# Hadronic Shower Substructure Reconstruction using Graph Neural Networks

**AHCAL main meeting**
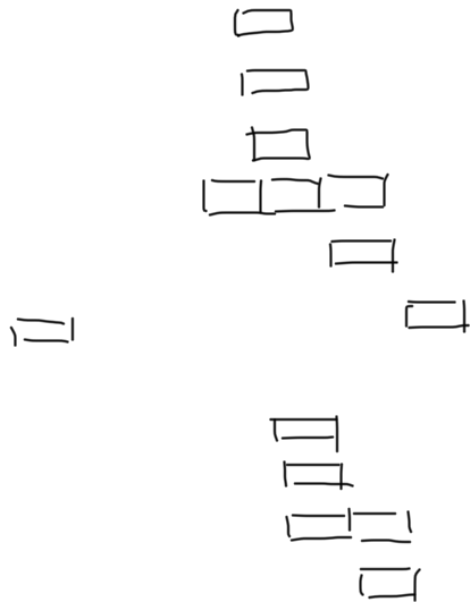
Vladimir Bocharnikov (DESY)

8 Dec 2021

Vladimir Bocharnikov (DESY)

# Calorimeter vision for hadronic showers

## Ultimate goal and general approach

**Set of hits in highly granular calorimeter**

**Particle interaction tree**

Per particle: ID, E, $\mathbf{p}$, $\mathbf{v}_{prod}$, $\mathbf{v}_{decay}$

**Potential applications of hit to secondary particle association:**

- Hadronic energy reconstruction

- Shower separation algorithms:

  - Recombination of secondaries between overlaid showers

- Validation of simulation performance:

  - Comparison of global physical distributions

  - Shower description on single event basis is possible
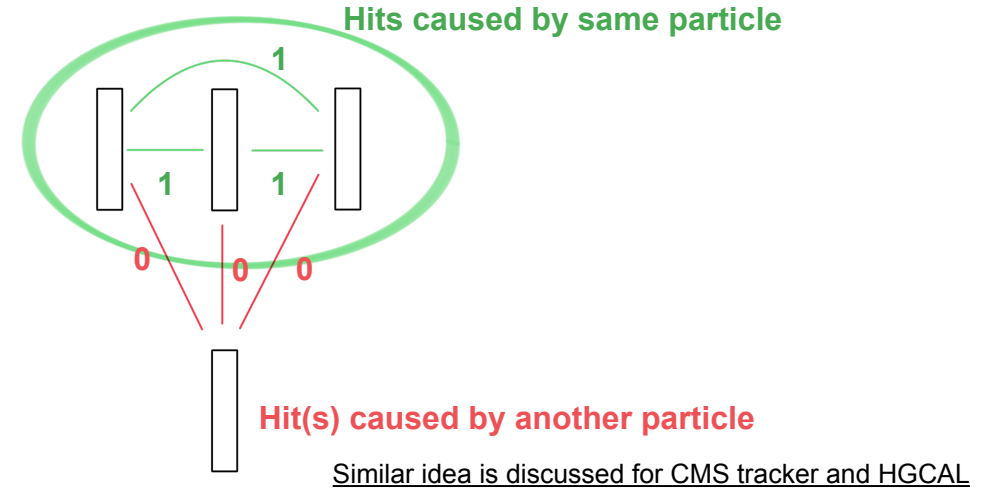
    ➡ essential for adversarial networks

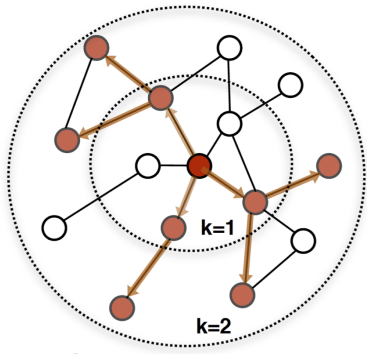# Graph representation of calorimeter event

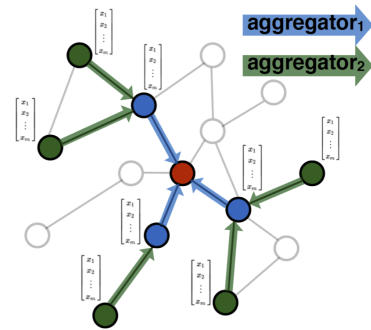## First steps

**Event graph:**

- ○ Nodes - calorimeter hits
- ○ Node features - position, energy, (time)
- — Edges - neighbours within distance < $R_{max}$ (Radius graph)
- — Edge weights - 1 if pair of hits belong to same **fundamental object** (e/m sub-shower, track), otherwise 0
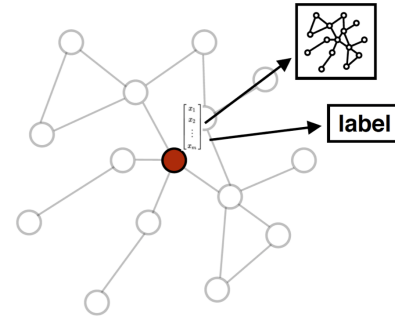- ○ ML **objective** - **predict edge weights** given the radius graph of event



Hits caused by same particle

Hit(s) caused by another particle

Similar idea is discussed for CMS tracker and HGCAL

**GraphSAGE** (**SA**mple and aggre**G**at**E**) architecture (Graph neural network model (GNN))**:**
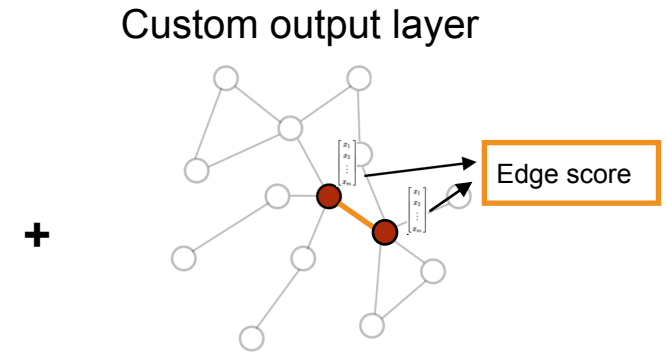


Sample neighbourhood of graph nodes

Aggregate feature information from neighbours

Get graph context embeddings for node using aggregated information

**+**

Custom output layer

Edge score

Predict edge score for each pair of connected nodes using embedded features

# Truth information from Monte-Carlo

## Algorithm to find truth e/m objects

**Simulations**

*Geant4 (v10.03.p02)* QGSP_BERT_HP using CALICE AHCAL geometry

Pure energy deposition in cells (before digitalisation and reconstruction)

**Truth electromagnetic sub-shower definition:**

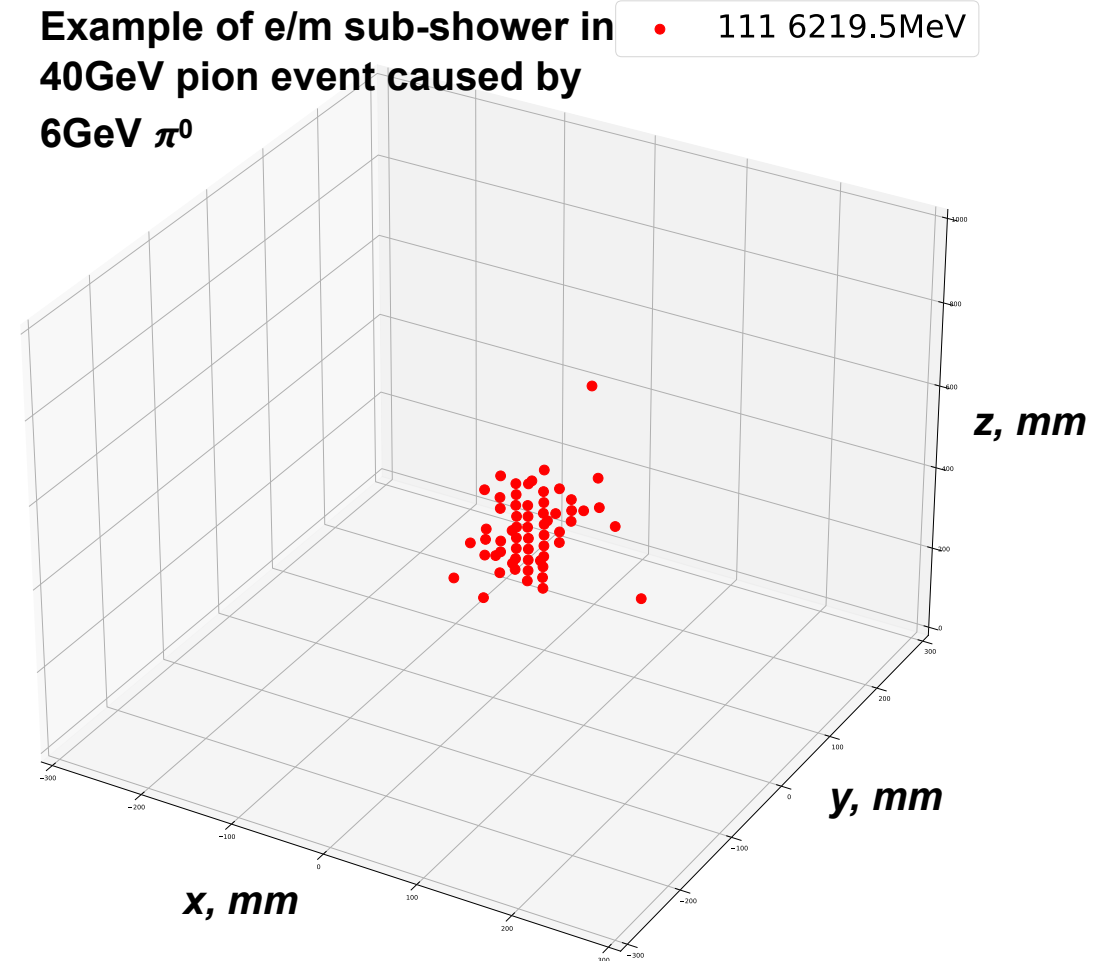"**Electromagnetic**" particles: $e^{\pm}, \gamma, \pi^0, \eta$

