

# Design and Evaluation of a GNN Algorithm for the ECL L1 Trigger Upgrade at Belle II

Thomas Lobmaier

Masterthesis

2nd March 2026

Institute of Experimental Particle Physics (ETP)

Reviewer:	Prof. Dr. Torben Ferber
Second Reviewer:	Prof. Dr. Markus Klute
Advisor:	Dr. Isabel Haide

Editing time: 1st March 2025 – 2nd March 2026



# Entwurf und Evaluation eines GNN Algorithmus für das ECL L1 Trigger Upgrade von Belle II

Thomas Lobmaier

Masterarbeit

2. März 2026

Institut für Experimentelle Teilchenphysik (ETP)

Referent: Prof. Dr. Torben Ferber  
Korreferent: Prof. Dr. Markus Klute  
Betreuer/in: Dr. Isabel Haide

Bearbeitungszeit: 1. März 2025 – 2. März 2026



---

Ich versichere wahrheitsgemäß, die Arbeit selbstständig angefertigt, alle benutzten Hilfsmittel vollständig und genau angegeben und alles kenntlich gemacht zu haben, was aus Arbeiten anderer unverändert oder mit Abänderungen entnommen wurde.

**Karlsruhe, 2. März 2026**

.....  
(**Thomas Lobmaier**)



# Disclaimer

Data analyses in high-energy physics, such as the topics presented in this master's thesis, are a collaborative effort. The SuperKEKB particle accelerator, which provides the particle beams essential for all studies at Belle II was built and is operated and maintained by the SuperKEKB accelerator group. The Belle II detector was built and is maintained and operated by the Belle II collaboration. The Belle II collaboration also creates the simulated and recorded data sets and maintains the computing infrastructure necessary to process them. The software environment necessary for studies with Belle II data plays an important role and was created and is maintained by the collaboration.

The GNN-ETM reconstruction algorithm, as well as the CaloClusterNet was designed by Isabel Haide and implemented on FPGA hardware by Marc Neu. The model used for evaluating the GNN-ETM is the original model, which was designed and trained by Isabel Haide. The hardware implementation of the final quantised model proposed in this thesis was implemented by Marc Neu, based on the previous work of Till Raedler and Fabio Papagno.

For this thesis, I use and adapt the training and evaluation framework created by Isabel Haide. This includes the software implementation of the GNN-ETM, the simulated training and evaluation samples, as well as the training and evaluation scripts. Own Monte Carlo samples are created using the simulation pipeline provided by Isabel Haide.

## Artificial Intelligence Tools

This thesis incorporates the use of Artificial Intelligence (AI) tools to check for grammatical or spelling mistakes, as well as program code creation. No rephrasing or stylistic adaptations were made based on AI tools.

Grammarly<sup>1</sup> is utilised throughout the thesis for spelling and grammar checks. I have approved all suggested changes.

GitHub Copilot<sup>2</sup> is used to aid the development of Python and HLS code, in particular, code restructuring and optimisation that does not constitute the core scientific work of this thesis. I have approved and tested all suggestions to provide robust and reliable results.

---

<sup>1</sup><https://www.grammarly.com/> (Access date: 21st February 2026)

<sup>2</sup><https://github.com/features/copilot> (Access date: 21st February 2026)



# Contents

<b>1. Introduction</b>	<b>1</b>
<b>2. Belle II</b>	<b>3</b>
2.1. SuperKEKB Collider and Belle II Detector . . . . .	3
2.2. Electromagnetic Calorimeter . . . . .	4
2.3. Beam Background . . . . .	7
<b>3. Trigger</b>	<b>9</b>
3.1. The Level 1 Trigger . . . . .	9
3.2. FPGAs as L1 Trigger Hardware . . . . .	10
3.3. The Electromagnetic Calorimeter L1 Trigger Pipeline . . . . .	11
3.4. Drawbacks of the ICN-ETM . . . . .	12
<b>4. Long Shutdown 2</b>	<b>19</b>
4.1. Planned SuperKEKB Upgrades . . . . .	19
4.2. Estimated Run Conditions after LS2 . . . . .	20
4.2.1. Beam Background Extrapolation . . . . .	20
4.2.2. L1 Trigger Rate Extrapolation . . . . .	20
4.3. Belle II Upgrade Plans . . . . .	22
4.4. Upgrades regarding the ECL L1 trigger pipeline . . . . .	22
<b>5. Network Design</b>	<b>25</b>
5.1. CaloClusterNet . . . . .	26
5.1.1. GravNet . . . . .	26
5.1.2. ObjectCondensation . . . . .	26
5.1.3. Architecture . . . . .	27
5.2. Integration as GNN-ETM . . . . .	28
5.2.1. Quantisation . . . . .	28
5.3. Adaptation for the LS2 Upgrade . . . . .	30
<b>6. Datasets and Metrics</b>	<b>33</b>
6.1. Preprocessing . . . . .	33
6.1.1. Flat Energy Cut . . . . .	34
6.1.2. Tower Cut . . . . .	34
6.2. Training Targets and Training Dataset . . . . .	35
6.2.1. Training Targets . . . . .	35

6.2.2.	Training Dataset . . . . .	37
6.3.	Metrics . . . . .	38
6.3.1.	Cluster-Finding Metrics . . . . .	38
6.3.2.	Predictive Performance Metrics . . . . .	39
6.3.3.	Signal Background Classifier . . . . .	41
6.4.	Evaluation Datasets . . . . .	42
6.4.1.	Single Photon Sample . . . . .	42
6.4.2.	Bhabha Sample . . . . .	43
6.4.3.	Dimuon Sample . . . . .	44
6.4.4.	$B\bar{B}$ Sample . . . . .	47
<b>7.</b>	<b>Evaluation</b>	<b>53</b>
7.1.	Model Training . . . . .	53
7.1.1.	Cluster Cut Selection . . . . .	54
7.2.	Single Photon Sample . . . . .	57
7.2.1.	Granularity Comparison . . . . .	57
7.2.2.	Quantisation Comparison . . . . .	59
7.3.	Bhabha Sample . . . . .	63
7.4.	Dimuon Sample . . . . .	70
7.5.	$B\bar{B}$ Sample . . . . .	70
7.5.1.	Granularity Comparison . . . . .	74
7.5.2.	Quantisation Comparison . . . . .	77
7.5.3.	Signal/Background Classification Performance . . . . .	78
<b>8.</b>	<b>Conclusion</b>	<b>83</b>
<b>A.</b>	<b>Appendix</b>	<b>89</b>
A.1.	Single Photon Sample . . . . .	89

# 1. Introduction

One of the main challenges of the Belle II experiment and the SuperKEKB collider is to further increase the instantaneous luminosity and the corresponding dataset. Only this enables the multitude of precision measurements at the forefront of particle physics.

Due to the low cross-section, an actual collision only takes place in a small subsection of bunch crossings. Furthermore, in the majority of collisions, Bhabha scattering occurs, which is not relevant for the physics goals of Belle II. To prevent the analysis and storage of empty or uninteresting events for the physics program, a filter is required. This filtering is called triggering and is performed in two steps: the hardware-based Level 1 Trigger (L1 trigger) followed by the software-based High Level Trigger (HLT). The trigger identifies interesting collisions and triggers the data acquisition and storage of events.

With regard to increasing the overall dataset, efforts are being made to further increase the instantaneous luminosity of the SuperKEKB collider. Namely, during the Long Shutdown 2 (LS2) planned for 2032, there are many hard- and software upgrades planned, both for the accelerator and the detector [1]. The expected increase in luminosity directly leads to an increase in trigger rates [2]. As the extrapolations suggest that the combined trigger rates of the L1 trigger might exceed the limit of 30 kHz, either the current trigger parameters have to be tightened, or the trigger pipeline has to be adapted. As the first option would directly harm the physics program, alternative ways of reducing the trigger rate have to be studied.

The L1 trigger runs optimised algorithms for a preliminary reconstruction of the events on Field Programmable Gate Array (FPGA)s and has to make a real-time decision to discard or save an event. The task of the L1 trigger is to reduce the event rate to the maximal processing rate of the HLT, while keeping the interesting events. The different sub-detectors of the experiments have individual trigger pipelines, which reconstruct high-level objects, which are then passed to a centralised decision logic, where they are combined, and the final trigger decision is made.

For the current Electromagnetic Calorimeter (ECL) L1 trigger, a fast pattern matching algorithm is deployed [3], which identifies connected energy depositions in the detector. In order to accelerate and simplify this, each group of  $4 \times 4$  crystals is combined into a so-called Trigger Cell (TC). This automatically reduces the complexity of the clustering task, at the cost of a reduced resolution of both the energy and the position. Additionally, the resolvability of nearby or overlapping clusters is reduced.

Recent efforts were made to replace the cluster finding algorithm with a Graph Neural Network (GNN) based clustering algorithm, called the CaloClusterNet [4, 5], which is based on the ObjectCondensation algorithm [6] and GravNet architecture [7].

In this thesis, I describe the design and evaluation of a high-granularity adaptation of this GNN approach, as a planned electronics upgrade enables the use of single crystal information for the ECL L1 trigger pipeline. I further present two input reduction method schemes and two corresponding versions of the proposed model and evaluate them on Monte Carlo simulation (MC)-samples against each other, the current clustering algorithms and the offline reconstruction. To prove the feasibility of this approach, I further show a model that is implemented on hardware.

In chapter 2, an overview of the Belle II experiment and the ECL is given, followed by a detailed description of the L1 trigger in chapter 3 and analysis of the existing ECL L1 trigger. Subsequently, in chapter 4, the plans for the upcoming LS2 upgrade are presented. In chapter 5, the existing CaloClusterNet approach and its implementation as part of the L1 trigger are presented. Furthermore, I present a fully implemented variation of the CaloClusterNet designed specifically for the changes anticipated in the LS2 upgrade. In chapter 6 I define the preprocessing, the proposed input reduction methods, the training samples and targets, as well as the evaluation metrics and samples, followed by the training and evaluation of the model in chapter 7.

## 2. Belle II

Belle II is a high-precision and intensity experiment located at the SuperKEKB electron-positron collider in Japan. For the overarching objective of collecting data for high-precision physics analyses, the accelerator is tuned to the  $\Upsilon(4s)$  resonance. Combined with the known initial state variables, this provides a very clean experimental environment, as the resonance predominantly decays into  $B\bar{B}$  pairs. The study of these  $B\bar{B}$  pairs enables high-precision physics analyses in the flavour and dark sector. This chapter outlines the experimental setup with a focus on the ECL and an overview of the sources of beam background.

### 2.1. SuperKEKB Collider and Belle II Detector

SuperKEKB accelerates electron and positron bunches to asymmetric energies of 7 and 4 GeV respectively and collides them at the interaction point within the Belle II detector. Both rings have a circumference of roughly 3 km and are called high energy ring and low energy ring, shown in Fig. 2.1.

The asymmetric energies are chosen in order for the resulting  $B\bar{B}$  pairs formed by the  $\Upsilon(4s)$  to be boosted in the laboratory frame of the detector. This boost leads to a spatial separation of the decay vertices, based on the otherwise unresolvable time-dependent evolution of the B-mesons. Due to the expected asymmetric distribution of the decay products, the Belle II detector is designed asymmetrically accordingly. Belle II is a cylindrical  $4\pi$  detector comprised of stacked sub-detector layers, shown in Fig. 2.2. The coordinate system is right-handed, with the z-axis pointing in the flight direction of the electron bunches. The x-axis points horizontally outwards while the y-axis points vertically upwards. The polar angle  $\theta$  is defined from  $0^\circ$  along the z-axis to  $180^\circ$ . The azimuthal angle  $\phi$  ranges from  $-180^\circ$  to  $180^\circ$ , with  $0^\circ$  being defined as the direction of the x-axis at  $y = 0$ .

The detector is designed for high precision particle detection, reconstruction and identification in mind. For this purpose, the particle response of the different sub-detectors has to be measured, combined and fundamental properties, like their position, energy, momentum and type of particle, have to be reconstructed. The innermost part of the detector is the silicon-based Pixel Detector (PXD) [9]. Combined with the enclosing strip detector, called the Silicon Vertex Detector (SVD) [10], these subdetectors are responsible for increasing the vertex position reconstruction. While the main particle tracking is performed by the gas-filled Central Drift Chamber (CDC) [11], which enables the reconstruction of charged particle tracks, their momenta and particle identification. Surrounding the CDC, the

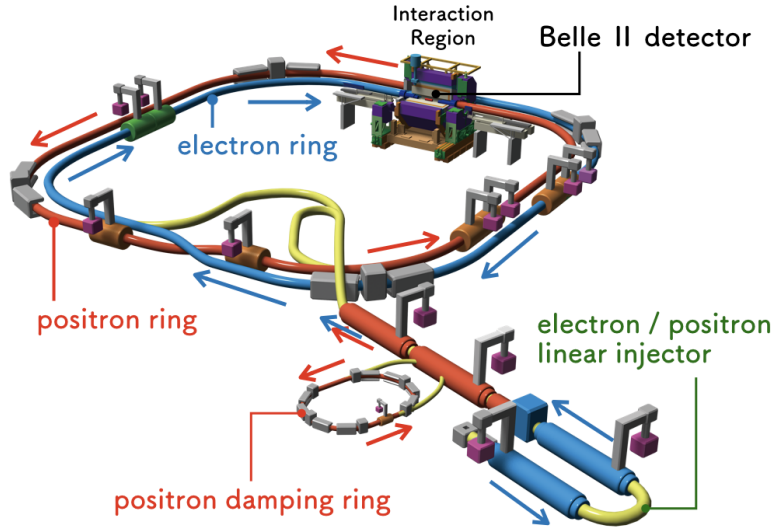


Figure 2.1.: Overview of the SuperKEKB Accelerator. Figure taken from [8].

particle identification system is placed. In the barrel, the Time-Of-Propagation (TOP) [12] detector and in the forward endcap the Aerogel Ring-Imaging Cherenkov (ARICH) [13] detector exploit the emittance of Cherenkov light to boost the particle identification (PID) capabilities. Encapsulating the PID system, the ECL [14] detects energy depositions in the scintillation crystals. Outside of the ECL and the enclosing superconducting solenoid, the  $K_L^0$  and  $\mu$  detector (KLM) [15] is located, which provides additional material for  $K_L^0$  to shower hadronically and act as magnetic flux return for the solenoid as well as for the detection and identification of muons [16, 17].

To ensure a consistent and optimised reconstruction of detector data, the centralised Belle II Analysis Framework (basf2) is provided by the collaboration. Apart from reconstruction, it also incorporates the necessary tools for MC simulation and analyses [18, 19].

## 2.2. Electromagnetic Calorimeter

The ECL serves multiple objectives. Primary goals are to detect and measure the energy of photons, identify electrons and other particles like neutral hadrons. Additional tasks are the luminosity determination, event triggering, determination of the rest of event variables and background suppression.

In order to determine the energy of the incident particles, the ECL is comprised of a high-density material that provokes an inelastic interaction with the primary particle. Due to the interaction, a particle shower is induced, resulting in numerous low-energy particles. These low-energy shower constituents can then cause the emission of detectable scintillation light. Depending on the particle type and energy, the deposition characteristics vary. Minimal ionising particles, e.g. muons, do not interact inelastically with the detector and therefore only generate a small and highly localised signal. The electromagnetic calorimeter in the

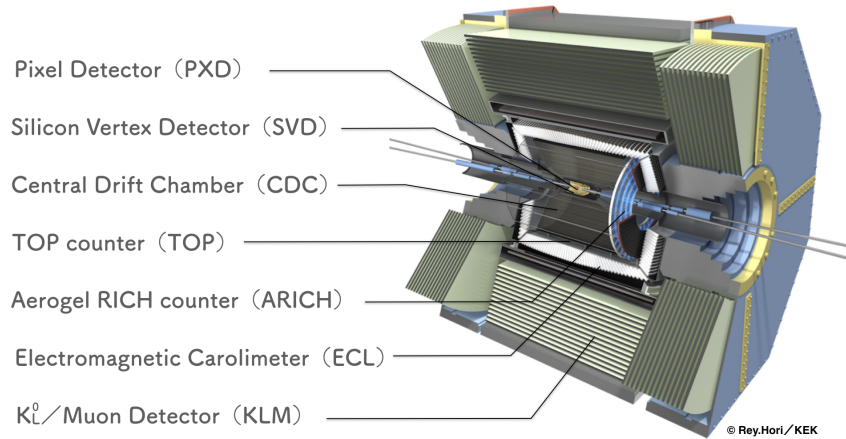


Figure 2.2.: Overview of the Belle II detector and its subdetectors. Figure taken from [8].

Belle II detector consists of 8736 CsI(Tl) crystals, distributed over a cylindrical surface. Hence, the detector can be divided into 3 detector regions: the barrel region and the two endcaps in the forward and backward direction, with respect to the electron direction. A schematic of the detector is shown in fig. 2.3. The holes in the endcaps accommodate the collimator magnets required to achieve the high instantaneous luminosity. Therefore, the angular coverage of the ECL is between  $12.4^\circ$  and  $155.1^\circ$ , with additional  $1^\circ$  gaps between the endcaps and the barrel. In order to achieve a consistent coverage with reduced gaps and dead material across the whole detector, the crystals have varying geometries. For all visualisation purposes, the fronts of the crystals are simplified to trapezoids. The crystals have an average cross-section of roughly  $6 \times 6 \text{ cm}^2$  and a uniform length of 30 cm [16].

Two photodiodes are attached to each crystal for an independent scintillation light readout. After an amplification and digitisation step, the two signal waveforms are summed and fitted on a subsequent shaper board to derive timing and energy information for each crystal. For crystals exceeding 30 MeV, the raw waveform is also saved, enabling an offline pulse-shape analysis. This pulse shape discrimination (PSD) is used for the discrimination between hadronic and electromagnetic showers [20].

As the particle showers induced by the primary particles are, in general, not contained within a single crystal, the signal of multiple crystals has to be combined to gain information on the underlying particle shower. Due to the possibility of overlapping showers and the additional complication of background contributions, the assignment of crystals to clusters is not straightforward. The cluster reconstruction used for the Belle II analyses is done by the central Belle II software framework basf2, and is described in more detail in [4]. In the first step, the algorithm combines interconnected crystals above a certain threshold to Connected Region (CR)s. Inside these connected regions, crystals are identified as Local Maximum (LM) if their energy surpasses that of all directly neighbouring crystals and an additional energy cut. For this thesis, this energy cut is set to 20 MeV. In order to attribute the energy of surrounding crystals to the LMs, an iterative reweighing within a CR is performed. Resulting in optimised weights for each LM and all associated crystals.

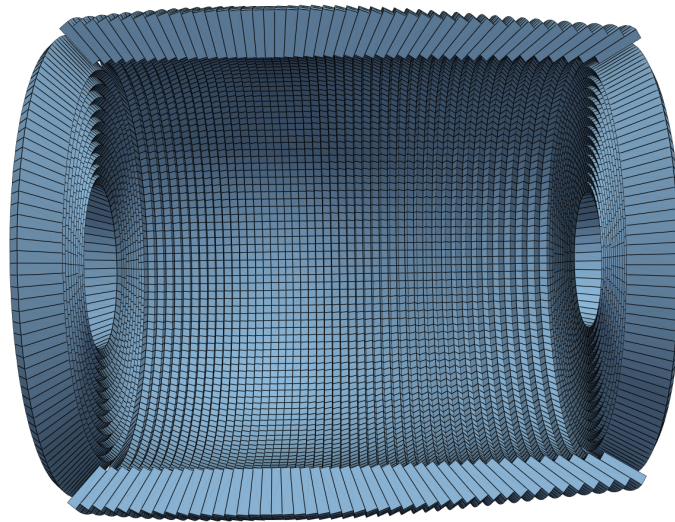


Figure 2.3.: 3D rendering of the Belle II electromagnetic calorimeter. On the left side, the backward endcap is visible. In the middle, the Barrel region is shown. On the right side, the forward endcap is depicted. To house cables and electronics as well as structural components, there are small gaps between the barrel and endcap regions.

For the energy estimation of the cluster, the weighted sum of crystals within a  $5 \times 5$  window is used. The position is the corrected position of the LM, and the timing of the cluster is determined by the timing of the highest energetic crystal. After the clustering is completed, the cluster energies are additionally corrected to account for longitudinal leakage, as the clusters are, in general, not fully contained within the crystal lengths. The whole algorithm is once performed by the HLT and a second time after data calibration, for the actual physics analyses. The thresholds and parameters of the clustering are tuned to the specific background conditions to ensure a consistent and stable cluster reconstruction [4, 16].

## 2.3. Beam Background

Instead of only detecting the desired signatures of particles from the collision, the detector also measures beam background. Beam Background is a collective term describing multiple sources of detector response which are caused by beam particles changing their direction or energy. Thus, leaving the nominal beam orbit and interacting with inactive detector or accelerator material, causing particle showers, which can reach the active detector volume. It is vital to identify the beam background, as it can mimic signatures from actual collisions and therefore lead to misreconstructions and a reduced overall resolution power.

