



# Graph Clustering: a graph-based clustering algorithm for the electromagnetic calorimeter in LHCb

Núria Valls Canudas<sup>a</sup> , Míriam Calvo Gómez<sup>b</sup> , Xavier Vilasis-Cardona<sup>c</sup> , Elisabet Golobardes Ribé<sup>d</sup> 

Smart Society Research Group, Engineering Department, La Salle Universitat Ramon Llull, Sant Joan de la Salle 42, 08022 Barcelona, Spain

Received: 22 December 2022 / Accepted: 10 February 2023 / Published online: 25 February 2023  
© The Author(s) 2023

**Abstract** The recent upgrade of the LHCb experiment pushes data processing rates up to 40 Tbit/s. Out of the whole reconstruction sequence, one of the most time consuming algorithms is the calorimeter data reconstruction. It aims at performing a clustering of the readout cells from the detector that belong to the same particle in order to measure its energy and position. This article presents a new algorithm for the calorimeter data reconstruction that makes use of graph data structures to optimise the clustering process, that will be denoted Graph Clustering. It outperforms the previously used method by 65.4% in terms of computational time on average, with an equivalent efficiency and resolution. The implementation of the Graph Clustering method is detailed in this article, together with its performance results inside the LHCb framework using simulation data.

## 1 Introduction

LHCb is one of the four main experiments at the LHC at CERN. It consists of a forward-arm spectrometer detector designed to measure the production and decay properties of charm and beauty hadrons with high precision [1, 2]. Starting in 2022, a major upgrade has taken place in order to adapt the luminosity rates of the experiment to the LHC conditions in Run 3. It implies an increment of the instantaneous luminosity by a factor of five to  $\mathcal{L} = 2 \times 10^{33} \text{cm}^{-2} \text{s}^{-1}$  [3] and a readout rate of 30 MHz, or a maximum data rate of 40 Tbit/s for all the subdetectors. In these conditions, collisions which are of interest for many LHCb analyses reach the MHz level in the detector's geometrical acceptance [4]. Therefore,

the event selection process is expected to provide an offline-quality reconstruction within the throughput requirements. With the vision of future upgrades implying even tighter time constraints, the LHCb reconstruction needs to be as optimal as possible.

Among the five most time-consuming algorithms in the High Level Trigger 2 (HLT2) sequence of LHCb, the calorimeter data reconstruction was the fourth one with around 15% of the total cost. To make a significant improvement, the data reconstruction process from this specific sub-detector has been a target to optimise. In this article, a new algorithm for calorimeter data reconstruction called Graph Clustering is presented. In terms of execution time, Graph Clustering outperforms the previous method by up to 65.4% with an equivalent efficiency. Overall, it provides an average throughput reduction of 9.8% in the whole HLT2. Furthermore, it is currently the default solution for calorimeter data reconstruction in the upcoming Run 3.

This article is aimed to give a background in other reconstruction methodologies used for similar problems, specifically in High Energy Physics, in Sect. 2. In Sect. 3, an introduction to the Electromagnetic Calorimeter (ECAL) of LHCb is given. Section 4 provides an extensive detail of the Graph Clustering implementation. Finally, in Sect. 5 a review of the performance of the algorithm is given, followed by a discussion and conclusions.

## 2 Background

Calorimeter data reconstruction can be understood as a clustering problem, as it aims to group the energy deposits from particles following a set of rules. Classical unsupervised clustering algorithms use extensive recursive functions to create clusters according to metrics related to distance or density [5]. Despite the cluster concept, the calorimeter data recon-

<sup>a</sup> e-mail: [nuria.valls@salle.url.edu](mailto:nuria.valls@salle.url.edu) (corresponding author)

<sup>b</sup> e-mail: [miriam.calvo@salle.url.edu](mailto:miriam.calvo@salle.url.edu)

<sup>c</sup> e-mail: [xavier.vilasis@salle.url.edu](mailto:xavier.vilasis@salle.url.edu)

<sup>d</sup> e-mail: [elisabet.golobardes@salle.url.edu](mailto:elisabet.golobardes@salle.url.edu)

struction strategy for LHCb has not much in common with classical clustering algorithms, due to the strong physics and execution time requirements.