Energy threshold - *0.1GeV* (arbitrary now)

If MC particle is "electromagnetic", all it's "electromagnetic" daughters compose e/m shower are removed from further consideration

Corresponding simulated hits compose sub-shower,

0.5MIP cut: $E_{hit} > 0.25MeV$

**Example of e/m sub-shower in 40GeV pion event caused by 6GeV $\pi^0$**



111 6219.5MeV

*z, mm*

*y, mm*

*x, mm*

MC history for **ionising particles** is more complicated to easily define individual objects (tracks). Work in progress

# Datasets and model parameters

## Edge score model

**Train&test** dataset:

- ~6000 MC event graphs (50/50 split)
  - Pure energy deposition in calorimeter cells (before digitalisation and reconstruction)
  - **10-100 GeV pion** samples
- ➡ Radius graphs with calorimeter hit nodes $(x, y, z, E_{hit})$ $R_{max}$ = **59 mm**
  - Electromagnetic relation between hits is encoded in edge scores (0/1)
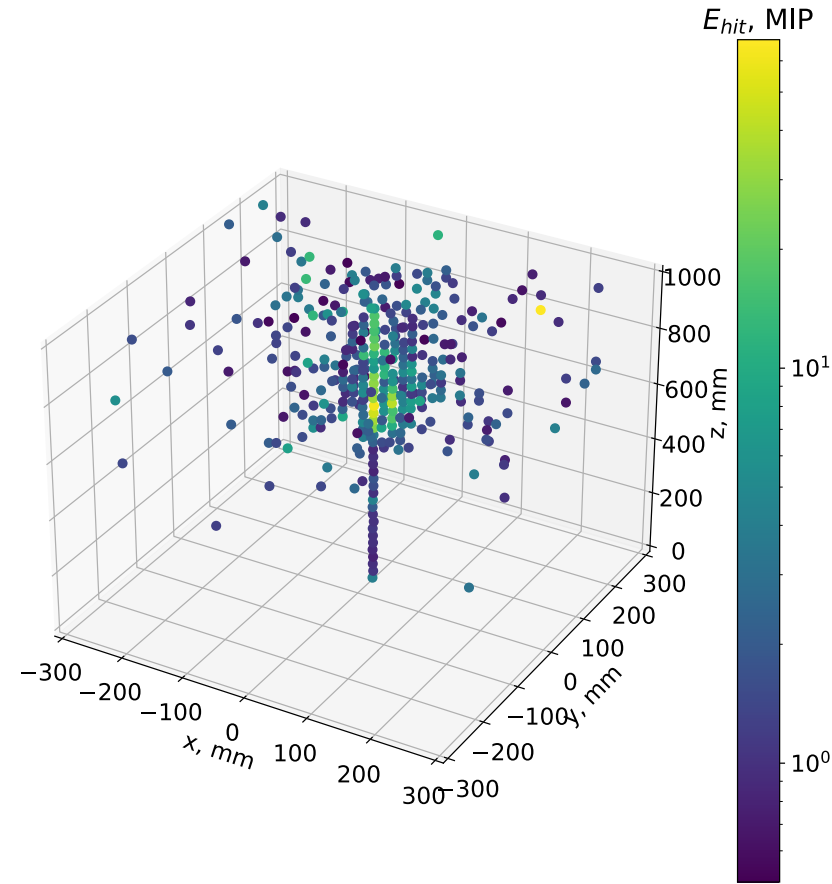
**Model:**

**GraphSAGE** GNN

8 layers with 16 hidden channels + 1 linear output layer to convert node embeddings to edge scores

Objective: prediction of edge scores

Loss: binary cross entropy

# Hadronic shower reconstruction with GNN
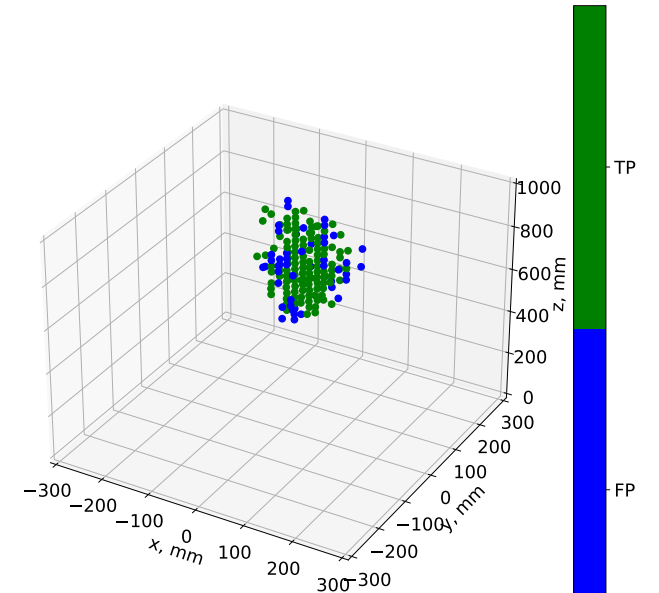
**Results for single test event.**

$E_{hit}$, MIP

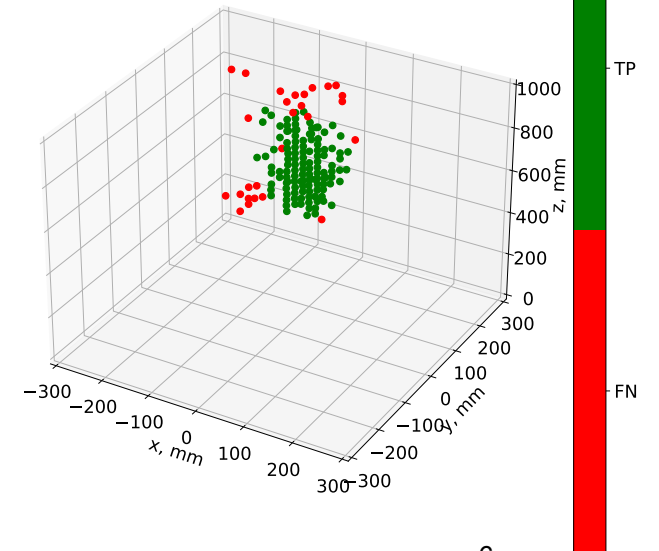*Electromagnetic relation between hits is encoded in graph edge weights:*
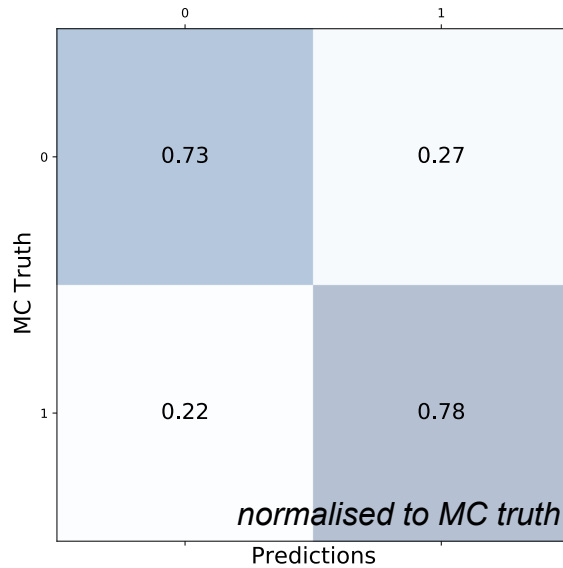
2650 graph edges

*Electromagnetic part of the shower:*
**Cut:** sum over all link attributes per hit > *0.5*
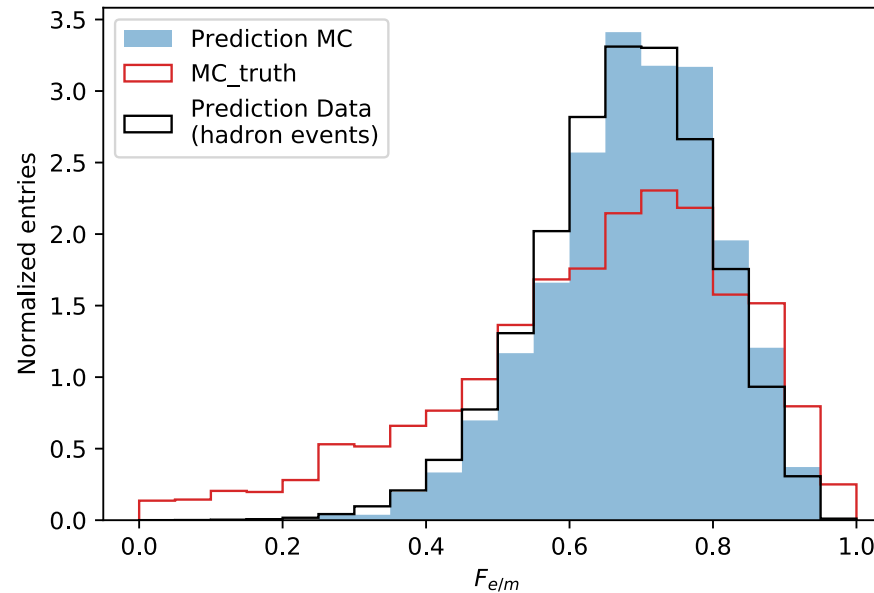
W o r k   i n   p r o g r e s s ...