Beam particles can Coulomb scatter with other particles in the same bunch, leading to the Touscheck effect, where the Due to the momentum transfer caused by this scattering, both scattering partners deviate from the nominal beam energy and path. The Touschek scattering rate is inversely proportional to the total number of bunches in the ring and their respective beam size, while being proportional to the beam current squared. The beam particles can also interact with the residual beam gas, either via Coulomb scattering or the emittance of Bremsstrahlung. Leading to a change in direction or a loss of energy. This beam-gas background is proportional to the residual beam gas pressure and the beam current. As the beam particles experience a constant perpendicular acceleration towards the centre of the ring, they emit synchrotron radiation. These resulting synchrotron photons predominantly affect the inner tracking detectors and are mainly caused by the higher energetic electron beam, as the power of the synchrotron radiation scales with the energy squared. In order to provide a stable beam current, new particles are injected in top-up injections. As these newly injected bunches are perturbed regarding the optimal beam orbit, they oscillate around the main beam until they stabilise. During this time, they cause much higher background rates in the detector, leading to a time dependent injection-background.

Apart from the background sources caused by the single beams, there is additional background directly generated by the beam collisions at the interaction point. Consequently, it is proportional to the luminosity and called luminosity background. Bhabha scattering,  $e^+e^- \rightarrow e^+e^-(\gamma)$ , is a major luminosity background channel. If a radiative photon is emitted, the energy of the electron and positron is reduced, and a fraction of those lower-energy particles is then lost in the detector. Even if no radiative photon is emitted, the scattered electron and positron can still interact with the beam pipe, causing electromagnetic showers visible at the edge of the detector. Superimposing that, the radiative photons along the beam line can produce neutrons in the accelerator magnets, which can backscatter towards the detector. Apart from the radiative Bhabha scattering, another prevalent background effect is caused by the pair production of low-energy electron-positron pairs,  $e^+e^- \rightarrow e^+e^-e^+e^-$ .

At the target luminosity of SuperKEKB, the luminosity-driven background is expected to dominate the background occupancy for the ECL [21].

As these background effects can mimic actual collision signatures and swamp the data acquisition system, it is vital to perform a filtering of the recorded detector response.

## 3. Trigger

Analogously to other collider experiments, Belle II deploys a trigger system to filter the recorded detector data prior to permanent storage, reducing unnecessary overhead. This event filtering is divided into two successive systems. The first is the L1 trigger, which applies highly accelerated reconstruction algorithms to filter the events. For event windows, triggered by the L1 trigger, the subsequent HLT runs a full reconstruction, applying the basf2 reconstruction algorithms. Only events passing both filtering steps are permanently stored, calibrated and reconstructed to be available for subsequent Belle II physics analyses. In this chapter, I give an overview of the L1 trigger, the deployed hardware, so-called FPGAs and the trigger pipeline dedicated to the ECL. As the current clustering algorithm exhibits irregularities, I perform a dedicated analysis regarding these.

### 3.1. The Level 1 Trigger

During operation, the detector data is constantly fed into a ring buffer, which stores the data until a decision is made by the L1 trigger to pass on an event to the HLT or the data is overwritten by new detector data. Due to the limited size of this buffer, this decision has to be made within  $5\ \mu\text{s}$  or the event is lost. For the L1 trigger, the detector data is gathered in time sections, so-called trigger windows. The detector response of the CDC, the ECL, the KLM and TOP detectors in each trigger window is digitised and fed into four independent trigger pipelines. In the case of the ECL, the trigger windows are 250 ns long. These individual trigger chains reconstruct high-level objects like charged particle tracks out of hits in the CDC or clusters out of energy depositions in the ECL. Subsequently, these objects are combined and assessed against set conditions to form so-called trigger input bits. Simple bits, requiring only information of a single sub-detector, are calculated within the trigger chain, while bits combining multiple sub-detectors are determined by the Global Reconstruction Logic (GRL). In the final step, the Global Decision Logic (GDL) gathers all input bits and combines them to trigger output bits by deploying a combinational logic. For trigger bits, with rather lenient requirements, an additional prescale can be applied. Effectively triggering only every  $n$ -th event, where this trigger bit would otherwise be active. Or else these trigger bits would exceed the maximum trigger rate, but are vital for some analyses and also for validating the trigger system. Whenever one of the prescaled trigger output bits is active, the event is passed on to the HLT [22, 23].

### 3.2. FPGAs as L1 Trigger Hardware

To make the trigger decision in time, performant algorithms are needed, which, as well, need to be implemented on the right computational hardware. As general-purpose devices like CPUs or GPUs do not reach the tight latency requirement for the L1 trigger and ASICs are not adaptable enough to be used, the L1 trigger is computed on FPGAs. FPGAs are standardised for the whole L1 trigger and are designated as Universal Trigger Board (UT). Currently, the third and fourth iterations of the UT are deployed with plans for a timely upgrade to the fifth generation in place. The UT4 features the Xilinx Ultrascale XCVU080/160 FPGA.

Contrary to CPUs or GPUs, FPGAs are inherently parallelised and do not process a sequential instruction list. Instead, they process all of the input simultaneously with logic functions. This removes the requirement of saving the data to a centralised memory, which causes a large latency overhead for CPUs or GPUs. FPGAs are based on the principle that any Boolean logic can be represented by a combination of a basic set of logic operations, also called logic gates. Examples for these are the AND and OR operations. On a hardware level, these gates are constructed by connecting multiple transistors. The FPGAs are an accumulation of a large number of basic logic gates that are linked by programmable connections. As different logical operations require different amounts of time, a periodic clock signal is also applied at which the logic gates are evaluated. Theoretically, a multitude of NOR or NAND gates would be sufficient to enable the encoding of any arbitrary function. However, this leads to a large overhead of required resources, which can be reduced by using more complex but still basic building blocks like flip-flops or look-up tables. A flip-flop is a single-state memory, which is reevaluated with each clock-cycle. Look-up table (LUT)s are more advanced, as they store a custom logic table, which encodes an arbitrary logic function. For the FPGAs relevant in this thesis, a multitude of LUTs and flip-flops are combined into configurable logic blocks. Apart from these logic blocks, there are additional dedicated memory banks, consisting of RAM cells and serving as buffers between calculations. To further boost the capabilities of the FPGA, especially with regard to arithmetic functions, there are additional Digital Signal Processor slices. These are predefined circuits to efficiently perform multiplications of two numbers. The programmability of the connections between the components makes FPGAs adaptable, compared to ASICs, ensuring that adaptations of the L1 trigger are possible.

Within recent years, AMD pushed towards the development of innovative computing architectures on FPGAs, resulting in the AMD Versal AI Core Series [24]. The series combines the traditional programmable logic with a processor system and novel AI Engine (AIE)s, all contained within a single System on Chip design. The different components are interconnected by a Network on Chip (NoC). Additionally, there exists a faster interface between the programmable logic and the AIEs. The AIEs are composed of individual AIE-Core, which are interconnected in a grid-like structure, forming the AIE-array. Each AIE-Core is a single instruction multiple data very long instruction word processor. This means that every operation is composed of multiple instructions and applied to multiple different inputs simultaneously. The available instructions contain move, load, store, scalar and vector instructions, enabling AIEs to efficiently calculate matrix multiplications. Furthermore, each AIE-Core has in total 4 attributed memory banks, which can be accessed

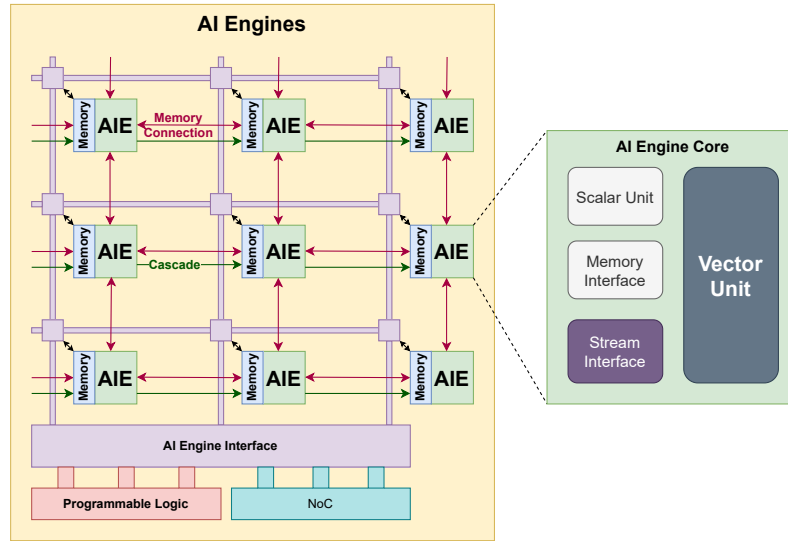


Figure 3.1.: Simplified visualisation of the inner structure of the AIEs. They consist of interconnected AIE-Cores, which are especially suited for efficient matrix multiplications.

by the associated and neighbouring AIE-Cores. A visualisation of the AIE structure is shown in fig. 3.1. For future upgrades of the L1 trigger hardware, the AMD Versal AI Core Series is a valid possibility for boosting the capability of the L1 trigger.

### 3.3. The Electromagnetic Calorimeter L1 Trigger Pipeline

The goal of the ECL Trigger pipeline is to find and reconstruct clusters and the total energy in the calorimeter. Based directly on the ECL trigger of Belle, the current ECL trigger is depicted in fig. 3.2. The readout signal of 16 neighbouring crystals is summed to form Trigger Cells. This reduction of input information is done to simplify and accelerate the trigger pipeline by reducing background contributions and pre-clustering the input data.

TCs surpassing an energy threshold of 100 MeV are passed to the ICN-ETM. The clustering is performed by applying a  $3 \times 3$  pattern-matching algorithm. Figure 3.3 displays the pattern logic. In the first step for each active TC, the pattern logic is evaluated. If the logic is fulfilled, the TC is classified as an ICN-hit. The logic is comprised of three conditions:

- The central TC (TC0) has to be active.
- The TC directly above (TC1) and the TC directly to the left (TC2) must not be active.
- At maximum, one of the lower left (TC3) and the central lower (TC4) TCs is active.

To summarise this, this logic yields in general one ICN-hit for TCs which are interconnected by shared edges. For up to 6 individual ICN-hits, an ICN-cluster is created by shifting the  $3 \times 3$  evaluation window to be centred around the previously highest energetic TC. All TC

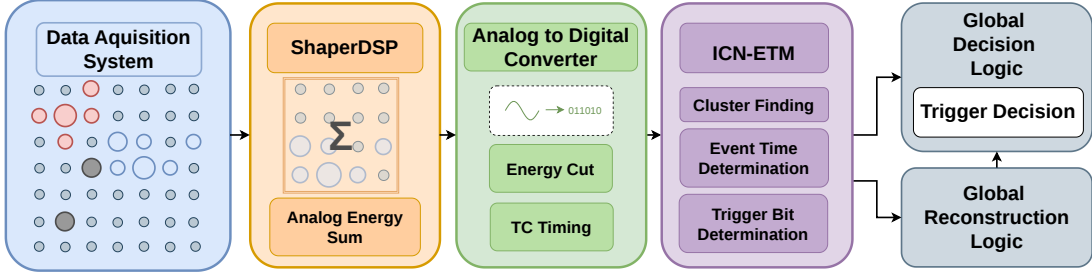


Figure 3.2.: Simplified overview over the ECL L1 trigger pipeline. The response of all crystals within a trigger window is analogously summed into TCs, and after digitization a 100 MeV cut is applied. The TCs of the different detector regions are merged and passed on to the Isolated Cluster Number ECL Trigger Module (ICN-ETM), where the clustering is performed. Additionally, trigger bits, which only depend on the ECL information, are formed. The resulting trigger bits are passed to the GDL, and all reconstructed clusters are passed to the GRL, where trigger bits based on combinations of the different sub-detectors are determined and relayed to the GDL. This figure is adapted from [4].

energies covered by this adjusted  $3 \times 3$  window are summed up, and the position of the new central TC is the ICN-cluster energy and position.

Out of these cluster predictions, the ECL stand-alone trigger bits can directly be determined, while matching trigger bits are determined by the GRL. The cluster bits are often limited to certain detector regions, represented by the  $\theta$ -ID, which encodes the  $\theta$ -value of a corresponding TC. A detailed definition of these can be found in [4].

### 3.4. Drawbacks of the ICN-ETM

The ICN-ETM proves quite robust and, due to the simplicity, also very fast. However, there are multiple limitations: First of all, a separation of interconnected TCs into multiple ICN-clusters is not possible. Secondly, the energy and position predictions are not consistent, as they are biased towards the lower left of interconnected TCs. Additionally, if the  $3 \times 3$  evaluation windows of multiple ICN-clusters overlap, they can share common TCs leading to an overall increase of the total predicted energy compared to the deposited energy. A special case of this occurs whenever two diagonally adjacent TCs are active  $\begin{smallmatrix} \blacksquare & & \blacksquare \\ & & \end{smallmatrix}$ . Both TCs fulfil the hit logic and therefore get shifted towards the same higher energetic TC. Effectively causing two identical ICN-clusters, both position and energy-wise. Lastly, the ICN-ETM is not robust under rotational transformations. More specifically, whenever the pattern of at least four active TCs contains the  $\begin{smallmatrix} \blacksquare & \blacksquare \\ \blacksquare & \blacksquare \end{smallmatrix}$  pattern, an additional ICN-hit will be created. However, if this pattern is rotated or mirrored, it does not yield the additional prediction.

In order to evaluate the extent of those effects, I perform dedicated studies on the occurrence of different TC patterns, in different MC-samples. The summarised result of this study is included in [5]. As the predictions of the ICN-ETM depend on the pattern of active TCs, these patterns will be used to analyse the algorithm. Due to the locality induced by the  $3 \times 3$  evaluation window, it is sufficient to study patterns, which are surrounded by one

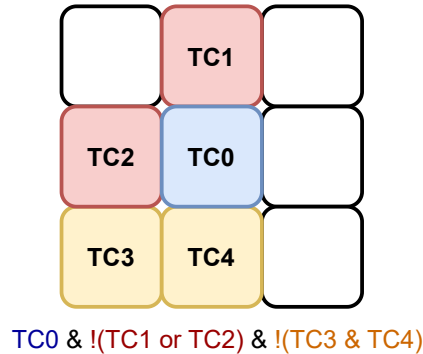





Figure 3.3.: ICN-ETM pattern logic. This is evaluated for all active TCs, and if fulfilled, this TC is classified as an ICN-hit.

layer of inactive TCs. To further reduce the complexity of this analysis, the shifting to the highest energetic TC is omitted in the first step, as it does not change the number of cluster predictions, but rather the amount of energy that is shared.

In a second consideration, the extent of different TC patterns has to be set. It is apparent that the occurrence of patterns drastically reduces with their size. Hence, the patterns studied for this analysis are limited to a total of four TCs. For each possible configuration, with up to four TCs, a mask is defined containing each possible rotational or mirrored version of the pattern. With this mask, the absolute occurrence of the patterns is studied. An additional point of interest is how many clusters are actually contained within the pattern. Ideally, the number of ICN-clusters matches the number of offline clusters obtained by the basf2 reconstruction. The number is obtained by counting the clusters, which deposit the main energy contribution within a crystal, that is summed into an active TCs. Due to large background contributions, this number can also be zero.

With multiple MC-samples, the behaviour of the ICN-ETM is probed. As the extrapolated beam background levels after LS2 are relevant for the future performance of the algorithm, the background levels for the MC-simulation are chosen as such. A detailed definition of the overlaid background files is given in section 4.2.1. The first sample is a generic  $e^+e^- \rightarrow B\bar{B}$  sample and the corresponding results are shown in table 3.1. For each pattern, the different orientations are enumerated, and the number of offline clusters depositing energy within the pattern is broken down accordingly. As stated above, a clear tendency towards small patterns is visible. Additionally, a dependence of the orientation is clearly visible for patterns like . This can be explained by the fact that due to the boost of the centre of mass frame, elongated energy depositions along the  $\theta$  direction are more probable compared to the  $\phi$  direction. Furthermore, it becomes apparent that for small patterns, the number of ICN-clusters matches well with the number of offline clusters. For larger patterns, however, this is not always the case. As the  pattern leads to only one ICN-cluster even though it is primarily produced by two offline clusters.

Even though the  pattern is not that common, it still occurs in roughly one third of the studied events. As this pattern results in two identical ICN-clusters with shared energy,

this is a major design flaw, leading to overestimations of energies. Whereas, the issue of the rotational differences of the  $\begin{smallmatrix} \blacksquare \\ \blacksquare \end{smallmatrix}$  pattern arises much less frequently. And therefore should not affect the L1 trigger in a notable way.

As a second MC-sample, a  $e^+e^- \rightarrow e^+e^-(\gamma)$  sample, also with high beam background conditions, is simulated and studied. This channel represents a control channel, as the detector is expected to be much less populated, while the energies of the expected clusters reach maximal energies. Overall, the same tendency towards smaller TC-patterns is visible. Mentionable is the discrepancy between the  $\begin{smallmatrix} \blacksquare \\ \blacksquare \end{smallmatrix}$  and  $\begin{smallmatrix} \blacksquare \\ \blacksquare \end{smallmatrix}$  pattern, especially compared to the  $B\bar{B}$ -sample. Additionally, the  $\begin{smallmatrix} \blacksquare \\ \blacksquare \end{smallmatrix}$  pattern is predominantly composed of only one offline cluster. This can be explained by the much higher expected cluster energies and, therefore, also cluster sizes, caused by the comparable high energetic electron and positron energies. Due to the drastically lower occurrence of the  $\begin{smallmatrix} \blacksquare \\ \blacksquare \end{smallmatrix}$  pattern, the energy sharing is less of a problem in the Bhabha events.

To study the effect of overlapping clusters, a set of MC-samples containing  $e^+e^- \rightarrow a(a \rightarrow \gamma\gamma)$  events is considered. The  $a$  denotes an Axion like particle (ALP), which decays into two photons. The lifetime of the ALP is inversely proportional to its mass. Therefore, a lighter ALP results in two clusters generally closer together in comparison to a heavier ALP. To sample different cluster distances, four different alp masses, namely 200 MeV, 250 MeV, 300 MeV and 400 MeV are simulated. Exemplary, the results of the 300 MeV sample are shown in table 3.3. As the samples created contain fewer events, the statistics are limited, and an absolute number comparison to the other samples is not possible. To note here, the relative occurrence of some larger patterns like  $\begin{smallmatrix} \blacksquare \\ \blacksquare \end{smallmatrix}$  and  $\begin{smallmatrix} \blacksquare \\ \blacksquare \end{smallmatrix}$  are increased, while other large patterns like  $\begin{smallmatrix} \blacksquare \\ \blacksquare \end{smallmatrix}$  are not. Due to the number of contained offline clusters, this can be explained by the expected separation of the two photon clusters of the alp, which seldom create a double diagonal pattern. The  $\begin{smallmatrix} \blacksquare \\ \blacksquare \end{smallmatrix}$  and  $\begin{smallmatrix} \blacksquare \\ \blacksquare \end{smallmatrix}$  pattern incorporates the case, where the two photon clusters are still very close to each other and are not separated by the ICN-ETM. The energy sharing due to the  $\begin{smallmatrix} \blacksquare \\ \blacksquare \end{smallmatrix}$  pattern is also more pronounced in this sample. Hence, the ICN-ETM will predict in these cases two clusters with the total ALP energy.

Summarising the above, the energy sharing is a more pronounced problem than the much rarer case of the additional predictions due to the  $\begin{smallmatrix} \blacksquare \\ \blacksquare \end{smallmatrix}$  pattern. This is mainly attributed to the suppression of large TC patterns. More importantly, this analysis underlines that a simple pattern-matching approach will never be able to consistently predict the correct number of clusters. Simply due to the fact that the correct number of clusters is strongly dependent on the event topologies and underlying physics processes, and not the TC patterns. This directly motivates the study and development of alternative clustering algorithms, which can separate interconnected TC patterns.

Table 3.1.: Absolute occurrence of different TC-patterns in a high background  $B\bar{B}$  MC-sample. The orientation denotes the different possible configurations by rotation and mirroring. The number of offline clusters denotes the number of basf2 clusters contained by the active TCs.






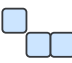
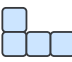


Pattern	Orientation	Number of all Offline Clusters			
		0	1	2	>2
	0	3099	66175	1304	13
	$\Sigma$	<b>3099</b>	<b>66175</b>	<b>1304</b>	<b>13</b>
	1	41	4300	3222	207
	2	17	2434	2045	117
	$\Sigma$	<b>58</b>	<b>6734</b>	<b>5267</b>	<b>324</b>
	1	11	208	1529	63
	2	5	231	1471	58
	$\Sigma$	<b>16</b>	<b>439</b>	<b>3000</b>	<b>121</b>
	1	2	81	227	119
	2	4	83	247	133
	3	3	76	210	95
	4	5	79	236	137
	$\Sigma$	<b>14</b>	<b>319</b>	<b>920</b>	<b>484</b>
	1	1	6	132	68
	2	3	11	261	160
	$\Sigma$	<b>4</b>	<b>17</b>	<b>393</b>	<b>228</b>
	1	0	5	44	54
	2	0	1	52	52
	3	0	2	57	40
	4	0	2	60	51
	5	0	9	88	83
	6	2	4	88	82
	7	0	7	99	88
	8	0	4	95	72
	$\Sigma$	<b>2</b>	<b>34</b>	<b>583</b>	<b>522</b>
	1	0	0	4	22
	2	0	0	9	20
	3	0	0	5	13
	4	0	0	2	12
	5	0	0	6	13
	6	0	0	10	20
	7	1	0	4	16
	8	0	1	7	18
$\Sigma$	<b>1</b>	<b>1</b>	<b>47</b>	<b>134</b>	
	1	0	0	9	43
	2	0	1	7	23
	3	0	0	10	31
	4	0	2	9	41
	$\Sigma$	<b>0</b>	<b>3</b>	<b>35</b>	<b>138</b>
	1	0	0	14	30
	2	0	0	7	35
	$\Sigma$	<b>0</b>	<b>0</b>	<b>21</b>	<b>65</b>

Table 3.2.: Absolute occurrence of different TC-patterns in a high background Bhabha MC-sample. The orientation denotes the different possible configurations by rotation and mirroring. The number of offline clusters denotes the number of basf2 clusters contained by the active TCs.



