# k4Clue: Empowering Future Collider Experiments with CLUE

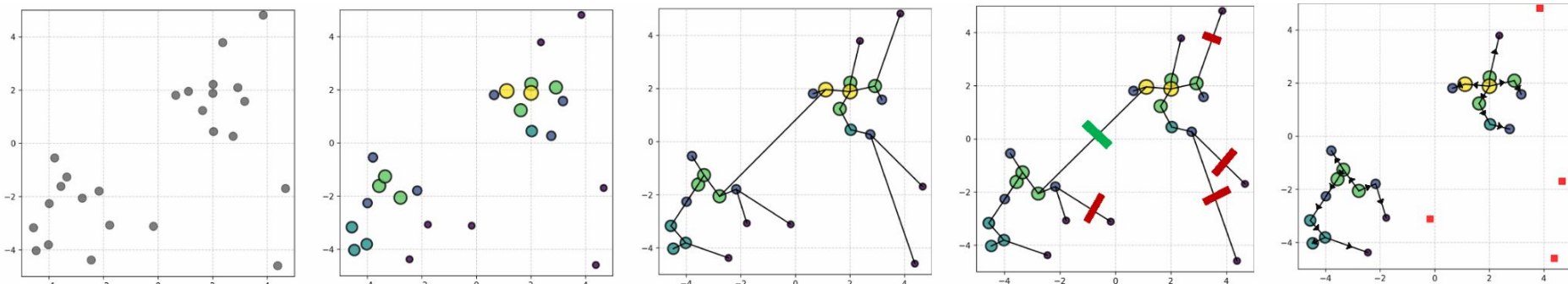**Erica Brondolin**,
Felice Pantaleo, Marco Rovere

# Introduction

# CLUstering of Energy

- CLUE (**CLUstering of Energy**) is a fast density-based clustering algorithm for the next generation of sampling calorimeter with high granularity in HEP
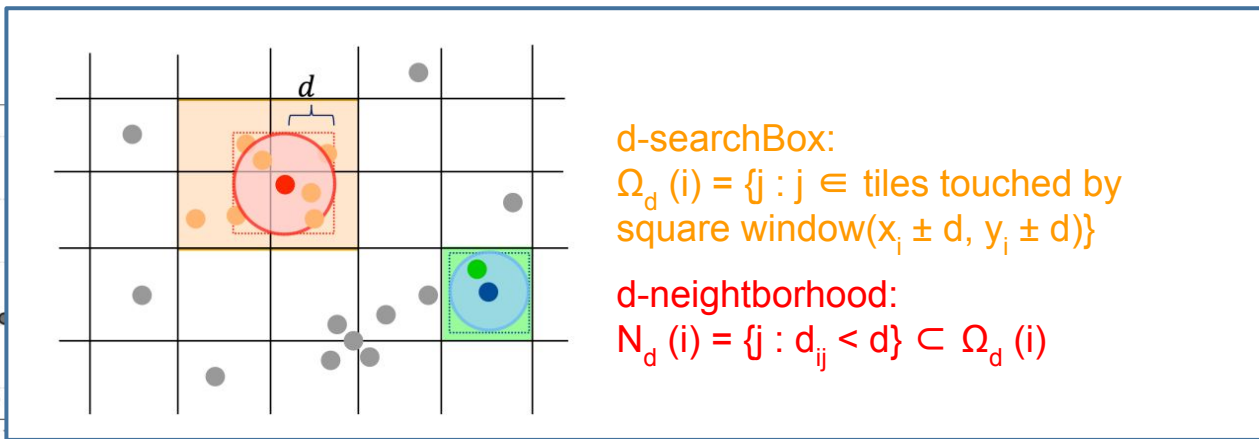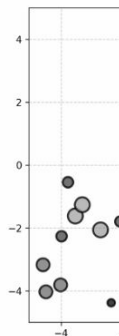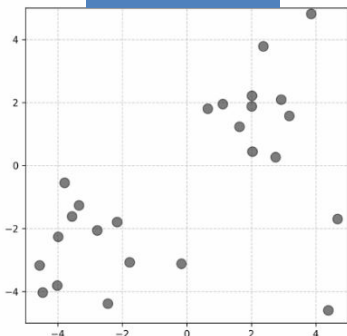
    **CLUE**  doi: 10.3389/fdata.2020.591315

- It uses **energy density** - rather than individual cell energy - to establish **seeds, outliers, and followers** in 2D planes.

- **GPU-friendly**, i.e. suitable for the upcoming era of heterogeneous computing in HEP

- Standalone repo:  CLUE - kalos
    gitlab.cern.ch

# Step 1: Building Data Structure

- Querying neighborhood is a frequent operation in density-based clustering → **fast!**

- Build **Fixed-Grid Spatial Index** for hits on each layer:

  - Each tile in the grid hosts indices of hits inside it and has a fixed length of memory to store the hosted indices. It is independent by the detector granularity.

- To find the neighborhood hits $N_d(i)$ of $i$-hit, we only need to loop over hits in $\Omega_d(i)$



build data structure

d-searchBox:
$\Omega_d(i) = \{j : j \in$ tiles touched by square window$(x_i \pm d, y_i \pm d)\}$

d-neightborhood:
$N_d(i) = \{j : d_{ij} < d\} \subset \Omega_d(i)$

# Step 2: Local energy density

- Calculate local energy density ($\rho_i$) in a distance ($d_c$)

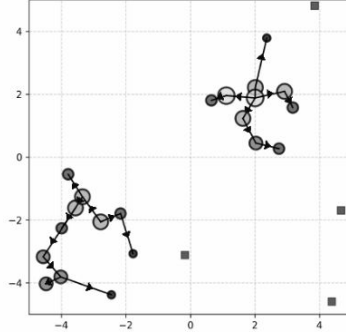  - Each hit $j$ weighted by the deposited energy ($E_j$)
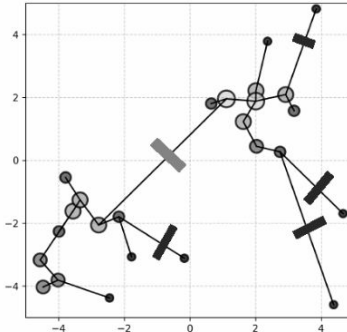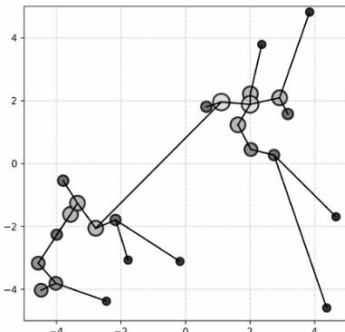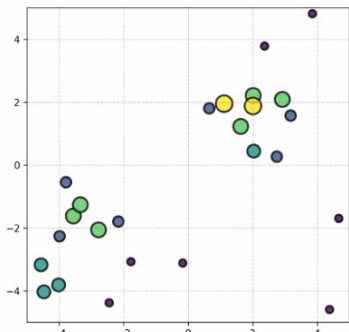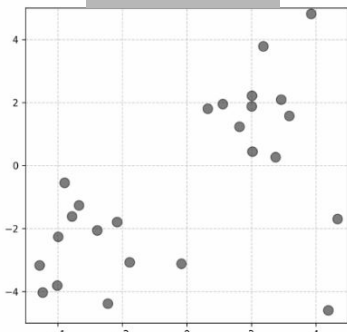
  - For each hit, calculate $\rho_i$

convolution kernel

k = 0.5

$$\rho_i = \sum_{j \in N_d(i)} E_j \times f(d_{ij}); \quad f(d_{ij}) = \begin{cases} 1, & \text{if } i = j \\ k, & \text{if } 0 < d_{ij} \leq d_c \\ 0, & \text{if } d_{ij} > d_c \end{cases}$$

**build data structure**

**density**

Erica Brondolin (erica.brondolin@cern.ch)

# Step 3: Find "closest higher hit"

- Calculate "Nearest-Higher" hit within $N_{dm}(i)$

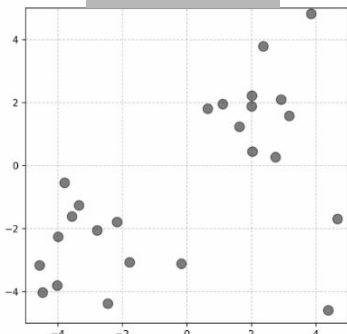  - Define $d_m = o_f * d_c$

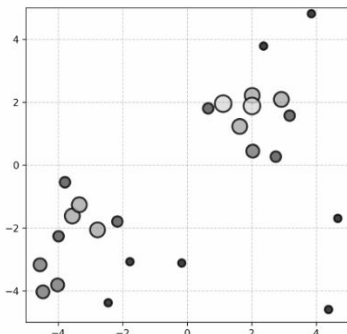  - Find the closest hit with higher local energy density, $nh_i$

$$nh_i = \begin{cases} argmin_{j \in \hat{N}_{d_m}(i)} d_{ij}, \text{if } |\hat{N}_{d_m}| \neq 0, \hat{N}_{d_m}(i) = \{j : j \in N_{d_m}(i), \rho_j > \rho_i\} \\ -1, \text{otherwise} \end{cases}$$
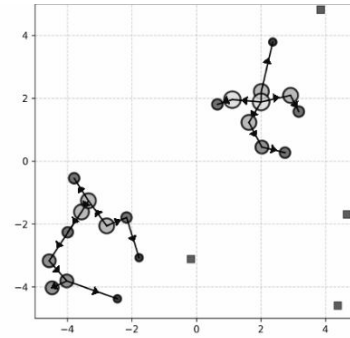
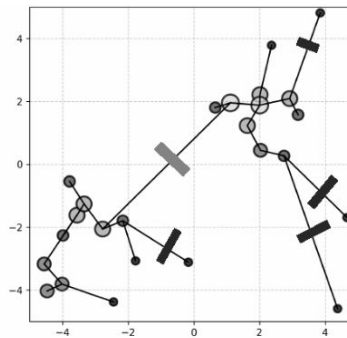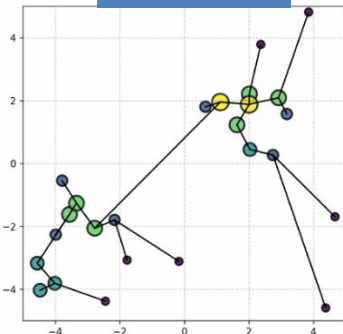  - Calculate the separation distance $\delta_i = dist(i, nh_i)$



