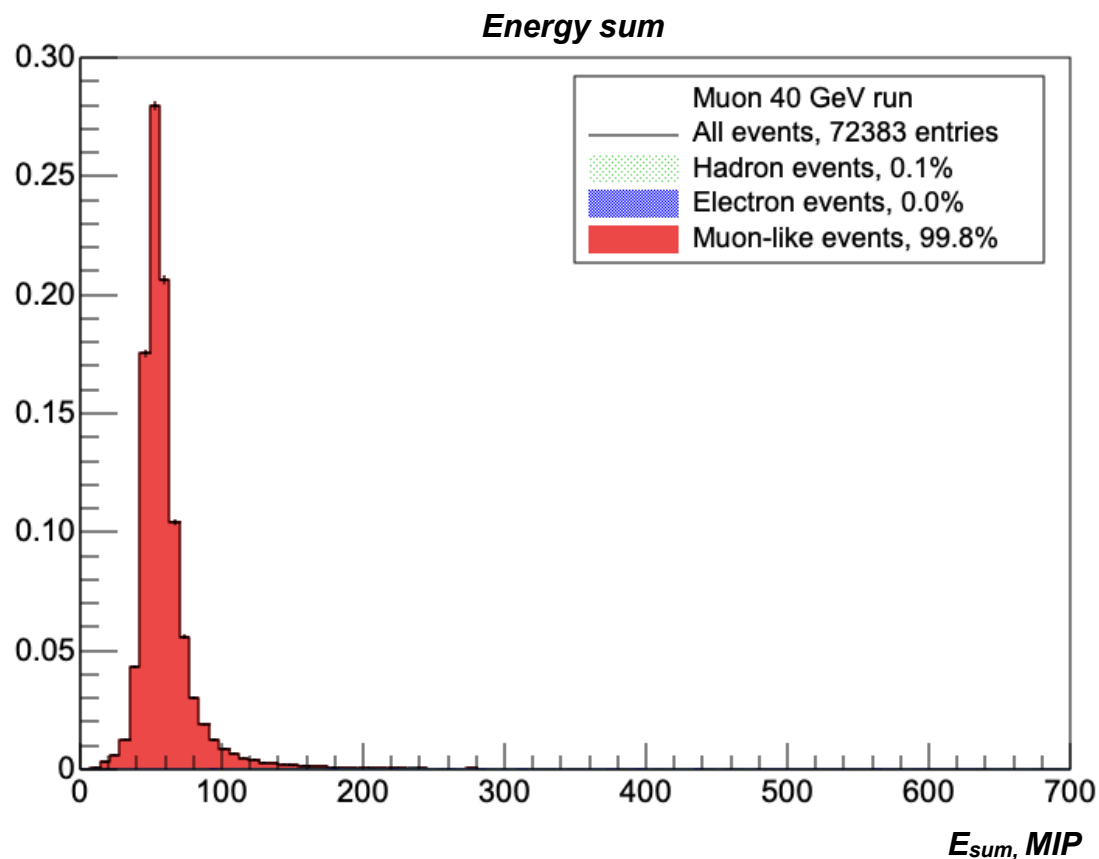
Energy sum distributions for 80GeV runs



- Energy expectation for electron events in pion run is close to real electron run
- Energy distribution of hadron events in 80GeV electron run looks very similar to actual 80GeV pion

Results on test beam data taken in June 2018

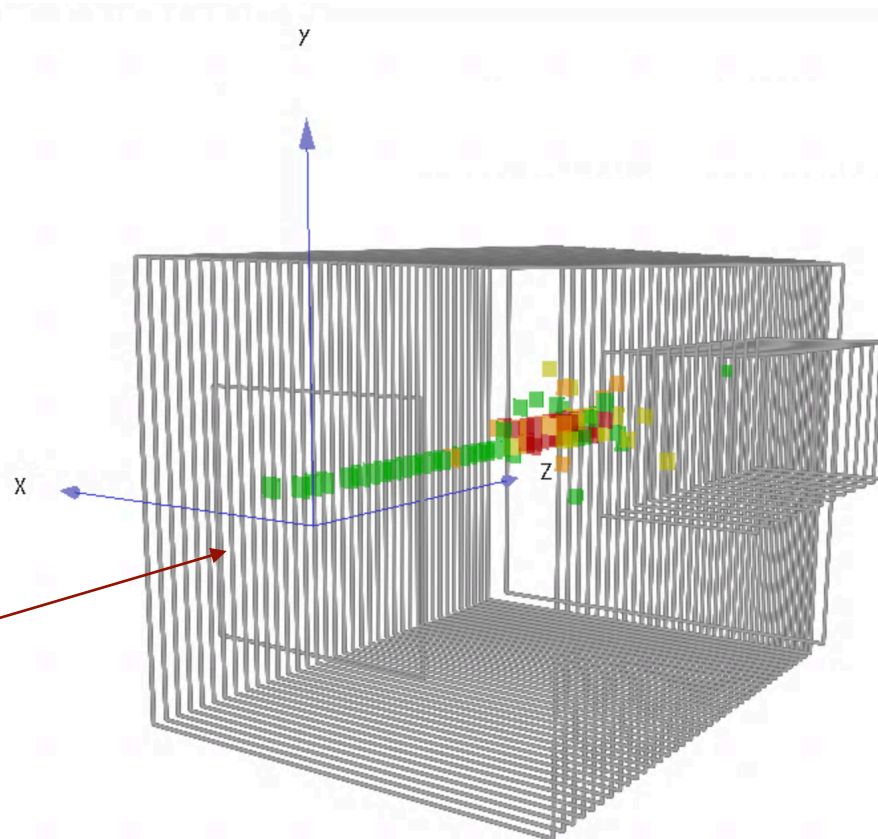
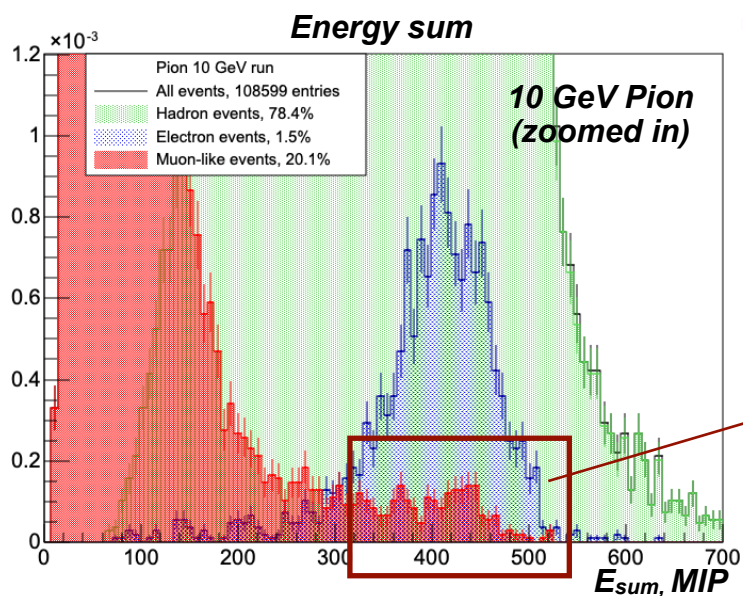
Energy sum distribution for 40GeV muon run



