

Study of residual artificial neural network for PID using the CEPC AHCAL Prototype

Siyuan SONG

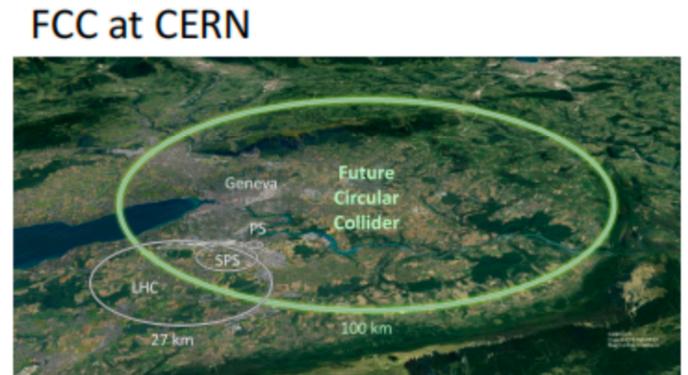
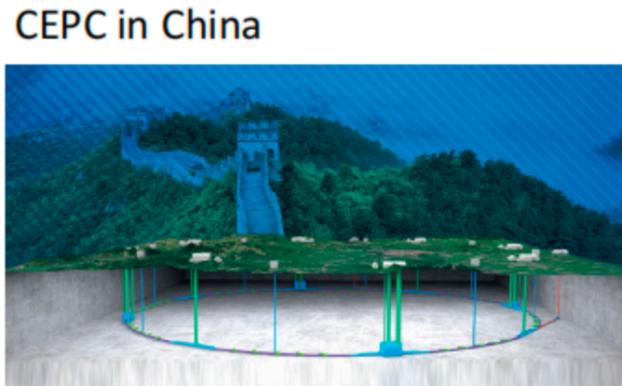
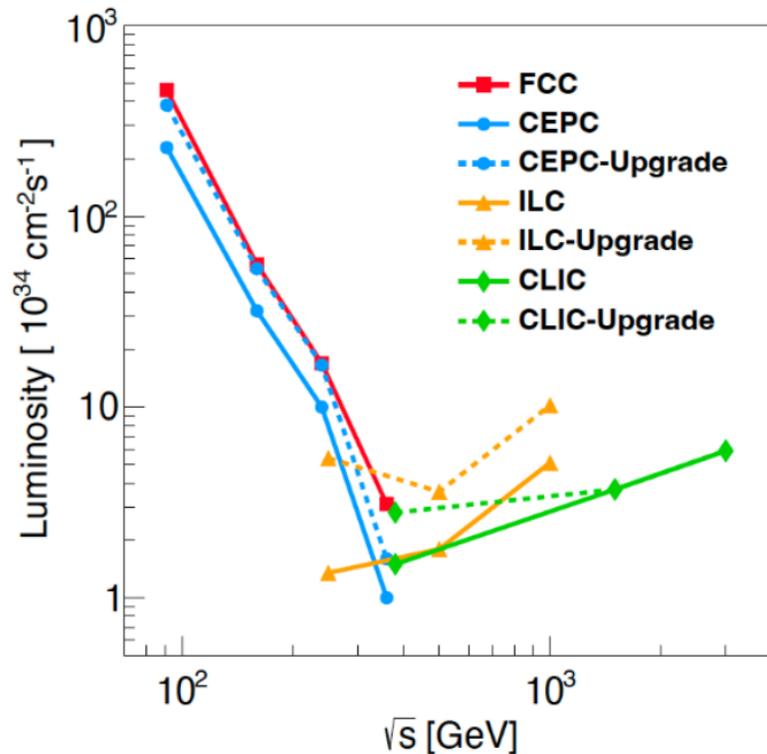
Shanghai Jiao Tong University

The 2023 International Workshop on the CEPC from Oct 23-27

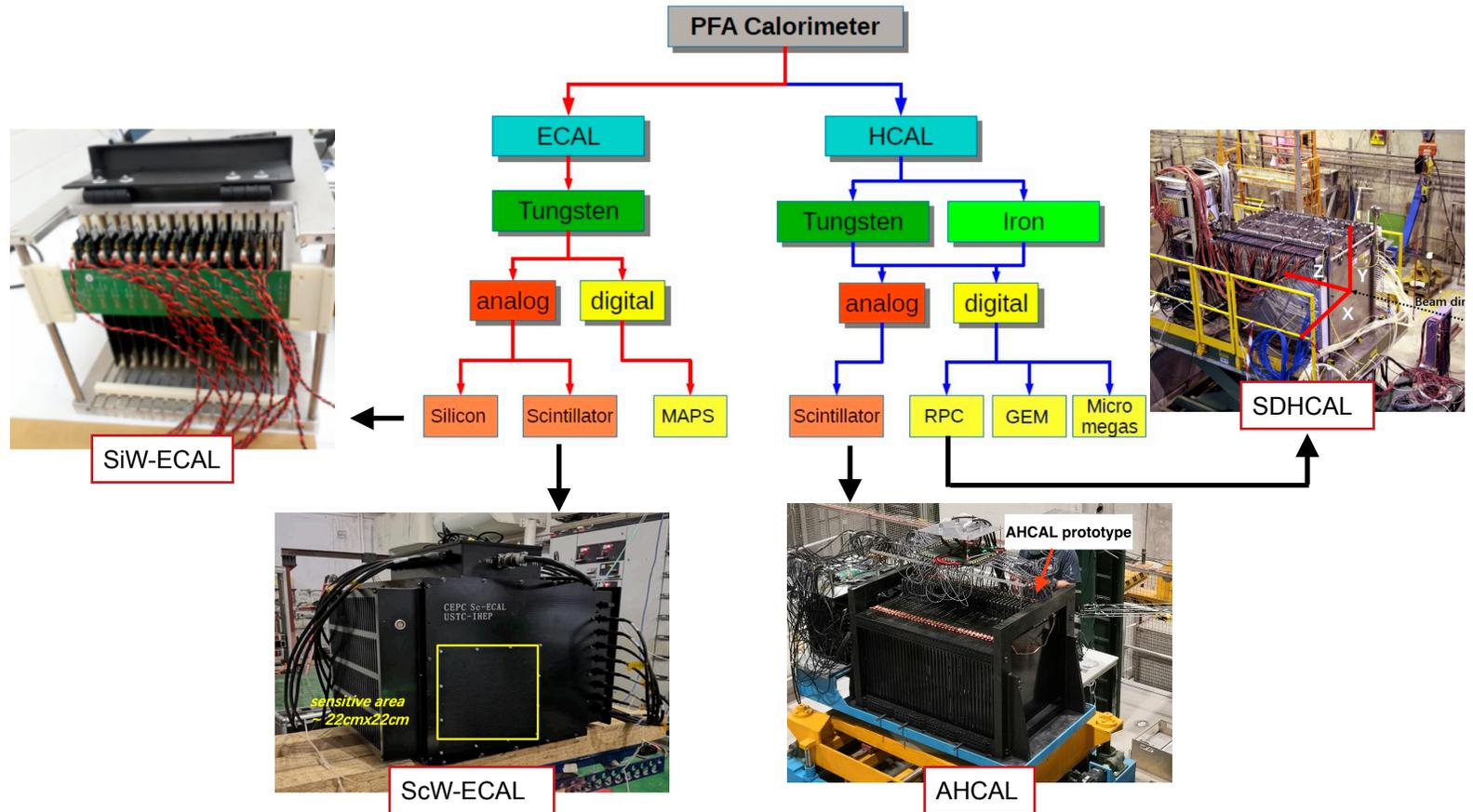
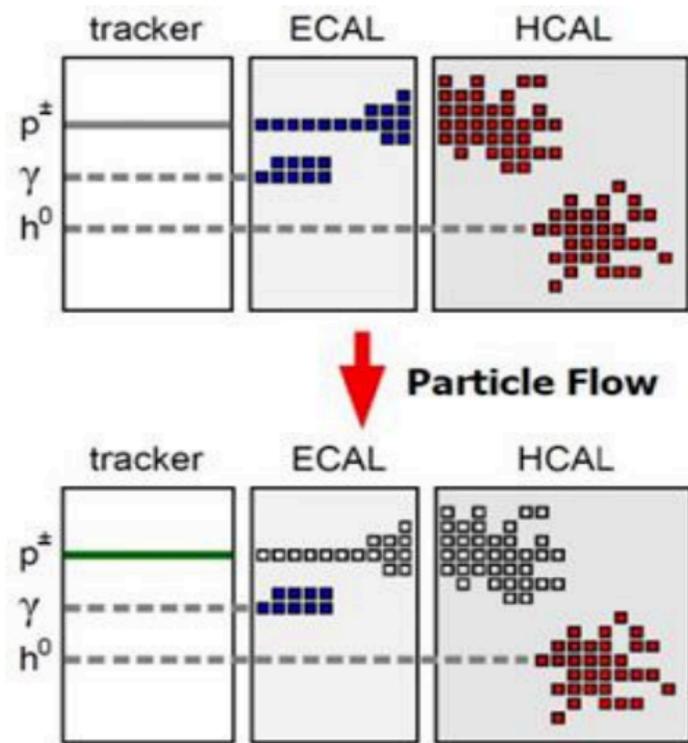


- **Motivation**
- **MC samples**
- **PID based on BDT**
- **PID based on ANN**
- **PID application on Beam data**

- **Future Higgs/W/Z/Top factories require high jet energy resolution**
 - Particle identification is essential for jet tagging and measurement.
 - Accurate differentiation between hadronic and electromagnetic showers is crucial for determining the energy and type of each particle.



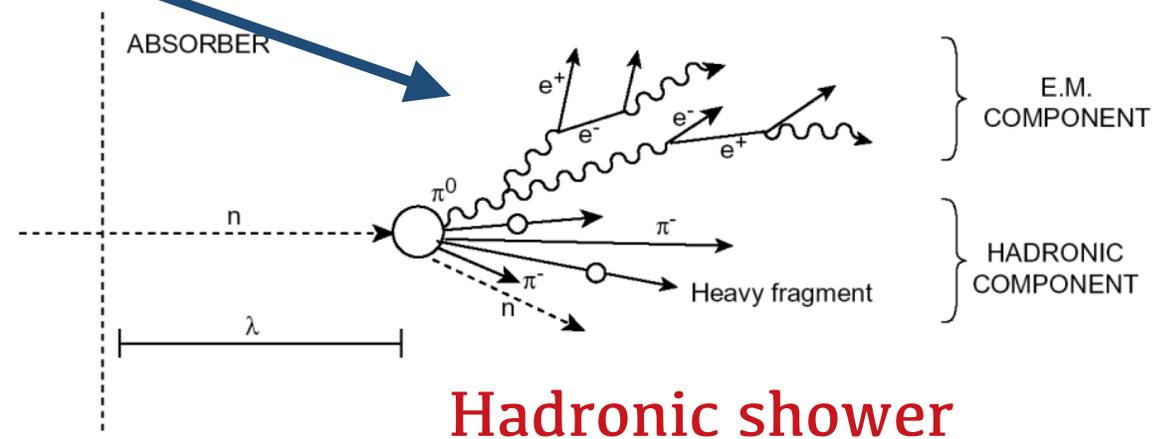
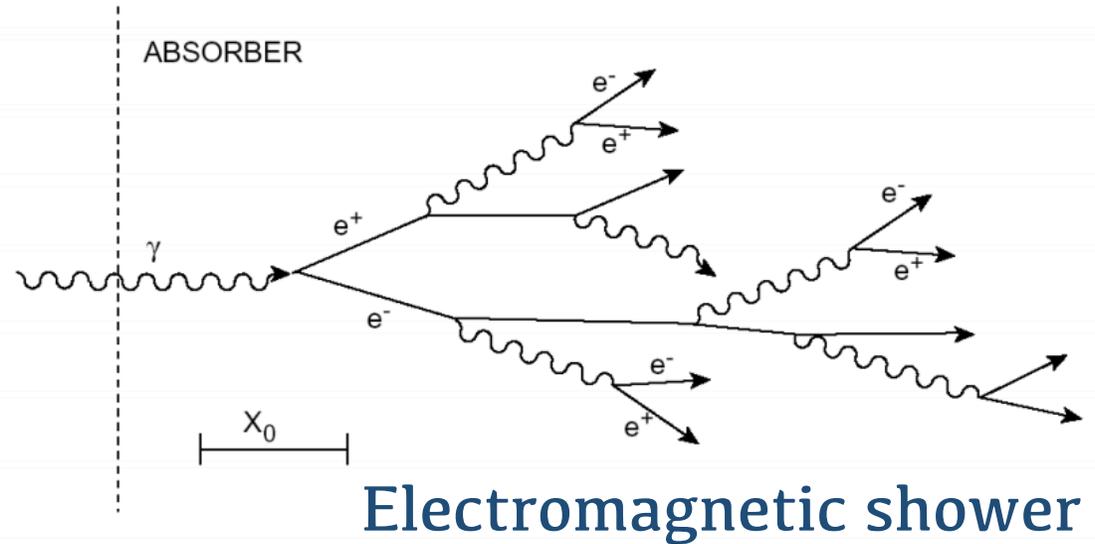
- **Exploration for ML PID methods in PFA-oriented calorimeters**
 - A major calorimetry option of future Higgs/W/Z/Top factories:
 - High granularity (imaging) -> A “camera”.



Overview of world-wide development of high-granularity calorimeters

• Challenge:

- Hadronic showers may undergo secondary interactions and exhibit complex development patterns.
- Hadronic showers also contain EM component.
- Develop effective shower topology variables.



• CEPC AHCAL prototype parameter

- Geometry

- 40 sampling layers.
- 72cm × 72cm in transversal plane.
- 120cm in longitudinal direction.

- Absorber

- 2 cm thickness/layer steel.

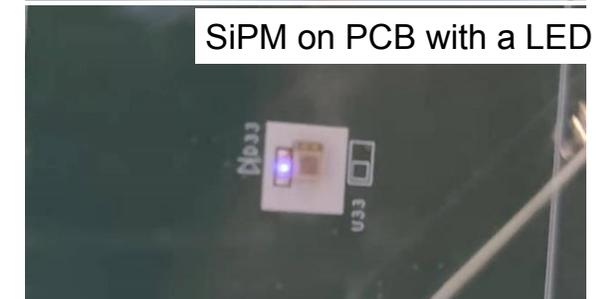
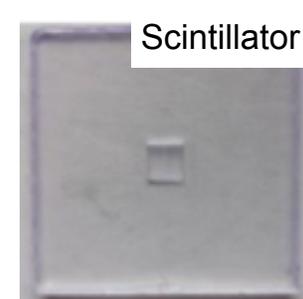
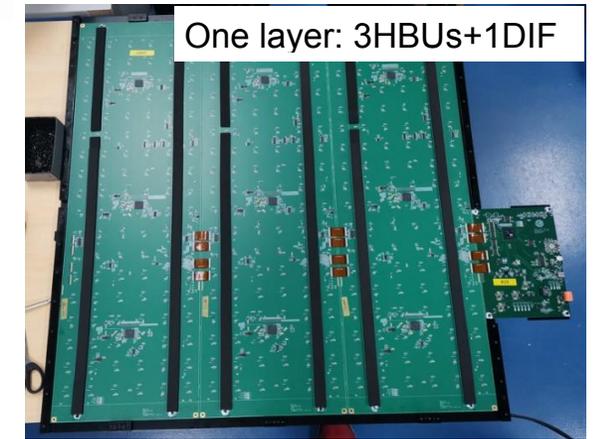
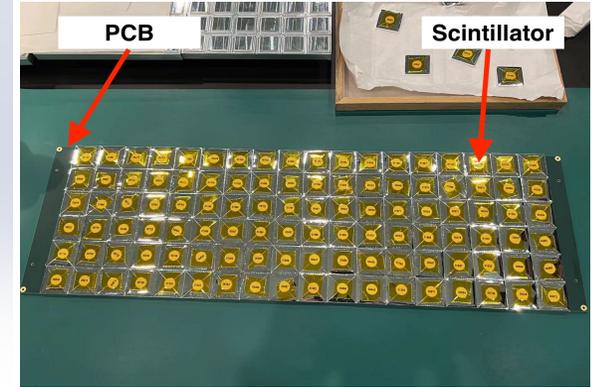
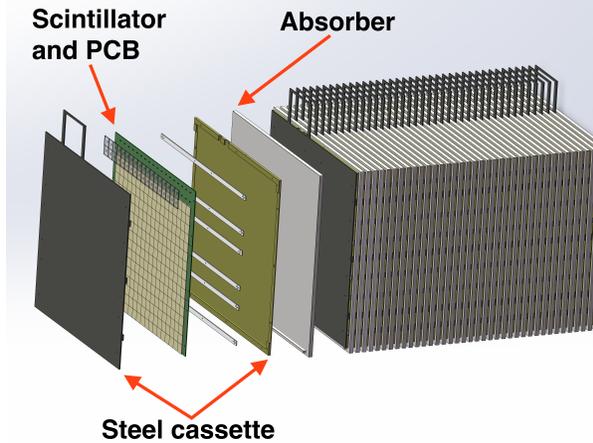
- Sensitive cells

- 40mm × 40mm × 3mm scintillator tile coupled with SiPM (SiPM-on-tile).
- **18 × 18 × 40** array.



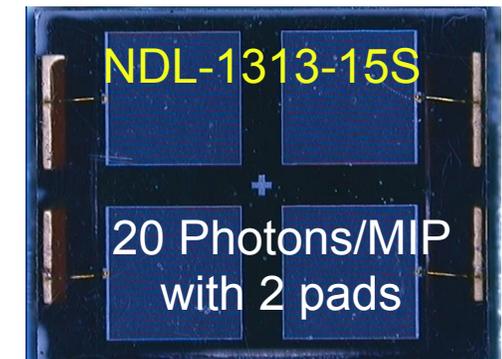
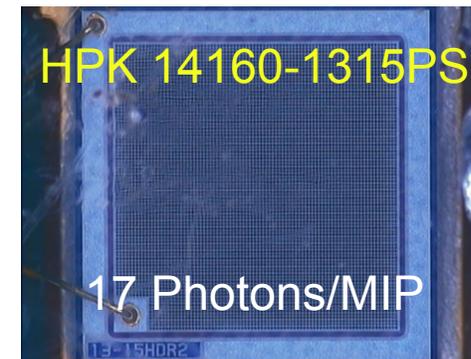
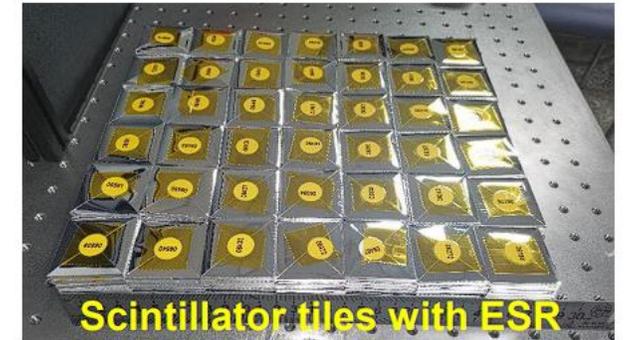
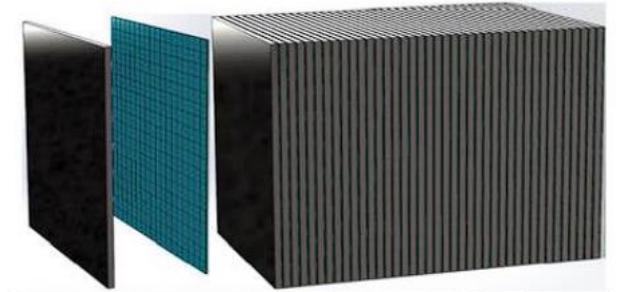
Spatial distribution of hits viewed as a 3D image

Journal of Instrumentation, 2021, 16(03): P03001.
Journal of Instrumentation, 2022, 17(11): P11034.