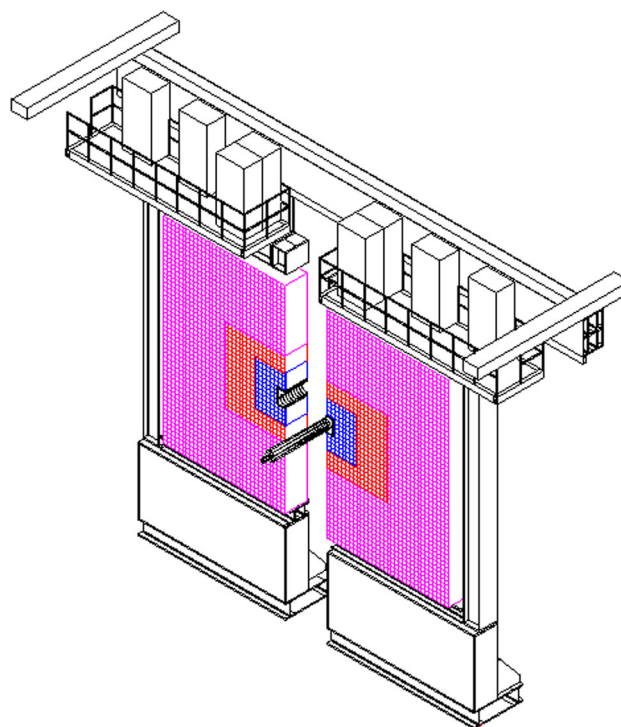
Focusing on the field of calorimetry in High Energy Physics, the Cellular Automaton has been a benchmark solution for many years [6]. LHCb has been using this strategy for Runs 1 and 2. In 2004 an approach using spanning trees was proposed, using this flexible data structures to exploit the neighbourhood definitions in general calorimeter data [7], but it does not consider the cluster separation needed in LHCb. Graph data structures started to appear in the field with increasing popularity of deep learning in the form of a neural model based on graphs [8]. Several approaches have used these graph neural networks on layered calorimeters [9, 10] showing promising results on clustering energy deposits in consecutive calorimeter layers. However, the LHCb calorimeter geometry is bi-dimensional. Within this context, other approaches have also used graph neural networks [11] and convolutional neural networks [12, 13] with similar conditions as ECAL in LHCb. That said, the inference of some deep learning models is still not mature enough to be incorporated in the LHCb software framework.

Graph structures have demonstrated to be suited for calorimeter data. Hence, the Graph Clustering algorithm stores the calorimeter digits into graphs and makes use of its flexible neighbourhood properties to define the clusters. Moreover, it follows the same reconstruction principles from the Cellular Automaton strategy, which has proved to give a good performance in terms of reconstruction efficiency.

### 3 Detail of the electromagnetic calorimeter

The LHCb experiment has a subset of eight dedicated detectors to acquire data from the particles generated in the LHC collisions. The electromagnetic calorimeter is one of them. Its main purpose is the identification of electrons and photons, and the measurement of their energies and positions [14]. The ECAL has a rectangular shape of  $7.8 \times 6.3 \text{ m}^2$  and is placed perpendicular to the accelerator beam pipe. The energy measurement area is segmented into individual square-shaped modules. Each module is made from lead absorber plates interspaced with scintillator tiles as active material. The general structure is segmented in three different rectangular shaped regions, as can be seen in Fig. 1.

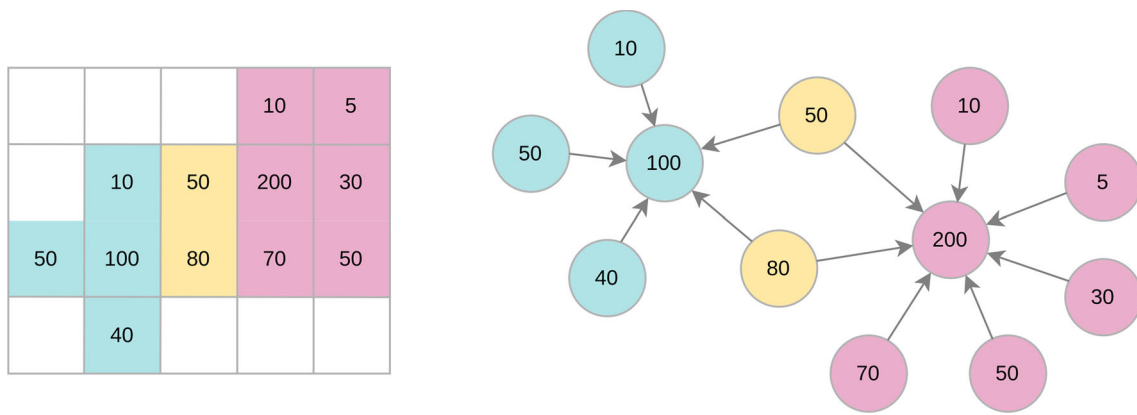
Although all modules have the same size of  $12 \times 12 \text{ cm}^2$ , the number of readout cells on a module depends on the region. The inner region is the closest to the beam pipe and has the highest occupancy of incident particles. Thus, it has the highest granularity among the three regions, with nine readout cells of  $4 \times 4 \text{ cm}^2$  per module. The middle region surrounds the inner one and has four readout cells of  $6 \times 6$



**Fig. 1** The electromagnetic calorimeter 3d view from behind the detector towards the interaction point [14]

$\text{cm}^2$  per module. The outer region has a single readout cell of  $12 \times 12 \text{ cm}^2$  per module.