- Very low admixture of other particles
- Little fraction of delta electrons can be classified as hadron event

Sources of confusion

From 10GeV pion run

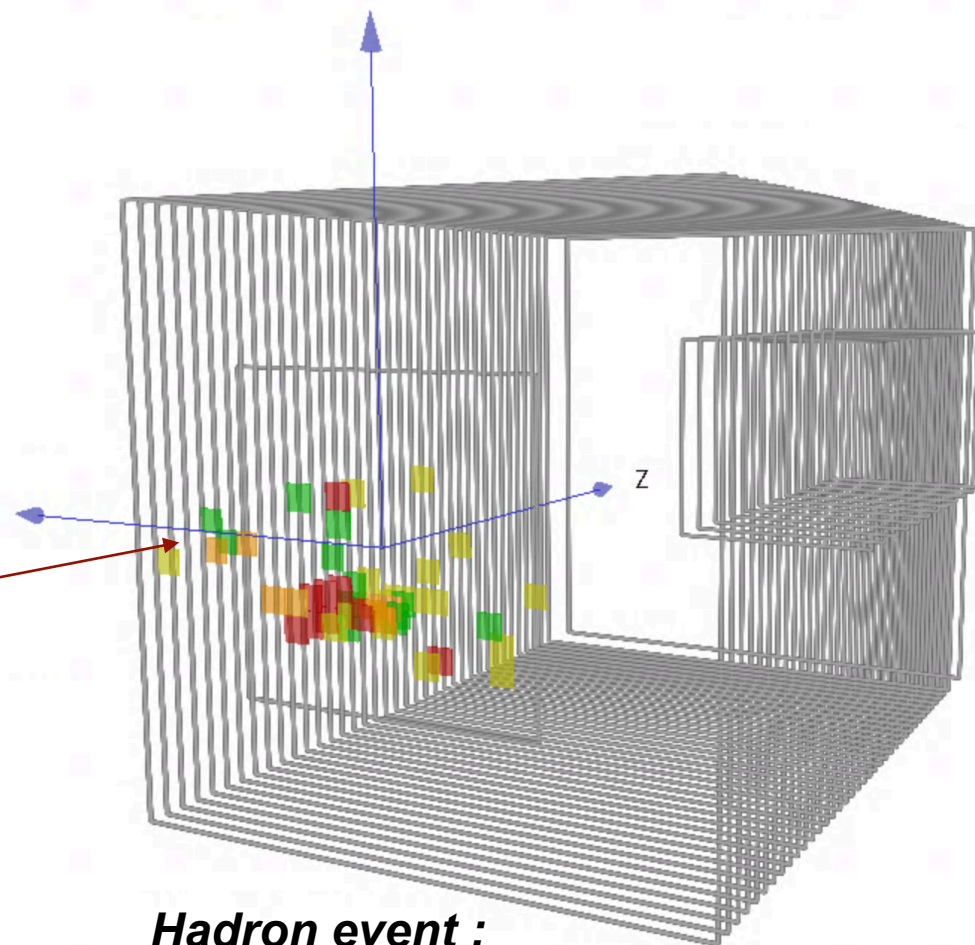
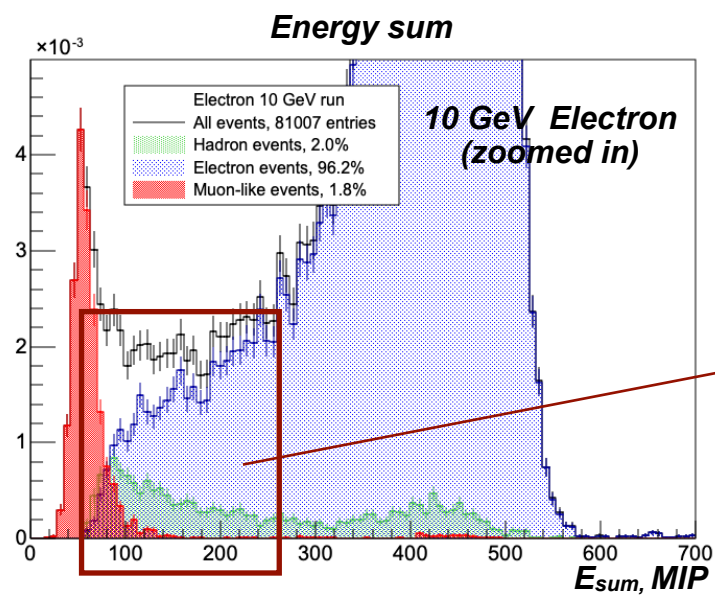