Simulation set up

- the Geant4 11.1.1 Toolkit with the QGSP_BERT physics list was employed.
- 2 mm plastic scintillator + 0.25 mm \times 2 ESR + 20 mm Steel.
- **Digitization:**
 - Photon statistics: Poisson distribution concerning #detected photons (light output).
 - SiPM saturation : $response = \#pixel \times e^{-\frac{photon}{\#pixel}}$.
 - ADC error: assume 0.02%, very low.
 - Energy cut: 0.5 MIP.
- **SiPM:**
 - S14160-1315PS for first 38 layers.
 - EQR15 22-1313D-S for last 2 layers.



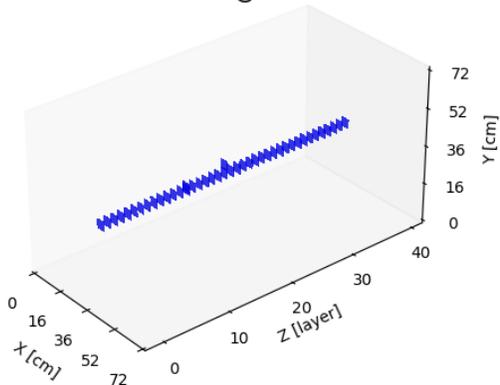
<https://indico.cern.ch/event/847884/contributions/4831207/>



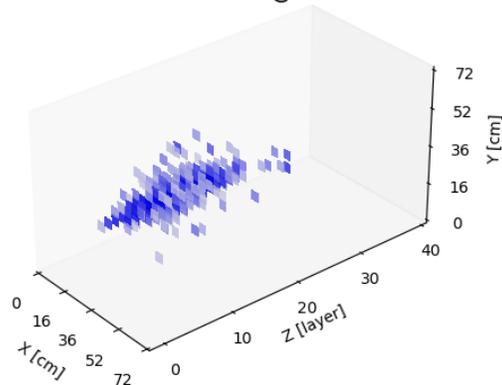
- **Generally, MC and Data are close in shower profile.**

- Data come from beam test at SPS-H2, CERN.
- Several shower topology variables are reconstructed:
 - **Shower density:** Mean hits number in a 3×3 cell.
 - **Shower length:** Distance between the start of the shower and the layer with maximum RMS of hit transverse coordinates.

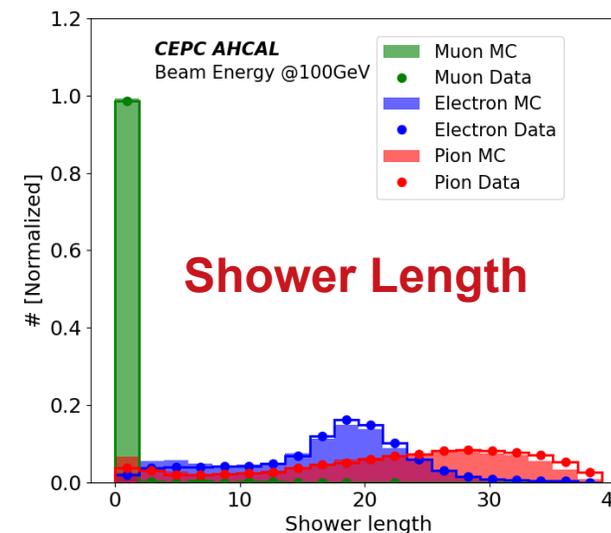
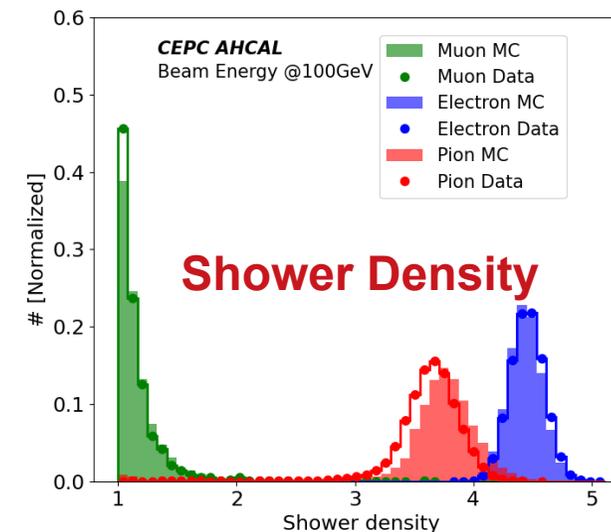
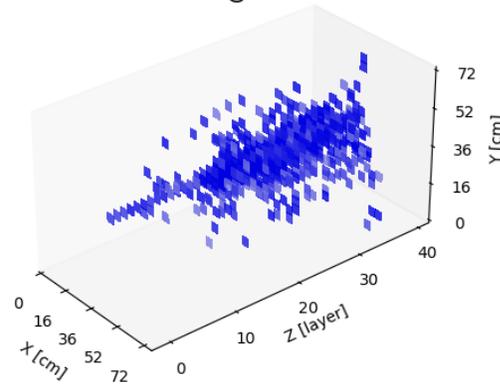
CEPC AHCAL
Muon Simulation @100GeV



CEPC AHCAL
Electron Simulation @100GeV



CEPC AHCAL
Pion Simulation @100GeV

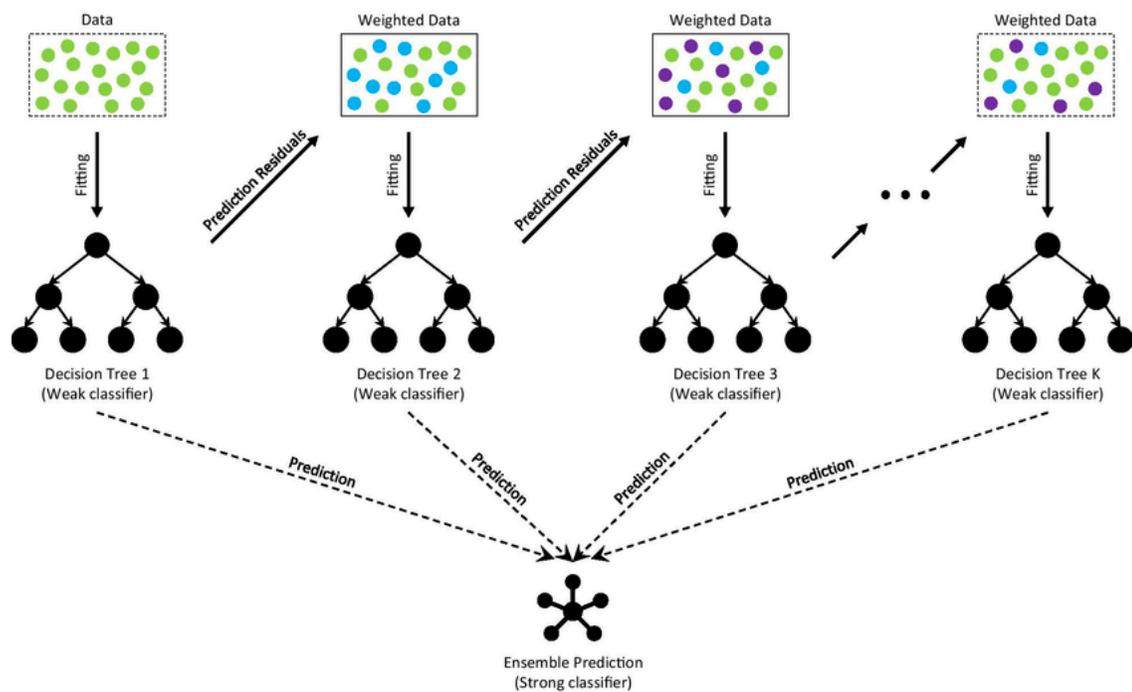


• Monte Carlo Samples

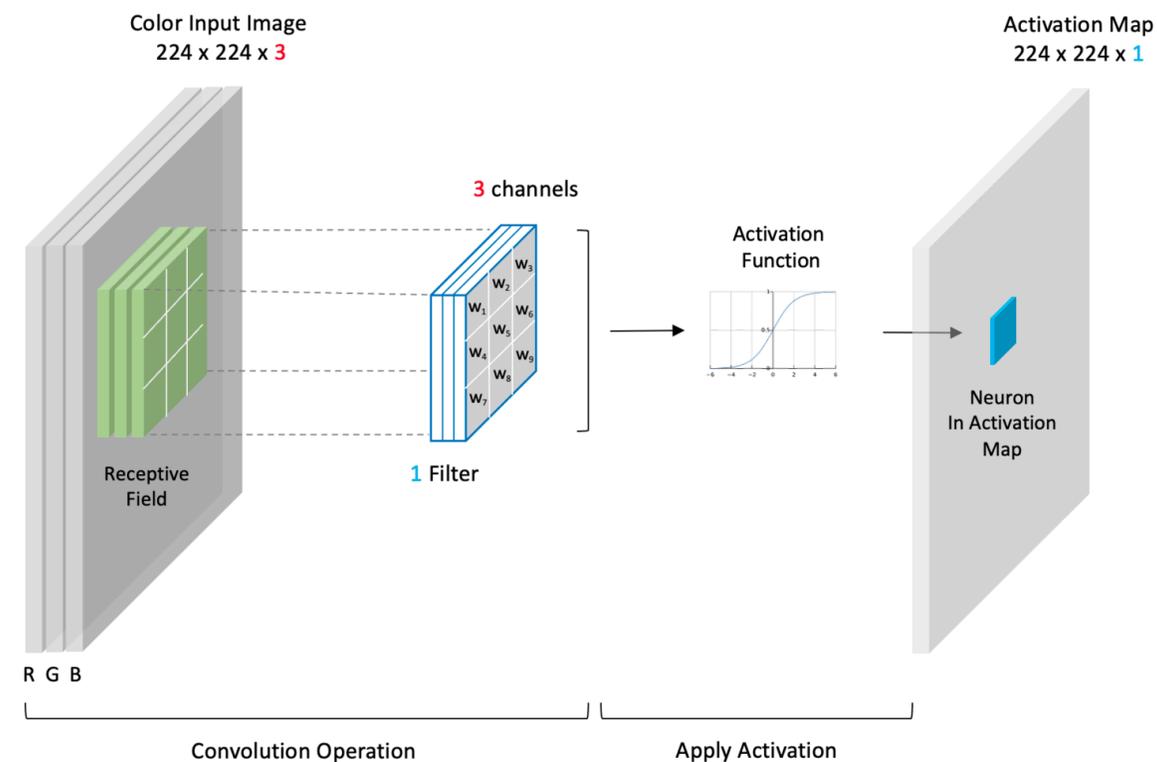
- To study the separation power in HAD showers and EM showers.
- Training set : Validation set : Test set = 5 : 1 : 4 .

Energy	5GeV	10GeV	30GeV	50 GeV	60GeV	80GeV	100GeV	120GeV
Electron	100k	100k	100k	100k	100k	100k	100k	100k
Pion-	100k	100k	100k	100k	100k	100k	100k	100k

- **Boosted decision tree (BDT)**
 - Need pre-reconstructed input.
- **Cell-based artificial neural networks (ANN)**
 - Treat spatial distribution of hits as images.



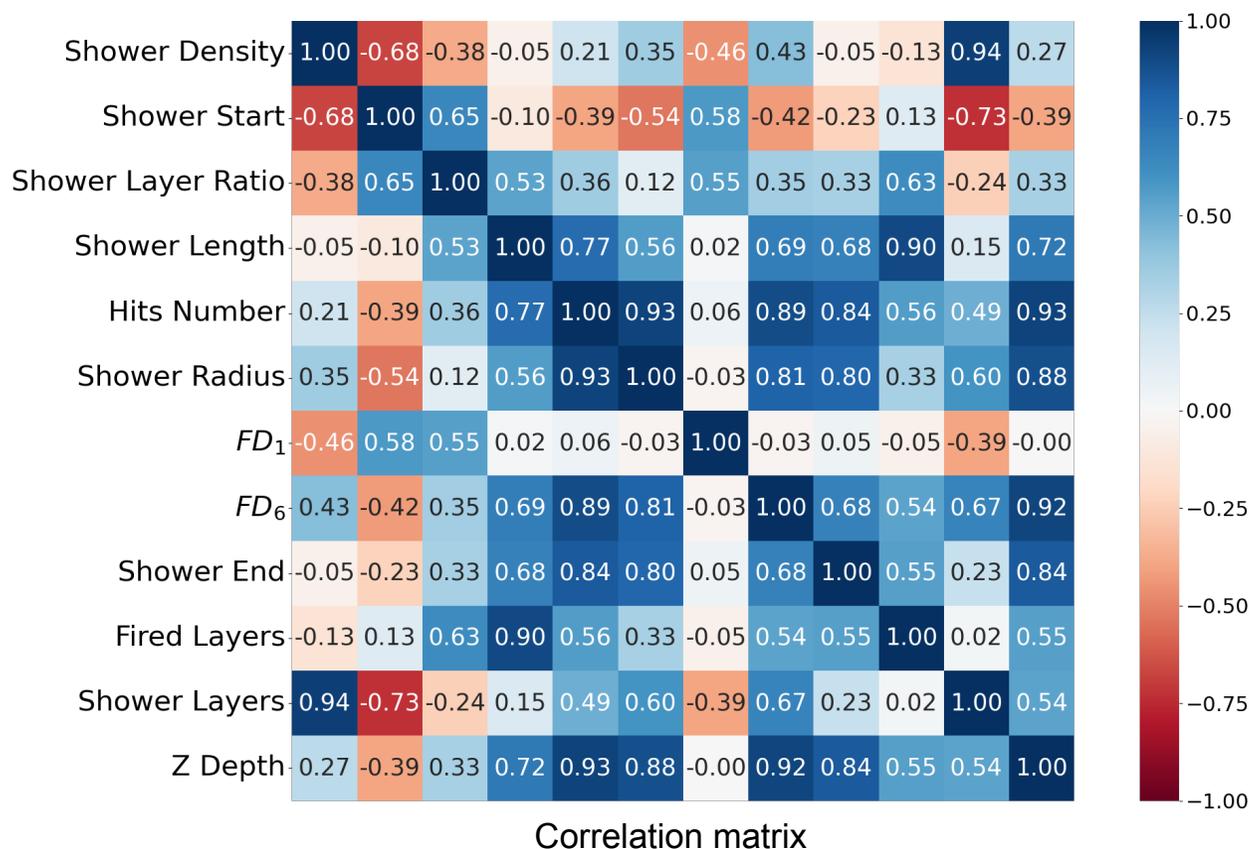
Boost decision tree



Neural networks

- **Apply Extreme Gradient Boosting (XGBoost).**

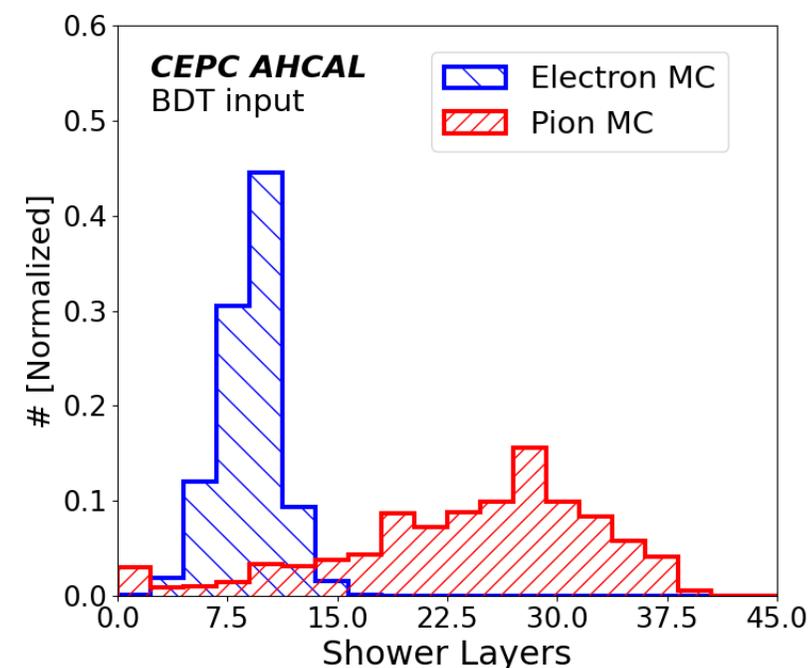
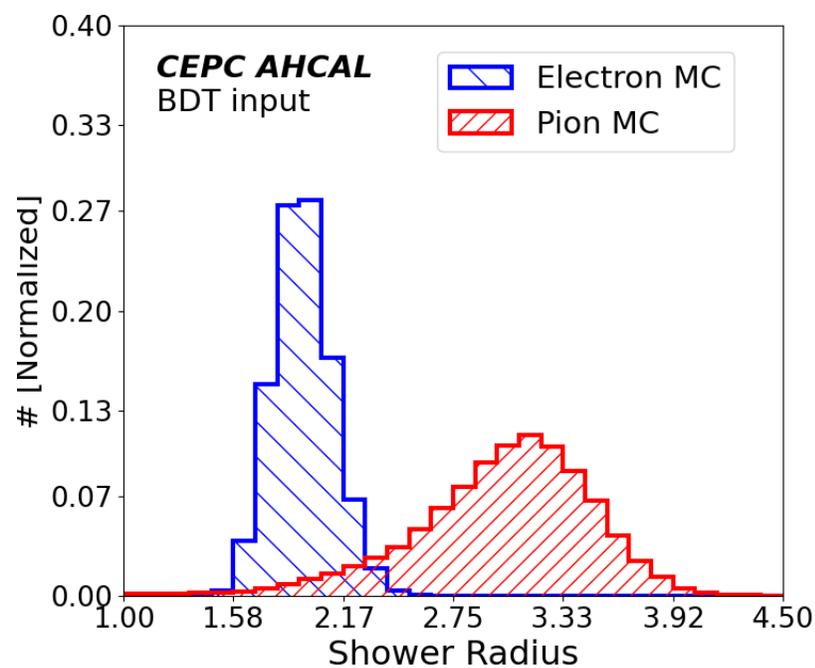
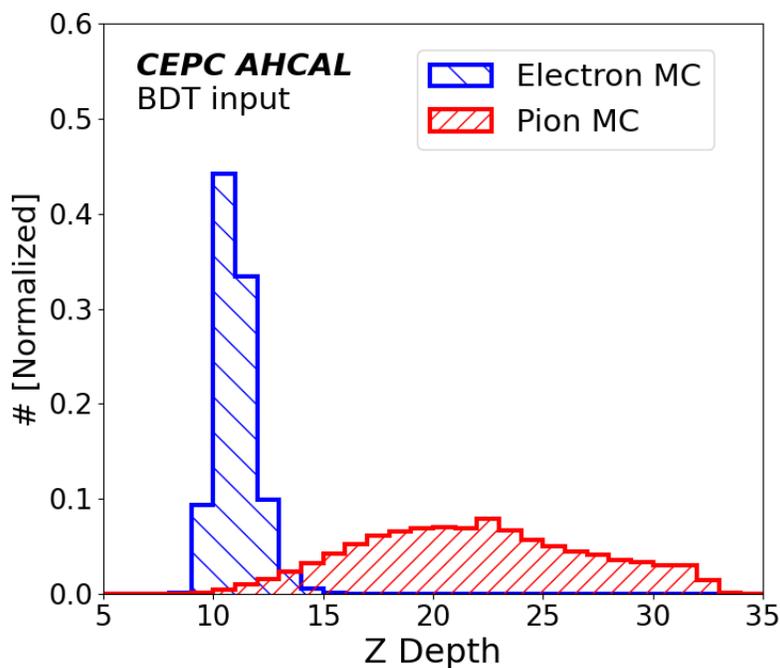
- 12 variables are reconstructed.
- Z depth, Shower radius, and Shower layers are top 3 variables.



