

Gradient Boosted Decision Tree for Particle Identification in the MPD

V. Papoyan^{1,3}

Coauthors: A. Aparin², A. Ayriyan^{1,3}, H. Grigorian^{1,3}, A. Korobitsin², A. Mudrokh²

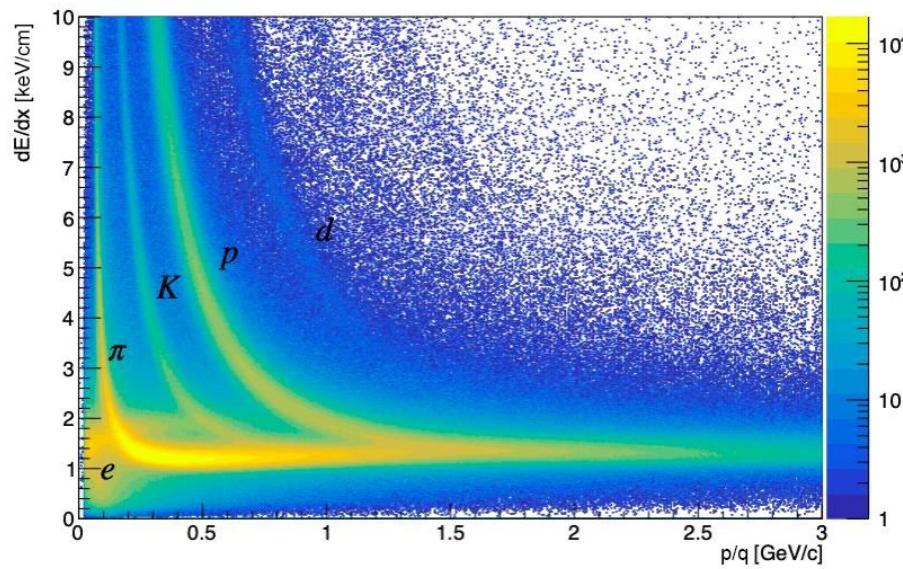
¹MLIT JINR, ²VBLHEP JINR, ³AANL (YerPhi)

This work was done with support from the Russian Science Foundation under Grant No. 22-72-10028

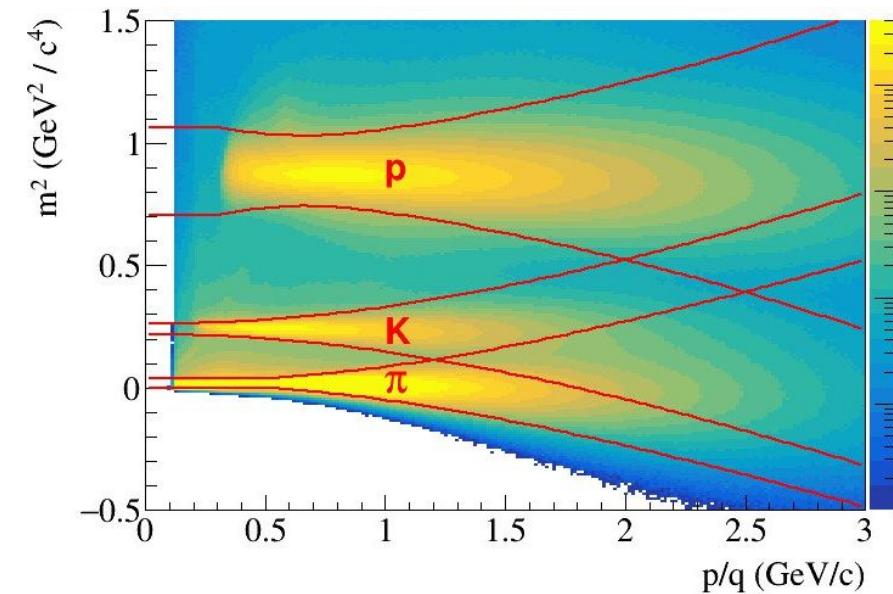
Particle Identification in MPD experiment

MPD particle identification (PID) based on **Time-Projection Chamber** (TPC) and **Time-of-Flight** (TOF).

A TPC can identify charged particles by measuring their specific ionization **energy losses** (dE/dx);



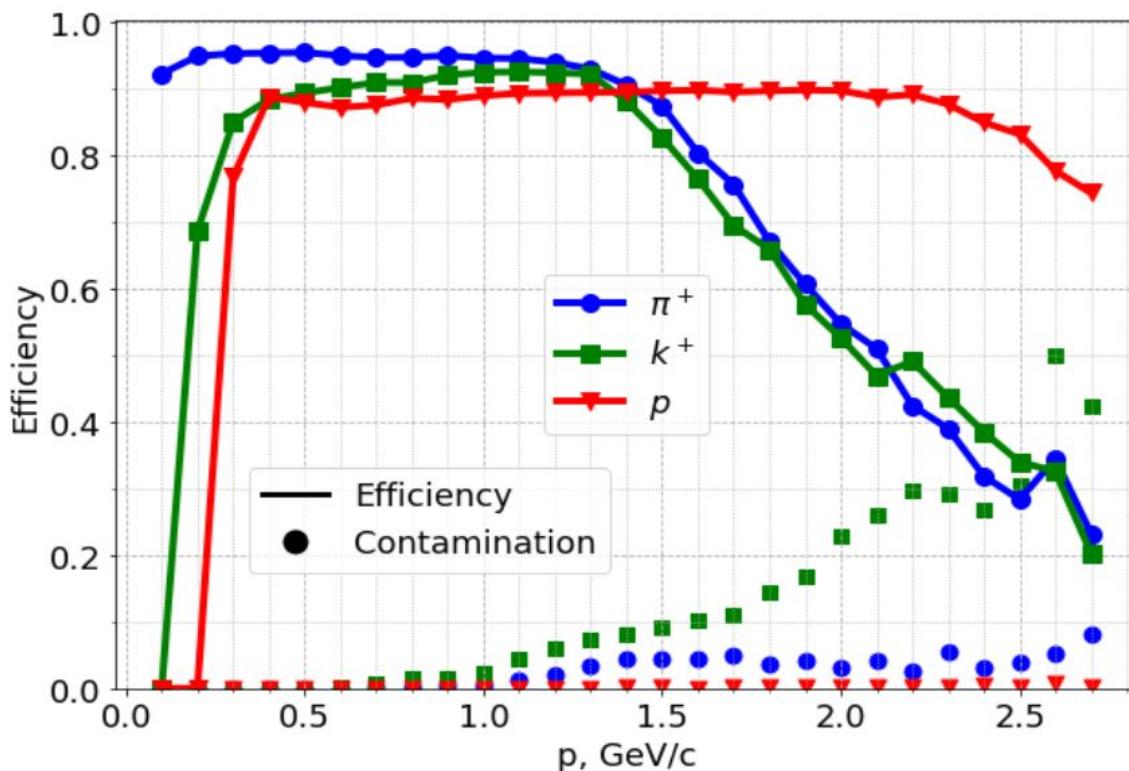
A TOF measures the particle flight **time** over a given **distance** along the track trajectory;