# Electromagnetic fraction of hadronic showers

## Results for 10,20,30,40,60,80 GeV pions

### Hit classification



### Electromagnetic fraction



### Prediction vs truth correlation



- ~75,5% hit classification accuracy
- Higher MPV for $F_{em}$ than expected
  - ➡ Non-e/m contributions to the hits are not taken into account
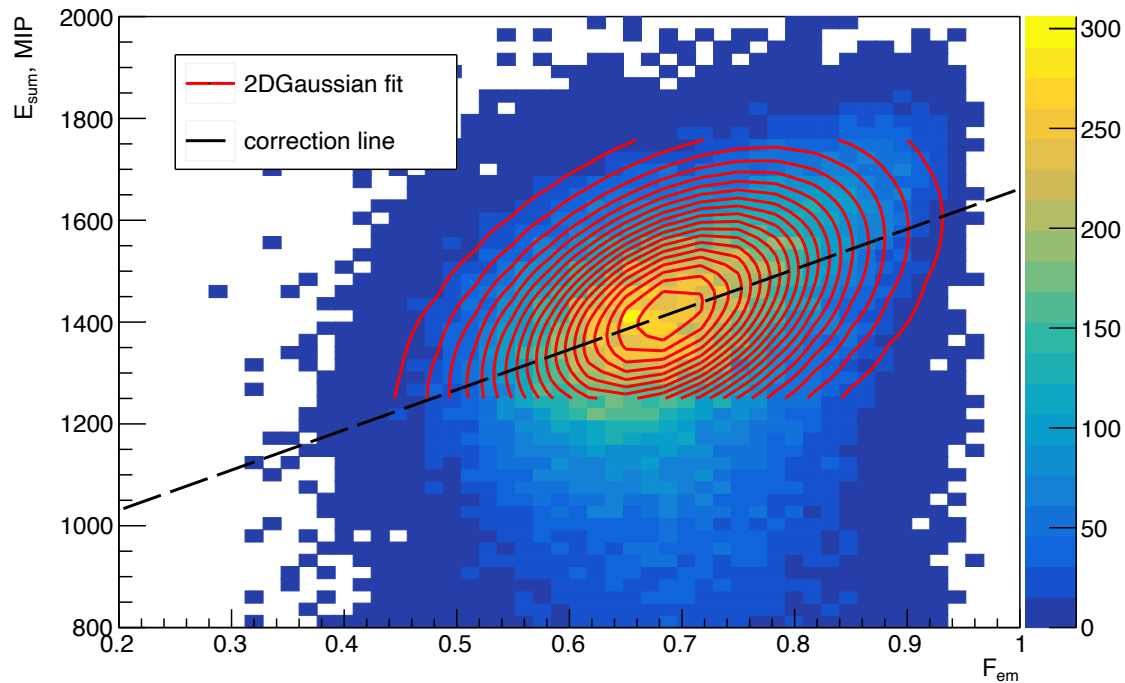- Less pronounced tails for $F_{em}$ prediction than for MC truth

- Reasonable correlation of predicted EM fraction with truth in MPV region

W o r k   i n   p r o g r e s s …
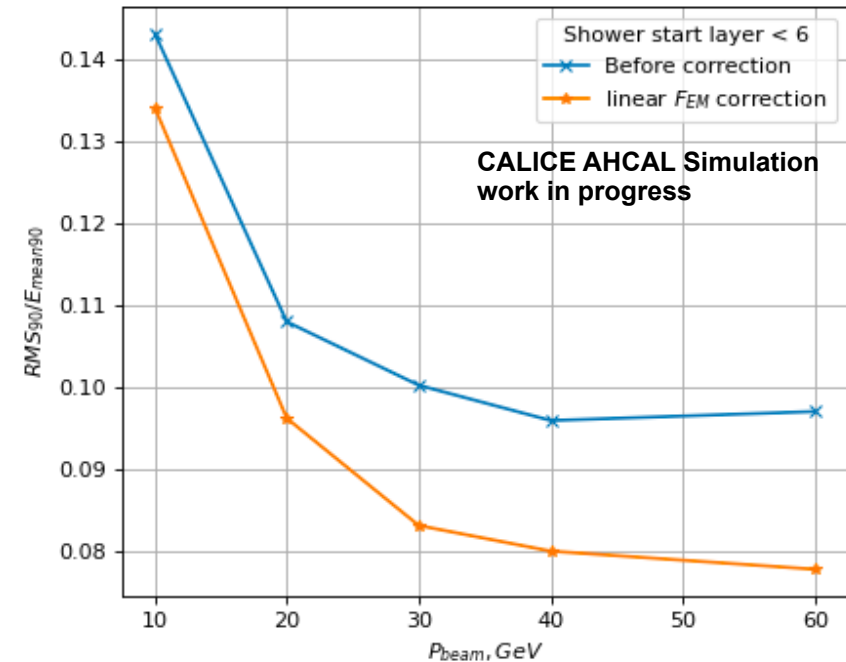
# Hadronic shower reconstruction with GNN

## Using reconstructed EM fraction for energy correction

### Correlation example for 40 GeV pion



### Energy resolution estimation



CALICE AHCAL Simulation
work in progress

- Well pronounced correlation between $E_{sum}$ and $F_{em}$ observed for all energies

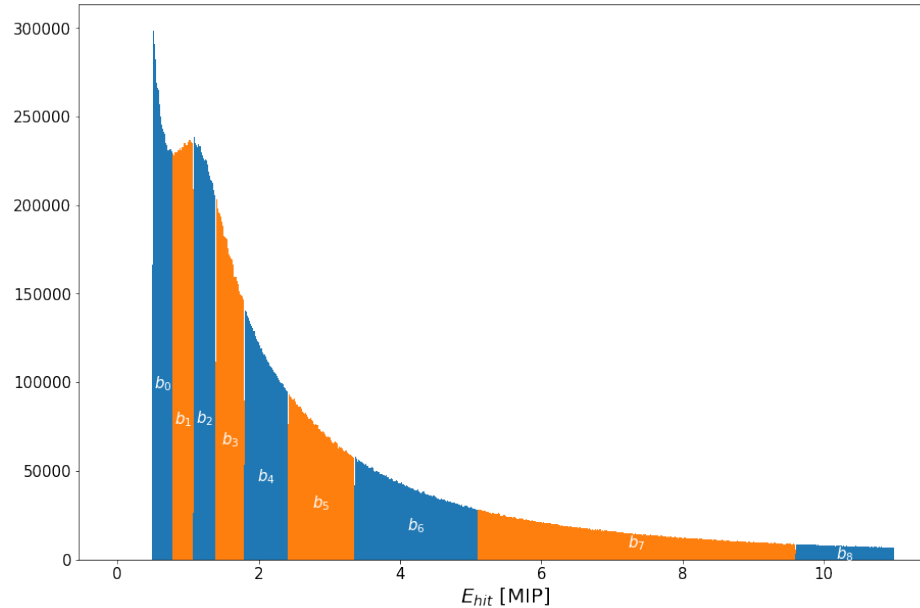- For each energy point simultaneous gaussian fit is performed to extract the correction line

- Simple linear correction gives resolution improvement of ~6-20%

- Promising resolution improvement, baseline for more complex compensation algorithms using reconstructed EM information

# "Standard" LSC

## Code provided by Jack (used as a reference)

> $E_{hit}$ **distribution** split into **bins** of **equal frequency probability**;

> i.e. **equal likelihood (on average!)** of **hits** falling into **each bin**.

> **Three weights defined**, per bin, using **Chebyshev Polynomial**;

> **Fraction of shower energy** falling into **each bin** is **weighted** according to the $E_{sum}$.



$$w_b = w_{b0} + w_{b1}\left(\frac{E_{sum}}{S}\right) + 2w_{b2}\left(\left(\frac{E_{sum}}{S}\right)^2 - 1\right) \quad (1)$$

$S$ is a normalization constant, 150 GeV

- Binning and weights are updated with latest available simulations

- 10-80 GeV range

  - 10K events before shower start cut:

    - 2 < st < 15

    - 28652 events in total

# Energy reconstruction using predicted EM information
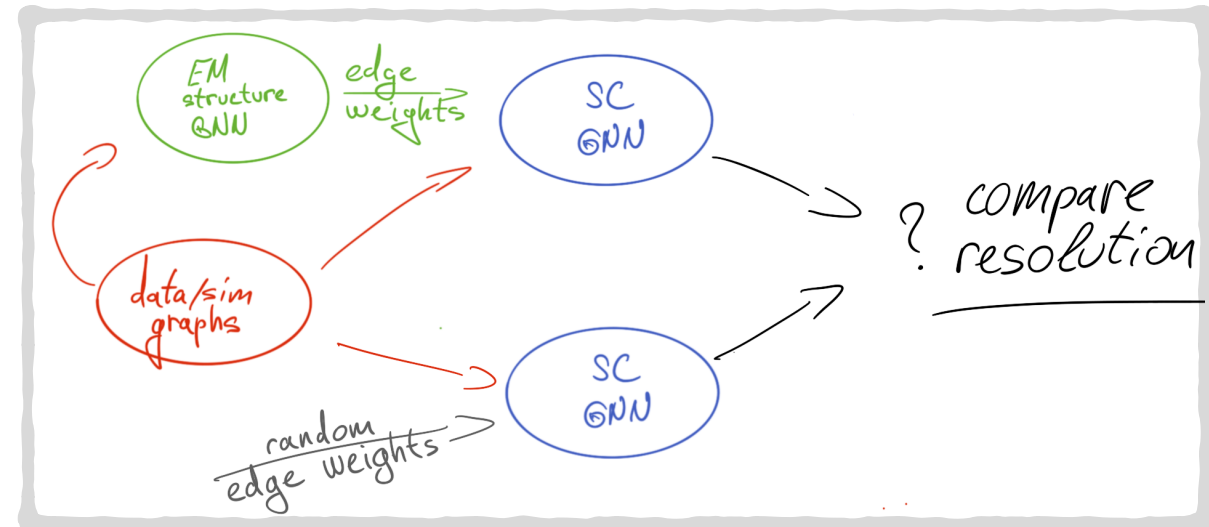
## SC experiment

- Test if use of predicted edge weights improves the energy resolution

- Almost same GNN as for EM structure prediction:

  - 1 GraphSAGE layer replaced with <u>ARMAConv</u> (capable to exploit edge attributes during message passing), output has shape [$N_{nodes}$]

  - Train using predicted EM edge weights

    - Simulations: 10,20,40 GeV, st<6, 30 Kevents

  - Compare resolution for the test sample using predicted EM attributes or random edge weights

    - Simulations: 10,20,30,40,60 GeV, st<6



Training:



1.1

Experiment:

# Resolution and linearity

## 10-60 GeV. Simulations only.



- SC_GNN gains some resolution performance by using reconstructed EM connections between hits

- Problems with LSC linearity are already visible at 60GeV (fit range was up to 80 GeV)
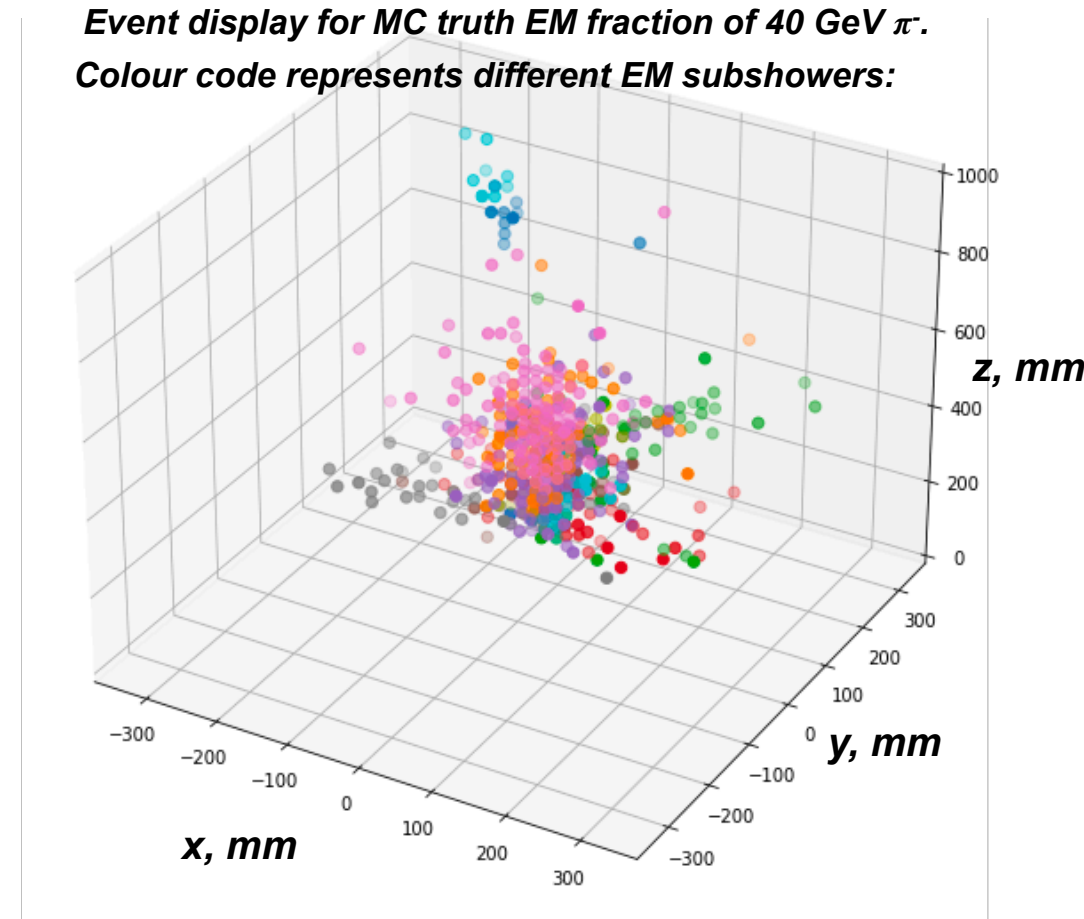
TODOs:

- Estimate leakage effect (check if methods are doing leakage correction in addition to SC) using tail catcher information

- Compare with TB data

# Towards distinct secondary particle reconstruction
## Outlook

**Motivation:**

- In HAD showers we can have many EM subshowers at first HAD interaction (overlaid) and later in the had cascade (displaced)

- Further look into the structure of EM fraction:

  - Reconstruct distinct particle components

    - No easy rule-based algorithm to merge overlaid subshowers on MC truth level ➡ go unsupervised!

    - Test Bayesian Gaussian Mixture model with Dirichlet process on point clouds from calorimeter events

      - SKlearn implementation is tested, own flexible Pyro implementation is planned

  - ➡ Tune training dataset for substructure GNN

    - e.g. energy thresholds (some EM sub showers have topology closer to ionising tracks)



*Event display for MC truth EM fraction of 40 GeV $\pi^-$.*
*Colour code represents different EM subshowers:*

# Applying Bayesian GM to EM component of had showers

## Truth EM component

- SKlearn implementation can handle only scatter plots

- To keep hit energy information, artificial scatter plot is produce:

  - 10 points per MIP

  - uniformly distribute within cell volume: ±15mm,±15mm,±1mm

  - Normalise coordinates: *(-0.36m,0.36m) (-0.36m,0.36m) (0m,1m)*

- Max number of components = 10,

- Object size can be optimised by modifying covariance prior
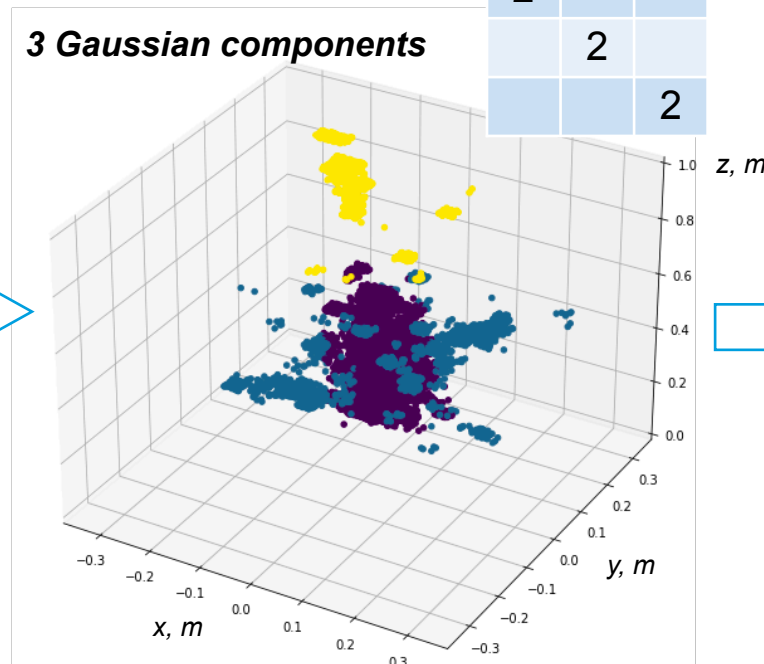
- Clusters can be filtered by likelihood and energy density

*MC truth EM fraction of 40 GeV MC π⁻.*

*Colour code represents different EM subshowers:*

*1st BGM iteration:*

*3 Gaussian components*

Cov prior 1st iter

| 2 | | |
|---|---|---|
| | 2 | |
| | | 2 |

*2nd BGM iteration:*

*7 components*

Cov prior 2d iter

| 0.5 | | |
|---|---|---|
| | 0.5 | |
| | | 0.5 |

# Applying Bayesian GM to EM component of had showers
## Truth vs reco EM component

**1st BGM iteration (truth EM):**

**3 Gaussian components**

Cov prior truth EM

| 2 | | |
|---|---|---|
| | 2 | |
| | | 2 |



**1st BGM iteration (predicted EM):**

**2 Gaussian components**

Cov prior 1st iter

| 2 | | |
|---|---|---|
| | 2 | |
| | | 2 |



**2nd BGM iteration (rest of event) :**

**20 components**

Cov prior 2d iter

| 0.5 | | |
|---|---|---|
| | 0.5 | |
| | | 0.5 |



- Visual similarity for main gaussian component
  - Hints of agreement for $E_{sum}$ and $E_{density}$ on several hundred events between truth and predicted EM fraction (see backup slides)
- Physical observables to be determined and compared with TB data
  - some examples of main GM component distributions in the backup

- Smaller clusters are more challenging

➡ Room for improvement

# Conclusion

- Reconstruction method for electromagnetic substructure of hadronic showers using Graph Neural Networks is presented

- Reconstructed electromagnetic structure can be used to improve hadronic energy resolution

  - GNN software compensation model is capable to exploit EM information

    - can extrapolate and interpolate to different energies

  - Better performance for "standard" local SC to be understood

- Gaussian Mixture model is a promising tool to reconstruct distinct particle contributions within hadronic showers

# Backup

# Single energy examples

# Hit energies
## GNN vs LSC. Simulations

### 10 GeV

### 30 GeV

### 60 GeV

# Truth vs reco EM

## 500 40GeV pion events

# Applying GM to larger dataset

## Some distributions for simulated 40 GeV pions. 10 Kevts.

- Reconstructed EM fraction.

