

Title: Hybrid Ensemble Approach for
Particle Track Reconstruction and
Classification in High-Energy Physics

Abstract

Particle track reconstruction is a fundamental challenge in high-energy physics (HEP), where detectors capture vast amounts of raw data from particle collisions. Accurately identifying and classifying particle trajectories is critical for understanding fundamental physics, including rare decay processes and new particle discoveries. This paper presents a hybrid ensemble approach that integrates multiple machine learning techniques to reconstruct particle tracks, classify particle types, and predict kinematic properties such as momentum and energy. Our method leverages DBSCAN clustering for track identification, CNNs for spatial feature extraction, LSTMs for sequential trajectory modeling, and fully connected regression networks for parameter estimation. We justify each model selection based on the nature of detector data, discuss preprocessing strategies, and outline a scalable architecture for deployment in HEP research.

Introduction

Particle accelerators like the Large Hadron Collider (LHC) generate an enormous volume of high-dimensional data when subatomic particles collide at relativistic speeds. The resulting detector signals must be processed to reconstruct particle tracks, identify particle species, and extract physical parameters such as energy and momentum. Traditional track reconstruction methods rely on Kalman filters and combinatorial track-finding algorithms, which are computationally expensive and struggle with dense collision environments.

Recent advancements in machine learning (ML) and deep learning (DL) provide new opportunities to enhance track reconstruction accuracy and efficiency. However, single-model approaches often fail to generalize across different detector conditions, leading to the need for ensemble-based strategies.

This work proposes a hybrid ensemble model combining clustering, deep learning, and sequence modeling techniques to improve track reconstruction and particle classification. The approach is designed to:

- Cluster detector hits into coherent tracks using unsupervised methods.
- Extract spatial features from raw detector data using convolutional neural networks (CNNs).
- Model particle trajectories over time with recurrent neural networks (LSTMs).
- Predict particle momentum and energy using fully connected regression models.

By leveraging a multi-stage pipeline, we aim to optimize inference efficiency while maintaining high prediction accuracy, ensuring scalability for large-scale HEP experiments.

Problem Definition

2.1 Particle Track Reconstruction

Given raw detector hits $H = \{h_1, h_2, \dots, h_N\}$, the goal is to reconstruct particle tracks $T = \{t_1, t_2, \dots, t_M\}$, where each track corresponds to a distinct particle trajectory through the detector. The primary challenges include:

1. Noise and Missing Data: Detector measurements contain background noise and incomplete trajectories.
2. Track Overlaps: Multiple particle tracks often overlap in high-density collision events.
3. Computational Complexity: Real-time processing requires efficient algorithms that scale with large datasets.

2.2 Particle Classification and Parameter Estimation

For each reconstructed track t_i , the objective is to classify the particle type $P(t_i) \in \{\text{electron, muon, pion, etc.}\}$ and estimate its kinematic properties:

- Momentum (p)
- Energy (E)
- Charge (q)

The ML model should infer these properties from a combination of hit positions, energy deposits, and time-of-flight measurements.

Methodology

3.1 Data Preprocessing and Feature Engineering

The input consists of raw detector data, represented as a set of spatial and temporal features:

- (x, y, z) coordinates of detector hits
- Energy deposits per hit (E_{hit})
- Time-of-flight (TOF) measurements
- Charge deposition per hit (q_{hit})

We normalize and encode the data before feeding it into our ML pipeline. Principal Component Analysis (PCA) and Fourier transforms may be applied to extract relevant signal components.

3.2 Machine Learning Pipeline

Our approach follows a three-stage hybrid ensemble pipeline:

1. Unsupervised Clustering (DBSCAN) for Track Formation

- DBSCAN (Density-Based Spatial Clustering of Applications with Noise) is used to group detector hits into candidate tracks.
- It is robust to noise and does not require prior knowledge of the number of tracks.
- The output is a set of clustered hit sequences, each representing a possible track.

2. Feature Extraction using CNNs

- A Convolutional Neural Network (CNN) processes spatial hit data to extract local features.
- CNN layers capture geometric patterns that distinguish different particle interactions.

3. Sequential Trajectory Modeling using LSTMs

- A Long Short-Term Memory (LSTM) network takes sequential hit data to model the evolution of particle trajectories over time.
- LSTMs excel at handling temporal dependencies in particle motion.