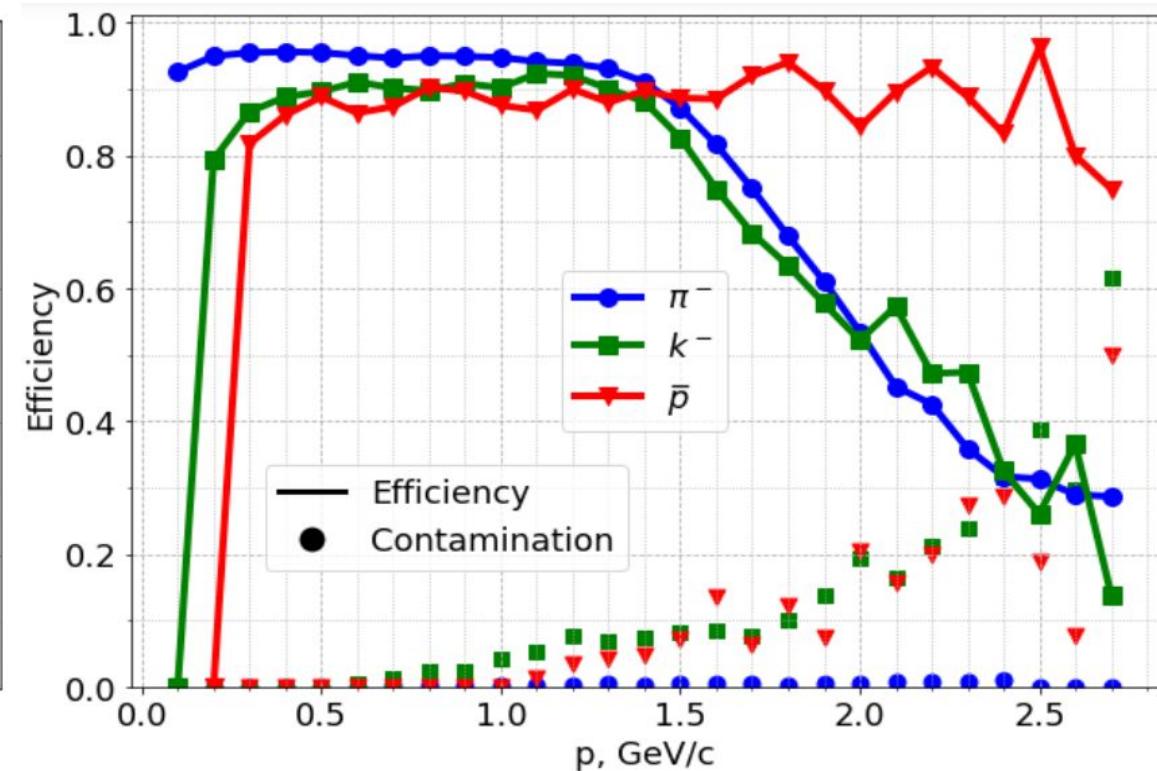
Knowing the particle **momentum** (from TPC) one obtains the **mass squared** and thus identity of the particle.

Baseline PID in MPD - N-sigma

$$E^s = \frac{N^s_{corr}}{N^s_{true}}$$



$$C^s = \frac{N^s_{incorr}}{N^s_{corr} + N^s_{incorr}}$$



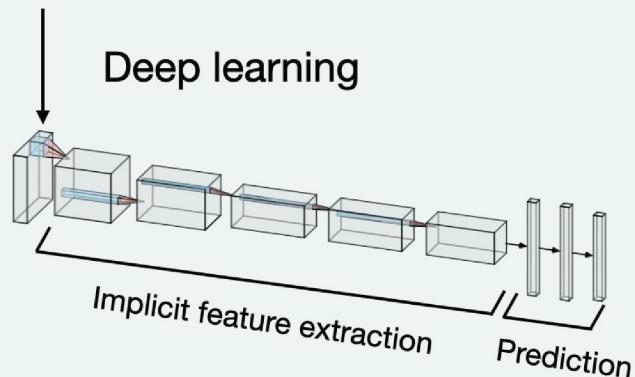
PID efficiency and contamination for positive (left) and negative (right) charged hadrons
in Bi+Bi collisions at 9.2 GeV

Tabular Data: Deep Learning vs Gradient Boosting

Unstructured data



Deep learning



Structured data

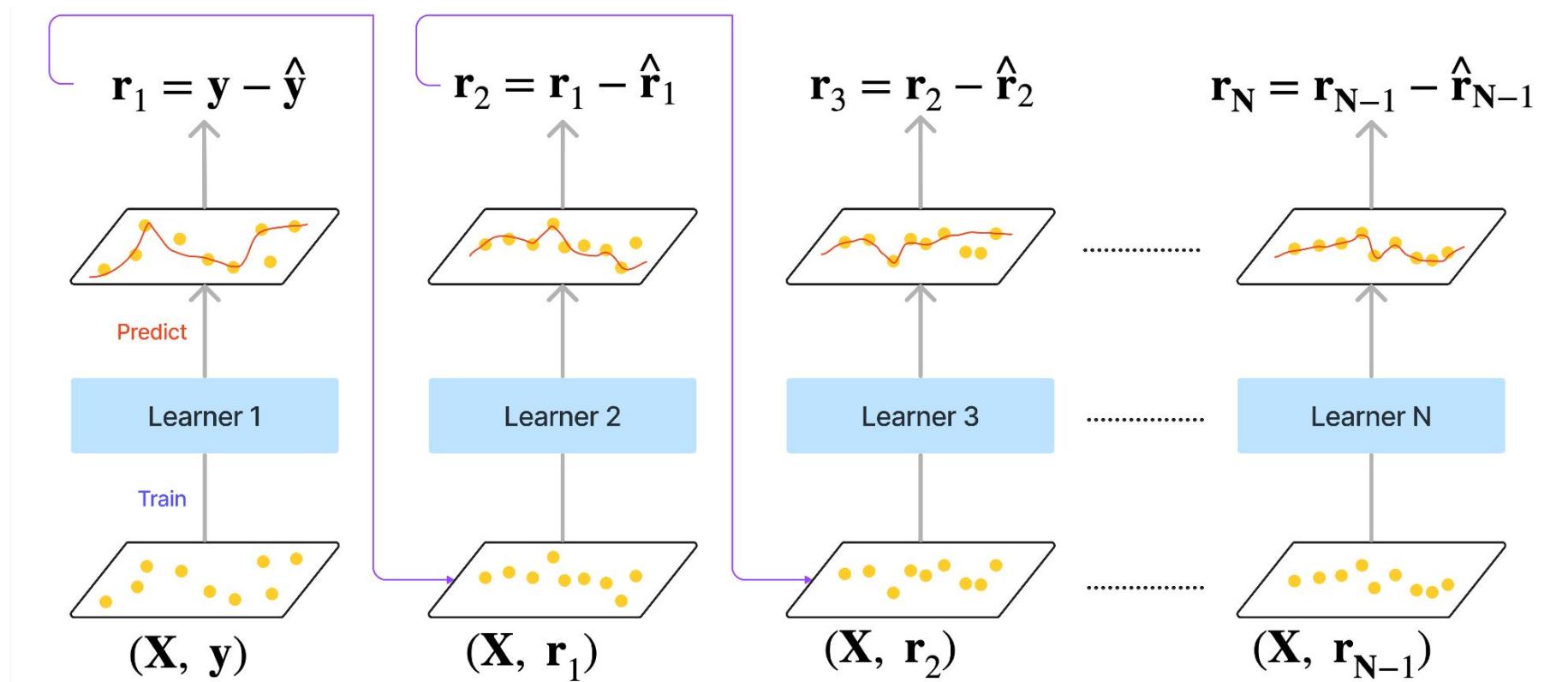
	Fuselage length	Wingspan
Boeing 707	44,07	39,9
Cessna 172	8,28	11
B-2 Spirit	20,90	52,12

