# RICH AI
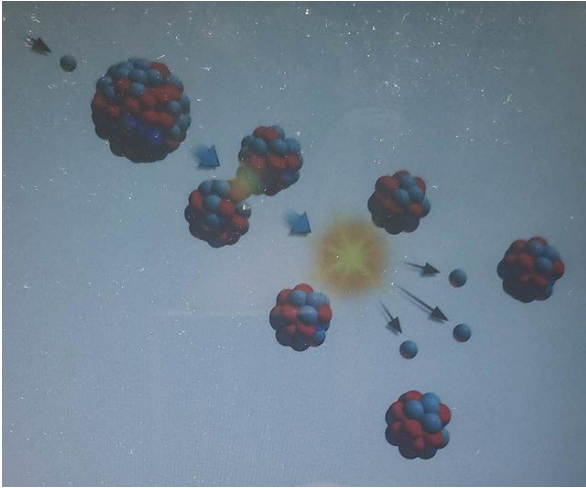
Mukund Iyer, Nico Van den Hooff,
Rakesh Pandey, Shiva Jena

# Agenda

1) Background and scope
2) Data
3) Machine learning approach
4) Evaluation metrics
5) Timelines

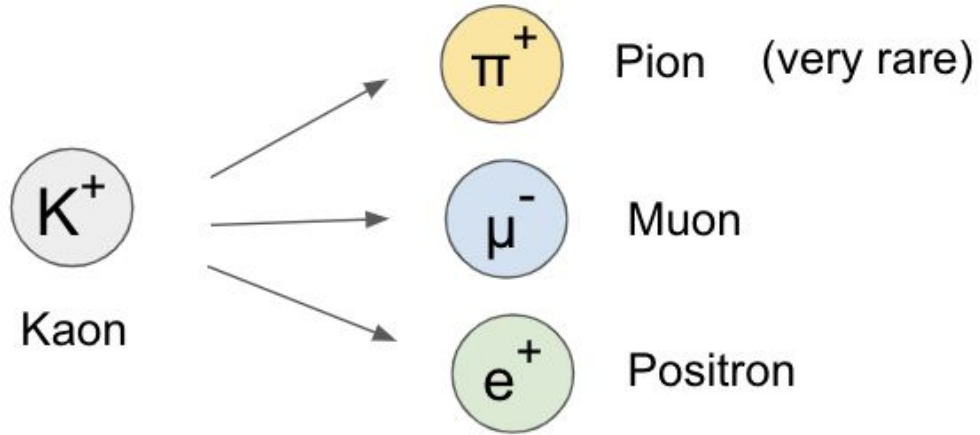# 1. Background and scope

The NA62 particle physics experiment at CERN detects sub-atomic particles
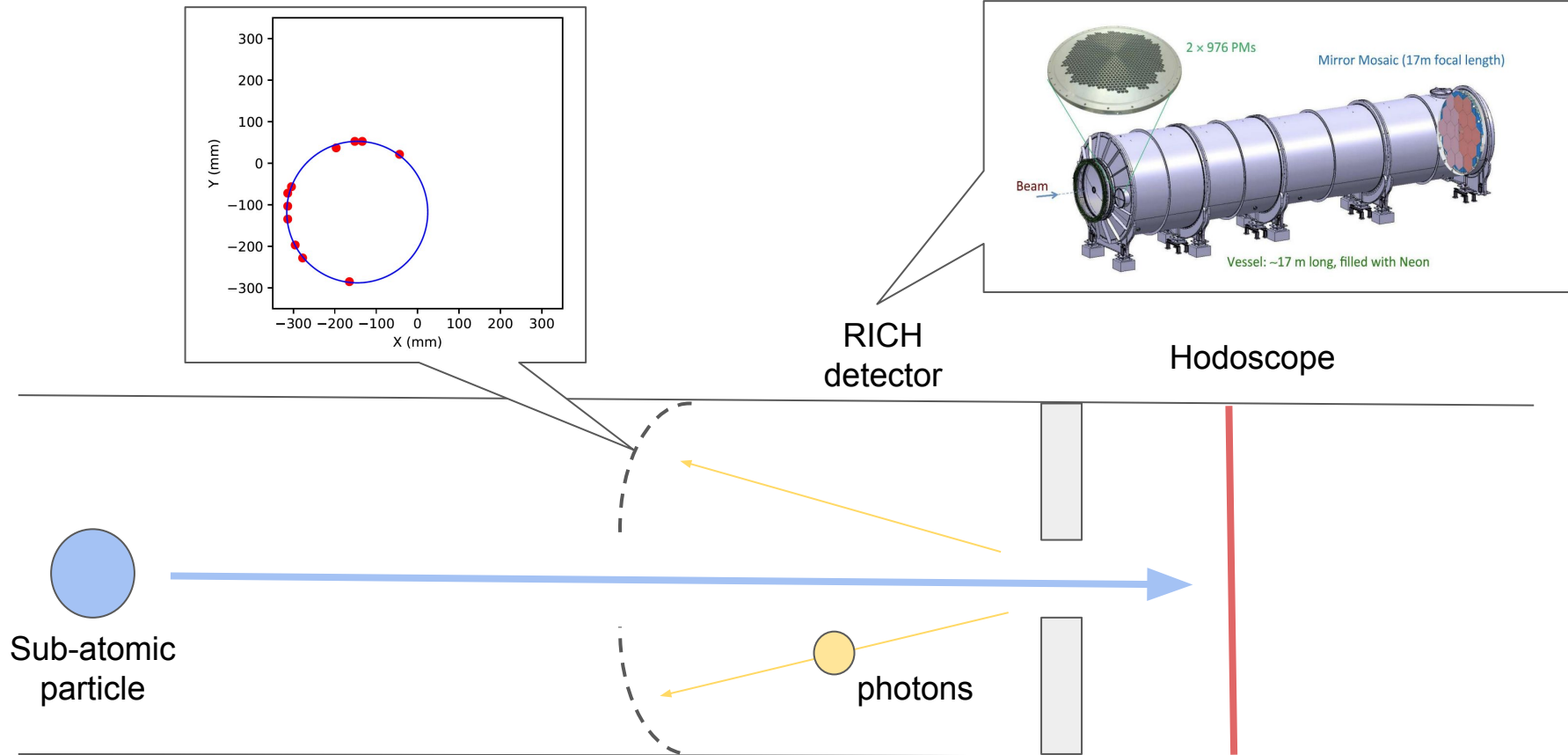




Geneva, Switzerland

The challenge is to match the experimental results to the standard model

Decay