The output data obtained from the ECAL modules are the values from each readout cell concerning the accumulated energy deposited by incident particles in a collision event with 12 bit precision. The digitized value is converted to MeV as the measure of energy. Another measure used in the data reconstruction is the transverse energy, which is calculated using the measured energy of a cell and its angular position in the ECAL. Since particles may deposit energy in more than one readout cell, the energy deposits need to be reconstructed and clustered together with the ones belonging to the same particle. This process is precisely the cluster reconstruction for the calorimeter. The current definition of a calorimeter cluster stands for a  $3 \times 3$  block of readout cells around an energy peak. Studies have been done regarding the cluster shapes [15] where a combination of  $2 \times 2$  and swiss-cross cluster shapes show promising performance for high luminosity, although the  $3 \times 3$  cluster is used as a base for masking other shapes on clusters. Hence, the definition of  $3 \times 3$  readout cell clusters is maintained through all the regions of the detector.



**Fig. 2** An example of two clusters with overlapping cells on the calorimeter on the left and its graph representation on the right

### 4 Method

The baseline idea behind the Graph Clustering algorithm is to use graphs as a data structure to store the event digits. It transforms the calorimeter digits into independent graph structures, where only relevant digits for a cluster are contained into isolated graphs. Following graph theory nomenclature, each energy digit from an event is represented as a vertex  $v$  in the graph, also called node. The relations between digits, representing links to the same cluster, are defined as directional edges  $(u, v)$  between the source digit node  $u$  and the target node  $v$ . By design, the target nodes of all edges in the graph are the seeds of the reconstructed clusters, where a seed is defined as a local maximum energy digit in the calorimeter grid over a threshold of 50 MeV in transverse energy. With this, the cluster seeds can be easily identified as nodes with only incoming edges. Furthermore, a node can be linked to more than one seed if it has energy deposits from more than one particle. These particular cases are called overlap cells. Overall, the graph derived from an event may contain structures like the example shown in Fig. 2.

The following subsections describe in detail the four steps needed in the Graph Clustering reconstruction process.

#### 4.1 Sorting

To achieve the mentioned representation of the digits, the algorithm needs to make an efficient insertion of the edges into the graph structure. Since all the edges are based on the cluster seeds, the initial key point is to identify seed candidates. As defined previously, a cell in the ECAL grid is considered a seed if it is a local maximum and has a minimum transverse energy value of 50 MeV. A local maximum in this context defines a cell that has the highest energy value among a  $3 \times 3$  cell area around it in the calorimeter grid. This definition is the same as the one used in the Cellular Automaton algorithm [6].

In order to process the seed candidates of an event in the first place, all the digits above 50 MeV need to be sorted by decreasing transverse energy value. In the proposed algorithm, the sorting is computed using Introspective Sorting [16], which is a hybrid sorting algorithm that combines three different methods to provide fast average performance and optimal worst-case performance.

#### 4.2 Insertion

The role of the insertion step is to build the graph edges between the event digits such that the graph structures of Fig. 2 are obtained. A pseudo-code notation of this process is stated in Algorithm 1.

---

#### Algorithm 1 Graph insertion

---

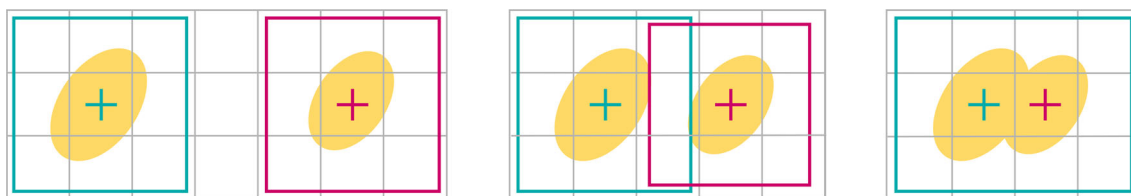
```

1:  $G \leftarrow$  directional weighted graph
2: for each  $energy, id \in sortedDigits$  do
3:   if  $id$  not inserted in  $G$  then
4:     if  $id$  is local maxima then
5:       add node  $id$  in  $G$ 
6:       for each  $n_{energy}, n_{id} \in$  neighbours of  $id$  do
7:         add node  $n_{id}$  and edge  $(n_{id}, id, w = 1)$  in  $G$ 
8:         if  $id$  is a merged  $\pi^0$  candidate then
9:           add  $id$  and  $n_{id}$  to merged  $\pi^0$ 
10:        end if
11:       end for
12:     end if
13:   else if  $id \in merged\pi^0$  then
14:      $seed =$  first seed from  $id$  in  $G$ 
15:     for each  $n_{energy}, n_{id} \in$  neighbours of  $id$  do
16:       if  $energy > n_{energy} \ \& \ n_{id}$  not in  $G$  then
17:         add node  $n_{id}$  and edge  $(n_{id}, id, w = 1)$  in  $G$ 
18:       end if
19:     end for
20:   end if
21: end for

```

---

It essentially iterates over each sorted digit. That digit may have already been inserted in the graph. If so, this means it



**Fig. 3** Diagram representation of  $\pi^0$  cluster cases on the calorimeter. From left to right: the two photons are separable and without overlap, it is a resolved  $\pi^0$ . The two photons are separable but have three over-

lapping digits, it is however a resolved  $\pi^0$ . The two photons are not separable, it is a merged  $\pi^0$  and is reconstructed as a single cluster bigger than  $3 \times 3$

is a neighbour of a more energetic digit. In that case, it cannot be a seed since there cannot be two adjacent maxima by construction, except for the case of merged  $\pi^0$ s, which is explained in Sect. 4.2.1. Therefore, that digit is not inserted. On the other hand, if the digit has not yet been inserted on the graph, it can be either a seed or a residual digit, meaning it is not a local maximum and does not have any seed on its neighbourhood. To distinguish between the two, the algorithm checks if that digit is a local maximum. If it is the case, the seed is inserted in the graph together with all its neighbour digits linking them with edges to the seed. The default weight value for all edges is one.