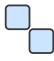


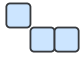
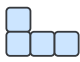
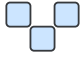






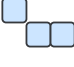
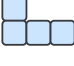


Pattern	Orientation	Number of all Offline Clusters			
		0	1	2	>2
	0	4337	97123	2640	3
	$\Sigma$	<b>4337</b>	<b>97123</b>	<b>2640</b>	<b>3</b>
	1	4	56081	16533	175
	2	2	35940	2181	15
	$\Sigma$	<b>6</b>	<b>92021</b>	<b>18714</b>	<b>190</b>
	1	4	65	231	16
	2	1	92	205	8
	$\Sigma$	<b>5</b>	<b>157</b>	<b>436</b>	<b>24</b>
	1	0	3648	1227	43
	2	0	3223	1135	48
	3	0	3539	1203	47
	4	0	3253	969	51
	$\Sigma$	<b>0</b>	<b>13663</b>	<b>4534</b>	<b>189</b>
	1	0	1	58	2
	2	1	16	2243	235
	$\Sigma$	<b>1</b>	<b>17</b>	<b>2301</b>	<b>237</b>
	1	0	0	22	1
	2	0	0	21	3
	3	0	0	18	1
	4	0	1	32	4
	5	0	3	43	6
	6	0	4	39	7
	7	0	0	50	7
	8	0	0	42	17
	$\Sigma$	<b>0</b>	<b>8</b>	<b>267</b>	<b>46</b>
	1	0	0	53	8
	2	0	0	58	14
	3	0	0	9	1
	4	0	0	1	3
	5	0	0	7	1
	6	0	0	55	17
	7	0	1	43	14
	8	0	0	8	4
$\Sigma$	<b>0</b>	<b>1</b>	<b>234</b>	<b>62</b>	
	1	0	0	0	4
	2	0	0	0	1
	3	0	0	0	1
	4	0	0	0	0
	$\Sigma$	<b>0</b>	<b>0</b>	<b>0</b>	<b>6</b>
	1	0	0	0	0
	2	0	0	0	0
	$\Sigma$	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>

Table 3.3.: Absolute occurrence of different TC-patterns in a high background  $e^+e^- \rightarrow a(a \rightarrow \gamma\gamma)$  MC-sample. The ALP mass is set to 300 MeV. The orientation denotes the different possible configurations by rotation and mirroring. The number of offline clusters denotes the number of basf2 clusters contained by the active TCs.

Pattern	Orientation	Number of all Offline Clusters			
		0	1	2	>2
	0	64	3783	89	0
	$\Sigma$	<b>64</b>	<b>3783</b>	<b>89</b>	<b>0</b>
	1	0	2151	1116	10
	2	0	1808	994	2
	$\Sigma$	<b>0</b>	<b>3959</b>	<b>2110</b>	<b>12</b>
	1	0	2	160	0
	2	0	1	153	1
	$\Sigma$	<b>0</b>	<b>3</b>	<b>313</b>	<b>1</b>
	1	0	123	503	8
	2	0	127	501	8
	3	0	131	496	6
	4	0	138	512	8
	$\Sigma$	<b>0</b>	<b>519</b>	<b>2012</b>	<b>30</b>
	1	0	1	115	0
	2	0	1	249	26
	$\Sigma$	<b>0</b>	<b>2</b>	<b>364</b>	<b>26</b>
	1	0	0	17	1
	2	0	0	16	0
	3	0	0	20	0
	4	0	0	16	0
	5	0	0	24	0
	6	0	1	29	1
	7	0	0	14	0
	8	0	0	23	2
	$\Sigma$	<b>0</b>	<b>1</b>	<b>159</b>	<b>4</b>
	1	0	0	6	1
	2	0	0	9	1
	3	0	0	6	0
	4	0	0	2	0
	5	0	0	5	0
	6	0	0	8	6
	7	0	0	11	2
	8	0	0	7	0
	$\Sigma$	<b>0</b>	<b>0</b>	<b>54</b>	<b>10</b>
	1	0	0	0	4
	2	0	0	0	0
	3	0	0	0	1
	4	0	0	0	0
	$\Sigma$	<b>0</b>	<b>0</b>	<b>0</b>	<b>5</b>
	1	0	0	0	0
	2	0	0	0	0
	$\Sigma$	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>



## 4. Long Shutdown 2

After the start of the Belle II operation in 2019, it became apparent that the targeted design luminosity of SuperKEKB was too ambitious, as the background conditions exceeded estimations. Therefore, the target peak luminosity was reduced to  $6 \cdot 10^{35} \text{ cm}^{-2} \cdot \text{s}^{-1}$  while pursuing a total integrated luminosity of  $50 \text{ fb}^{-1}$ . For this goal, a roadmap was formulated, establishing two upgrade windows Long Shutdown 1 (LS1) and LS2 in which both the accelerator and the detector could be exchanged, upgraded and maintained. The first upgrade was conducted from 2022 to 2023, while the LS2 is scheduled for the time between 2032 and 2033. As the different upgrade proposals are still under study, no final decisions have been made. In the following section, I therefore present an overview of currently discussed upgrades for the Belle II detector and the SuperKEKB accelerator. The summarised information is mainly based on the references [1, 25]. As the proposed upgrades are vast, a selection is made with this thesis and a general interest in mind. Additionally, the target parameters of SuperKEKB can be used to extrapolate the beam background conditions after LS2, described in full detail in [25]. As these play a vital role in the MC simulations used in this thesis, a short overview of these extrapolations is given. Furthermore, an extrapolation of actual L1 trigger rates was done in [2], which is summarised as well.

### 4.1. Planned SuperKEKB Upgrades

To increase the instantaneous luminosity, the SuperKEKB main upgrade plan consists of a redesign of the interaction region (IR). There are two proposed approaches, which both assist in more accurately controlling the beams close to the point of interaction. The first part of the proposal suggests relocating the final focusing magnet closer to the actual interaction point. The second part consists of adding a compensation solenoid coil along the beam pipe, further extending the focusing capabilities towards the interaction point. With these adaptations, the cryostat of the focusing apparatus will interfere with the Vertex Detector (VXD), therefore making a redesign of the inner Belle II detector layers unavoidable. Further SuperKEKB upgrades regard the accelerator system, promising higher possible beam currents and improved beam stability. Additionally, the LS2 will be used for many replacement and maintenance tasks, in order to counteract the ageing of the system.

## 4.2. Estimated Run Conditions after LS2

With the planned SuperKEKB upgrades and the associated changes in beam parameters, the impact on the beam background levels can be extrapolated. Especially for this thesis, these background extrapolations are vital, as they provide background overlay files, which can be used for MC simulations.

### 4.2.1. Beam Background Extrapolation

The extrapolation of the beam background was performed without a final decision on the final detector design; these extrapolations are only rough estimations based on the current detector design and, therefore, with large uncertainties attached. In order to simulate the beam backgrounds, in the first step, the actual beam bunches are simulated, revolving around the accelerator. This simulation includes beam-gas and Touschek scattering and determines the coordinates of particles hitting the beam pipe or collimators in dependence on the ring pressure and beam current. Particles lost in the vicinity of the Belle II detector are then simulated with basf2, in order to yield the induced particle showers and the resulting detector response. For the luminosity background generated at the beam collisions, the beam bunch simulation is not required, and the basf2 simulation is used directly. With current measurements, the individual background components are scaled with measured average data/MC ratios. As the detector design will remain as is, the beam background levels can be simulated for the time before LS2. However, as there are many possible but unknown changes to the detector and IR, the simulation of the beam bunches and subsequent basf2 simulation is not possible for the time after LS2. Therefore, the simulated single-beam backgrounds for the time before LS2 have to be manually scaled to the target beam parameters after LS2. For this, the dependency of the single-beam backgrounds of the beam parameters is fitted heuristically, and an extrapolation to the target beam parameters is performed.

Additionally to these scaling factors, the single-beam backgrounds are scaled with an additional factor to incorporate for unexpected additional increases in beam background. There are 3 different scaling factors used, resulting in the following three Scenarios, as defined in [25]:

1.  $\times 2$  *optimistic* Scenario 1: Only the scaling factors derived from the fits is used. Resulting in roughly a doubling of single-beam backgrounds.
2.  $\times 5$  *intermediate* Scenario 2: An additional factor of 2.5 is applied, to account for unexpected effects.
3.  $\times 10$  *conservative* Scenario 3: Pessimistic scaling, with background estimations of an order larger than before LS2.

The luminosity backgrounds are simulated for the target luminosity after LS2 and only scaled with the measured average data/MC ratios. For the studies conducted in this thesis, the Scenario 2 background overlay files are chosen for the MC-production.

### 4.2.2. L1 Trigger Rate Extrapolation

A similar approach is used to extrapolate the L1 trigger rates. For this, the dependency of each individual trigger bit is assumed to be a linear combination of the single beam

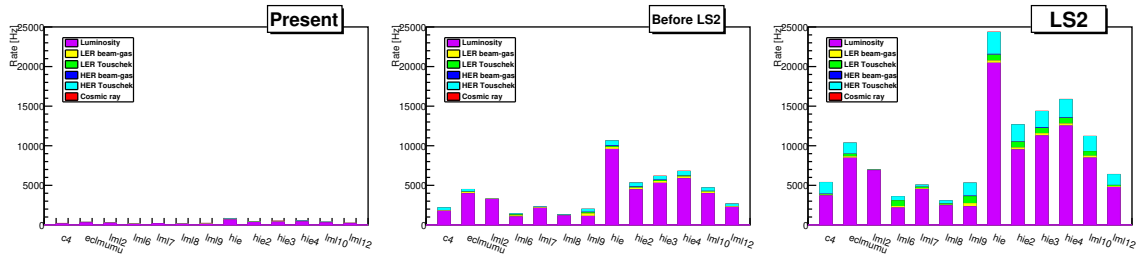


Figure 4.1.: The L1 trigger rate extrapolation for trigger-bits, that only rely on ECL information. An increase in trigger rates, caused by the rising beam background levels, is visible. The most dominating background component is the luminosity-driven background. LER and HER denote the low and high-energy ring components, respectively. This figure is taken from [2].

components and luminosity components. The single beam component is a combination of the trigger rate caused by beam-gas scattering  $TR_{\text{beam-gas}}$ , Touschek scattering  $TR_{\text{Touschek}}$  and a constant rate caused by cosmic rays  $TR_{\text{CR}}$ . The  $TR_{\text{CR}}$  is expected to be constant and is determined by dedicated cosmic runs. The trigger rate induced by Touschek scattering scales like,

$$TR_{\text{Touschek}} \propto \frac{I^2}{n_B \sigma_x \sigma_y \sigma_z}. \quad (4.1)$$

Thereby,  $I$  denotes the beam current,  $\sigma_i$  the bunch sizes in the  $i$ -th coordinate and  $n_B$  the number of bunches in each ring. The trigger rate caused by beam-gas scattering is expected to be proportional to the beam current and the effective residual gas pressure  $P_{\text{eff}}$ . By combining all single-beam components for both beams and assuming a linear behaviour for the luminosity induced trigger rate, a formula for the overall trigger rate in dependence of the beam parameters arises.

In order to fit the individual trigger bits with this formula, the trigger rates have to be determined in dependency of varying beam parameters. For this, a dedicated beam run was conducted, in which the beam parameters were continuously varied, to cover a larger beam parameter space.

The L1 trigger rate extrapolation for the trigger-bits, which only rely on ECL information, is shown in fig. 4.1. The individual background components are colour-coded. A drastic increase in L1 trigger rates is observable, with the "hie" trigger bit reaching almost 25 kHz, which is close to the trigger rate limit of the HLT. It is worth mentioning that these extrapolations also do not incorporate the possible new inner collider and detector design. Additionally, the individual trigger bits are not exclusive and heavily correlated; a simple addition of trigger rates is not possible.

A similar trend of rising trigger rates is also observable in the CDC and KLM. However, there, the proportion between the different background sources differs. As a conclusion of this analysis, the authors suggest studying background reduction methods to cope with the drastically increasing background conditions.

### 4.3. Belle II Upgrade Plans

As mentioned above, a new design of the inner part of Belle II will be required to accommodate the adapted focusing system, if this gets upgraded. Primarily, a new innermost detector will be deployed, replacing both PXD and SVD, called Vertex Detector (VTX). There are three different proposed layout designs, combined with a newly developed sensor and readout. Furthermore, the buffer size of the VTX will be drastically increased, which will extend the latency requirement for the total L1 trigger chain from 5  $\mu\text{s}$  to approximately 15  $\mu\text{s}$ .

For the detector, ageing is also a major problem, especially for the CDC. During the operation, a major tracking efficiency loss was observed, which emerged from the irradiation and the consequent degradation of the CDC. Studies on the extent and reversibility of this degradation are still ongoing. Possible propositions for the LS2 propose a complete drift gas replacement, a removal of the innermost part of the detector, or even building a completely new CDC. If only the inner layers of the CDC are removed, an additional detector system between the VTX and CDC is under debate. This new proposed Inner Tracking and Timing detector (ITT) will be a silicon strip detector, with two layer types. A silicon strip tracker, analogous to the current SVD and a Fast Timing Layer (FTL). Motivating the FTL is the goal of improving low-momentum particle identification through a time-of-flight measurement. Similarly, a time-of-flight measurement of the KLM could boost the  $K_L$  identification. For this, a replacement of the current resistive plate chambers by scintillators with SiPM light readout in the barrel region is studied. Additionally, this would remove the efficiency loss of the KLM with increasing background rates. As currently the KLM is operating in streamer mode, the signal yield is increased; each KLM hit leads to an extended space charge, impeding further particle detection. Alternatively to the replacement with scintillators, this efficiency loss could also be reduced by changing to an avalanche operating mode. Nevertheless, a replacement of the readout electronics would still be necessary. The TOP detector also struggles with degrading hardware, as well as with data acquisition stalls, especially after injections, where the background rate overstrains the detector readout. To handle both issues a upgrade of the used control and readout electronics is required.

### 4.4. Upgrades regarding the ECL L1 trigger pipeline

With increasing background conditions, the problem of pileup noise arises for the ECL. Pileup noise describes the effect that, due to the long decay time in the CsI crystals, the scintillation light curve of one event is still visible in the following event, and overlays the new upcoming signal curve. Due to this pileup noise, the energy and time resolution degrade for higher background rates, especially with regard to the extrapolated background scenarios. In order to counteract this, a new shaper-digitizer (ShaperDSP) is commissioned, featuring a faster shaper and Analog-to-Digital converter. The ShaperDSP serves the purpose to filter, integrate and digitise the signal from the photodiodes and subsequently fit a waveform to the data points. Upgrading the ShaperDSP enables a higher sampling frequency and a more sophisticated fitting of the signal curve, increasing the separation between the signal and the pile-up.

The new ShaperDSP is fast enough that the digitisation for the data acquisition can also be used for the L1 trigger pipeline. In the current system, the digitisation for the data acquisition was slower but more precise than the digitisation for the L1 trigger. As the separate signal shaping for the trigger is therefore obsolete, the analog-sum to create the TCs is also no longer required in this step. Moreover, the digital signal can be transformed to arbitrary TC shapes. With this opportunity, the question arises whether an adaptation of the TC shape can help to improve the accuracy and hence reduce the trigger rate of the ECL L1 trigger. Especially since an adaptation of the TC shape directly entails the adaptation of the whole ECL L1 trigger chain. The alternative way to tackle the increasing trigger rates is to adapt the threshold values applied at the different steps of the pipeline. For example, the 100 MeV threshold applied to each TC could be increased, or the number of clusters required for certain trigger bits could be increased. However, such approaches will always lead to a loss in trigger efficiency, especially for low multiplicity event topologies, harming physics analyses, and should therefore be omitted.

In favour of the change of the TC shape are several theoretical advantages, which are visualised in fig. 4.2.

- As the input information is more highly granulated, the cluster predictions are also expected to have a better position resolution. Especially for single active TC, the position resolution is expected to be much improved.
- Following up on this, also the resolvability of close-by clusters is drastically improved.
- The identification of particles can be improved, especially with regard to muon signatures. This is especially relevant, for example, for the proposed chiral Belle II program.
- Most importantly, the higher granular input enables a shower shape analysis to distinguish between clusters caused by the interaction and clusters caused by beam background, explained in section 2.3

Based on these considerations and the extrapolated trigger rates, the study of the effect of alternative clustering algorithms and smaller TC geometries is required.

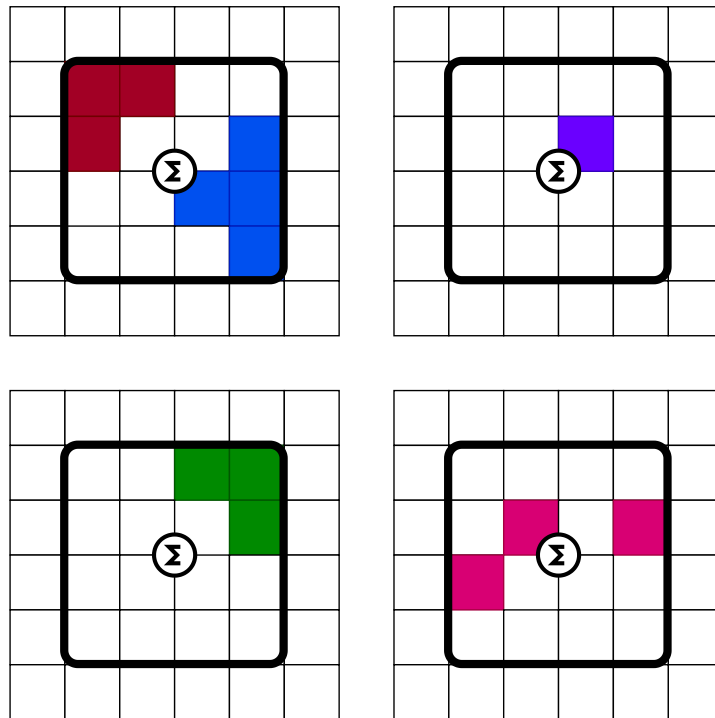


Figure 4.2.: Motivation for a TC shape adaptation. On the upper left, the resolvability of overlapping clusters is visualised. On the upper right, the drastically different signature of muons. On the lower left, the expected improvement regarding position resolution and on the lower right, the possibility to apply a shower shape analysis to distinguish signal from background clusters.

## 5. Network Design

To replace the current clustering algorithm in the ECL L1 trigger, a new algorithm must be devised. The algorithm must reliably find and reconstruct cluster objects. Ideally, this model further incorporates means to reduce the impact of background contributions, to be robust against the rising background levels anticipated for the LS2 upgrade. The surrounding L1 trigger chain imposes strict requirements regarding the latency and size of the algorithm. Currently, the clustering of each event has to be performed within 1  $\mu$ s. The current size restriction is given by the size and capabilities of the currently deployed UT-4 FPGA. For the current L1 trigger chain, an alternative approach to the existing ICN-ETM is already developed and studied in detail in [4, 5]. This model is called the CaloClusterNet and represents a GNN-based approach. It utilises a combination of the GravNet [7] message passing and the ObjectCondensation loss [6]. In the first part of this chapter, the CaloClusterNet approach and architecture are explained.

In order to be integrated into the L1 trigger, a hardware implementation of the model on the current UT-4 board is required. This hardware implementation with the required pre- and post-processing is the so-called GNN-ETM. To comply with the hardware and timing restrictions, the CaloClusterNet must be quantised and transformed into a fixed-point representation. In the second part of this chapter, an overview of quantising a model and the GNN-ETM is given, which is based on [5].

With the proposed LS2 upgrade, the latency constraint will be relaxed. The size constraint of the algorithm will also be altered, as both an upgrade and the usage of multiple FPGAs is possible. Due to these changes, as well as the ShaperDSP upgrade, novel possibilities arise for the clustering algorithm. As clustering on the single crystal information promises to improve the resolutions as well as the background rejection capabilities. Based directly on the architecture of the CaloClusterNet, I deduce an adapted model, which uses the single crystal information instead of the  $4 \times 4$  TCs. To ensure a realistic assessment of the size and complexity of a possible model, a complete implementation of the model on a FPGA is demanded. Therefore, the proposed adaptation of the CaloClusterNet is completely quantised, and a working implementation on an AMD Versal SoCs with AIE is provided in collaboration with the electrical engineering department. In the last part of this chapter, I present my fully quantised adaptation of the CaloClusterNet.