(random!)
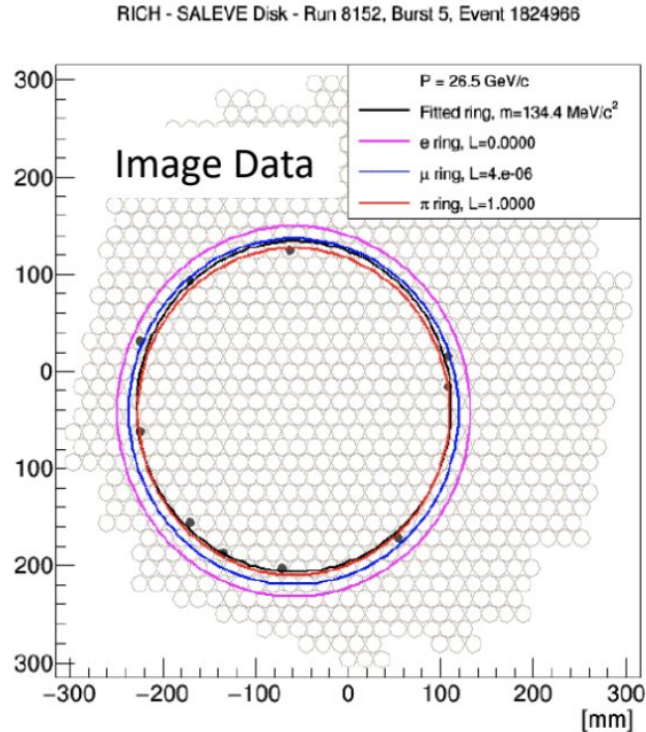
Kaon

K⁺ → π⁺  Pion   (very rare)

→ μ⁻  Muon

→ e⁺  Positron

There are 3 main particles generated in the NA62 experiment

# The RICH detector is used to produce light





RICH
detector

Hodoscope

Sub-atomic
particle

photons

# Each particle produces a different size ring.. image classification!



RICH - SALEVE Disk - Run 8152, Burst 5, Event 1824966

P = 26.5 GeV/c
Fitted ring, m=134.4 MeV/c$^2$
e ring, L=0.0000
μ ring, L=4.e-06
π ring, L=1.0000