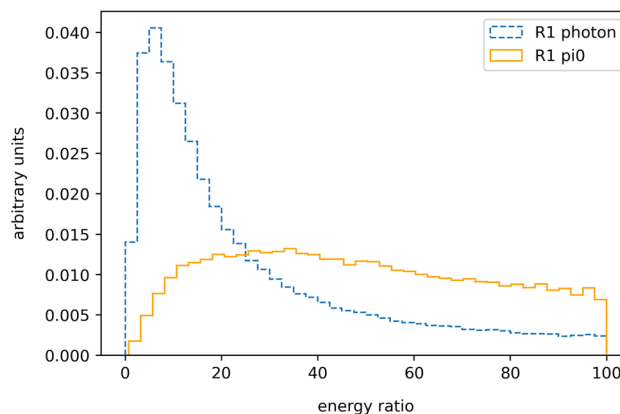
Additionally, if a merged  $\pi^0$  candidate is identified following the metrics described in Sect. 4.2.1, there is a second seed added to the same cluster. The neighbours of the second seed are also linked with an edge to the first seed if they are not already inserted and if its energy deposit is lower than the energy of the first seed.

At the end of the insertion step, the clusters are already grouped in the graph. However, the overlap cases still need to be processed to adjust the weight of the overlap edges.

#### 4.2.1 Merged $\pi^0$ case

One of the reconstruction requirements for the LHCb calorimeter is the correct identification of neutral pions,  $\pi^0$ , which decay into two photons before reaching the calorimeter. Depending on the energy and momentum of the  $\pi^0$ , the two photons arrive at the calorimeter with a certain separation.

If the seeds of the two photons are distanced more than one cell, they will be reconstructed as separate clusters. This case is called a resolved  $\pi^0$ . Otherwise, the two photons may travel very close to each other and reach the calorimeter at one cell distance or less. In that case, the reconstruction is done as a single cluster, since the definition of maxima does not allow two adjacent cluster seeds. When photons are not separable, it is then called merged  $\pi^0$  case. Hence, the super-cluster from a merged  $\pi^0$  can be bigger than the  $3 \times 3$  window around the seed as can be seen in Fig. 3.



**Fig. 4** Normalized histograms of the energy ratio between the second most energetic digit and the cluster seed for photon samples and  $\pi^0$  samples

There are other dedicated algorithms in the LHCb sequence that use the output of the calorimeter data reconstruction, together with other detector data, to properly identify and classify  $\pi^0$  particles. The cluster shape used in these cases is a mask of  $5 \times 5$  cells around the seed [17]. Therefore, the residual energy outside the  $3 \times 3$  window of a merged  $\pi^0$  is crucial, as it contains part of the energy from the second photon. That is why the Graph Clustering algorithm adapts the shape of potential merged  $\pi^0$  candidates, expanding the cluster up to the neighbours of the second most energetic digit in the cluster.

To avoid adding complexity to the data reconstruction algorithm, the energy deposits of the  $3 \times 3$  cluster are the only source of information used to find a metric that can provide a soft selection filter for merged  $\pi^0$  candidates at run time. Therefore, we have studied the relation between the two most energetic digits as a ratio labeled  $R1$ . Figure 4 shows a normalized histogram of the  $R1$  ratio for over 46,000 samples of single  $\pi^0$  deposits from  $B^0 \rightarrow \pi^+\pi^-\pi^0$  decays simulated using Run 3 conditions compared to photon samples from  $B^0 \rightarrow K^*\gamma$  decays also simulated using Run 3 conditions. The main difference between the two distributions is that the majority of photons have an energy ratio between 0 and 25 whereas  $\pi^0$  tend to have higher energy ratios in most

cases. Therefore, the algorithm sets a threshold of value 25 in  $R1$  to determine if a cluster needs to be expanded more than  $3 \times 3$ . This value has been optimized and ensures that the residual energy left outside the cluster is less than 9% for the studied  $\pi^0$  samples and that the cluster expansion affects an average of 8.2% of the clusters in an event. Further studies have determined that small variations around 10% of the selected threshold value do not significantly change the time complexity of the algorithm nor the  $\pi^0$  resolution.

Moreover, given that only high energetic  $\pi^0$ s will be merged, a second threshold is added cutting the seed candidates under 1 GeV in the merged  $\pi^0$  candidate selection. This value has been chosen according to the  $\pi^0$  samples studied since it is the minimum seed value for a  $\pi^0$  to be merged and not resolved.

