

# Machine-learning-based Particle Identification in using CEPC AHCAL Prototype

Siyuan SONG, Jiyuan CHEN

Advisor: Prof. Haijun YANG

2023-09-06

饮水思源 · 爱国荣校<sup>1</sup>

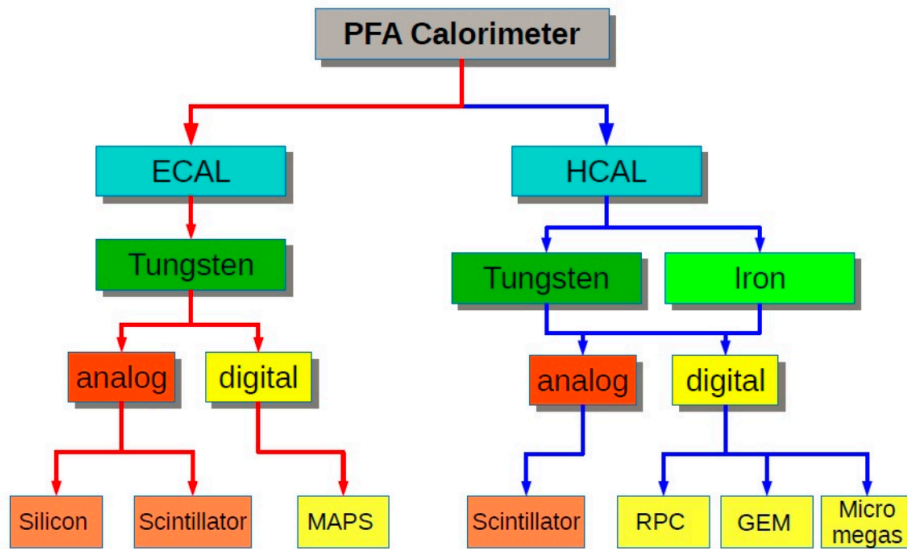


1. Introduction.
2. Monte Carlo Samples & Test Beam Samples.
3. PID based on BDT.
4. PID based on ANN.
5. Purity of 2022 AHCAL Test Beam.
6. Summary.



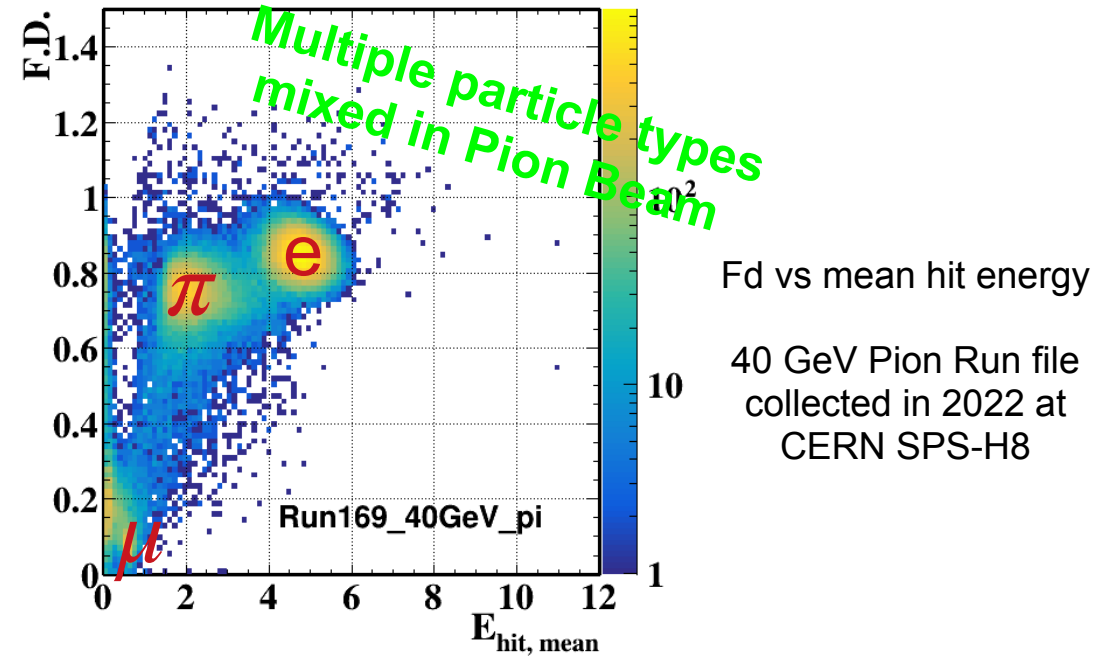
- **Exploration for multiple PID methods in high-granularity calorimeters**
- **CEPC AHCAL prototype test beam data require PID**

- Several PFA oriented high-granularity calorimeters have been developed.

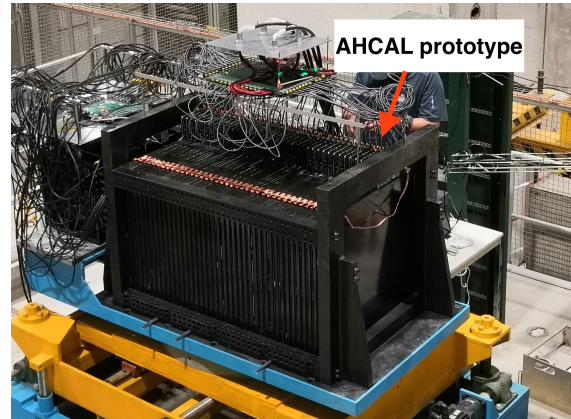
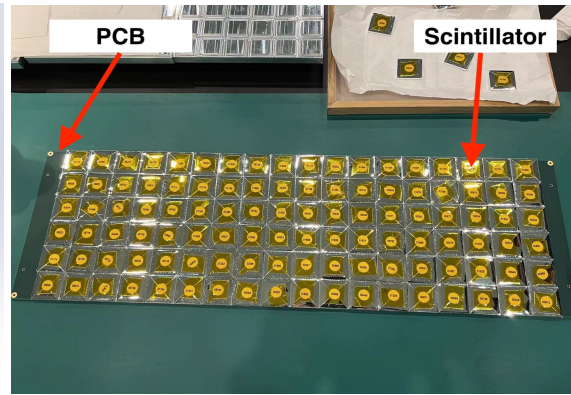
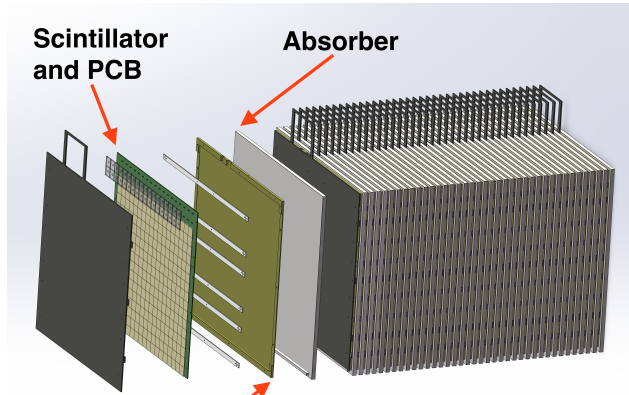


Overview of world-wide development of high-granularity calorimeters

- Contamination in test beam data collected in 2022 at CERN SPS-H8.
- Main task is purifying Pion beams.



- **Utilize High-granularity CEPC AHCAL prototype in this PID research**



## 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.
- Electronics readout channels
  - 12960 (18 × 18 × 40).



- **Monte Carlo Samples:** Employ Geant4 11.1.1 Toolkit with the QGSP<sub>BERT</sub> physics list.

Energy point	5 GeV		10 GeV		30 GeV		50 GeV		60 GeV		80 GeV		100 GeV		120 GeV	
	#	Source	#	Source	#	Source	#	Source	#	Source	#	Source	#	Source	#	Source
Muon	10k	MC	10k	MC	10k	MC	10k	MC	10k	MC	10k	MC	10k	MC	10k	MC
Electron	10k	MC	10k	MC	10k	MC	10k	MC	10k	MC	10k	MC	10k	MC	10k	MC
Pion	10k	MC	10k	MC	10k	MC	10k	MC	10k	MC	10k	MC	10k	MC	10k	MC

- **Test Beam Samples:** Pre-processed purer 2023 CERN SPS-H2 & PS-T9 test beam data.

Energy point	5 GeV		10 GeV		30 GeV		50 GeV		60 GeV		80 GeV		100 GeV		120 GeV	
	#	Source	#	Source	#	Source	#	Source	#	Source	#	Source	#	Source	#	Source
Muon	-	-	40k	Data	-	-	-	-	-	-	-	-	40k	Data	-	-
Electron	10k	Data	10k	Data	10k	Data	10k	Data	10k	Data	10k	Data	10k	Data	10k	Data
Pion	10k	Data	10k	Data	10k	Data	10k	Data	10k	Data	10k	Data	10k	Data	10k	Data

Each sample set is split to a Train set and a Test set in a ration of 3:2.

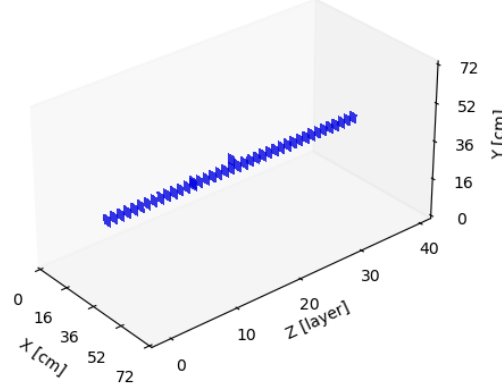
Each sample set would be utilized to build classifiers.

## Event Display

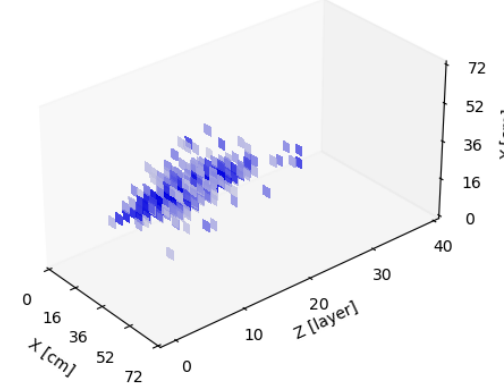
- PID depends the shower topology.
- Shower type:
  - Muon: Non-showering track.
  - Electron: Electromagnetic shower.
  - Pion: Hadronic shower.
- Shower topology of the same particle type is similar between MC and Data.

## Monte Carlo Samples

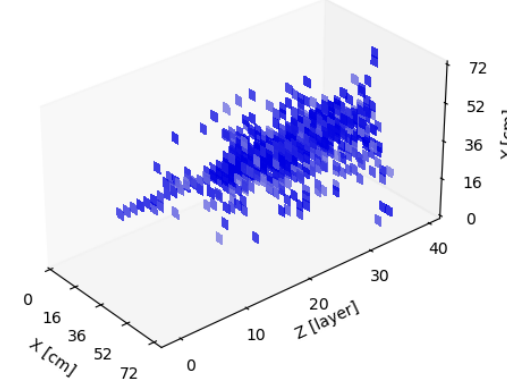
**CEPC AHCAL**  
Muon Simulation @100GeV



**CEPC AHCAL**  
Electron Simulation @100GeV

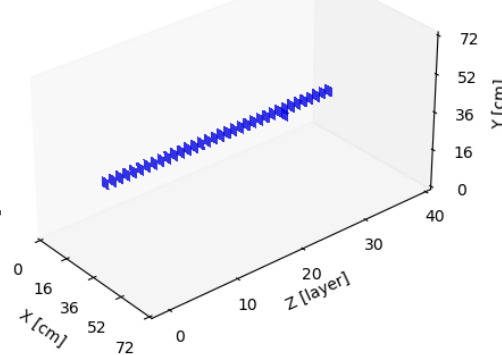


**CEPC AHCAL**  
Pion Simulation @100GeV

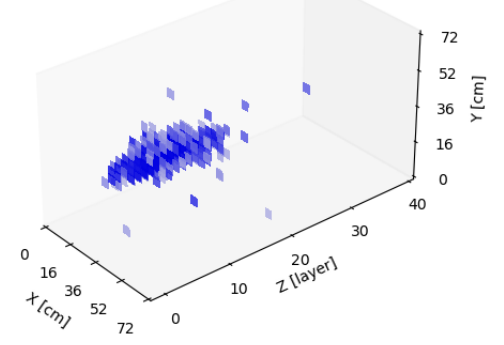


## Test Beam Samples

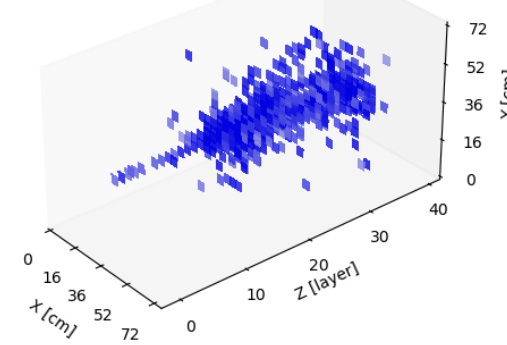
**CEPC AHCAL**  
CERN SPS Test Beam Muon @100GeV  
2023-04-27 21:01:02 CEST



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

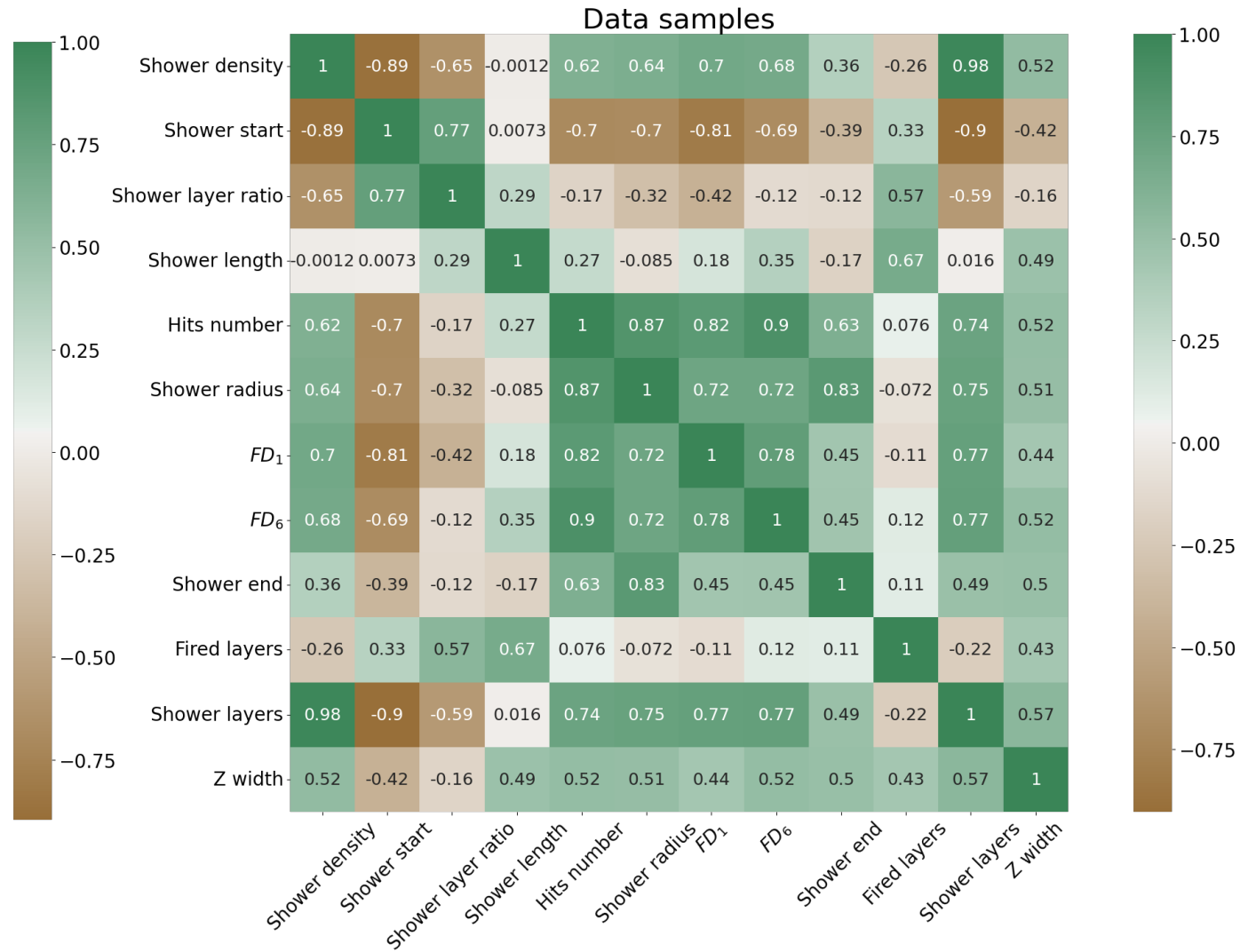
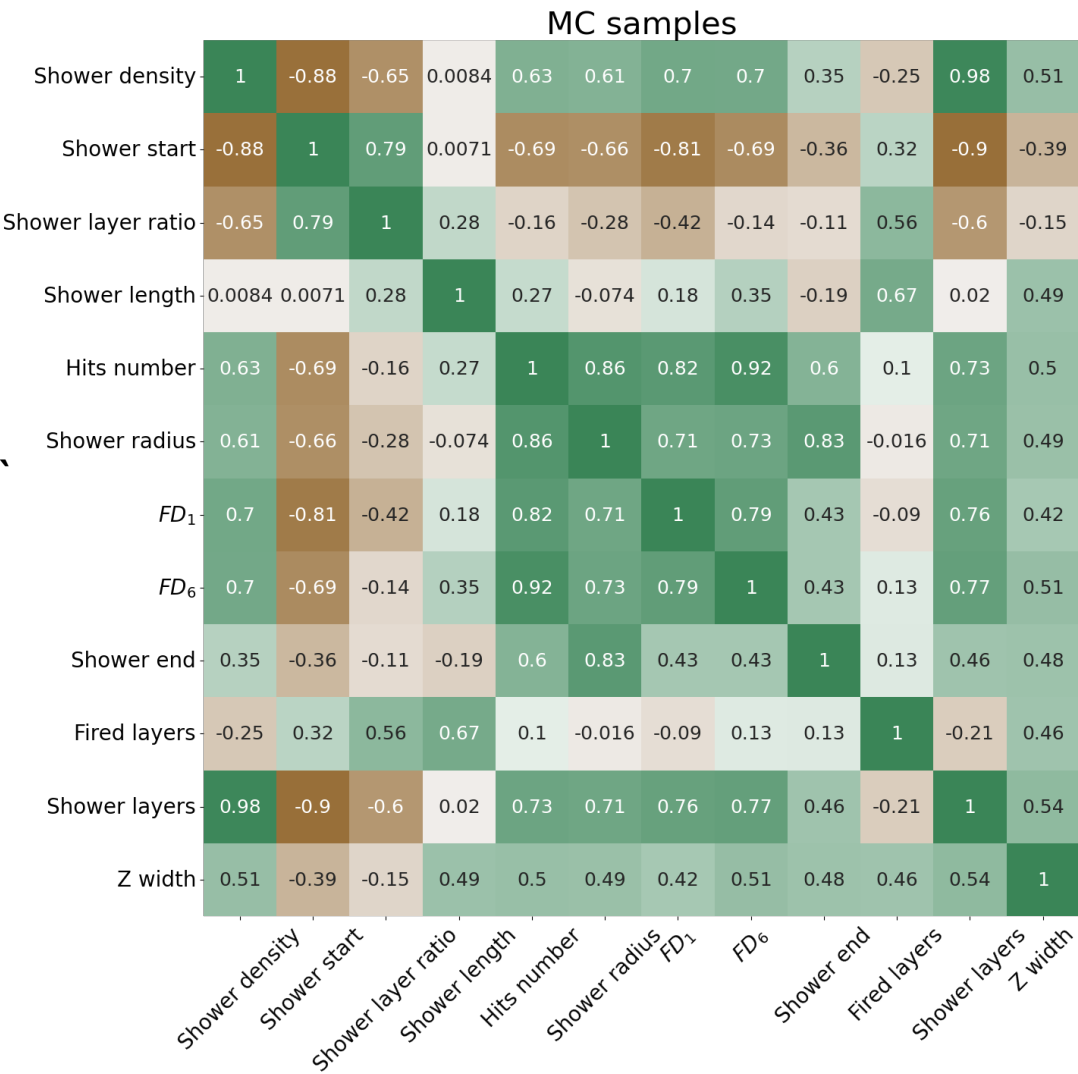