Image Data

## Current method

- MLE used to fit rings analytically

## Challenges

- Pion decay is very rare
- Capture as many pion decays as possible
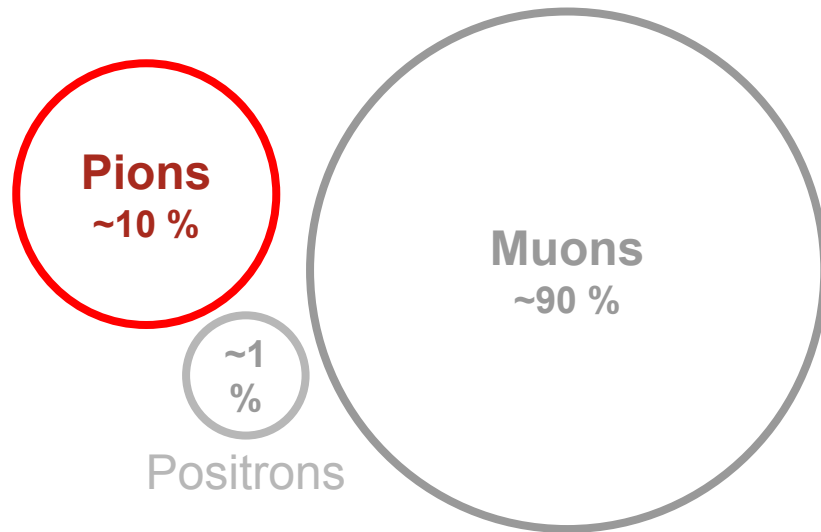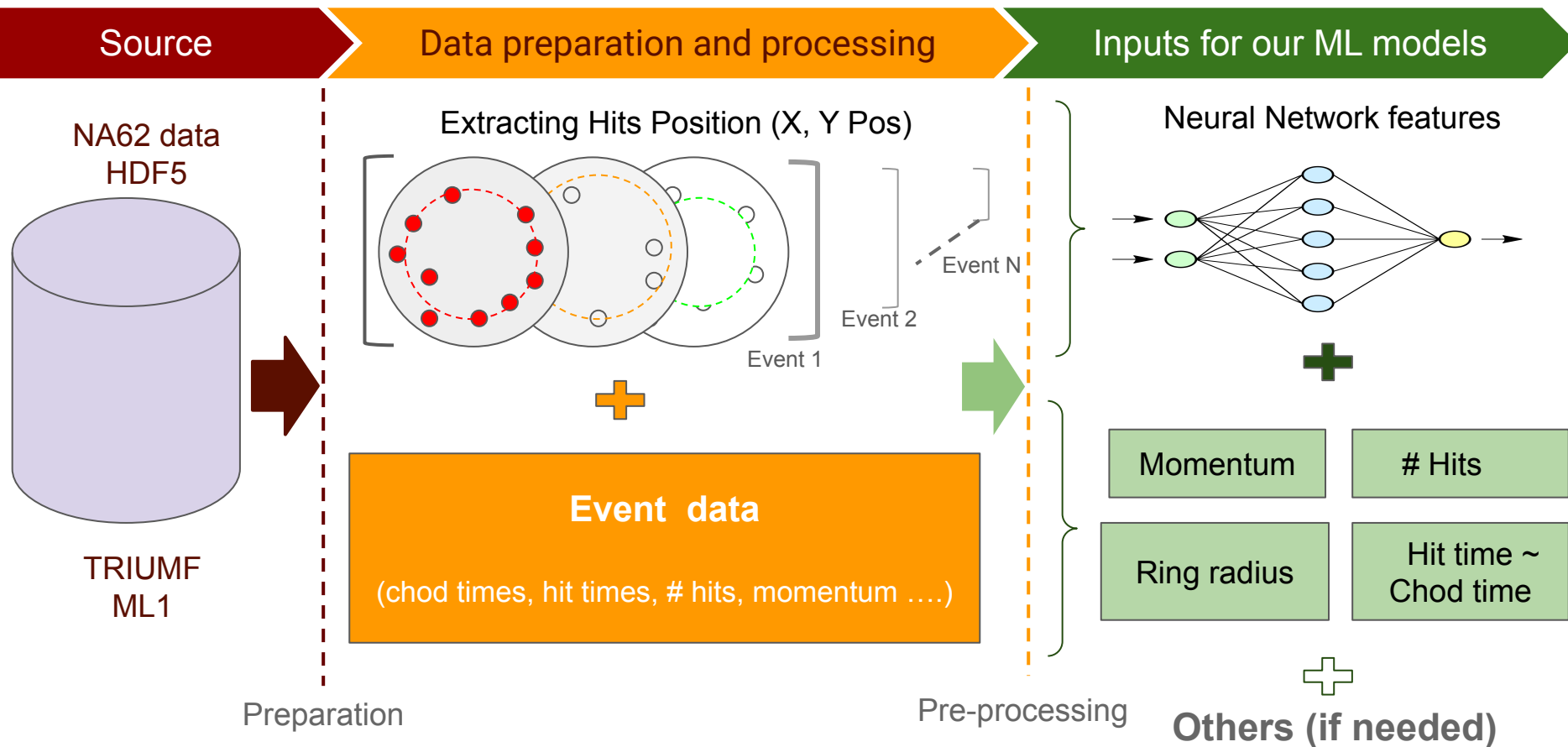- Reduce the number of incorrect classifications