Muon-like event :
Mu-like score is 0.51
Had score is 0.48

- Compact pion showers with late shower start can be classified as muons
- Additional variables can improve identification
- Fraction $\ll 1\%$

Sources of confusion

From 10GeV electron run

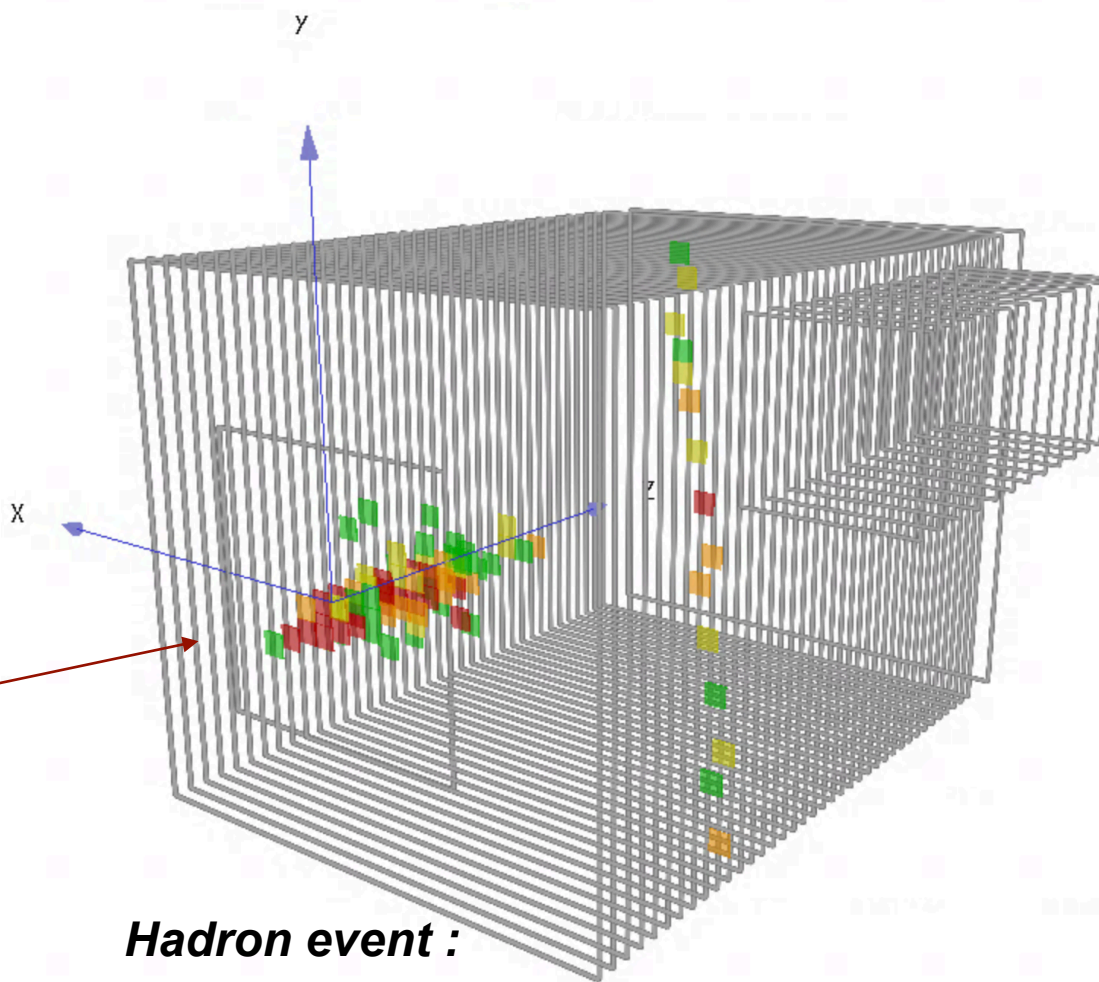
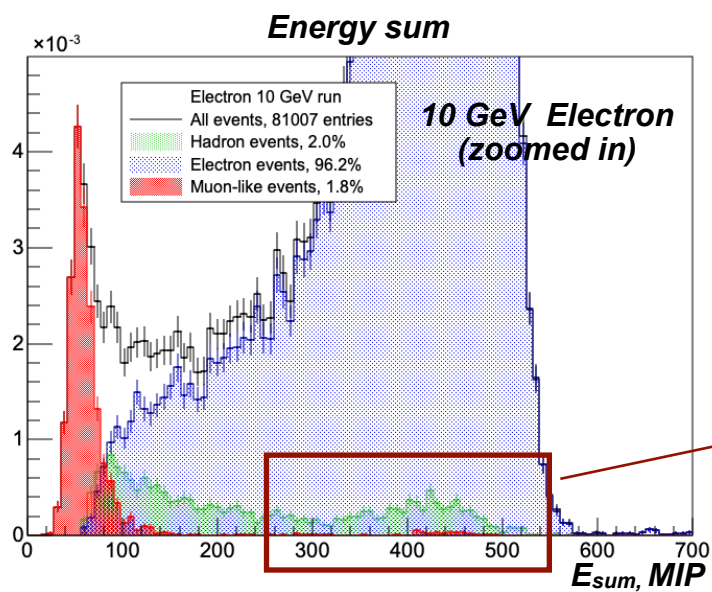