Rank: Variable	Variable weight
1: Z depth	0.532
2: Shower radius	0.186
3: Shower layers	0.073
4: Fired layers	0.065
5: Shower density	0.370
6: Shower start	0.026
7: Shower layer ratio	0.022
8: FD ₁	0.018
9: Hits number	0.013
10: FD ₆	0.012
11: Shower end	0.009
12: Shower length	0.006

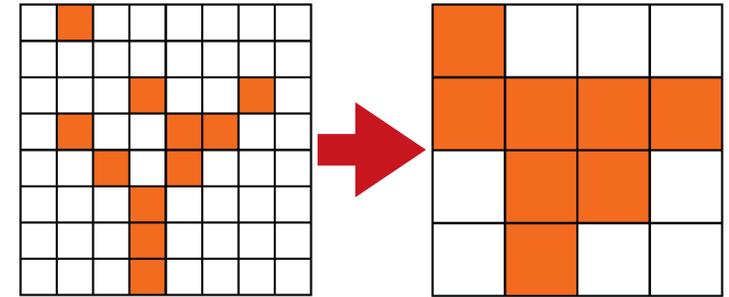
- **Top 3 variables in separating between EM showers and HAD showers.**

- **Z depth:** The RMS of the z-axis coordinates.
- **Shower Radius:** The RMS of the distance with respect to the z-axis.
- **Shower layers:** The number of layers in which the RMS of positions in the x-y plane exceeds 4 cm.



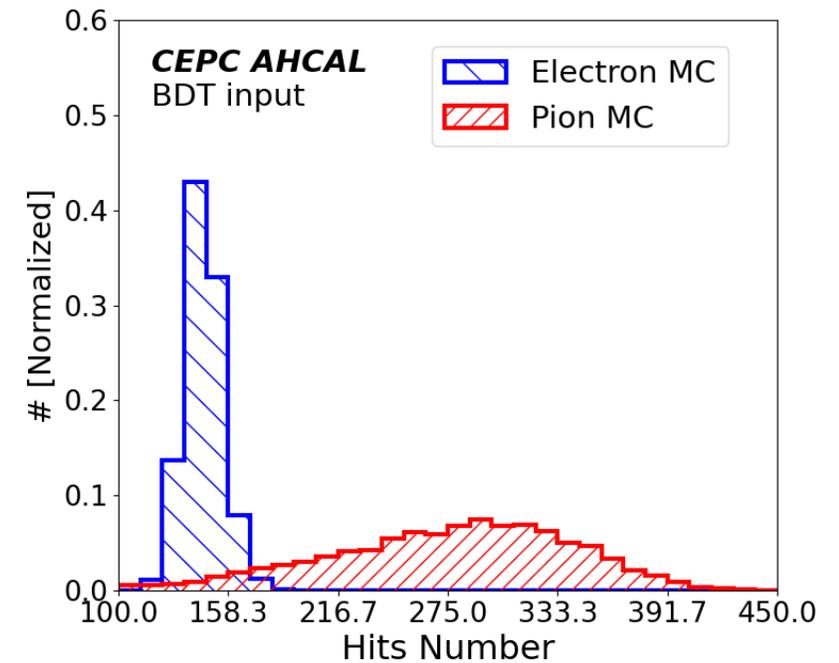
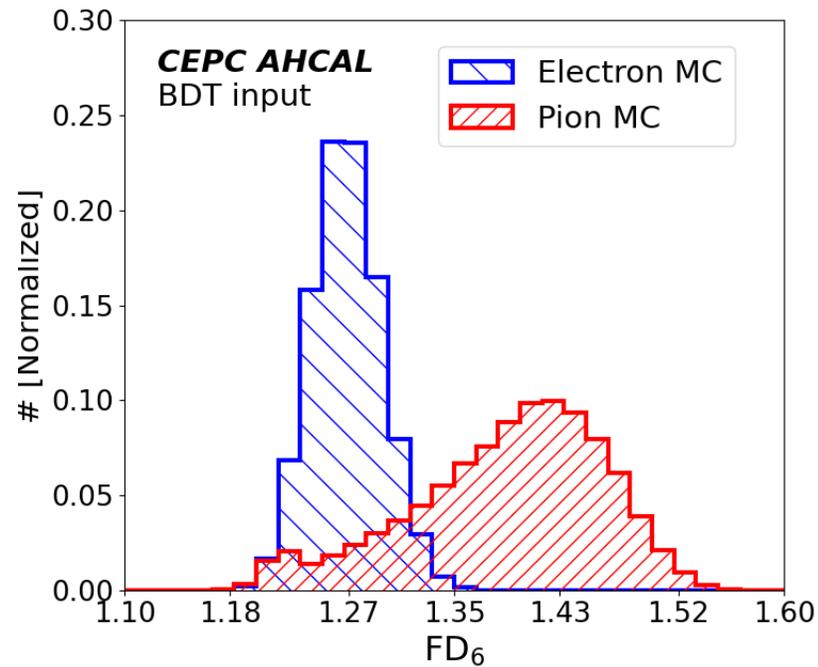
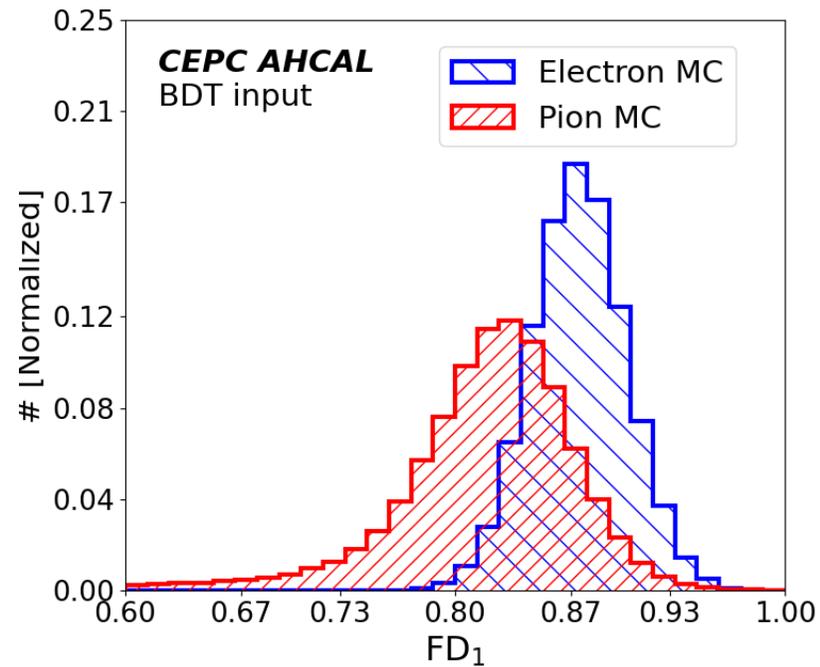
- **Fractal dimension:** $FD_\beta = \left\langle \frac{\log(R_{\alpha,\beta})}{\log(\alpha)} \right\rangle + 1$, where $R_{\alpha,\beta} = N_\beta/N_\alpha$.
- **Hits number:** The number of hits.

- N_α : number of hits scaled by α .



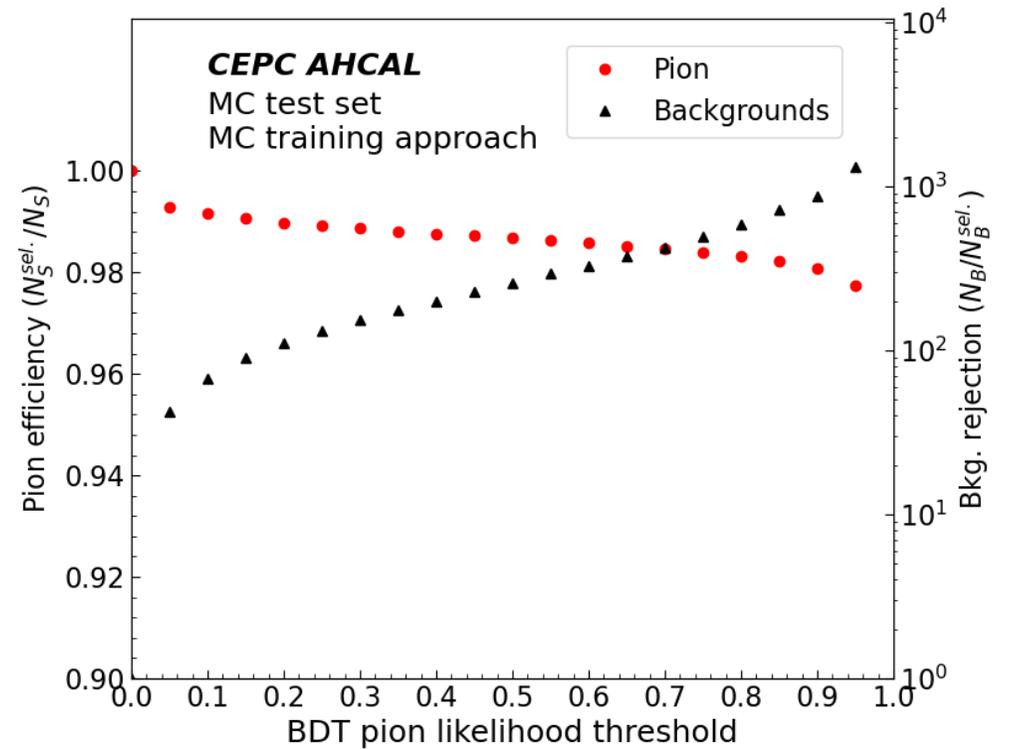
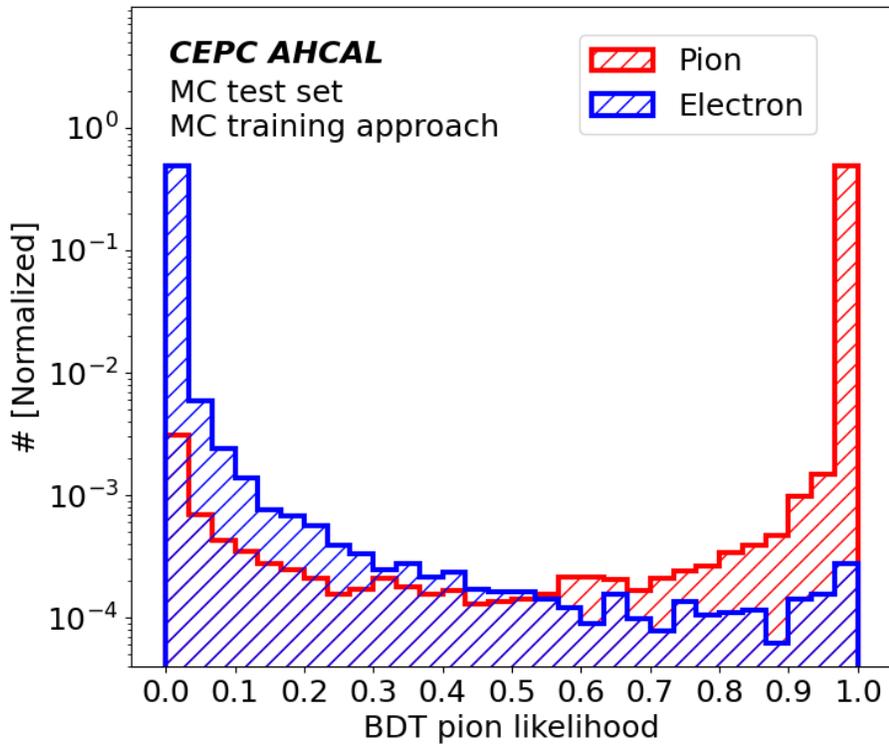
Take $\alpha = 2$ as an example

FD Ref: PhysRevLett.112.012001

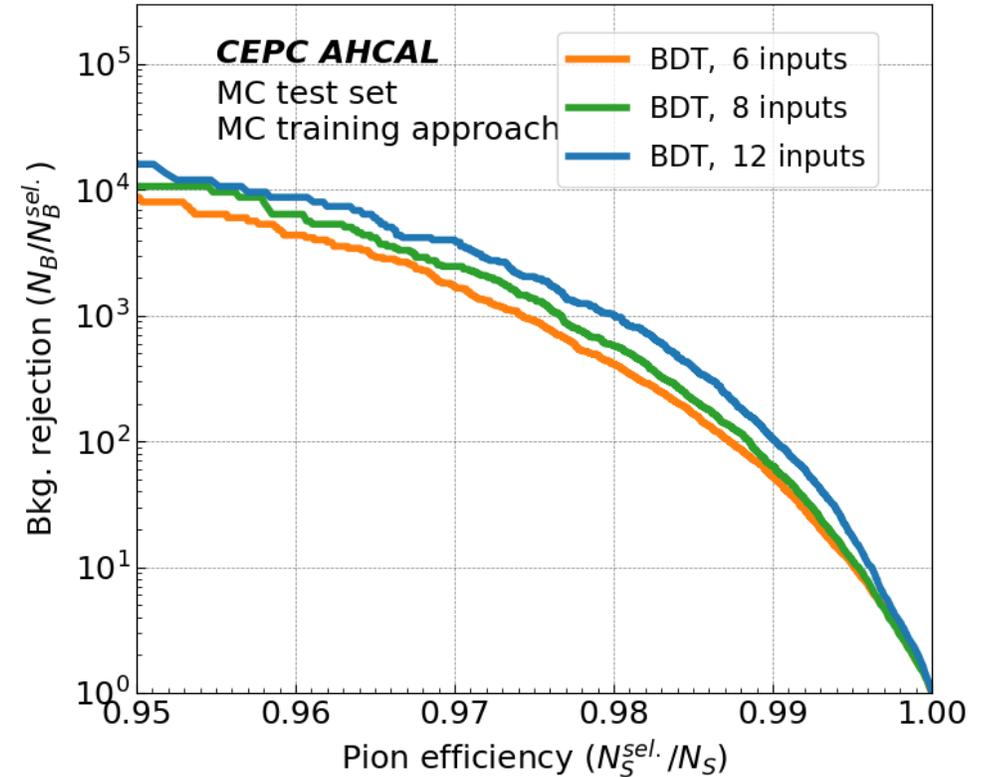


BDT classifier performance

Pion efficiency	95%	96%	97%	98%	99%
Electron rejection (1/e efficiency)	16012.7	8734.2	3843.0	970.5	105.3

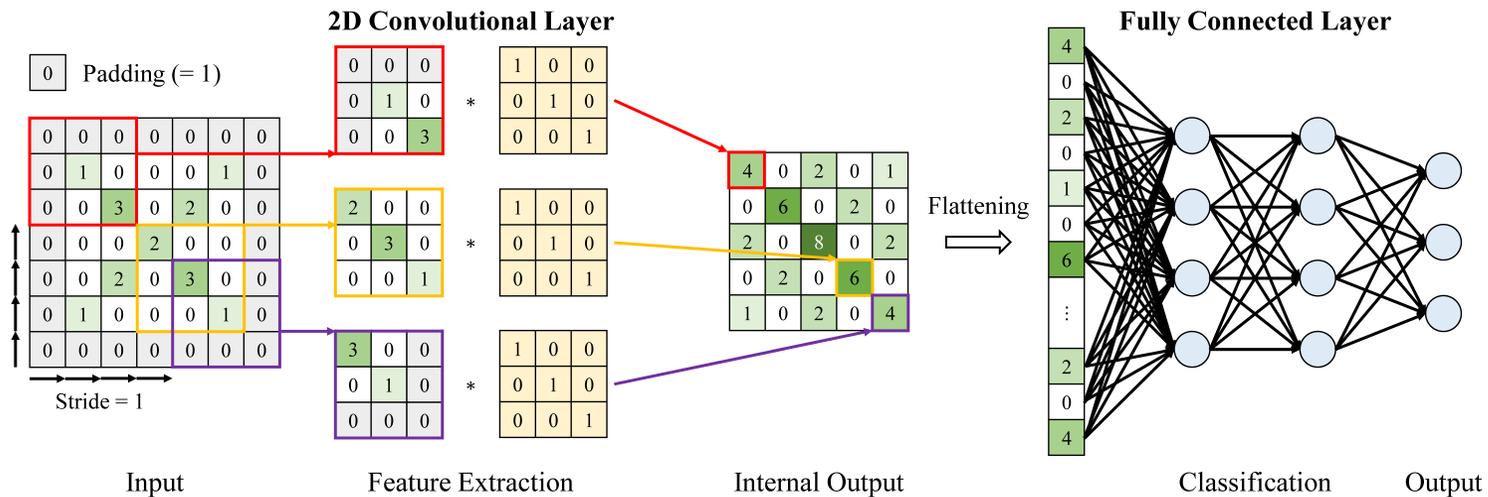


- **We observe dependence of BDT performance on input variables**
 - Remove Shower End, Shower Layers, Fired Layers, and Z Depth to build BDT with 8 inputs.
 - Further remove FD_1 and FD_6 to build BDT with 6 inputs.
- **Feature engineering can be sometimes tricky.**



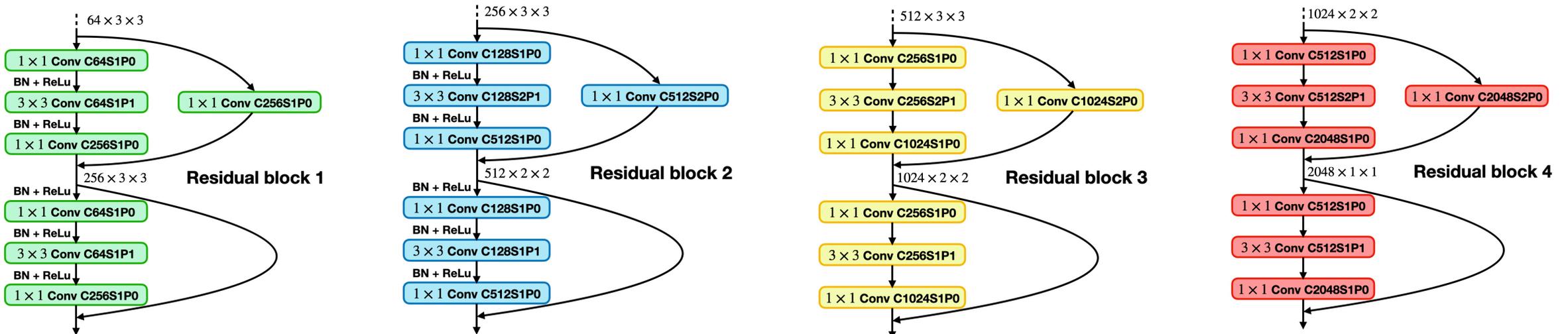
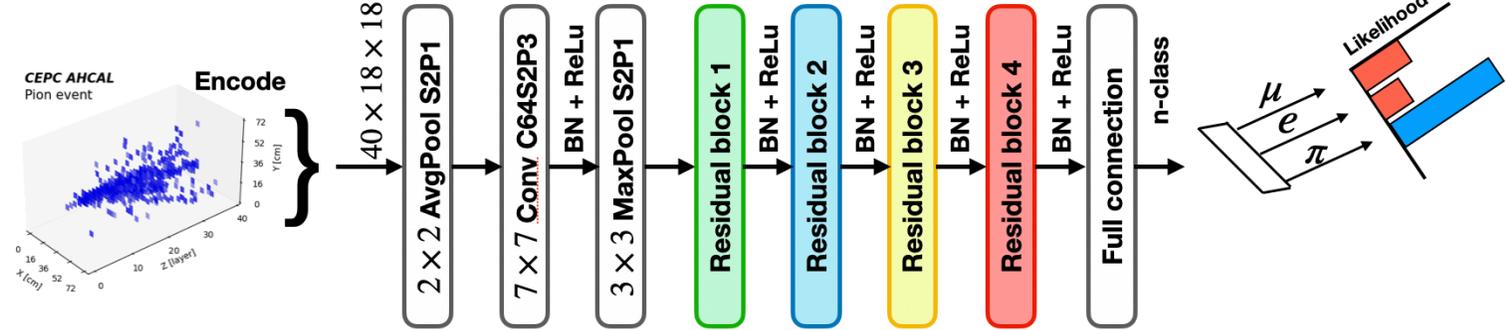
- **BDT optimization is still on going.**

- **Cell-based Artificial Neural Networks (ANN) make full use of high-dimensional input.**
 - Compile layers to extract features.
 - **Input:** Variable stands for events (Spatial distribution of hits).
 - **Output:** the likelihood of each particle type candidate.
 - After iteration, the mapping: **Input**->**Output** close to truth.