build data structure
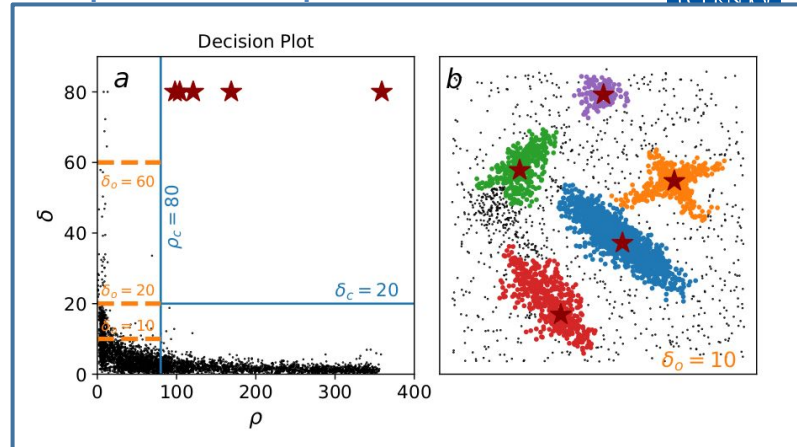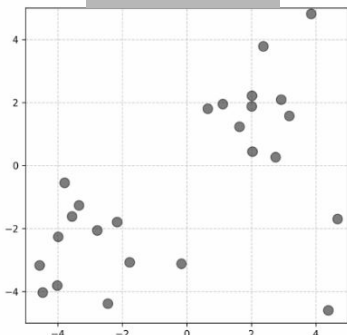
density

nearest higher

# Step 4: Classify hits

- Promote as **seed** if $\rho_i > \rho_c$ , $\delta_i > d_c$

- Demote as **outlier** if $\rho_i < \rho_c$ , $\delta_i > o_f * d_c$

- Assign unique, progressive cluster ID to each cluster

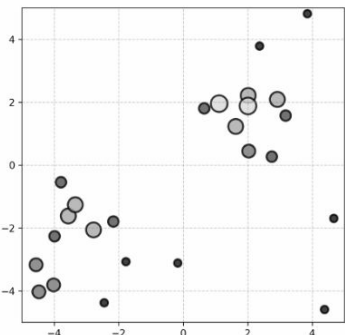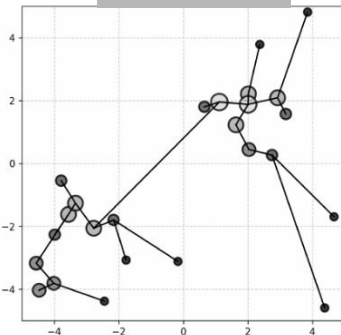  - **Followers** are defined and associated to their closest seed
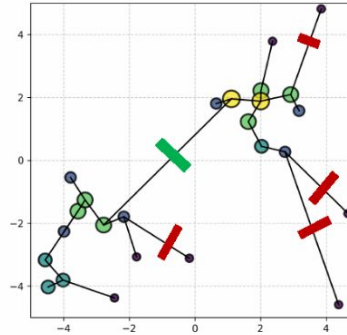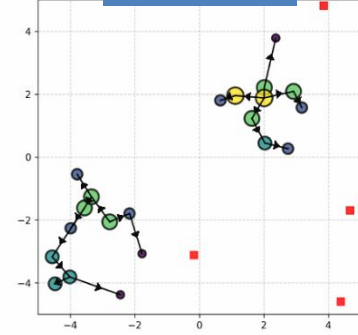


Decision Plot





build data structure

density

nearest higher

find seed

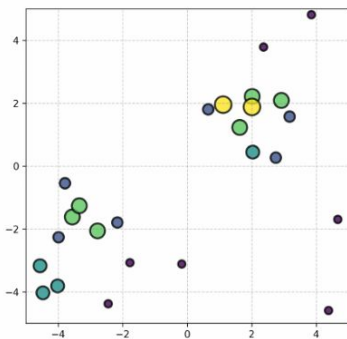assign clusters

# Clustering procedure recap

**build data structure**

Query the neighborhood of a point by looping over the points in $N_d$ in the bins touched by the tiles intersected by $d_c$
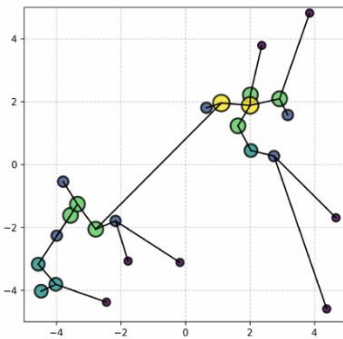
**density**

Hit position and energy used to calculate the hit's local energy density $\rho_i$ and its distance $\delta_i$ to the nearest hit with higher local density
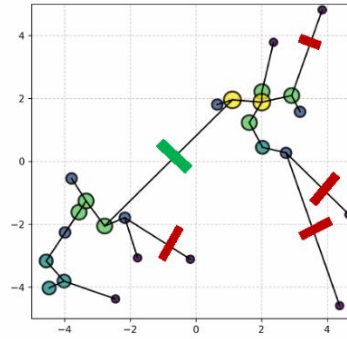
**nearest higher**

Define the nearest-higher of each hit as the hit with the local energy density higher then the hits itself and within a distance of $d_m = o_f \times d_c$
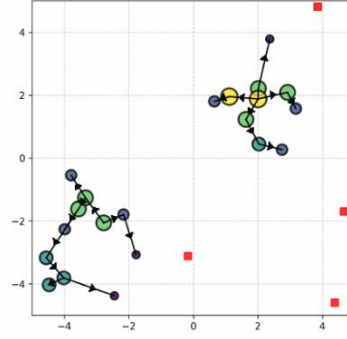
**find seed**

Use following criteria:
• seed: $\rho_i \geq \rho_c$ and $\delta_i \geq d_c$;
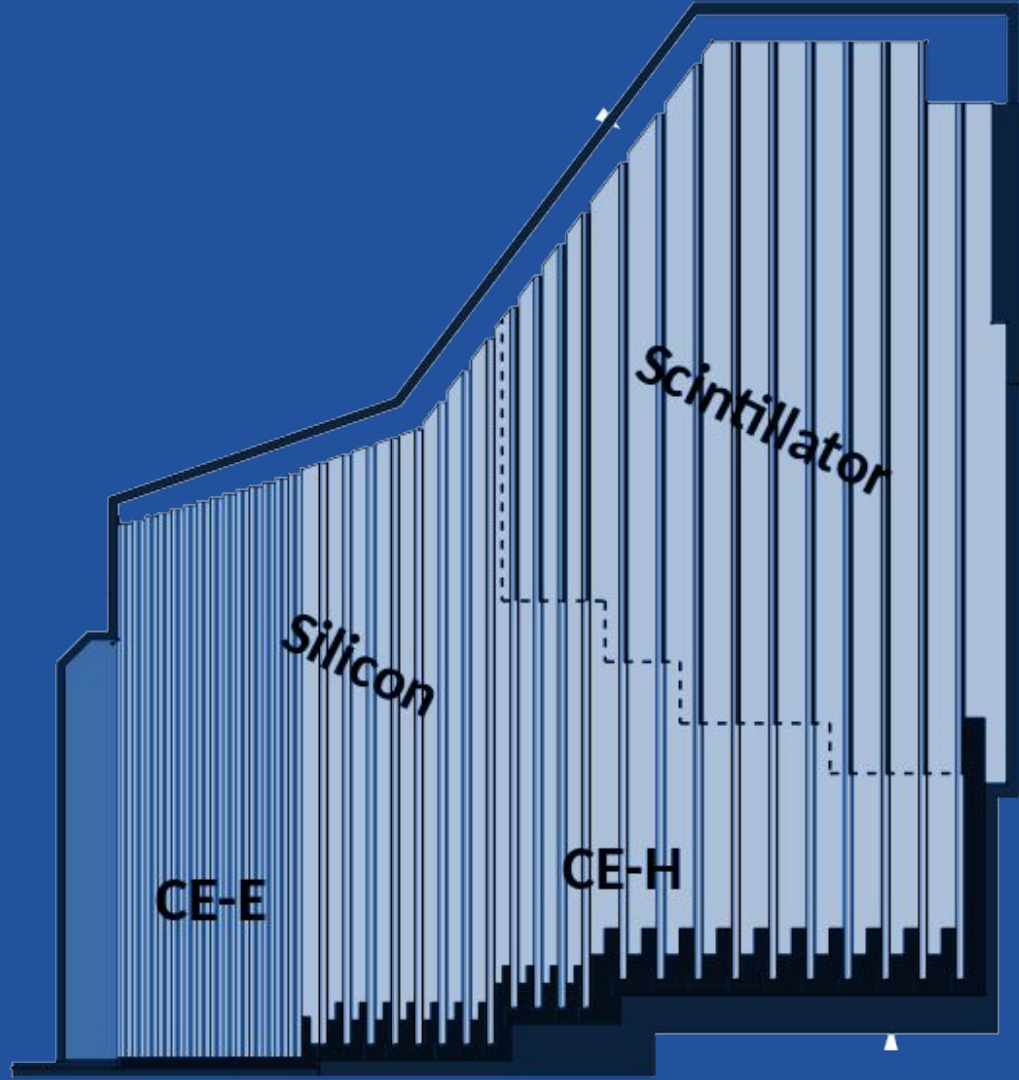• outlier: $\rho_i < \rho_c$ and
$\delta_i \geq (o_f \times d_c)$

**assign clusters**

Register each remaining point as a follower to its nearest-higher

CLUE in the HGCAL reconstruction

Silicon

Scintillator

CE-E

CE-H

# CMS High Granularity Calorimeter
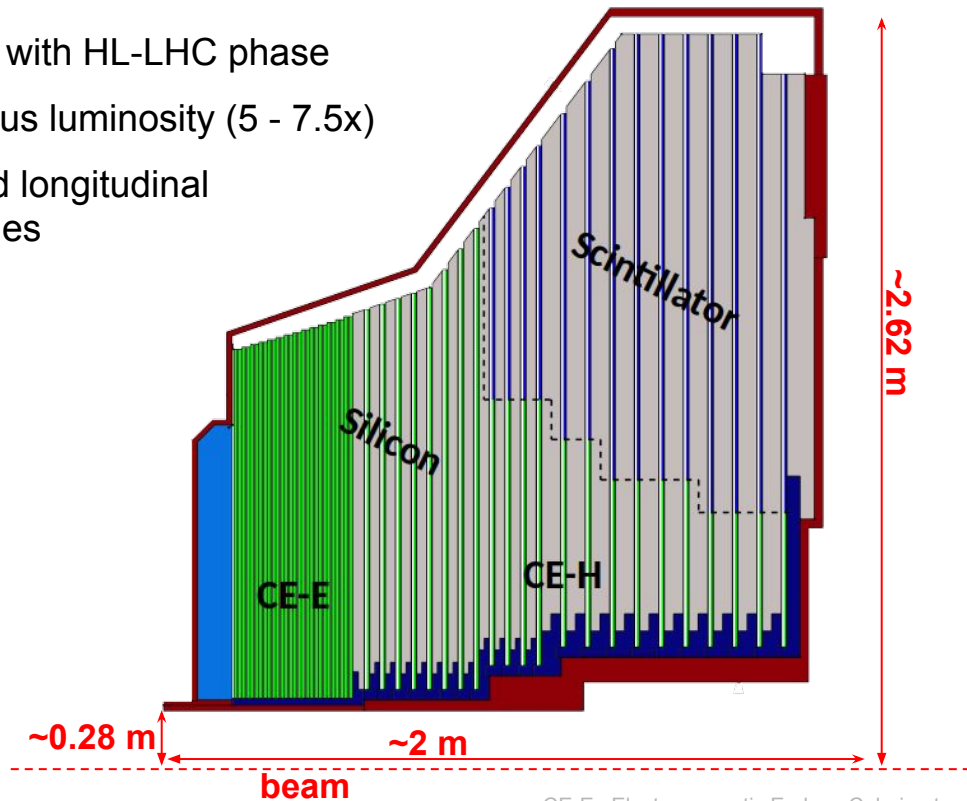
- Phase-2 upgrade of CMS is needed to cope with HL-LHC phase
  - A significant increase in the instantaneous luminosity (5 - 7.5x)
- Imaging calorimeter with very fine lateral and longitudinal segmentation, and precision timing capabilities
  - Covering 1.5 < η < 3.0