Hadron event :
Had score is ~0.9

- Multi-particle/upstream shower events with small fragments can be classified as hadron events
- Multi-particle events can be partly filtered out using timing information

Sources of confusion

From 10GeV electron run



Hadron event :
Had score is ~ 0.98

- Some events are contaminated with cosmic muons
- Multi-particle events can be partly filtered out using timing information

Summary and outlook

AHCAL Particle ID using BDTs

- ☒ BDT particle ID method in the AHCAL was discussed

 - ☒ Method shows good performance

 - ☒ Similar response on data and MC

 - ☐ Feature importance study* is planned as next step

*sort input observables by importance to drop less useful ones

- ☐ More advanced event filtering for data is needed

 - ☐ Timing analysis

Backup

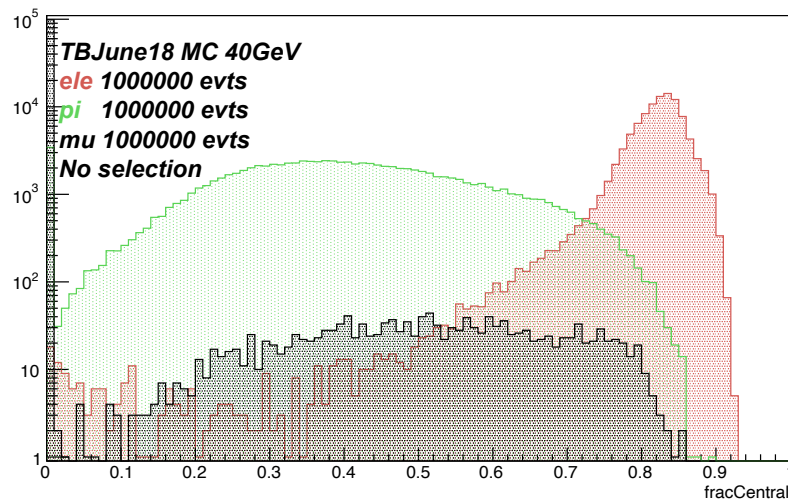
Disadvantages of cut-based method

Towards BDT ID

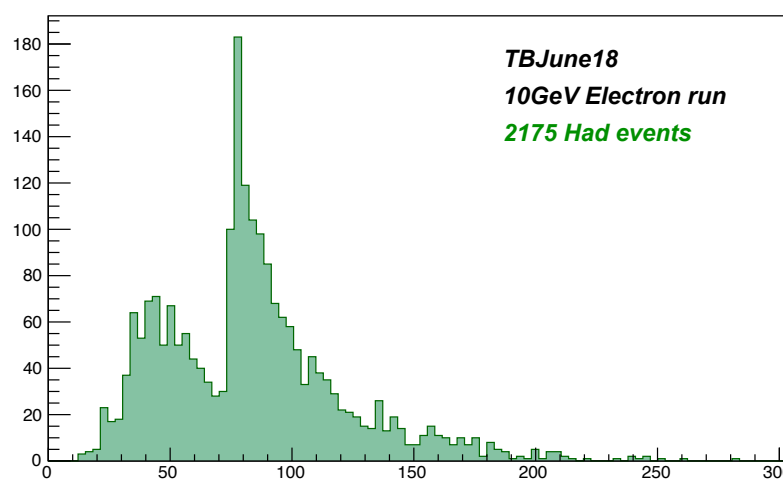
Cut-based method:

- > 10 steering parameters for each energy
- Asymmetric distributions/ long tails can be problematic
- Cut artefacts

Central energy fraction



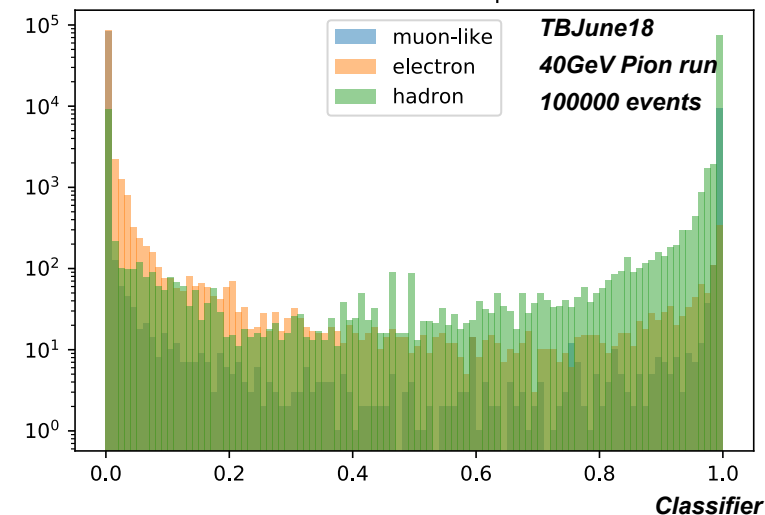
Shower radius



Multivariate methods:

- Can provide probabilistic classifier trained on given distributions of observables
- One model can be used for whole dataset