• ANN-based PID: Taking the advantage of ResNet

- Input: energy depositions in AHCAL (Input tensor size: $40 \times 18 \times 18$).
- Output: likelihood of each particle type candidate.



ResNet Ref: He K, Zhang X, Ren S, et al. Deep residual learning for image recognition[C]//Proceedings of the IEEE conference on computer vision and pattern recognition. 2016: 770-778.

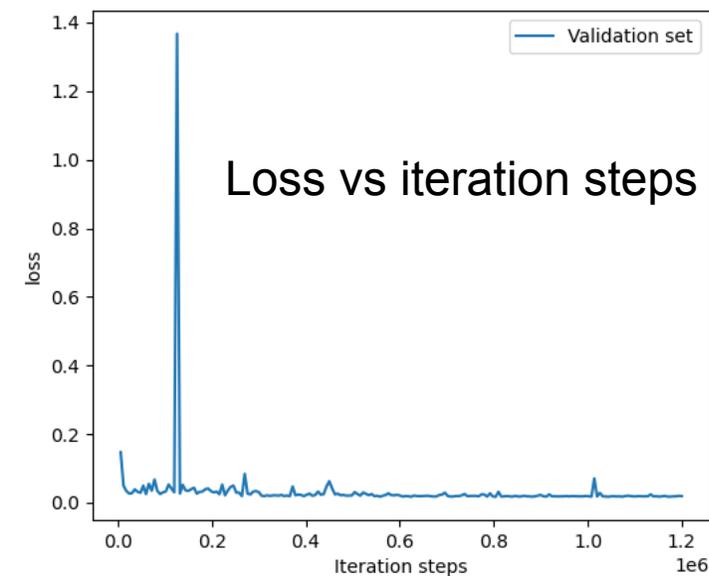
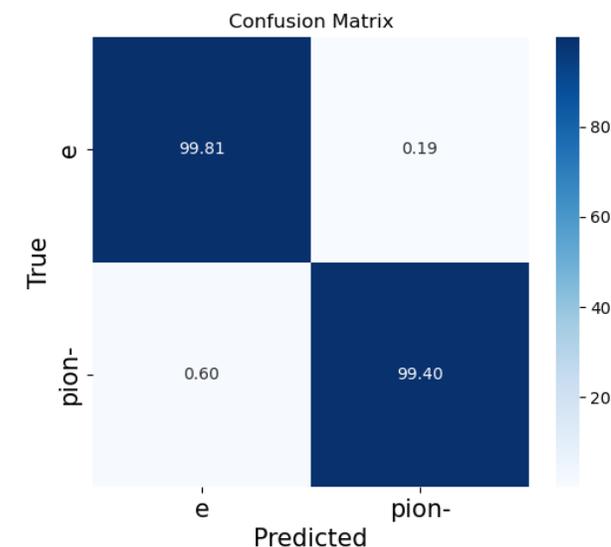


- Executed on an NVIDIA V100 NVLink GPU.
- Loss: Cross Entropy.
- Hyper-parameter
 - Batch size= 64
 - lr=0.0001
 - epoch = 200

Algorithm 1 Artificial Neural Network.

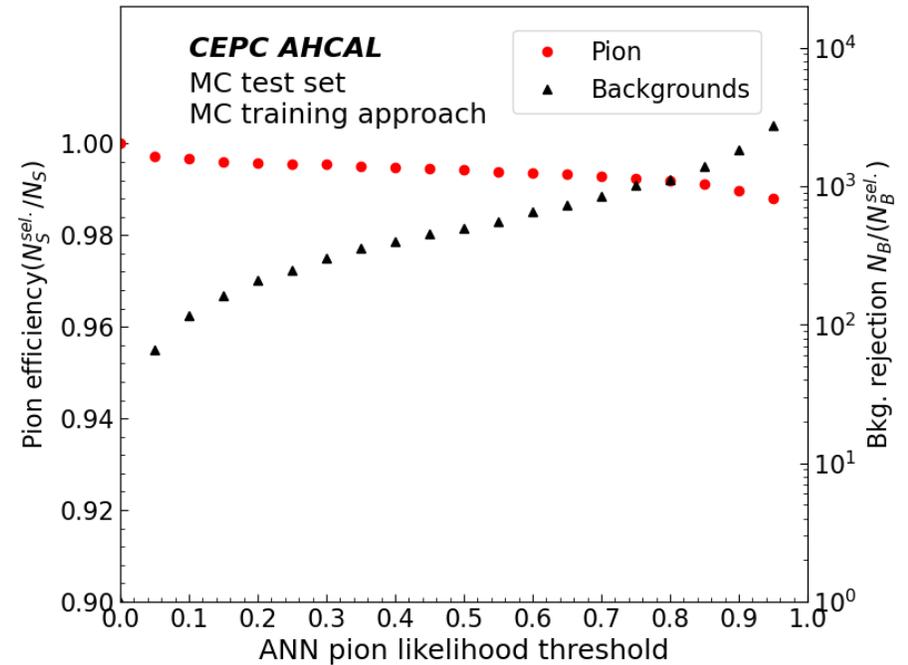
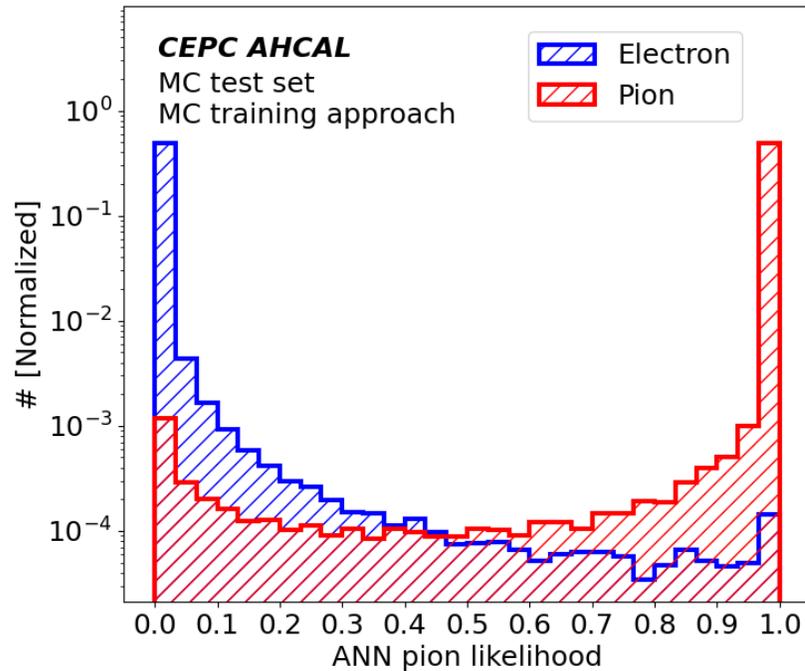
Require: The batch size m , the epoch number n , initial learning rate lr , initial net parameters θ_0 .

- 1: Assign corresponding label y to data x .
 - 2: **for** $t = 1, \dots, k$ iteration steps **do**
 - 3: **for** $i = 1, \dots, m$ **do**
 - 4: $\hat{y} \leftarrow \text{Net}(x, \theta)$
 - 5: $\text{Loss}(y_i, \hat{y}_i)_\theta^{(i)} \leftarrow (-\log(\hat{y}_i))$
 - 6: **end for**
 - 7: $\theta \leftarrow \text{SGD} \left(\nabla_\theta \frac{1}{m} \sum_{i=1}^m \text{Loss}_\theta^{(i)} \right)$
 - 8: **end for**
-



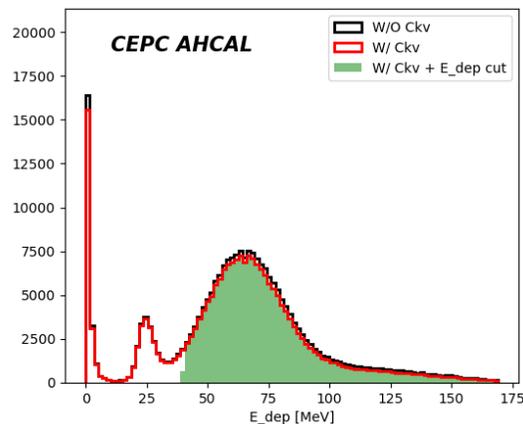
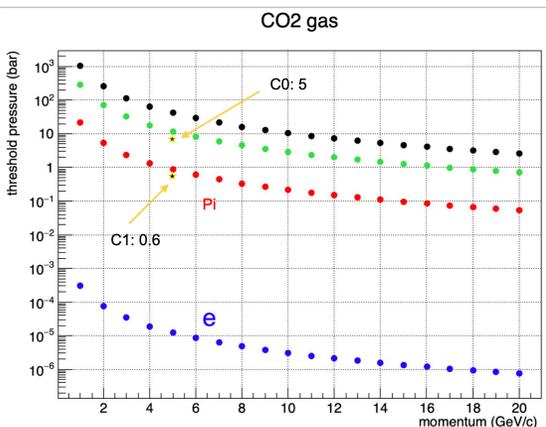
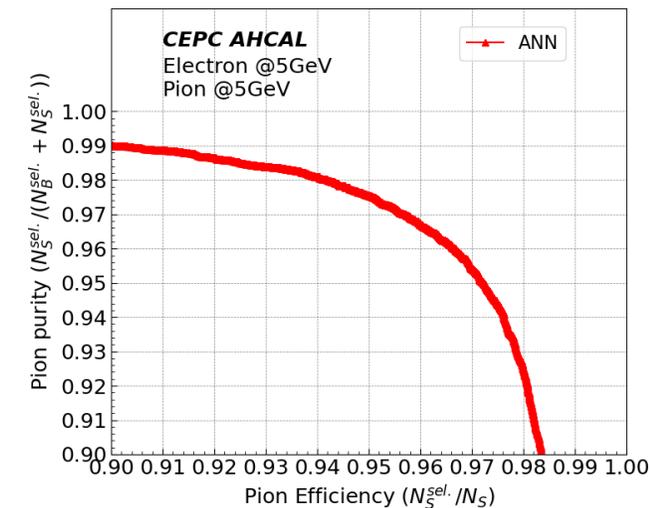
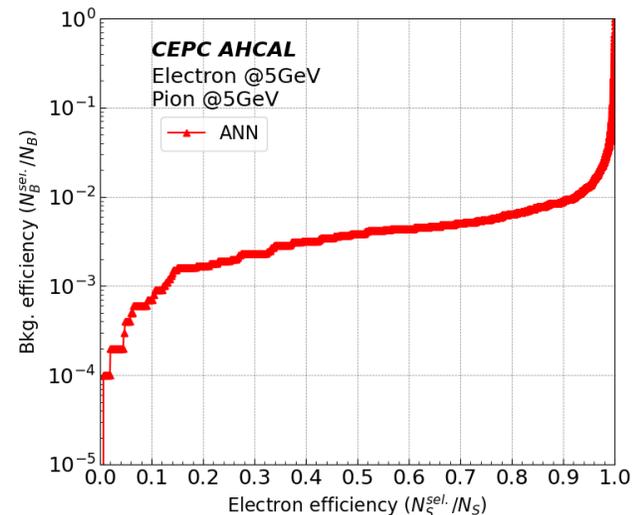
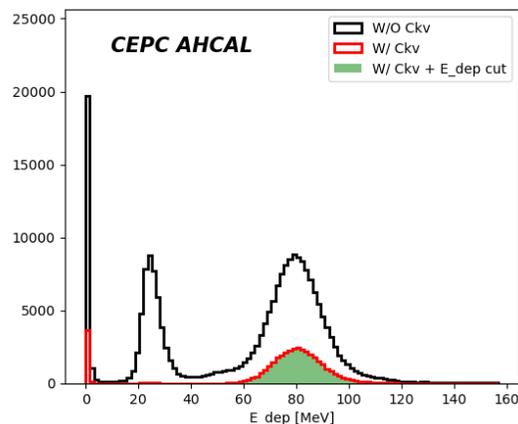
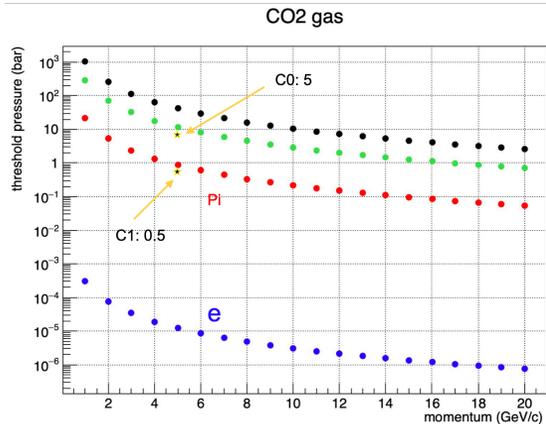
ANN classifier performance

Pion efficiency	95%	96%	97%	98%	99%
Electron rejection (ANN)	56012.7	48010.9	19769.2	9083.1	2154.3
Electron rejection (BDT)	16012.7	8734.2	3843.0	970.5	105.3
Improvement	199.8%	249.8%	393.8%	752.1%	1945.0%



ANN PID cross-check using Cherenkov detectors

- Two CO2 Cherenkov detectors are available at PS (<15 GeV).
- 20,000 Electron and 20,000 Pion samples are selected as truth.



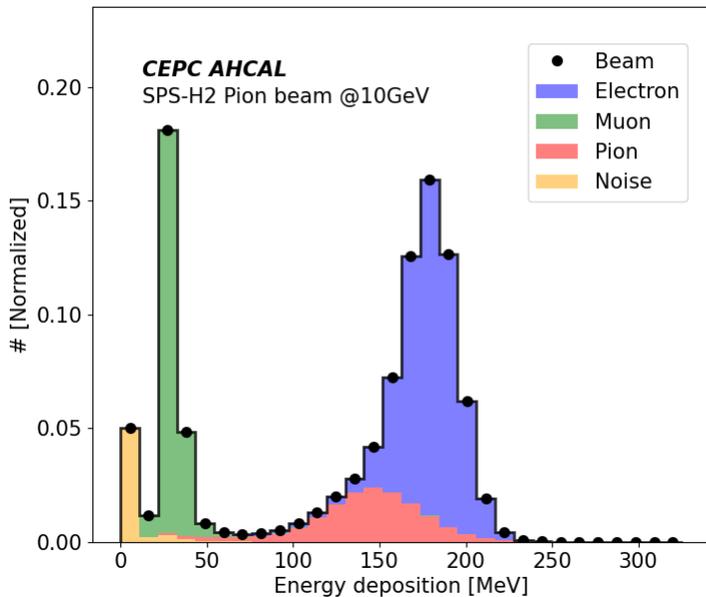
- Achieve 90% Pion efficiency and 99% Pion purity at the same time.

Cherenkov detector setting strategy.

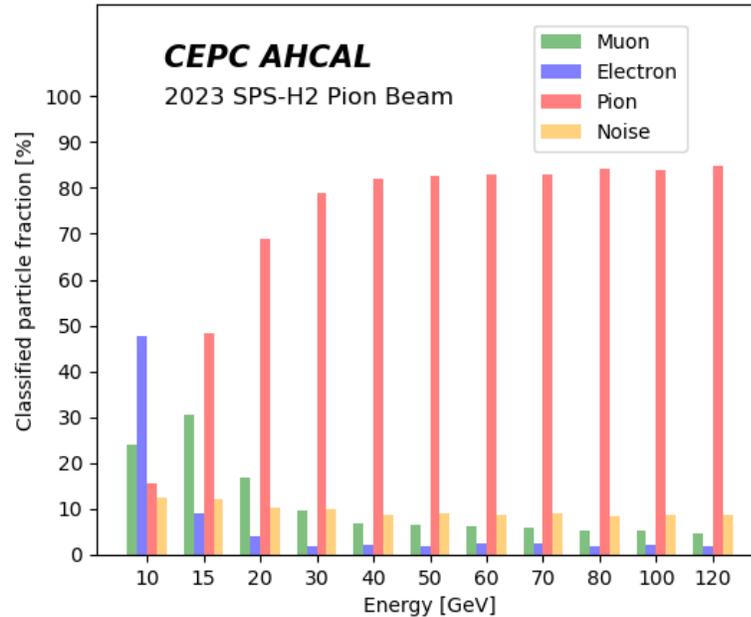
Samples after selection.

- **Beam composition is given by ANN classifier**

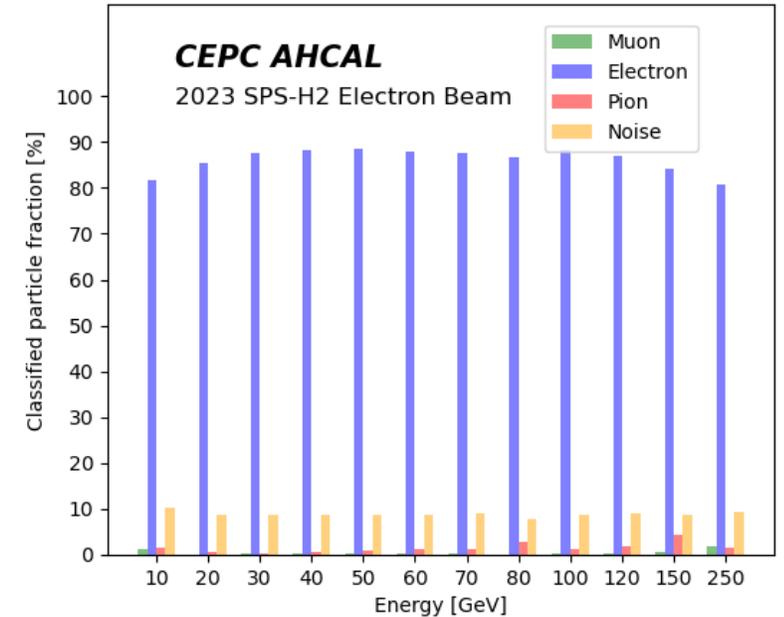
- Trained on pre-selected data also tagged by Cherenkov detector.
- Pion beam purity: around 80% when beam energy is over 30 GeV.
- Electron beam purity: over 80% at each energy point.