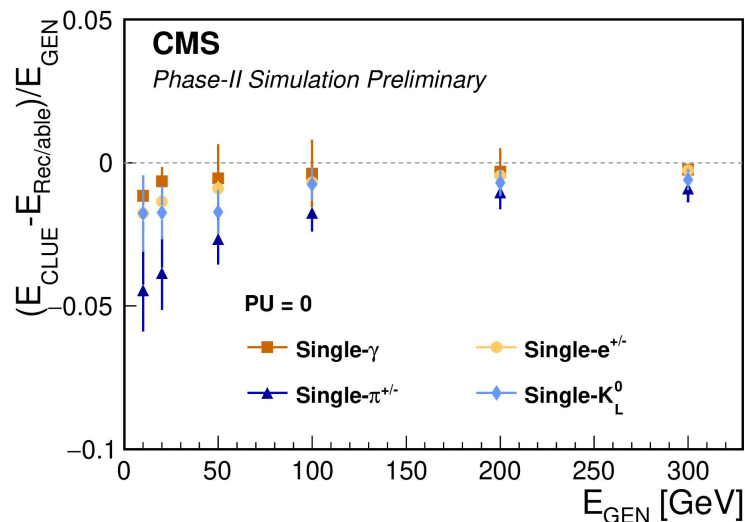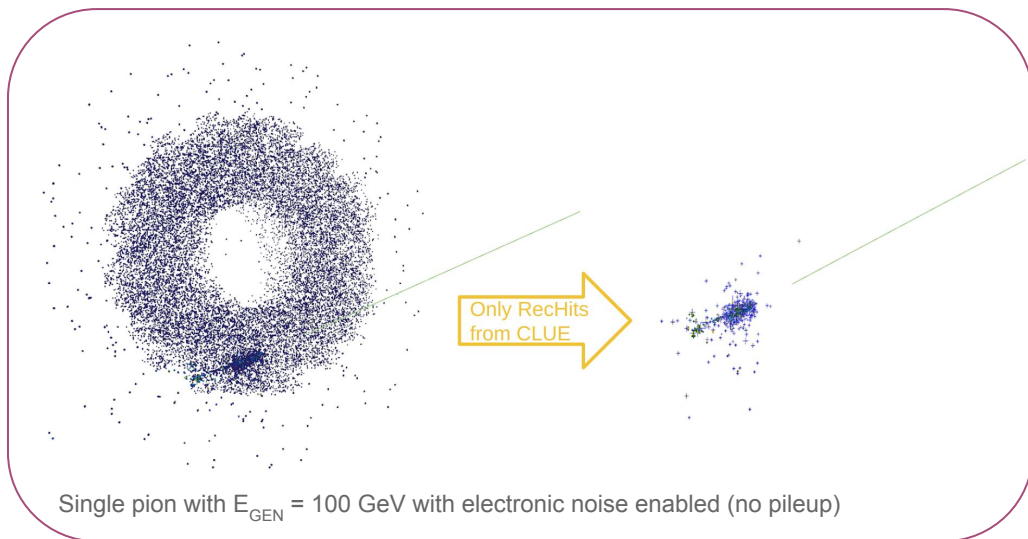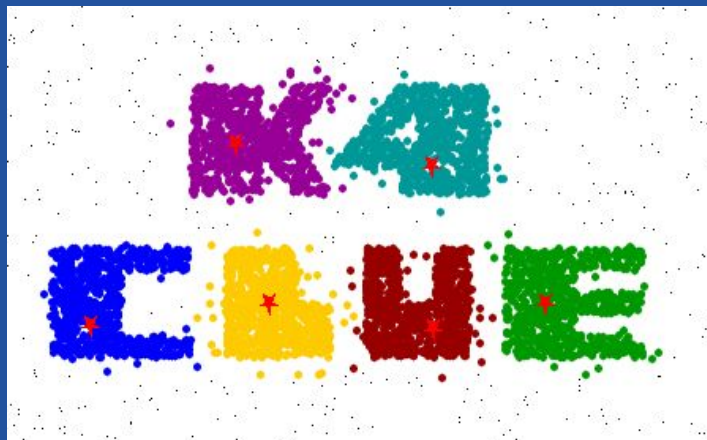| Both endcaps | Silicon | Scintillators |
|---|---|---|
| Area | ~620 m$^2$ | ~400 m$^2$ |
| Channel size | 0.5 - 1 cm$^2$ | 4 - 30 cm$^2$ |
| #Modules | ~30'000 | ~4'000 |
| #Channels | ~6 M | 240 k |
| Op. temp. | -30 °C | -30 °C |

Ref.



CE-E : Electromagnetic Endcap Calorimeter
CE-H : Hadronic Endcap Calorimeter

# HGCAL Software Reconstruction

- The HGCAL reconstruction framework is **TICL (The Iterative Clustering)**

- It starts by calibrating deposited energy in individual cells, also called RecHits → an order of **$10^5$ RecHits** in the HGCAL detector for events @ 200 pileup

- CLUE clusters the RecHits in the same layer to produce **Layer Clusters (LCs)**
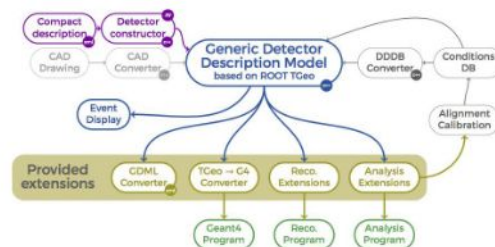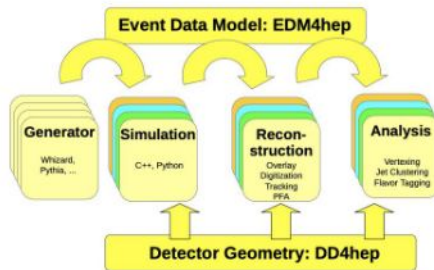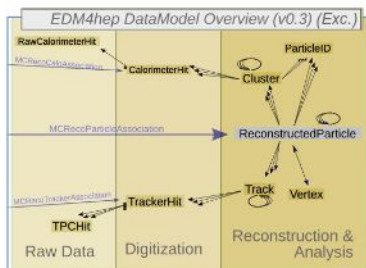


Single pion with $E_{GEN}$ = 100 GeV with electronic noise enabled (no pileup)

Only RecHits from CLUE



**CMS**

*Phase-II Simulation Preliminary*

PU = 0

- Single-$\gamma$
- Single-e$^{+/-}$
- Single-$\pi^{+/-}$
- Single-K$_L^0$

$(E_{CLUE}-E_{Rec/able})/E_{GEN}$

$E_{GEN}$ [GeV]

The package

# Key4hep in a nutshell

# Integrating CLUE in Key4hep

- **k4Clue v01-00** (doi: 10.5281/zenodo.7851995)
  - It's adapted to the common event data model, `EDM4hep`
  - It includes a wrapper class to run in the Gaudi software framework
  - It's included in the new Key4hep releases managed by Spack

# Additional features w.r.t. [kalos/Clue](#)

- **Cluster hits in the entire 4π detector region**
  - Definition of the tessellated space (`LayerTile`) in the standalone version defines coordinates and searches only in the transverse plane
  - Modified basic structure of the `LayerTile` and the search algorithm to to allow for the definition of a cylindrical surface

    $x \rightarrow r\Phi$         $y \rightarrow z$
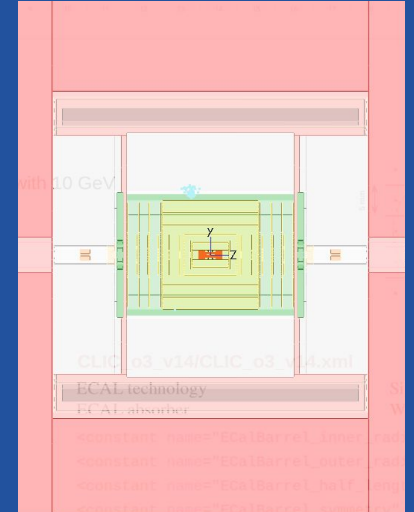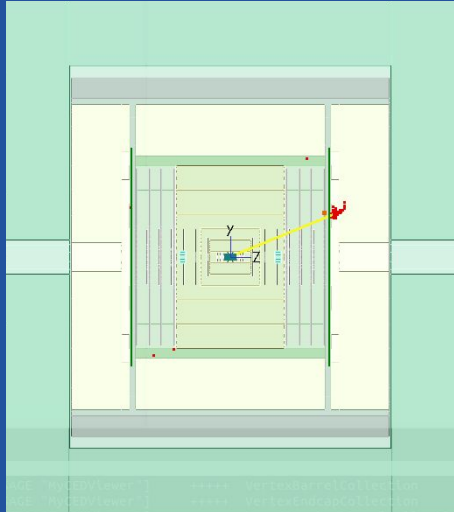
- **Template CLUE algorithm classes**
  - To allow the possibility of defining several different calorimeter layouts
  - A dedicated documentation page in the package ([`include/readme.md`](#)) allows the user to follow a simple but detailed step-by-step procedure to introduce and test the preferred layout.