Gradient
Boosting



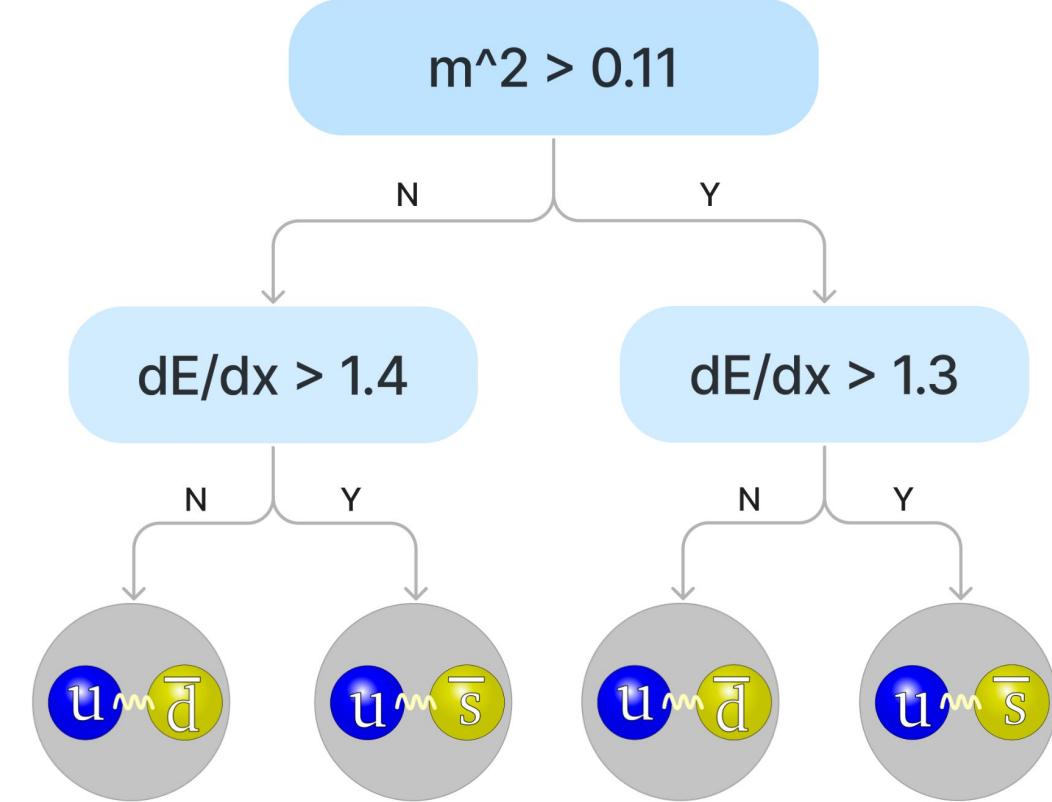
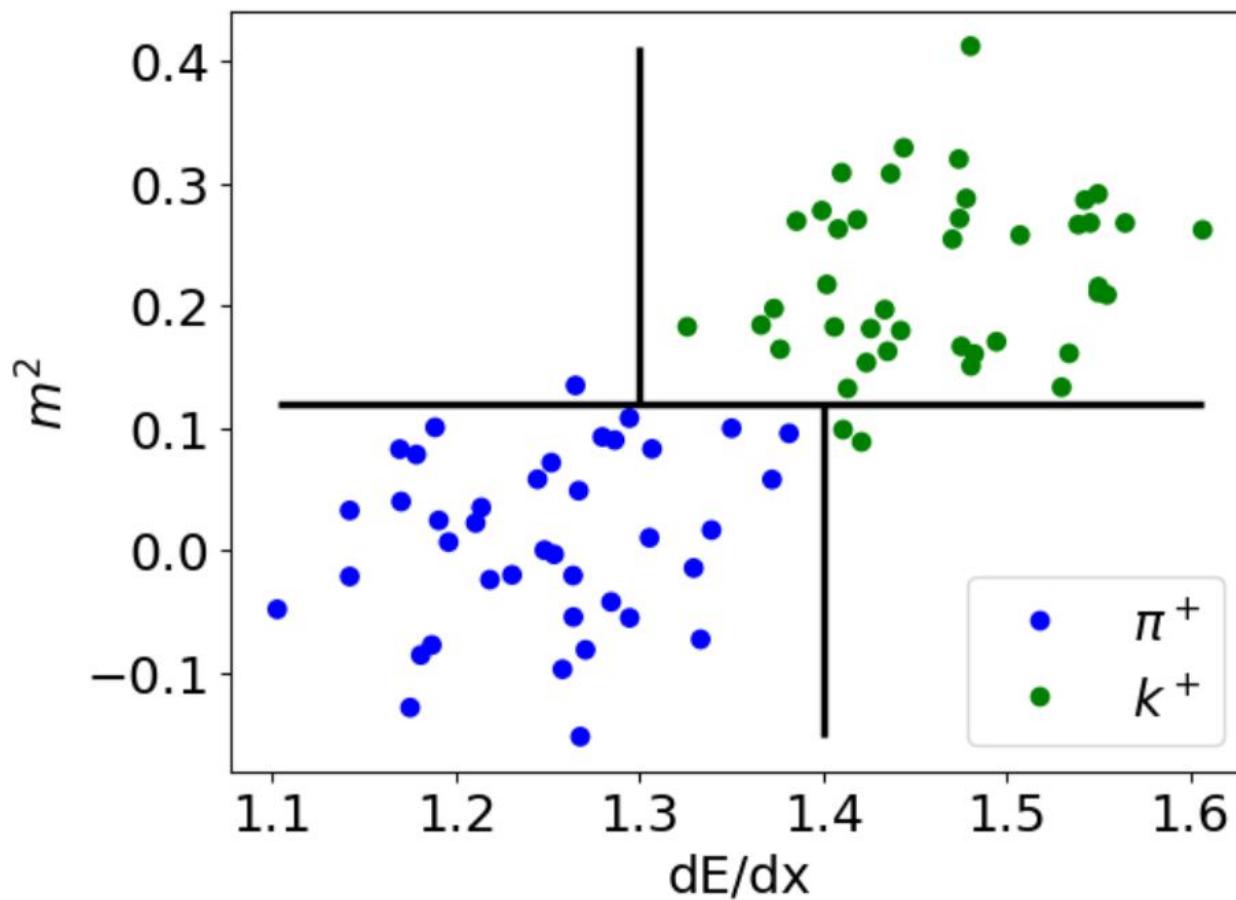
Gradient Boosting

Gradient boosting is a machine learning technique which combines weak learners into a single strong learner in an iterative fashion



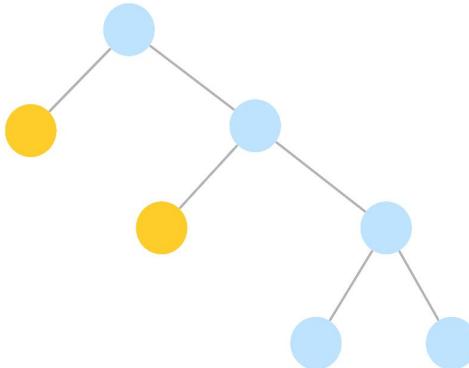
Gradient Boosted Decision Tree

Gradient Boosted Decision Tree (GBDT) uses decision trees as weak learner. They can be considered as automated multilevel **cut-based** analysis

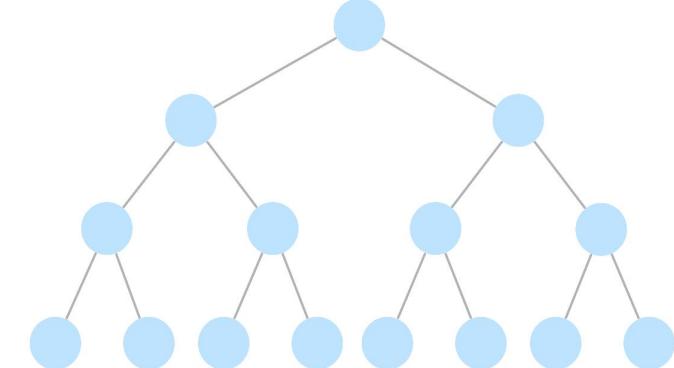


XGBoost vs LightGBM vs CatBoost vs SketchBoost

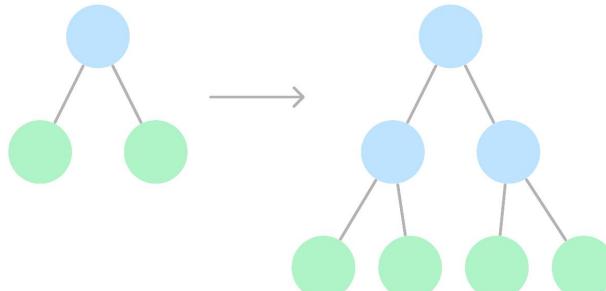
Asymmetric Tree (XGB, LGBM)



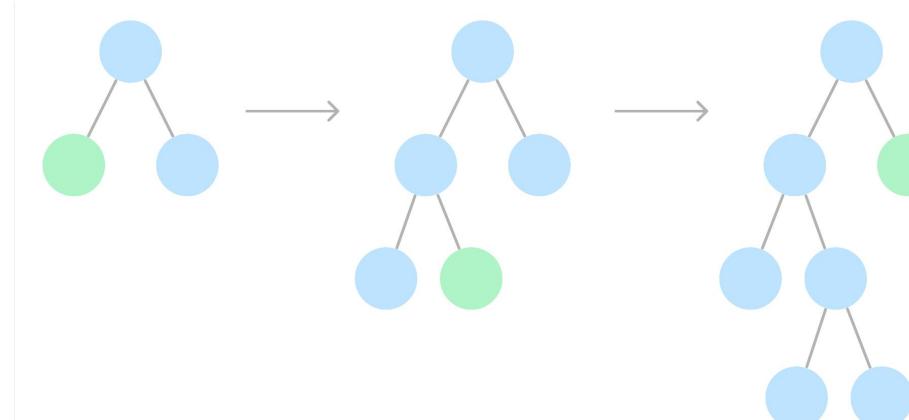
Symmetric Tree (CatBoost, SketchBoost)



Level-wise Tree Growth (XGB)