Beam composition in Energy deposition spectrum (10 GeV)



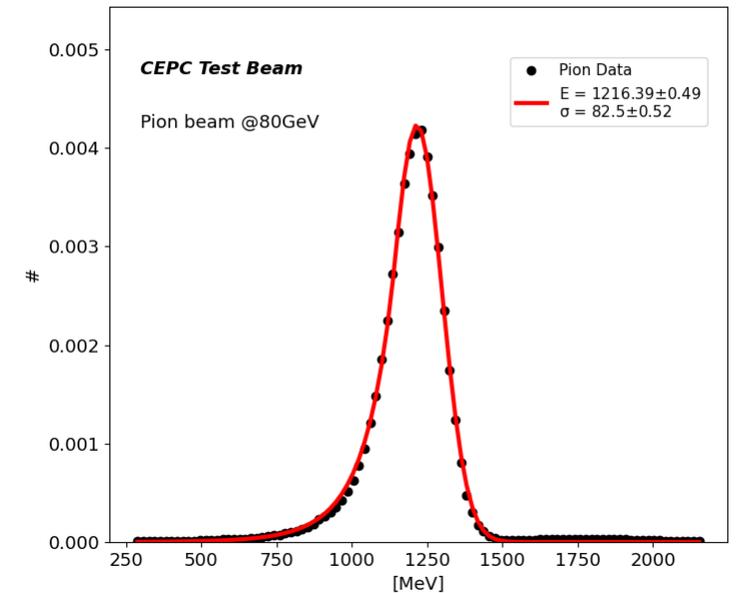
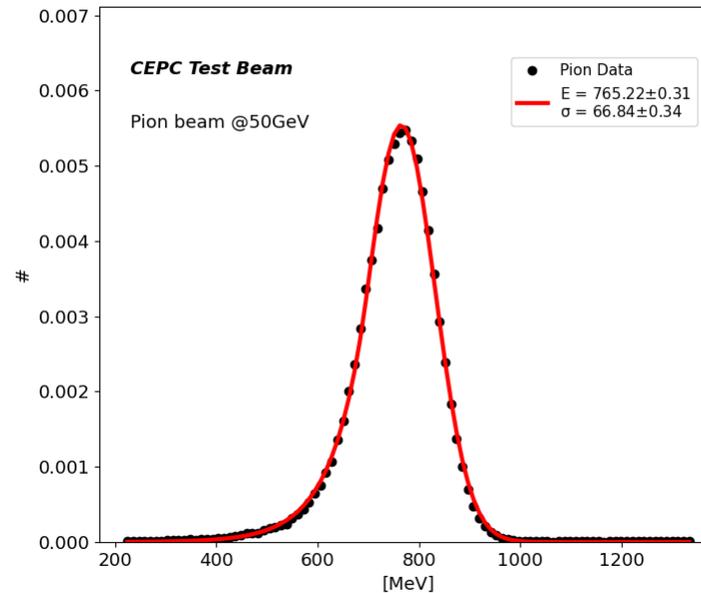
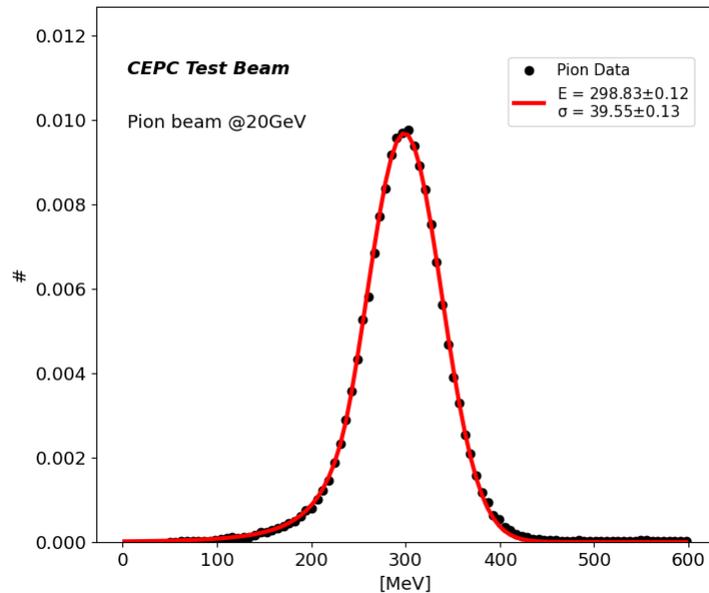
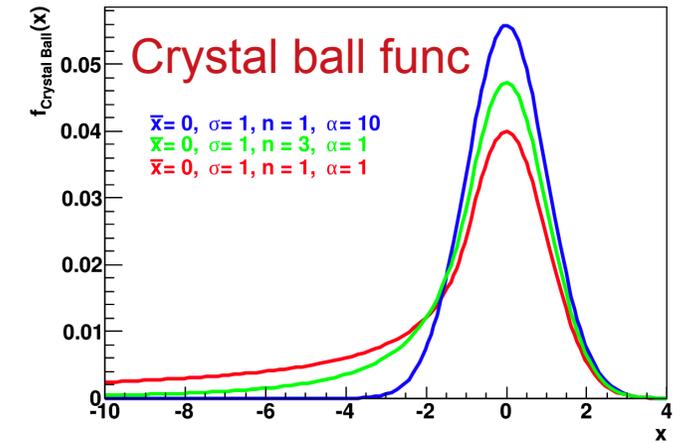
Pion Beam composition



Electron Beam composition

• Purified beam data fitting

- ANN classier is first used for purifying the Pion beam.
- Crystal ball function is then used for fitting purified Pion Data.



The energy of 20, 50, 80 GeV.

- **ANN(ResNet) outperforms BDT in our cases.**
 - **Automatic feature extraction:** This allows ANN to make full use of all input information, and potentially uncover hidden patterns in the data that may be missed by BDT, which relies on limited reconstructed features.
 - **Effective in handling large and high-dimensional inputs:** ANN is well-equipped to handle high-dimensional data and capture complex patterns within it.
 - **Non-linearity:** ANN can model complex non-linear relationships in data more effectively than BDT.
- **More validation is still on going.**

Thank you



感谢聆听

饮水思源 爱国荣校²⁴



Beam Test at CERN in 2022 & 2023



• SPS-H8: Oct 19 - Nov 2, 2022:

- μ^+ : 160 GeV (for calibration)
- π^+ : 10, 20, 30, 40, 50, 60, 70, 80, 90, 100, 120 GeV (~1M events each point)
- e^+ : 10, 20, 30, 40, 50, 60, 70, 80, 90, 100, 120 GeV (~0.3M events each point)

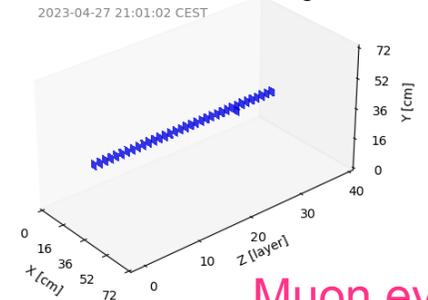
• SPS-H2: Apr 26 - May 10, 2023:

- μ^- : 100 GeV (for calibration)
- π^- : 10, 15, 20, 30, 40, 50, 60, 70, 80, 100, 120, 350 GeV (~1M events each point)
- e^- : 10, 20, 30, 40, 50, 60, 70, 80, 100, 120, 150, 250 GeV (~0.3M events each point)

• PS-T9: Apr 17 - May 31, 2023:

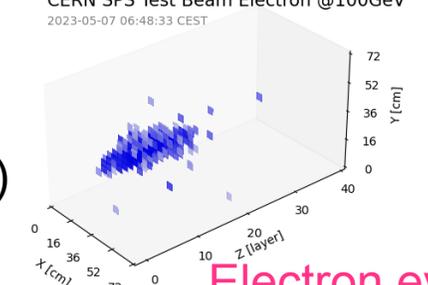
- μ^- : 10 GeV (for calibration)
- π^- : 1, 2, 3, 4, 5, 6, 7, 8, 10, 12, 15 GeV (~0.5M events each point)
- e^- : 1, 2, 3, 4, 5 GeV (~50K events each point)

CEPC AHCAL
CERN SPS Test Beam Muon @100GeV
2023-04-27 21:01:02 CEST



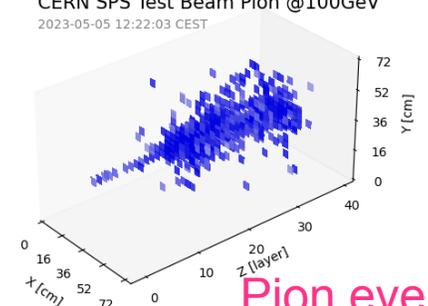
Muon event

CEPC AHCAL
CERN SPS Test Beam Electron @100GeV
2023-05-07 06:48:33 CEST



Electron event

CEPC AHCAL
CERN SPS Test Beam Pion @100GeV
2023-05-05 12:22:03 CEST



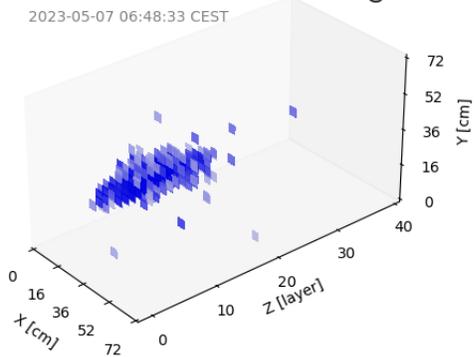
Pion event

- **Successful beam test for AHCAL and ScW-ECAL was conducted.**

- Beam test at CERN SPS & PS in 2022, 2023.
- Beam: Muon, Electron, Pion, Proton (1-350GeV).

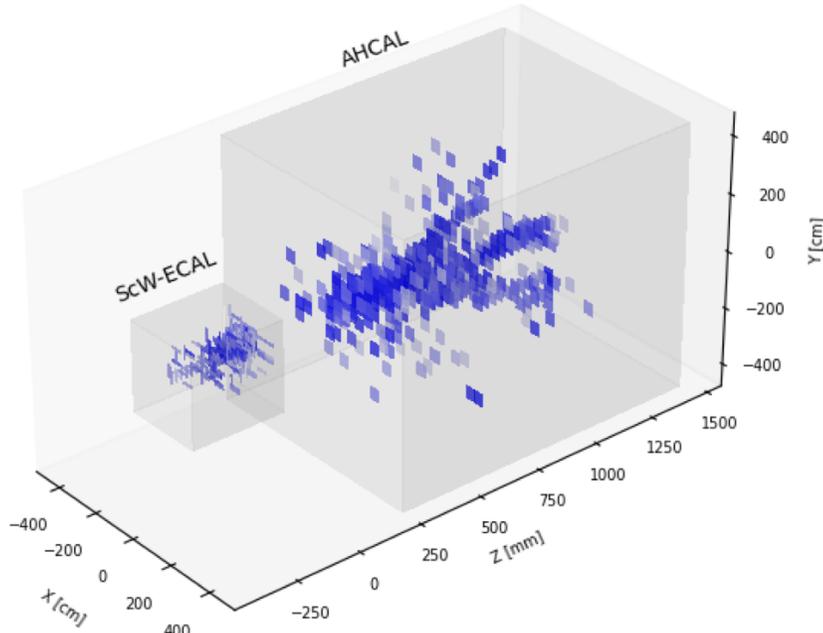
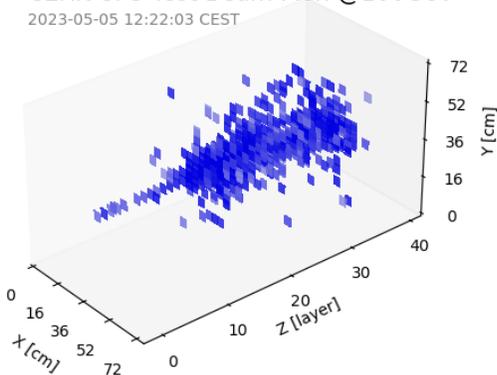
CEPC AHCAL

CERN SPS Test Beam Electron @100GeV
2023-05-07 06:48:33 CEST



CEPC AHCAL

CERN SPS Test Beam Pion @100GeV
2023-05-05 12:22:03 CEST

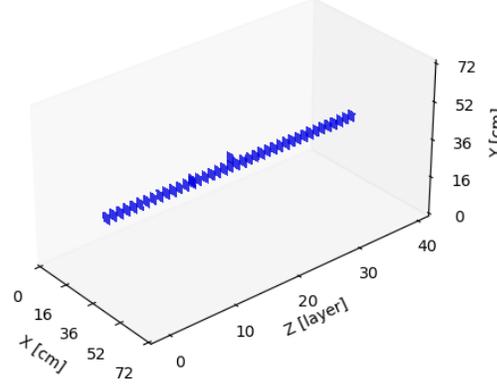


Event Display

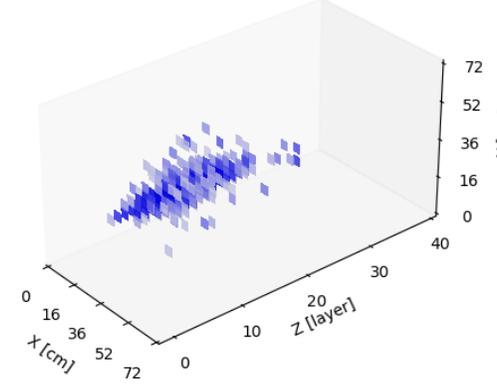
- AHCAL only.
- Shower topology of the same particle type is similar between MC and Data.
- Shower type:
 - Muon: Non-showering track.
 - Electron: Electromagnetic shower.
 - Pion: Hadronic shower.

Monte Carlo Samples

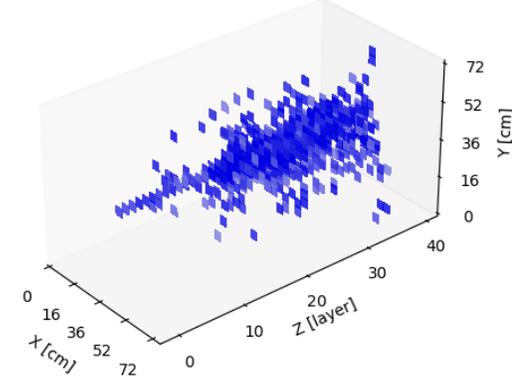
CEPC AHCAL
Muon Simulation @100GeV



CEPC AHCAL
Electron Simulation @100GeV

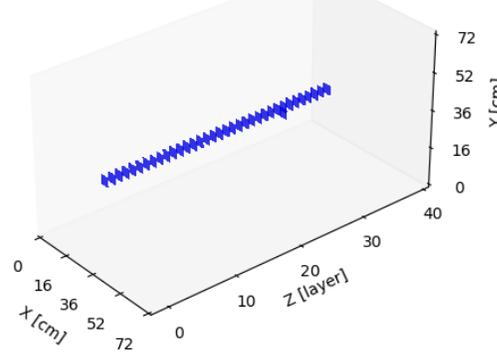


CEPC AHCAL
Pion Simulation @100GeV

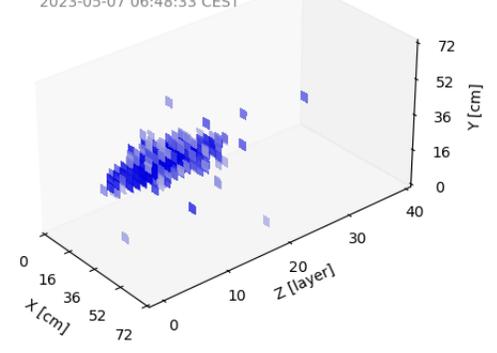


Test Beam Samples

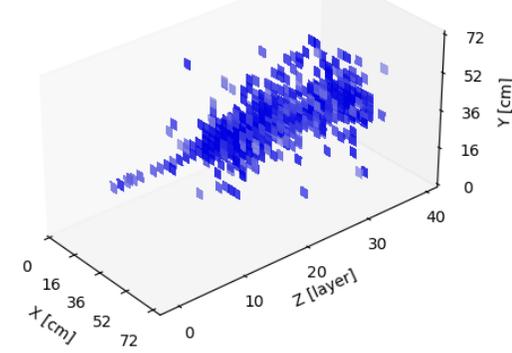
CEPC AHCAL
CERN SPS Test Beam Muon @100GeV
2023-04-27 21:01:02 CEST

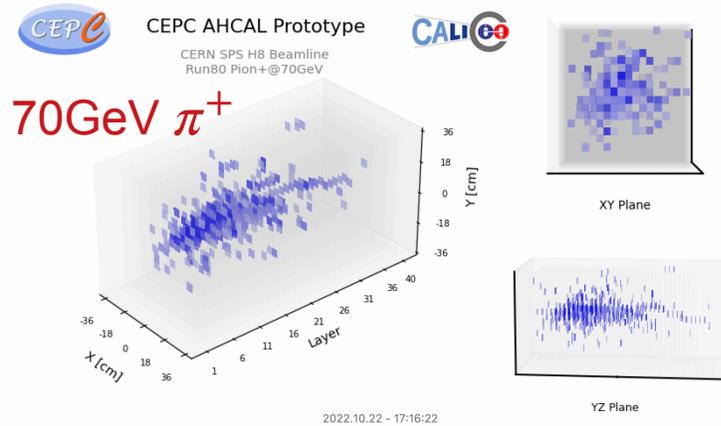
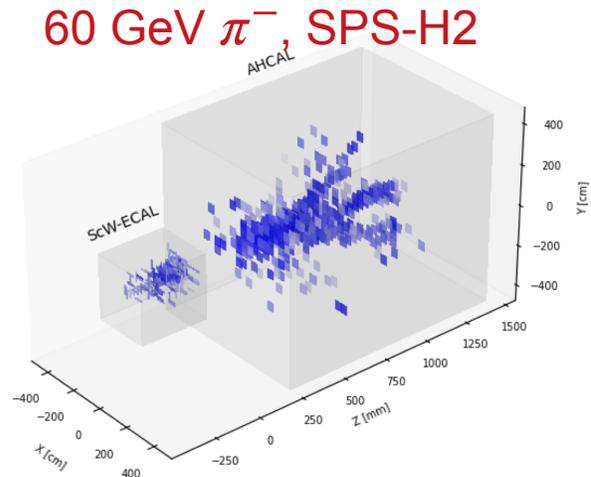
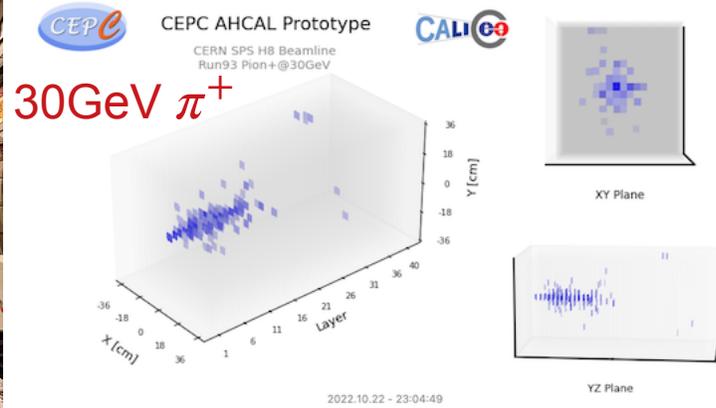


CEPC AHCAL
CERN SPS Test Beam Electron @100GeV
2023-05-07 06:48:33 CEST

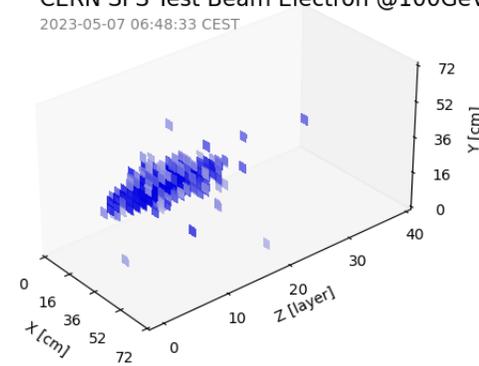


CEPC AHCAL
CERN SPS Test Beam Pion @100GeV
2023-05-05 12:22:03 CEST

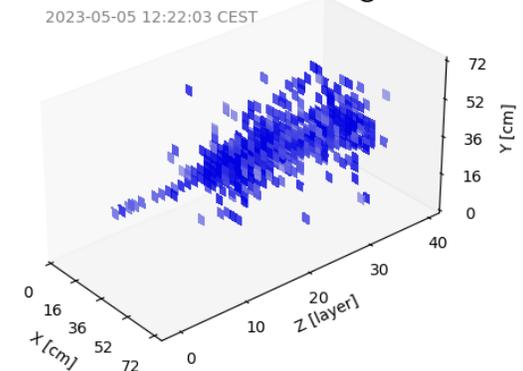




CEPC AHCAL
CERN SPS Test Beam Electron @100GeV
2023-05-07 06:48:33 CEST



CEPC AHCAL
CERN SPS Test Beam Pion @100GeV
2023-05-05 12:22:03 CEST



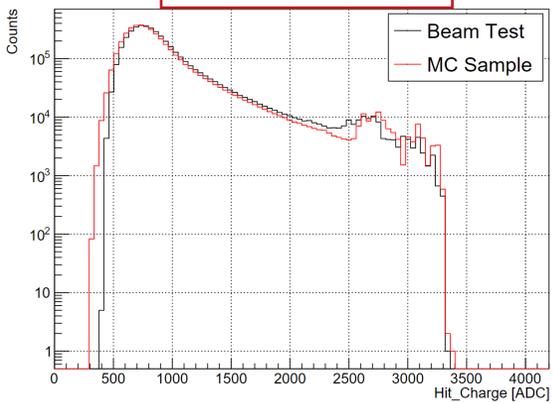
- First beam test at CERN SPS-H8: Oct-Nov, 2022.
- Beam: Muon, Electron, Pion (10-160GeV).
- Notice beam contamination issue.