- **GitHub CI & EDM4hep Validation**
  - `edm4hep:CLUECalorimeterHit : CalorimeterHit` class with specific methods related to the CLUE algorithm
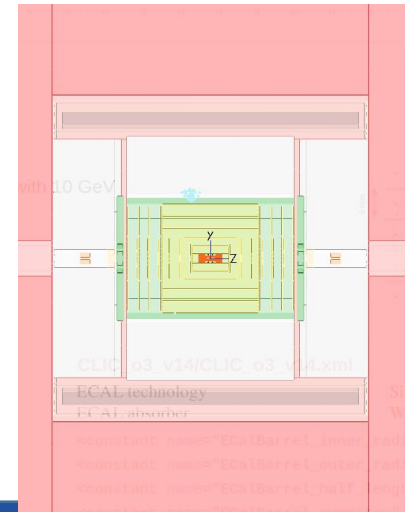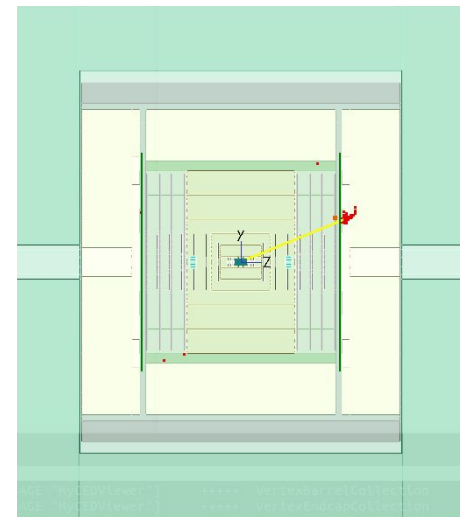
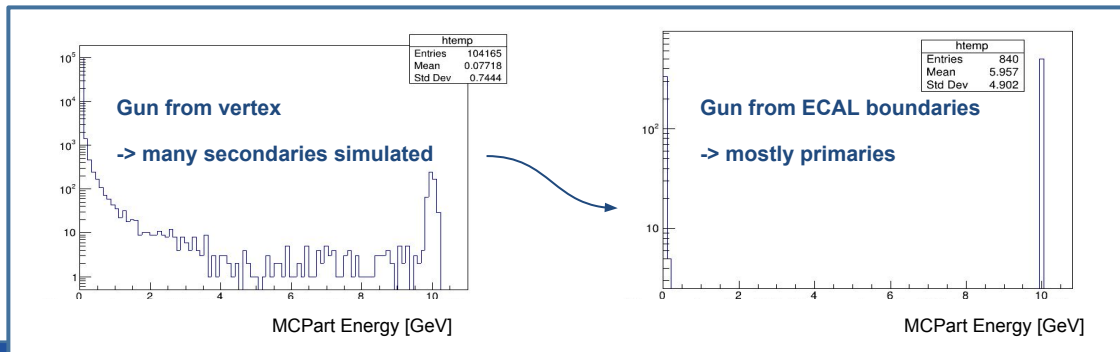# Performance evaluation
## - CLD & CLICdet -

# ECAL of CLICdet & CLD

- 40 layers of 5x5 mm$^2$ Silicon cells & W

- The main difference between the two calorimeters lies in the layout parameters → To compensate for a lower detector solenoid field, the **CLD design** starts from a larger radius both in the barrel and in the endcap region w.r.t. **CLICdet**.
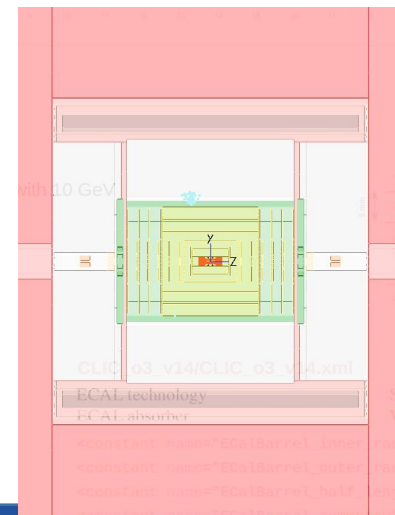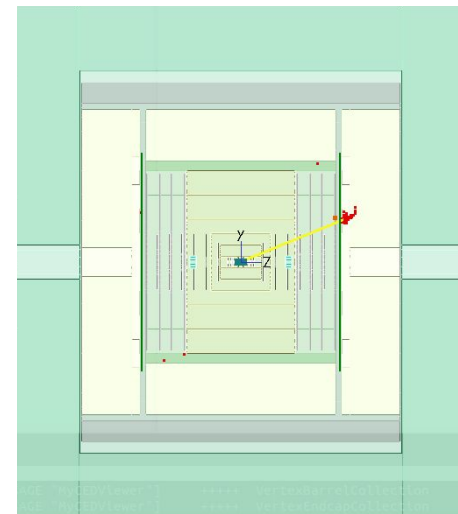  - Further details in backup

# ECAL of CLICdet & CLD

- 40 layers of 5x5 mm$^2$ Silicon cells & W

- The main difference between the two calorimeters lies in the layout parameters → To compensate for a lower detector solenoid field, the **CLD design** starts from a larger radius both in the barrel and in the endcap region w.r.t. **CLICdet**.
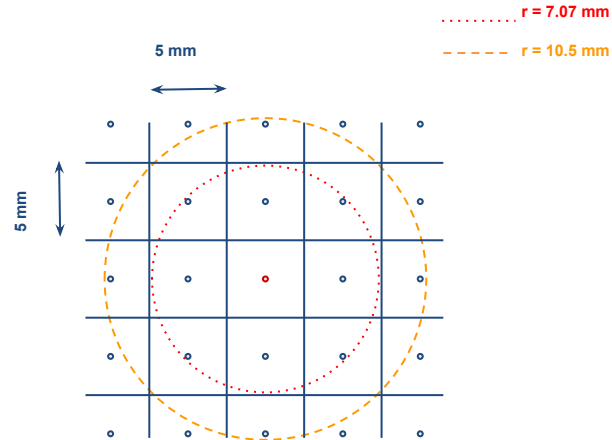  - Further details in backup

- 500 events of single gamma at 10 GeV generated perpendicular to the surface with Geant4 General Particle Source
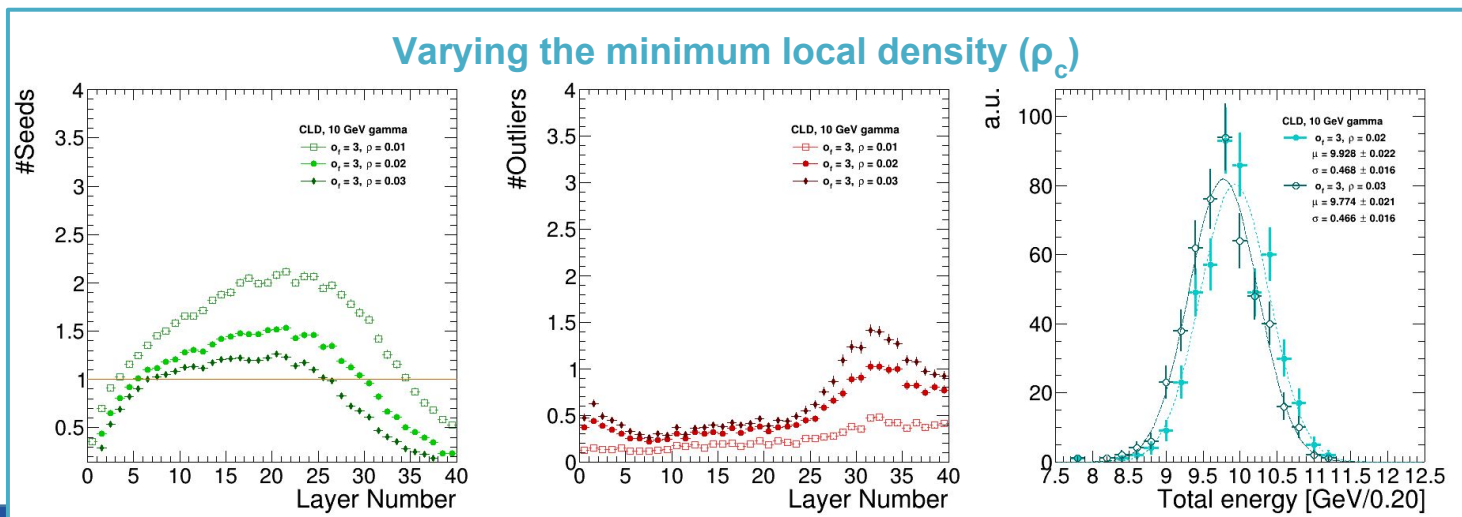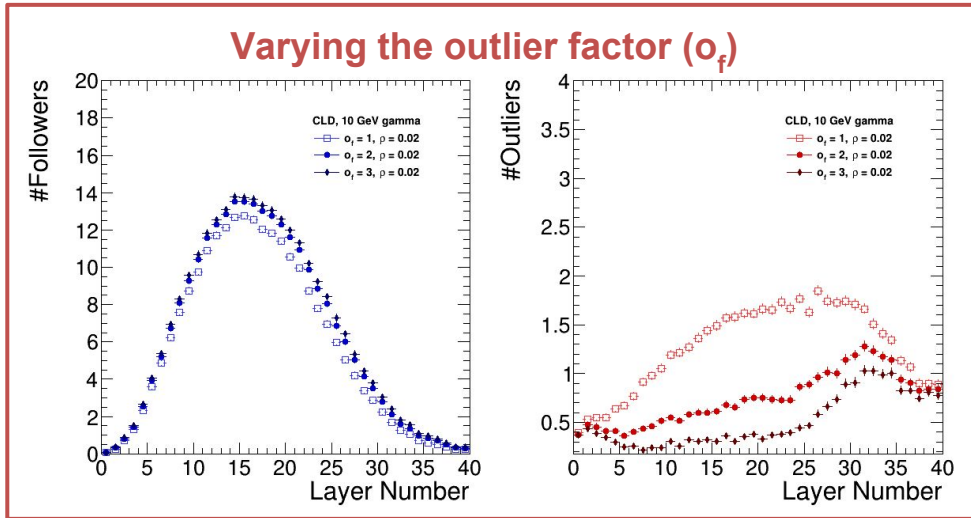  - Main reason: no conversion in the tracker volume

# Parameters tuning
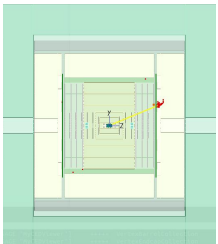
- Input parameters tuned for CLD

- Same ones tested also for CLICdet **(similar geometry, same granularity)**

- Critical Distance ($d_c$) is established by geometry granularity to contain (minimum) the close neighbors cells:

  - $d_c$ = 15 mm

# Parameters tuning



**Varying the outlier factor ($o_f$)**



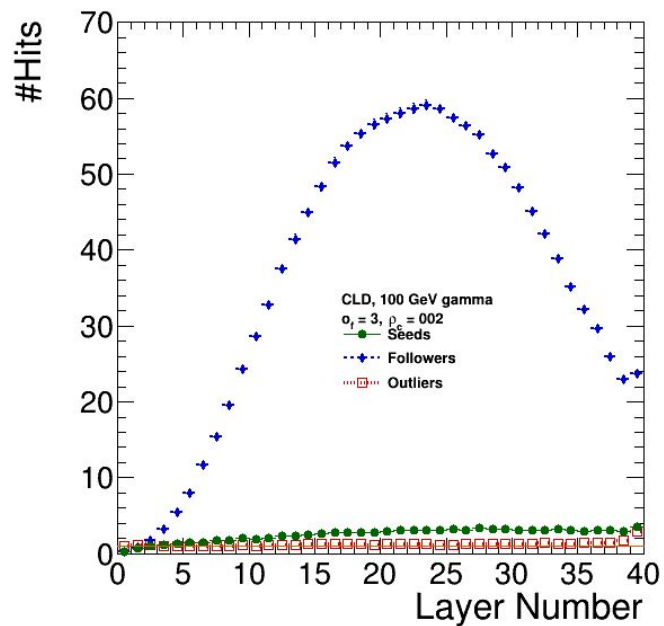**Varying the minimum local density ($\rho_c$)**