## 5.1. CaloClusterNet

The CaloClusterNet is a GNN-based clustering approach introduced and optimised in [4]. This algorithm is implemented in the GNN-ETM, published in [5], and currently tested in the data taking supplementary to the ICN-ETM. The CaloClusterNet architecture uses the TC information on energy, timing, and position as input. The target is to find and predict all clusters in the event correctly. For this, it predicts the cluster position and energy, as well as an additional signal-background classifier. In order to apply a GNN, a graph-based representation of this input is required. A graph consists of nodes, which are connected by edges. In a GNN, both the nodes and edges contain information. CaloClusterNet relies on a dynamic graph building, so for each event, a new graph is constructed. The model can be separated into two parts. In the first part, it utilises the GravNet [7] message passing to enable an information exchange between the individual inputs. For this, a graph representation of the input nodes is constructed and based on the relation between them, an information exchange is performed. The second part of the model performs the actual clustering and cluster prediction by employing the ObjectCondensation [6] approach. For this, an additional latent space representation is learnt.

### 5.1.1. GravNet

To build the graph representation of the input required for the message passing of GravNet [7], each TC is transformed by multiple dense layers. This provides multiple learnable outputs for each input. A subset of these is interpreted as coordinates in a latent space, representing the position of the graph nodes. The rest of the outputs are learnt features attached to these nodes. The edges of the graph are introduced as connections between each node and its nearest neighbours. The information assigned to each edge is the inverse distance-weighted node feature of the neighbour. Hence, a trainable and dynamic graph representation of the input is constructed. The limitation to only nearest neighbours, as well as the inverse distance weighting, is required to reduce the size and complexity of the graph and is motivated by the locality of the underlying clustering task. The message passing is performed on each node of the graph by aggregating the connected edge information. The used aggregation functions are arbitrary; for the CaloClusterNet, the average and maximum are used.

### 5.1.2. ObjectCondensation

The actual clustering step is performed by the ObjectCondensation [6] and is performed subsequent to a message passing between the individual inputs. In clustering, multiple inputs have to be combined to form a cluster object. In the ObjectCondensation stage, each transformed input is processed by dense layers to predict a latent space position.

The underlying concept of ObjectCondensation assumes that each input can be assigned to an underlying target cluster object or background. For the training of the model, a loss is applied to the predicted latent space positions, which is inspired by physical potentials. All inputs which are assigned to the same underlying target cluster are subject to an attractive potential in the latent space. Inputs assigned to different target clusters or labelled as background are instead subject to a repulsive force. With this approach, all inputs that

have the same target cluster attached are clustered together in the learnt latent space, while simultaneously being separated from other clusters and background contributions.

In this learnt latent space representation, the clustering is performed. Out of each group of clumped-up points, a cluster prediction has to be formed. For this task, the condensation point selection is used, which is also introduced by [6]. Instead of performing a manual aggregation of cluster information in the latent space, this ansatz exploits the fact that the model can be trained to aggregate the information itself. Due to the message passing and assignment to a target cluster, a training target with associated cluster properties exists for each input. By expanding the output of the model by a prediction of these properties, each point in the latent space also predicts the properties of the target cluster. Hence, the clustering task is now simplified from the combination of multiple inputs to the selection of one prediction per cluster.

In principle, this selection could be done randomly, as each point carries a prediction of the underlying cluster. However, this could lead to bad predictions, especially for ambiguous inputs. As an example of this, the case of overlapping clusters in the calorimeter might be considered. In this case, inputs are comprised of energy depositions of multiple clusters. To reduce this ambiguity, the model additionally predicts a  $\beta$ -value. This  $\beta$ -value indicates the certainty of the model, whether the input and the corresponding prediction are a good representative of the underlying cluster object, or not. In the condensation point selection, the nodes are evaluated in descending  $\beta$ -order, to ensure that the most representative nodes are selected first. If a node surpasses a set  $\beta$ -threshold, it is selected as a condensation point, and the attached learnt cluster features are interpreted as a cluster. All other nodes within the vicinity around the condensation point are interpreted as related nodes, representing the same cluster object. Therefore, they are excluded from becoming a condensation point themselves. This algorithm ensures that one representative cluster prediction is performed for each cluster and that background contributions are suppressed due to the  $\beta$ -threshold.

### 5.1.3. Architecture

The combined CaloClusterNet design is shown in fig. 5.1. The inputs are the individual  $4 \times 4$  TCs active in each trigger window. It harnesses the position, energy and timing information of each TC. These inputs are passed to two subsequent GravNet blocks that enable information exchange between the individual inputs. The skip connections around the GravNet block are used to retain the original node information corresponding to each TC and not only the aggregated edge information. In the final dense layer, a position, energy and signal prediction is performed for each input, as well as the position in the condensation point latent space, and the attributed  $\beta$ -value. For stability reasons, instead of predicting the absolute value of the cluster energy, the model predicts a scale factor, which is applied to the input energy to obtain the cluster energy. Only the predictions passing the condensation point selection are treated as cluster predictions.

For the training, additional loss terms apart from the required Object Condensation loss are used for each of the learnable cluster features. For the energy scale factor and position predictions, the respective losses are computed as the mean absolute errors. The goal of the separate binary classifier is to suppress clusters caused by beam background. To ensure a uniform scaling of all loss components, with regard to the Object Condensation loss, they

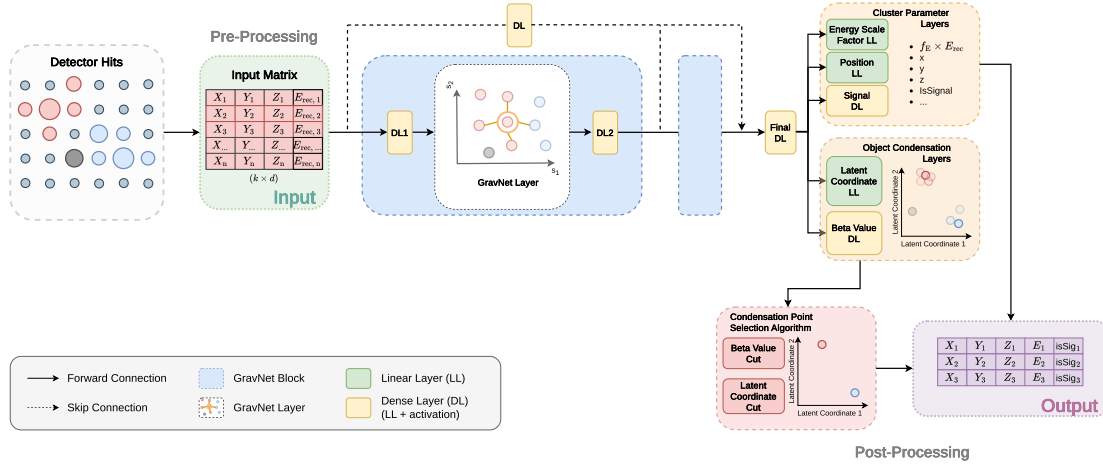


Figure 5.1.: Overview of the CaloClusterNet design. The input is provided in a matrix and passed on to the network. The network is comprised of two subsequent GravNet blocks, which enable information exchange between the individual inputs. The output is predicted per node, and in the inference step, the condensation point selection is performed to gain the final cluster predictions. Figure taken from [5].

are scaled with the magnitude of the learnt potential, which directly corresponds to the  $\beta$ -values [4].

## 5.2. Integration as GNN-ETM

With the architecture of the GNN in place, the model has to be implemented on the current UT-4 board and comply with all the required timing and size constraints. This integration into the L1 trigger is called the GNN-ETM. For the current data-taking period, the GNN-ETM is running in parallel to the ICN-ETM; however, it is still being tested and not yet included in the active trigger decision.

In the following, the adaptations and limitations for the GNN-ETM are introduced. These also hold in the same or slightly adapted ways for the model I propose in the last part of this chapter. The GNN-ETM supports up to 32 active TCs as input, per trigger window. If fewer TCs are active, the input is zero-padded to guarantee a consistent input size. Additionally, the input is pre-scaled to give each input feature approximately the same weight. The most demanding calculations, both utilisation as well as timing-wise, are sorting operations. Therefore, the number of nearest neighbours used in the GravNet layers is limited to eight. The distances calculated in both the GravNet and the ObjectCondensation parts of the model are simplified by employing the Manhattan distance instead of the Euclidean distance. The most important adjustment applied to the model is the quantisation from a float to a fixed point type representation.

### 5.2.1. Quantisation

The encoding of numbers in computers is a well-known problem, and for different tasks, different representations are better suited. The standard floating-point representation

is adaptable and powerful enough to handle a vast range of possible numbers. This representation of numbers is so powerful that it is used in almost every part of modern computing. In the memory of a computer, a fixed number of bits is used to represent a number. For a floating-point number  $n$ , this is usually set to 32 or even 64 bits. In the case of 32 bits, the first bit  $sb$  denotes the sign of the number, followed by the 8 bits encoding the exponent  $exp$  and lastly the 23 bits representing the fraction  $frac$ . Out of this bit sequence, the represented number is calculated as

$$n = (-1)^{sb} \cdot (frac) \cdot 2^{exp}. \quad (5.1)$$

This enables the bit sequence to represent values between  $\pm 1.175e - 38$  to  $\pm 3.4e38$ . The multiplication of two floating-point numbers  $n_1$  and  $n_2$  can be computed as

$$n_1 \cdot n_2 = (-1)^{(sb_1+sb_2)} \cdot (frac_1 \cdot frac_2) \cdot 2^{(exp_1+exp_2)}, \quad (5.2)$$

where the new indices denote the floating-point numbers. However, this result has to be rounded and normalised again to be representable in the same way as  $n_1$  and  $n_2$  are. For the use on FPGAs, this representation is not economical, as the multiplication of the two fractions already has to be implemented as a  $24 \times 24$  bit multiplication. Followed by multiple operations required for the rounding and normalisation, requiring even more time and resources.

To reduce the hardware utilisation and increase the processing speed of multiplications on FPGAs, a fixed point representation is embraced. The fixed point representation also uses a sign bit  $sb$  and is further composed of dedicated integer bits  $int$  and fraction bits  $frac$ . The bit width of all constituents can be chosen arbitrarily, and the data type is then often denoted as  $Qk.l$ .  $Q$  denotes that the sign bit exists;  $UQ$  is used if only positive numbers are represented.  $k$  denotes the number of bits representing the integer part and  $l$  the number of bits encoding the fraction. An arbitrary number  $n$  is represented as,

$$n = (-1)^{sb} \cdot [int + (frac \cdot 2^{-l})]. \quad (5.3)$$

As the multiplication of two fixed-point numbers is realised as a single multiplication on FPGAs, this reduces the required time to only one clock cycle. The total bit-length of both numbers, however, dictates the bit-length of the required multiplication; the required resources directly scale with the chosen fixed-point type. As in multiplications, the precision of the result is larger than the precision of the multiplied numbers; the precision of the result is often reduced to the input type by discarding the least significant bits. The fixed-point representation proves to be more economical if the range of expected values is known beforehand or can be adjusted. However, a number outside of the allowed value range leads to information loss. This can happen if either the value range is exceeded or the precision of the represented number is larger than the precision of the representation. This has to be considered in each calculation step.

For the deployment of the CaloClusterNet as part of the GNN-ETM, the full model is transformed into a 16-bit fixed-point representation. However, all major machine learning

software frameworks are designed to handle and train models in a floating-point representation. To result in a quantised model, either a model in the floating-point representation is quantised after the training, or the quantisation is included already during the training, to enable the model of self correcting the introduced precision loss. For the quantisation-aware training of the CaloClusterNet for the GNN-ETM, as well as for the rest of this work, the Qkeras framework [26] is used. In Qkeras the model is implemented in a floating-point representation. For each weight and bias, a quantisation is introduced, which emulates the precision loss caused by the quantisation. This requires a manual decision on each layer on which data-type to use. As this quantisation is non-differentiable, which would make training the model impossible, Qkeras employs the straight-through estimator. This approximates the gradient for each non-differentiable operation to 1, virtually not affecting the gradient. So internally the model is saved in a floating-point representation, and only for the inference, each floating-point number is quantised to a fixed-point representation.

### 5.3. Adaptation for the LS2 Upgrade

Regarding the LS2 upgrade, a new implementation of the CaloClusterNet is required. As the upgrade of the ShaperDSP enables the use of single crystal information for the L1 trigger, the amount of information accessible for the trigger decision increases. Hence, a more detailed reconstruction can be performed by the L1 trigger. As the CaloClusterNet is not dependent on the geometry of the input, it can handle the single crystal information with no need for adaptations. Though the current implementation in the GNN-ETM can not be used on the single crystal input, as the number of crystals per trigger window exceeds the maximum 32 inputs of the GNN-ETM. For the extrapolated background conditions, the number of inputs without any applied cuts reaches a few hundred per trigger window consistently. The GNN-ETM utilises a majority of the available resources on the current hardware; therefore, a larger FPGA is required to increase the maximum input size. With the size and performance increases of classical FPGAs in recent years, a new FPGA with the required size to implement a scaled-up version of the CaloClusterNet is not realistic. Instead of deploying a regular FPGA, a FPGA with an AMD Versal AI Core Series chip would be better suited, as the AIEs are optimised for matrix multiplications. Similar to most neural networks, major parts of the CaloClusterNet are basic matrix multiplications, which can be implemented much more efficiently on integrated AIE.

To ensure a realistic approach regarding the size and parameters for the CaloClusterNet, a specific and currently available FPGA is selected. This is the AMD VCK190 test board [24]. The board features 400 AIE with a programmable logic, which is slightly smaller compared to the UT-4 board. Hence, this board is not a realistic option to be considered for the LS2 upgrade. But by comparing it to the UT-4 board, it probes the potential of the AIEs. With this specific chipset, the hardware restrictions limit the implementation of the CaloClusterNet to an 8-bit representation, while increasing the maximum input size to 128. Furthermore, the number of nearest neighbours in the GravNet layers is limited to 8. Based on these restrictions, I present a new and optimised CaloClusterNet model, designed to be used on the single crystal input. The implementation of the model on hardware is fully functioning and performed by the electrical engineering department in close collaboration, namely Marc Neu, Fabio Papagno and Till Raedler. An overview of the final model and the

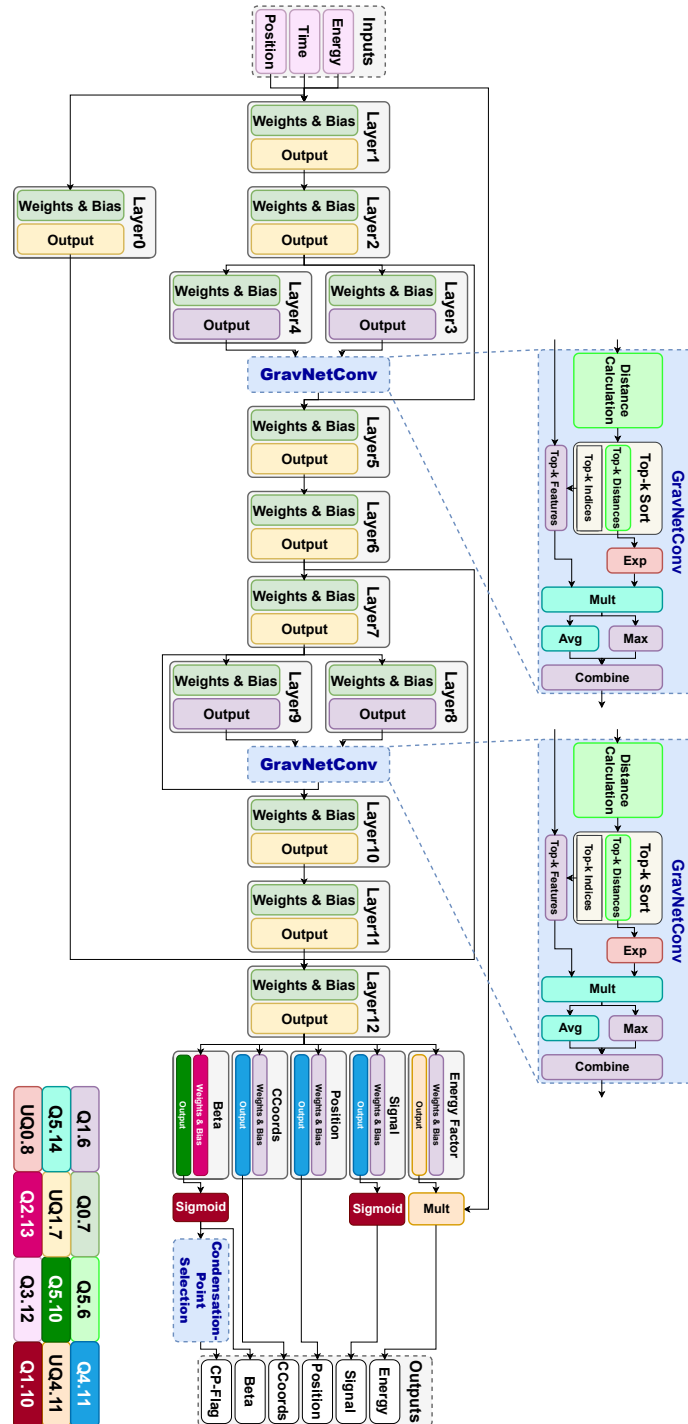


Figure 5.2.: Overview of the fully quantised model. The quantisations are colour-coded.

Table 5.1.: Layer Partitioning and size for the individual dense layers of the model. The Partition denotes whether the layer is implemented on the programmable Logic PL or the AI engines AIE.

Identifier	Partition	Input Dimension	Output Dimension
Layer 0	PL	(5)	(8)
Layer 1	PL	(5)	(8)
Layer 2	AIE	(8)	(16)
Layer 3	AIE	(16)	(6)
Layer 4	AIE	(16)	(8)
Layer 5	AIE	(32)	(32)
Layer 6	AIE	(32)	(16)
Layer 7	AIE	(16)	(16)
Layer 8	AIE	(16)	(6)
Layer 9	AIE	(16)	(8)
Layer 10	AIE	(32)	(32)
Layer 11	AIE	(32)	(16)
Layer 12	AIE	(40)	(16)
Energy Factor	PL	(16)	(1)
Signal	PL	(16)	(1)
Position	PL	(16)	(3)
CCoords	PL	(16)	(3)
Beta	PL	(16)	(1)

chosen quantisations is shown in fig. 5.2. By keeping the quantisation of the model uniform, bottlenecks and unnecessary inflations of the information flow are reduced. With the overall resource and latency contingency, the quantisation of the dense layers is set to an 8-bit fixed-point representation. For the GravNet layers, the quantisations are chosen to be larger to account for the exponential distance scaling, which heavily compresses the available range of parameters. Since the hardware limitations are strict, the width of the deployed dense layers is kept minimal. Especially as the crosstalk between the different inputs is restricted to the GravNet layers, the dense layers are primarily required to transform the respective input features into new representations. Hence, the focus of the architecture lies within the depth of the model. The widths of the dense layers are listed in table 5.1. In order to keep a consistent gradient across the network, multiple skip connections are implemented. The skip connections framing the GravNet blocks are especially relevant as they assure unhindered information propagation past the applied aggregation functions. Since the clustering task is expected to be heavily localised in the real space, the single dense layer "Layer 0" serves as a dedicated way of preserving this input information for the final layers.

Due to the performance gain in matrix multiplications, the dense layers and distance calculations are computed on the AIEs. As the sorting of distances required in the GravNet layers is not well suited to be calculated on AIEs, this part of the algorithm is performed on the programmable logic. Since the in- and output of the board are only accessible via the programmable logic, the first and last layers are also implemented on the programmable logic of the FPGA. For each transition between the programmable logic and the AIEs, a reordering of the data is performed. The limitation of the number of nearest neighbours in the GravNet layers is rooted in the required sorting algorithm, whose hardware requirements and computation time scales quadratically and can not be accelerated by the utilisation of AIEs.

## 6. Datasets and Metrics

As the currently postulated hardware limitations limit the maximum amount of input per event to 128 a input reduction method has to be devised. In order to gain the maximal resolution improvement, this thesis focuses on the single-crystal approach. In this chapter, I propose two different input reduction methods. Since the existing trigger simulation does not feature a single crystal input, to emulate the reduced trigger timing window, a realistic timing cut based on the existing  $4 \times 4$  trigger simulation is chosen. Furthermore, I define the training targets and sample. Followed by the relevant metrics and the samples used to evaluate the clustering and predictive performance of the trained models.

### 6.1. Preprocessing

As the timing windows of the ECL L1 trigger are 250 ns long, this timing selection has to be emulated by the trigger simulation. Since the current ECL L1 trigger is designed for  $4 \times 4$  TCs, the emulated timing cut is only implemented for these TCs, and is described in the following. The detector response for the ECL is simulated in a  $\pm 4 \mu\text{s}$  timing window around the simulated collision, and is binned into 200 ns timing segments. For each segment, the  $4 \times 4$  TCs are determined, and the timing is calculated by an energy-weighted average of the crystal timings. The energy of the TC is determined by an approximated digitisation and simulated fitting of the TC. The simulated TC information is then divided into halfway overlapping 250 ns windows and the window 125 ns before the simulated collision, and the subsequent two windows are selected for further investigation. In order to emulate the offset regarding the collision time, the three windows are shifted by a common random offset between 0 and  $-125$  ns. In the last step, the 250 ns window is selected, in which the largest energy deposition is visible in all TCs. To conserve a similar and comparable timing selection for the single crystal input, the selected timing window of the trigger simulation is used analogously for this work, by applying it to the crystal time stamps.