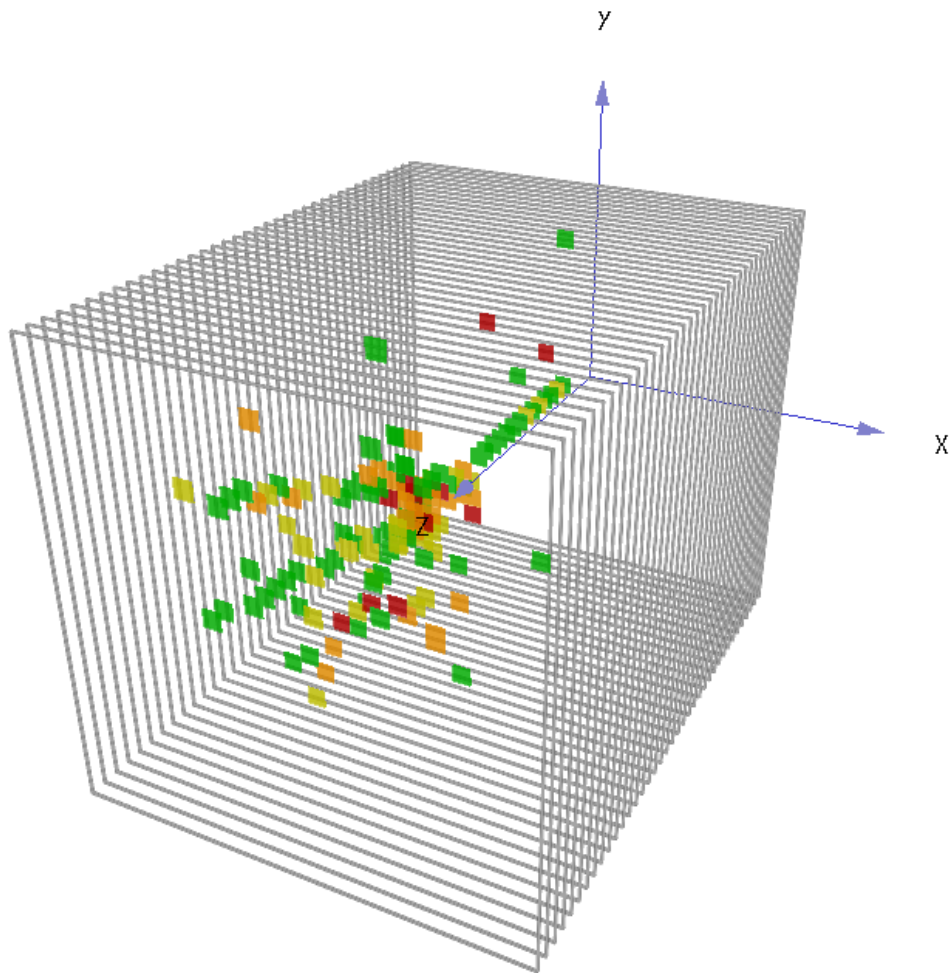
Predictors for 40GeV pion run



Will be discussed during one of the upcoming HGCal meetings

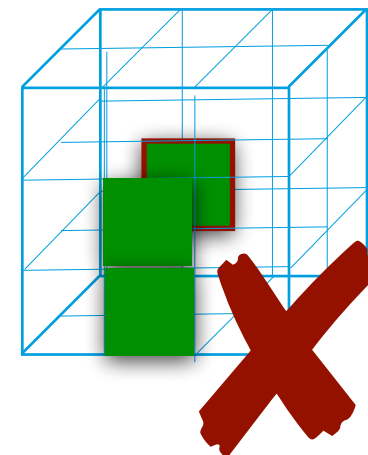
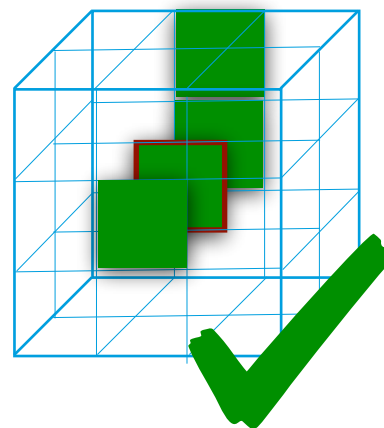
Track finding

Important tool for shower characterisation,
Can be used for particle ID



Track candidates:

2/3 neighbours in surrounding volume. 2 of them on different sides

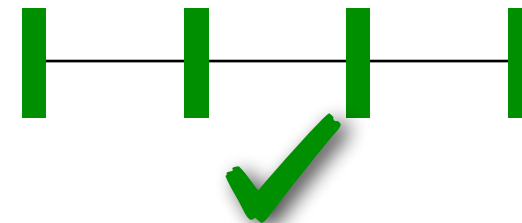
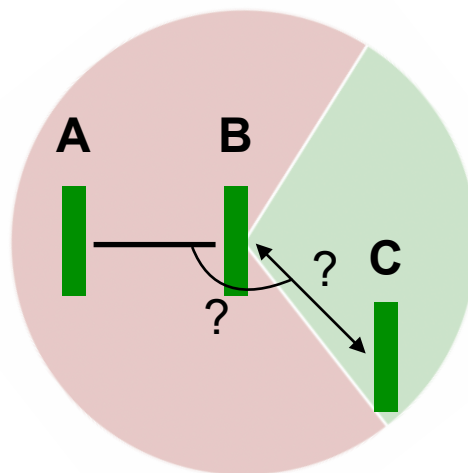
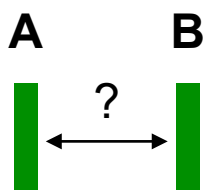
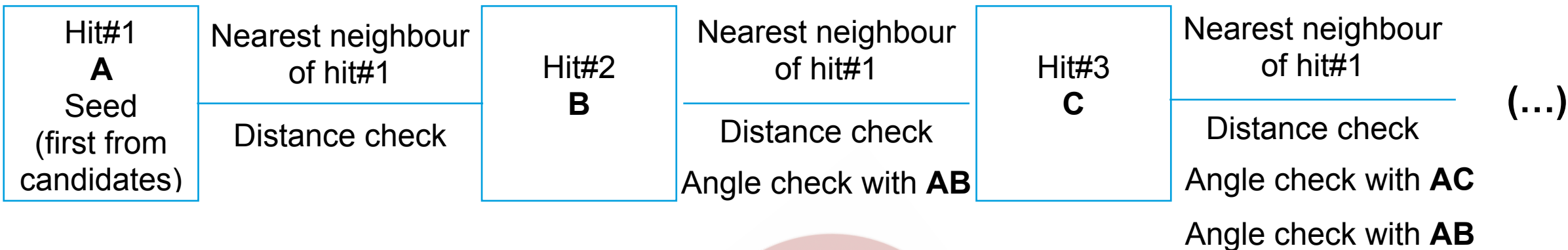