**CEPC AHCAL**  
CERN SPS Test Beam Pion @100GeV  
2023-05-05 12:22:03 CEST





## • Correlation Matrix of 12 input variables.





- **Apply XGBoost**
- **Variable Ranking in Pion identification**

- Shower radius,
- Shower layers ,
- Hits number are important.
- (Signal:  $\pi$ , Background:  $e$  &  $\mu$  )

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

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 <sub>6</sub>	0.017
9: FD <sub>1</sub>	0.015
10: Shower layer ratio	0.014
11: Shower end	0.006
12: Shower length	0.006

MC samples

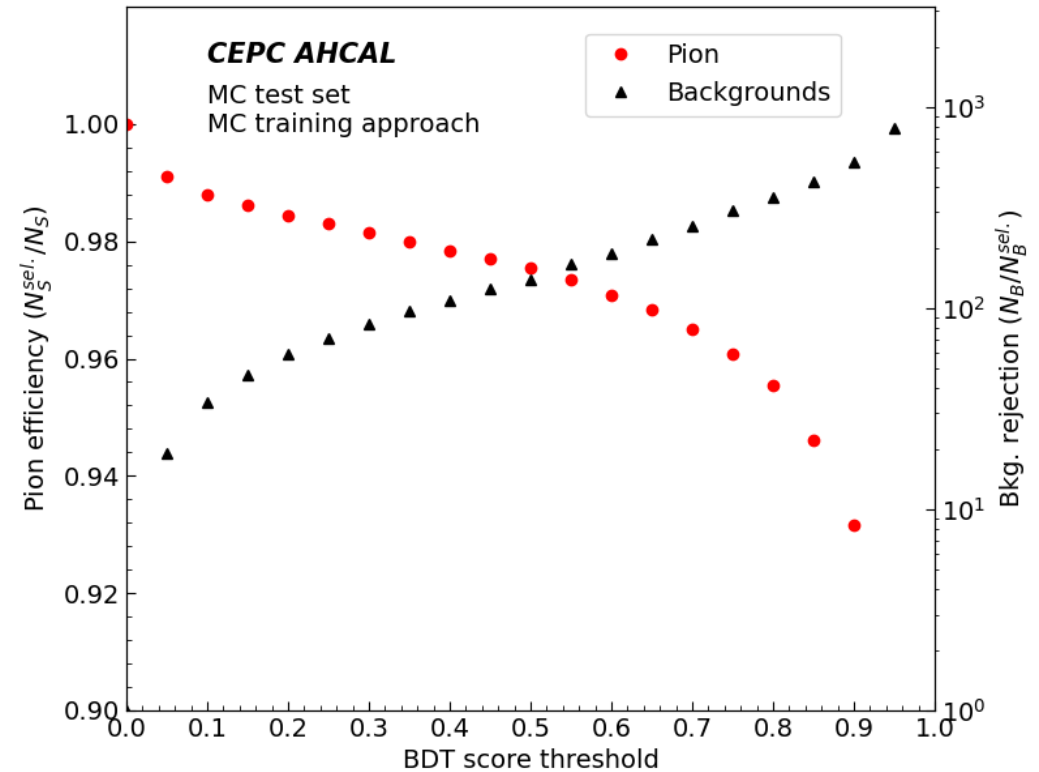
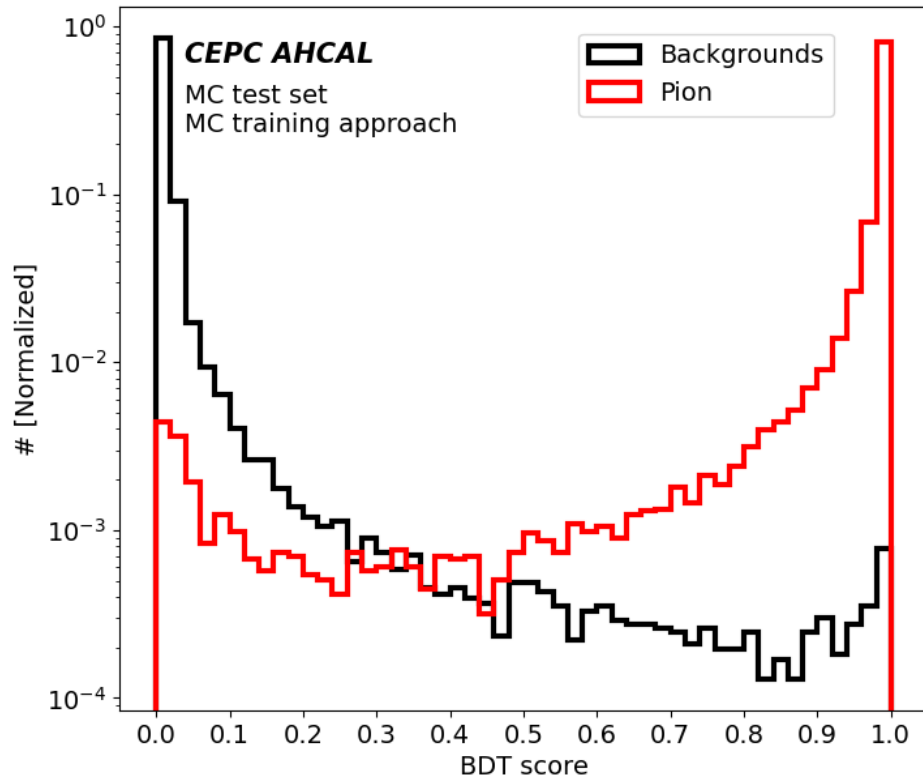
Rank: Variable	Variable weight
1: Shower radius	0.379
2: Shower layers	0.228
3: Hits number	0.133
4: Shower density	0.058
5: Fired layers	0.058
6: Z width	0.042
7: Shower start	0.039
8: FD <sub>6</sub>	0.019
9: FD <sub>1</sub>	0.016
10: Shower layer ratio	0.010
11: Shower length	0.010
12: Shower end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

Data samples



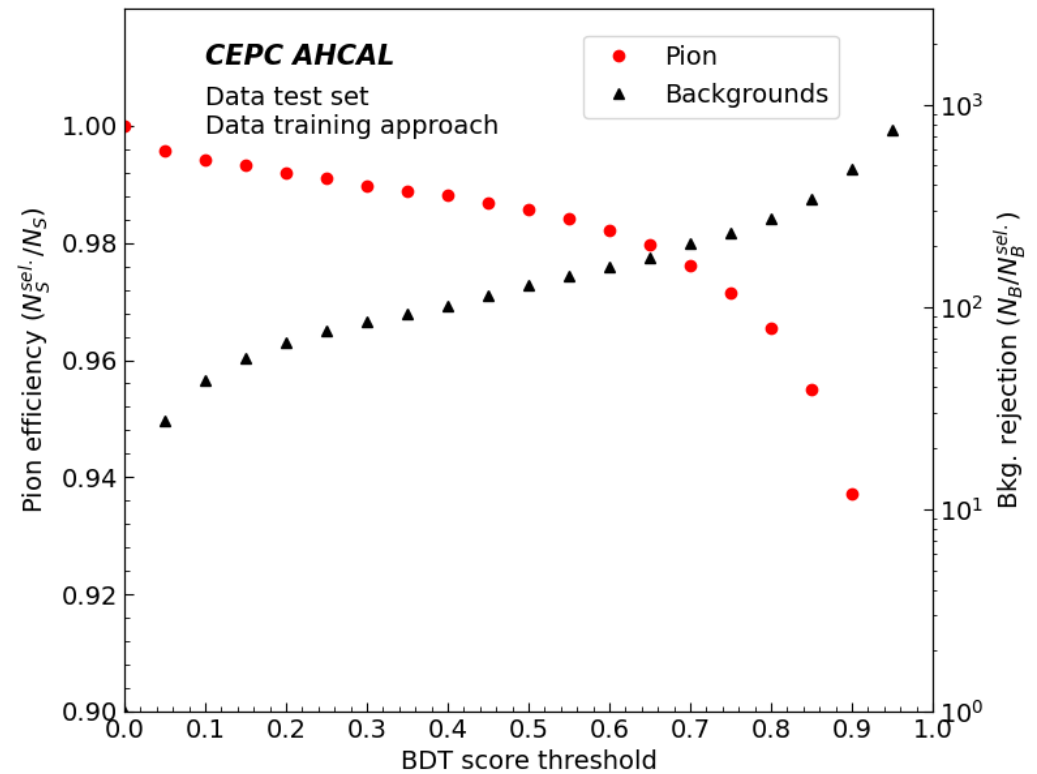
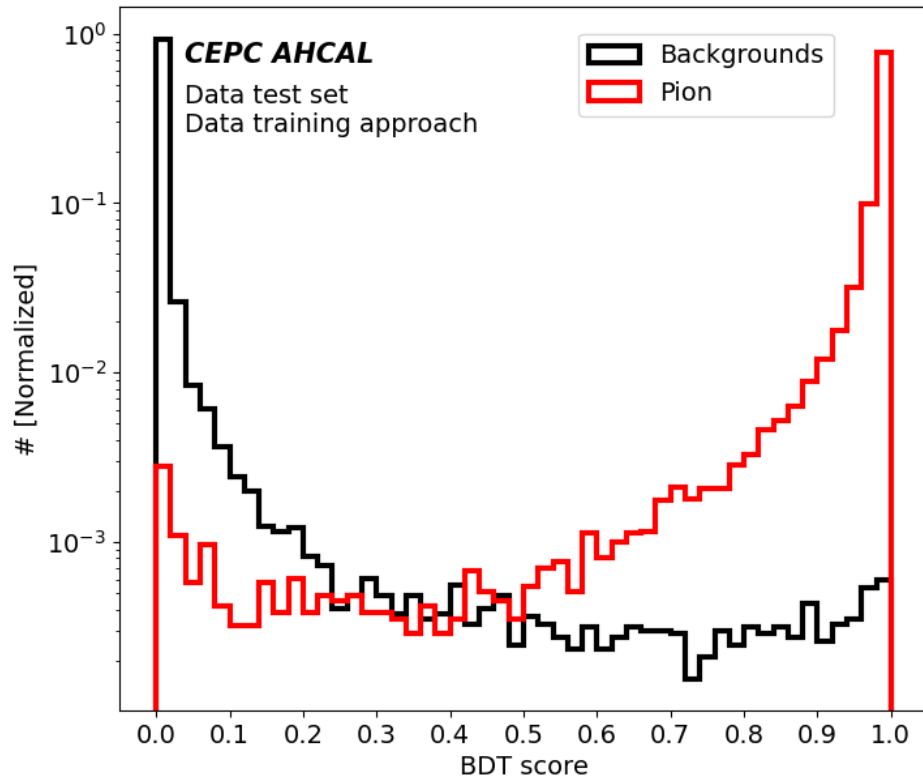
## MC training approach

- At 99% pion signal efficiency, Bkg. Rejection is 29.6 ( $N_{\text{Bkg.}}/N_{\text{Bkg.}}^{\text{sel.}}$ )



## Data training approach

- At 99% pion signal efficiency, Bkg. Rejection is 143.0 ( $N_{\text{Bkg.}}/N_{\text{Bkg.}}^{\text{sel.}}$ )

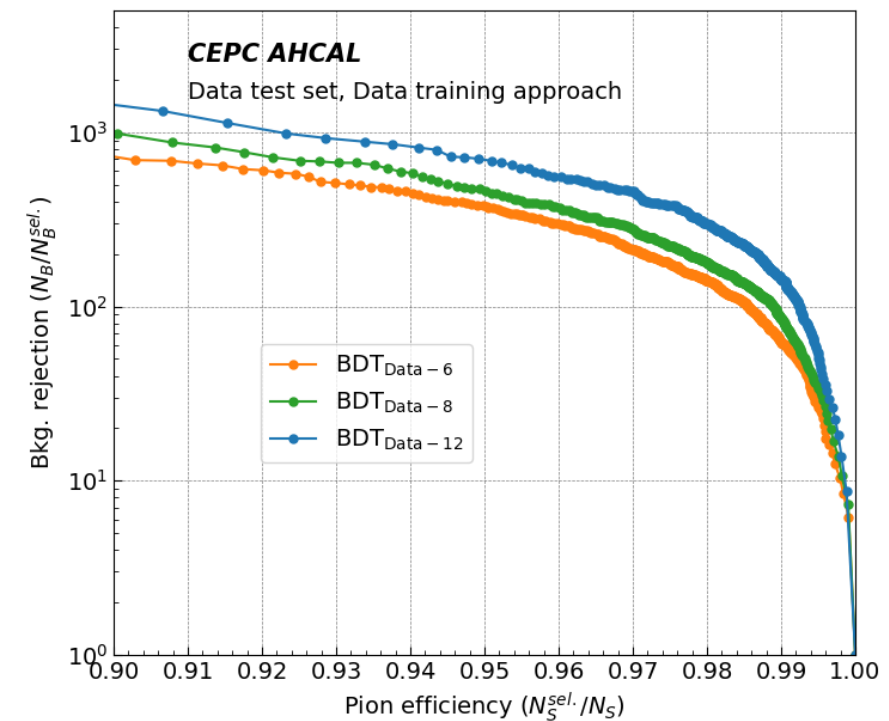
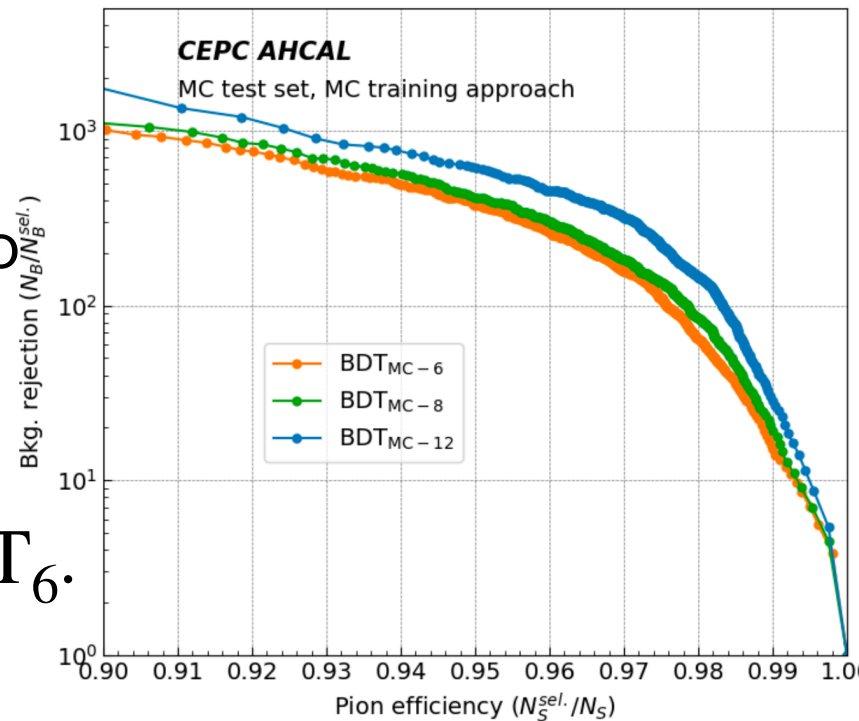




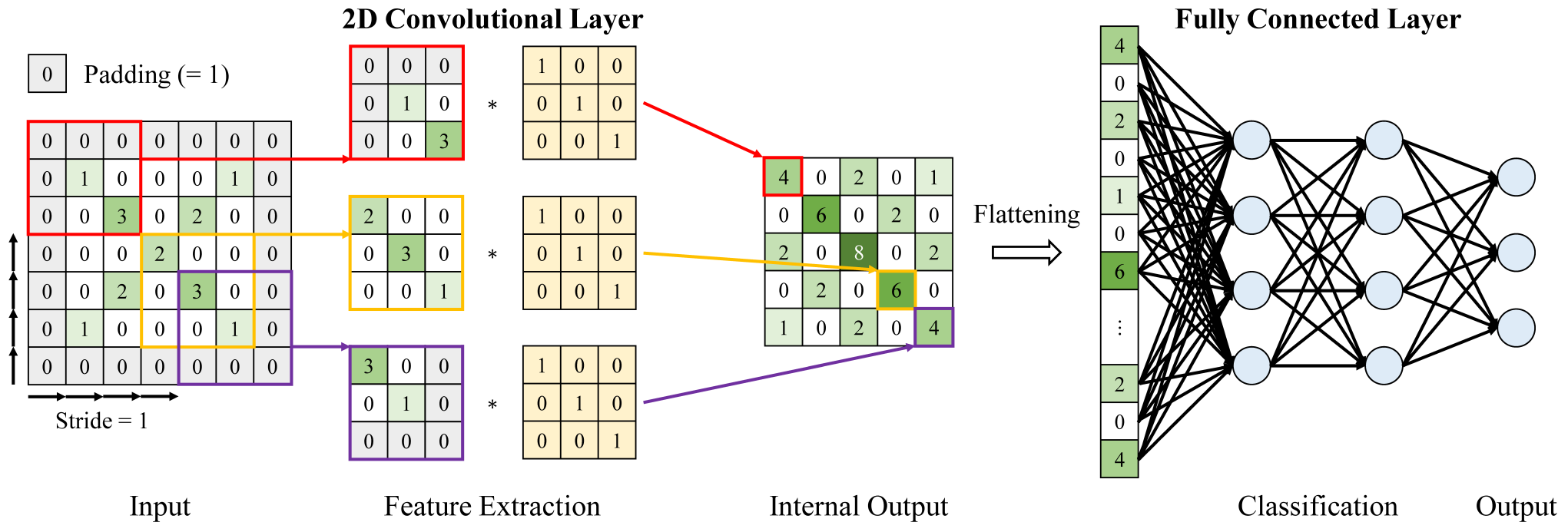
- **Dependence of BDT performance on input variables**

- Remove Shower End, Shower Layers, Fired Layers, and Z Width to build  $BDT_8$ .
- Further remove  $FD_1$  and  $FD_6$  to build  $BDT_6$ .