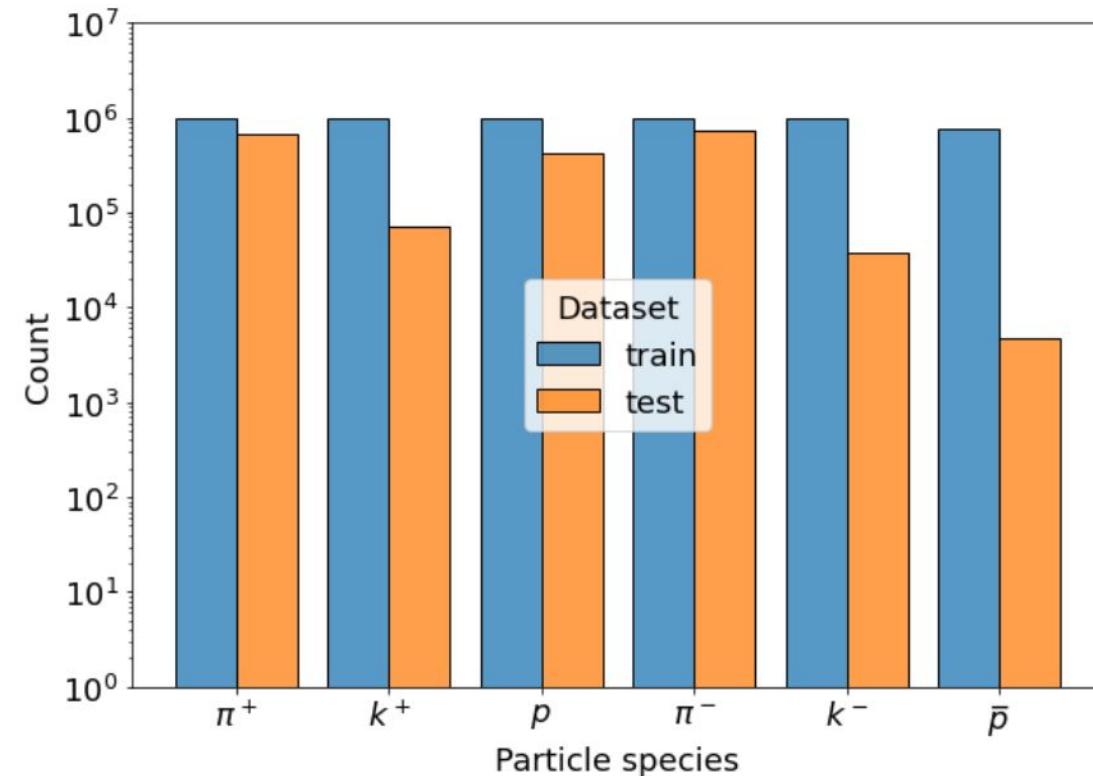
Leaf-wise Tree Growth (LGBM)



Datasets

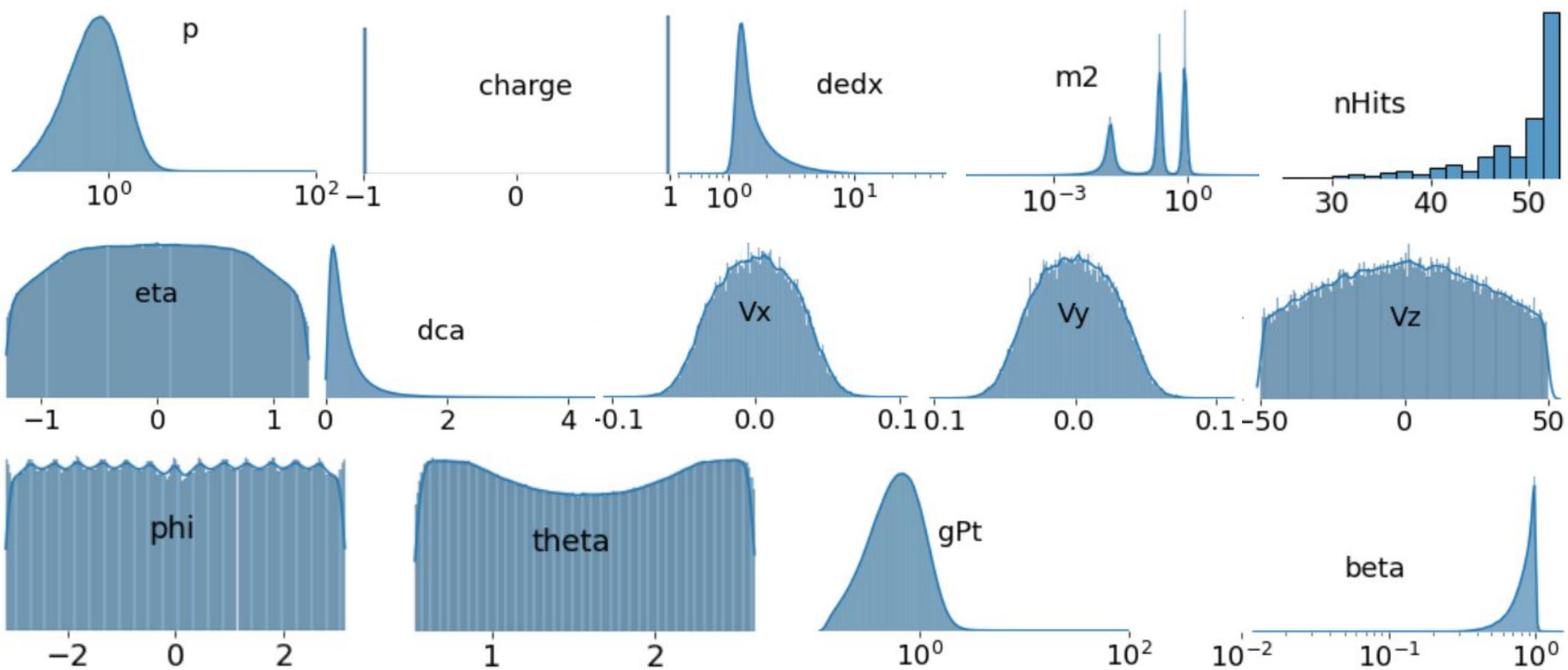
Subsamples of the two MPD Monte-Carlo productions have been used (Request 25 & Request 29)

	prod05	prod06
Event generator	UrQMD	PHQMD
Transport	Geant 4	Geant 4
Impact parameter ranges	0-16 fm (mb)	0-12 fm
Smear Vertex XY	0.1 cm	0.1 cm
Smear Vertex Z	50 cm	50 cm
Colliding system	Bi+Bi	Bi+Bi
Energy	9.2 GeV	9.2 GeV