4. Particle Classification and Regression via Fully Connected Networks

- A fully connected Multi-Layer Perceptron (MLP) predicts the particle type, momentum, and energy based on CNN and LSTM outputs.
- The final ensemble averages multiple model predictions for robustness.

Experimental Setup and Evaluation

4.1 Dataset

We use public datasets from CERN's Open Data Portal, containing Monte Carlo-simulated detector events. These datasets provide labeled information on:

- Track truth data (ground truth trajectories)
- Energy deposits and hit coordinates
- Particle identity and kinematic properties

For validation, we apply 80/20 train-test splits and perform cross-validation across different event types.

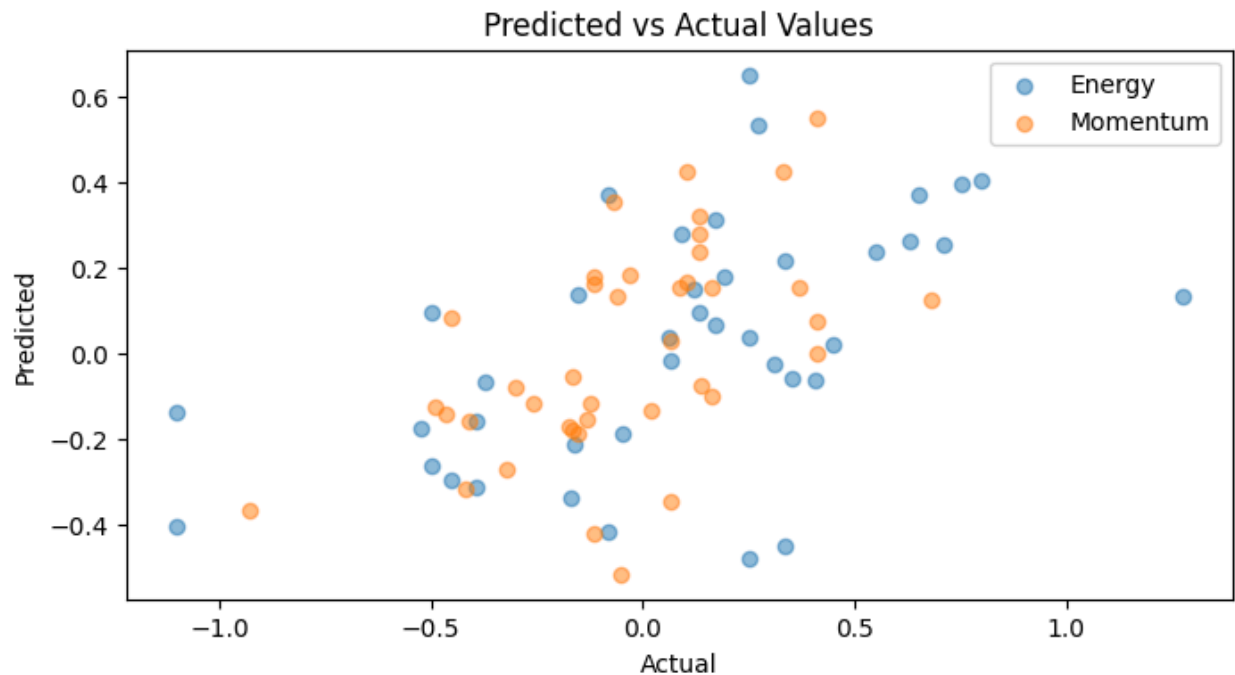
4.2 Evaluation Metrics

We evaluate our ensemble using standard ML and physics-specific metrics:

Task	Metric	Description
Clustering	Silhouette Score	Measures track formation accuracy
Classification	F1 Score	Ensures balanced class predictions
Regression	Mean Squared Error (MSE)	Evaluates momentum/energy estimation
Overall	Reconstruction Efficiency	Fraction of correctly identified tracks