Erica Brondolin (erica.brondolin@cern.ch)

# Higher energies
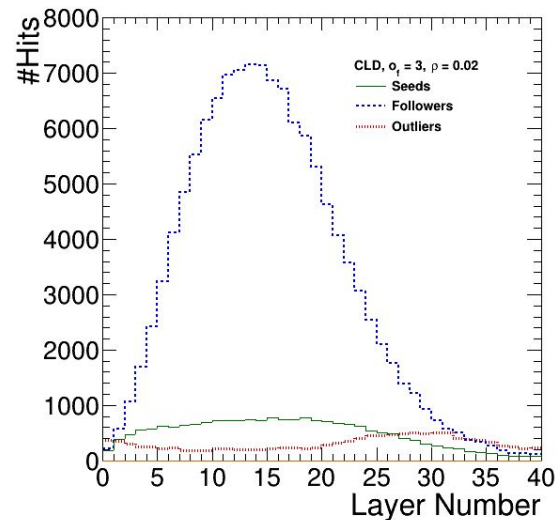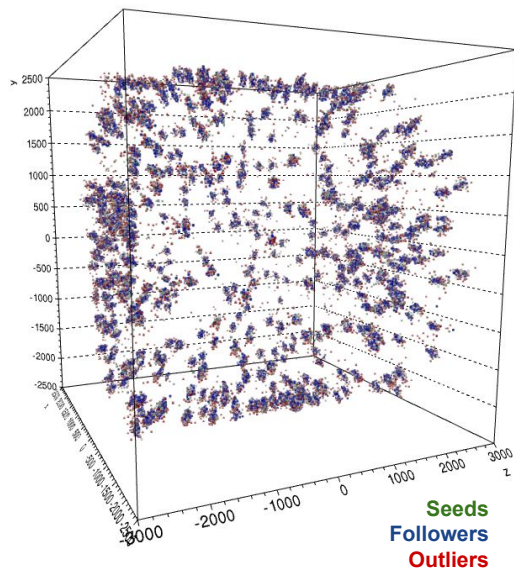
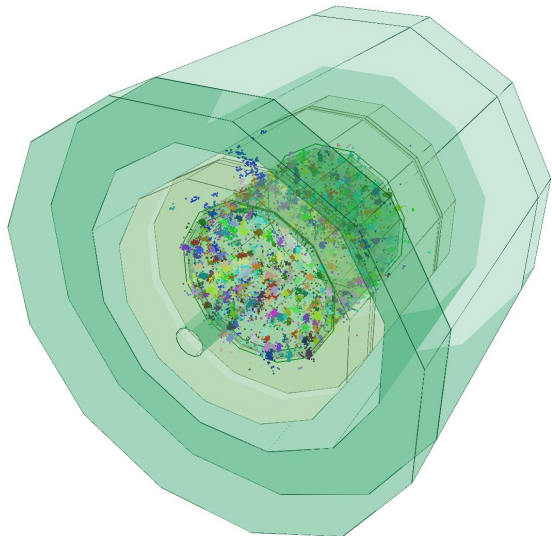- 500 events with single gamma (from ECAL surface) at <mark>100 GeV</mark>

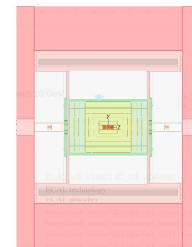- $d_c$ = 15.00, $\rho_c$ = 0.02, $o_f$ = 3.0
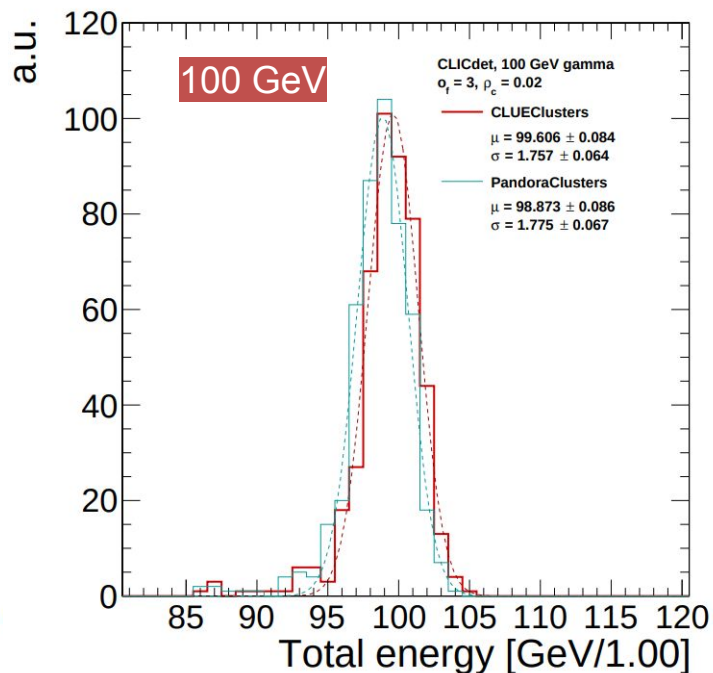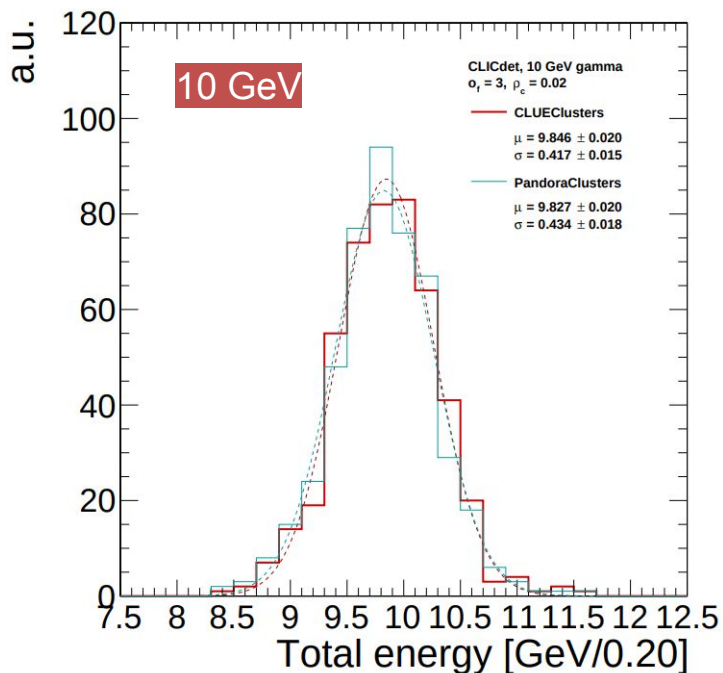
# Multiple gamma event

- Produced with normal gun, i.e. particles generated from vertex

- 1 event with 500 single gammas each produced with 10 GeV
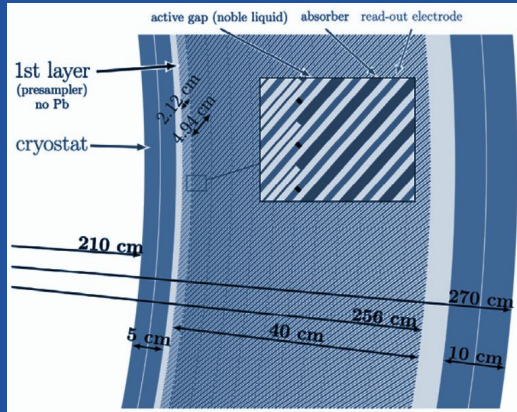  Only simulated calorimeter hits are shown



Seeds
Followers
Outliers



CLD, $o_t = 3$, $\rho = 0.02$
— Seeds
···· Followers
···· Outliers

# CLICdet results

- Using same input parameters selected for CLD



CLICdet, 10 GeV gamma
$o_r = 3$, $\rho_c = 0.02$

CLUEClusters
$\mu = 9.846 \pm 0.020$
$\sigma = 0.417 \pm 0.015$

PandoraClusters
$\mu = 9.827 \pm 0.020$
$\sigma = 0.434 \pm 0.018$

CLICdet, 100 GeV gamma
$o_r = 3$, $\rho_c = 0.02$

CLUEClusters
$\mu = 99.606 \pm 0.084$
$\sigma = 1.757 \pm 0.064$

PandoraClusters
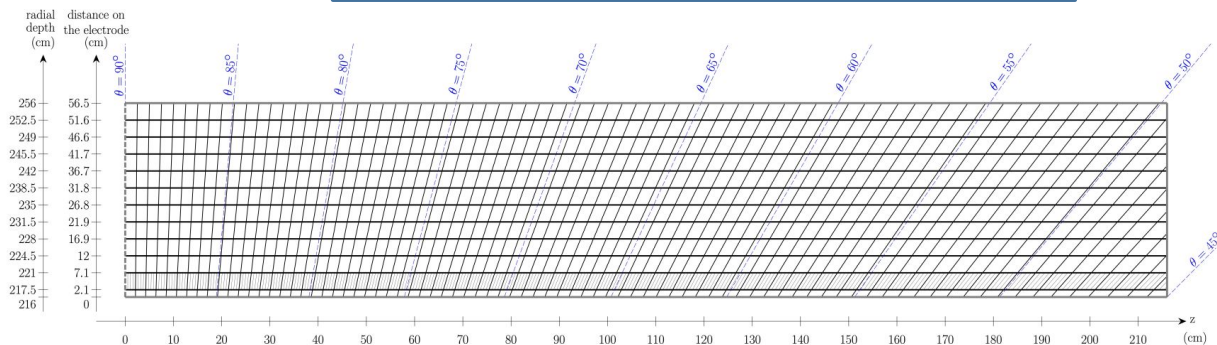$\mu = 98.873 \pm 0.086$
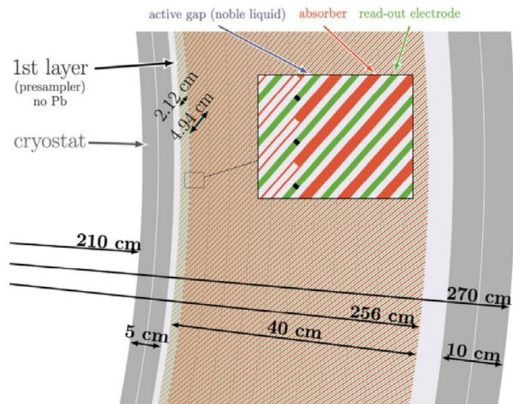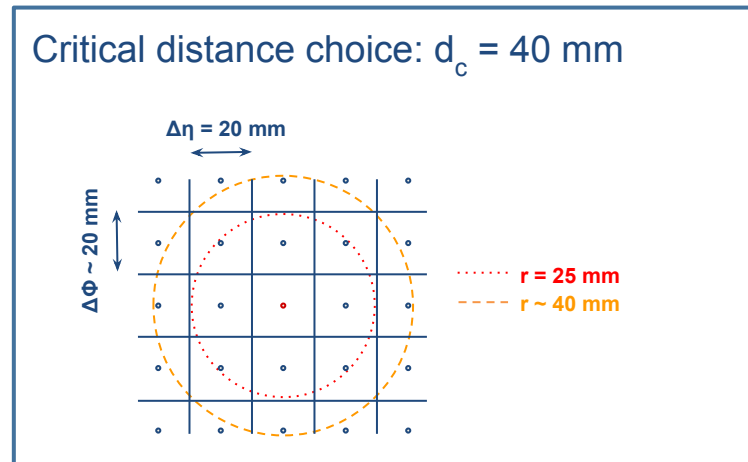$\sigma = 1.775 \pm 0.067$