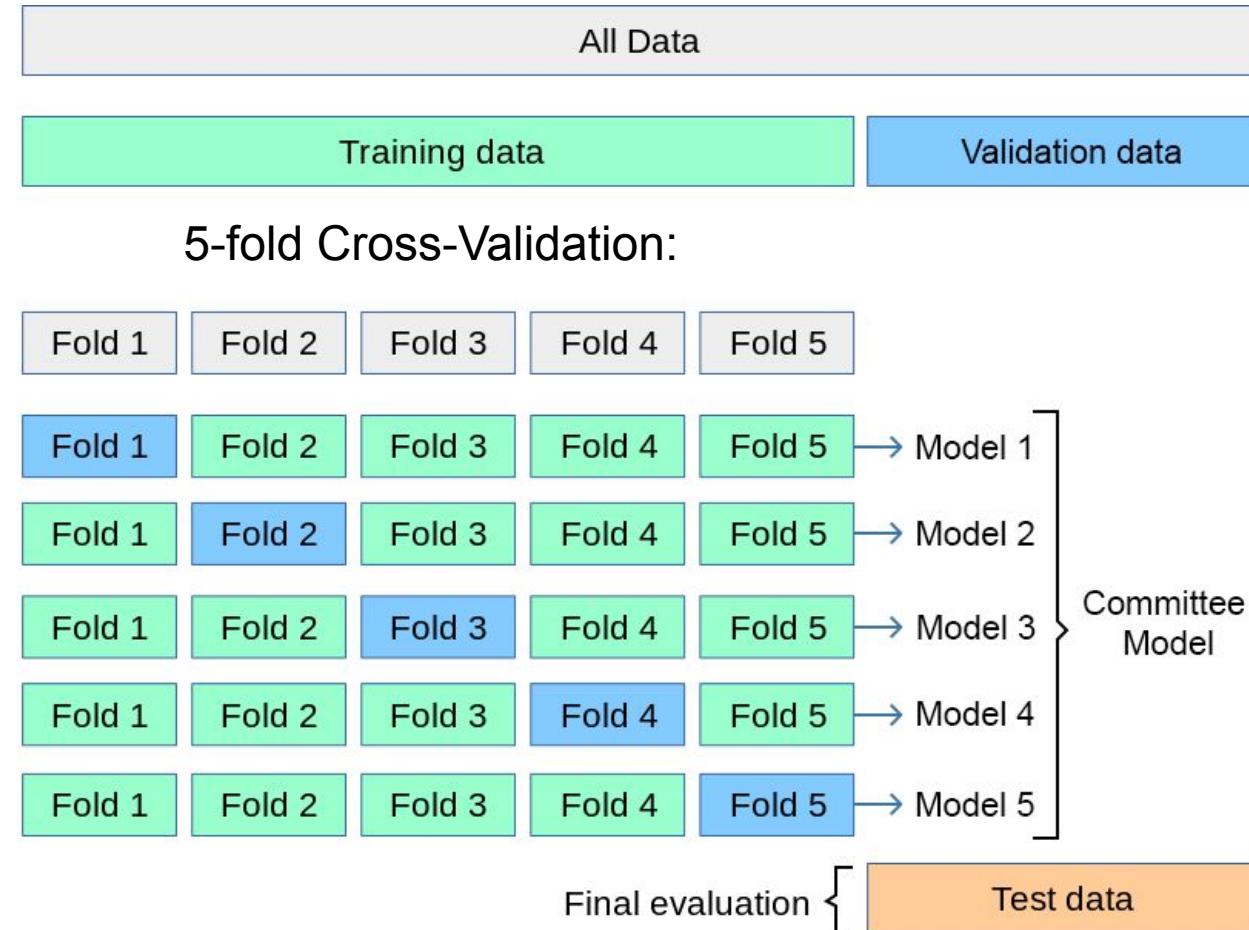


track selection criteria: ($p < 100$) & ($|m^2| < 100$) & ($n\text{Hits} > 15$) & ($|\eta| < 1.5$) & ($dca < 5$) & ($|\nabla z| < 100$)

Data description



Experiment design



All classifiers have been trained using the Nvidia Tesla V100-SXM2 NVLink 32GB HBM2 within the ecosystem for tasks of machine learning, deep learning, and data analysis at **HybriLIT** platform

Two stages of the experiments

Some parameters for the tuning and model evaluation stages

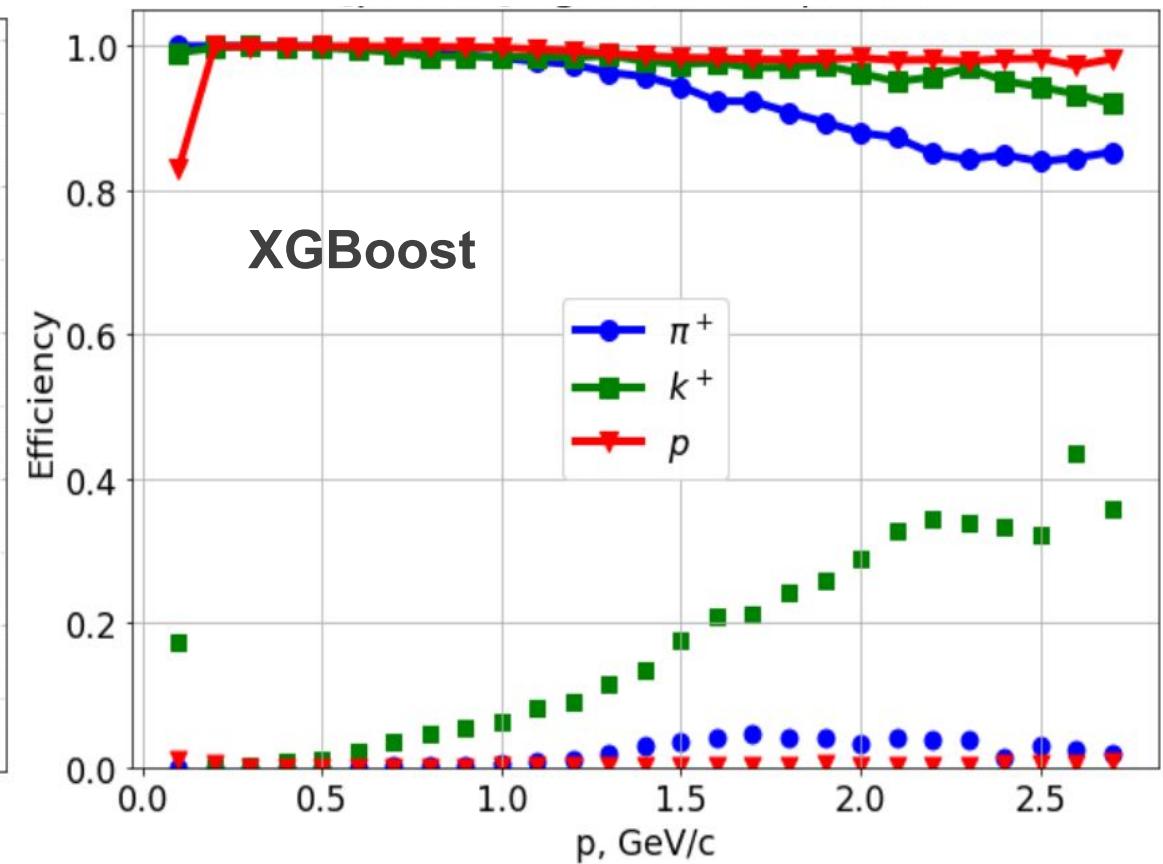
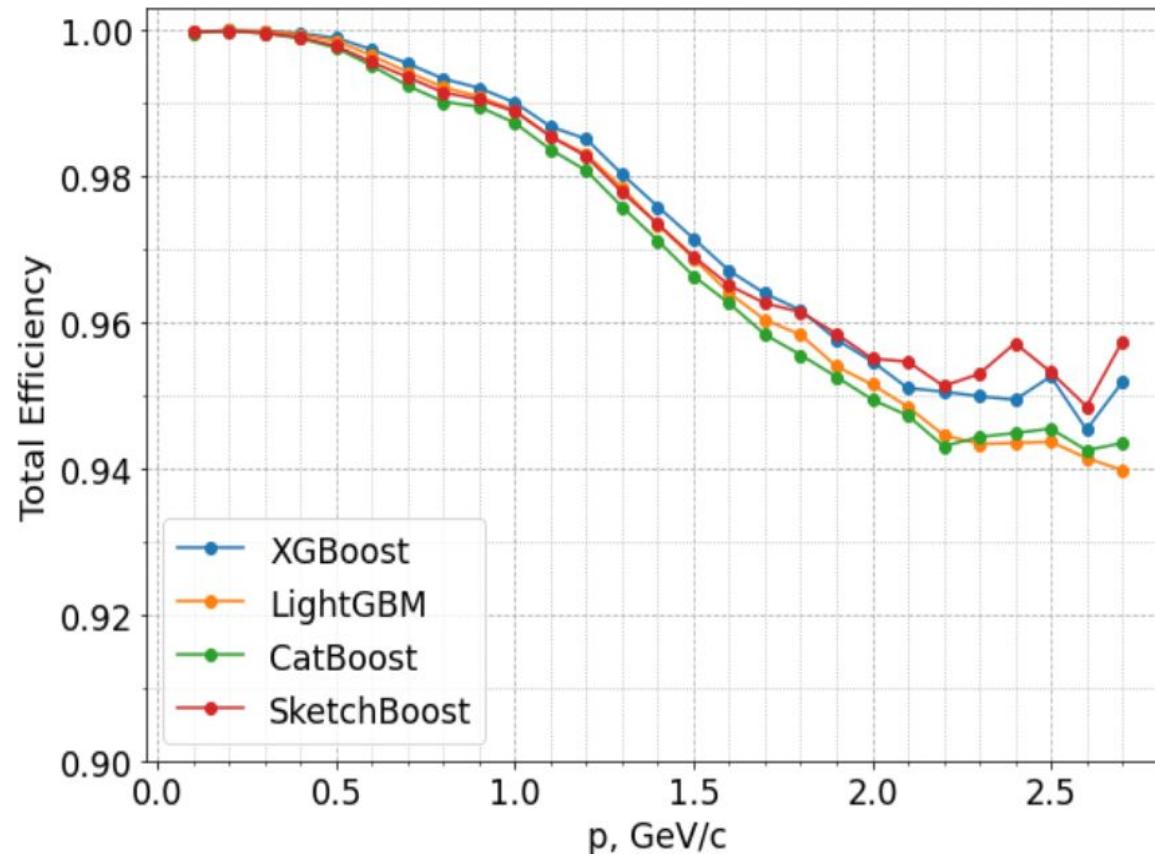
Stage	Learning Rate	Max Number of Iterations	Early Stopping
Tuning	0.05	5 000	200
Model Evaluation	0.015	20 000	500