### 4.3 Connected components

Once the insertion is finished, the graph structure contains all the relevant energy digits as nodes linked with the elements of each cluster and other overlapping clusters, if any. From this point on, the algorithm needs to process each cluster or group of overlapping clusters separately. Using graph theory terminology, a subset of nodes from a graph connected by some path is called a weakly connected component. Therefore, to retrieve the list of nodes that belong to the same cluster or group of overlapping clusters, the algorithm needs to find all the weakly connected components of the graph. In the proposed algorithm, this process is implemented as a depth-first search [18], which explores an entire graph exploring all its branches as far as possible before backtracking. Its time complexity is  $O(|V| + |E|)$  [19] where  $V$  is the number of vertices or nodes in the graph and  $E$  is the number of edges. Once all the vertices of the graph are visited, the nodes and edges on each weakly connected component are obtained.

### 4.4 Analysis of clusters

The final step of the reconstruction is to analyze each weakly connected component to resolve the overlap cases if any and transform the graph clusters into the output cluster format. The processing of a weakly connected component can be done independently of the others, since each one contains only the relevant nodes and edges for a cluster.

The analysis of clusters consists on iterating through the list of weakly connected components following the pseudocode in Algorithm 2. Only connected components with more than one node are considered as reconstructed clusters. Any isolated node is likely to be a residual energy deposit from a cluster and should not be considered a reconstructed cluster itself.

If there is more than one seed in a connected component, there is at least one cell overlapping between two clusters. In

### Algorithm 2 Analysis of connected components

```

1: for each  $wcc \in weaklyConnectedComponents$  do
2:   if  $wcc.size() > 1$  then
3:     calculate overlap weights (Algorithm 3)
4:     for each  $id \in wcc$  do
5:       if  $id$  in-edges  $> 1$  &  $id$  out-edges  $== 0$  then
6:         add  $id$  as a cluster seed to  $clusters$ 
7:         for each  $vertex$  connected to  $id$  do
8:           add  $vertex$  as entry of  $id$  in  $clusters$ 
9:         end for
10:      end if
11:    end for
12:  end if
13: end for

```

that case the overlap resolution, defined in Algorithm 3, consists in assigning a fraction of the energy of the overlapping cell to each of the seeds linked to it. The fraction is calculated as a function of the energy of the clusters and is stored as the weight of an edge.

### Algorithm 3 Calculate overlap weights

```

1:  $clusterEnergy \leftarrow$  empty map
2: for each  $vertex \in wcc$  do
3:   if  $vertex$  out-edges  $\geq 2$  then
4:     for each  $end\_vertex \in vertex$  out-edges do
5:        $energy =$  accumulate energy from the nodes linked to
         $end\_vertex$ .
6:        $energy+ = end\_vertex$  energy/num out-edges.
7:       store  $energy$  to  $clusterEnergy$ 
8:     end for
9:      $totalEnergy =$  accumulate  $clusterEnergy$  energies with
        entries  $\in vertex$  out-edges
10:    for each  $end\_vertex \in vertex$  out edges do
11:       $weight = \frac{clusterEnergy\ at\ end\_vertex}{totalEnergy}$ 
12:      set edge ( $vertex, end\_vertex, w = weight$ )
13:    end for
14:  end if
15: end for

```

Entering in more detail, this algorithm iterates through all the vertices in a connected component. It searches for overlap vertices, identified by having two or more output edges, and accumulates the energy of all the connected nodes on all the clusters involved in the overlap. The energy of the overlap node is equally fractioned among the number of involved clusters to avoid accounting it more than once. Then, the weight of every overlapping edge is computed as the fraction between the energy of the target cluster and the sum of all the clusters involved in the overlap.

## 5 Results

All the algorithm tests have been done within the GAUDI framework [20,21]. For comparison purposes, this paper

**Table 1** Efficiency results in number of reconstructed versus reconstructible clusters from 80,000  $B^0 \rightarrow K^*\gamma$  events

Algorithm	Reconstructible	Reconstructed	Efficiency (%)
Graph clustering	43234	35313	$81.68 \pm 0.19$
Cellular automaton	43234	34872	$80.66 \pm 0.19$

evaluates the performance of the Graph Clustering algorithm and the Cellular Automaton algorithm as it has been a benchmark solution until now. Both are tested with the same Monte Carlo data from  $B^0 \rightarrow K^*\gamma$  simulations using Run 3 conditions.