- **Feature engineering matters in BDT**
  - Increasing variable number can improve BDT



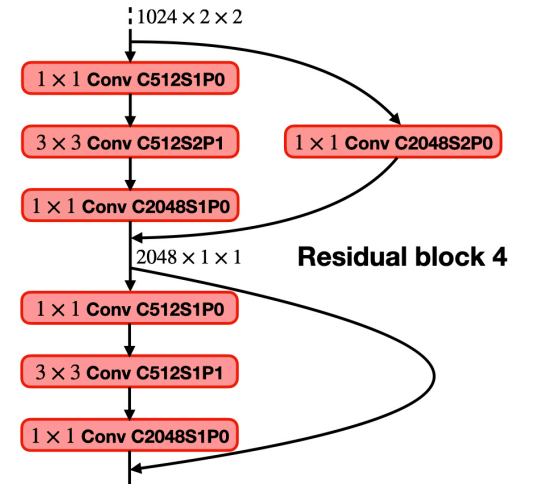
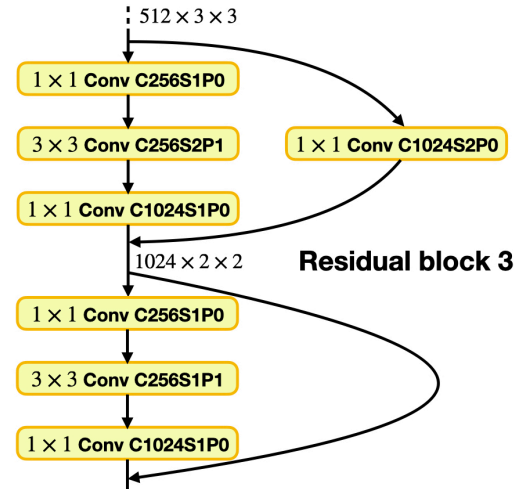
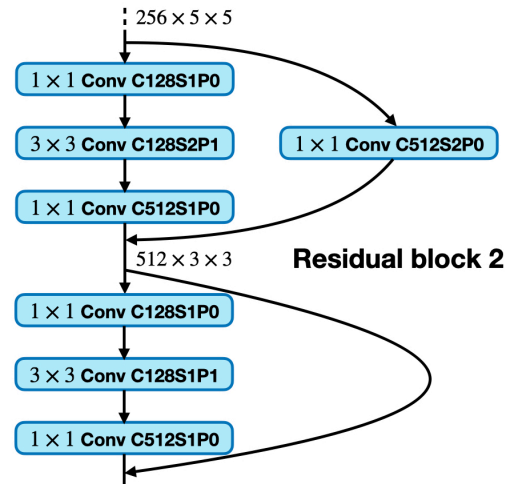
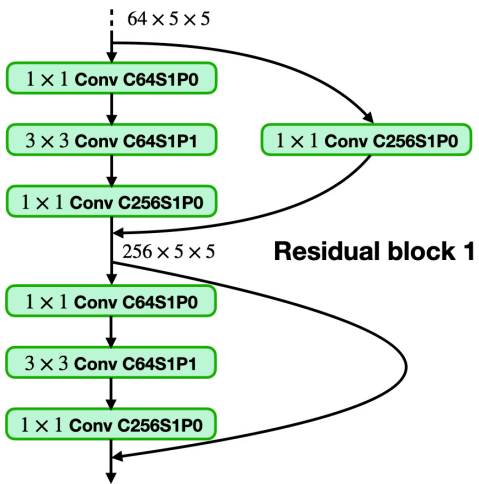
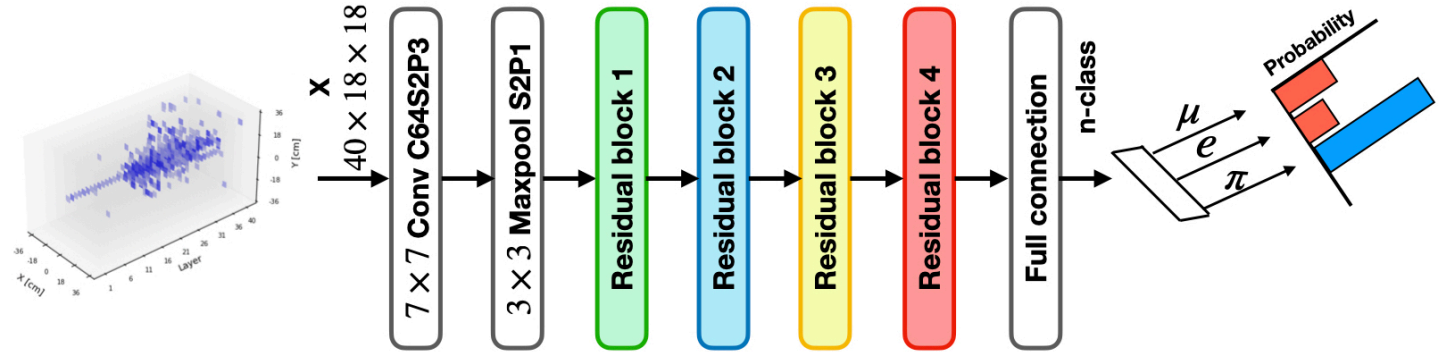
- **Cell-based Artificial Neural Networks (ANN)** make full use of **high-dimensional input** ( $18 \times 18 \times 40$ ).
  - Compile layers to extract features.
  - Output is the probability of each particle type candidate.



- Architecture: take the advantage of the Residual Block**

Input: energy deposits in AHCAL ( $18 \times 18 \times 40$ ).

output: probability of each particle type candidate.



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.

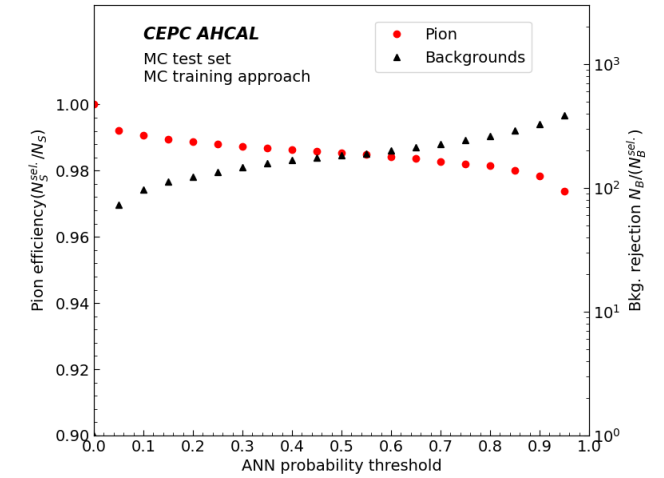
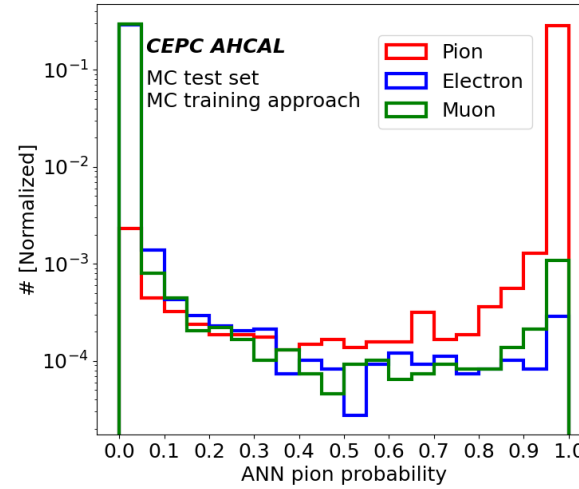


# PID based on ANN



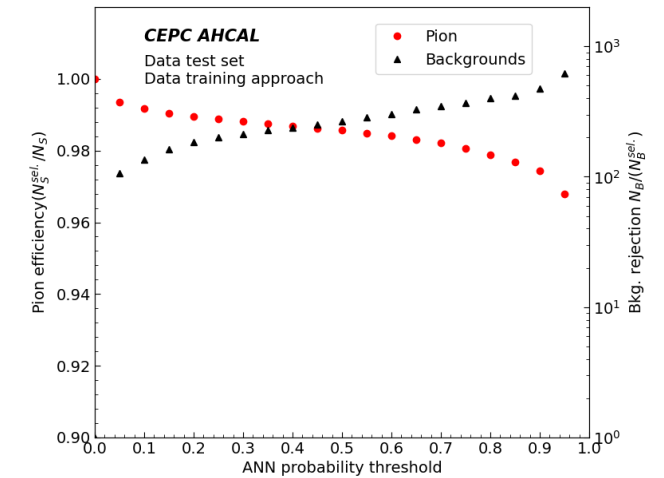
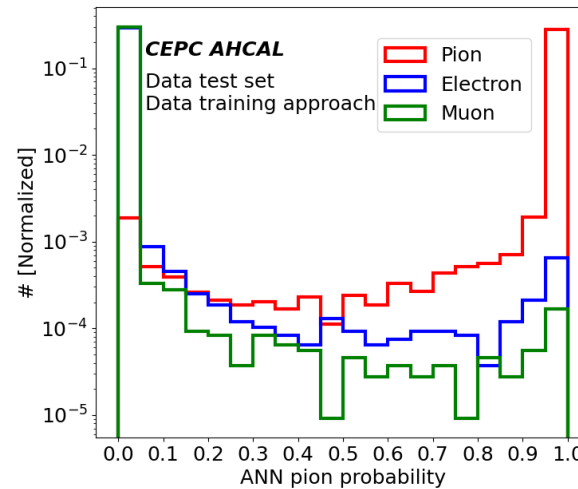
- **ANN<sub>MC</sub> is trained on MC samples.**

- At 99% pion efficiency, background rejection is 103.9, pion purity is 98.0%



- **ANN<sub>Data</sub> is trained on Test Beam samples.**

- At 99% pion efficiency, background rejection is 187.8, pion purity is 98.9%

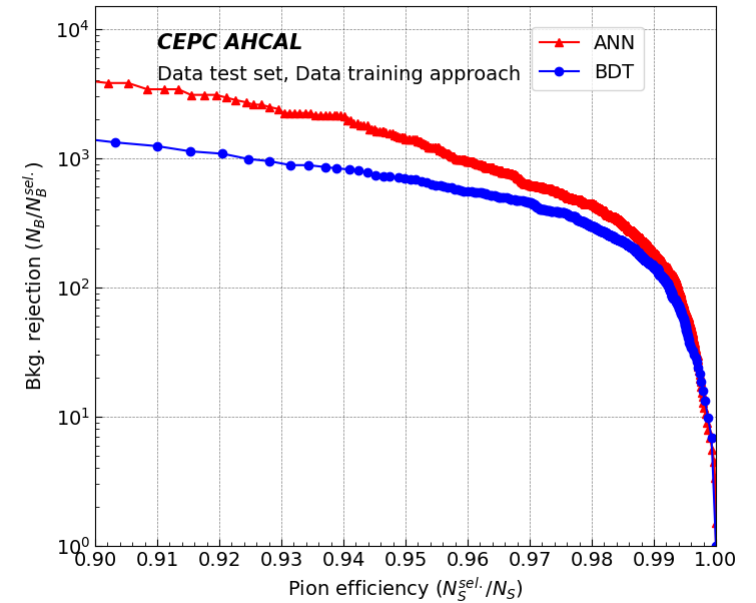
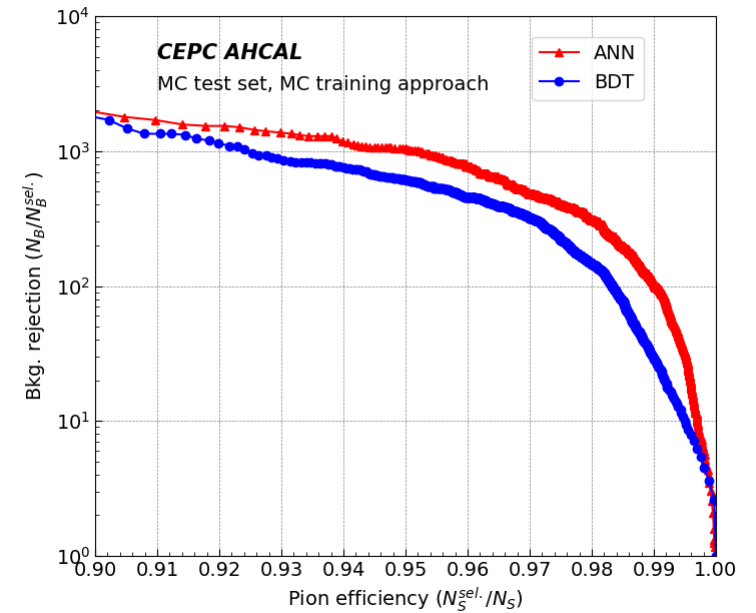




# Comparison between ANN and BDT

- We observe better performance of ANN in terms of Background rejection and Pion purity.

	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	<b>2015.7</b>	<b>3811.2</b>	<b>1040.3</b>	<b>1408.5</b>	<b>103.9</b>	<b>187.8</b>
Improvement	↑ 18.49%	↑ 163.12%	↑ 68.51%	↑ 103.65%	↑ 251.14%	↑ 31.37%
BDT pion purity	0.998	0.998	0.996	0.996	0.923	0.983
ANN pion purity	<b>0.999</b>	<b>0.999</b>	<b>0.998</b>	<b>0.998</b>	<b>0.980</b>	<b>0.989</b>
Improvement	↑ 0.05%	↑ 0.13%	↑ 0.21%	↑ 0.22%	↑ 6.20%	↑ 0.61%



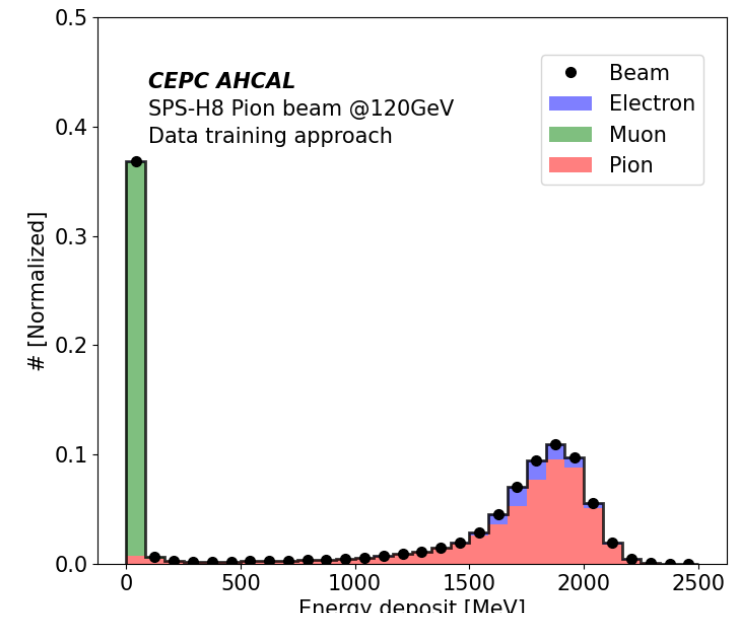
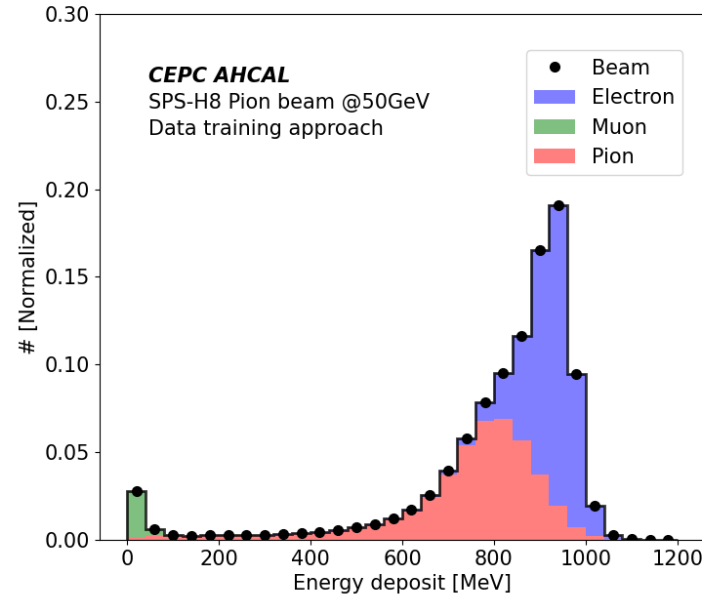
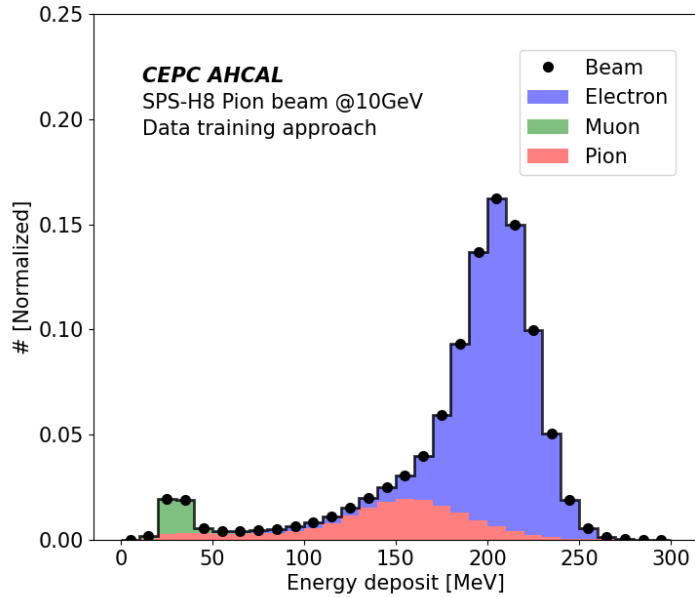


- **ANN outperforms BDT in our cases.**
  - **Automatic feature extraction:** This allows ANN to make full use of all input information, and potentially uncover hidden patterns in the data that may be missed by BDT, which relies on limited reconstructed features.
  - **Effective in handling large and 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.





