

Machine learning based optimizations of the ALICE TPC dE/dx response

HADRON Conference 2023 Tuba Gündem ALICE Collaboration 05.06.2023

Outline



- ALICE TPC
- Specific Energy Loss (dE/dx) Measurement
- Particle Identification
- Machine Learning

ALICE TPC





- The primary particle tracking and identification device in ALICE
 - Inner and outer field cages, high-voltage central electrode, endplates (18 sectors per side)

ALICE TPC particle track outer field cage readout chambers





inner field cage

• A charged particle loses energy by ionizing the gas

5 m

endplate





• The drifting electrons are projected onto readout planes





- The inner and outer readout chambers are denoted by IROC and OROCs, respectively
 - Different pad lengths and different numbers of pad rows

ALICE TPC projected particle track Pad plane central electrode outer field cage OROC3 OROC2 OROC1 IROC 5 m endplate GEM stack inner field cage

ALICE

• For the gas gain, GEM stacks are used



- GEM transparency affects the energy loss resolution
- The amplified charge then induces a signal on the pad plane
- The total charge of the signal is proportional to the specific energy loss of this particle



- The fluctuations of the signal are described by the Landau distribution
- Clustering in pad and time direction is performed
- The specific energy loss (dE/dx) is calculated using a truncated mean method on the cluster charge (Q_{tot} or Q_{max})
- The truncation value is 60% in the ALICE TPC

Particle Identification (PID)





• The average behaviour of the dE/dx of a charged particle is described by the Bethe-Bloch formula

• ALEPH parametrization:

 $\langle dE/dx \rangle_{ALEPH} = p1 (p2 - log(p3 + (\beta\gamma)^{-p5})/\beta^{p4} - 1) z^{fz} f_{MIP}$

• Momentum:

 $p = \beta \gamma . m$

Particle Identification (PID)

N σ -bands are defined around the average for PID •

d*E*/dx (arb. units)

N^{TPC}(e)





Optimizing the PID Performance

ALICE

- The Bethe-Bloch parameterization is not know a priori
 - A fit to data is required
- The Bethe-Bloch fit is not sufficient to describe all particles species in full phase space
 - A multidimensional correction is needed
- Resolution (sigma) parametrization requires a multidimensional description
- The dE/dx estimator might be improved by optimizing the truncation range separately on readout regions
- Different Machine Learning (ML) techniques can be used to carry out these tasks



Machine Learning

PID Response Parameterization

- 1. Update the Bethe-Bloch parameters
 - a. using Optuna (full dataset, no skimming)
 - b. using ROOT (skimmed clean samples available)
- 2. Skim data with optimal Bethe-Bloch parameters
- 3. Train Neural Networks (NN) for mean and sigma estimation
- 4. Mean+Sigma NN (relatively small, to optimize CPU usage) is stored
- 5. Application on analysis with the Open Neural Network Exchange (ONNX)

Full chain (1-5) is run per data taking period (several weeks of data taking)



Data Preparation

- Clean samples of pions, protons and electrons
 - \circ V0 selection based on
 - invariant mass
 - pointing angle (θ_{PA})
 - Armenteros-Podolanski selections
 - excluding overlapping regions





05.06.2023 | HADRON Conference 2023

Tuba Gündem | Goethe Universität

Data Preparation



- Skimming
 - Random downsampling of data using Tsalis/Hagedorn fit along the transverse momentum p_{T}

- In order to reduce computational complexity skimmed data is used for the PID response parameterization
 - PetaBytes -> GigaBytes



Bethe-Bloch Parameterization with Optuna

- Hyperparameter Optimization (HPO) framework Optuna
 - Designed for the automation and the acceleration of the optimization studies
 - How does it work?
 - Define a function to be minimized
 - Error or score functions
 - Define number of trials for the optimization
 - In each trial Optuna will generate the hyperparameters in user defined ranges
 - Optuna will give the best parameters depending on the error or score values
- The parameters of the Bethe-Bloch fit are considered as hyperparameters





Bethe-Bloch Parameterization with Optuna

- Full dataset and range for Bethe-Bloch parameters are given to Optuna •
- For each set of parameters
 - PID is assigned based on proximity to the nearest Bethe-Bloch curve (black line) Ο
 - For each particle species, the mean or density maxima (grey points) is evaluated Ο along the $\beta \gamma$ binning
 - The score function is calculated Ο
 - Mean squared sum of mean or max values to current Bethe-Bloch curve
- Optuna tries to optimize this score
- This runs on CPU (several CPU hours)



 10^{-1}



100

p [GeV/c]



Bethe-Bloch fit : p1 (p2 - log(p3 + $(\beta\gamma)^{-p5})/\beta^{p4}$ - 1)

101

Bethe-Bloch Parameterization Performance





- dE/dx_{exp}: Bethe-Bloch fit from ROOT (clean samples) or HPO (no clean samples)
- Fits are in good agreement with both mean and max estimations
- Residual deviations at low momenta
 - Further correction is needed

Mean Correction to Bethe-Bloch Parameterization

- Mean correction using Neural Networks in 6D phase space
- Network training is done on GPU using Pytorch
- Fully connected, 6 input layers, 3 hidden layer, 8 neurons per layer
- Loss function: Mean Squared Error (MSE)
- Activation function: tanh
- $dE/dx_{corr} = dE/dx_{exp} \cdot NN_{mean}$

NN learns the ratio



Deep neural network Input layer Multiple hidden layers Output layer m : features a_{1n} $a_{11} a_{12} a_{13}$ momentum a_{2n} inclination angle a ... a₂₂ a₂₃ ... q/p_{T} a_{3n} a ... a₃₂ a₃₃ mass multiplicity number of clusters a_{mn} a_{m1} $a_{m2} a_{m3}$. . .