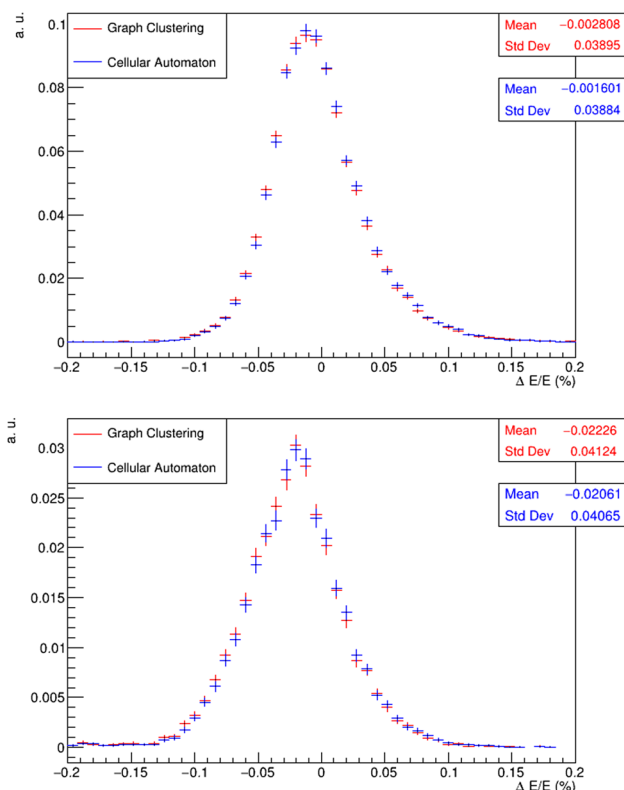
The quality of the reconstruction in calorimeter algorithms in LHCb is evaluated using metrics of efficiency, energy resolution and position resolution. The efficiency is defined as the fraction between reconstructed particles over reconstructible particles in a set of events. Reconstructible particles are photons that have deposited at least 90% of their energy in the calorimeter cells. On the other hand, reconstructed particles are reconstructible particles matching a cluster from which at least 90% of their energy belong to that particle. This ratio is later referred to as match fraction. Table 1 shows that Graph Clustering has a higher efficiency than the Cellular Automaton, with 1.02% more reconstructed clusters.

On the other hand, the energy resolution and position resolution metrics aim to measure the difference in energy and position between the reconstructed clusters and the associated Monte Carlo particles. Resolutions are evaluated for  $\gamma$  and  $\pi^0$  particles. For both cases, we evaluate the difference in position on the  $X$  and  $Y$  axis and the difference in energy as a percentage. For  $\gamma$  resolution, a total of 80,000 simulation samples of  $B^0 \rightarrow K^*\gamma$  decays have been used, and another 80,000 samples of  $B^0 \rightarrow \pi^+\pi^-\pi^0$  decays have been used for  $\pi^0$  resolution. The study accounts for all the clusters with a match fraction higher than 0.9 since it is the standard match threshold for a cluster to be considered reconstructed in terms of efficiency.

Figure 5 shows the energy distribution for both methods, before any corrections are applied [22], where  $\Delta E$  stands for the difference in reconstructed energy and truth energy of a cluster. It can be seen that for both  $\gamma$  and  $\pi^0$  samples the two distributions look very alike. For energy resolution, Graph Clustering is slightly more shifted to negative values, but overall it can be said that the resolution in energy is equivalent to the Cellular Automaton one.

Regarding the position resolution, Fig. 6 shows that the  $X$  and  $Y$  distributions have again an equivalent behavior for both methods. For simplicity, only the  $\pi^0$  resolutions are shown for position, since the differences with  $\gamma$  samples are minimal.

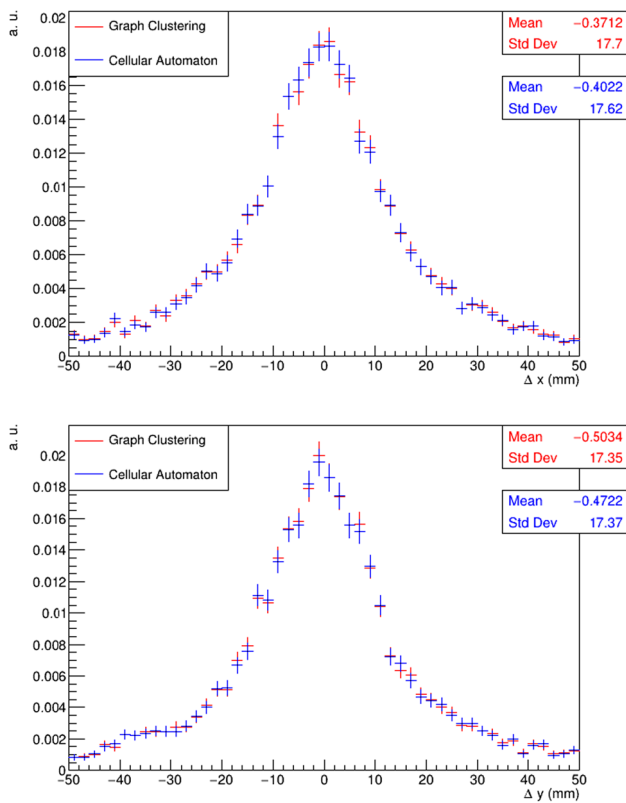
Regarding the execution time, it is defined as the time elapsed between the first and the last lines executed in an