Apply the ANN classifier to SPS-H8 beam collected in 2022.



Energy [GeV]	10	20	30	40	50	60	70	80	90	100	120
Muon percentage [%]	3.6	4.3	3.2	3.3	3.1	3.4	3.1	3.7	5.5	28.4	36.3
Electron percentage [%]	78.7	61.7	64.2	58.3	51.7	41.9	41.2	36.0	31.1	20.3	7.9
<b>Pion percentage [%]</b>	<b>17.7</b>	<b>34.0</b>	<b>32.6</b>	<b>38.5</b>	<b>45.3</b>	<b>54.8</b>	<b>55.7</b>	<b>60.3</b>	<b>63.4</b>	<b>51.3</b>	<b>55.8</b>



1. BDT and ANN based PID methods are all developed.
2. The preliminary purity of CEPC AHCAL test beam data on the type of incident particles is given.
3. This research promises the application prospect of cell-based ANN classifier in high-granularity calorimeters.



感谢聆听

饮水思源 爱国荣校<sup>19</sup>



# Backup

饮水思源 爱国荣校<sup>20</sup>

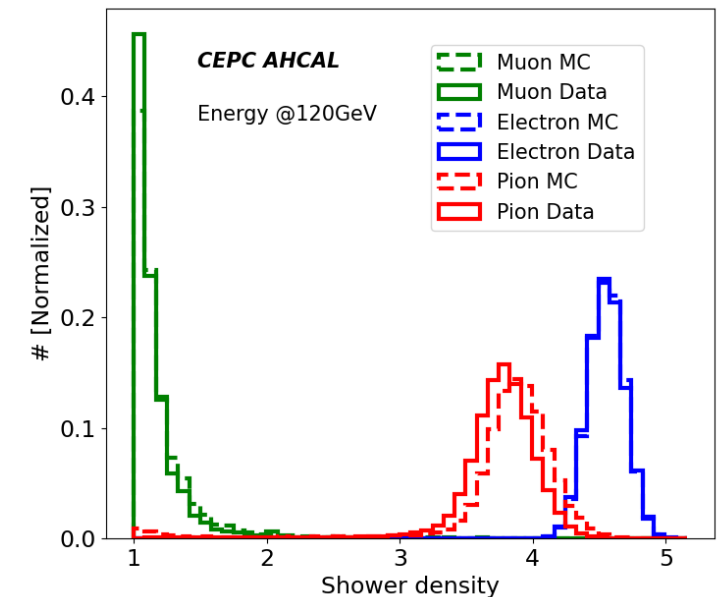
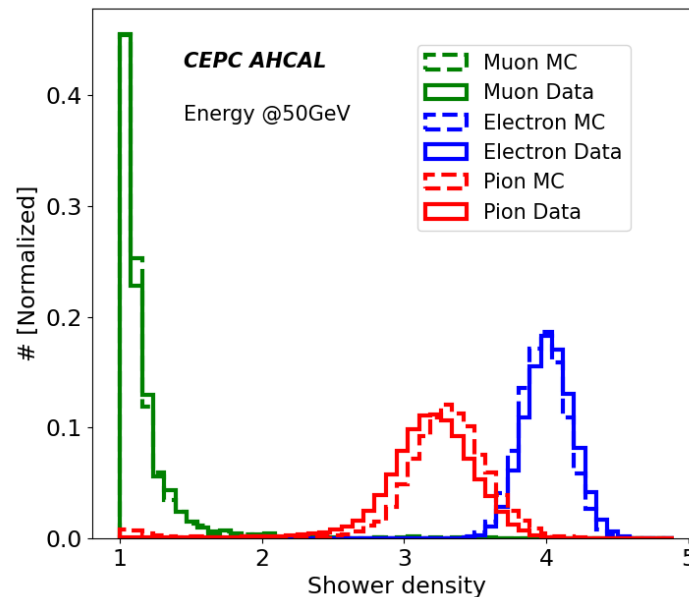
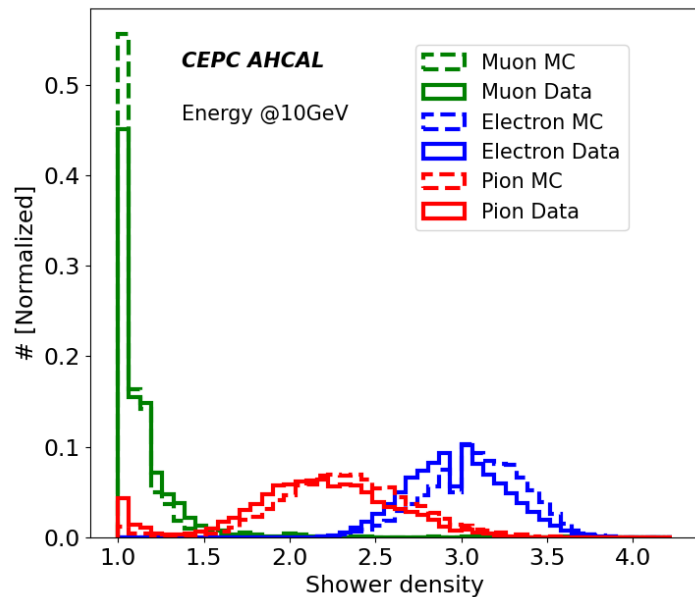


# MC & Data comparison



## Shower density

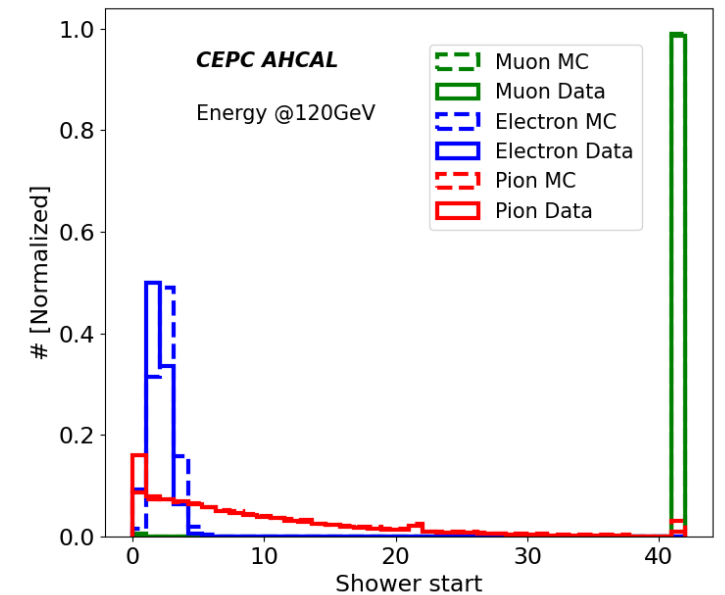
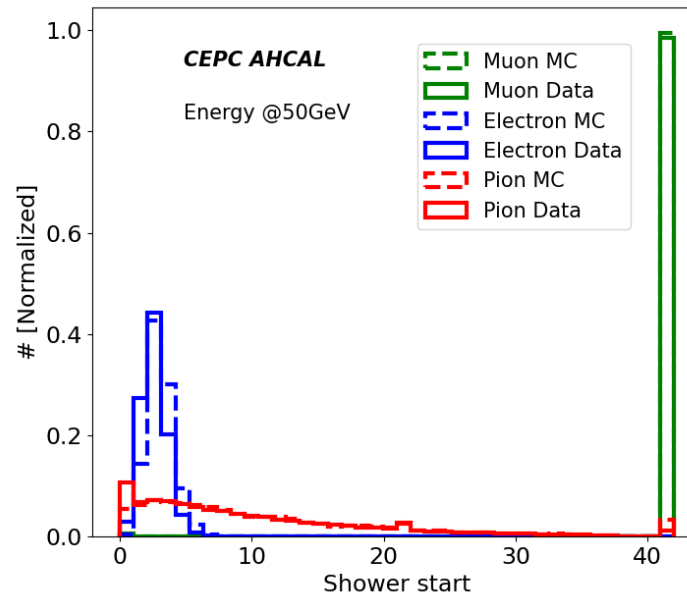
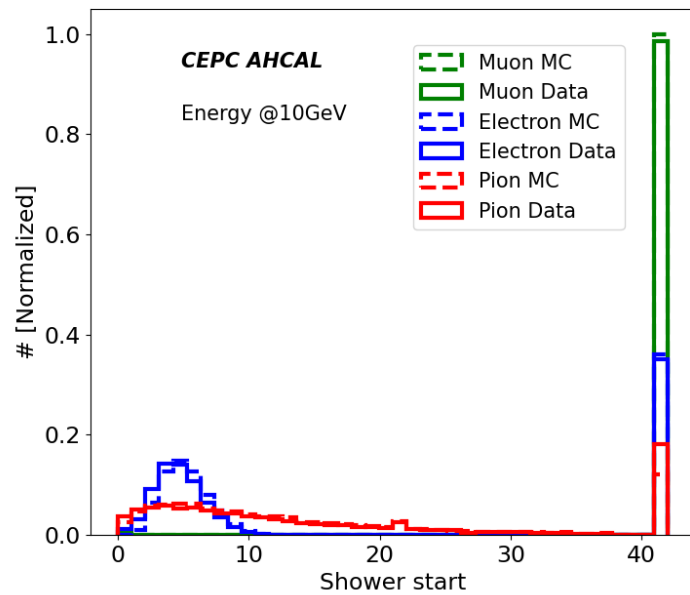
- The average number of the neighbouring hits located in the 3×3 cell around one of the hits including the hit itself.
- Data come from 2023's Data, MC samples come from: [run20230515\\_AHCAL\\_Shaping150ns\\_Window4us](#)
- The same 0.5 Mip hit energy threshold is pre-applied on both.
- They generally fit.





## Shower layer start

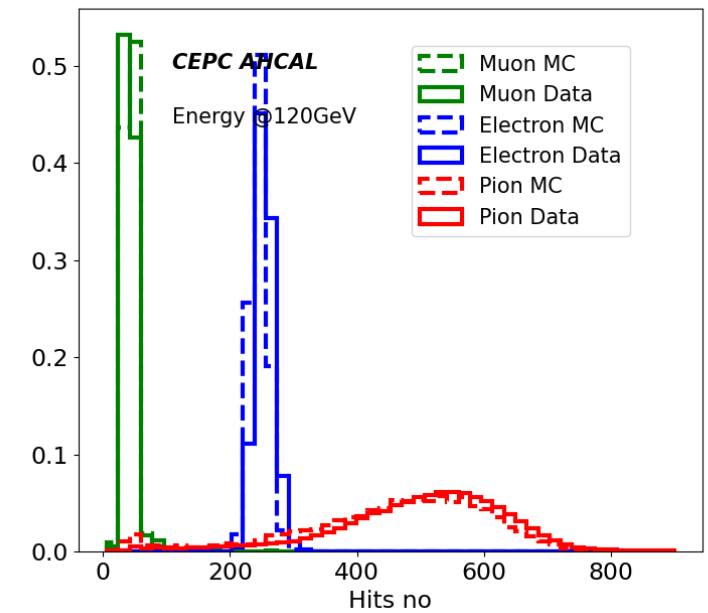
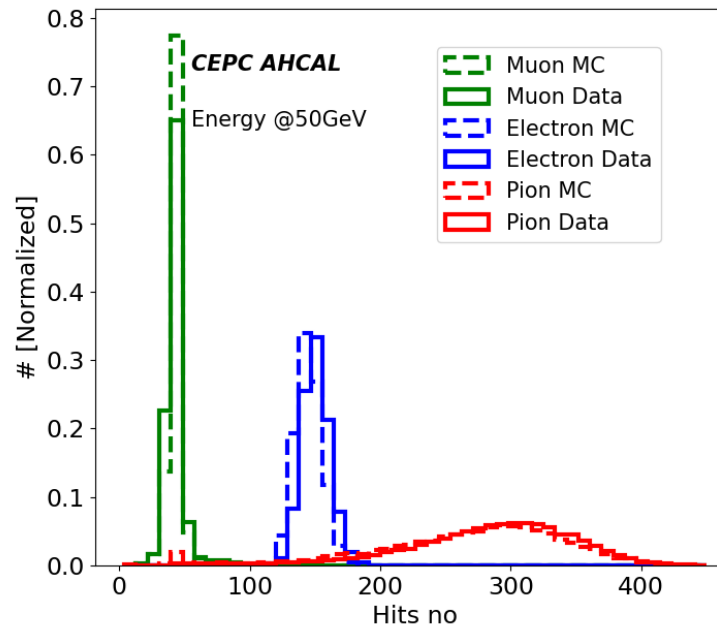
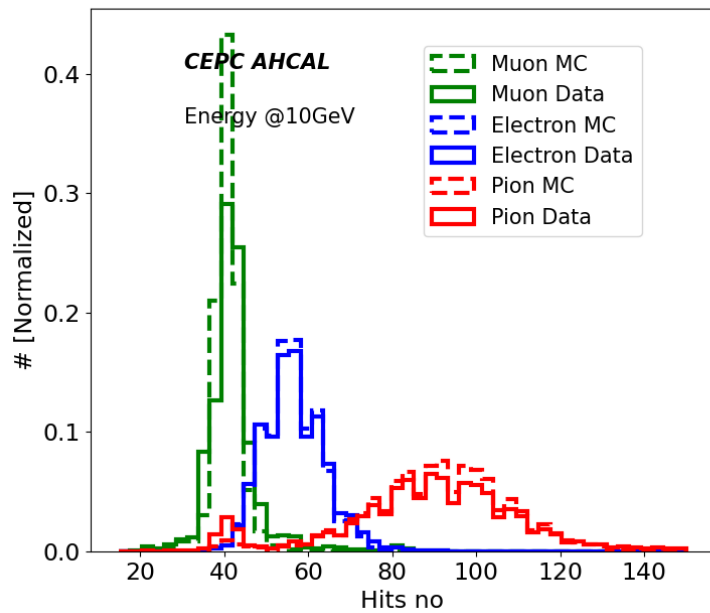
- The first layer of 3 consecutive layers with at least 5 hits. If no shower, it would be set as 42.
- Data come from 2023's Data, MC samples come from: [run20230515\\_AHCAL\\_Shaping150ns\\_Window4us](#)
- The same 0.5 Mip hit energy threshold is pre-applied on both.
- They generally fit.