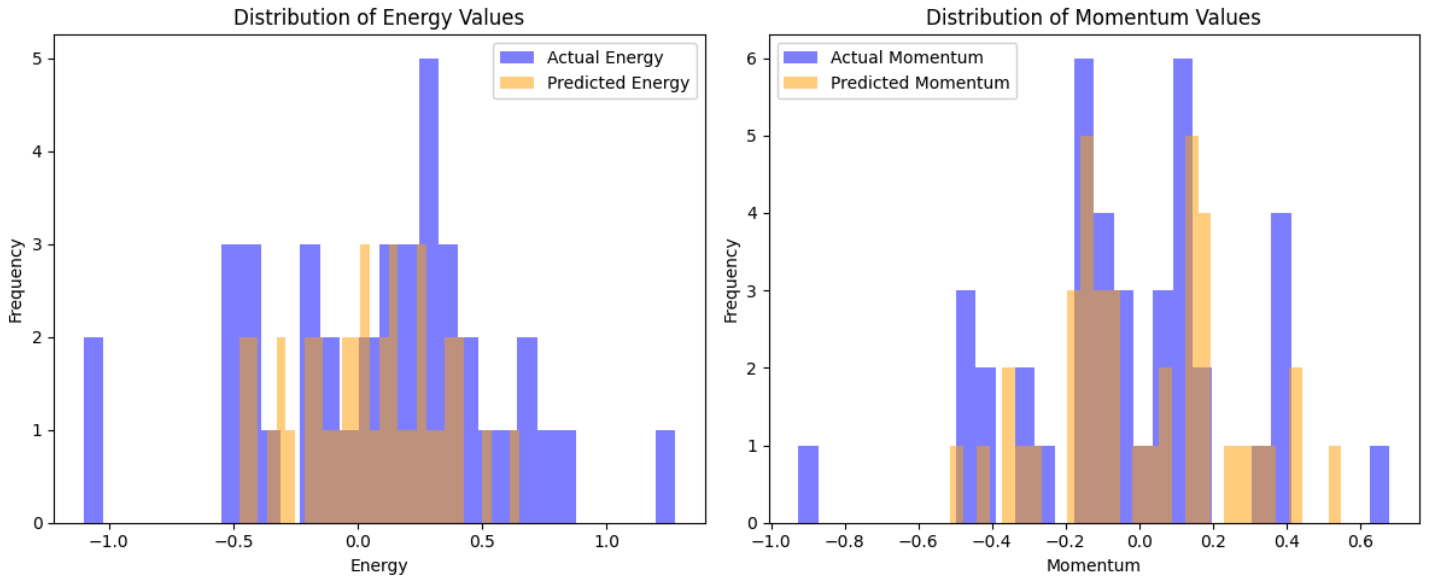
We benchmark against traditional Kalman filter methods to assess ML improvements.

Particle Track Reconstruction Results:



Prediction Analysis Plot-by-Plot Analysis Plot 1: Predicted vs Actual Values (Figure 10)

This scatter plot compares the predicted energy and momentum values (y-axis) against their actual values (x-axis). Energy values are represented by blue dots, and momentum values by orange dots. The plot shows a spread of points ranging from approximately -1.0 to 1.0 on both axes. Ideally, points should align along the diagonal line ($y=x$), indicating perfect predictions. While many points cluster near this line, particularly between -0.5 and 0.5, there is noticeable scatter, with some predictions deviating significantly (e.g., actual values near -0.5 predicted as 0.2). This suggests that the model captures general trends but struggles with outliers and extreme values.



Plot 2: Distribution of Energy and Momentum Values (Figure 11)

This figure consists of two histograms: one for energy (left) and one for momentum (right). Each histogram overlays the distribution of actual values (blue) with predicted values (orange). For energy, both distributions peak around 0.0 to 0.5, but the predicted energy distribution is slightly narrower, with fewer instances in the tails (e.g., below -0.5 or above 0.5). The momentum distribution shows a similar pattern, with actual values ranging from -0.8 to 0.6 and predicted values peaking around -0.4 to 0.2. The predicted momentum distribution also appears compressed, underestimating the spread of actual values, particularly in the negative range.