Comparison with **Pandora Clusters** not completely equitable comparison (it includes a dedicated calibration procedure), but comparable results in terms of energy linearity and resolution

# Performance evaluation
## - Noble Liquid Calo -

# Noble Liquid ECAL for FCC-ee

- 12 layers, only barrel considered
  - cell size in Φ: 17.9 mm - 20.7 mm
  - cell size in η: ~ 20 mm
- Sample (if not stated otherwise):
  - 500 single gamma at 10 GeV
  - $\theta_{[min, max]}$ = [50, 130]

Critical distance choice: $d_c$ = 40 mm
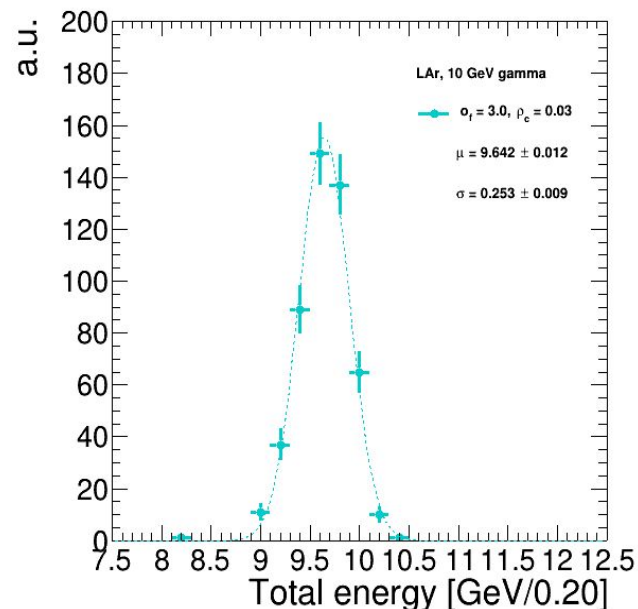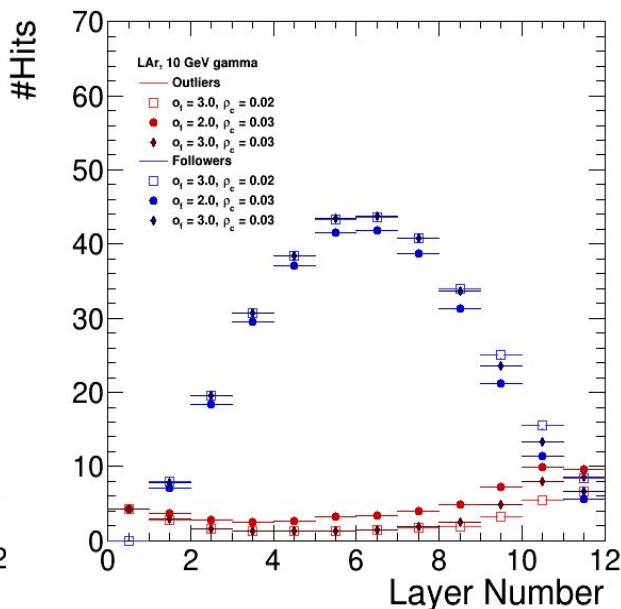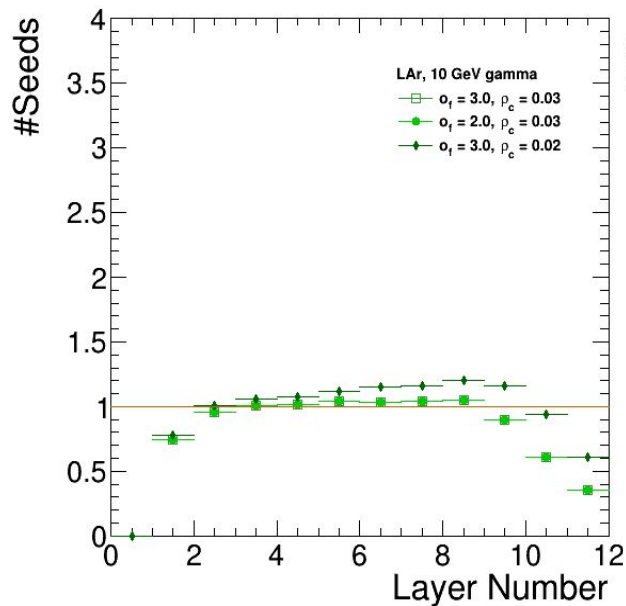
Δη = 20 mm

ΔΦ ~ 20 mm

r = 25 mm
r ~ 40 mm

active gap (noble liquid)    absorber    read-out electrode

1st layer
(presampler)
no Pb

cryostat

2.12 cm
4.91 cm

210 cm

270 cm

5 cm    40 cm    256 cm    10 cm

radial   distance on
depth    the electrode
(cm)     (cm)

256     56.5
252.5   51.6
249     46.6
245.5   41.7
242     36.7
238.5   31.8
235     26.8
231.5   21.9
228     16.9
224.5   12
221     7.1
217.5   2.1
216     0

θ = 90°   θ = 85°   θ = 80°   θ = 75°   θ = 70°   θ = 65°   θ = 60°   θ = 55°   θ = 50°

θ = 45°

0   10   20   30   40   50   60   70   80   90   100   110   120   130   140   150   160   170   180   190   200   210   (cm)
z

# Parameters tuning

- 500 events with single gamma (from vertex) at **10 GeV**

- **$d_c$ = 40.00, $\rho_c$ = 0.03, $o_f$ = 3.0**

# Comparison with other cluster algorithms

- **Sliding window**: It considers the calorimeter as a two-dimensional grid in η-φ space, neglecting the longitudinal segmentation of the calorimeter.

- **Topological clustering**: It starts with a seed cell and then adds topologically connected calorimeter cells
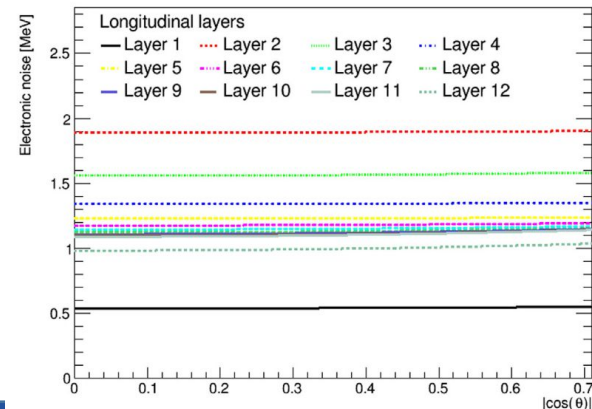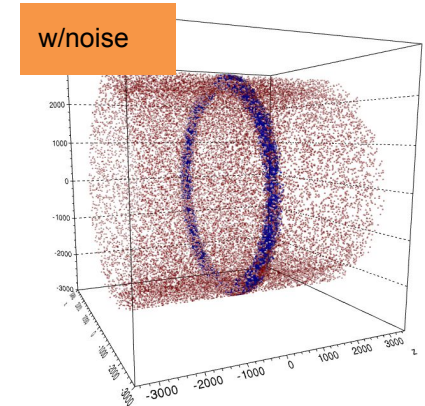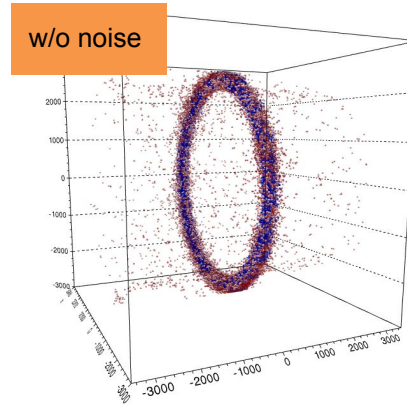


CLUE creates about ~10 clusters per event (up to few GeV per layer)
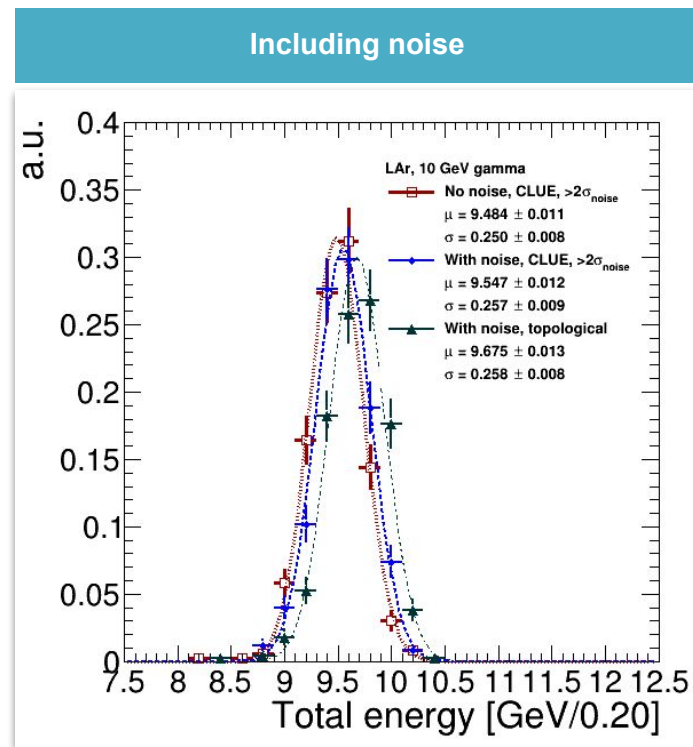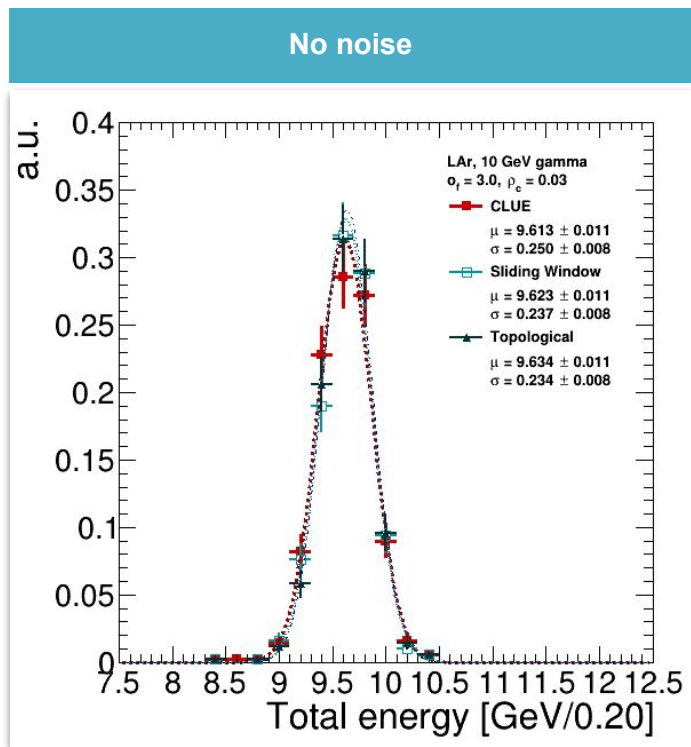
# Noise in Liquid Argon Calorimeter