## Hits number

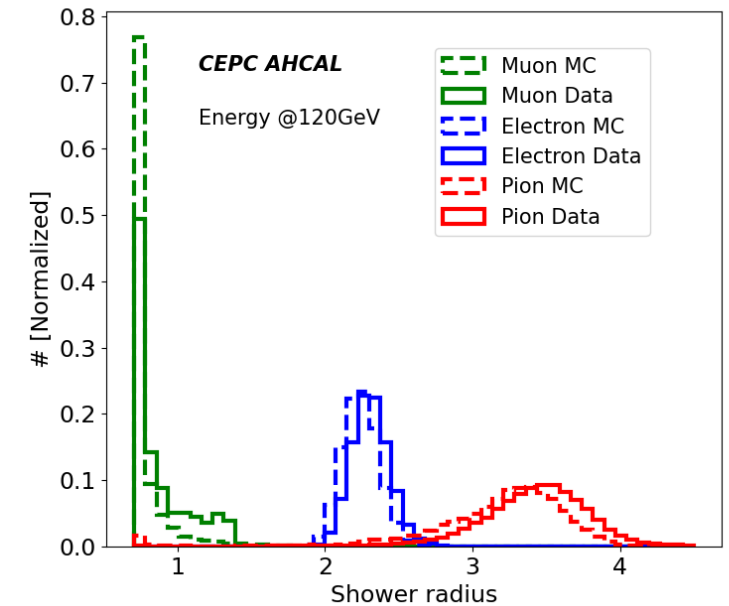
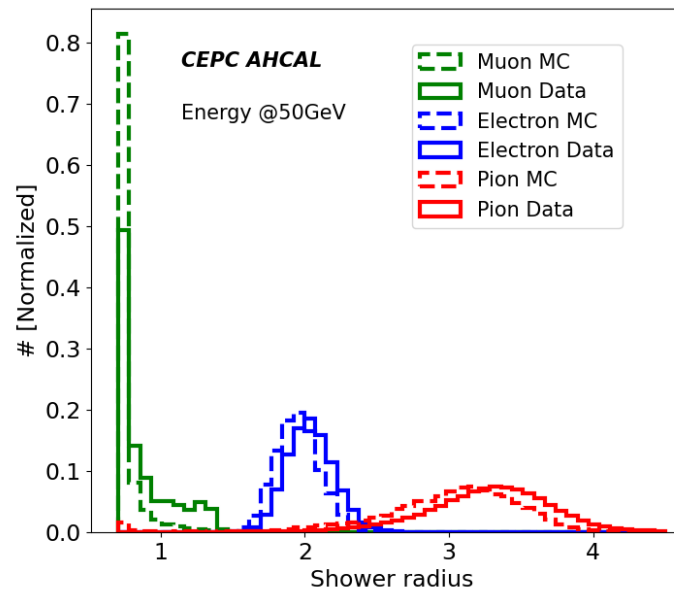
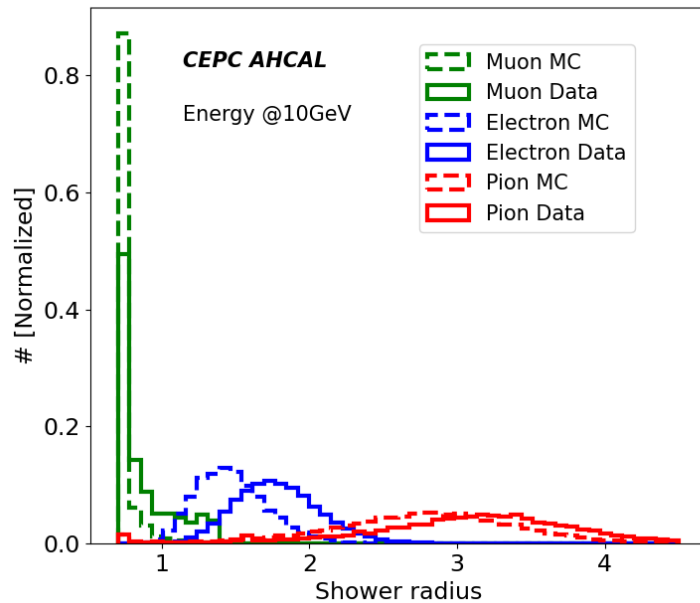
- The same 0.5 Mip hit energy threshold is pre-applied on both.
- They generally fit.





## Shower radius

- The average rms of the hits distance to the longitudinal axis of AHCAL (cross the center).
- Data come from 2023's Data, MC samples come from: [run20230515\\_AHCAL\\_Shaping150ns\\_Window4us](#)
- The same 0.5 Mip hit energy threshold is pre-applied on both.
- Deviation of peak value is observed.

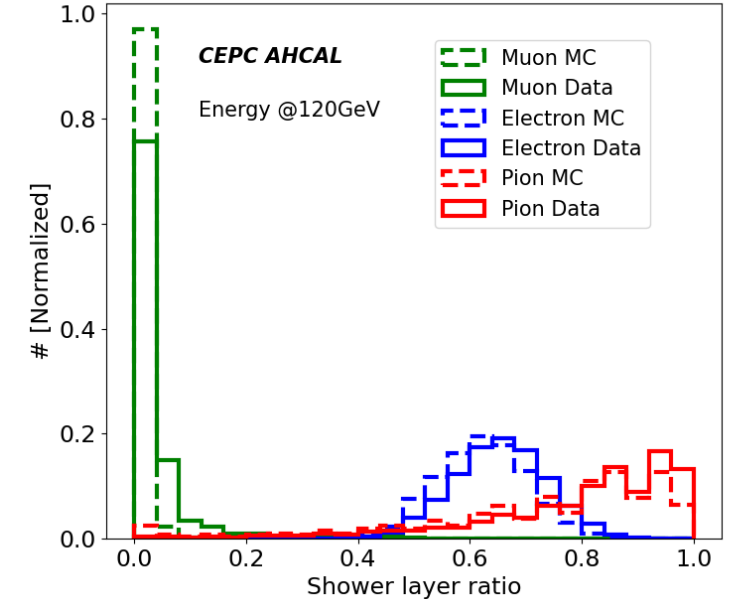
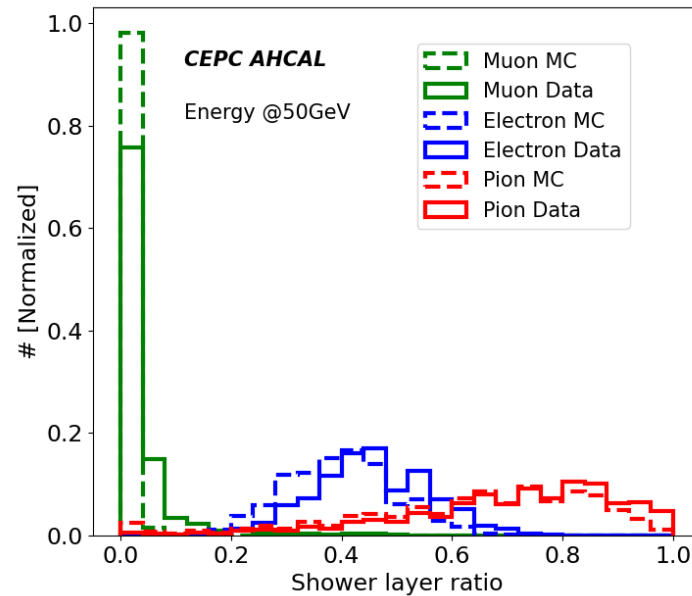
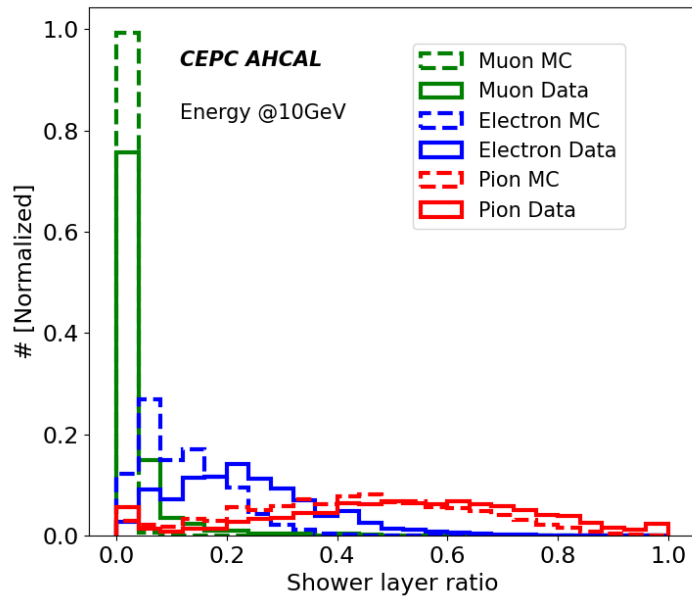






## Ratio of shower layers over total hit layers

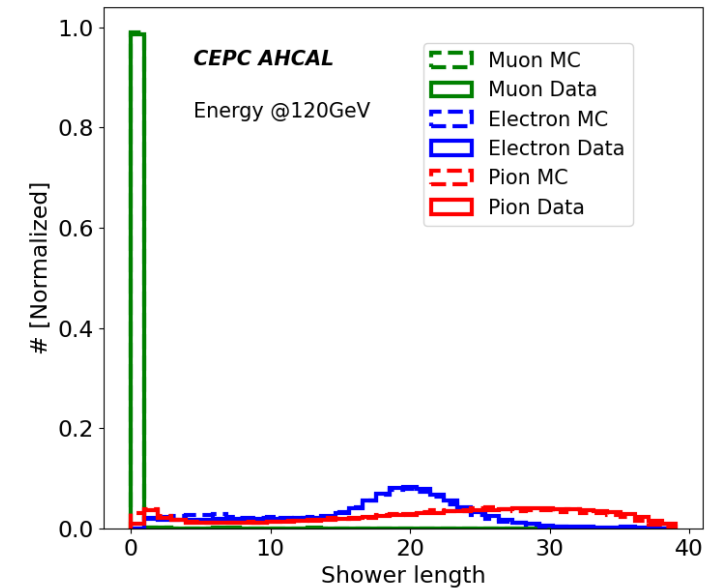
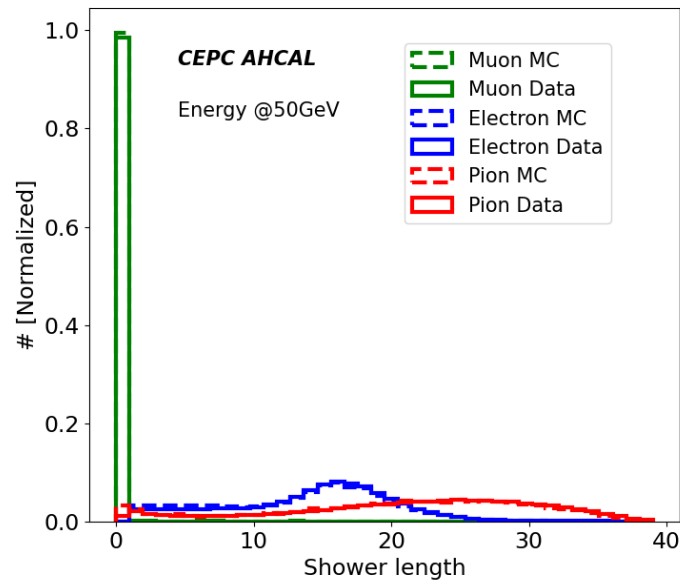
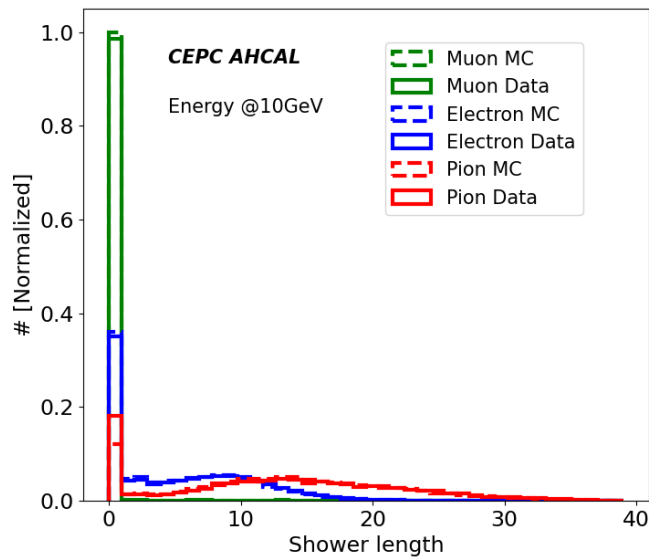
- The ratio between the number of layers in which the Root Mean Square (RMS) of the hits' position in the x-y plane exceeds 4 cm in both x and y directions and the total number of layers with at least one fired cell.
- Data come from 2023's Data, MC samples come from: [run20230515\\_AHCAL\\_Shaping150ns\\_Window4us](#)
- The same 0.5 Mip hit energy threshold is pre-applied on both.
- They generally fit.





## Shower length

- This is the distance between the start of the shower and the layer where the maximum RMS of hit transverse coordinates with respect to the z-axis occurs.
- Data come from 2023's Data, MC samples come from: [run20230515\\_AHCAL\\_Shaping150ns\\_Window4us](#)
- The same 0.5 Mip hit energy threshold is pre-applied on both.
- They generally fit.



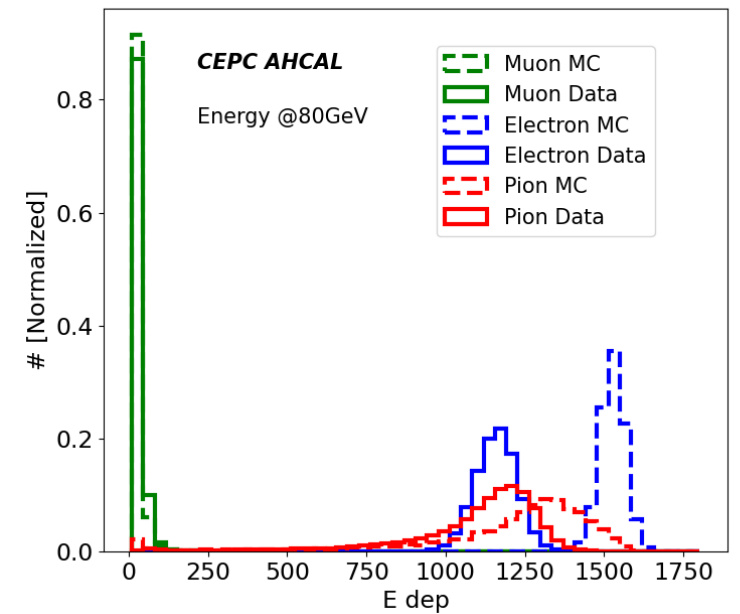
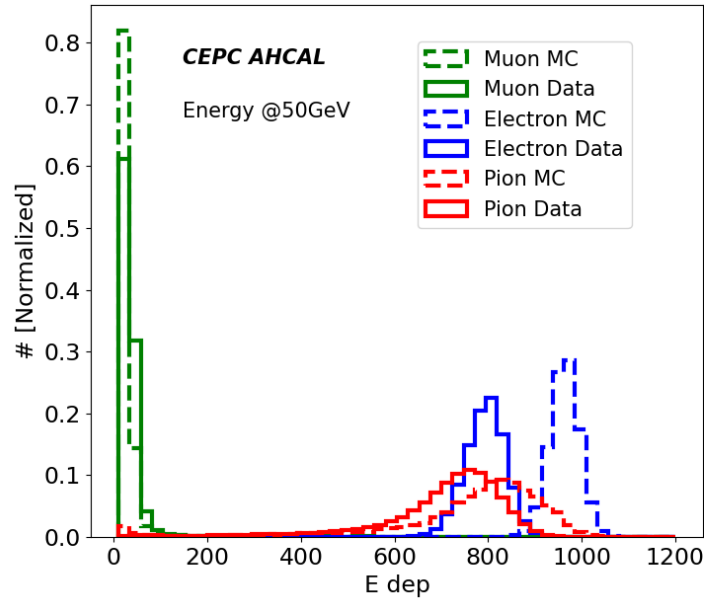
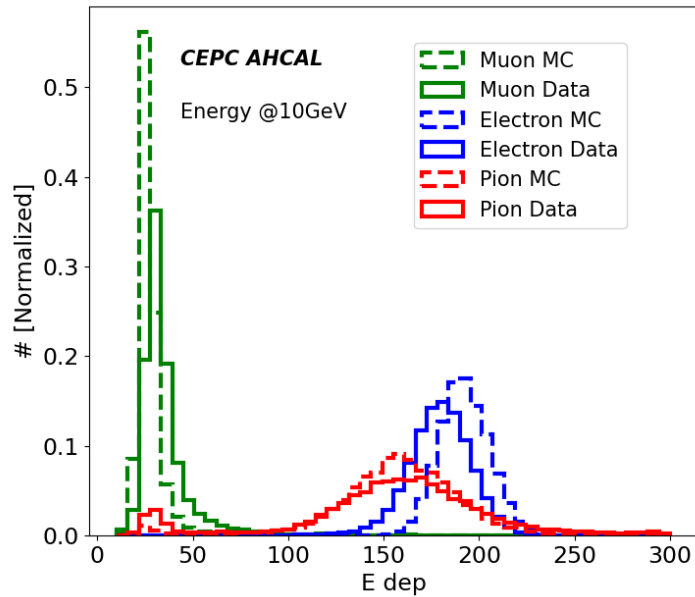


# MC & Data comparison



## Total energy deposit.

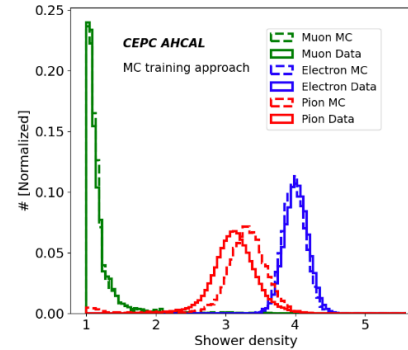
- Data come from 2023's Data, MC samples come from: [run20230515\\_AHCAL\\_Shaping150ns\\_Window4us](#)
- The same 0.5 Mip hit energy threshold is pre-applied on both.
- Obvious discrepancy appeared, especially on electron events.



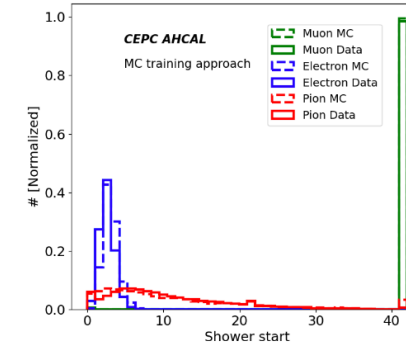


# 2022 SPS-H8 beam composition

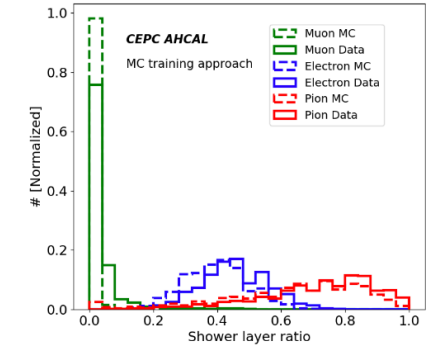
- The Pion samples are classified by using ANN trained on MC data set and are compared with MC.
- Not large disagreement observed.



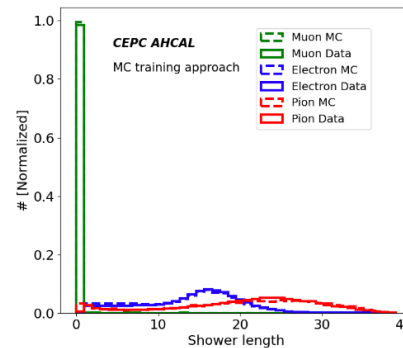
(a) Distribution of the shower density.



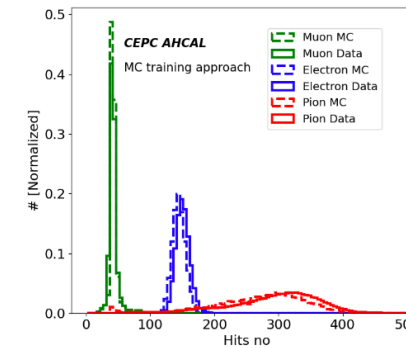
(b) Distribution of the layer that the shower starts.