Candidates ordered:

- z-coordinate
- Distance to (0,0,z) in same layer

Track finding

Grouping candidates into tracks

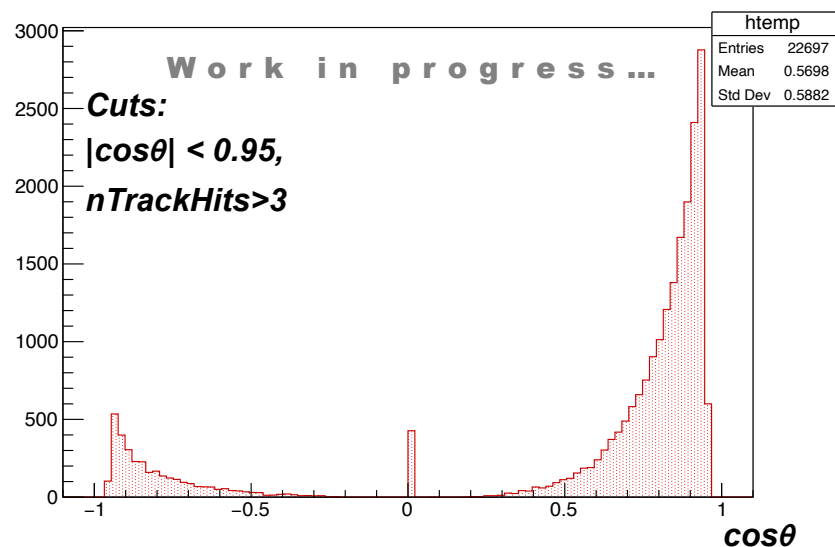
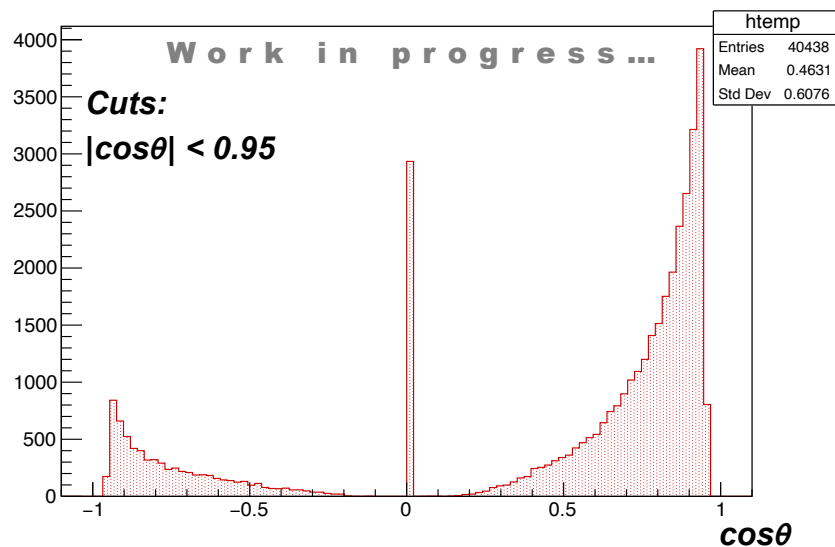