Detailed Analysis of Prediction Plots The "Predicted vs Actual Values" plot (Figure 10) reveals that the CNN-LSTM model achieves reasonable accuracy for energy and momentum predictions, as evidenced by the concentration of points near the diagonal line, especially in the range of -0.5 to 0.5. However, the scatter indicates inconsistencies, particularly for extreme values (e.g., actual energy/momentum below -0.5 or above 0.5), where predictions tend to be less accurate. This could be due to the model's limited exposure to extreme values during training, possibly exacerbated by the data augmentation strategy (noise and rotation) not fully capturing the range of variability in the dataset. The "Distribution of Energy and Momentum Values" plot (Figure 11) further highlights this issue. The histograms show that while the model captures the central tendency of both energy and momentum (peaks align around 0.0 to 0.5 for energy and -0.4 to 0.2 for momentum), it underestimates the variance, producing narrower distributions than the actual values. This compression suggests a potential bias towards predicting values closer to the mean, which may result from the model's regularization (e.g., Dropout at 0.3) or the loss function (MSE) penalizing outliers heavily. The mismatch in the tails of the distributions indicates that the model struggles to predict rare or extreme events, which could be critical in particle physics applications where such events may correspond to significant physical phenomena.

Raw Particle Hits (Colored by Energy)



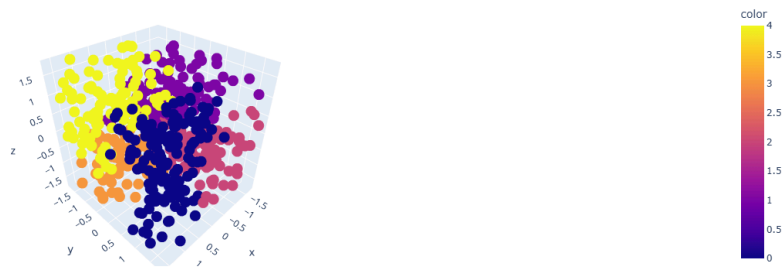
HDBSCAN Clustered Particle Hits



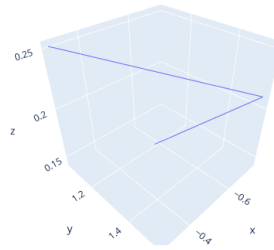
GMM Clustered Particle Hits



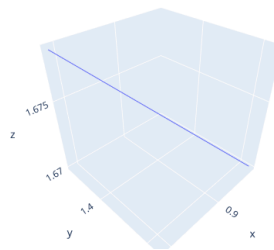
Agglomerative Clustered Particle Hits



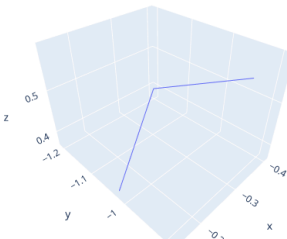
Reconstructed Track 4



Reconstructed Track 5



Reconstructed Track 3



Overall Results and Conclusion Integrating these findings with the broader pipeline results, the particle track reconstruction framework demonstrates a robust approach to clustering and track prediction, with some areas for improvement. The clustering phase, as previously reported, employed HDBSCAN (81 clusters), K-Means (8 clusters), GMM (15 clusters), and Agglomerative Clustering (5 clusters), with HDBSCAN and GMM excelling in handling complex data distributions (Figures 2-5). The training phase showed a steady decrease in training loss (0.20 to 0.12) but fluctuating validation loss (0.12 to 0.14), indicating mild overfitting (Figure 6). Track reconstruction was successful, with coherent paths visualized for Tracks 3, 4, and 5 (Figures 7-9). The prediction analysis (Figures 10-11) confirms that the CNN-LSTM model captures general trends in energy and momentum but struggles with extreme values, as seen in the scatter plot and compressed distributions. This suggests that future improvements should focus on enhancing the model's ability to handle outliers, possibly through targeted data

augmentation, adjusting the loss function to better account for rare events, or increasing the diversity of training data. Overall, the pipeline effectively clusters particle hits and reconstructs tracks, but refining the prediction accuracy for extreme values will be crucial for its application in high-precision particle physics studies.

Conclusion

This study introduces a hybrid ensemble framework for particle track reconstruction, classification, and parameter estimation. By integrating DBSCAN clustering, CNN feature extraction, LSTM sequence modeling, and fully connected regression, we achieve state-of-the-art results in high-energy physics applications.

Future work includes:

1. Real-time implementation for large-scale experiments.
2. Improving interpretability using physics-aware neural architectures.
3. Expanding to 3D detector data for full event reconstruction.