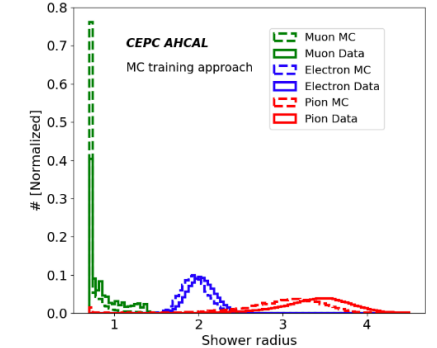
(c) Distribution of the ratio of shower layers over total hit layers.



(d) Distribution of the shower length.



(e) Distribution of the total number of hits.

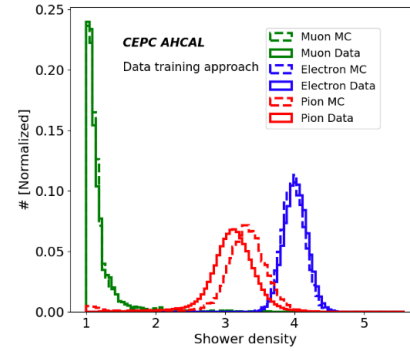


(f) Distribution of the shower radius.

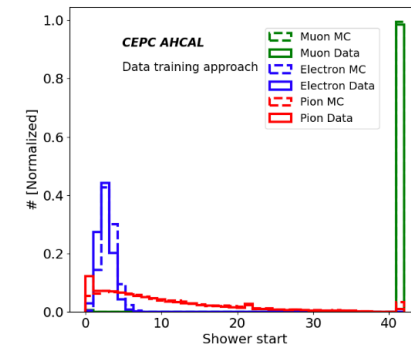


# 2022 SPS-H8 beam composition

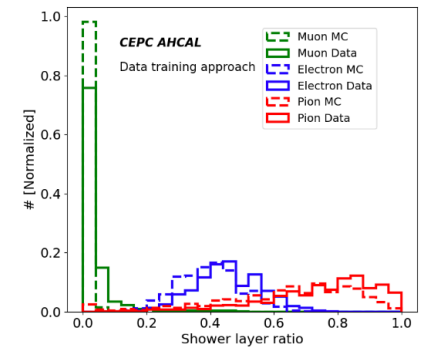
- The Pion samples are classified by using ANN trained on TB data set and are compared with MC.
- Not large disagreement observed.



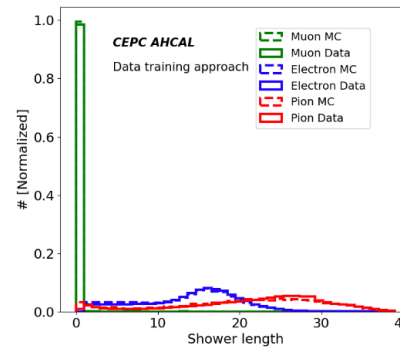
(a) Distribution of the shower density.



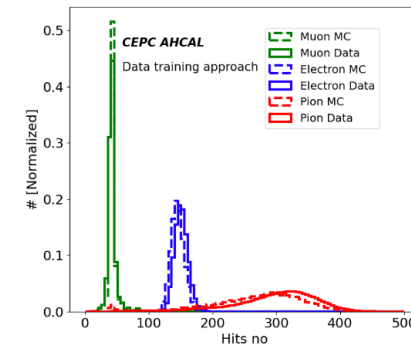
(b) Distribution of the layer that the shower starts.



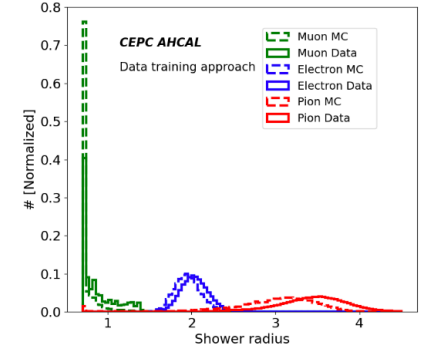
(c) Distribution of the ratio of shower layers over total hit layers.



(d) Distribution of the shower length.

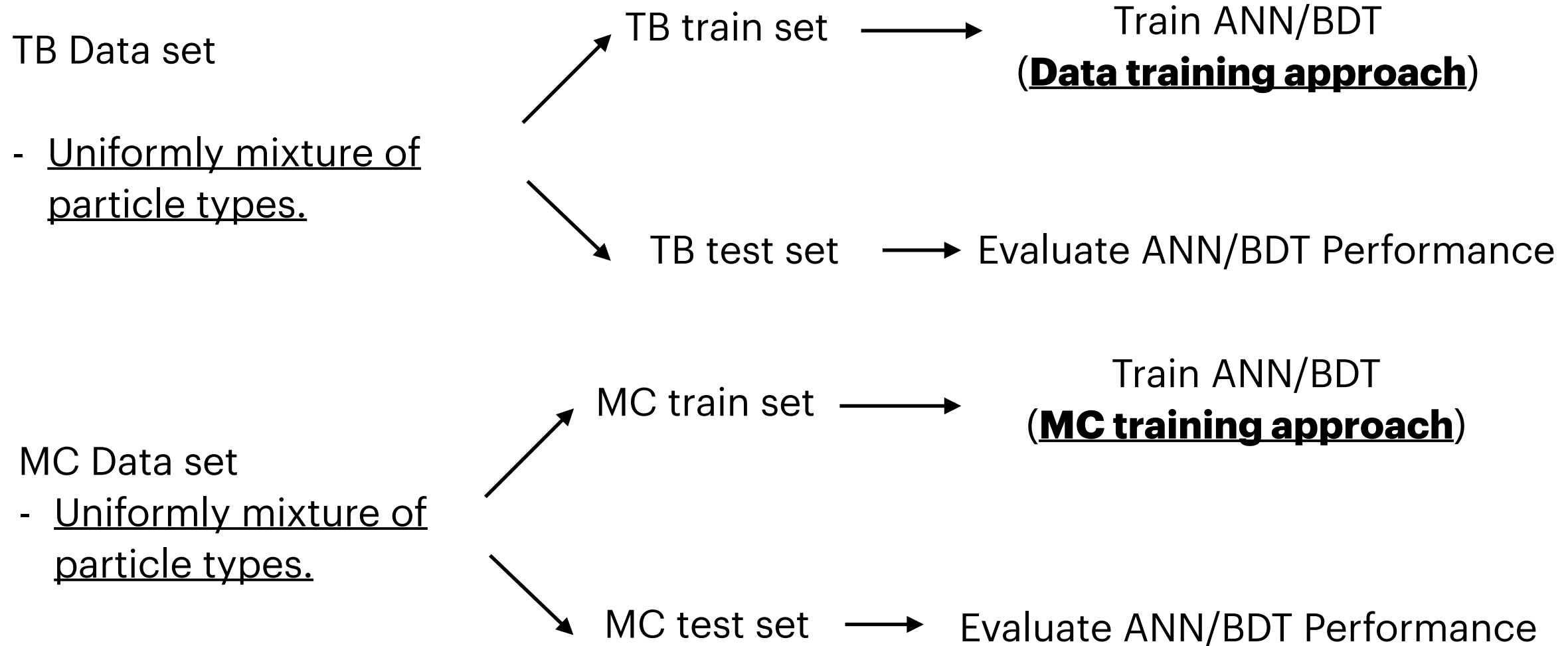


(e) Distribution of the total number of hits.



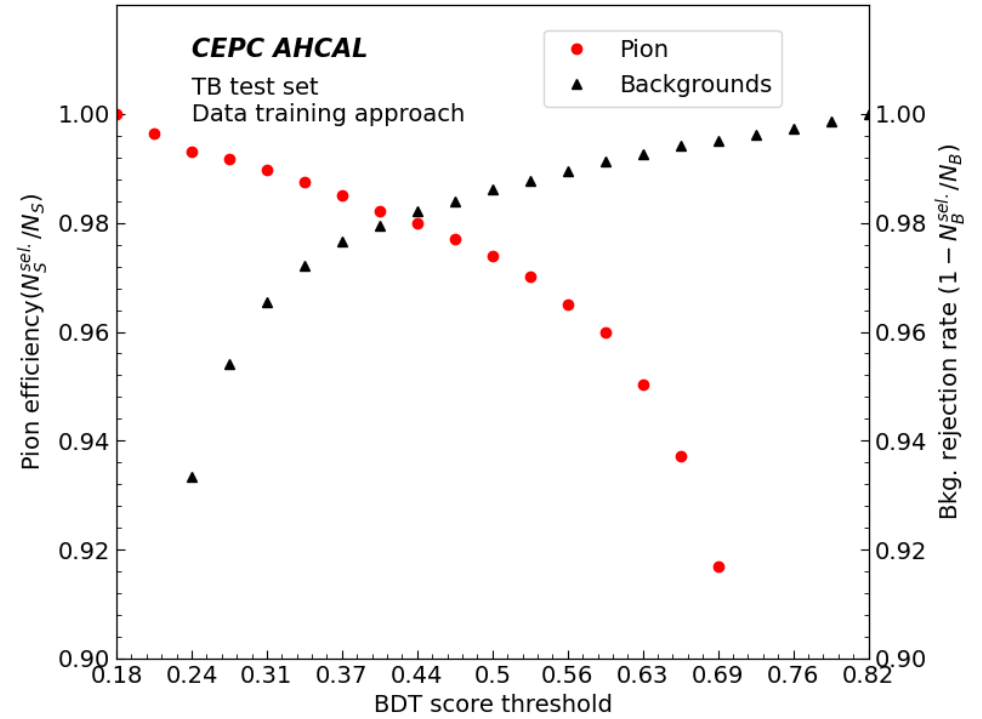
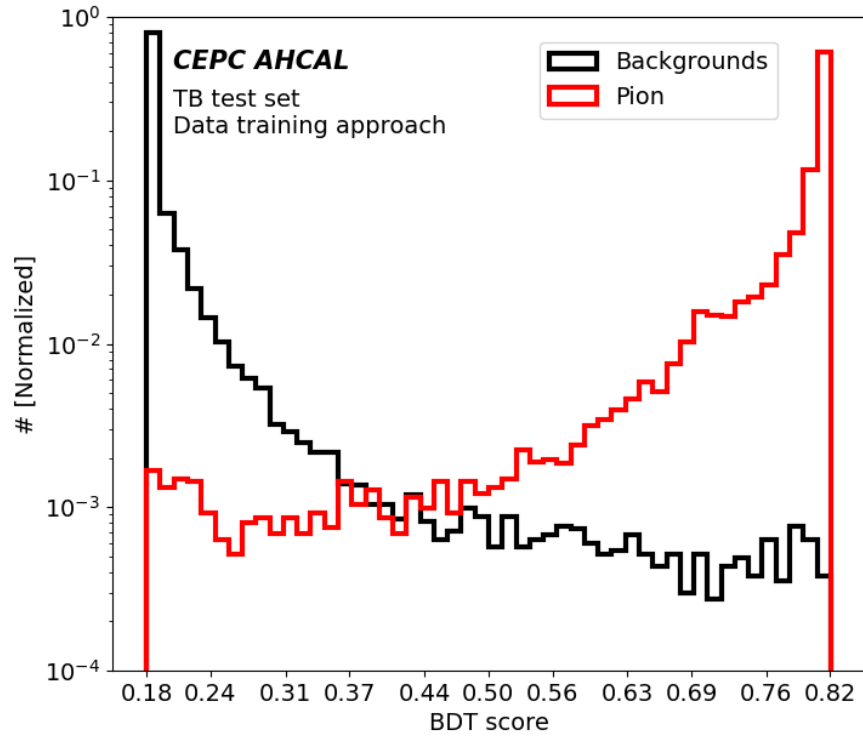
(f) Distribution of the shower radius.

# Two Data sets are prepared due to discrepancy between MC and Data





# BDT Output



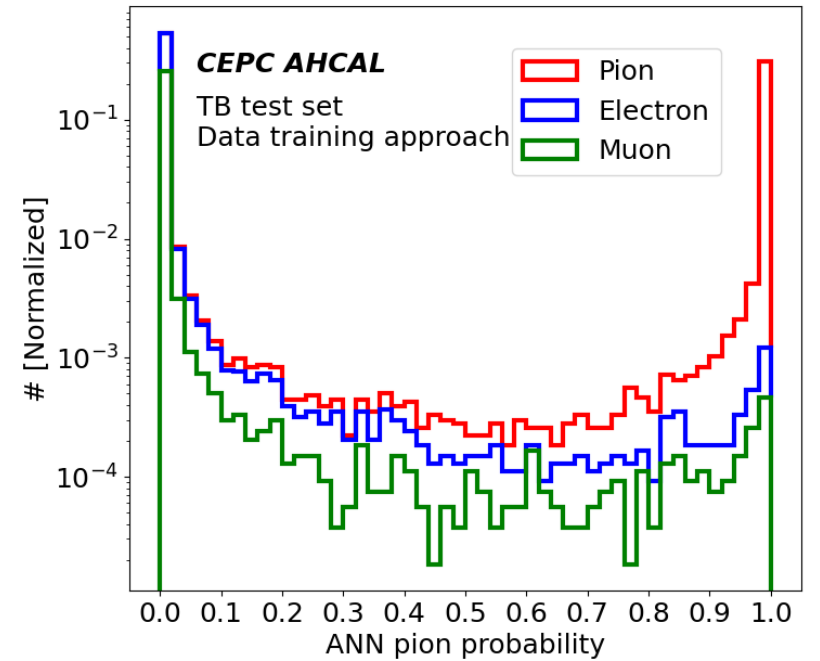
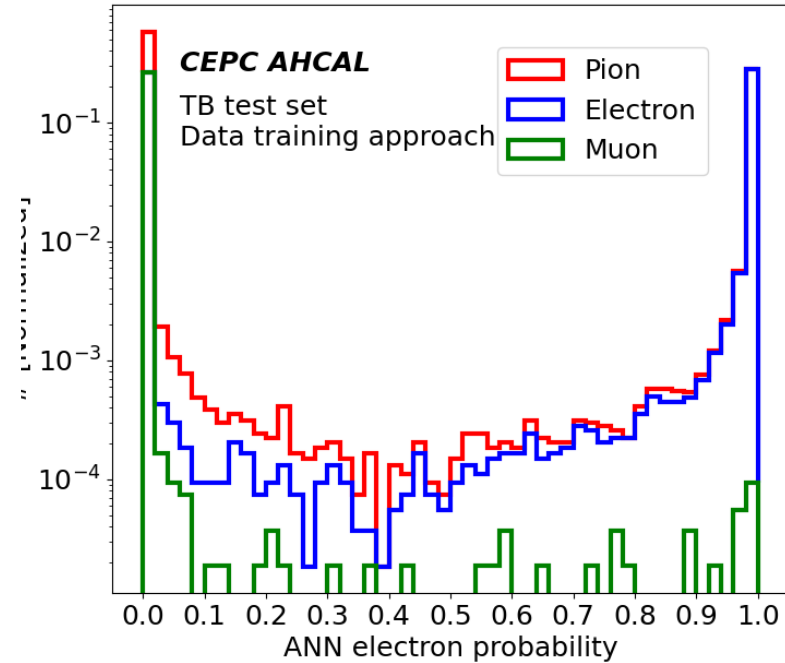
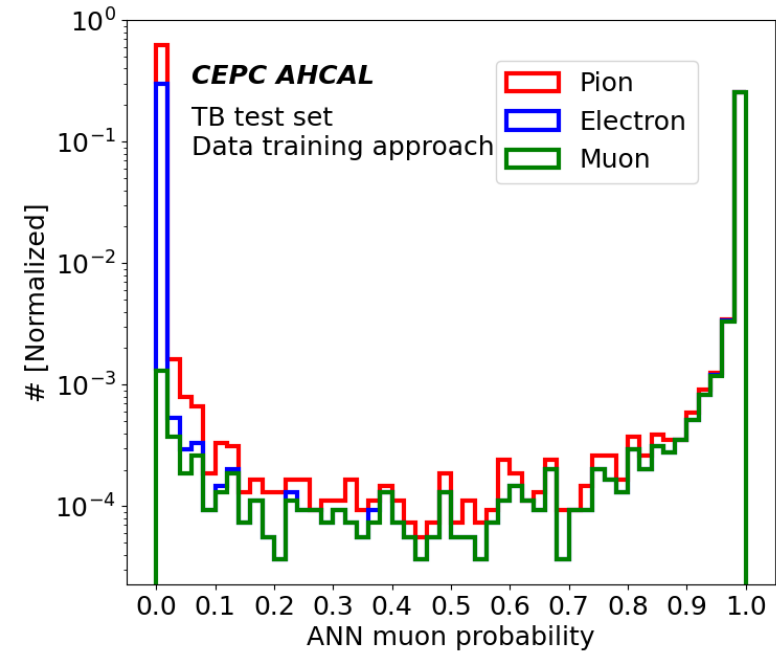
Backgrounds :Signal =2:1



# ANN Output



Separation power achieved. e.g. Pions get higher probability classified as pions.