In order to probe the resulting input size per event, a test sample containing  $B\bar{B}$  events is simulated.  $B\bar{B}$  events are high occupancy events, therefore being ideally suited to evaluate the input reduction methods. Considering the LS2 upgrade, the background files for the extrapolated Scenario 2, as mentioned in section 4.2.1, are used. With this timing cut applied to the simulated crystal information, the number of crystals above 1 MeV still consistently exceeds 800. Hence, a further input reduction method has to be devised and applied to keep the input size within the 128 crystals per event. I present two approaches which both lead to the desired reduction. For both approaches, the resulting number of

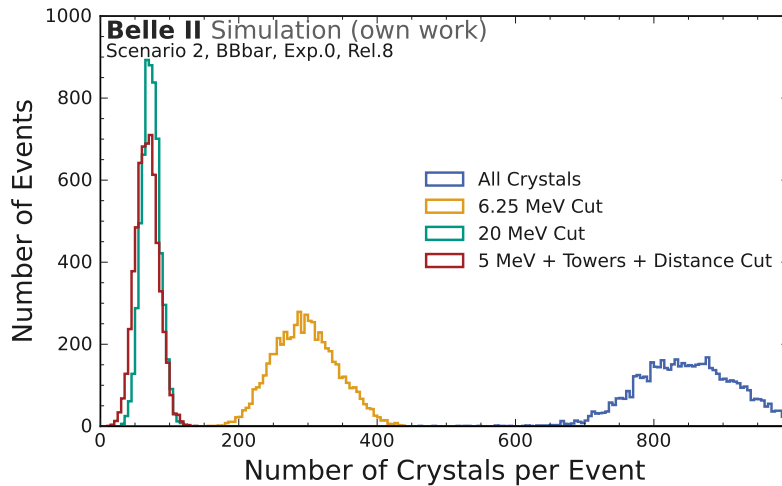


Figure 6.1.: Visualising the number of crystals per event, after applying the timing cut and the two proposed input reduction methods. As a reference, the total amount and the amount of crystals surpassing a naive 6.25 MeV cut are also shown. Both input reduction methods are chosen to comply with the 128-input limit.

crystals per event is shown in fig. 6.1, as well as the comparison to a no input reduction and a naive scaling of the current energy threshold.

### 6.1.1. Flat Energy Cut

The first approach consists of a uniform energy cut applied to all crystals. With the  $4 \times 4$  TC energy cut being set at 100 MeV, the naive scaling of the crystal energy cut would result in  $100 \text{ MeV}/16 = 6.25 \text{ MeV}$ . However, a flat energy cut of 20 MeV is required to stay within the 128 limit.

### 6.1.2. Tower Cut

The second input reduction method is the so-called trigger tower approach. Since clusters are expected to contain high-energy crystals in the centre, these high-energy tower crystals indicate possible cluster positions in the detector. By applying a high-energy cut, these trigger towers are obtained. With this cut, much information regarding the shower shape and also the energy of the underlying cluster is lost. To retain this information, a region of interest is devised around each trigger tower, including the lower energetic crystals at the edges of the cluster. With this approach, sufficiently high energetic clusters are still included, and additionally, more refined information on the crystals in the vicinity is provided. This approach is complexer as more free parameters need to be fixed. The most important parameter is the exact value of the high-energy cut. Depending on this cut, clusters may be completely lost if they do not contain a high-energy crystal. In fig. 6.2, the fraction of retained clusters in dependence on the cluster energy and the tower cut is shown. For this thesis, the tower cut is fixed at 50 MeV, as the fraction of retained 200 MeV clusters is 99% and is additionally marked. Most trigger decisions made by the ECL are driven by high-energy clusters. Especially as for the extrapolated conditions after LS2, this energy

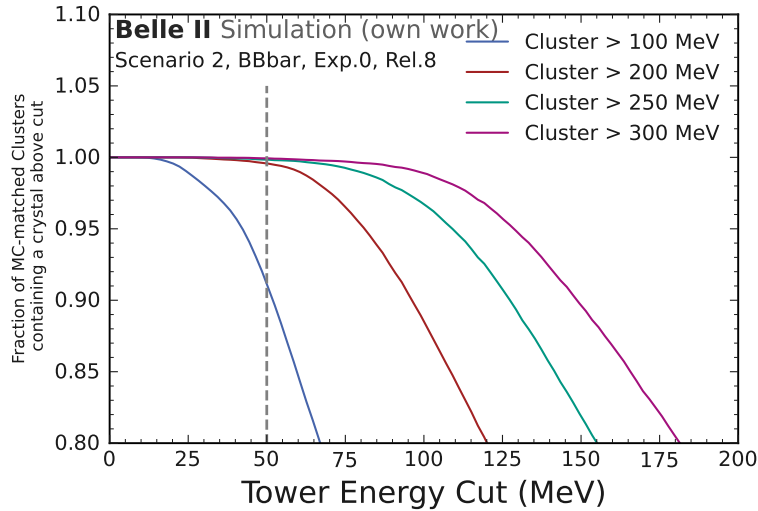


Figure 6.2.: Fraction of MC-matched offline clusters containing a crystal with an energy above the Tower energy cut. Visualised are different Cluster Energies in different colours, as well as the chosen tower cut of 50 MeV. Evaluated on a generic Scenario 2  $B\bar{B}$  MC sample.

region is dominated by background clusters, and the inherent efficiency loss of 100 MeV clusters is justifiable.

The size of the area of interest is also a free parameter, and could additionally differ in  $\phi$  and  $\theta$  direction. For the barrel, the chosen area of interest is 2 crystals wide in each direction. Resulting in a  $5 \times 5$  area of interest for a single tower crystal. For the endcaps for each tower, the two neighbouring crystals within the same  $\theta$ -ring are selected. The central position of the two outermost crystals is then used to define the  $\phi$ -coverage. The crystals in the adjacent two  $\theta$ -rings are selected if their central position is within the  $\phi$ -coverage, leading to a wedge shape. To visualise this, multiple exemplary tower crystals and the resulting regions of interest in the forward endcap are shown in fig. 6.3.

This area of interest can also span across the detector gaps, between the barrel and the endcaps. As the information gain of a crystal in the reconstruction of a cluster is heavily dependent on its energy, an additional 5 MeV energy cut on the crystals in the area of interest is applied. Both the size and low-energy cut are chosen to fulfil the 128 crystals per event requirement.

## 6.2. Training Targets and Training Dataset

With the input features and their prospective preprocessing defined, a universal training sample needs to be created, and the training targets have to be defined.

### 6.2.1. Training Targets

As each input of the model requires a target value for each prediction, each input crystal has to be assigned to a target. In order to remove the dependency of MC truth values and

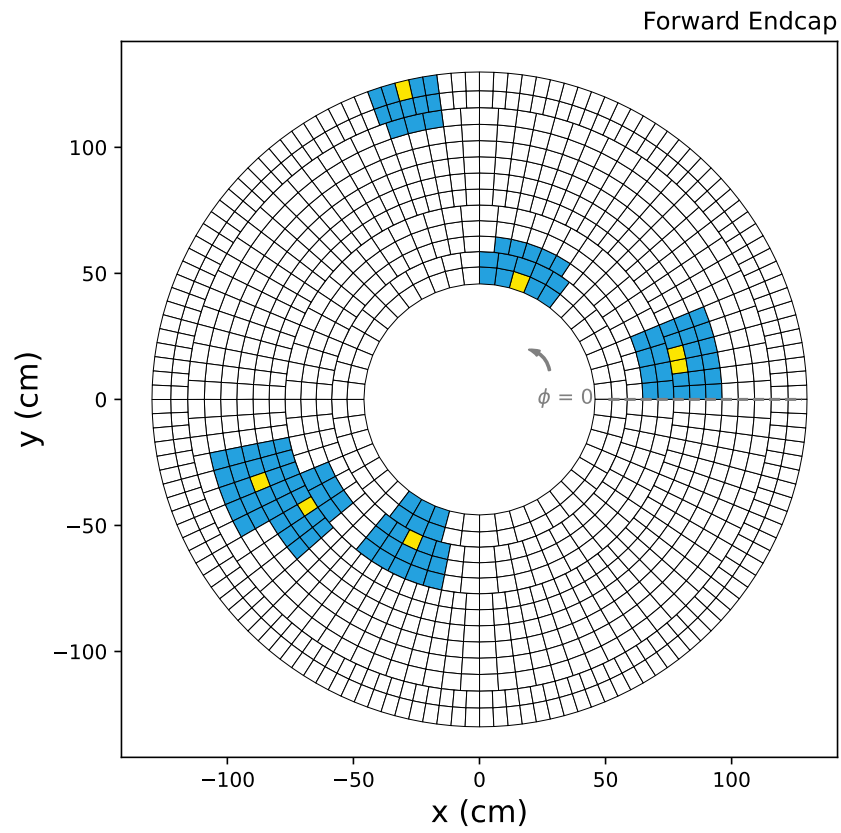


Figure 6.3.: Multiple exemplary tower crystals, marked in yellow, and their corresponding region of interest around them marked in blue, for the forward endcap of the detector. The region of interest of towers close to the detector gap is extended to the adjacent detector part.

the risk of associated mismodeling, the offline ECL clusters reconstructed by basf2 are used (see section 2.3). This theoretically also enables the training of the model on real data. Each crystal can contain energy from multiple offline clusters and an additional background contribution. The assignment of each crystal to a target cluster is performed by selecting the cluster which deposited the most energy in this crystal. If the largest energy deposition in the crystal is caused by background contributions, this crystal is labelled as background, and no target cluster information is provided. The training targets for each crystal are then the reconstructed energy and position of the assigned target cluster.

For the training of the signal background classifier, an additional signal label has to be defined. This label should represent whether a cluster is caused by particles that were produced in the original collision, or a decay product of those, or if the cluster was caused by beam background. In order to evaluate this, the cluster has to be matched with a MC particle. For this, the ratio of energy deposited by a MC particle over the total cluster energy has to exceed 0.2.

### 6.2.2. Training Dataset

In order to avoid an introduction of biases or artefacts within the training sample, a technical MC sample is introduced. As stated in [27], the use of physically constrained training samples can hinder the model from generalising to arbitrary event topologies, as it leads to overfitting to physical correlations. This sample is adopted from [4]. Therefore, no underlying collision is simulated, but the resulting particles are sampled out of simple and smooth distributions. As the focus lies on electromagnetic showers, only photons are simulated as primary particles. The sample consists of 1-6 photons, drawn from a uniform energy distribution between 0 and 7 GeV. While  $\theta$  is sampled from a uniform distribution between  $5^\circ$ - $175^\circ$ ,  $\phi$  is sampled uniformly across the full detector. The  $\theta$ -range is chosen to be larger than the acceptance of the ECL, to also incorporate particle showers created outside of the detector by interactions with passive material. The MC sample is overlaid with low beam background files. Which represent beam background conditions with lower levels of beam background, compared to the current operating conditions. In order to reduce the proportion of beam background clusters in the low energy regime, the sample is enriched with additional low energetic MC clusters. The number of clusters  $N_\gamma$  is approximated with a Poisson distribution, and the energy  $E_\gamma$  distribution is modelled by an exponential. By fitting dedicated background events, the parameters for both distributions are obtained. Resulting in the following distributions:

$$P(N_\gamma) = \frac{3^{N_\gamma}}{N_\gamma!} e^{-3}, \quad P(E_\gamma) = e^{5.046 - 32.621 E_\gamma} \quad (6.1)$$

In addition to that, another focus lies on improving the separation of overlaying clusters. For this, the same distribution of up to 6 high-energy photons is used to produce a further training sample. To this sample, a pair of photons is added, with the same angular and energy distribution as the other photons and an opening angle between  $2.8^\circ$  and  $11.2^\circ$ . The combined training sample is designed to induce as little bias as possible into the model.

### 6.3. Metrics

Assessing the trained models' performance requires a clearly defined set of metrics, describing both the cluster-finding performance as well as the predictive performance regarding the cluster parameters. For the training, the crystals are matched to the offline clusters as described above. The same procedure could also be applied to the evaluation of the model. However, if a crystal leads to a prediction of a close-by cluster and not the cluster it was assigned to by the offline reconstruction, this would artificially worsen the cluster finding metrics. This issue occurs especially in the case of overlapping clusters. Hence, a custom matching between the target offline clusters and the predicted clusters is performed.

The main parameters to match the clusters with the predictions are the predicted position and reconstructed energy. A prediction is matched to a target offline cluster if they are less than 10 cm apart, and if the energy ratio satisfies

$$\frac{E_{\text{predicted}}}{E_{\text{offline}}} \in [0.01, 2]. \quad (6.2)$$

For the ICN-ETM and GNN-ETM, which use the  $4 \times 4$  TC input, the distance requirement is relaxed to 40 cm, as argued in [4]. To confirm the reasonability of the 10 cm distance matching cut, the distance between predictions of a quantised and non-quantised high granularity model and the respective target clusters is shown in fig. 6.4. This is determined on a sample consisting of single-cluster events. For each event, a singular photon is simulated, with no beam background, and only events with one target and one predicted cluster are retained. The energy of the photon is sampled from a uniform distribution between 0 and 7 GeV. More details on the used sample are given in section 6.4.1. In the plot, it is clearly visible that the matching distance is loose enough not to cut into the main distribution and therefore reduce the cluster finding metrics artificially. For the quantised model, 0.6% of predictions across all detector and energy regions are not distance-matched to the corresponding target cluster. For the non-quantised model, this fraction is reduced to 0.08%.

If multiple offline clusters meet the matching requirement for a single prediction, only the closest one is selected. To resolve the ambiguity of multiple predictions being matched to the same offline cluster, the prediction with the energy ratio closest to 1 is selected in these cases. As the simulated offline clusters can have varying timings, even far away from the 250 ns trigger window and therefore invisible to the trigger, only the offline clusters to which at least one crystal is assigned are regarded as target clusters. This assignment is identically as for the training target determination described above.

#### 6.3.1. Cluster-Finding Metrics

The following cluster-finding metrics are evaluated with the described matching applied. Each offline cluster which is matched to a prediction is considered as found, and each prediction matched to an offline cluster is considered as correct. The cluster-finding efficiency is determined as

$$\text{Efficiency} = \frac{N(\text{found target clusters})}{N(\text{all target clusters})}, \quad (6.3)$$

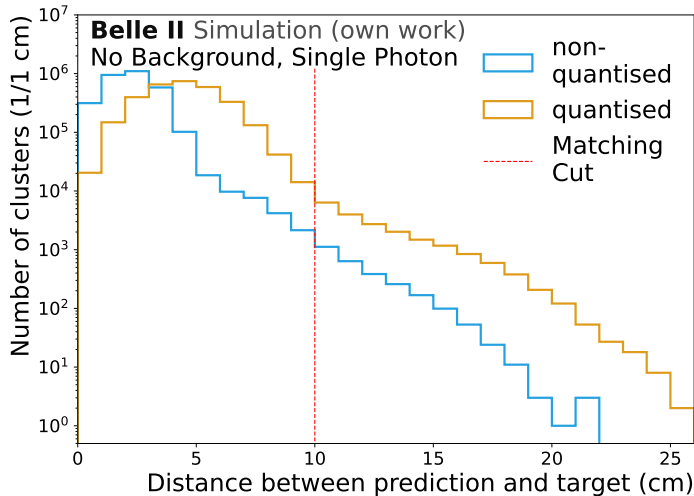


Figure 6.4.: The Euclidean distance between predictions and the target clusters of two high granularity models. The matching distance of 10 cm is shown in red. The sample consists of events containing a single simulated photon, with no beam background, and only events with one target and one predicted cluster are selected. More details on the used sample are given in section 6.4.1.

while the purity is calculated as

$$\text{Purity} = \frac{N(\text{correct predictions})}{N(\text{all predictions})}. \quad (6.4)$$

These metrics are determined for different energy and detector regions. For efficiency, the position and energy of the offline clusters are used, while for the purity, the parameters of the predictions are used to perform the mapping into the detector and energy regions. As the number of entries varies across different samples, energies and detector regions, the statistical uncertainty of each entry varies. Under the assumption that the efficiency and purity of the model are statistically independent for each target, the uncertainty is determined by using a binomial proportion confidence interval [28]. As both metrics are affected by the matching, which is dependent on the predicted energy and position, the corresponding resolution can affect the efficiency and purity and has to be evaluated alongside.

### 6.3.2. Predictive Performance Metrics

The position resolution is defined by the distance between the predicted cluster and the target in each Cartesian coordinate. As a target, either the matched offline cluster or the information of the associated MC particle is used. The MC particle position is only used for photons, as for these the extrapolation to the detector surface is less prone to mistakes. To express the width of the resulting distribution within a scalar value, the difference between the 10<sup>th</sup> and 90<sup>th</sup> percentile is calculated, and in the following referred to

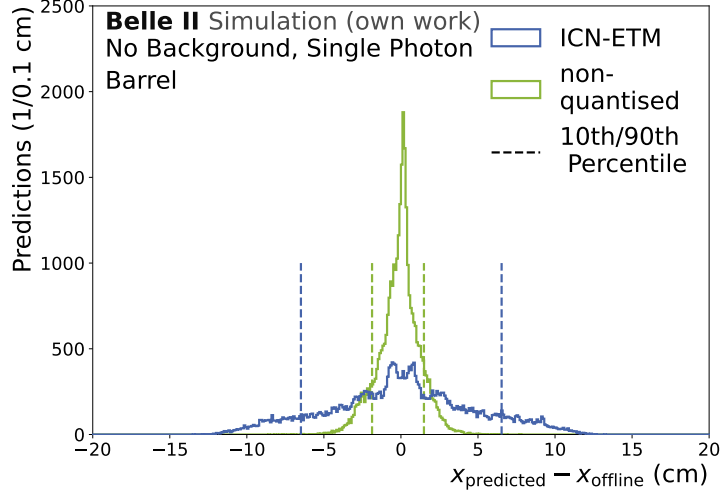


Figure 6.5.: Difference between the predicted cluster x-coordinate and the x-coordinate of the target offline cluster. Shown for the ICN-ETM and an exemplary non-quantised model with single crystal input. Additionally the extracted 10<sup>th</sup> and 90<sup>th</sup> percentile are visualised. The distribution is extracted on a sample with exactly one target cluster per event, and only the barrel region of the detector is shown. The distribution of both models shows steps and the ICN-ETM distribution, even a double peak structure. This directly necessitates the use of percentiles, instead of fits, to extract the position resolution.

as position resolution. This is done as the underlying distributions are heavily affected by the granularity of the input. Hence, they do not follow analytical distributions, especially for the comparison with the current  $4 \times 4$  TCs. A dedicated uncertainty estimation is not performed. Exemplary in fig. 6.5, the underlying distribution of a non-quantised high-granularity model and the currently deployed ICN-ETM, in the barrel region for targets between 100 and 200 MeV, is shown. Additionally the extracted 10<sup>th</sup> and 90<sup>th</sup> percentile are visualised. The distributions are extracted on a sample consisting of single-cluster events. The same as is used for fig. 6.4. More details on the used sample are given section 6.4.1. In both distributions, a step-like structure is visible; additionally, the ICN-ETM features two clearly separated peaks, making a consistent fitting procedure impossible.

The energy resolution is more well-behaved, as the granularity of the input affects the visible energy less. The energy resolution

$$\frac{E_{\text{predicted}} - E_{\text{target}}}{E_{\text{target}}} \quad (6.5)$$

peaks around 1 for a functional model. The predicted energy  $E_{\text{predicted}}$  is compared to the target energy  $E_{\text{target}}$ , which can either be the energy of the matched offline cluster or the associated MC particle energy. The MC particle energy, however, is only used for electrons, positrons and photons. As these particles develop an electromagnetic shower, it is expected that the total particle energy is deposited in the electromagnetic calorimeter. The resulting

distribution is fitted with a double-sided Crystal-ball function, and the full width at half maximum (FWHM) is extracted.