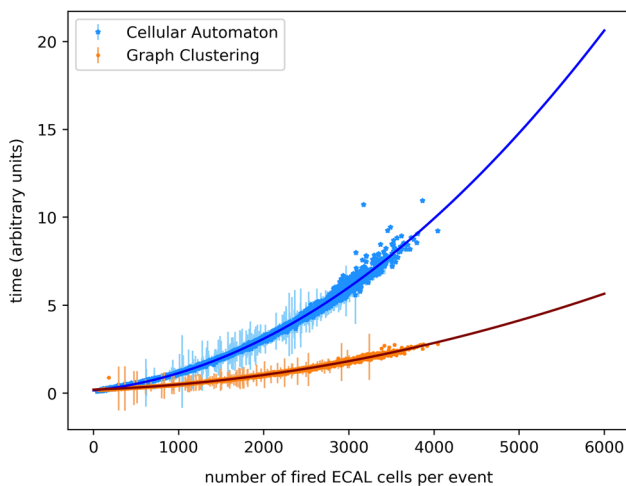
**Fig. 5** Normalized histograms of the energy resolution with no corrections for clusters with a match fraction over 0.9 using  $\gamma$  samples in the top plot and  $\pi^0$  samples in the bottom plot

algorithm. Figure 7 shows a plot of the execution time in arbitrary units as a function of the number of digits per event.

The plotted time measurements are obtained as the average measured time from all the events with the same number of digits, from a total of 100,000 events from  $B^0 \rightarrow K^*\gamma$  simulation. The figure also includes error bars from the standard deviation of the samples with the same number of digits. This error reflects the small variation of complexity that the algorithm may have according to the distribution of the digits in the event as well as, more significantly, the variations in the available resources from the distributed computing environment where the tests have been executed. Taking as a reference the fitted curves from the plot, for events with less than 150 digits, the Cellular Automaton is faster. However, from that point on, Graph Clustering outstands the benchmark algorithm showing a flatter complexity curve. Furthermore, the average number of digits per event from the analysed samples is 1520 digits. At that complexity level, Graph Clustering is 65.4% faster than Cellular Automaton on average.



**Fig. 6** Normalized histograms of the X axis resolution at the top and the Y axis resolution at the bottom. Both using  $\pi^0$  samples and clusters with a match fraction over 0.9 with no corrections



**Fig. 7** Execution time measured in arbitrary units as a function of the number of digits per event for the Cellular Automaton algorithm and the Graph Clustering algorithm. On top of them, a fitted curve for every algorithm is shown

## 6 Discussion and conclusions

Graph Clustering has shown to improve the computational complexity of the calorimeter data reconstruction in LHCb. Furthermore, it is the default reconstruction solution for the

ongoing Run 3 data taking period. The baseline of the algorithm is to reproduce the same reconstruction steps as in the previously used algorithm, the Cellular Automaton, but with an optimized codification using graph data structures. Hence, it is expected and observed to have similar results compared to the benchmark in terms of efficiency and resolution. The observed efficiency is consistent with the efficiencies in Run 1 and Run 2. It is considered good in terms of performance since the definition of a reconstructible particle does not take into account noise or other fully overlapping particles, known as the pileup effect. Hence, the data reconstruction efficiency is not expected to reach 100% but gives an overall idea of the algorithms performance.

Graphs have demonstrated to be suited for calorimeter data reconstruction. Within the proposed implementation, such data structures also provide a flexible interpretation of the neighbour cells in the calorimeter grid. This could also be used to adapt the shape of the clusters to an optimized pattern depending on the region at reconstruction time and significantly accelerate its execution. Currently, the definition of an optimal cluster shape for ECAL clusters is being studied considering pileup and overlap effects as well as precision.

Within the steps of the presented Graph Clustering, as mentioned in Sect. 4.4, the analysis of each connected component is completely independent of the rest of the graph. Although it is not the most time consuming part of the algorithm, it represents a 27.3% of the total algorithm's execution time, which could benefit from parallel execution. In the context of the first level of the trigger system (HLT1) ran in GPUs, calorimeter data reconstruction is at a preliminary stage. The current implementation builds simplified clusters with lower efficiency and resolution than the benchmark. In that direction, there is currently work in progress on adapting the presented Graph Clustering logic to a CUDA algorithm optimized for parallel architectures.

As a final conclusion, the complexity curve that Graph Clustering exhibits makes it a useful alternative for other calorimeters with higher occupancy. Furthermore, the vision of future upgrades in the LHCb calorimeter is a challenging opportunity to test the limits of this algorithm.

**Acknowledgements** The authors would like to thank the LHCb computing and simulation teams for their support and for producing the simulated LHCb samples used in the paper. Specially the RTA team for their help with code optimization and the integration into the LHCb framework. This research was funded by Ministerio de Ciencia e Innovación grant number PID2019-106448GB-C32.

**Data Availability Statement** This manuscript has no associated data or the data will not be deposited. [Authors' comment: All LHCb scientific output is published in journals, with preliminary results made available in Conference Reports. All are Open Access, without restriction on use beyond the standard conditions agreed by CERN. Data associated to the plots in this publication as well as in supplementary materials are made available on the CERN document server at <https://cds.cern.ch/record/2846012>.]

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

Funded by SCOAP<sup>3</sup>. SCOAP<sup>3</sup> supports the goals of the International Year of Basic Sciences for Sustainable Development.