- Shower start found

- Quality metrics (optimised on several events)

  - likelihood > 2 (first guess)

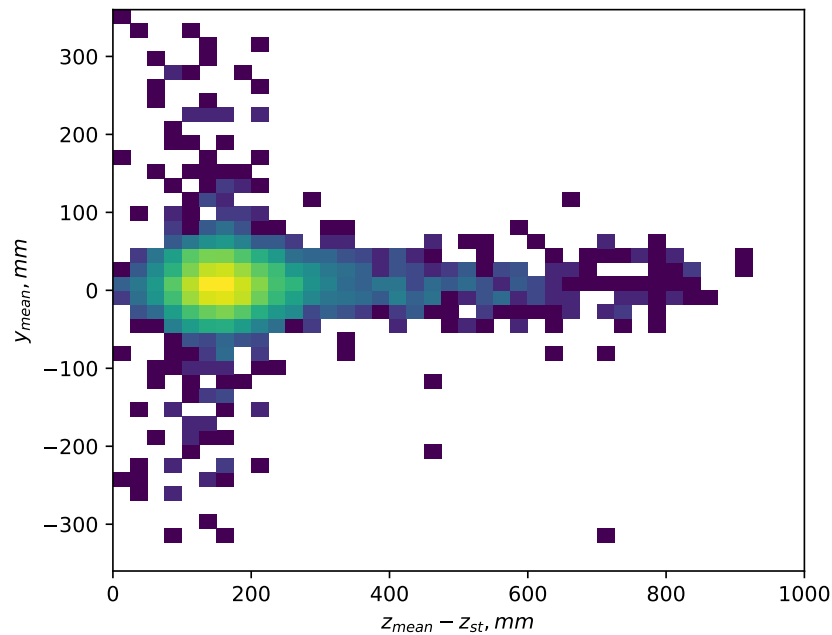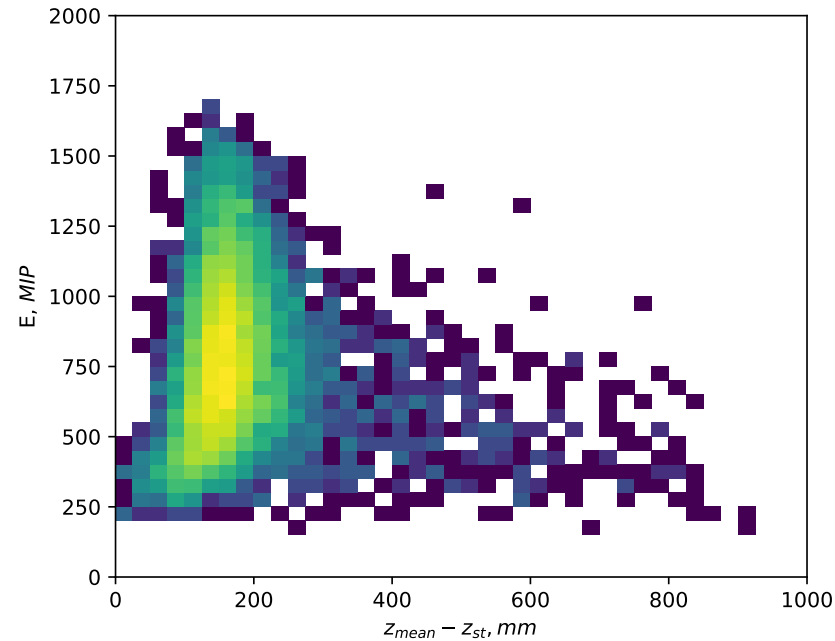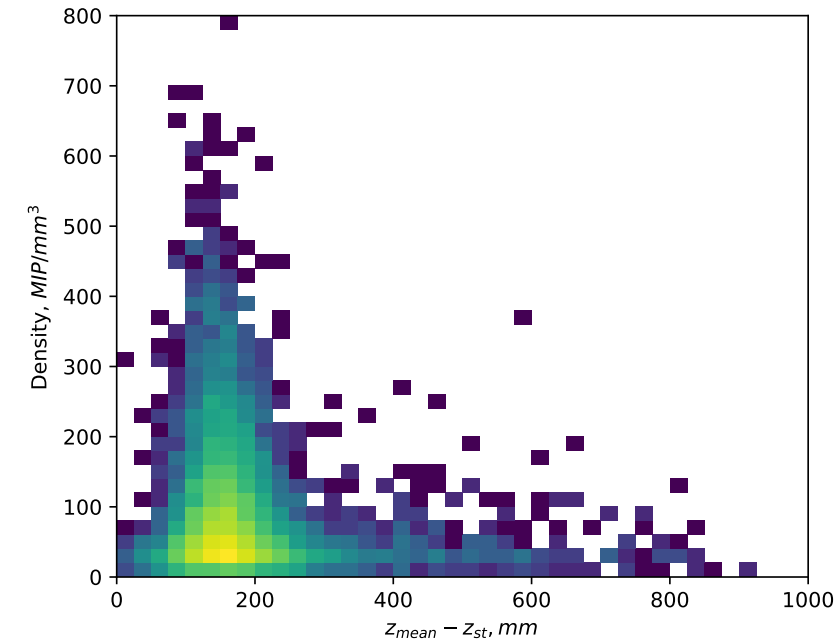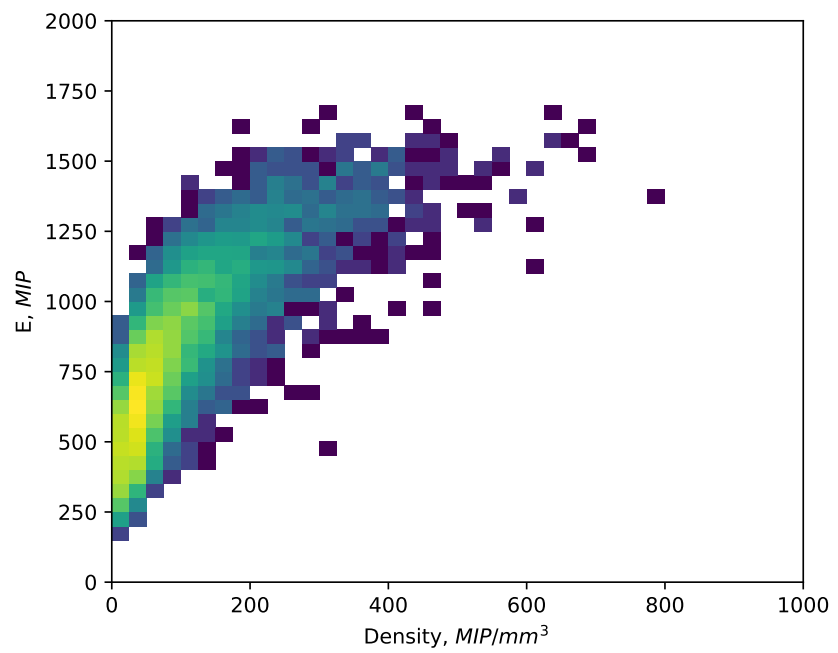  - energy density in ellipsoid [MIP/mm$^3$] > 20 (first guess)

# Main gaussian component (shower core)

## Some distributions for simulated 40 GeV pions. 10 Kevts.

- Reconstructed EM fraction.