An unbinned fit with zfit [29] is performed. The used density function is implemented in zfit as,

$$f(x; \mu, \sigma, \alpha_L, n_L, \alpha_R, n_R) = \begin{cases} A_L \cdot (B_L - \frac{x-\mu}{\sigma})^{-n_L}, & \text{for } \frac{x-\mu}{\sigma} < -\alpha_L \\ \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right), & \text{for } -\alpha_L \leq \frac{x-\mu}{\sigma} \leq \alpha_R \\ A_R \cdot (B_R + \frac{x-\mu}{\sigma})^{-n_R}, & \text{for } \frac{x-\mu}{\sigma} > \alpha_R \end{cases} \quad (6.6)$$

with

$$A_{L/R} = \left( \frac{n_{L/R}}{|\alpha_{L/R}|} \right)^{n_{L/R}} \exp\left(-\frac{|\alpha_{L/R}|^2}{2}\right), \quad (6.7)$$

$$B_{L/R} = \frac{n_{L/R}}{|\alpha_{L/R}|} - |\alpha_{L/R}|.$$

$\mu, \sigma, \alpha_L, n_L, \alpha_R, n_R$  denote free fit parameters and  $x$  represents equation (6.5). Even though the energy resolution is more well-behaved, there are still artefacts remaining in the resulting distributions. Most dominantly, a second peak can occur towards -1 if the model separates a cluster into multiple lower-energy cluster predictions. As the interested metric is given by the properties of the main peak around 1, the fitting range is adapted accordingly. Regardless of the selected target, the resulting distribution is expected to exhibit a bias, caused by energy leakage. To correct for this, a multiplicative correction is applied. The correction factor is computed as a function of the fitted mean  $\mu$ , resulting in a bias-corrected distribution, given by

$$\frac{(E_{\text{predicted}} \cdot [1 - \mu]) - E_{\text{target}}}{E_{\text{target}}}. \quad (6.8)$$

As an example a uncorrected distribution with the initial and the corrected distribution with the final fit and the resulting FWHM is shown in fig. 6.6. The uncertainties of the fit are propagated onto the FWHM. As the underlying distributions are affected by the matching criteria as well as the granularity and inherent flaws of the models, for example, double predictions of a singular target, the resulting distributions do not always follow a double-sided Crystal-ball function. Therefore, some fits exhibit much larger uncertainties.

### 6.3.3. Signal Background Classifier

For each cluster, the signal classifier is predicted, in addition to the physical properties of the cluster. Each cluster surpassing a classifier threshold is designated as a signal prediction, while all other clusters are interpreted as background clusters. By tuning this threshold value, the ratio of falsely predicted clusters changes. In order to evaluate the performance of the classifier, the signal retention rate is defined as

$$R_S = \frac{N(\text{matched signal predictions})}{N(\text{signal offline clusters})}. \quad (6.9)$$

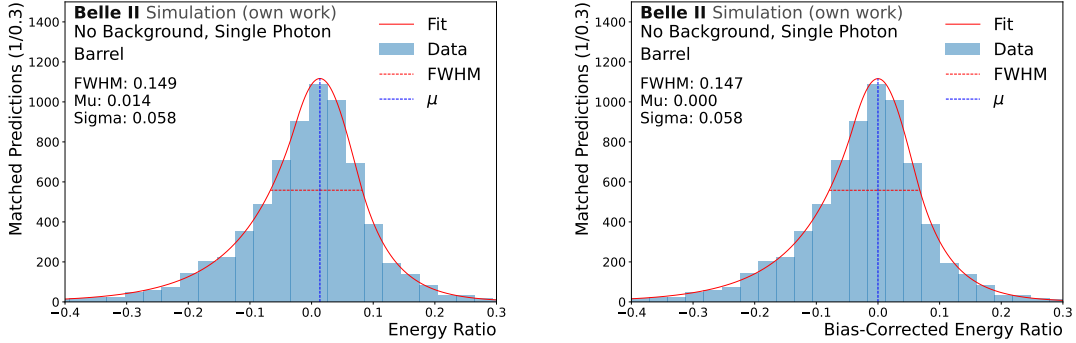


Figure 6.6.: Example of the bias correction and FWHM extraction. On the left, the uncorrected energy resolution is shown, with a double-sided Crystal-ball fit to extract the centre of the distribution. With this value, the energy resolution is now multiplicatively corrected, resulting in the distribution on the right-hand side. On the corrected distribution, a second fit is performed to determine the bias-corrected FWHM.

This corresponds to the true positive rate of the classifier. Additionally, the background rejection rate is defined as

$$R_B = \frac{N(\text{matched background predictions})}{N(\text{background offline clusters})}, \quad (6.10)$$

corresponding to the true negative rate. The binary label of the predictions is dependent on the threshold value, while the label of the target offline clusters is fixed. By varying the threshold value and comparing both metrics, the Receiver Operating Characteristic (ROC) is obtained. Even though for the deployment of the model a distinct threshold value has to be chosen, the behaviour of the can also ROC give insight into the distinctive power of the classifier, as well as its robustness.

## 6.4. Evaluation Datasets

Since the trigger algorithm will face all possible events, the evaluation datasets need to span a large range of topologies. The samples are chosen to cover the whole detector region and energy range, as well as being relevant to the L1 trigger trigger task.

### 6.4.1. Single Photon Sample

The first sample is also a technical photon sample. It is composed of only a single photon with no simulated background. The photon is sampled from a uniform energy and angular distribution. With energies ranging from 0 to 7 GeV.  $\theta$  ranges from  $5^\circ$  to  $175^\circ$ , while  $\phi$  covers the total angular range. This sample defines the simplest case of event topology. It is used to probe the clustering with the smallest influence of the underlying physics as possible. Even though only one particle is simulated, multiple clusters can occur as the photon undergoes pair conversion and produces an electron-positron pair. The resulting

clusters range from overlapping to clearly separable. To mitigate this effect, only events are considered in which exactly one offline cluster is present.

The distribution of the simulated photon energies, the resulting cluster energies of the resulting targets, as well as their  $\theta$  and  $\phi$  distribution are shown in fig. 6.7. In all plots, the two different input reduction methods as well as the current  $4 \times 4$  TCs are shown. A small deviation between the different granularities is visible in the lower cluster energy region and in the backward endcap of the detector. This deviation is expected, as the smaller granular input can retain lower energetic contributions, compared to the 100 MeV energy cut of the  $4 \times 4$  TCs. This deviation will affect the higher granular models, as the to predicted clusters have less visible energy compared to clusters retained by the  $4 \times 4$  TCs.

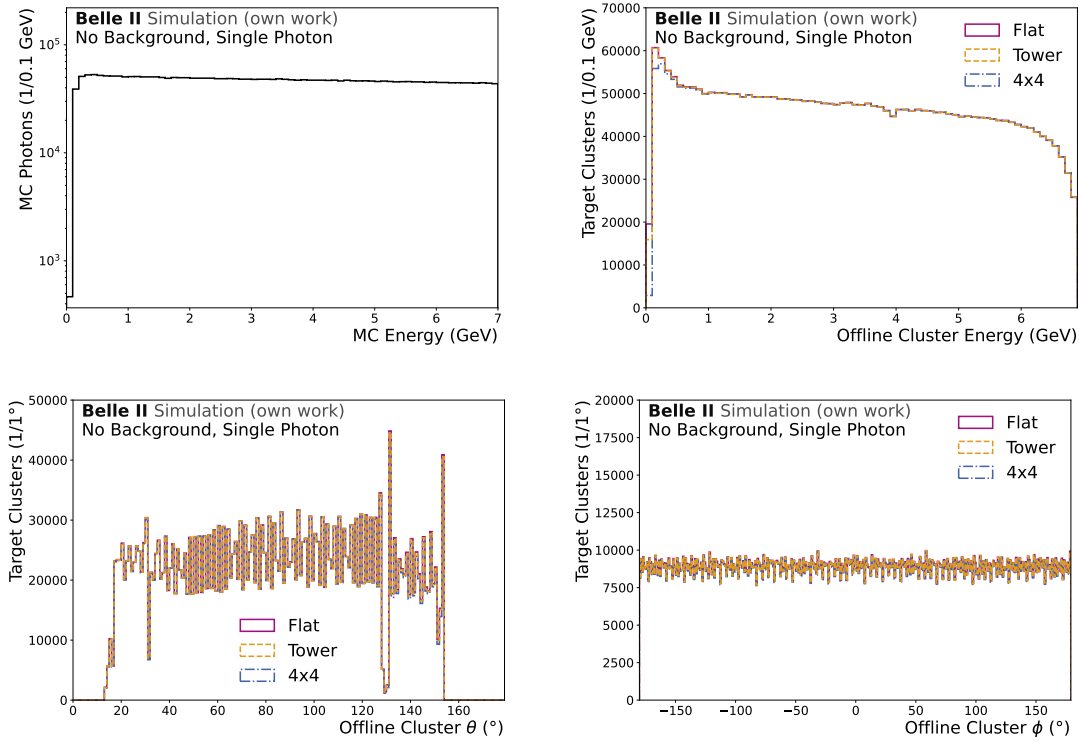


Figure 6.7.: Simulated photon energies and target parameter distributions, for the different preprocessing methods, evaluated on the single photon sample without background. A minute deviation between the distributions is visible in the lower cluster energy region and in the backward endcap of the detector.

### 6.4.2. Bhabha Sample

The most common event signature, the L1 trigger encounters, is caused by Bhabha scattering. Therefore, the performance of the trained models on Bhabha events is essential, and a MC sample is used, with the Scenario 2 overlay introduced in section 4.2.1. For this sample, only the comparison between the different high-granularity models is conducted. The focus of this comparison lies within the performance on the high-energy clusters caused by the

primary electron and positron. To obtain a clean sample of  $e^+e^- \rightarrow e^+e^-$ , only events containing two tracks with a large opening angle and high energetic matched ECL clusters are selected. The following selections are applied to the charged particle candidates:

- The transversal momentum  $p_t$  of each track has to supersede  $p_t > 0.2 \text{ GeV} \cdot c^{-1}$ .
- The centre of mass (CM) energy  $E_{\text{c.m.}}$  of each particle has to be  $2.5 < E_{\text{c.m.}} < 5.82 \text{ GeV}$ .
- The to the tracked matched cluster energy  $E_{\text{cluster}}$  has to be above 1 GeV.

For each event, the two candidates with the highest CM momentum are selected. If they are oppositely charged and deviate from a back-to-back configuration in the CM system by a maximum of  $5^\circ$ , both in  $\theta$  and  $\phi$ , the event is retained.

The resulting distributions of the targets are shown in fig. 6.8. Both in the MC energy as well as the offline cluster energy histograms, the expected double peak is visible at 4 and 7 GeV. In the cluster energies, the contribution of background clusters is clearly visible towards lower energies. Furthermore, a large discrepancy between the flat and tower preprocessing is visible, which is mainly located towards the backwards direction of the detector. This large discrepancy is mainly caused by background clusters. For low-energy background clusters, often only a single crystal surpasses the energy threshold of 20 MeV, while no crystal overcomes the tower cut. Therefore, the flat preprocessing contains a much larger fraction of low-energy background clusters. At higher energies, a similar effect can occur if background clusters are at the edges of the emulated timing window. This leads to individual crystals within the timing window of the L1 trigger, while the rest of the cluster is invisible to the trigger.

If only signal cluster targets are considered, as shown in fig. 6.9, this effect is not visible. To further probe the extent of the partial visibility of cluster energy, the visibility is defined as

$$\text{Visibility} = \frac{\sum E_{\text{crystals}}}{E_{\text{cluster}}}. \quad (6.11)$$

$\sum E_{\text{crystals}}$  denotes the energy sum of all crystals, which have the same offline cluster as the target assigned, while  $E_{\text{cluster}}$  denotes the energy of said cluster. In the visibility distribution, a clear double peak is visible for the flat preprocessing, where only half of the cluster energy is visible. This is caused by both the effect at the edges of the timing window, as well as the high 20 MeV energy cut, and visualises the drawbacks of this input reduction method.

### 6.4.3. Dimuon Sample

Dimuon events have a similar event topology as the Bhabha events, with two charged particles being back-to-back emitted in the centre-of-mass frame. Since the muons are minimising ionising particles, they deposited much less energy in the calorimeter, hence a lower cluster energy is required for the candidates, making the selection similar but orthogonal to the selection used for the Bhabha events. The following selections are applied to the charged particle candidates:

- The transversal momentum  $p_t$  of each track has to supersede  $p_t > 0.2 \text{ GeV} \cdot c^{-1}$ .

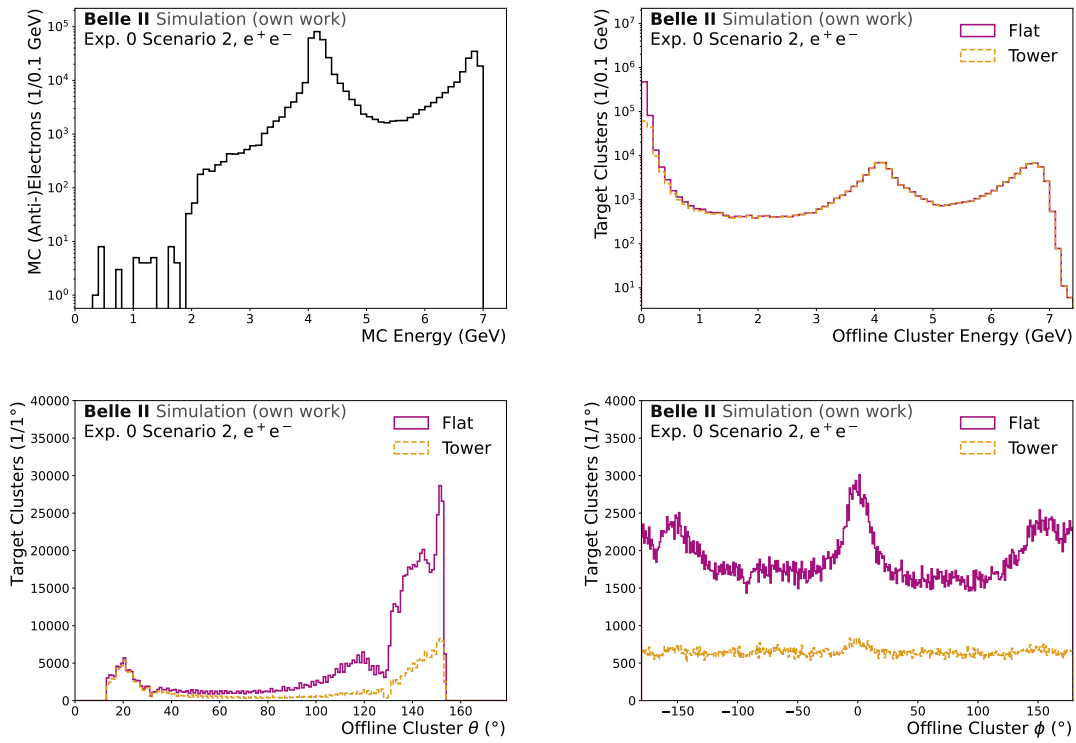


Figure 6.8.: Simulated (anti-)electron energies and target parameter distributions, for the different preprocessing methods, evaluated on the Bhabha sample with Scenario 2 background. A large deviation between the distributions is visible in the lower cluster energy region and in the backward endcap of the detector, as well as in the  $\phi$  distribution. These deviations are mainly caused by background clusters.

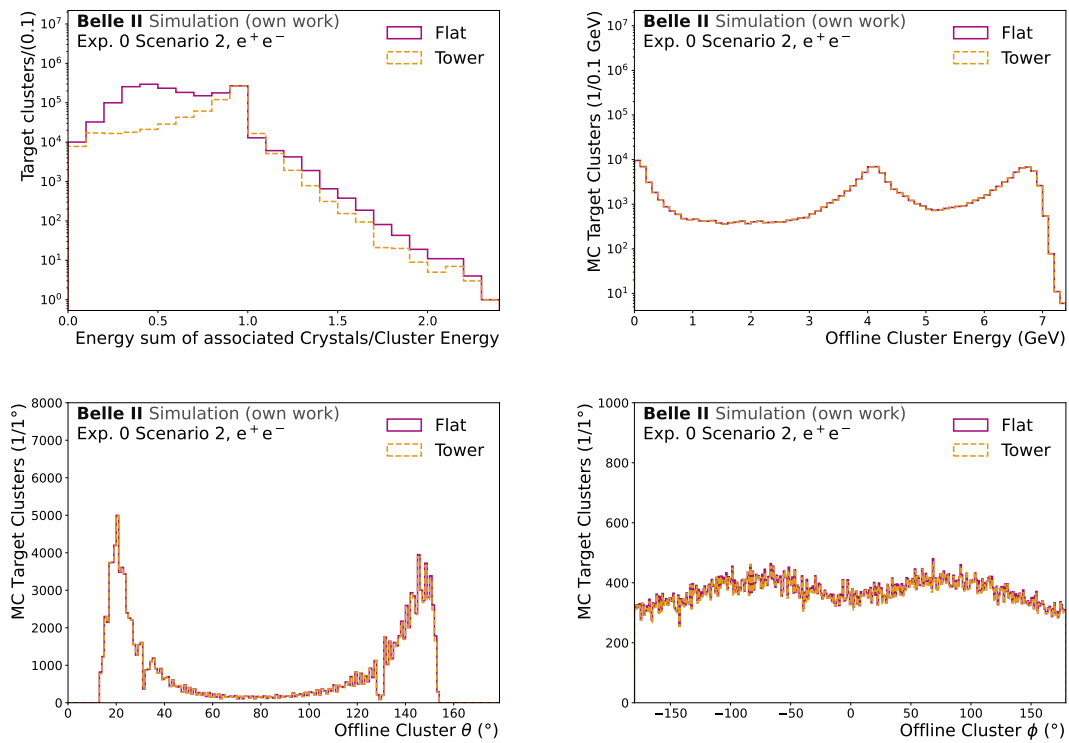


Figure 6.9.: Visibility of the target clusters as well as the distribution of MC-matched clusters, evaluated on the Bhabha sample with Scenario 2 background. Compared to all targets, these distributions exhibit the underlying physics features more clearly.

- The CM energy  $E_{\text{c.m.}}$  of each particle has to be  $2.5 < E_{\text{c.m.}} < 5.82$  GeV.
- The to the tracked matched cluster energy  $E_{\text{cluster}}$  has to be below 0.5 GeV.
- The momentum in the CM system has to exceed  $1.5 \text{ GeV} \cdot \text{c}^{-1}$ .
- The track has to approach the origin in the  $r - \phi$  plane with a maximum distance of 2 cm and in the  $z$ -direction with a maximum distance of 4 cm.

For each event, the two candidates with the highest CM momentum are selected. The candidates are required to be oppositely charged, and the combined dimuon system has to have a reconstructed energy of at least 9 GeV. Furthermore, the candidates can only deviate from a back-to-back configuration in the CM system by a maximum of  $5^\circ$ , both in  $\theta$  and  $\phi$ . The resulting distributions visualising the MC energies, the offline cluster energies and positions are shown in fig. 6.10. Also in this sample, the Scenario 2 background overlay was used. Similar to the Bhabha sample, the double peak structure at 4 GeV and 7 GeV is visible in the MC energies of the (anti-) muons. The offline cluster energies exhibit a strong decrease with rising energies. Compared to the tower preprocessing, the flat preprocessing is dominated by background contributions at lower energies. Concealing the visible peak in the tower sample at 200 MeV. As for the Bhabha sample, this background is concentrated towards the backward part of the detector.

In comparison, the distribution of signal clusters is even across the different parts of the detector. Furthermore, the expected peak of the (anti-)muons is clearly visible in the offline cluster energy distribution.

#### 6.4.4. $B\bar{B}$ Sample

Regarding the filtering process performed by the L1 trigger trigger and the physics motivation of Belle II, the performance on  $B\bar{B}$  events is essential. As the decay of these particles entails many clusters, it is also a good probe for high-occupancy events. Furthermore, it also includes more diversified cluster shapes, as hadronic and muonic signatures are inherent. The MC sample uses the Scenario 2 extrapolation overlay files, introduced in section 4.2.1.

The resolution of the models is only determined on the MC-matched photons, which show a declining distribution in fig. 6.12. The energy distribution of the target clusters features a similar distribution, with deviations between the preprocessing samples at low energies. These targets are located mainly in the backward endcap and are caused primarily by background clusters, as the distributions in fig. 6.13 show. In the visibility plot, the effect of the raised background in the  $4 \times 4$  TCs is clearly visible.

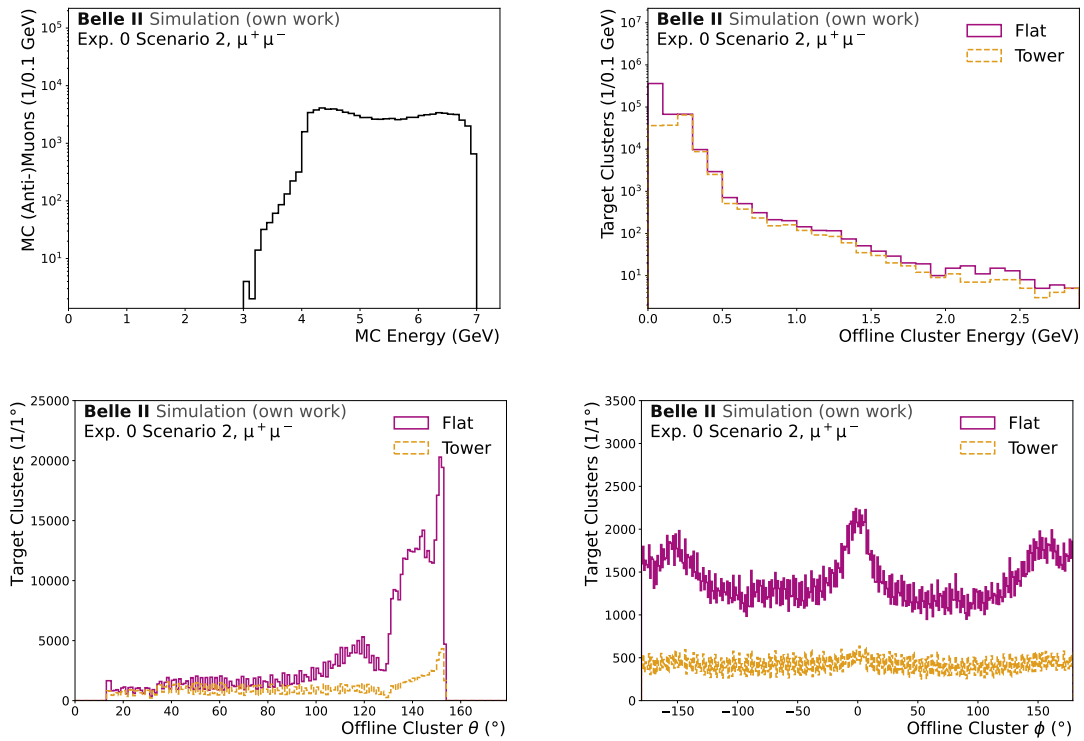


Figure 6.10.: Simulated (anti-)muon energies and target parameter distributions, for the different preprocessing methods, evaluated on the muon pair sample with Scenario 2 background. A large deviation between the distributions is visible in the lower cluster energy region and in the backward endcap of the detector, as well as in the  $\phi$  distribution. These deviations are mainly caused by background clusters.