- High level of noise in the detector

- In the topoclustering, there is no filter directly at the beginning for the noise, but this is done using cuts in the algorithm itself

- The main observable is the cell significance $\xi_{cell}$ which is defined as the absolute value of the ratio of the cell signal to the expected noise in this cell

$$\xi_{cell} = \left| \frac{E_{cell}}{\sigma_{cell}^{noise}} \right|$$

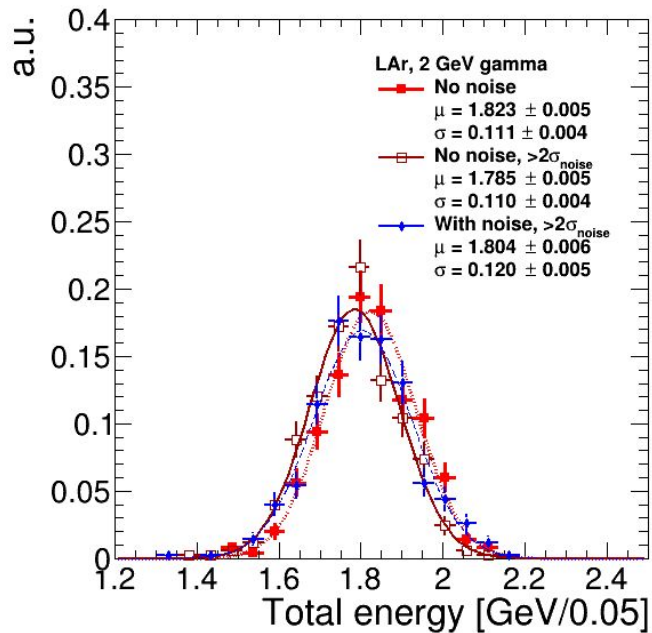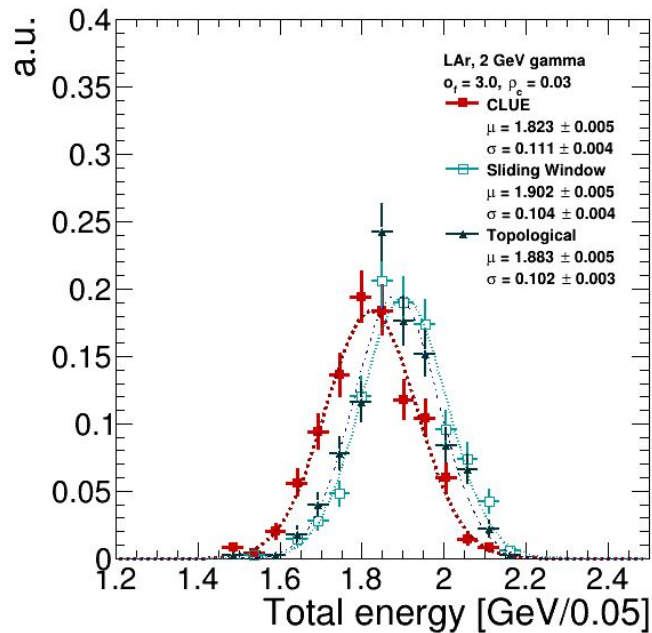- CLUE hits w/noise **selected with filter of > 2σ$_{noise}$**

Signal produced only with θ ~ 90.25°



w/o noise

w/noise

# Comparison with other cluster algorithms

**No noise**

**Including noise**

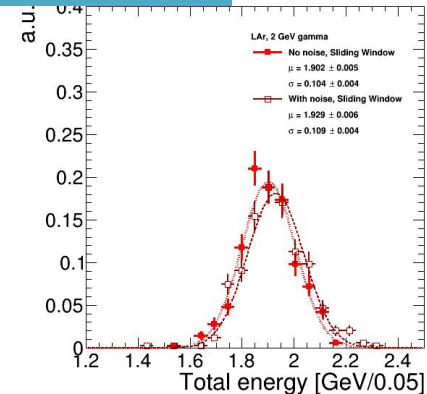# Low(er) energy
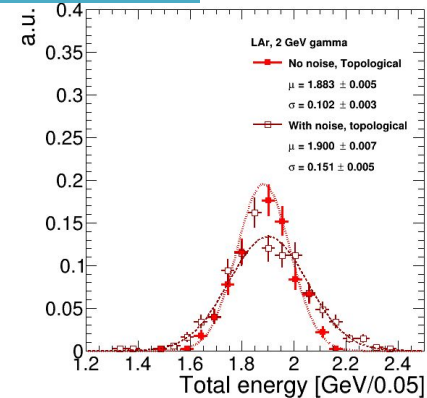
- 500 events with single gamma (from vertex) at **2 GeV**
  Motivated by flavor physics searches at Z peak


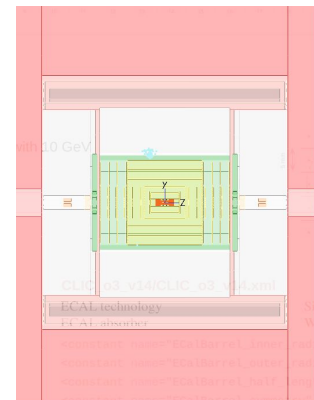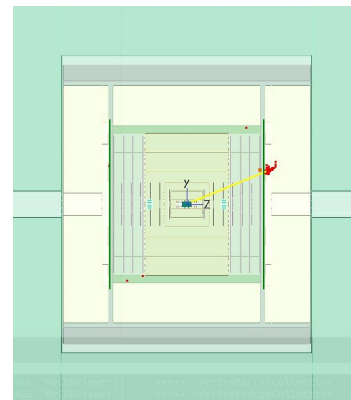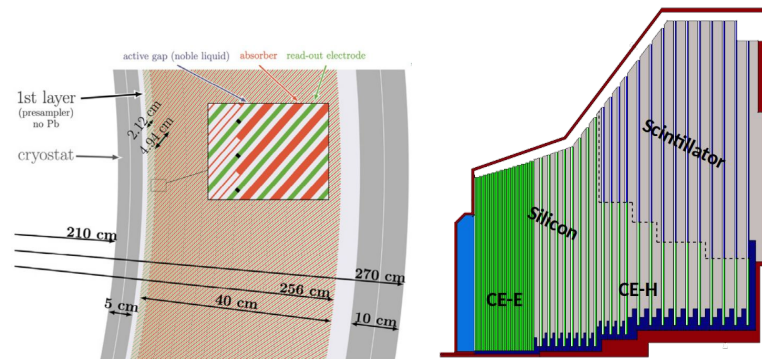
Topological

# Conclusions & Outlook

# Conclusions

- k4Clue package ([v01-00](#)) has improved upon the standalone CLUE
  - Run on the full detector (barrel & endcap)
  - Adapted for different types of calorimeters

- Analysis on three different future calorimeters has demonstrated the good performance for single gamma events
  - Good performance even in the presence of noise
  - Compared favorably to other baseline algorithms

This work highlights the adaptability and versatility of the CLUE algorithm for a wide range of experiments and detectors, as well as its potential for future high-energy physics experiments beyond CMS

  - Improvements from k4clue also under discussion to use the developments also in CMS (Phase-2 barrel region)

Erica Brondolin (erica.brondolin@cern.ch)

# Conclusions

- Final article summarizing k4clue and its performance for future collider detectors almost ready
  - Computing time under study

- This research was supported by the CERN Strategic R&D Programme on Technologies for Future Experiments

- Special thanks go to the Key4hep team and the FCC-ee liquid calorimeter software experts for the support
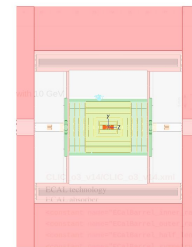
**EP-R&D**

**Programme on Technologies for Future Experiments**
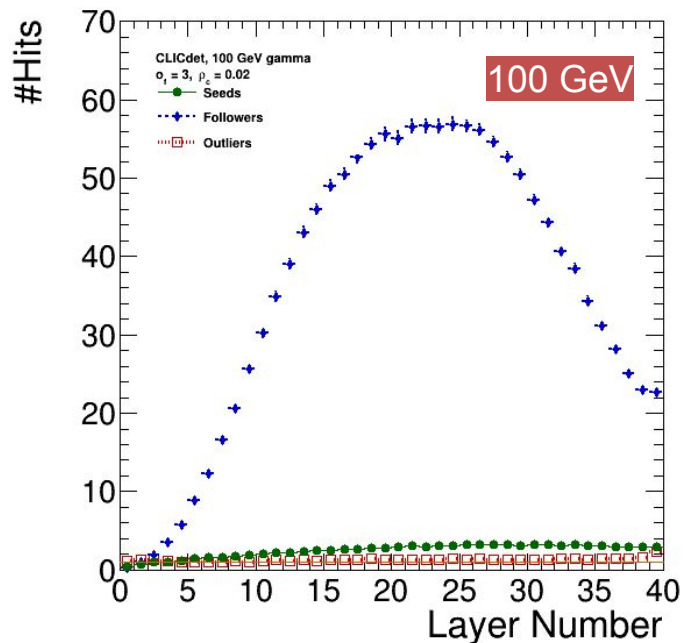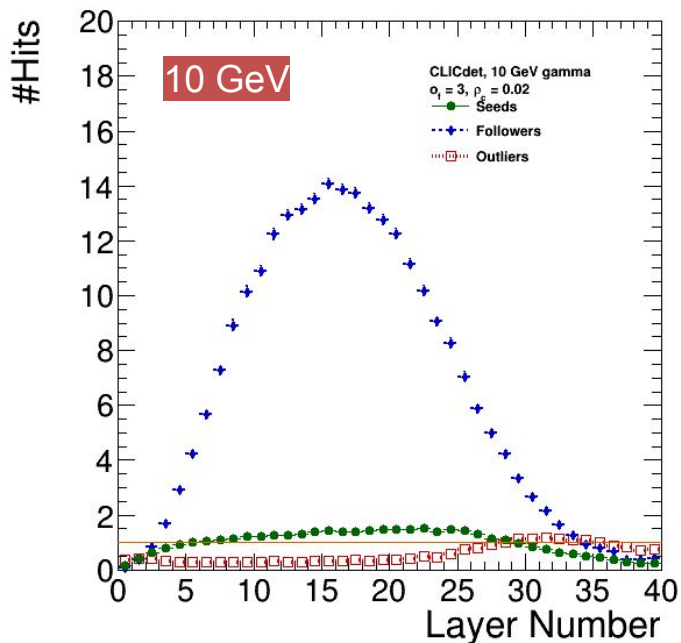
# Backup