- Shower start found

- Quality metrics (optimised on several events)

  - likelihood > 2 (first guess)

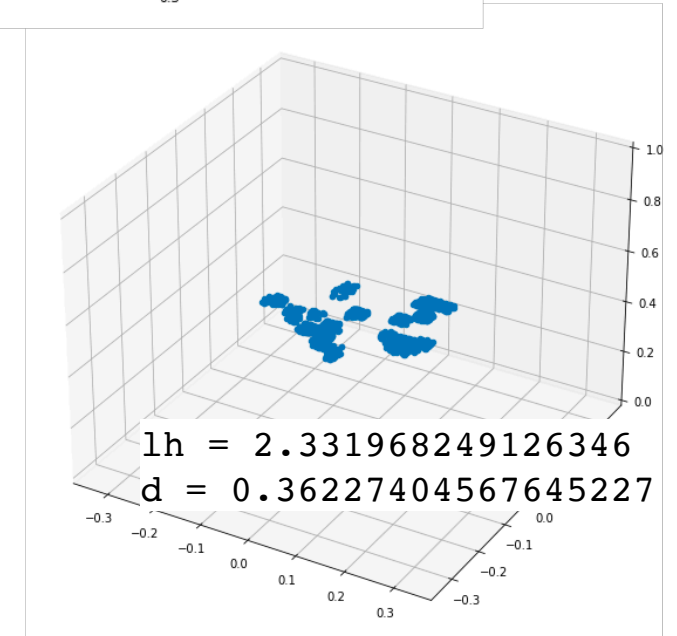  - energy density in ellipsoid [MIP/mm$^3$] > 20 (first guess)
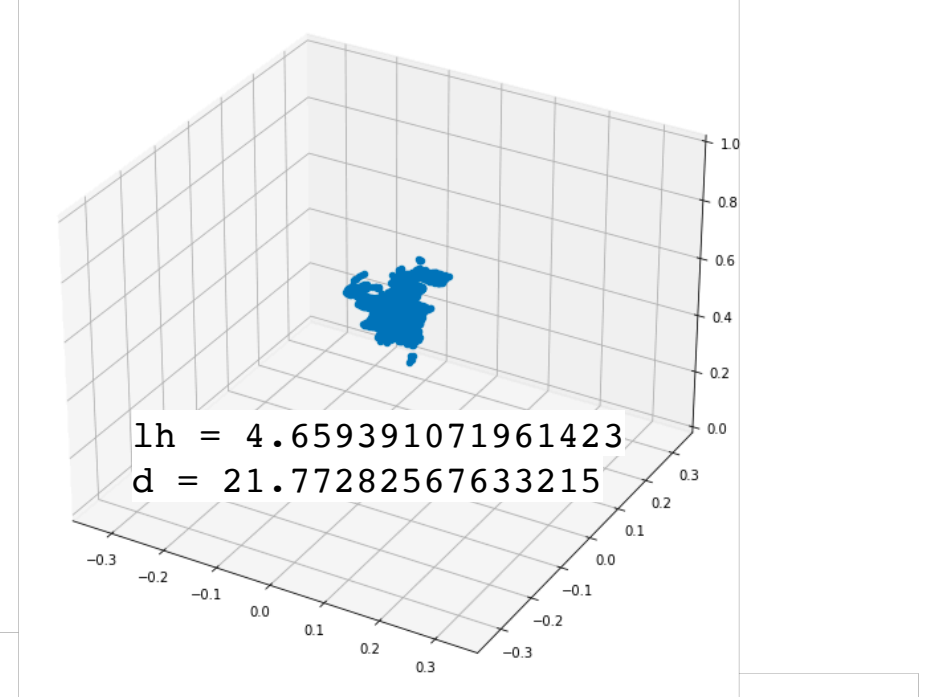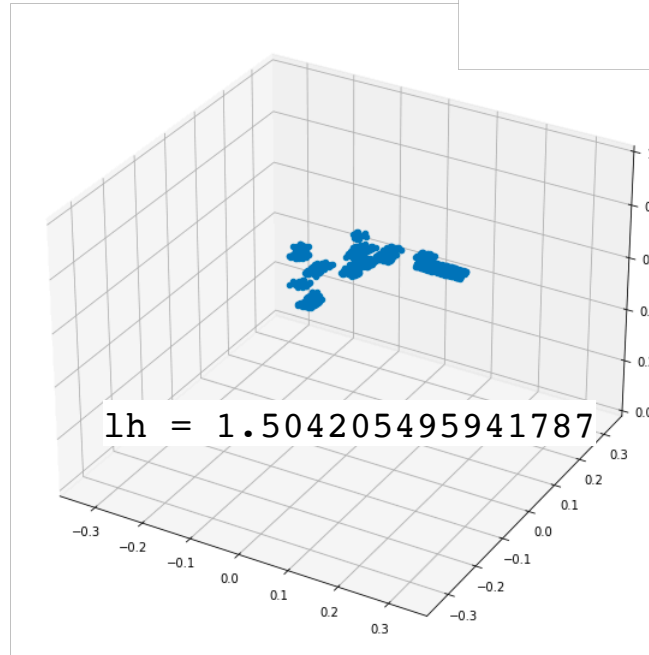


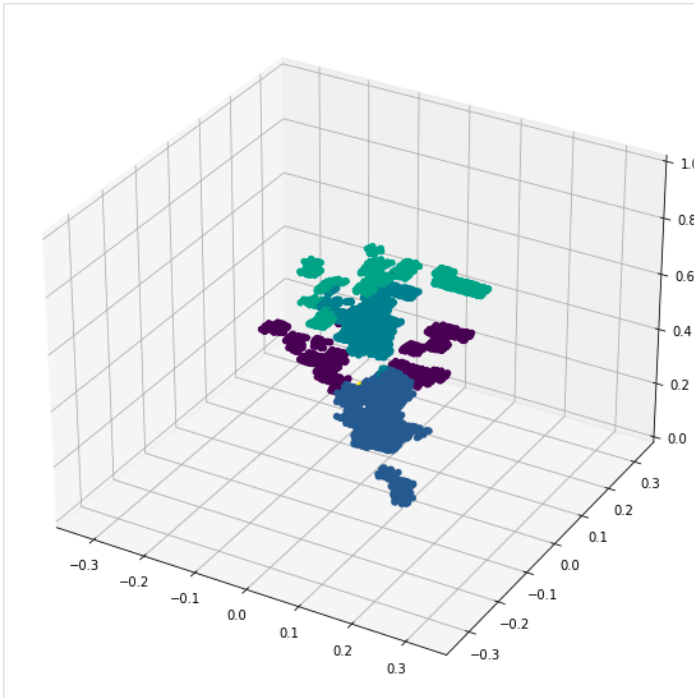Mean position in XY plane



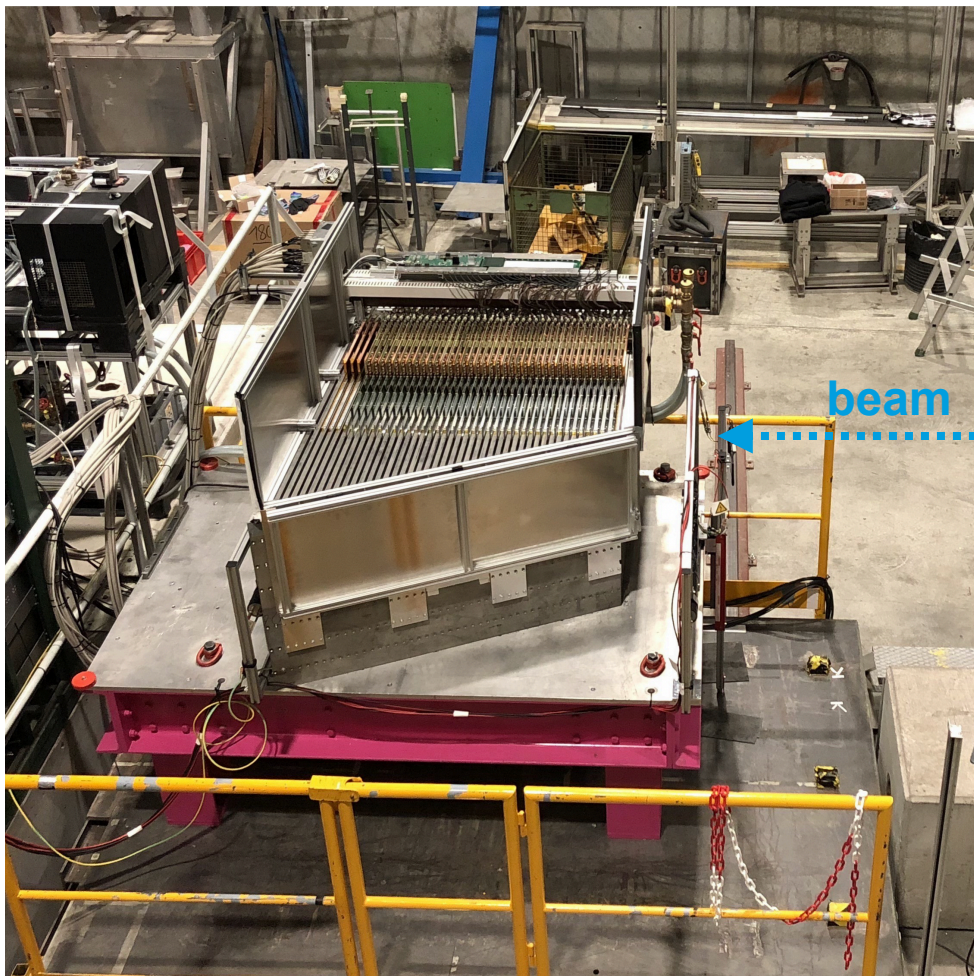Cluster energy vs Z mean



Cluster energy density vs Z mean

# Dealing with background clusters

- Quality metrics (optimised on several events)

  - likelihood > 2 (first guess)