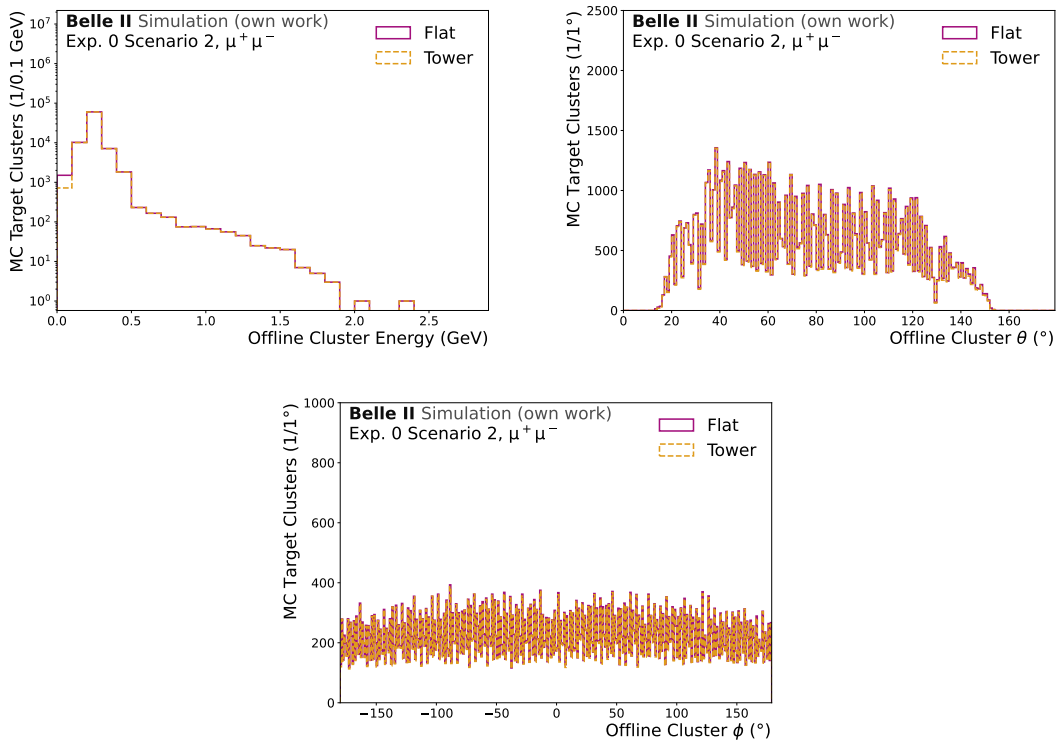


Figure 6.11.: Visibility of the target clusters as well as the distribution of MC-matched clusters, evaluated on the Bhabha sample with Scenario 2 background. Compared to all targets, these distributions exhibit the underlying physics features more clearly.

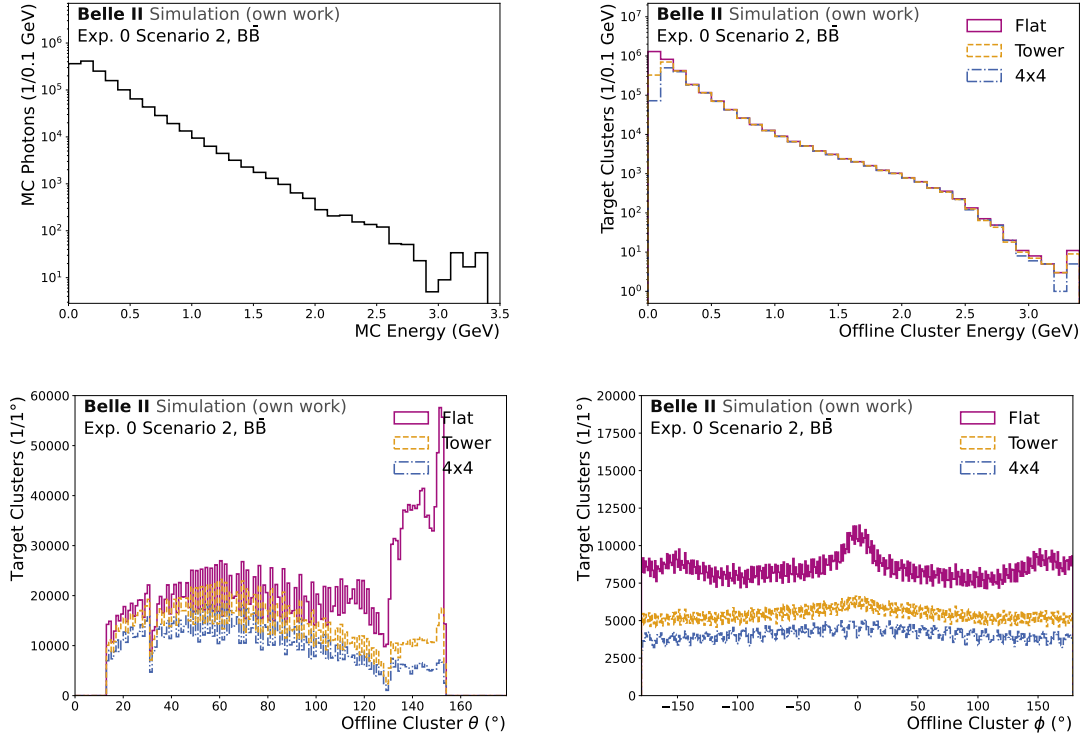


Figure 6.12.: Simulated photon energies and target parameter distributions, for the different preprocessing methods, evaluated on the  $B\bar{B}$  sample with Scenario 2 background. A large deviation between the distributions is visible in the lower cluster energy region and in the backward endcap of the detector, as well as in the  $\phi$  distribution. These deviations are mainly caused by background clusters.

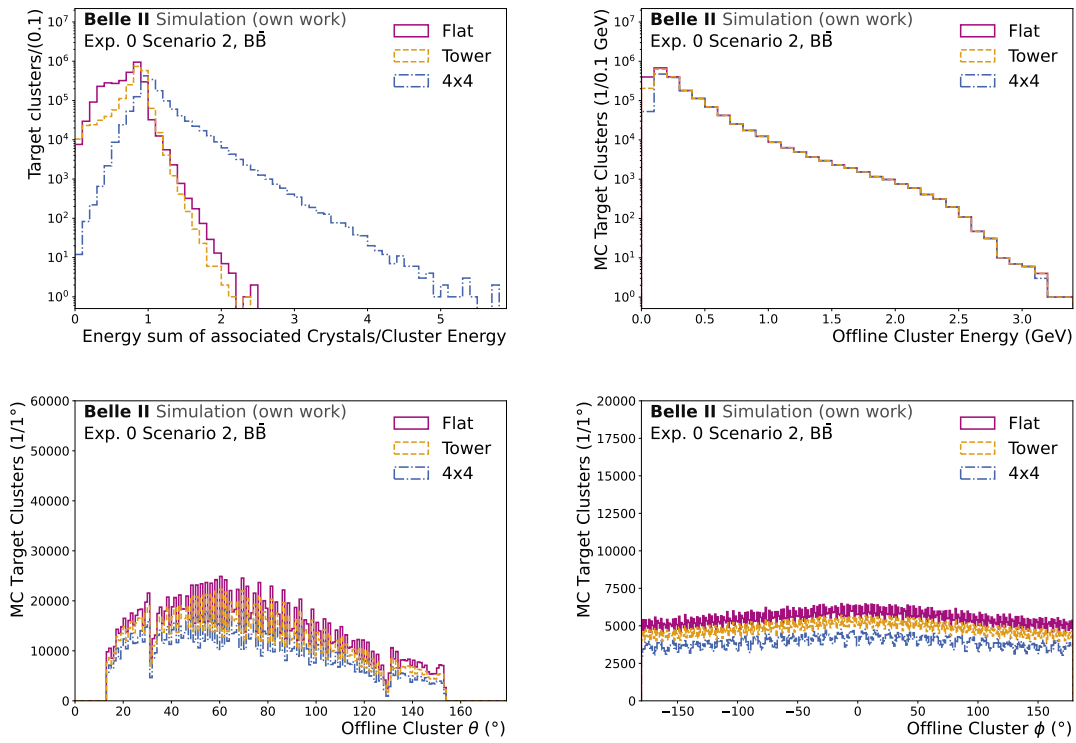


Figure 6.13.: Visibility of the target clusters as well as the distribution of MC-matched clusters, evaluated on the  $B\bar{B}$  sample with Scenario 2 background. Compared to all targets, these distributions exhibit a more even distribution.



## 7. Evaluation

With the proposed upgrade of the ECL L1 trigger pipeline and adaptation of the clustering algorithm, two questions arise. First of all, whether the higher-granularity input can increase the model’s performance in comparison to the current ICN-ETM and GNN-ETM. Furthermore, how well does the model perform compared to the offline reconstruction, and the extent to which this performance increase is dependent on the chosen input reduction method. The second question is focused more towards the feasibility of an actual implementation and what limitations such an implementation might entail. For this, I present the compared models and evaluate them in the following on the evaluation samples defined in section 6.4.

### 7.1. Model Training

To probe the impact of the input granularity, a non-quantised model is trained, but still features the overall model architecture presented in section 5.3. This non-quantised model is trained once on the flat energy cut and once on the tower cut, defined in section 6.1, to probe the impact of the input reduction. To further enhance the information aggregation and extract the best possible performance out of these models, the number of nearest neighbours in the GravNet layers is increased to 32. These models serve as baseline models for the high-granularity input and are compared to the low-granularity ICN-ETM and GNN-ETM. In the following, these models are referenced as **Flat 20** and **Tower 5x5** models.

As the Flat 20 and Tower 5x5 model can not be implemented on currently obtainable FPGA hardware, a third model is introduced. This **Quantised** model features the exact same specifications as described in section 5.3. Including the limit of 8 nearest neighbours in the GravNet layers. This model is fully implemented on the AMD VCK190 board and probes the performance on a currently available FPGA. As an input reduction method for the Quantised model, the flat energy cut is selected. This preprocessing is selected, as there is no further adaptation of the L1 trigger pipeline required. The energy cut could be applied to each crystal, directly at the ShaperDSP. For the tower cut, neighbouring crystals, even across different detector regions, influence each other. Hence, to apply the tower cut, an interconnection between different readout boards and detector regions has to be realised, or it has to be implemented on the FPGA running the clustering, effectively reducing its available resources.

All models are trained for 600 epochs on the same training sample, described in section 6.2.2, consisting of 80000 events. As a validation sample, a similarly designed sample as the training sample is used. For the training of the Quantised model, qnoise is introduced. For an arbitrary float value  $x$ , the specified target quantisation results in  $x_{\text{quantised}}$ . Instead of expecting the model to converge on the heavily quantised model weights, the training is started without any quantisation applied. Once the model is already trained to a certain degree, the qnoise is added. This qnoise scales the degree of quantisation as a function of the current training epoch  $n$ , until the target quantisation is reached. For the training of the Quantised model, the chosen function is given by

$$q(n) \begin{cases} 0 & \text{if } n \leq 300 \\ 1 - \left(\frac{500-n}{200}\right)^3 & 300 < n \leq 500 . \\ 1 & \text{if } n > 500 \end{cases} \quad (7.1)$$

With this qnoise-factor  $q$  the current quantised value of  $x$  is determined by

$$x_{\text{qnoise}} = x + x_{\text{quantised}}. \quad (7.2)$$

As this qnoise is introduced gradually, while the model is already pretrained, the model is already close to a local minimum and can adapt its values to compensate for the qnoise, and therefore the quantisation. Compared to a fully quantised training, the applied qnoise results in a performance and stability increase of the model.

Alternative scalings of the individual loss components, compared to the GNN-ETM, are studied, with no performance increase observed. Instead, this leads to increased training instabilities. For fine-tuning the models, both a learning rate scheduling and an early stopping procedure are used. Additionally, a pruning of 40% is applied to reduce the model size [30]. In pruning, the smallest absolute weights of each layer are masked with zero, as they are the least relevant to the model prediction. As this is already introduced in the training, the model can adapt to the loss of these weights, and the size of the model can be reduced without a loss in performance. The resulting validation losses of the three models are shown in fig. 7.1. To reduce the impact of fluctuations between trainings, each model setup is trained 5 times, and the overall best-performing models are selected for the analysis in this thesis.

### 7.1.1. Cluster Cut Selection

Once the models are fully trained, the hyperparameters of the Condensation point have to be determined. The hyperparameters are the  $\beta$ -threshold value and the size of the latent space around each condensation point, in which all other predictions are removed from becoming a condensation point themselves. This so-called cluster cut is determined on an overlapping photon MC-sample. This sample is simulated specifically for this reason and is sampled from the same angular and energy distribution as the dedicated overlapping clusters used in the training sample. Only events with exactly two target clusters are selected. For this purpose, no background is overlaid to enable the following procedure, additionally visualised in fig. 7.2.

First, the size of the two clusters and their distance are determined in the cluster space. The centre of the cluster is determined by the position of the node with the highest  $\beta$ -value

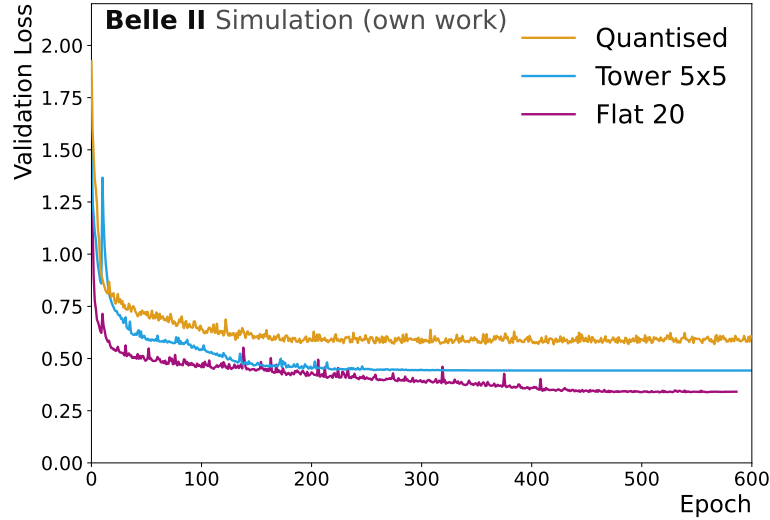


Figure 7.1.: Validation loss curves of the three models evaluated in this thesis. A converging behaviour is visible with only small noise contributions. For the Flat 20 model, the early stopping is visible. Across the multiple training iterations, the loss of the Quantised model is overall larger than that of comparable non-quantised models.

prediction. So the node, which leads to the condensation point representing the cluster. To ensure that every output node associated with the same target cluster is included within the cluster cut, the size of the cluster is determined by the largest distance of the centre to one of those output nodes. The distance between the two clusters in the event is determined by the distance of the cluster centres. To optimise the cluster cut, the fraction of resolved cluster pairs is maximised by varying the cluster cut. These resolved cluster pairs are defined as both clusters having a size smaller than the cluster cut and a respective distance larger than the cluster cut. If the cluster cut is too small, the individual clusters are not contained, leading to multiple predictions; if the cluster cut is too large, the clusters are not separated properly, leading to a missing prediction for one of the clusters. This definition ensures that for a resolved cluster pair, both clusters lead to exactly one prediction, caused by their respective highest  $\beta$ -value. With the cluster cut fixed, the  $\beta$ -cut is determined on the validation sample by simultaneously optimising for an optimal efficiency and purity for high energetic clusters. The validation sample is the sample used for the validation loss, which is simulated following the same underlying distributions as the training sample. The optimisation of the  $\beta$ -cut is heavily dependent on the sample it is optimised for, and leads to a trade-off between purity and efficiency. By tightening the  $\beta$ -cut, fewer clusters are predicted, leading to a lower efficiency, but also fewer multiple predictions or predictions on background crystals occur, leading to an increased purity.

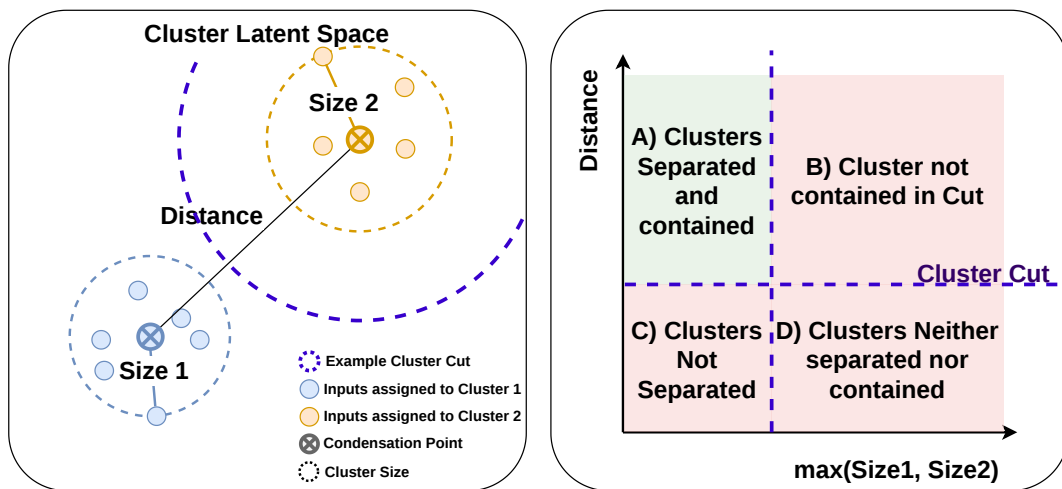


Figure 7.2.: Visualisation of the automatised cluster cut selection on an overlapping two-cluster sample with no overlaid background. For each event, the size of the two clusters and their distance are determined in the clustering space. The cluster cut is optimised by maximising the fraction of events in which the two clusters are correctly separated, visualised in quadrant A) on the right side. This is given if the cluster cut is smaller than the distance between the clusters and larger than both of the clusters. If the cluster cut is too small, the individual clusters are not contained, as seen in quadrants B) and D), leading to multiple predictions; if the cluster cut is too large, the clusters are not separated properly, leading to a missing prediction for one of the clusters, represented by the quadrants C) and D).

## 7.2. Single Photon Sample

The first sample to evaluate the models on is the single photon sample. It aims to probe the clustering performance and predictive performance in a clean environment. This helps to identify and understand the underlying effects that are also present in subsequent samples. For all metrics and models evaluated on this sample, only events that contain exactly one target offline cluster are evaluated.

### 7.2.1. Granularity Comparison

In fig. 7.3, the cluster-finding efficiency and purity are shown in dependency of the detector region and energy. For each model, both metrics are determined on all respective target clusters. Both metrics show an overall approximately constant behaviour across all bins. However, both metrics exhibit a drop at low energies. This loss is accentuated in the endcaps, especially in the backward region.

Two effects occur that cause a loss in efficiency. Either the model does not predict a cluster, or the position and energy prediction exhibit a discrepancy large enough that the prediction is not matched to the target cluster. For the tower model, the case that no prediction is made occurs in 0.01% of the events, and the case of not matching a prediction to the target cluster occurs in 0.003% of events. For the GNN-ETM, no prediction is made at all in 0.004% of events, and ten times more seldom, the prediction is not matched. For the Flat 20 model, in only 27 of the over 3 million events a efficiency loss is reported due to no cluster prediction, and in 0.0008% the prediction is not matched to the target cluster. As expected and per the design of the algorithm and the dataset, the ICN-ETM predicts every cluster, and the matching criteria are loose enough to match all ICN-ETM predictions to the respective target clusters.