Backgrounds :Signal =2:1

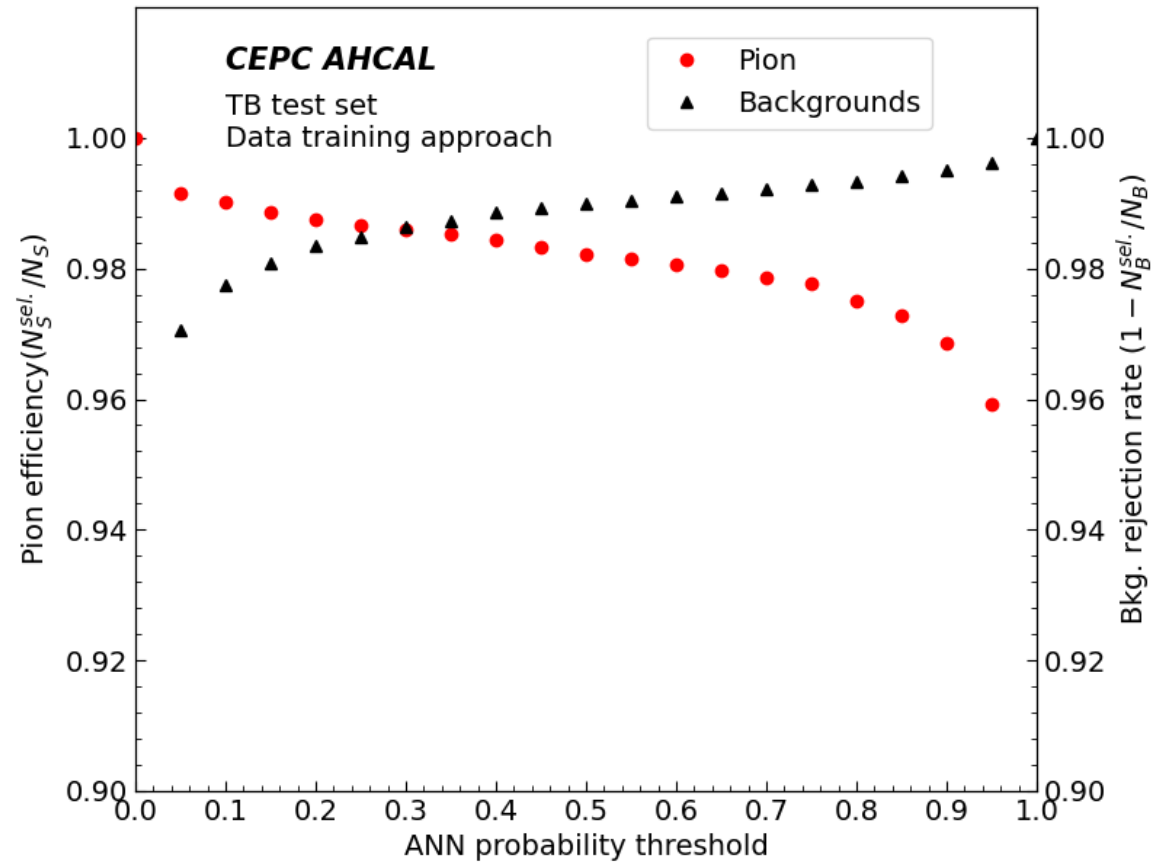




# ANN Output



- Separation power is confirmed.
  - Pion efficiency and background rejection rate can be both over 98%



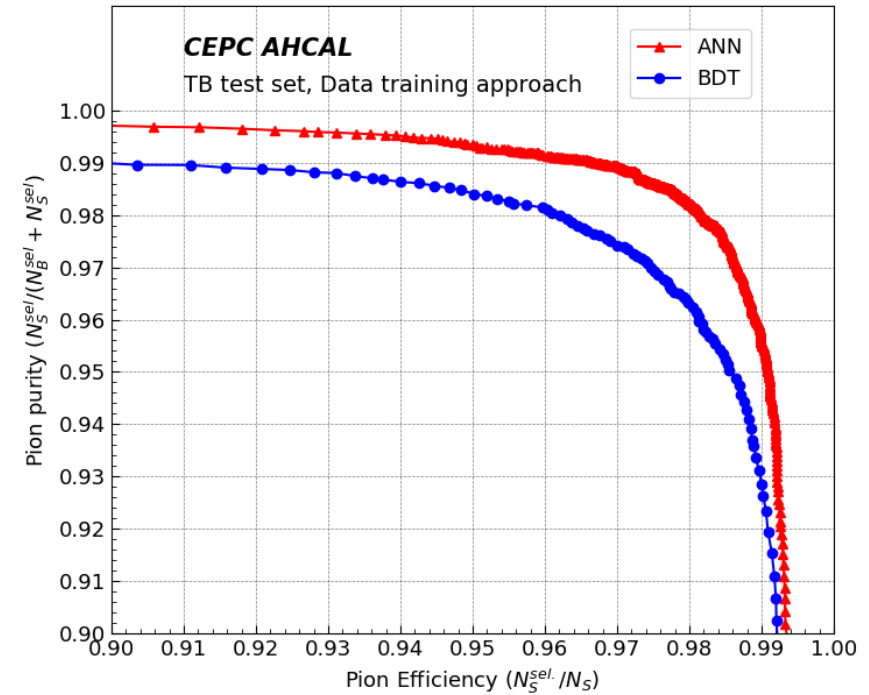
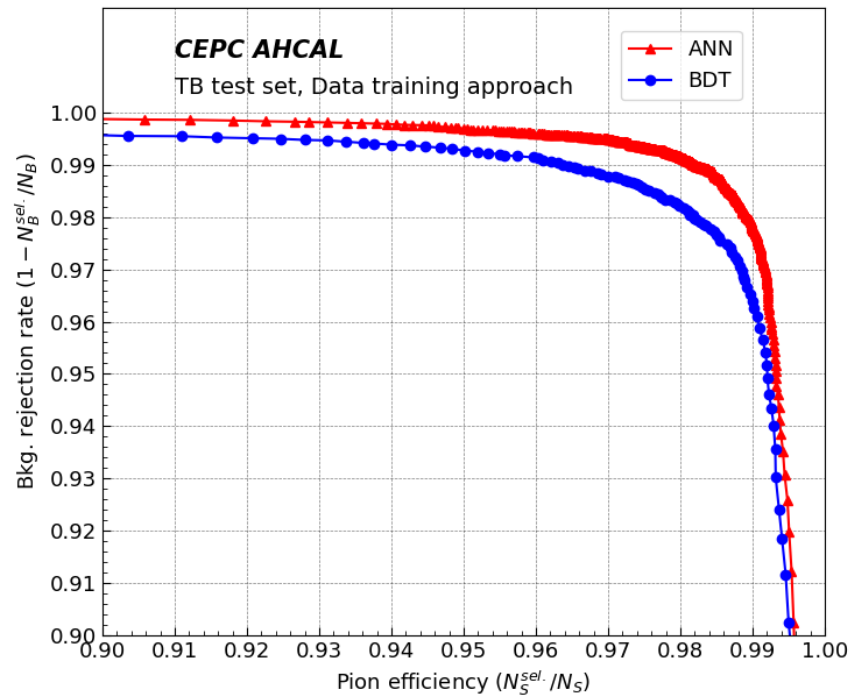


# Comparison with ANN PID method



Pion efficiency	0.90	0.92	0.94	0.96	0.98	0.99
BDT <sub>TB</sub> bkg. rej. rate	0.996	0.995	0.994	0.991	0.982	0.964
ANN <sub>TB</sub> bkg. rej. rate	<b>0.999</b>	<b>0.998</b>	<b>0.998</b>	<b>0.996</b>	<b>0.991</b>	<b>0.978</b>
	↑ <b>0.3%</b>	↑ <b>0.3%</b>	↑ <b>0.4%</b>	↑ <b>0.5%</b>	↑ <b>0.9%</b>	↑ <b>1.5%</b>
BDT <sub>TB</sub> pion purity	0.99	0.989	0.986	0.981	0.963	0.929
ANN <sub>TB</sub> pion purity	<b>0.997</b>	<b>0.996</b>	<b>0.995</b>	<b>0.991</b>	<b>0.982</b>	<b>0.956</b>
	↑ <b>0.7%</b>	↑ <b>0.7%</b>	↑ <b>0.9%</b>	↑ <b>1.0%</b>	↑ <b>2.0%</b>	↑ <b>2.9%</b>

TABLE V. The background rejection rate and the pion purity of the BDT<sub>TB</sub> and the ANN<sub>TB</sub>. The results highlight the improvement of the ANN<sub>TB</sub> compared to the BDT<sub>TB</sub>.

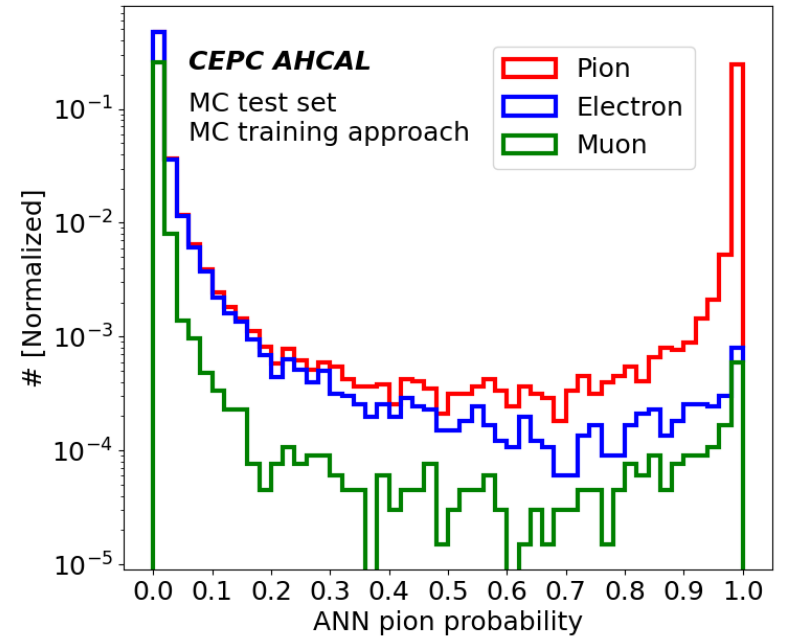
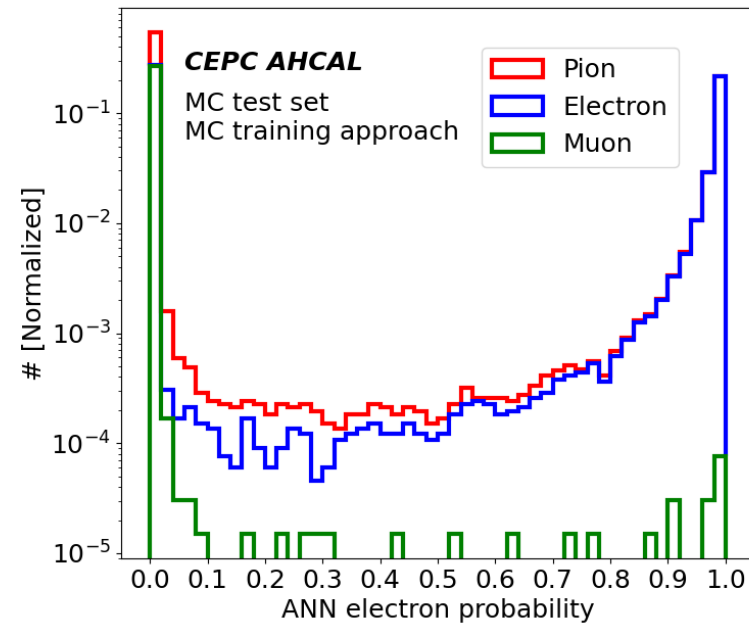
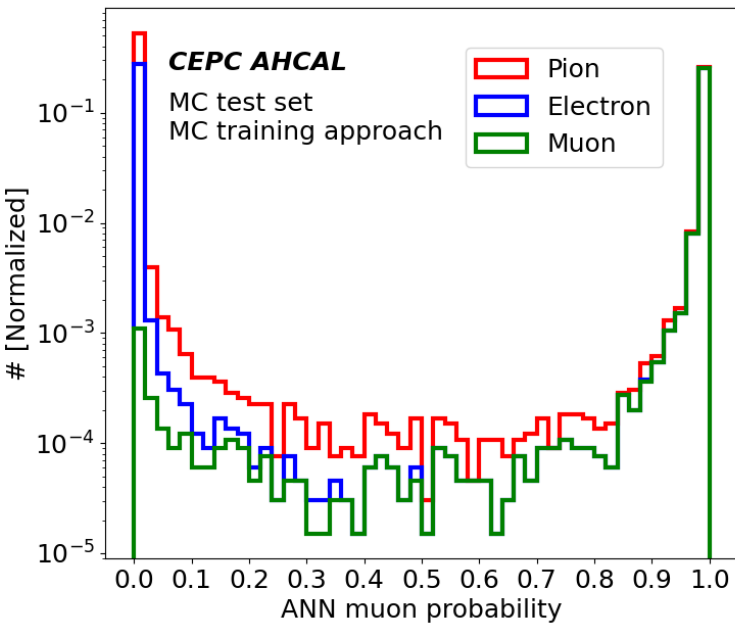




# ANN Output



Separation power achieved. e.g. Pions get higher probability classified as pions.



Backgrounds :Signal =2:1



- **Apply Extreme Gradient Boosting**
- **Reconstruct 12 input variables**

- **Shower density**
- **Shower start**
- **Shower length**
- **Hits number**
- **Shower radius**
- **Fractal dimension**
  - $FD_1, FD_6$
- **Fired layers**
- **Shower layers**
- **Shower layer ratio**
- **Z width**

MC samples

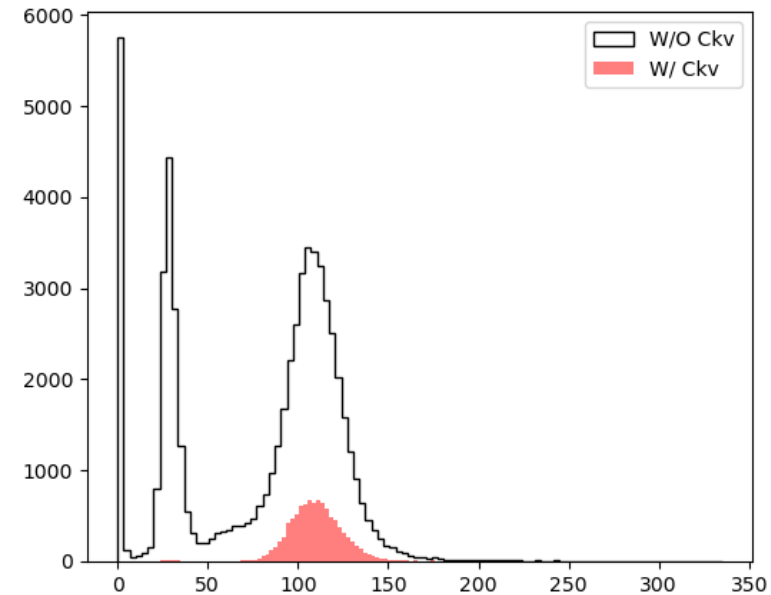
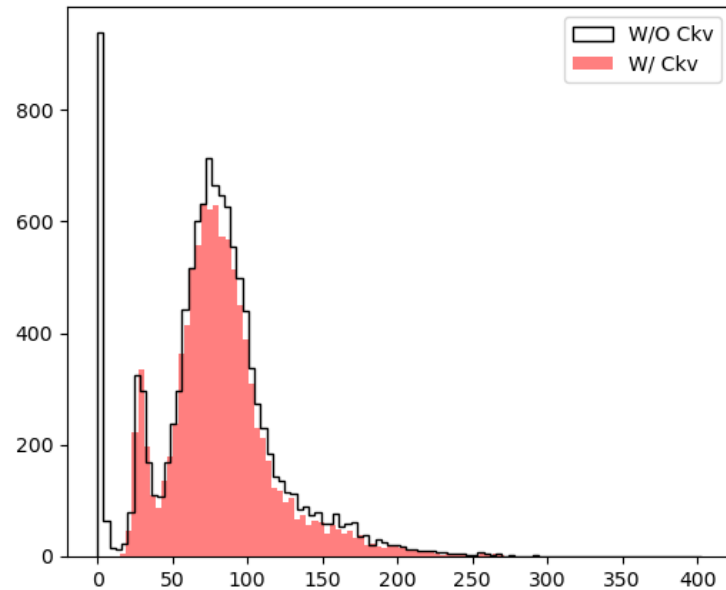
Rank: Variable	Variable weight
1: Shower radius	0.377
2: Shower layers	0.232
3: Hits number	0.088
4: Fired layers	0.083
5: Shower start	0.080
6: Shower density	0.049
7: Z width	0.034
8: $FD_6$	0.017
9: $FD_1$	0.015
10: Shower layer ratio	0.014
11: Shower end	0.006
12: Shower length	0.006

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

- The sensitivity of the Cherenkov detector in the low energy range (<30 GeV) could allow the collection of electron and pion samples.

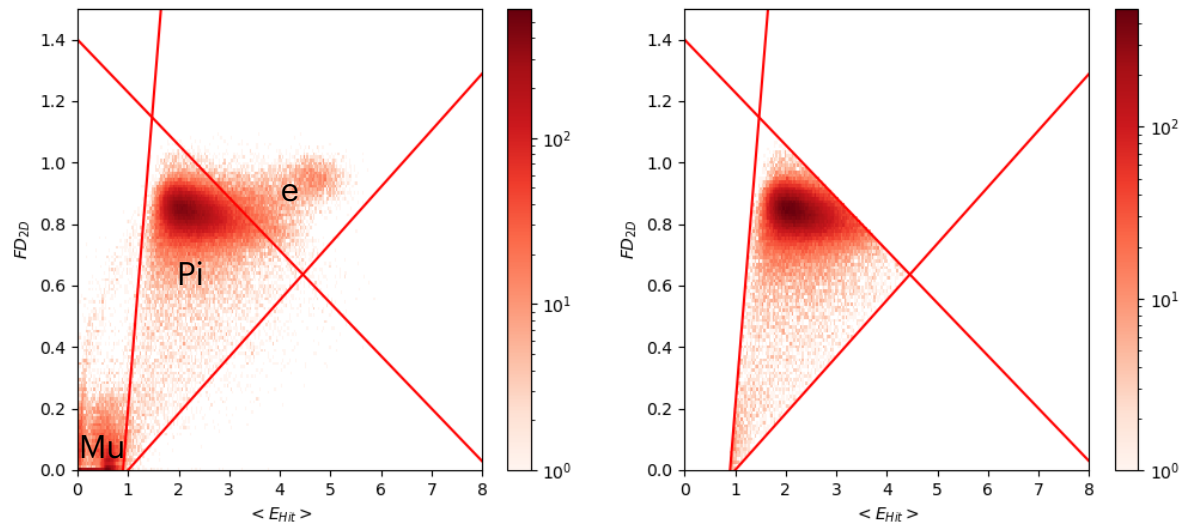


- Cherenkov cut. 5 GeV pion and 5 GeV electron TB data.