After grouping, track angle is obtained using MSE linear regression

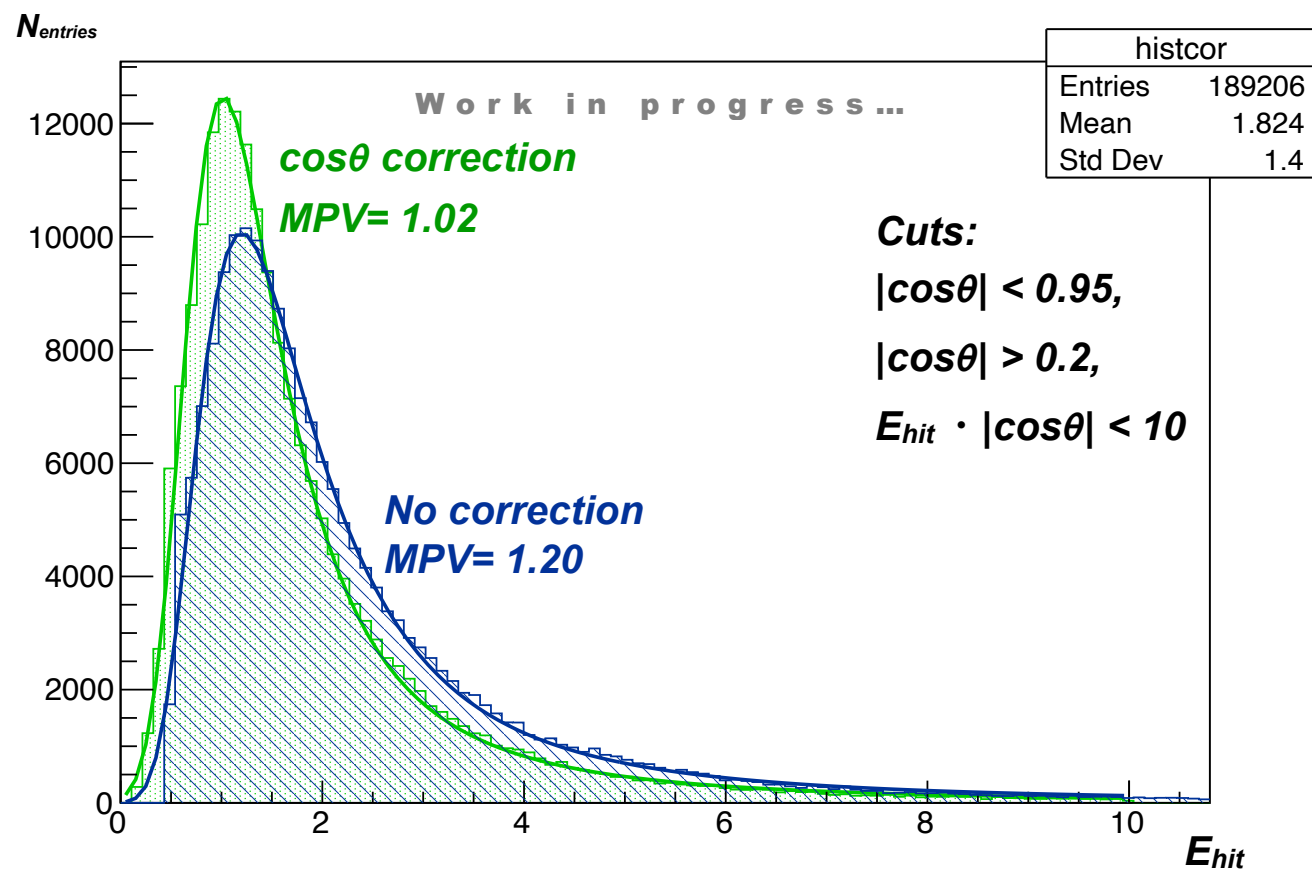
**** Procedure repeated iteratively ****

Tracking quality check

TBMay18 10GeV pion run. 50039 events.