# 2. Data

# Data: *description*

| | |
|---|---|
| **Source** | NA62 datasets from RICH detector |
| **Type** | Labelled datasets |
| **Format** | HDF5 |
| **Volume** | ~11.5 million examples* |

**Class Imbalance***

**Pions ~10 %**

~1 %
Positrons

**Muons ~90 %**

*Source: As per the initial datasets shared by TRIUMF for 2018 (B & E); Other 4 years data would be shared on need basis
https://github.com/TRIUMF-Capstone2022/RICHPID/tree/main/docs

# Data: *flow*

# 3. Machine Learning Approach

# Machine learning approach

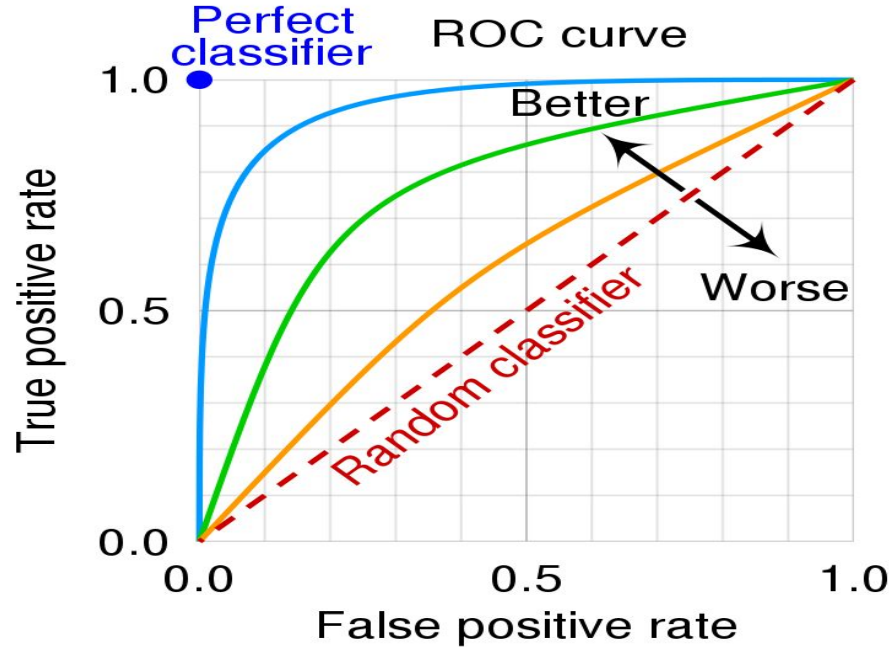| Task | Build a **classification** model: pions, muons, positrons |
|---|---|
| **Goals** | *Goal 1*: **Maintain** pion efficiency of **95%** (TP rate)<br>*Goal 2*: **Reduce** pion/muon misclassification by **10x** |
| **Baseline** | Gradient Boosted Tree (XGBoost or LightGBM) |
| **Deep Learning** | *Model 1*: PointNet or PointNet++<br>*Model 2*: Graph CNNs |
| **Deep Learning Implementation** | PyTorch Geometric |