- Beam test at CERN SPS-H2 & PS-T9: Apr-May, 2023.
- Beam: Muon, Electron, Pion, Proton (1-350GeV).
- Beam purity is better than 2022's.
- Available Cherenkov detectors (effective: $E_{Beam} < 30\text{GeV}$).

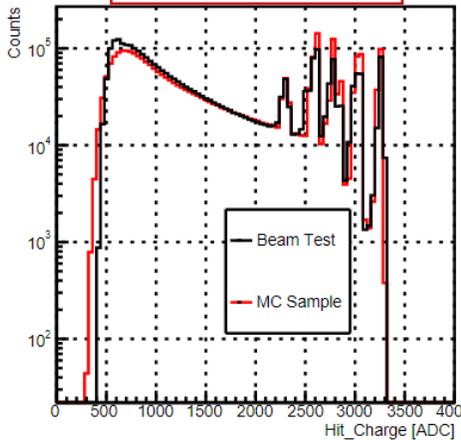


- **Validation on High/Low Gain - 0.5 MIP energy threshold**
 - Generally fit in High Gain; Need optimization in Low gain saturation correction.

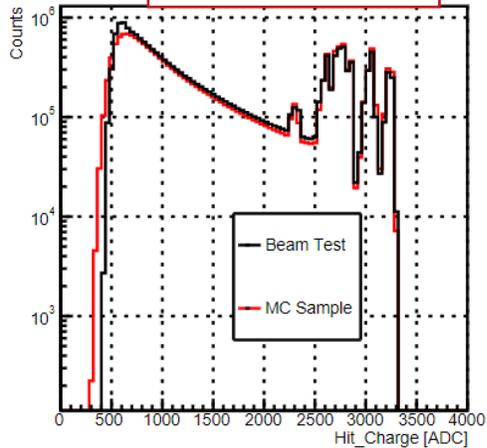
HG 100GeV μ^-



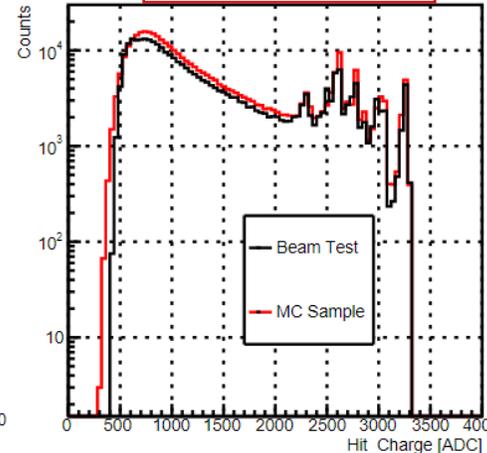
HG 10GeV e^-



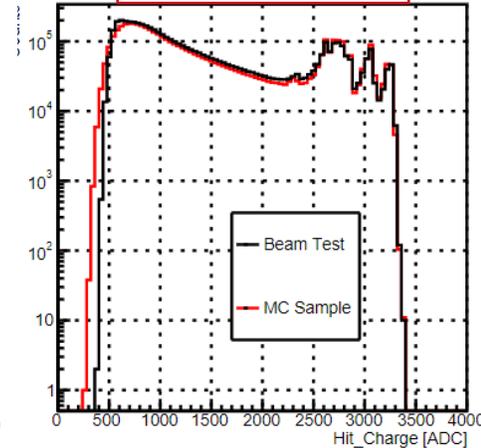
HG 80GeV e^-



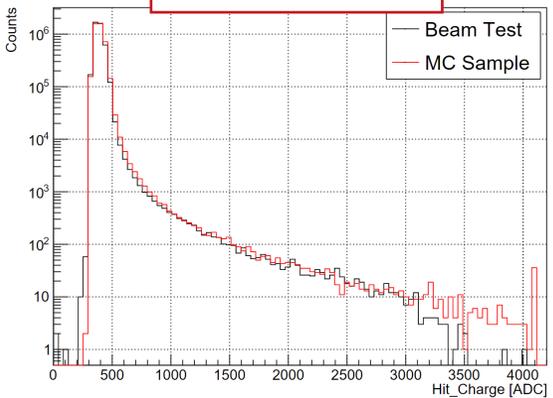
HG 10GeV π^-



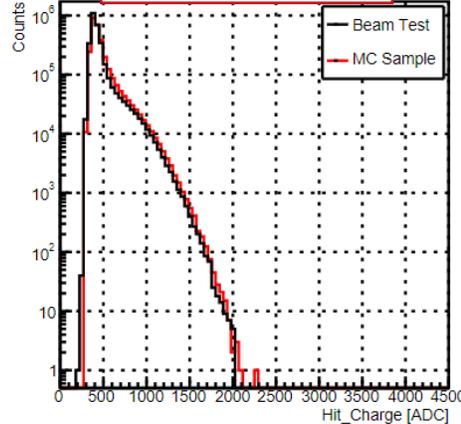
HG 80GeV π^-



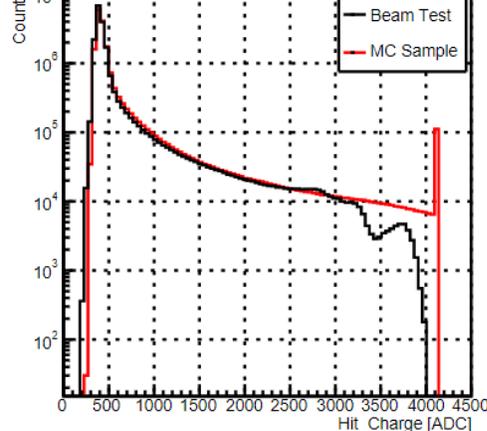
LG 100GeV μ^-



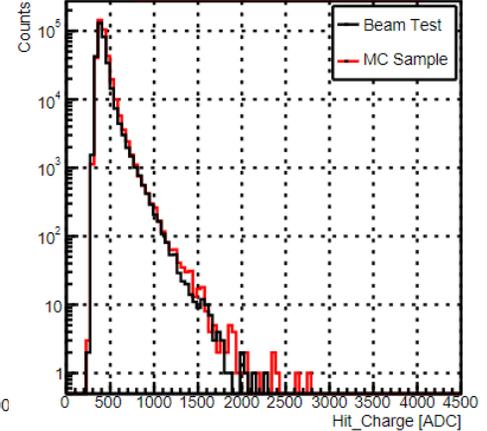
LG 10GeV e^-



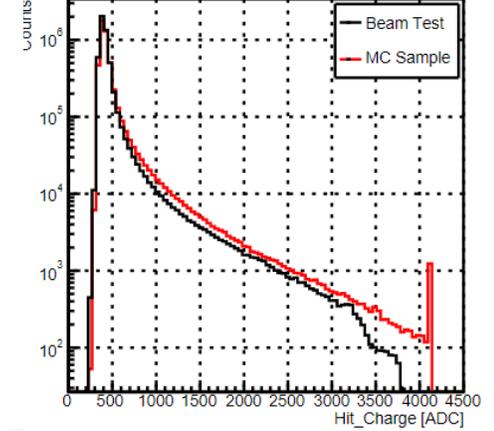
LG 80GeV e^-



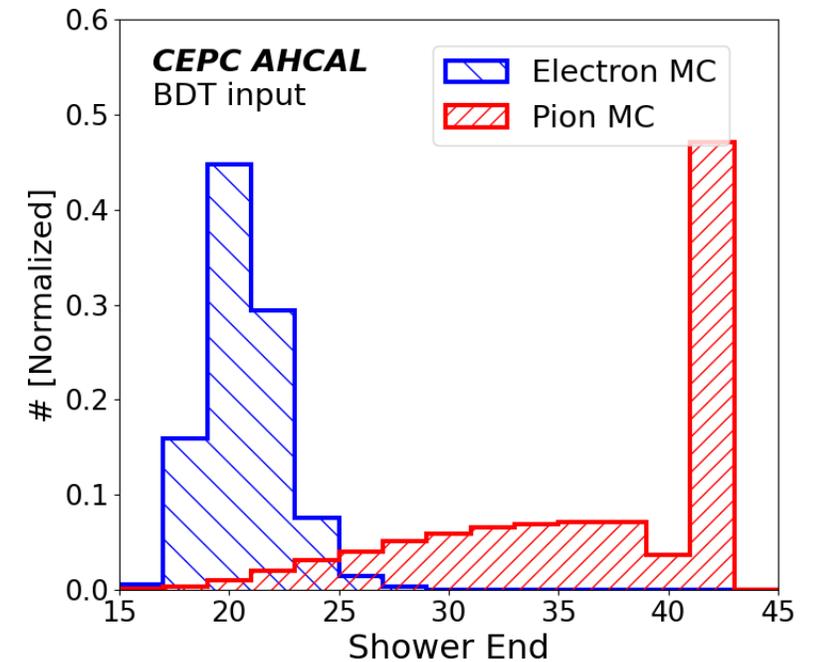
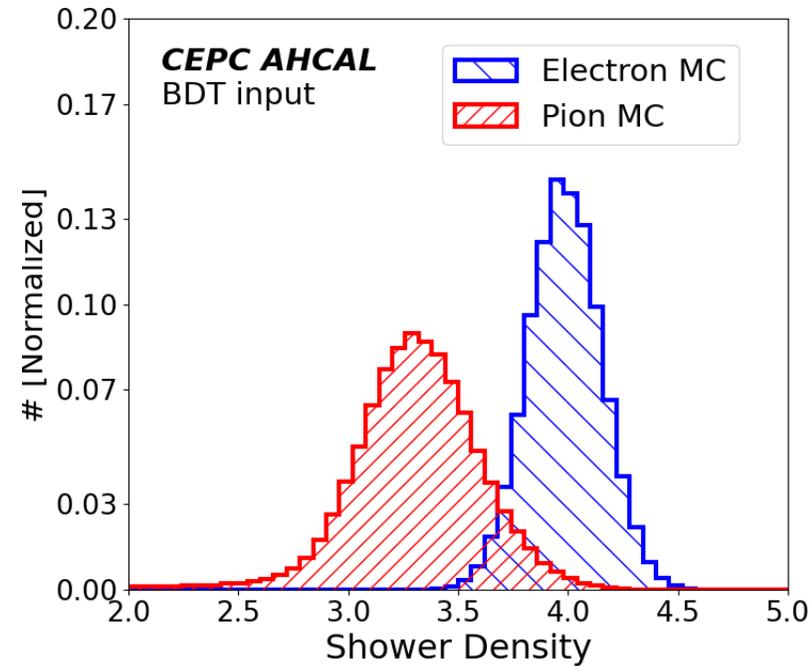
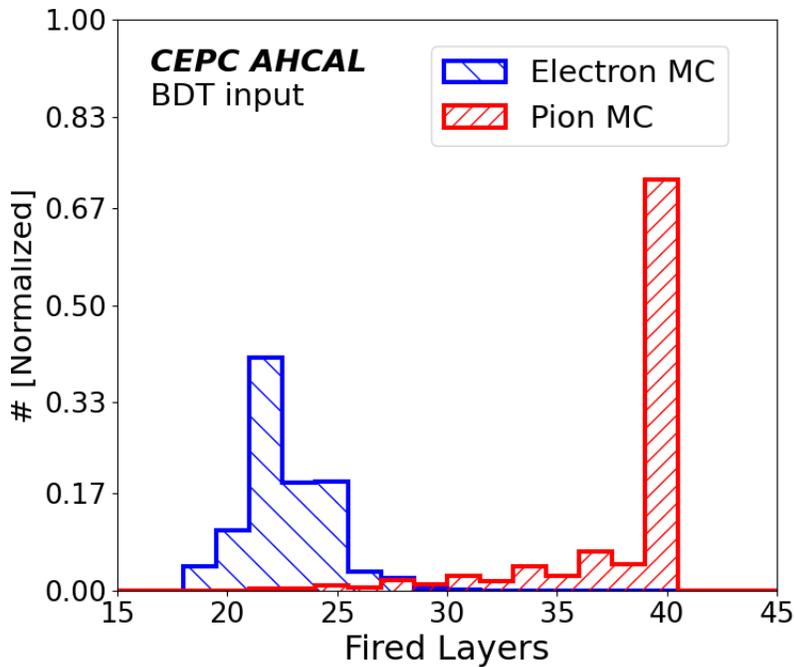
LG 10GeV π^-



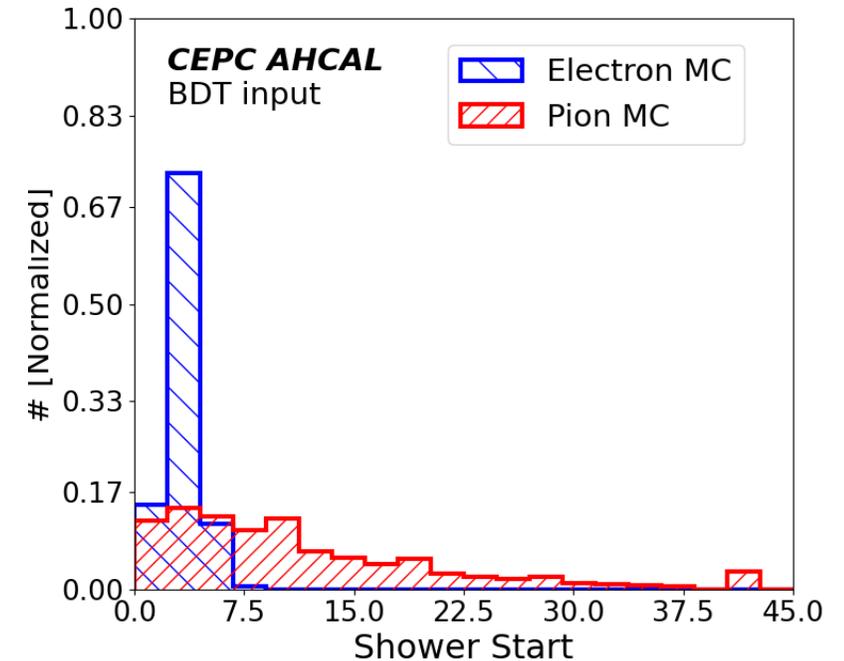
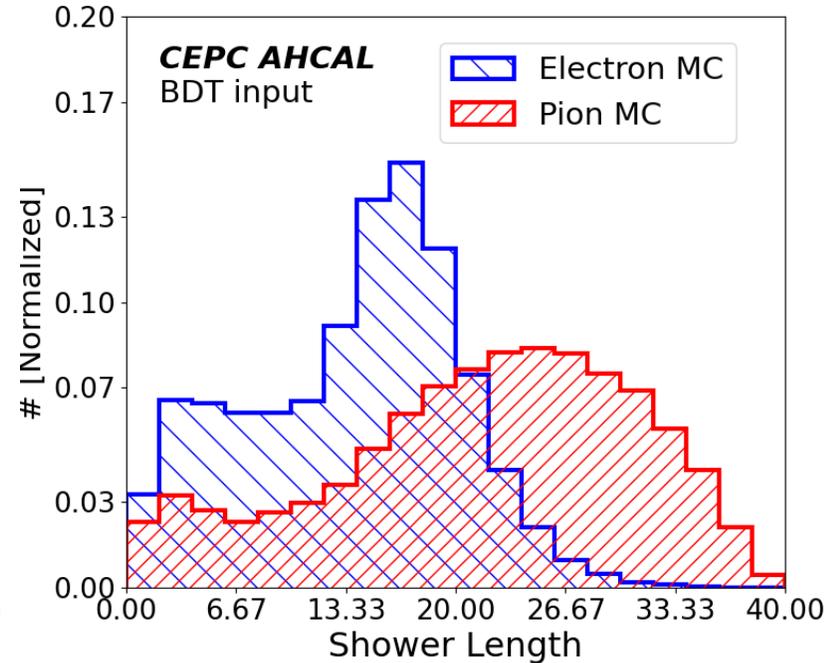
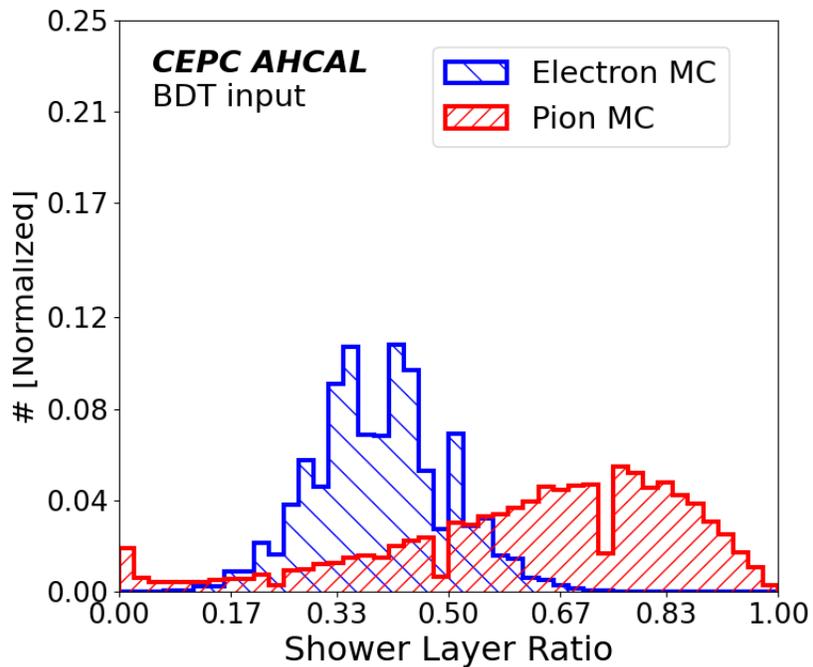
LG 80GeV π^-