Results for hyperparameter tuning (after **30 iterations** of the TPE algorithm for each GBDT)

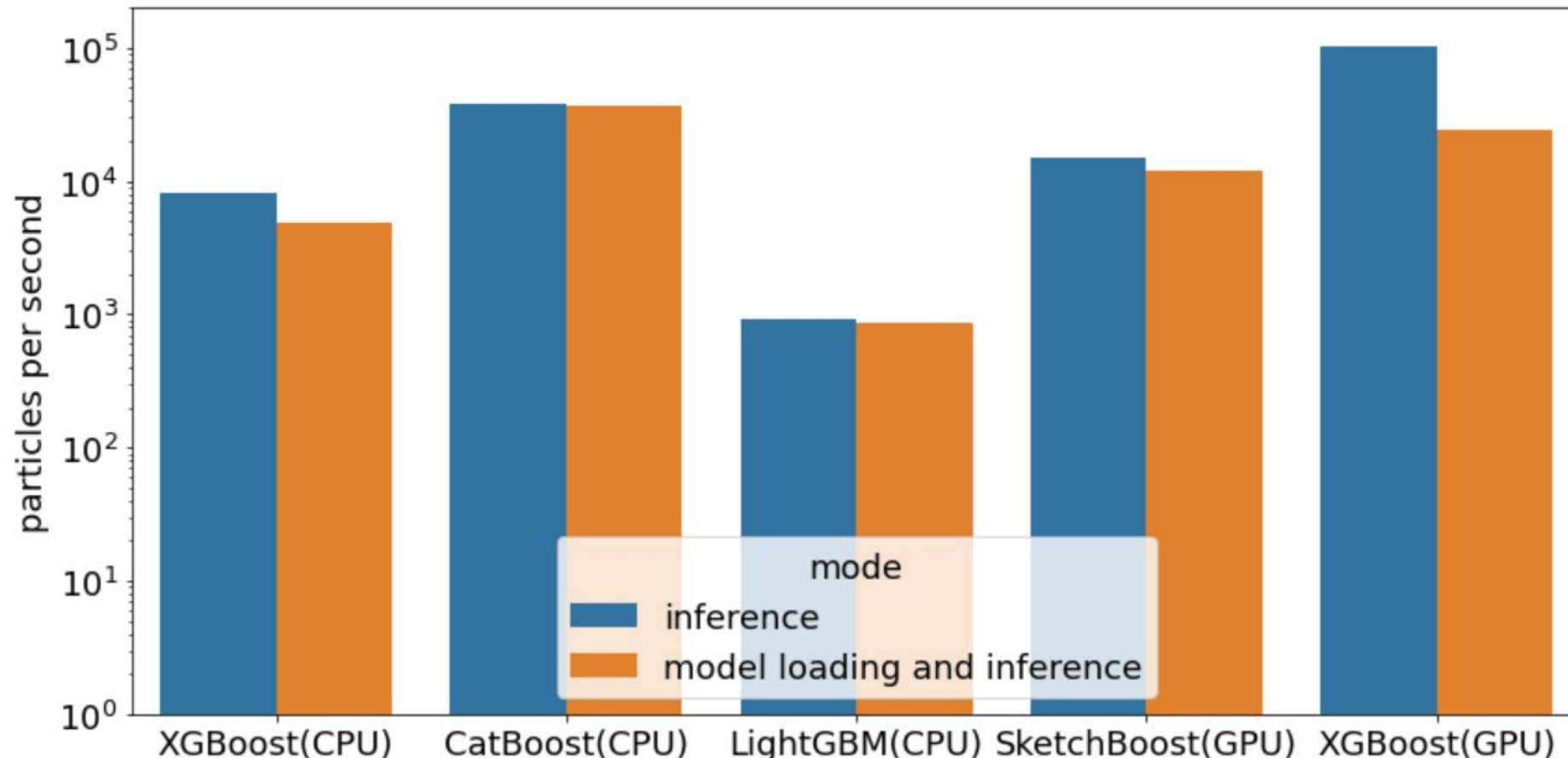
Framework	Max. Depth	L2 leaf reg.	Min. data in leaf size	Rows sampling rate
XGBoost	8	2.3	0.00234	0.942
LightGBM	12	0.1	4	0.981
CatBoost	8	3.0	5	0.99
SketchBoost	8	3.0	5	0.99

Comparative analysis of the algorithms. Efficiency

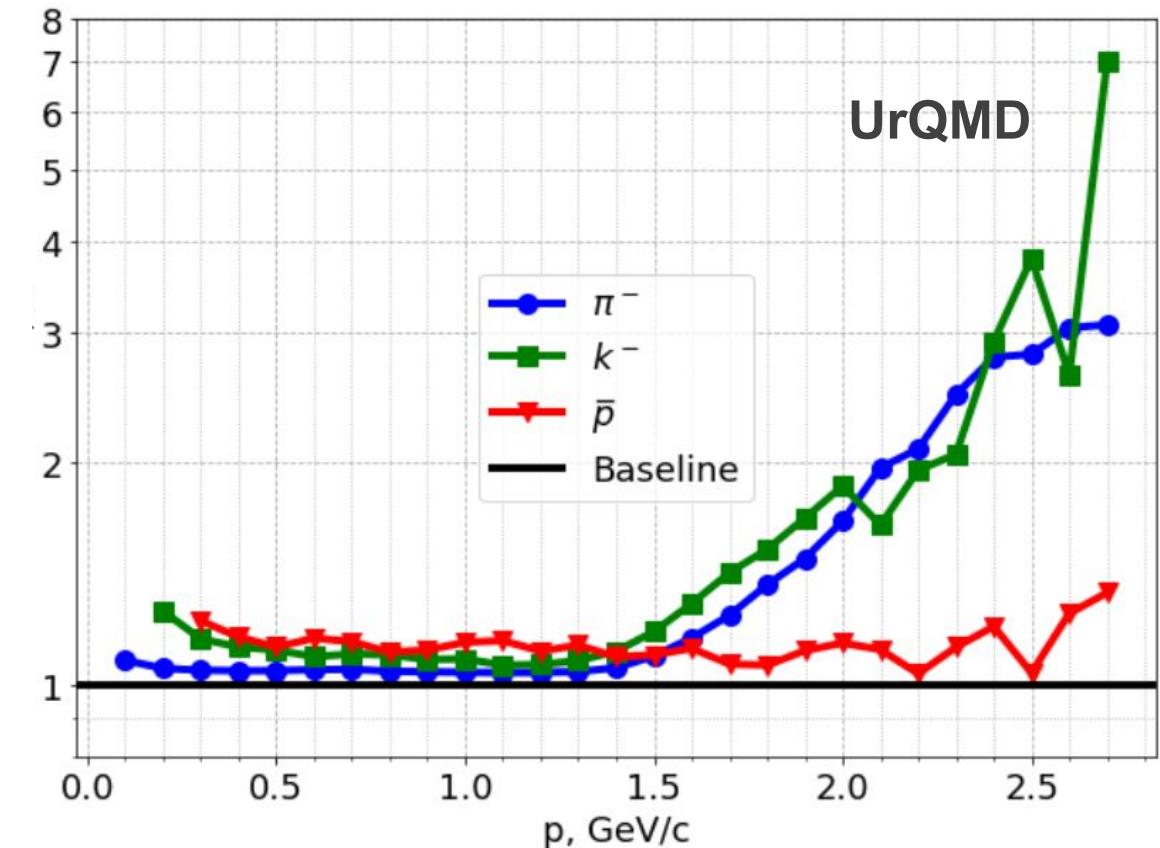
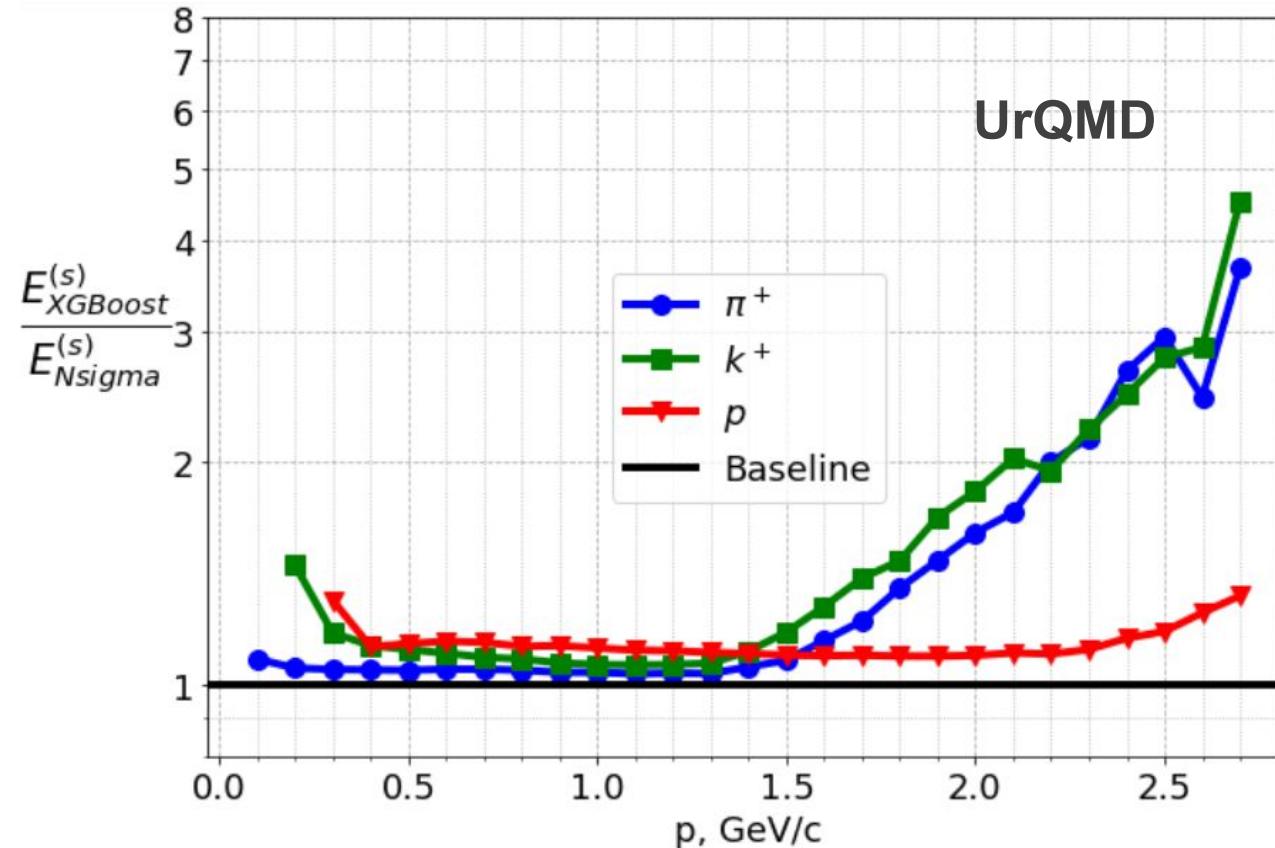
	XGBoost	LightGBM	CatBoost	SketchBoost
Total Efficiency	0.99327	0.99235	0.99138	0.99239