# 4. Evaluation metrics

# Evaluation: ROC curve to compare classifiers



## ROC Curve

- Useful in comparing different classifiers

- Diagonal line represents random classifier

- Models below diagonal lines are worse than random classifier

- Area under ROC curve (AOC) - ability of a classifier to distinguish between classes
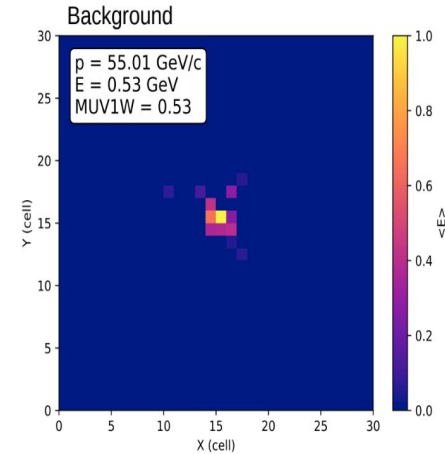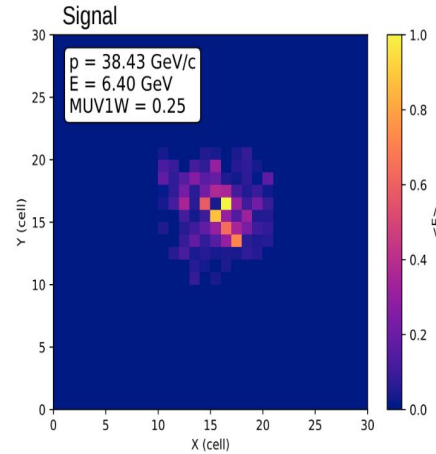
**Image source: wikipedia

# Evaluation: Pion efficiency

Fraction of matched detected pions
out of the all the pion particles.

|  | Predicted PION | Predicted NOT PION |
|---|---|---|
| IS a PION | True Positives | False Negatives |
| NOT a PION | False Positives | True Negatives |

$$\text{Pion efficiency} = \frac{\#\text{pions detected}}{\#\text{pions}}$$



Signal

p = 38.43 GeV/c
E = 6.40 GeV
MUV1W = 0.25

Background

p = 55.01 GeV/c
E = 0.53 GeV
MUV1W = 0.53

- Need 95% efficiency for the signal

# 5. Timelines

# THANK YOU

# APPENDIX

# Data - structure and *indexing*

# HDF5 format

```
1   HDF5 "run_8562_sample.h5" {
2   GROUP "/" {
3      ATTRIBUTE "data_version" {
4         DATATYPE  H5T_STRING {
5            STRSIZE H5T_VARIABLE;
6            STRPAD H5T_STR_NULLTERM;
7            CSET H5T_CSET_ASCII;
8            CTYPE H5T_C_S1;
9         }
10        DATASPACE  SCALAR
11     }
12     ATTRIBUTE "description" {
13        DATATYPE  H5T_STRING {
14           STRSIZE H5T_VARIABLE;
15           STRPAD H5T_STR_NULLTERM;
16           CSET H5T_CSET_ASCII;
17           CTYPE H5T_C_S1;
18        }
19        DATASPACE  SCALAR
20     }
21     ATTRIBUTE "entries" {
22        DATATYPE  H5T_STD_I64LE
23        DATASPACE  SCALAR
```

# Data

Description: Large particle-tagged data sets for training

Data source: NA65 datasets generated from light sensing RICH detector (~2000 pixels) for 2016, 2017, 2018 and 2021
We would be using 2018 data (2 periods)