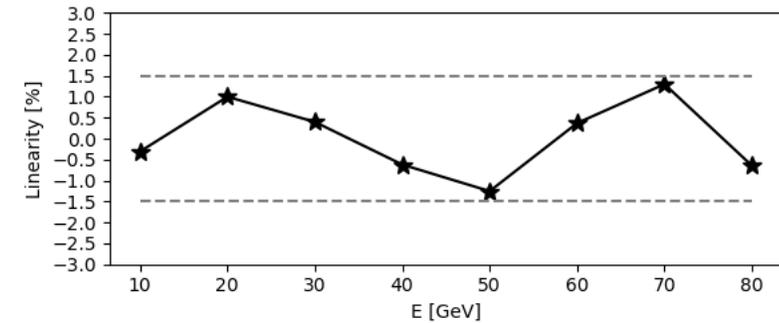
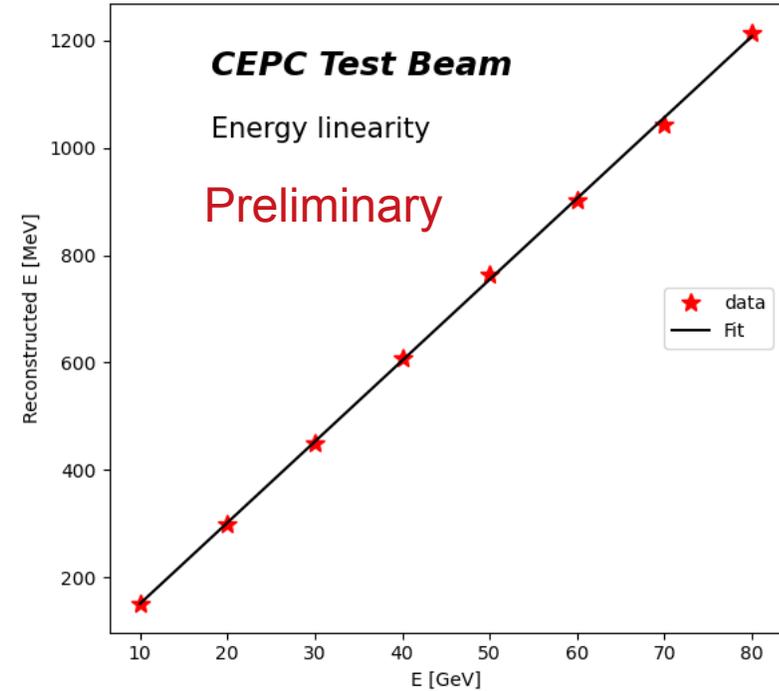
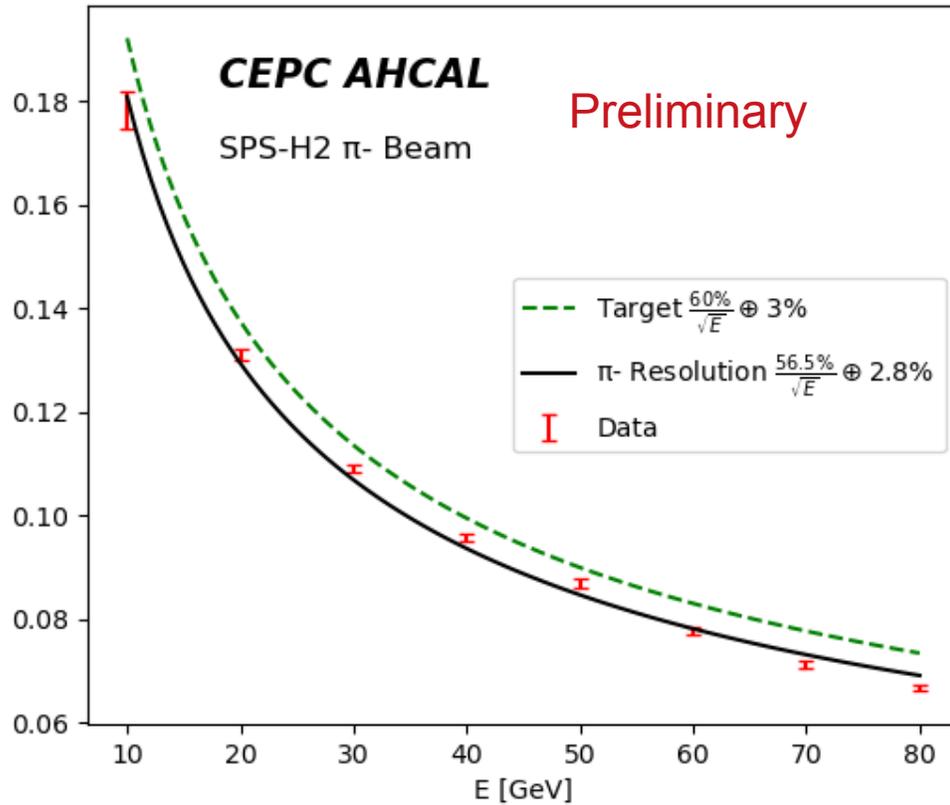
- **Fired layers:** The number of layer with hits.
- **Shower density:** The average number of neighboring hits around one hit, including the hit itself, in a 3×3 cell area in a given event is calculated.
- **Shower end:** After the shower starts, the first layer of two consecutive layers without 2 hits. If no shower is formed, it is set to 42.



- **Shower layer ratio:** Ratio of shower layers over fired layers.
- **Shower length:** The distance between the start of the shower and the layer where the maximum Root Mean Square (RMS) of hit transverse coordinates with respect to the z-axis occurs.
- **Shower start:** The first layer of the first three consecutive layers with at least 5 hits. For events without showers, the shower start layer is set to 42.



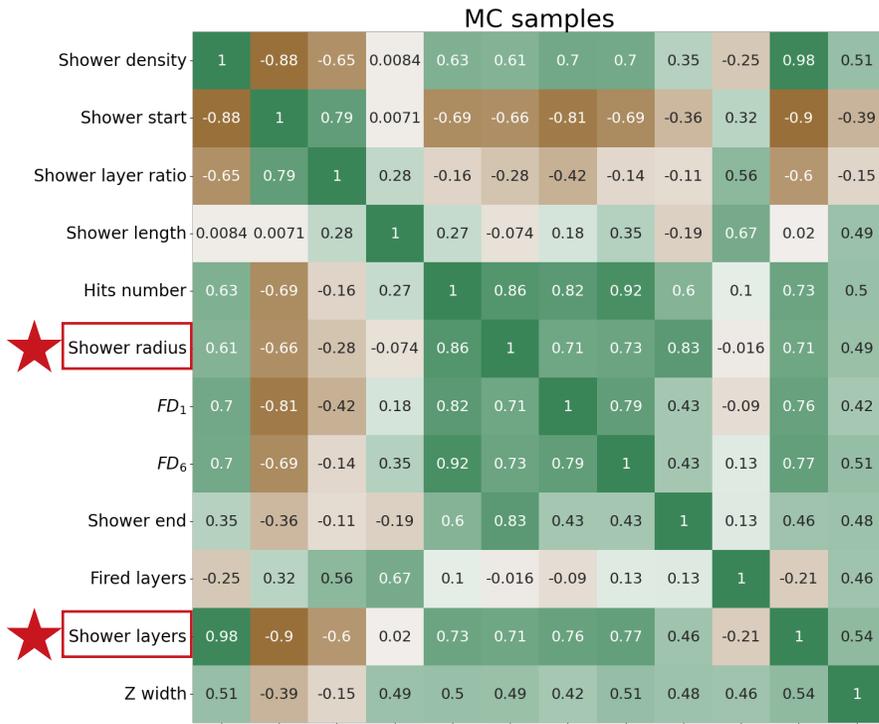
- The Energy resolution is $\frac{56.5\%}{\sqrt{E}} \oplus 2.8\%$.
- The Energy linearity is $\pm 1.5\%$.



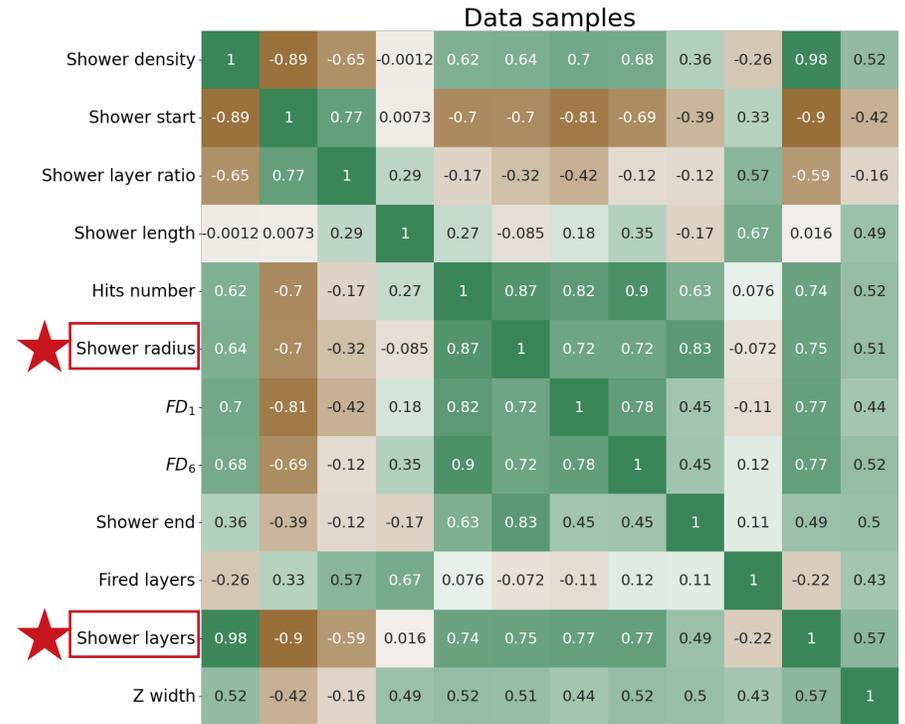
- **Apply Extreme Gradient Boosting (XGBoost)**

- 12 variables are reconstructed (Signal: π , Bkg.: e, μ)

- MC samples to build BDT_{MC-12} , Data samples to build $BDT_{Data-12}$



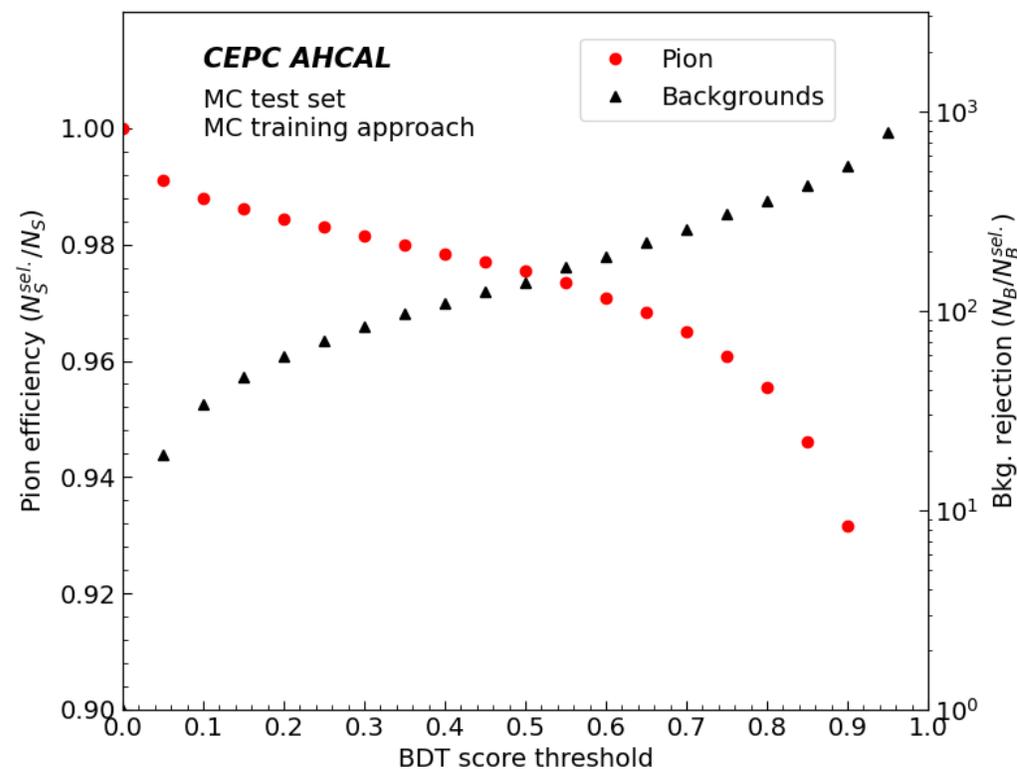
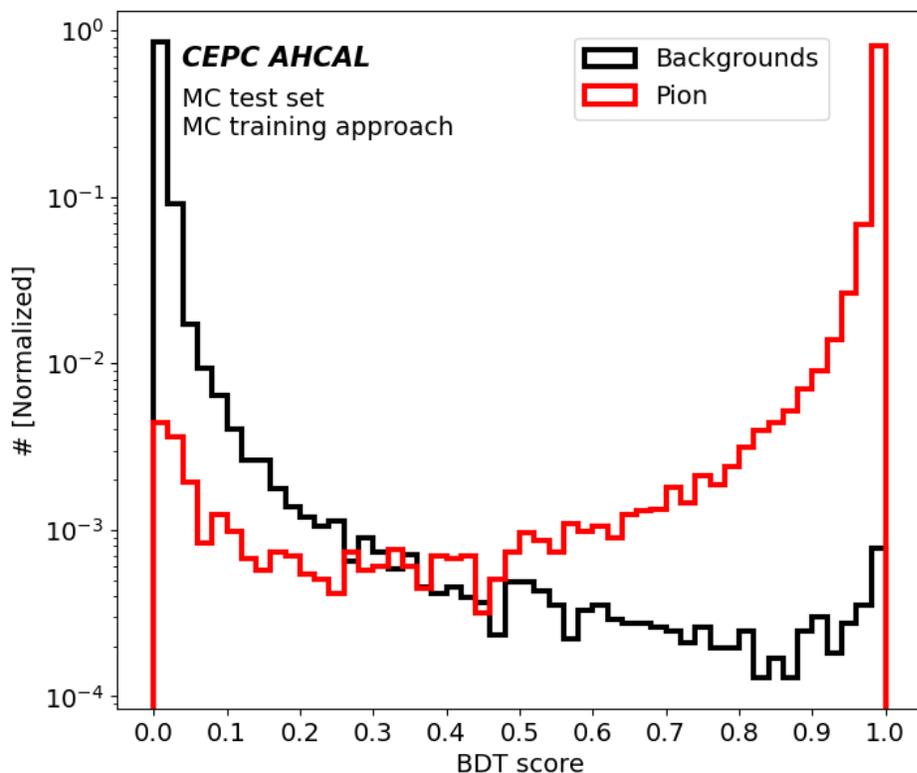
Correlation matrix



Correlation matrix

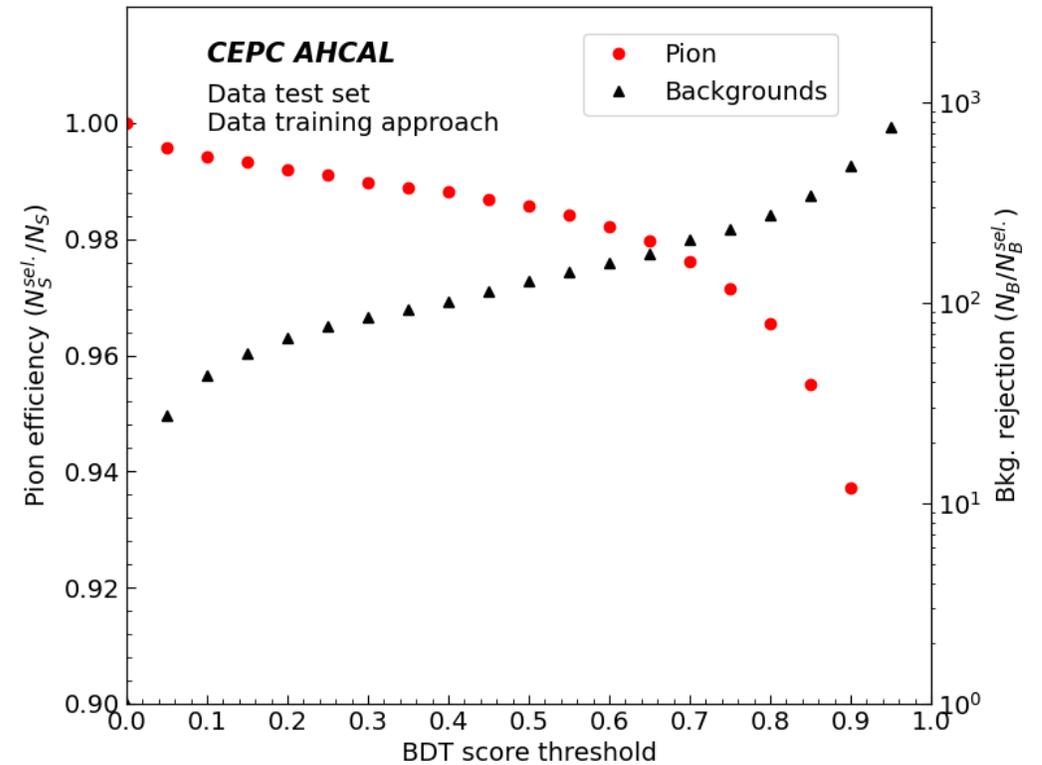
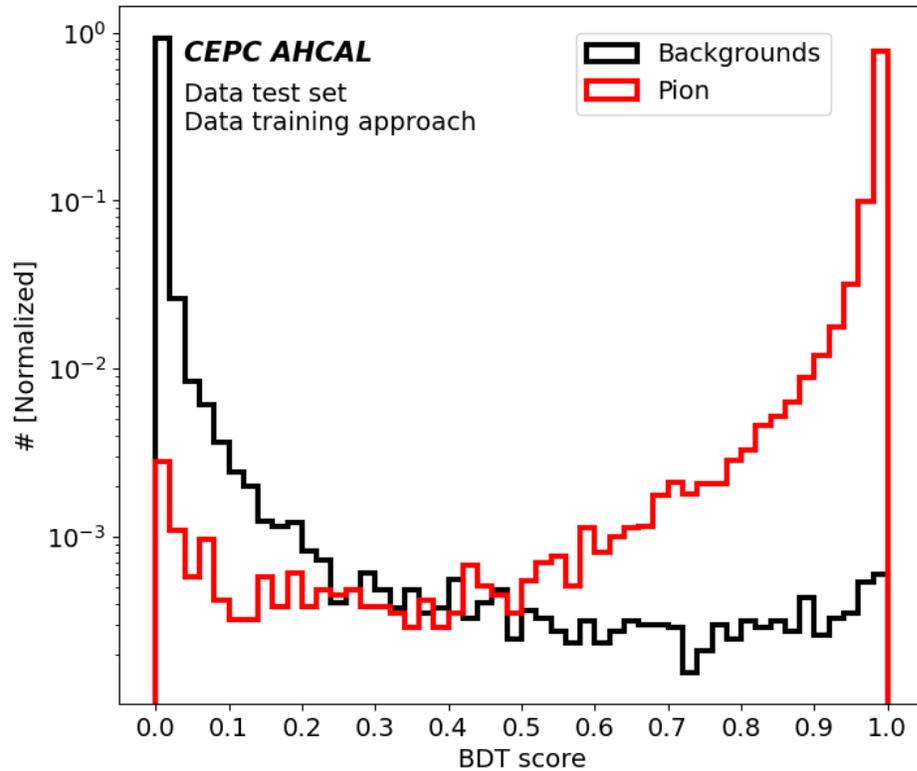
MC training approach