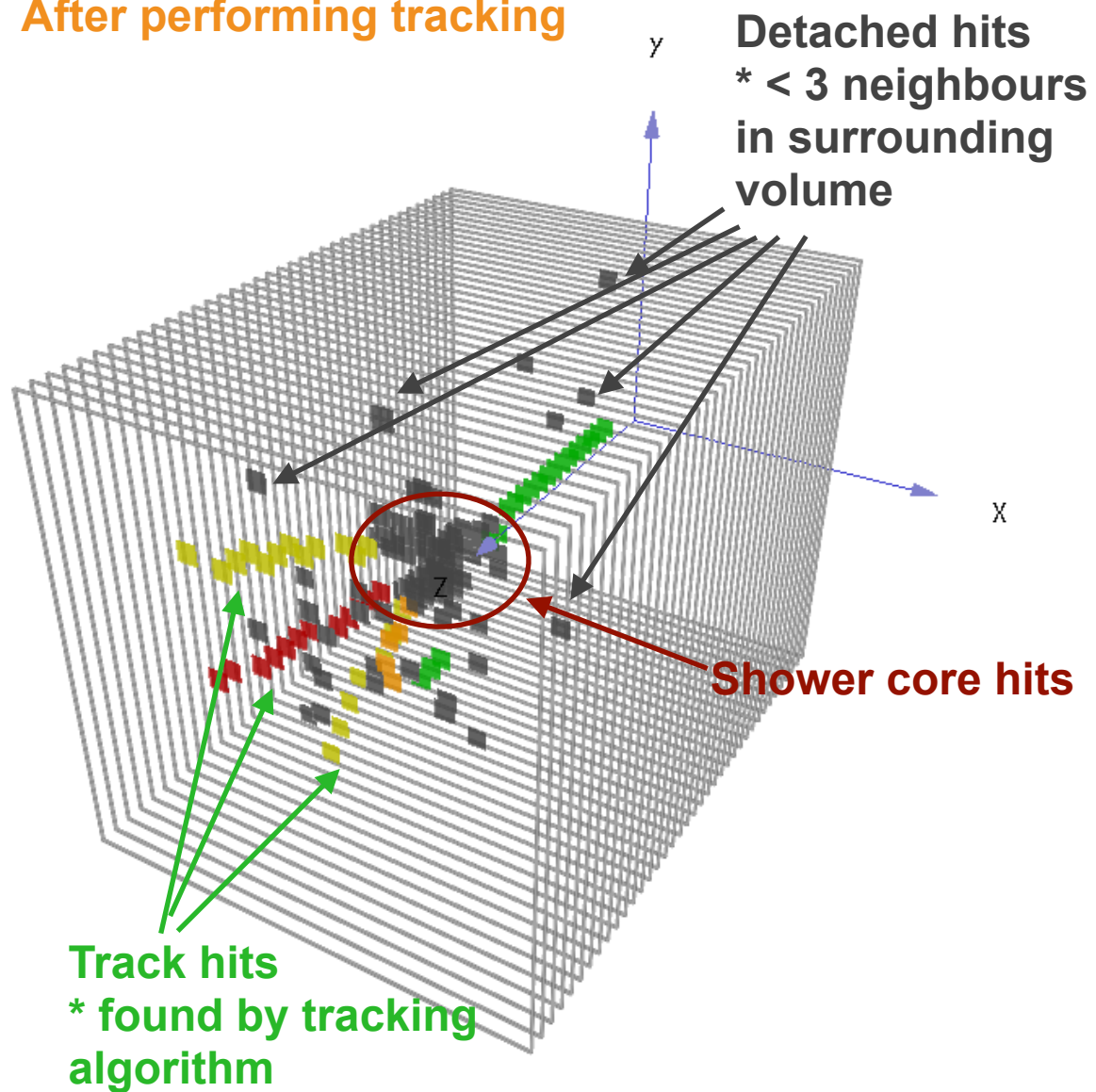


Scintillator path length correction for track hits



Resulting ID variables

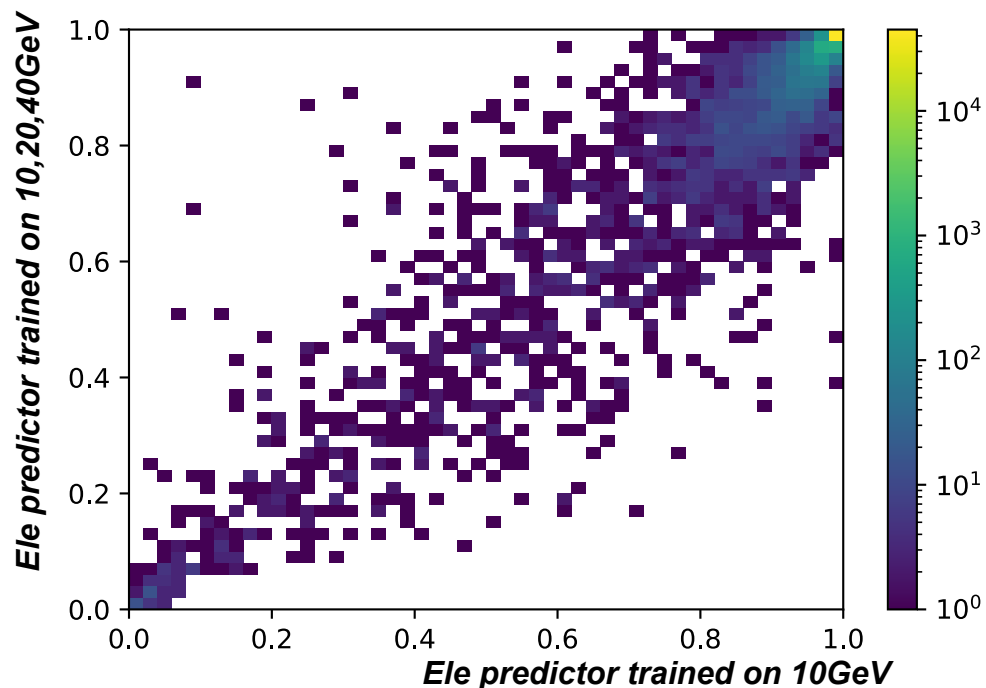
After performing tracking



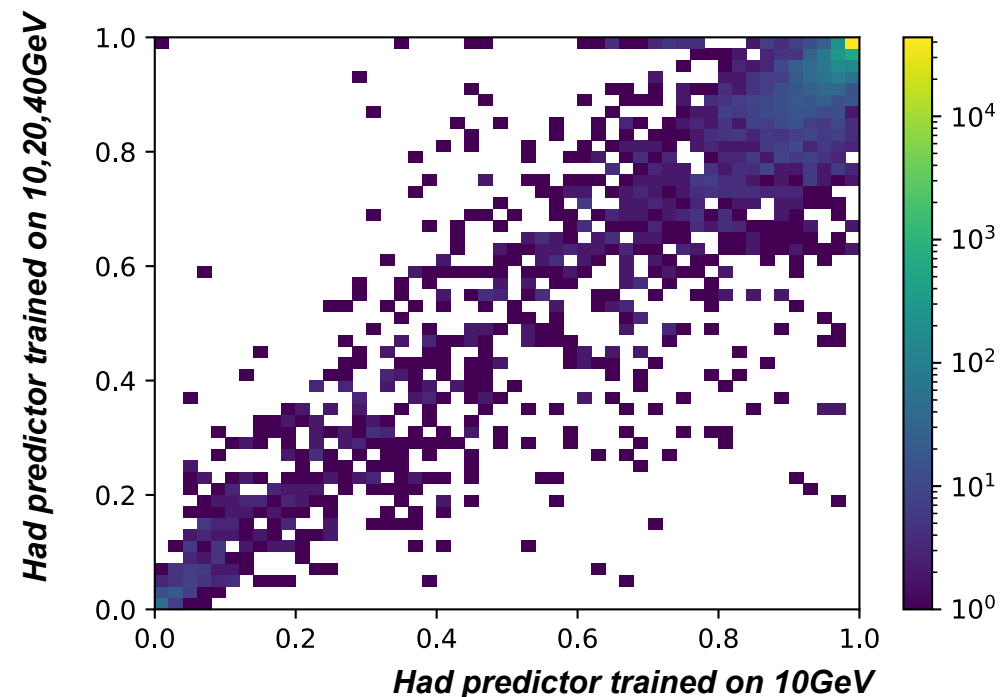
BDT output

Comparison with separate model trained only on 10GeV particles.

10GeV MC electron test sample
50000 events



10GeV MC pion test sample
50000 events



Application on electron data

Of trained BDT model

Electron events: $\text{classifier}_{\text{ele}} > 0.5$