Comparative analysis of the algorithms. Inference time

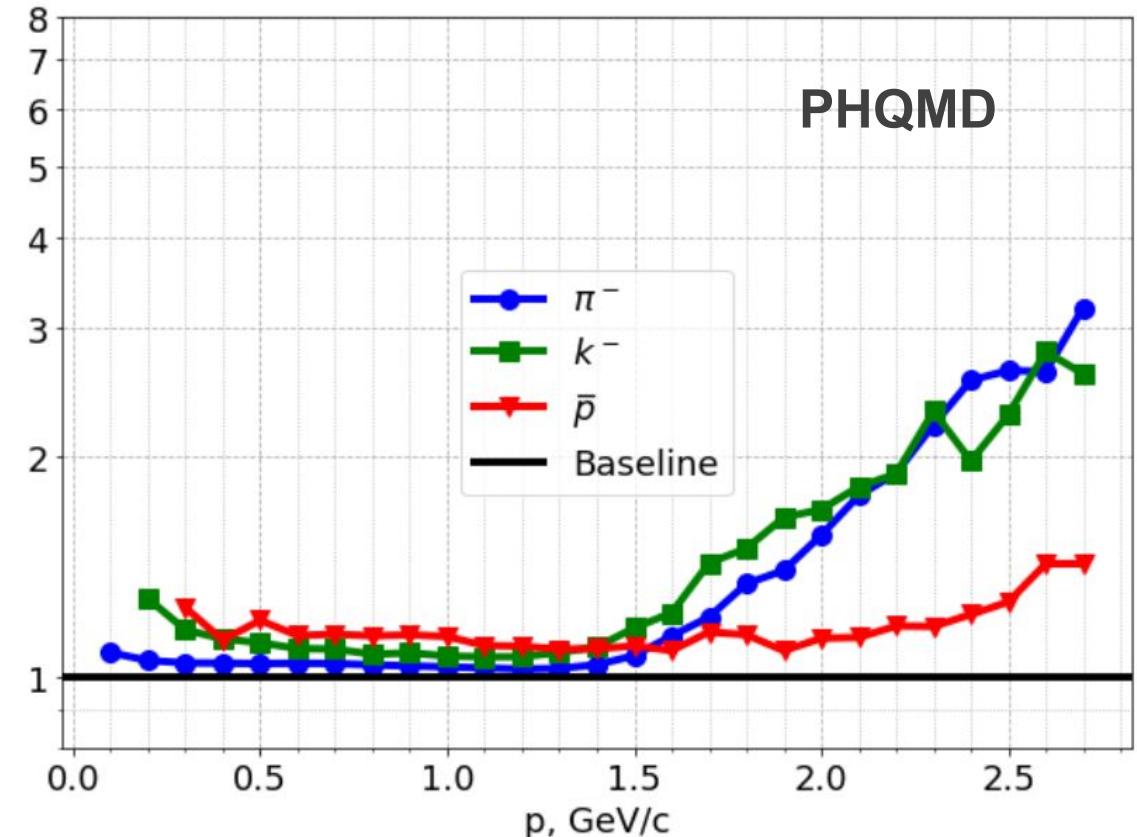
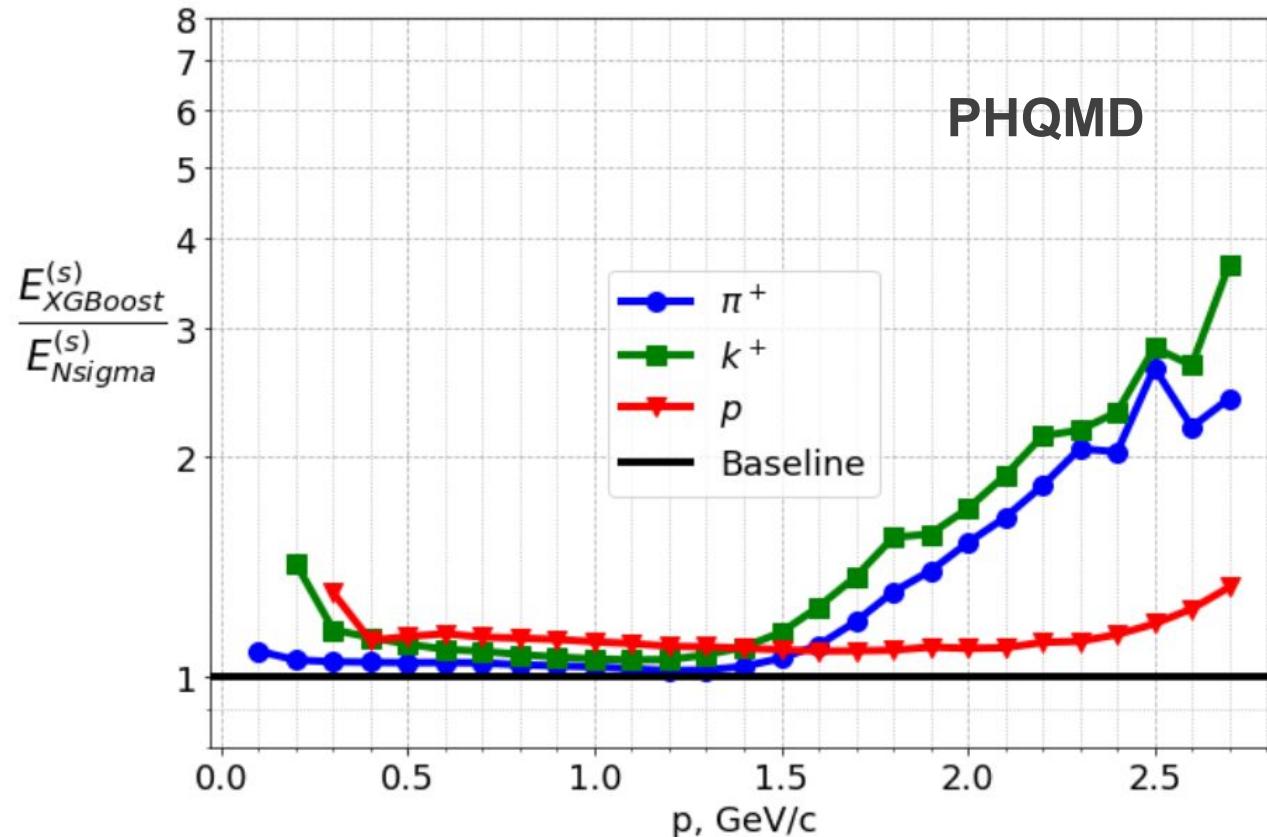