Pion efficiency	90%	95%	99%
Bkg. rejection (1/Bkg. efficiency)	1701.2	617.4	29.6

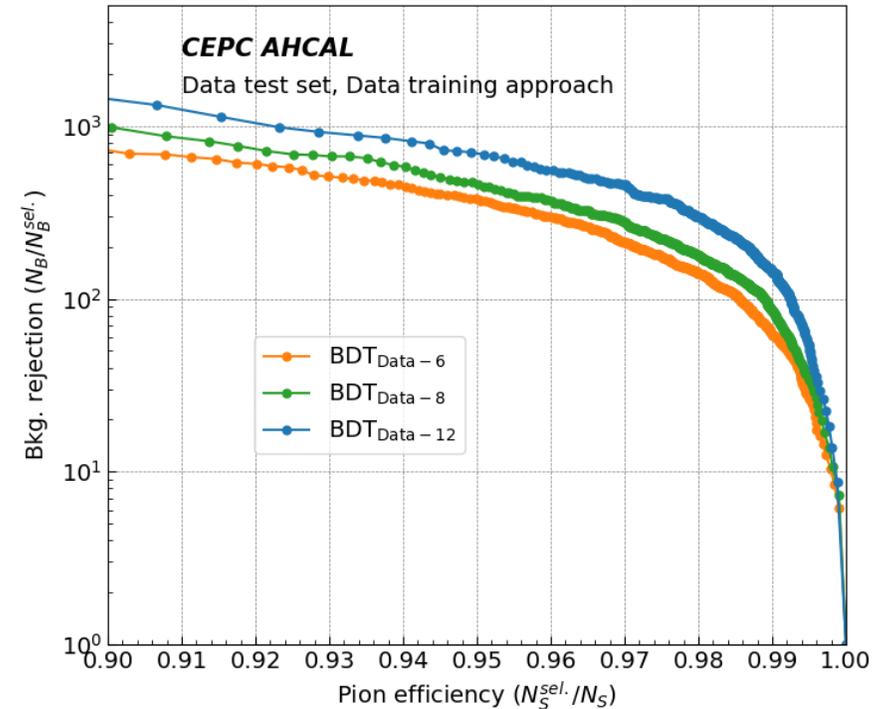
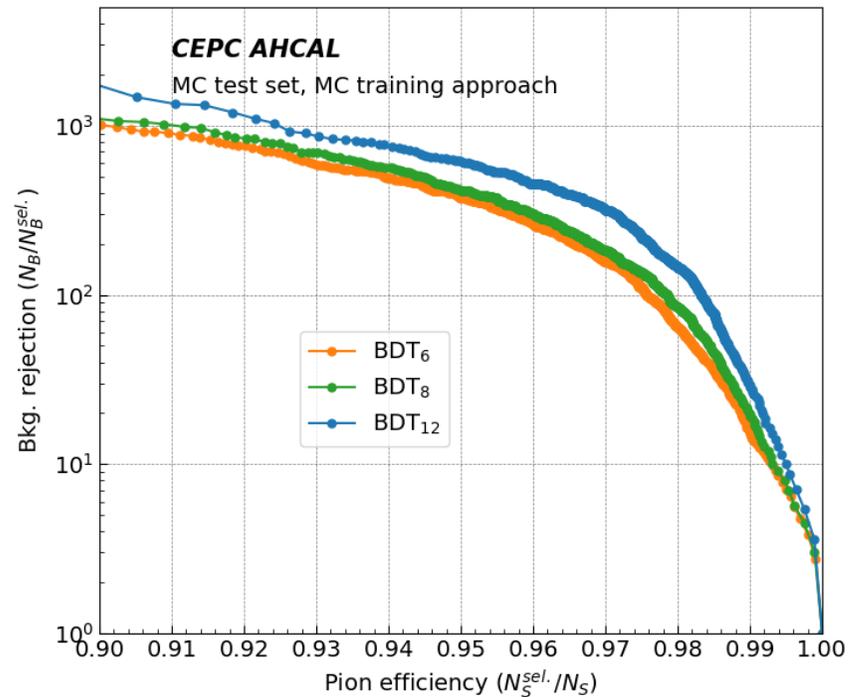


Data training approach

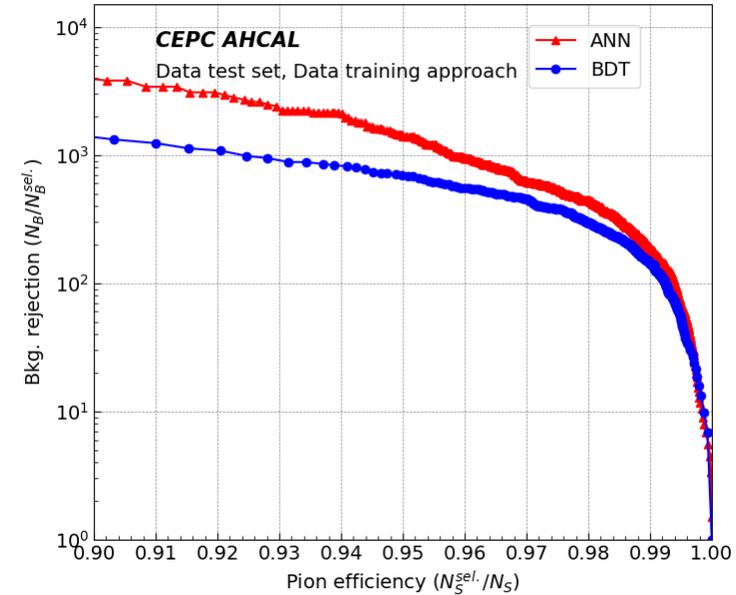
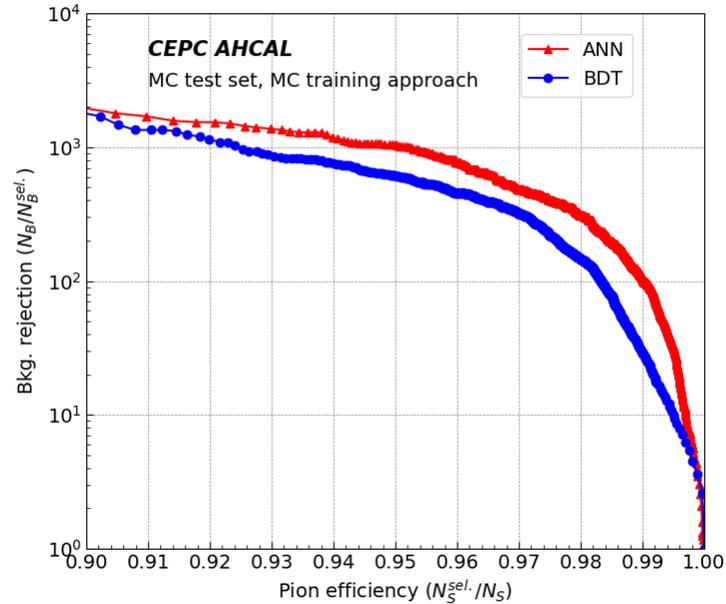
Pion efficiency	90%	95%	99%
Bkg. rejection (1/Bkg. efficiency)	1448.5	691.6	143.0



- **We observe dependence of BDT performance on input variables**
 - Remove Shower End, Shower Layers, Fired Layers, and Z Width to build BDT with 8 inputs.
 - Further remove FD_1 and FD_6 to build BDT with 6 inputs.
- **Feature engineering matters in BDT**
 - BDT optimization is on going.



- We observe obvious improvement in terms of Background rejection when tested on two sets of samples.



Pion efficiency	90%		95%		99%	
	MC	Data	MC	Data	MC	Data
BDT bkg. rejection	1701.2	1448.5	617.4	691.6	29.6	143.0
ANN bkg. rejection	2015.7	3811.2	1040.3	1408.5	103.9	187.8
Improvement	18.49%	163.12%	68.51%	103.65%	251.14%	31.37%

Quick view on shower topology variables

- Data samples are pre-selected by Cherenkov detectors and rough F.D. cut.
- MC and Data are close in Shower topology term .
- **Obvious discrepancy in energy term.**

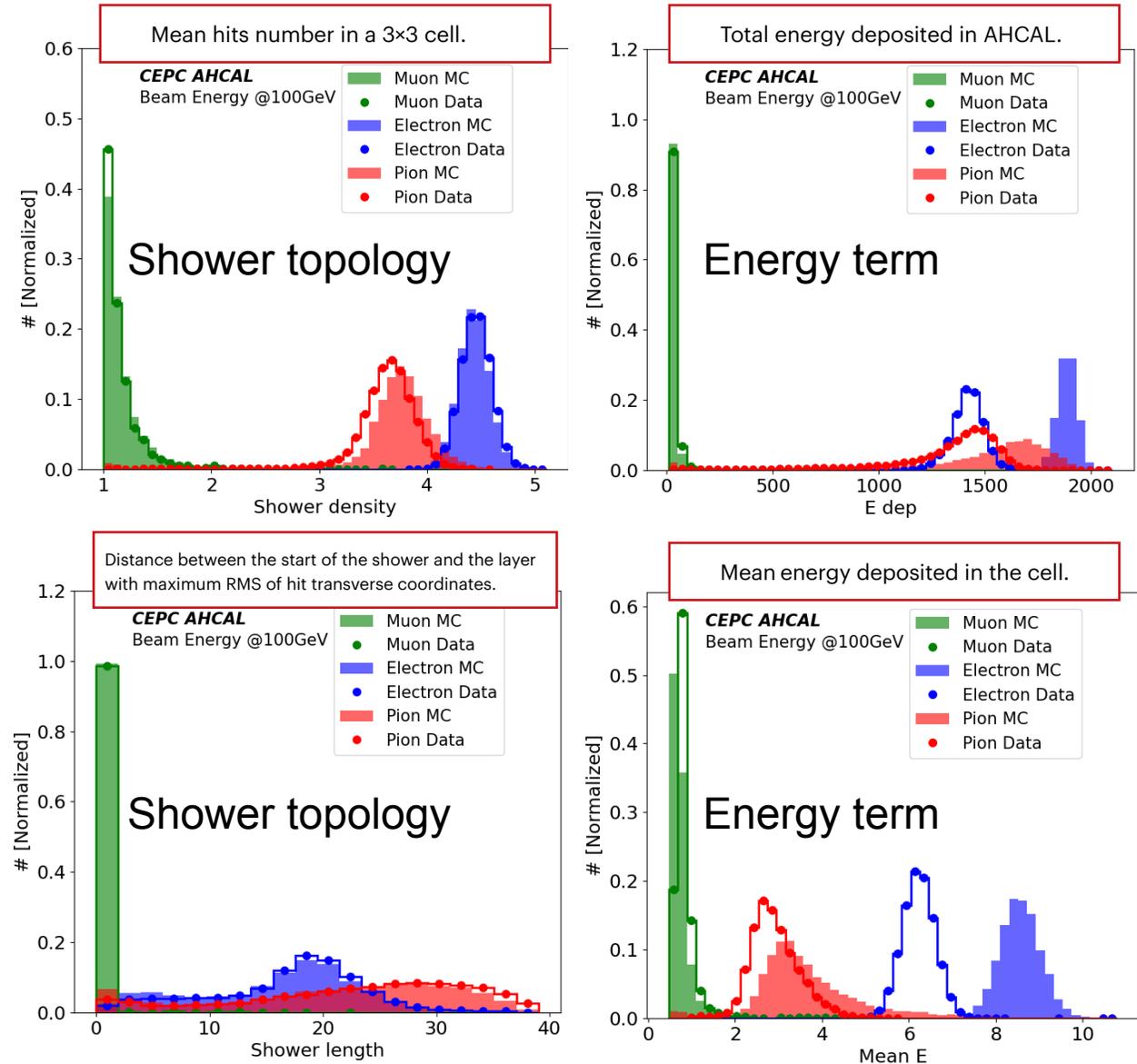
Currently

- ✓ Do Research on PID approach.
 - Separate Data training and MC training.
- ❖ Trained on MC and then applied on Data.
 - Discrepancy between Data and MC.
 - Bridging gap is on going.

Muon data: SPS_2023/100GeV_mu_Run25.

Electron data: SPS_2023/100GeV_e_run267.

Pion data: SPS/100GeV_pi_run230.



- **Apply Extreme Gradient Boosting (XGBoost)**

- 12 variables are reconstructed (Signal: π , Bkg.: e, μ)

- MC samples to build BDT_{MC-12} , Data samples to build $BDT_{Data-12}$

MC samples

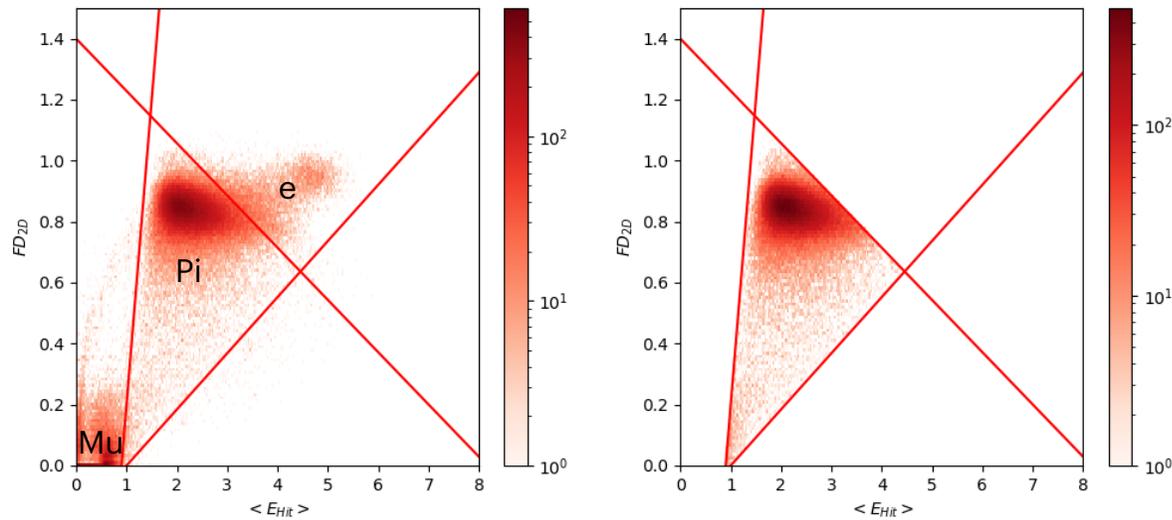
Rank: Variable	Variable weight
1: Shower radius	0.377
2: Shower layers	0.232
3: Hits number	0.088
4: Fired layers	0.083
5: Shower start	0.080
6: Shower density	0.049
7: Z width	0.034
8: FD_6	0.017
9: FD_1	0.015
10: Shower layer ratio	0.014
11: Shower end	0.006
12: Shower length	0.006

Data samples

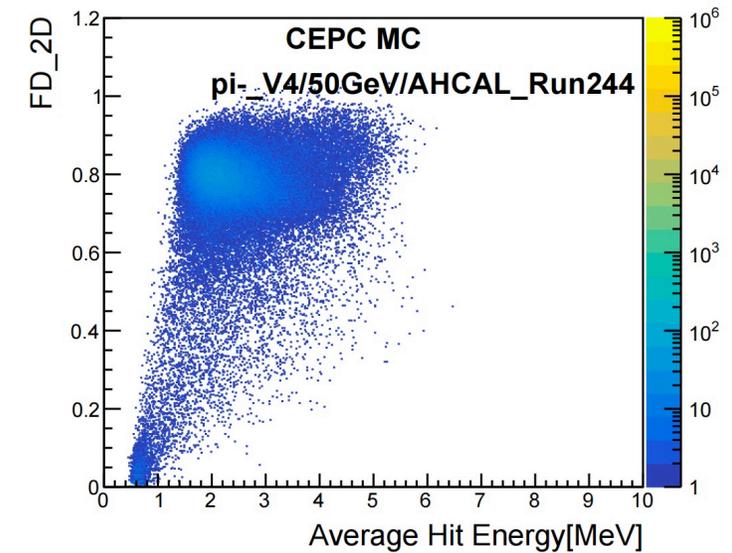
Rank: Variable	Variable weight
1: Shower radius	0.379
2: Shower layers	0.228
3: Hits number	0.133
4: Shower density	0.058
5: Fired layers	0.058
6: Z width	0.042
7: Shower start	0.039
8: FD_6	0.019
9: FD_1	0.016
10: Shower layer ratio	0.010
11: Shower length	0.010
12: Shower end	0.008

Data pre-selection

- Collect pion samples in 20pion run files.
- Cut approach is guided by MC.



FD cut

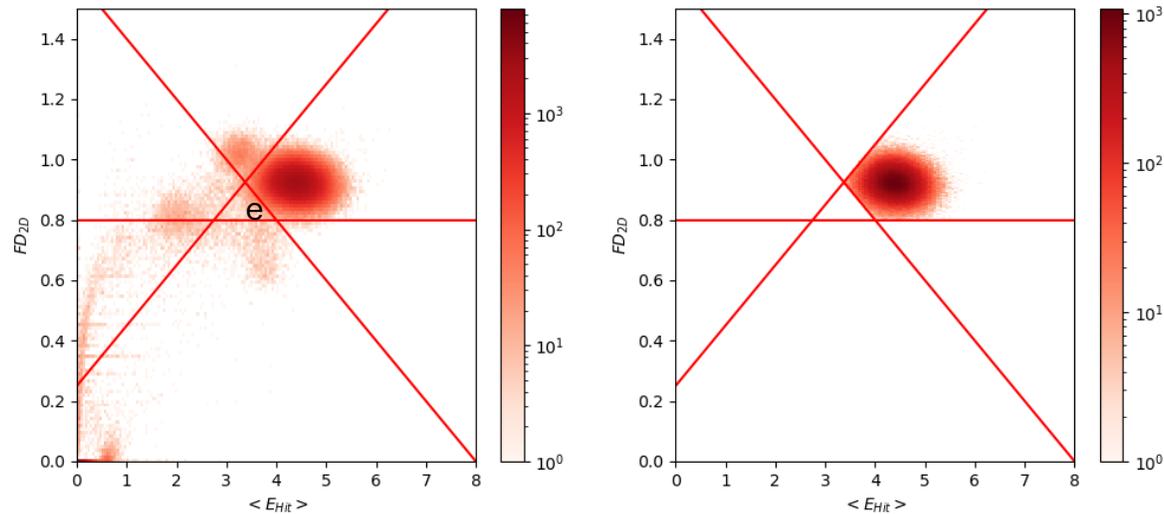


MC

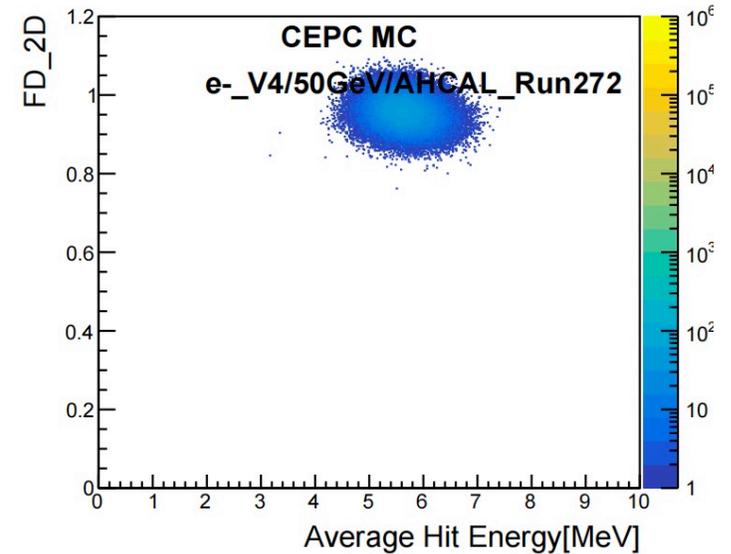
For data collected in 2023, SPS and PS

Data pre-selection

- Collect e samples in e run files.
- Cut approach is guided by MC.



FD cut



MC

For data collected in 2023, SPS and PS