  - energy density in ellipsoid [MIP/mm$^3$] > 20 (first guess)



lh = 4.659391071961423
d = 21.77282567633215





lh = 1.504205495941787



lh = 2.331968249126346
d = 0.36227404567645227

# CALICE AHCAL

## Test beam prototype.

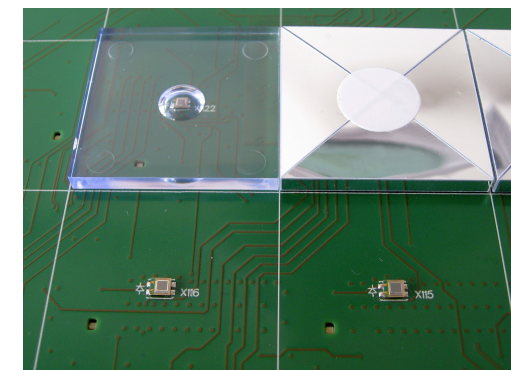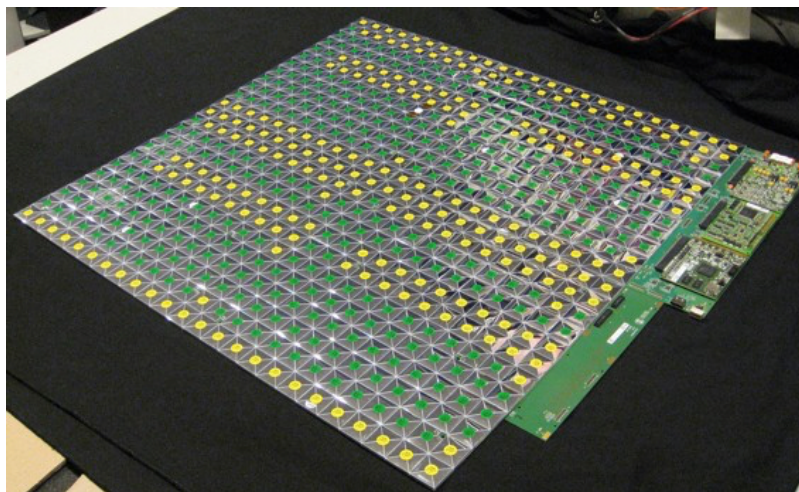

**39 active layers** of *24x24* scintillator tiles (*3x3 cm² each*) with individual SiPM readout. Active layers alternate with *~2 cm* steel absorber.

In total: **~22000 channels** (*<1‰ dead channels*), **~4 λ,  ~38X0**

Beam particles: muons, electrons, **pions**

Energy range: *10-200 GeV in 10-40 GeV steps*
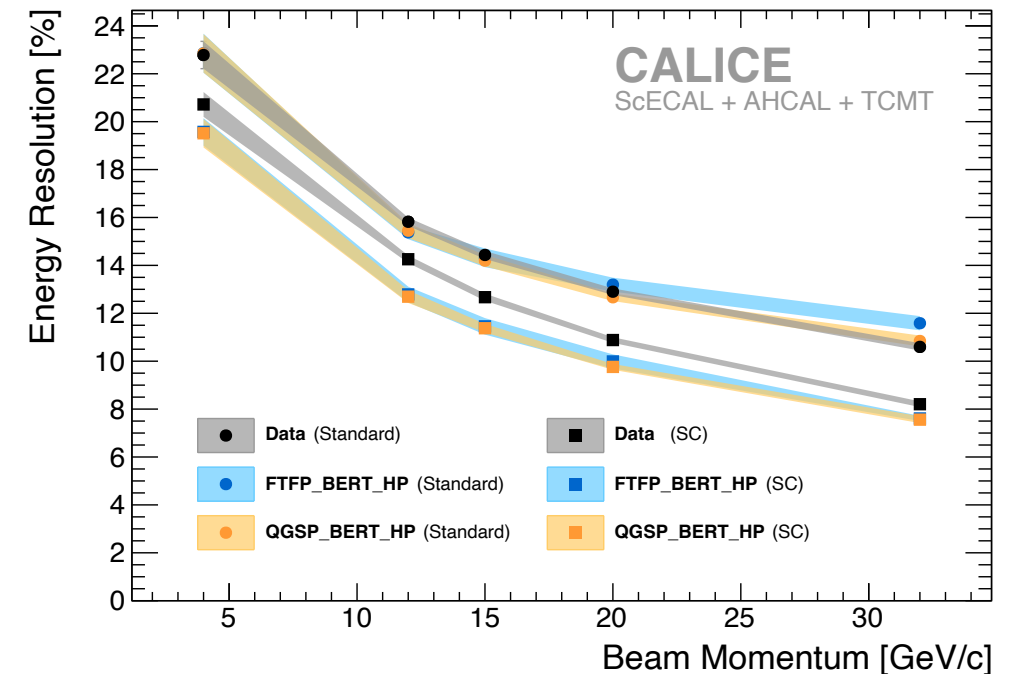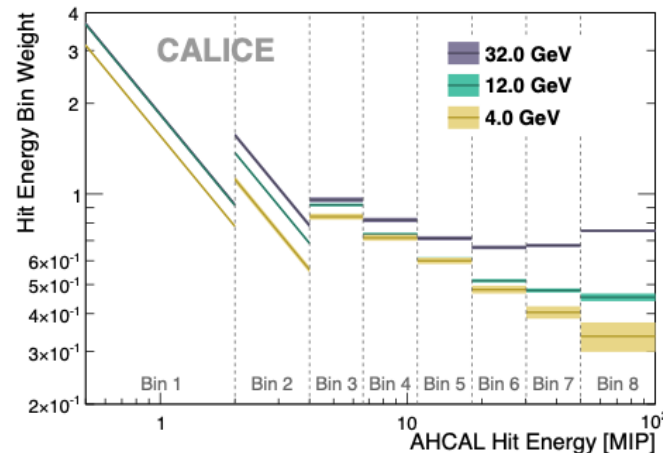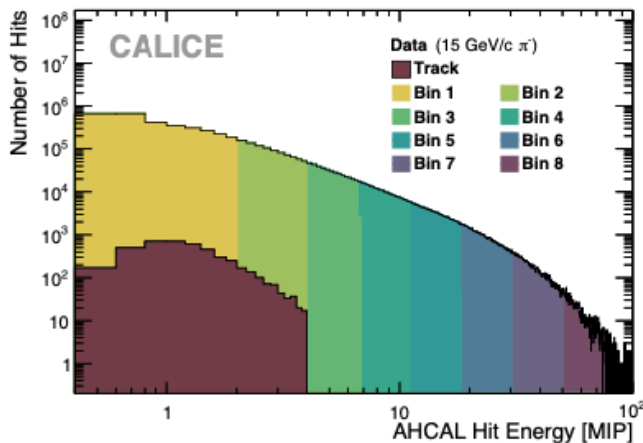
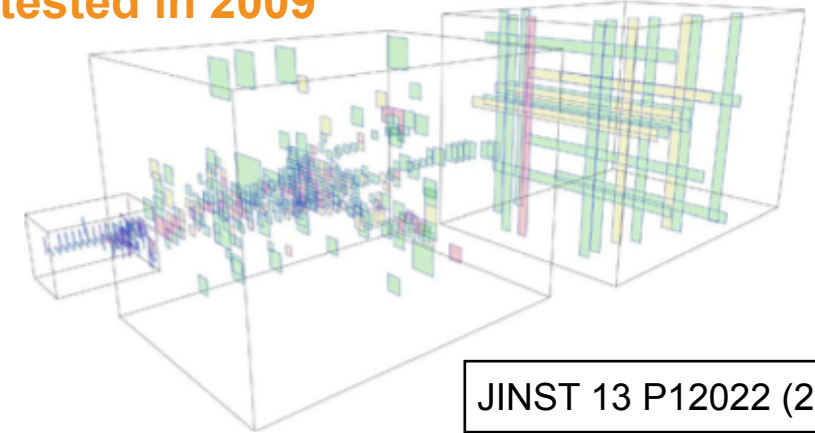*O(1M) hadron events per energy point*

# Software compensation method

## Example for CALICE combined setup ECAL+AHCAL+Tailcatcher tested in 2009
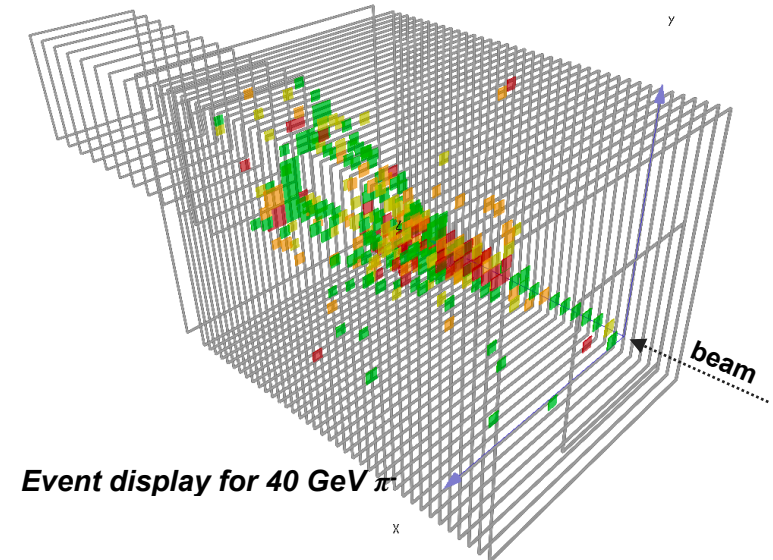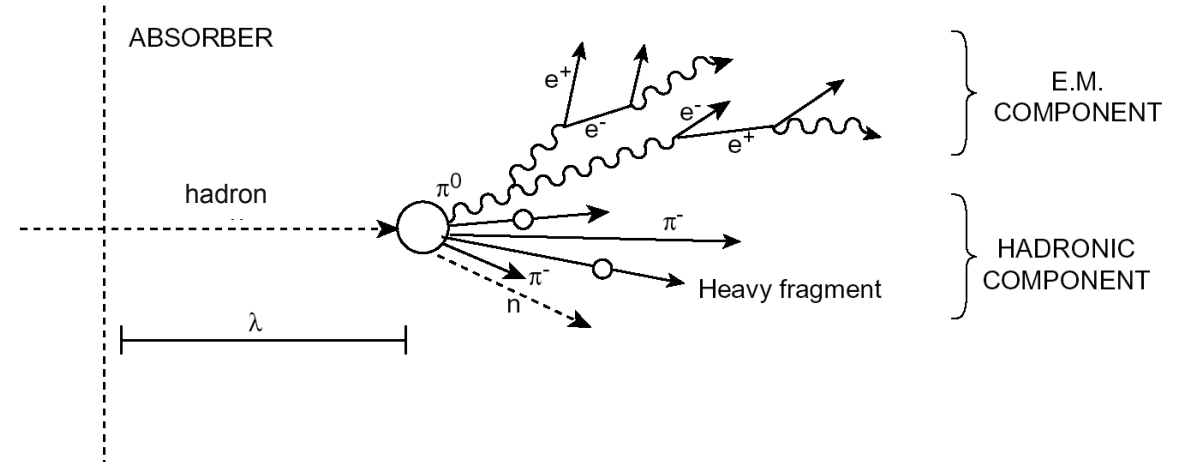
- h/e response compensation by assigning energy-dependent weights to hit energies (⇒local energy density)