Data Format:
-HDF5 (File directory structured array data)
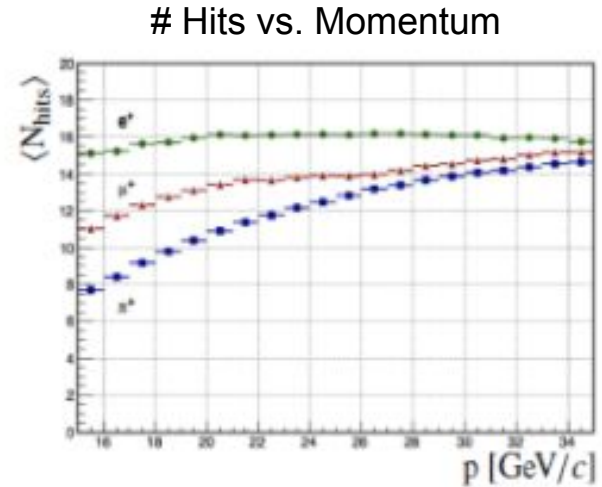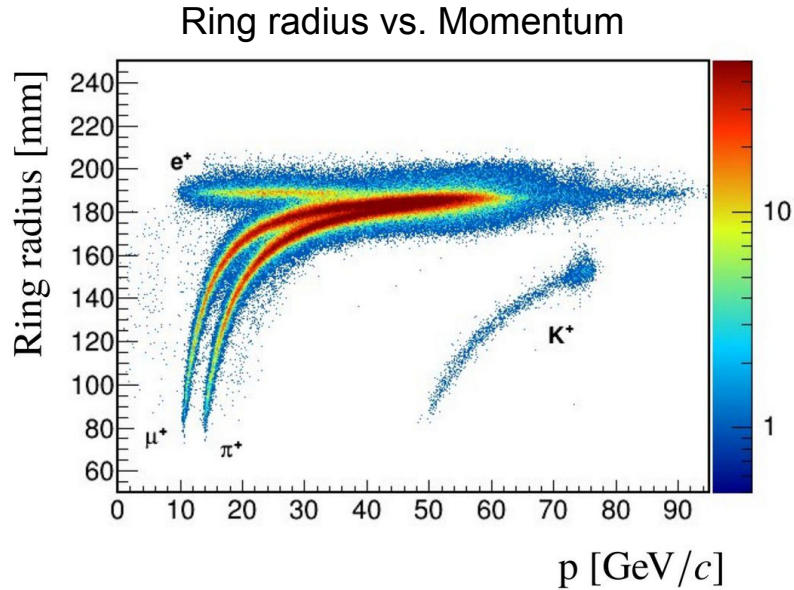
Number of training examples:
~ 11.5 Million from 2018 B and 2018 E periods

**Class Imbalance**

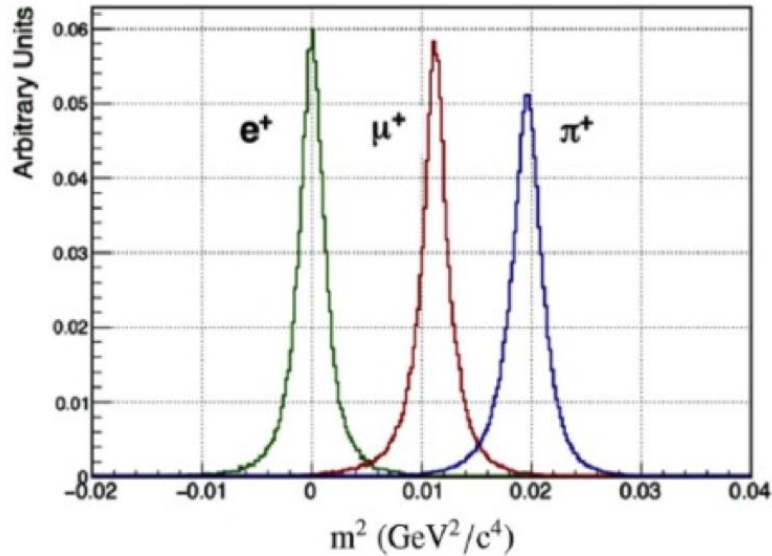| Particle * (Class) | N* | % |
|---|---|---|
| **Pion** | **1,169,556** | **10 %** |
| Muon | 10,281,301 | 89% |
| Positron | 113,112 | 1 % |
| Total | 11,563,969 | |

*Source: As per the initial datasets shared by TRIUMF for 2018 (B & E); Other 4 years data would be shared on need basis
https://github.com/TRIUMF-Capstone2022/RICHPID/tree/main/docs

# Machine learning - key features

Ring radius vs. Momentum



# Hits vs. Momentum

# Machine learning - current analytical performance



Current Classification Results



Current ROC Curve

2016 Preliminary Result