# Integrating CLUE in Key4hep

- **GitHub CI** to ensure that the modifications or additions to the software do not break the clusterization process
    - In the latest release was modified to focus on C++ code and `EDM4hep` data
    - Test on both current key4hep release and nightlies
- **Validation**
    - `edm4hep:CLUECalorimeterHit`
        - CalorimeterHit class with specific methods related to the CLUE algorithm
    - `CLUEHistograms` class to produce ntuples

# CLICdet results

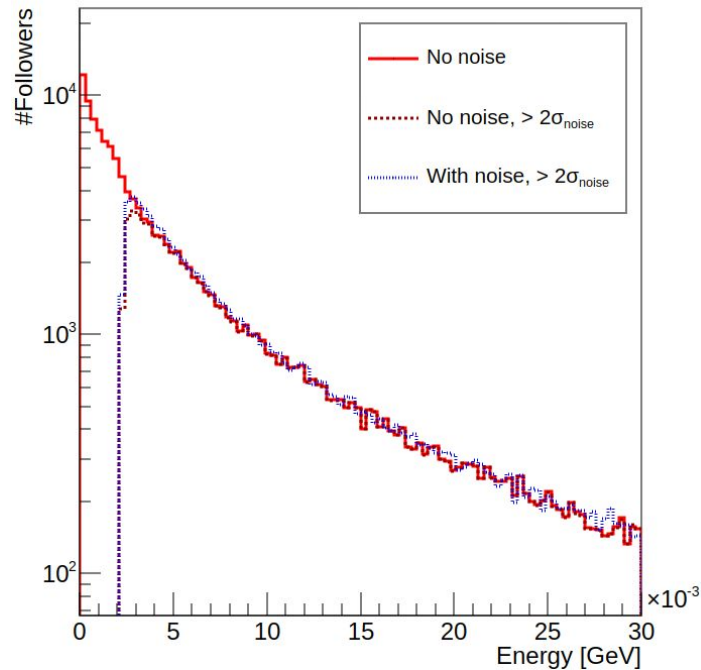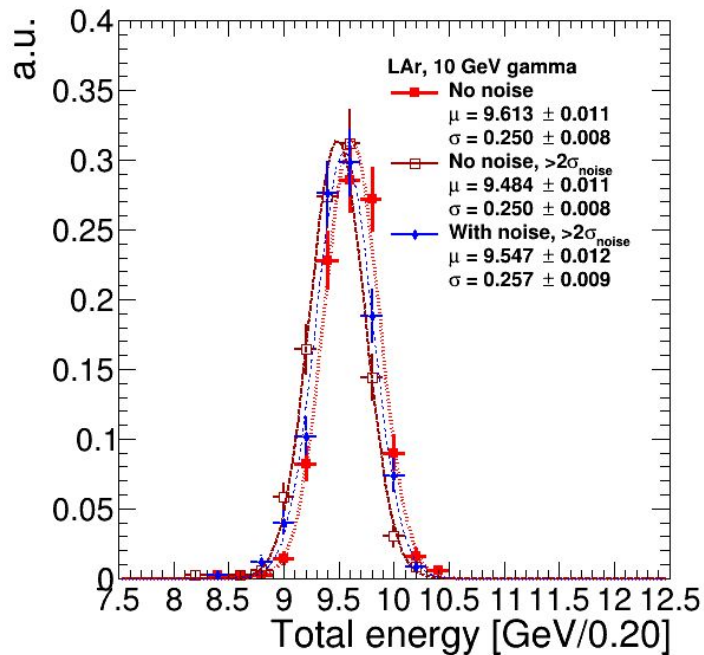- Using same input parameters selected for CLD



Comparison with **Pandora Clusters** not completely equitable comparison (it includes a dedicated calibration procedure), but comparable results in terms of energy linearity and resolution
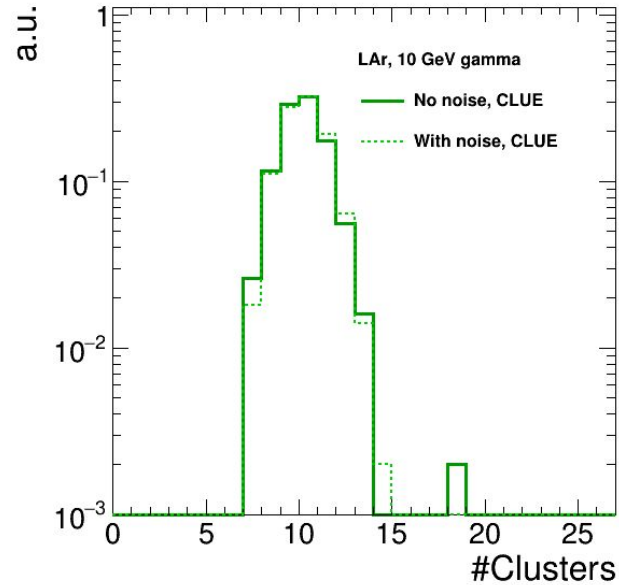
# Noble Liquid Calo

Pre-filtering

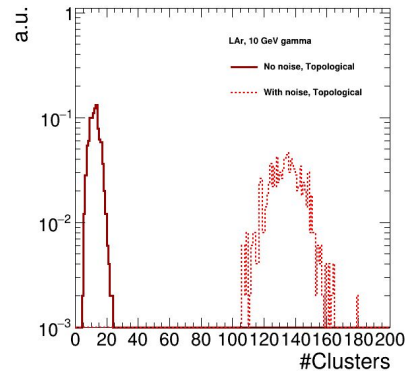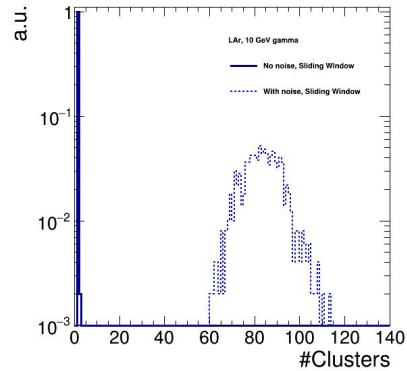- CLUE hits w/noise selected with filter of > $2\sigma_{noise}$
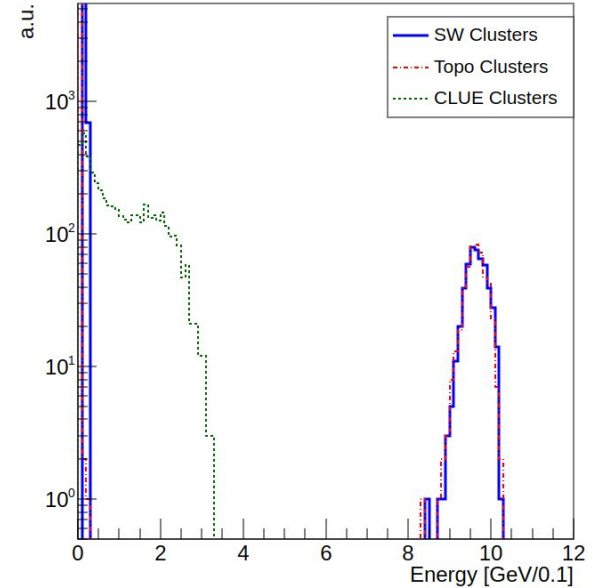
# Noble Liquid Calo

Comparison with other cluster algorithms



No significant effect on CLUE clusters - about ~10 per event (one per layer)
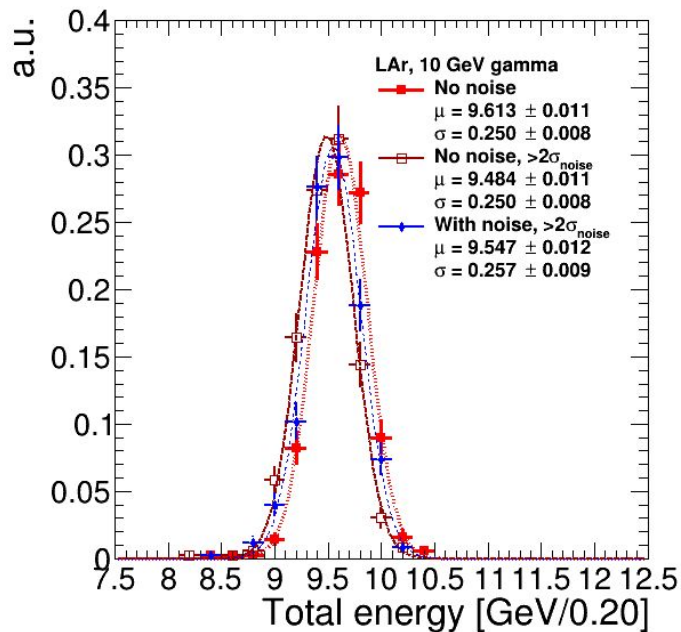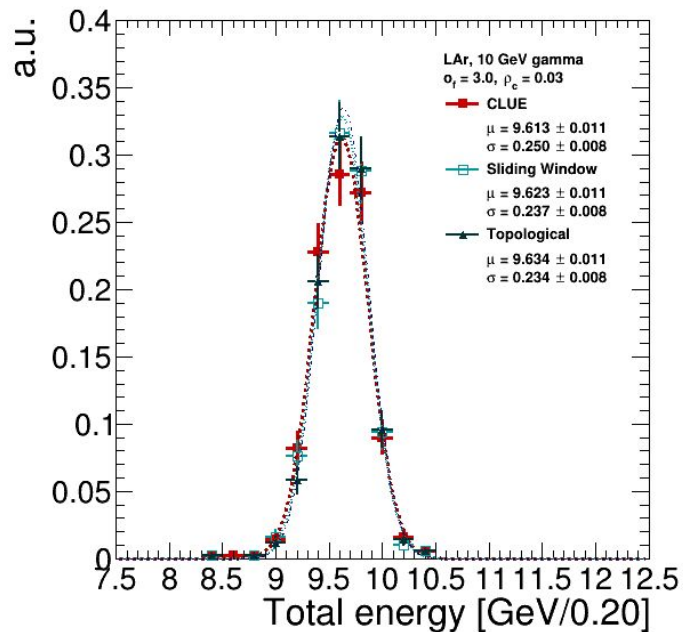
w/noise

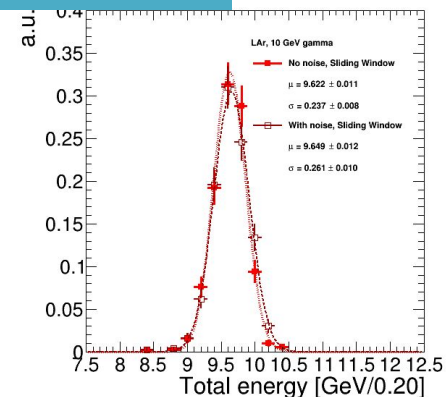**More than one SW and Topo cluster per event, but most of them with low energy**

# Noble Liquid Calo

Summary for 10 GeV gammas