  - Higher weights for **low energy hits** - dominated by **HAD** component
  - Lower weights for **high energy hits** - dominated by **EM** component

- 8 bins for hit energies

  - Polynomial fit to get energy dependent weight for each bin

- ➡ Energy resolution improvement 10-20%

- Disadvantages: limited to fit energy range, polynomial dependence has no physics motivation, additional topological information of hit context is not used
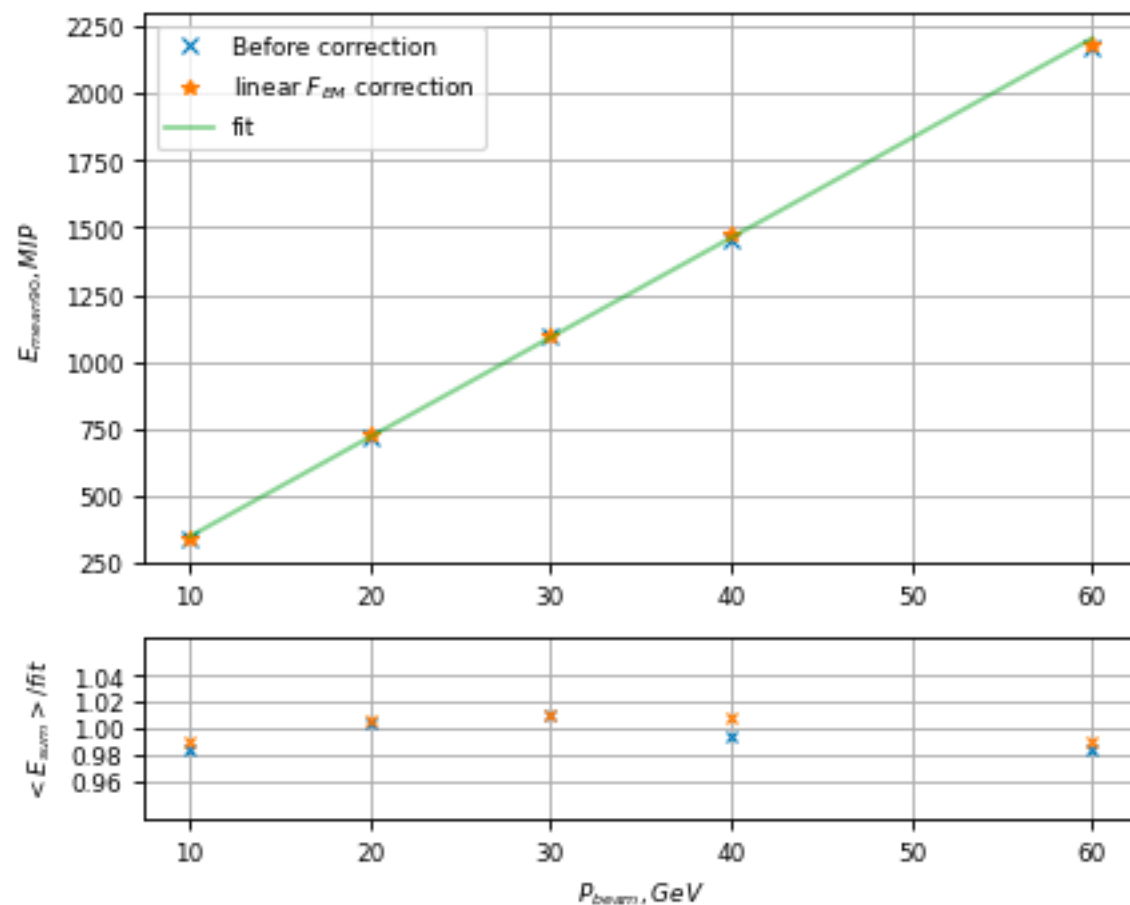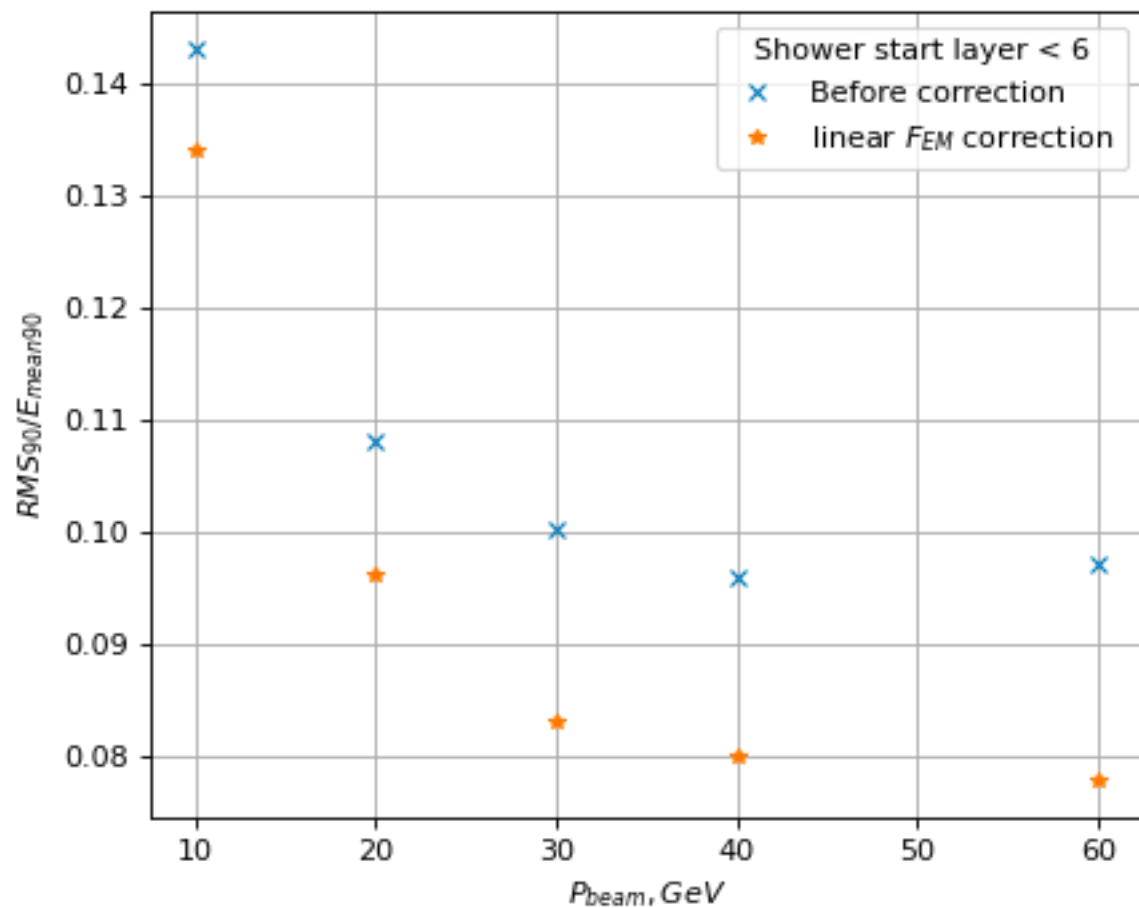
JINST 13 P12022 (2018)

# Hadronic showers

## General properties

- Hadronic shower development is rather complex:

  - Narrow EM core component from $\pi^0/\eta$

  - Surrounding halo dominated by charged hadrons

  - Large event-by-event fluctuation of EM/HAD ratio

  - Response to EM and HAD components is different in non-compensating calorimeters

  - Invisible energy as binding energy, nuclear recoil, neutrinos + late component

  ➡ Limited hadronic energy resolution

  ➡ Detailed simulation is challenging

- Highly granular calorimeter prototypes

  - Imaging capabilities provide detailed calorimetric images

  - Real test beam data for crosschecks and development of data-driven algorithms





*Event display for 40 GeV $\pi^-$*

- Linear $E_{sum}(F_{em})$ correction:

  $E_{cor} = C*E_{sum}$

  $C = <F_{em}>/(p_1 \cdot F_{em} + p_0)$

# Unified correction

## Getting $P_{beam}$-independent correction

Correction parameters as a function of $<E_{sum}>$:



- $p_0, p_1$ and $<F_{em}>$ are calculated for each event from the observed energy using resulting fits

  - More energy points need to be included to check the overfitting

  - Parameter uncertainties are not taken into account

  - Performance decrease for resolution ~3%