n : observation

Tuba Gündem | Goethe Universität

AITCE

Resolution Parameterization



- The data distribution for each particle species can be well approximated by a Gaussian distribution in dE/dx for every slice in momentum
- NN can learn the mean description of the data
 - For the resolution (sigma) parameterization the data is transformed in a way that the mean description is is proportional to sigma of dE/dx



blue curve: original distribution

red curve: transformed distribution

orange line: NN estimation which is proportional to original sigma

 $\mu = \operatorname{sqrt}(2/\pi) \sigma$

Neural Network Corrections Performance

- NN corrections via ONNX for the analysis code for mean and sigma
- Comparison between No plots



ALICE

05.06.2023 | HADRON Conference 2023

Tuba Gündem | Goethe Universität

05.06.2023 | HADRON Conference 2023

Optimization of *dE***/***dx* **Estimator**

- Toy Monte Carlo to simulate signal on the pad plane of the TPC readout
- Truncated mean calculations with truncation range from 30-100% for different regions
- Use Random Forest Regression model to estimate input of the simulation from truncated mean values
- Estimate best truncation combination from different pad regions
- Ongoing study





Toy Monte Carlo

- Fast microscopic simulation •
- Sequence of processes •

Simulating 1 track



Input	Primary ionization	Total ionization
(dN_{prim,in}/dx) : Mean number of primary electrons per cm, per track	N_{prim}: Number of primary electrons along the particle trajectory	$N_{tot} \sum^{Nprim} N_{sec}$
Studied in range [10-510] (1/cm)	Poisson distribution	N_{sec}: Number of secondary electrons (power law) generated for each primary electron along the particle trajectory
region track pad length (cm) OROC3 1.5 OROC2 1.2	k 30 25 20 30 25 20 30	track 35 30^{-} 25^{-} $\frac{9}{220^{-}}$ $\frac{9}{220^{-}}$ $\frac{9}{220^{-}}$ $\frac{1}{220^{-}}$ $\frac{1}{220^{-}}$ $\frac{1}{220^{-}}$ $\frac{1}{220^{-}}$ $\frac{1}{220^{-}}$ $\frac{1}{220^{-}}$
OROC1 1 IROC 0.75	10- 5-	10- 5-
$\langle dN_{prim,in}/dx \rangle = 14 (1/cm)$	im 0 100 200 300 400 500 dN _{prim} /dx (1/cm)	$N_{\text{prim}} N_{\text{sec}} = 0$ 100 200 300 400 500 dN_{tot}/dx (1/cm)
5.06.2023 HADRON Conference 2023	Tuba Gündem Goethe Universität	

Toy Monte Carlo

- Fast microscopic simulation
- Sequence of processes



Total ionization	GEM transparency	Amplification
N_{tot} : $\Sigma^{Nprim} N_{sec}$	GEM transparency : Studied in range [50%-100%]	Q : \sum^{Ntot} transparency fluctuations + gas gain fluctuations
N _{sec} : Number of secondary electrons (power law) generated for each primary electron along the particle trajectory	electrons passing through GEM holes electrons approaching	Gas gain fluctuations: exponential, mean=1
track 35 Entries = 152 Mean = 57.5	Transparency fluctuations: uniform, mean=1	Sat _{on/off} : Saturation cut
30- 25-	e ⁻	25- 20-
	GEM s	
	pad plane	5-
prim N _{sec} 0 100 200 300 400 500 dN _{tot} /dx (1/cm)	GEM transparency = 50%	0 100 200 300 400 500 d/dx (ADC/cm)

05.06.2023 | HADRON Conference 2023

Tuba Gündem | Goethe Universität

Truncation

- ALICE
- $dQ/dx (dQ/dx \approx dE/dx)$ values are sorted for each track and truncated by each percentage value:
 - o 30, 40, 50, 60, 70, 80, 90, 100 %
- Mean is calculated for different regions:
 - IROC, OROC1, OROC2, OROC3 and all



- Ensemble of decision trees
- Each tree is created from a randomly sampled subset of data
- Each tree makes its own prediction
 - Var_i: truncated mean values
 - Prediction target: $\langle dN_{prim.in}/dx \rangle = 14$ (for a MIP)
- These predictions are averaged to produce a single result
- Number of trees and depth of a tree are optimized using Grid Search



- Prediction of $\langle dN_{\text{prim.in}}/dx \rangle$ with a single estimator
 - Estimator: truncated mean (percentage, region)
 - Best single estimator + truncated mean (60%, all)





- Prediction of $\langle dN_{prim.in}/dx \rangle$ with a combined estimator
 - Estimator_i: truncated mean[(percentage_i, IROC), (percentage_i, OROC1), (percentage_i, OROC2), (percentage_i, OROC3)]





^{05.06.2023 |} HADRON Conference 2023

- Prediction of $\langle dN_{prim.in}/dx \rangle$ with a combined estimator
 - Estimator_i: truncated mean[(percentage_i, IROC), (percentage_i, OROC1), (percentage_i, OROC2), (percentage_i, OROC3)]





^{05.06.2023 |} HADRON Conference 2023

Summary & Outlook

- PID response parameterization
 - Bethe-Bloch parameterization is obtained from Optuna or ROOT fit (clean samples)
 - Mean and sigma are estimated in multidimensional space from real data using NNs
 - Application on analysis with ONNX
- Optimization of dE/dx estimator using Random Forest Regression
 - Effects of different truncation ranges on pad regions studied
 - Toy Monte Carlo More detector effects to be added
 - Combinations of different truncated mean (percentage, region) values studied





Thank you for your attention