Comparison with N-sigma



Efficiency ratio of XGBoost and n-sigma method

Comparison with N-sigma



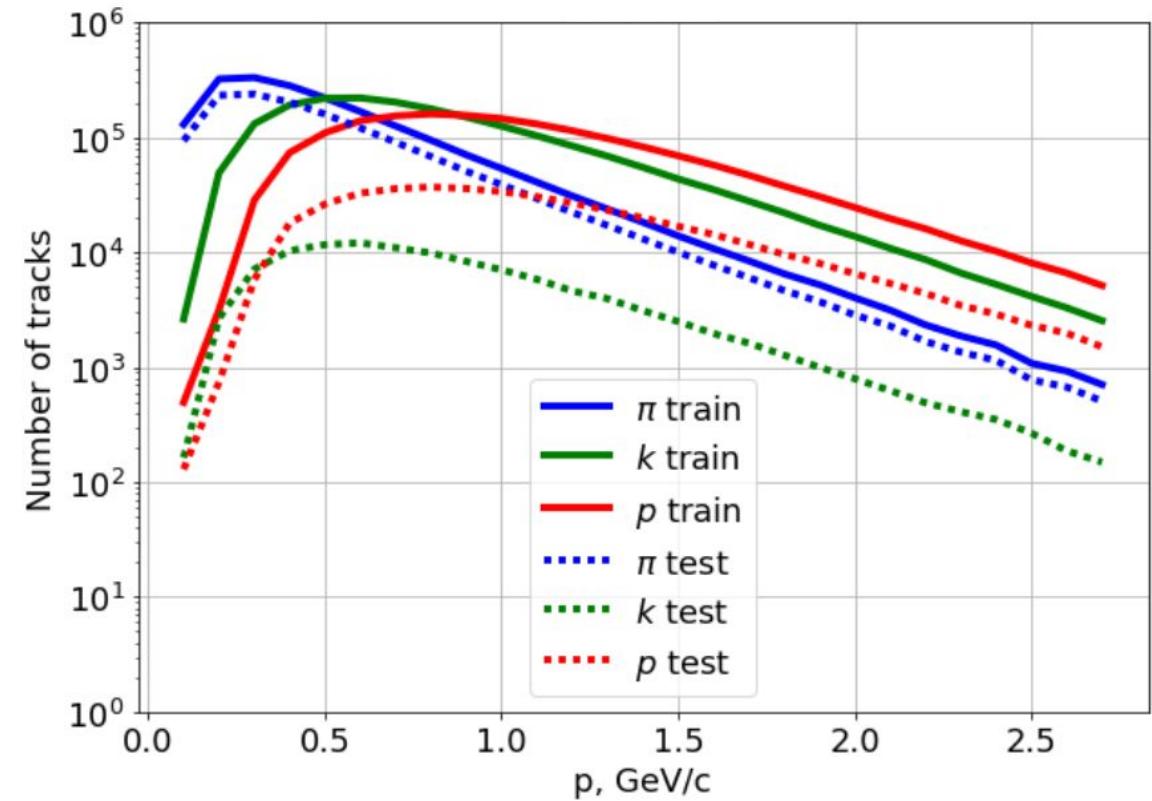
Efficiency ratio of XGBoost and n-sigma method

Conclusion and Outlooks

In general XGBoost has been demonstrated highest PID efficiency in comparison with considered algorithms of GBDT.

Next we are going to do additional testing to characterize identification stability of the model on data produced with different initial parameters of generated MC tracks at the MPD detector;

Also we are going to analyse the nature of the misclassifications and investigate the class imbalance problem.



Backup

Formulas

$$m^2 = \frac{p^2}{c^2} \left[\frac{t^2 c^2}{L^2} - 1 \right] \quad \beta = \frac{L}{ct}$$

$$-\left(\frac{dT}{dx}\right) = \frac{4\pi n_e z^2 e^4}{m_e v^2} \left[\ln \frac{2m_e v^2}{I} - \ln(1 - \beta^2) - \beta^2 - \delta - U \right],$$

Classification of Charged Particles

In Machine Learning terms PID can be considered as **classification** task (**Supervised** learning).

Let

X - is the input space (particle characteristics such as: dE/dx , m^2 , β , q , etc)

Y - is the output space (particle species such as: π , k , p , etc)

Unknown mapping exists

$$m : X \rightarrow Y,$$

for values which known only on objects from the finite training set

$$X^n = (x_1, y_1), \dots, (x_n, y_n),$$

Goal is to find an algorithm **a** that classifies an arbitrary new object $x \in X$

$$a : X \rightarrow Y.$$

Data description

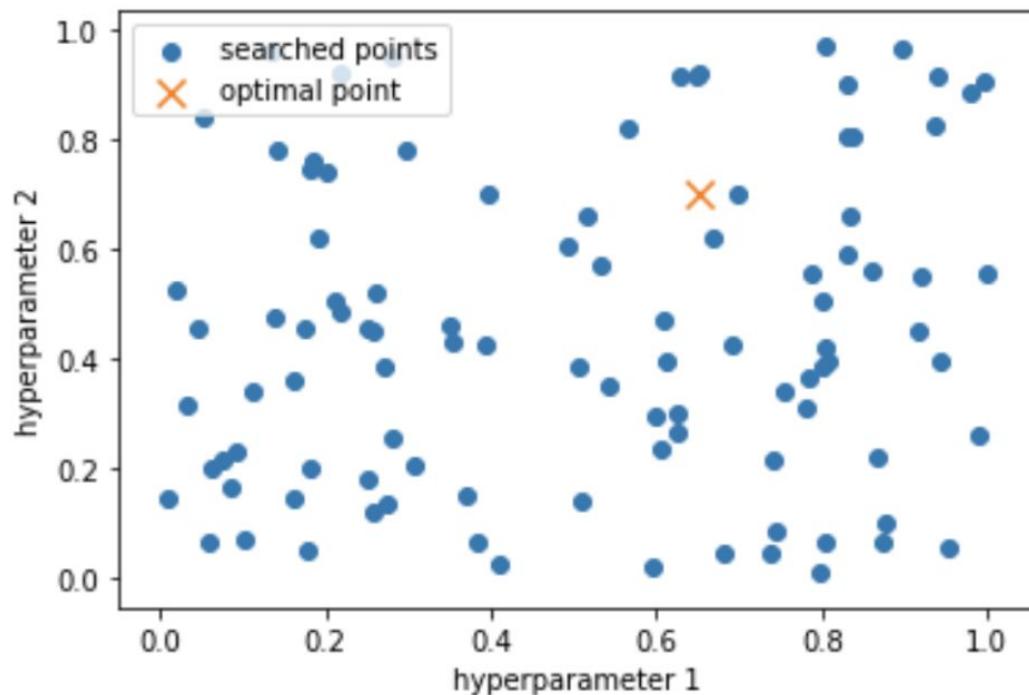
feature	values range
p	(0.1, 100)
q	{-1, 1}
dedx	(0, 72)
m2	(-100, 100)
nHits	[20, 53]
eta	[-1.3, 1.3]
dca	(0, 5)

feature	values range
Vx	(-0.106, 0.106)
Vy	(-0.103, 0.112)
Vz	(-50, 54.1)
phi	(-3.1415, 3.1415)
theta	(0.53, 2.61)
gPt	(0.106, 98)
beta	[0.012, 1.564]

Hyperparameters tuning

Tree-structured Parzen Estimator (TPE) was used to find the optimal hyperparameters;
TPE is a form of Bayesian Optimization.

Random search



TPE search