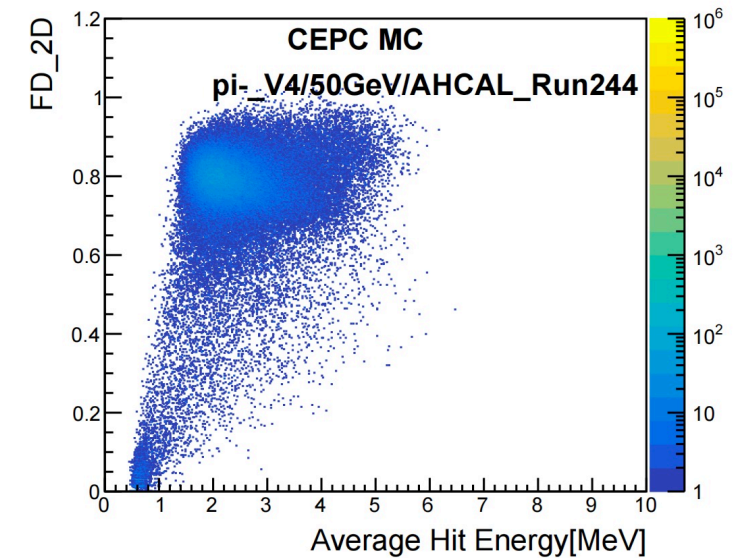
For data collected in 2023, SPS and PS

# Data pre-selection

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



**FD cut**

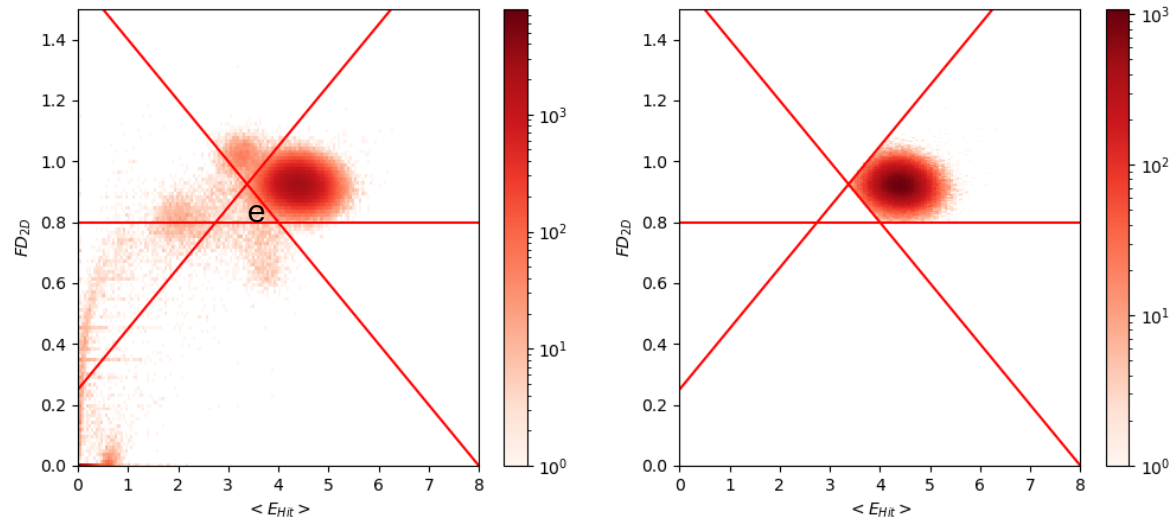


**MC**

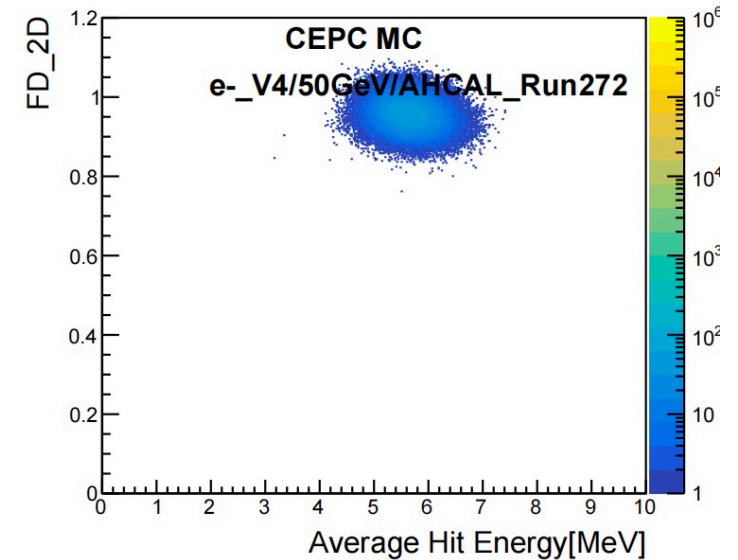
For data collected in 2023, SPS and PS

# Data pre-selection

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



**FD cut**

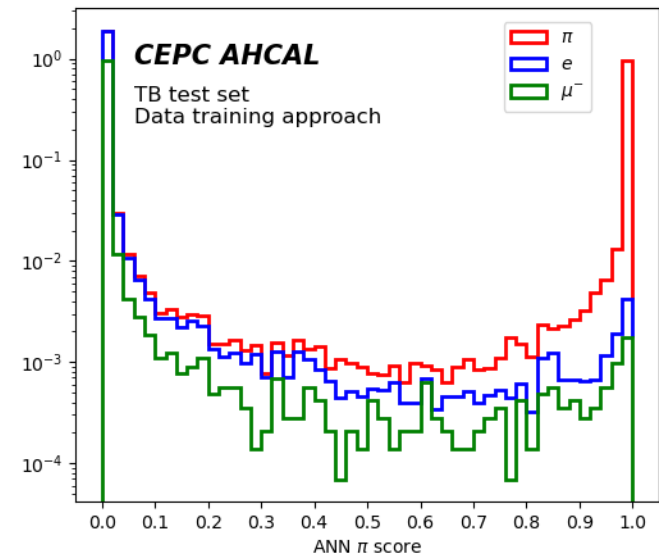
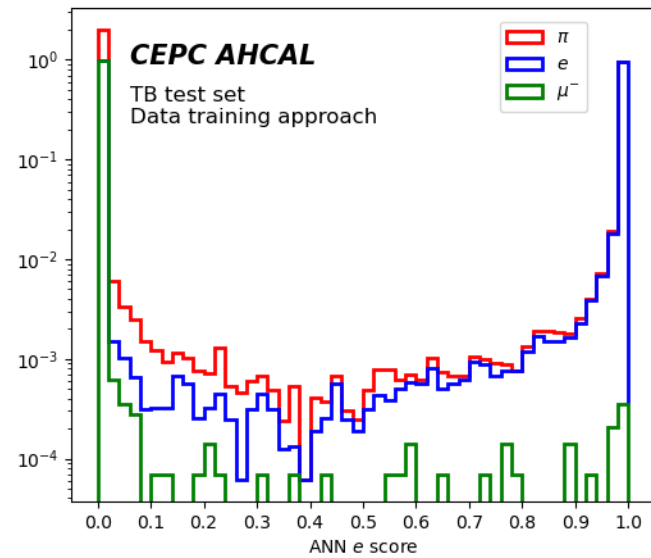
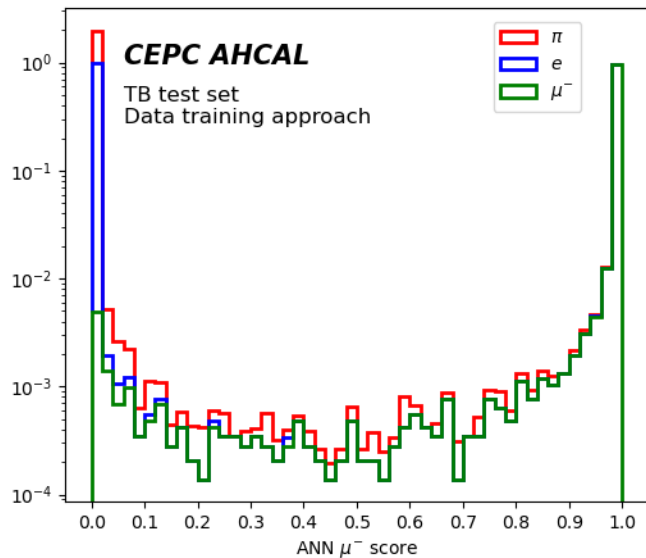


**MC**

For data collected in 2023, SPS and PS

# Evaluated on TB test set.

- Signals would get higher scores (closer to 1).
  - E.g. a pion would get higher ANN Pion score.
- An additional threshold ANN threshold cut would help to reject backgrounds.

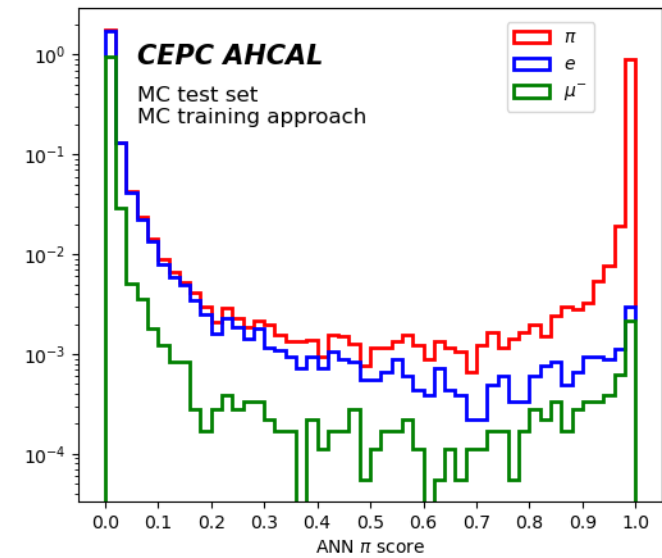
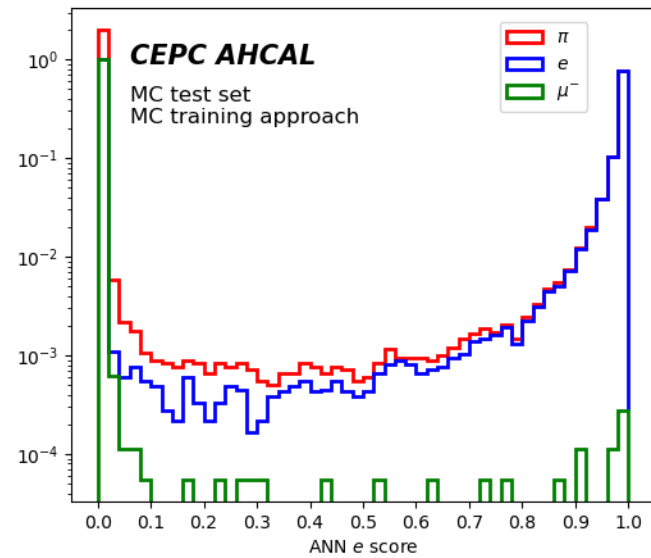
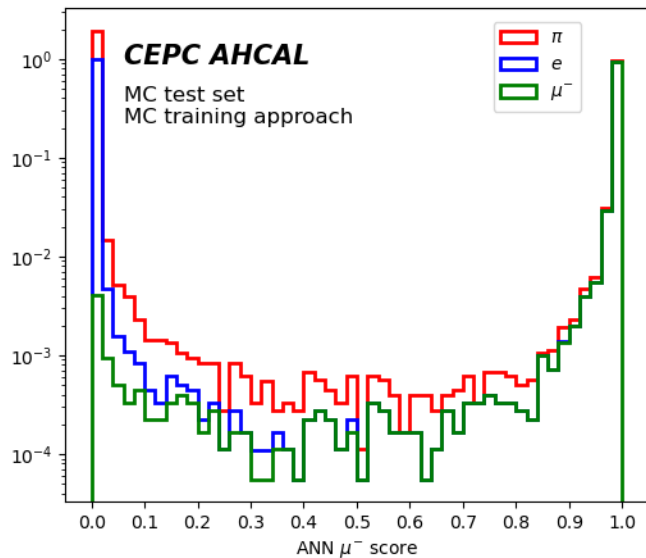


Output of ANN (data training approach).



# Evaluated on MC test set.

- Signals would get higher scores (closer to 1).
  - E.g. a pion would get higher ANN Pion score.
- An additional threshold ANN threshold cut would help to reject backgrounds.



Output of ANN (data training approach).

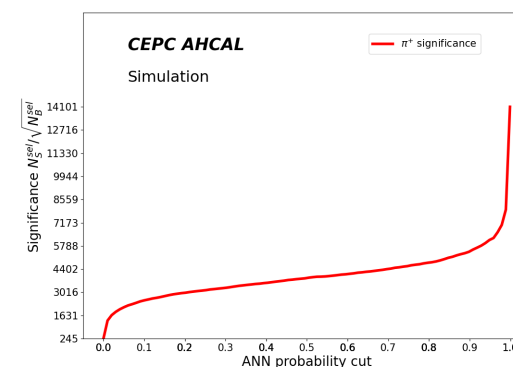
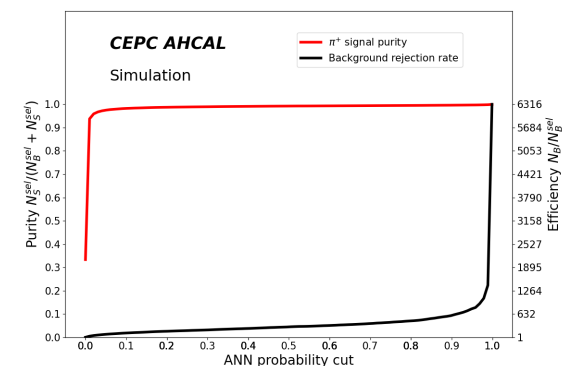
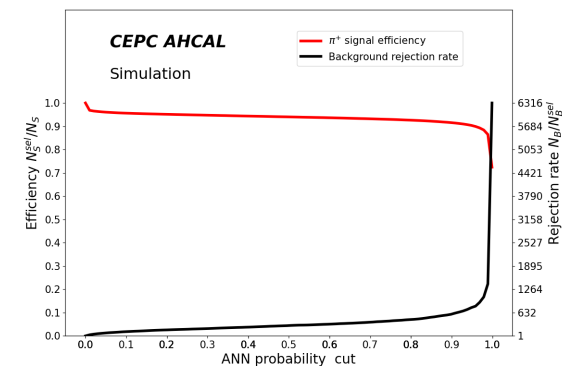
# Evaluated on MC test set

## Pion signal VS Backgrounds

Tested on Mixed MC Test Set.

- 120k Pi + 120k mu + 120k positron

ANN Score Cut	Pi Purity $(N_S^{sel} / (N_B^{sel} + N_S^{sel}))$	Pi Efficiency $(N_S^{sel} / N_S)$	Background Rejection Rate $(N_B / N_B^{sel})$
0.1	0.98	0.96	117
0.3	0.99	0.95	200
0.7	0.994	0.93	376
0.9	0.996	0.91	618



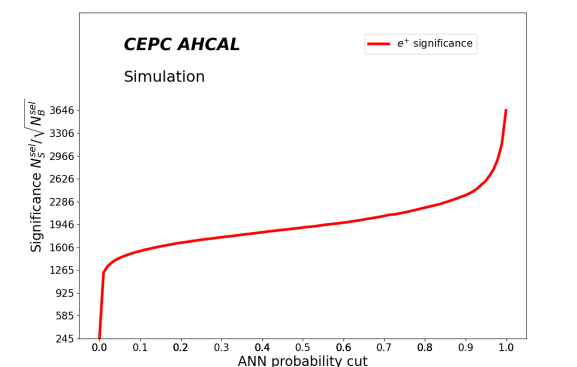
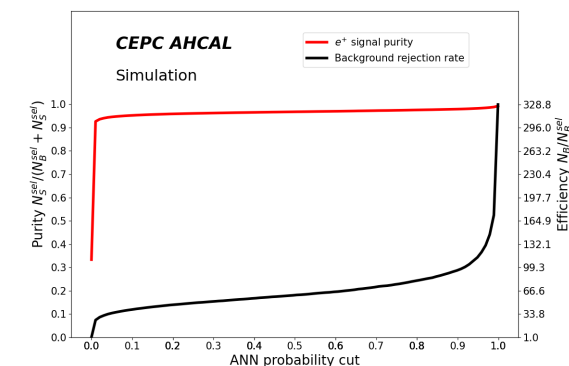
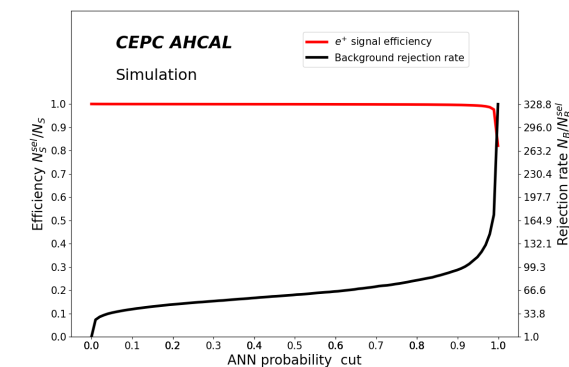
# Evaluated on MC test set

## Electron signal VS Backgrounds

Tested on Mixed MC Test Set.

- 120k Pi + 120k mu + 120k positron

ANN Score Cut	E Purity $(N_S^{sel} / (N_B^{sel} + N_S^{sel}))$	E Efficiency $(N_S^{sel} / N_S)$	Background Rejection Rate $(N_B / N_B^{sel})$
0.1	0.95	0.999	40
0.3	0.96	0.999	52
0.7	0.97	0.998	72
0.9	0.98	0.996	97



# Evaluated on MC test set

## Muon signal VS Backgrounds

Tested on Mixed MC Test Set.

- 120k Pi + 120k mu + 120k positron

ANN Score Cut	Mu Purity ( $N_S^{sel} / (N_B^{sel} + N_S^{sel})$ )	Mu Efficiency ( $N_S^{sel} / N_S$ )	Background Rejection Rate ( $N_B / N_B^{sel}$ )
0.1	0.97	0.99	59
0.3	0.97	0.99	69
0.7	0.98	0.99	80
0.9	0.98	0.98	86