## References

1. A.A. Alves Jr., L. Andrade Filho, A. Barbosa, I. Bediaga, G. Cernicchiaro, G. Guerrer, H. Lima Jr., A. Machado, J. Magnin, F. Marujo et al., The LHCb detector at the LHC. *J. Instrum.* **3**, S08005 (2008)
2. LHCb Collaboration, LHCb detector performance. *Int. J. Mod. Phys. A* **30**(07), 1530022 (2015)
3. I. Bediaga, H. Chanal, P. Hopchev, S. Cadeddu, S. Stoica, M. Calvo Gomez, S. T' Jampens, I.V. Machikhiliyan, Z. Guzik, A.A. Alves Jr. et al., *Framework TDR for the LHCb Upgrade: Technical Design Report; Technical Report; LHCb-TDR-012* (CERN, Geneva, 2012)
4. C. Fitzpatrick, V.V. Gligorov, *Anatomy of an upgrade event in the upgrade era, and implications for the LHCb trigger. Technical report; LHCb-PUB-2014-027, CERN-LHCb-PUB-2014-027* (CERN, Geneva, 2014)
5. J. Han, J. Pei, H. Tong, *Data Mining: Concepts and Techniques* (Morgan Kaufmann, Cambridge, 2022), p.2022
6. V. Breton, N. Brun, P. Perret, *A Clustering Algorithm for the LHCb Electromagnetic Calorimeter Using a Cellular Automaton; Technical Report; CERN-LHCb-2001-123* (CERN, Geneva, 2001)
7. G. Mavromanolakis, Calorimeter clustering with minimal spanning trees (2004). [arXiv:physics/0409039](https://arxiv.org/abs/physics/0409039)
8. F. Scarselli, M. Gori, A. Tsoi, M. Hagenbuchner, G. Monfardini, The graph neural network model. *IEEE Trans. Neural Netw.* **20**(1), 61–80 (2008)
9. X. Ju, S. Farrell, P. Calafiura, D. Murnane, L. Gray, T. Klijnsma, K. Pedro, G. Cerati, J. Kowalkowski, G. Perdue et al., Graph neural networks for particle reconstruction in high energy physics detectors (2020). [arXiv:2003.11603](https://arxiv.org/abs/2003.11603)
10. S.R. Qasim, K. Long, J. Kieseler, M. Pierini, R. Nawaz, Multi-particle reconstruction in the High Granularity Calorimeter using object condensation and graph neural networks. *EPJ Web Conf.* **251**, 03072 (2021)
11. S.R. Qasim, J. Kieseler, Y. Iiyama, M. Pierini, Learning representations of irregular particle-detector geometry with distance-weighted graph networks. *Eur. Phys. J. C* **79**(7), 1–11 (2019)
12. M. Mazurek, Deep learning solutions for 2D calorimetric cluster reconstruction at LHCb. *4th Inter-experiment Machine Learning Workshop, Talk; LHCb-TALK-2020-178* (2020)
13. N. Valls Canudas, M. Calvo Gómez, E. Golobardes Ribé, X. Vilasis-Cardona, Use of deep learning to improve the computational complexity of reconstruction algorithms in high energy physics. *Appl. Sci.* **11**(23), 11467 (2021)
14. O. Omelaenko, P. Dalpiaz, Z. Guzik, E. Spiridenkov, P. Jarron, V. Semenov, J. Ocariz, A. Khan, P. Perret, O. Schneider et al., *LHCb Calorimeters: Technical Design Report; Technical Report; LHCb-TDR-002* (CERN, Geneva, 2000)
15. A. Abba, F. Caponio, A. Cusimano, A. Geraci, LHCb, *LHCb Particle Identification Upgrade: Technical Design Report* (CERN, Geneva, 2013)
16. D.R. Musser, Introspective sorting and selection algorithms. *Softw. Pract. Exp.* **27**(8), 983–993 (1997)
17. O. Deschamps, I. Belyaev, F.P. Machefert, G. Pakhlova, M.H. Schune, *Photon and Neutral Pion Reconstruction. Technical report (CERN-LHCb-2003-091, Geneva, 2020)*
18. M. Ginsberg, *Essentials of artificial intelligence. Newnes* (2012)
19. T. Cormen, C. Leiserson, R. Rivest, C. Stein, *Introduction to Algorithms* (MIT Press, Cambridge, 2022)
20. The LHCb Collaboration, *Upgrade Software and Computing. Technical report, ERN-LHCC-2018-007, LHCb-TDR-017* (CERN, Geneva, 2018)
21. G. Barrand, I. Belyaev, P. Binko, M. Cattaneo, R. Chytracsek, G. Corti, M. Frank, G. Gracia, J. Harvey, E. van Herwijnen, P. Maley, P. Mato, S. Probst, F. Ranjard, GAUDI—a software architecture and framework for building HEP data processing applications. *Comput. Phys. Commun.* **140**(1), 45–55 (2001)
22. A. Vallier, Measurement of the CKM angle  $\gamma$  in the  $B^0 \rightarrow DK^{*0}$  decays using the Dalitz method in the LHCb experiment at CERN and photon reconstruction optimisation for the LHCb detector upgrade. PhD Thesis, Université Paris Sud-Paris XI (2015)