The purity loss is caused by either a deviating position or energy prediction, or by multiple predictions for only one target cluster. Hence, only one prediction is matched correctly to the target cluster, while the other prediction is considered false. In general, for the GNN-based models, this also results in two lower energetic predictions, as the model separates the cluster and its energy into two. The purity drop in all displayed models is dominated by the effect caused by multiple predictions per target cluster. Even so, the purity drop off the GNN-based models extends to higher predicted energies; the ICN-ETM is the model with the most double predictions. The double predictions of the ICN-ETM are mainly located at the gaps between the endcaps and the barrel region. If a TC is active in both the endcap and the barrel, the ICN-ETM predicts two clusters, from which only one is matched to the target offline cluster. In roughly 1.5% of events, the ICN-ETM logic leads to double predictions for the one target cluster, followed by the GNN-ETM at 0.6%, the Tower 5x5 model at 0.5% and the Flat 20 model at 0.3%. For the GNN-based models, the purity loss is extended to higher energies. The models separate the high-energy target cluster, leading to lower-energy cluster predictions. Due to the loose energy matching criteria, the target cluster is still considered found by one of those two predictions; therefore, not result in a bad efficiency. The second prediction then leads to a reduced purity. This is especially visible for the GNN-ETM and in the backward endcap, also for the Tower 5x5 model. Overall, the Flat 20 model performs most reliably across all energy and detector regions, with the largest instabilities in the backward detector.

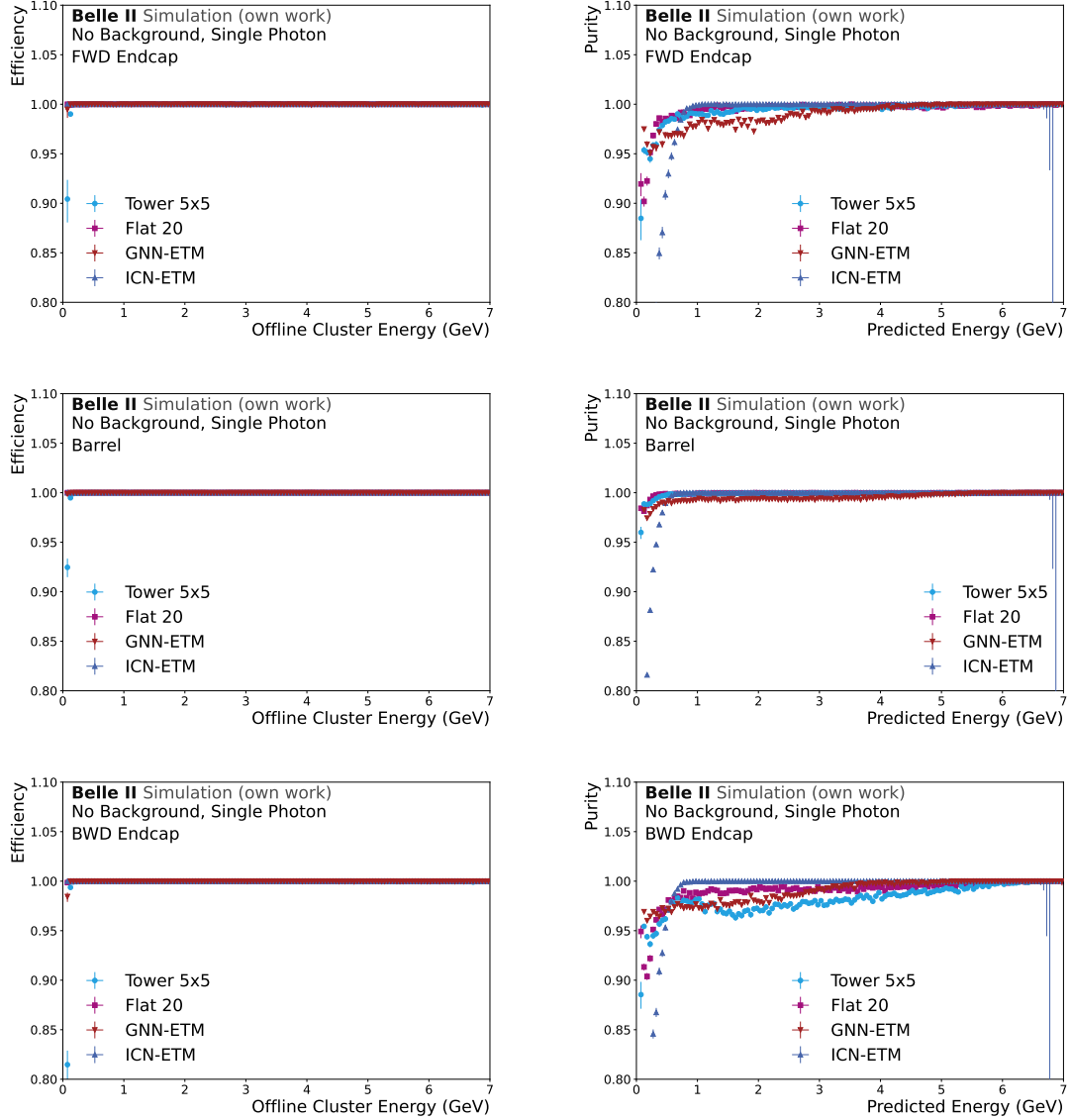


Figure 7.3.: Cluster finding efficiency and purity for the different detector regions in 50 MeV energy bins. Determined for the high-granularity Tower 5x5 and Flat 20 and the low granularity ICN-ETM and GNN-ETM. For each model, the respective targets in the single photon sample are used. A selection of only events with one target cluster is applied.

In this and all following resolution determinations, the subset of MC-matched target clusters found by all models, displayed in the respective plot, is used to determine the resolution metrics. The position resolution with respect to the MC position is displayed in fig. 7.4, which exhibits a clear improvement of the Tower 5x5 and Flat 20 model compared to the GNN-ETM and ICN-ETM. This improvement is roughly a 4 times improvement, which is consistent with the granularity change of the input. The offline reconstruction, however, still exhibits a clear performance advantage. The energy resolution binned in the detector regions and for varying MC energies is shown in fig. 7.5. The overall trend exhibits a reduced energy resolution at lower cluster energies. In general, the GNN-ETM performs worst, while the Tower 5x5, Flat 20 and ICN-ETM do not show a consistent hierarchy. The Offline reconstruction performs best, with the largest performance gain at low energies. Comparing the two input reduction methods, at lower energies in the endcaps, the Tower 5x5 model performs better, while for larger energies, especially the Flat 20 model excels.

### 7.2.2. Quantisation Comparison

Comparing the non-quantised Flat 20 and Tower 5x5 models to the Quantised model in fig. 7.6, both cluster-finding metrics behave similarly as above. A constant behaviour close to 1 is observed across all energy and detector regions, with a drop off towards lower energies. Furthermore, a small dip in efficiency can be observed for the Quantised model in the area around 0.5 GeV and a small efficiency drop off toward high energies. The efficiency reduction of the Quantised model is primarily caused by the reduced resolution. Therefore, the prediction is matched more seldom to the target cluster. While the fraction of events, for which no prediction was made, is similar to the Tower 5x5 model with roughly 0.01%, the proportion of events in which the matching criteria was not met is significantly larger with 0.18%.

Regarding the purity, the Tower 5x5 model exhibits the accentuated purity loss in the backward detector, which is caused by the multiple predictions. The Flat 20 model exhibits only a slight decrease towards lower predicted energies, which is also present in the Quantised model. The Quantised model furthermore exhibits a dip at around 1 GeV, which is present in all detector regions, most visibly in the barrel region. This drop is caused by double predictions on clusters above 4 GeV, for which multiple predictions are made. The higher energetic predictions are overall close to the target energy, while the distribution of the lower energetic predictions peaks around zero and also around 1 GeV. This effect is also visible for the Tower 5x5 model, however, shifted to higher energies. For both models, the effect can not be easily reduced by a higher  $\beta$ -cut. As this leads to an increasing efficiency loss.

Regarding the position resolution, shown in fig. 7.7, a performance loss caused by the quantisation of the model is apparent. Only the backward endcap is shown here, as the forward endcap and barrel region look similar, as can be seen in fig. A.1. This performance decrease varies between 10% and 30%. The overall shape, between the different models, stays consistent across the different detector regions. Clearly visible is a performance increase of the Flat 20 model compared to the Tower 5x5 model, however, only in the x-coordinate prediction. This effect is only visible in this specific training of the model and spans across the different evaluated samples.

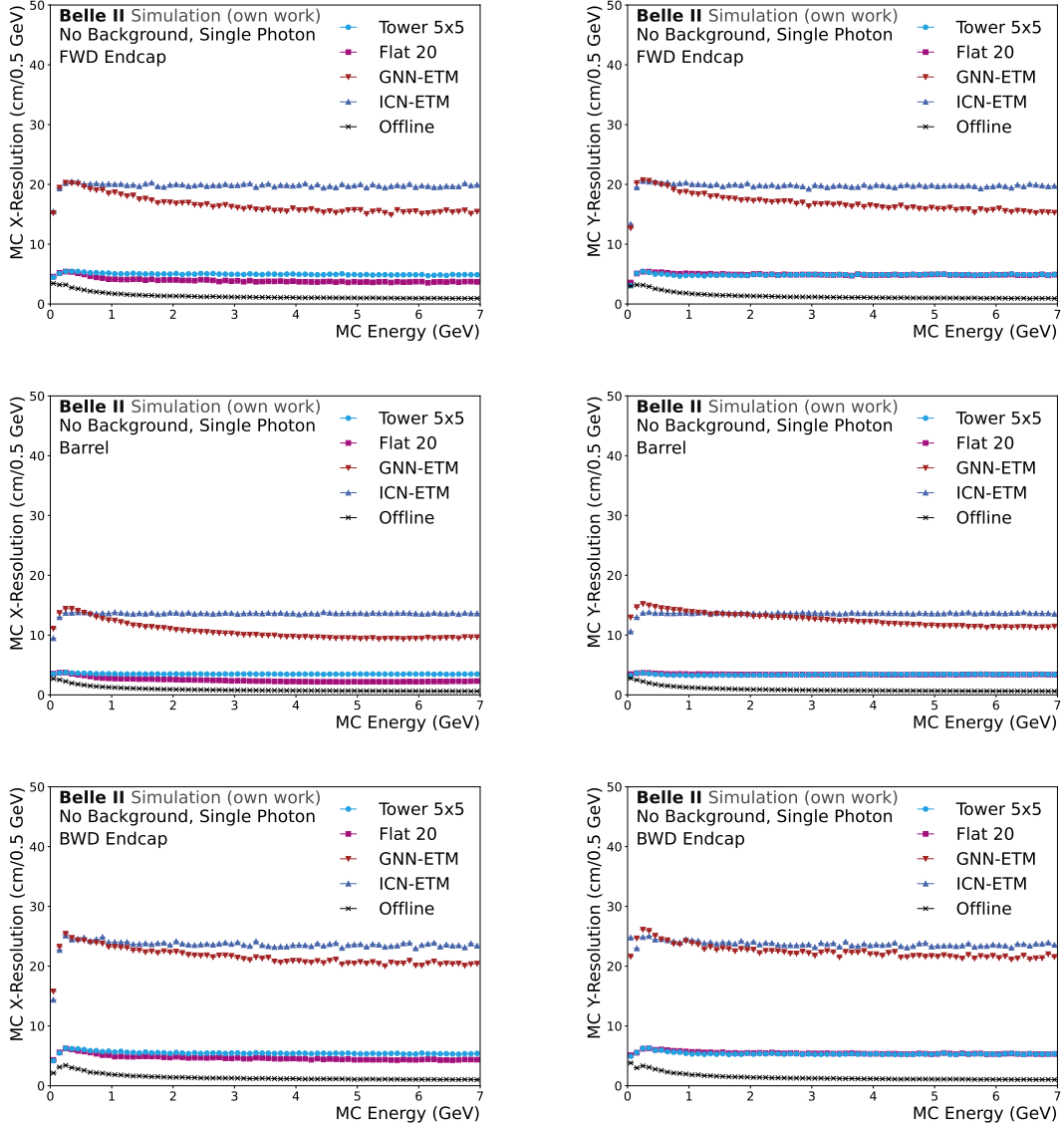


Figure 7.4.: Position Resolution in x and y direction for the different detector regions. Determined for the high-granularity Tower 5x5 and Flat 20 and the low granularity ICN-ETM and GNN-ETM on the single photon sample. A selection of only events with one target cluster is applied. Furthermore, the resolution is determined on the subset of MC-matched target clusters found by all models. As the target position, the extrapolated MC-information of the simulated photon is used. For comparison, the performance of the offline reconstruction on the same subset of clusters is shown in black.

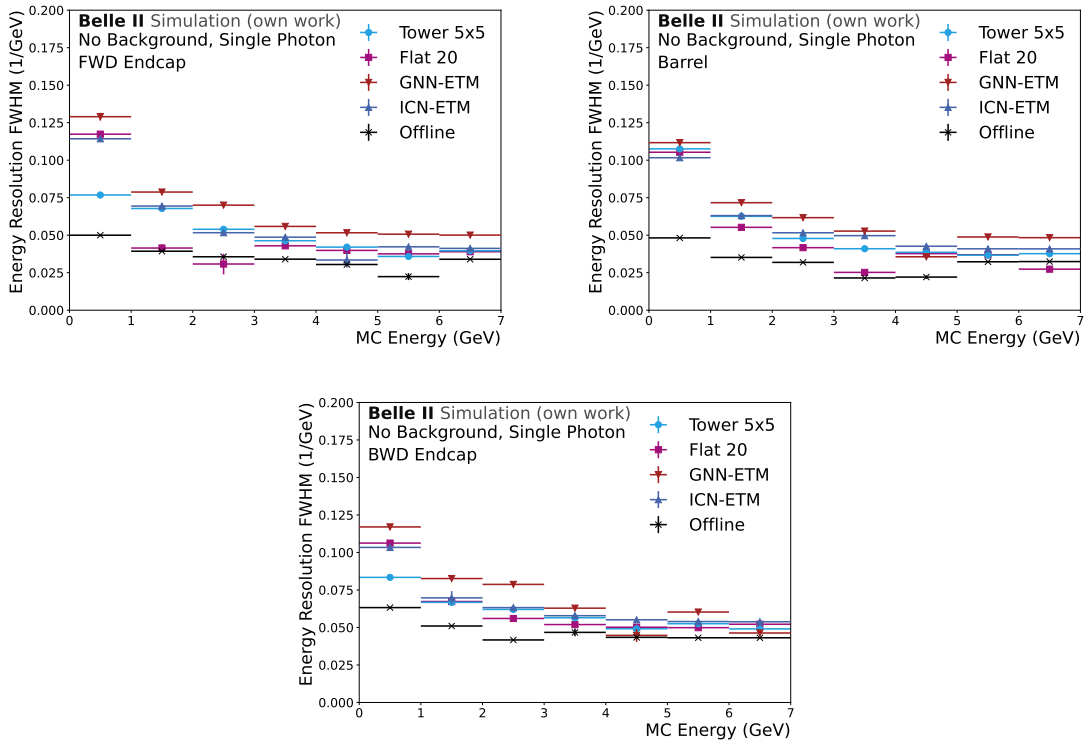


Figure 7.5.: Energy Resolution in the different detector regions. Determined for the high-granularity Tower 5x5 and Flat 20 and the low granularity ICN-ETM and GNN-ETM in the single photon sample. A selection of only events with one target cluster is applied. Furthermore, the resolution is determined on the subset of MC-matched target clusters found by all models. As the target energy, the MC-information of the simulated photon is used. For comparison, the performance of the offline reconstruction on the same subset of clusters is shown in black.

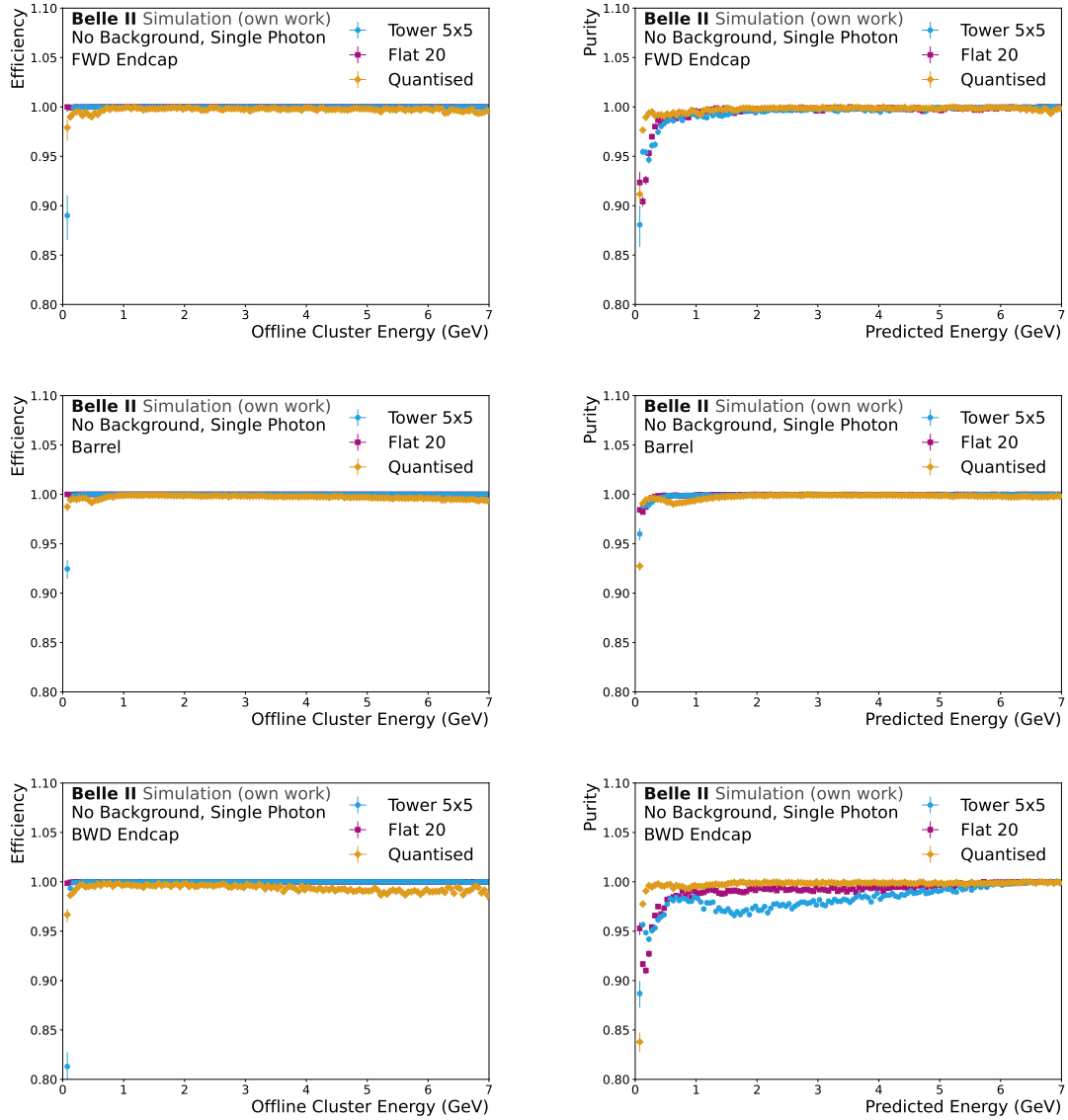


Figure 7.6.: Cluster finding efficiency and purity for the different detector regions in 50 MeV energy bins. Determined for the two non-quantised Tower 5x5 and Flat 20 models and the Quantised model. For each model, the respective targets in the single photon sample are used. A selection of only events with one target cluster is applied.

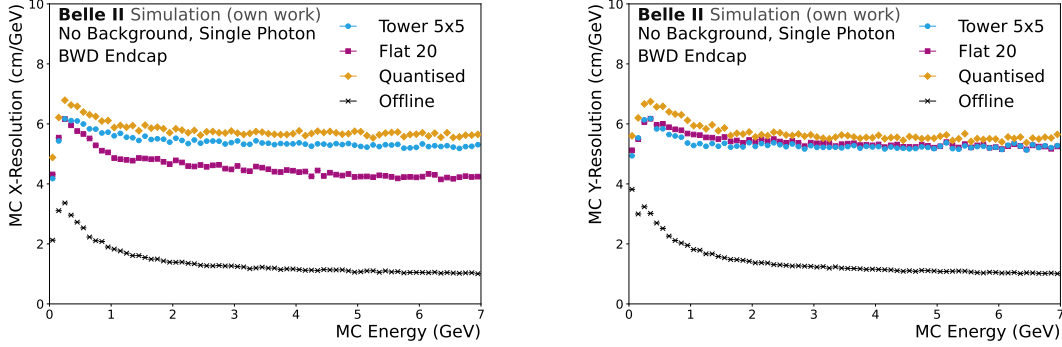


Figure 7.7.: Position Resolution in x and y direction for the backward endcap region. Determined for the two non-quantised Tower 5x5 and Flat 20 models and the Quantised model. A selection of only events with one target cluster is applied. Furthermore, the resolution is determined on the subset of MC-matched target clusters found by all models. As the target position, the extrapolated MC-information of the simulated photon is used. For comparison, the performance of the offline reconstruction on the same subset of clusters is shown in black.

For the energy resolution, seen in fig. 7.8, a more distinct performance decrease can be observed. The energy resolution of the Quantised model is up to twice as broad as that of the other two models. This loss in performance is traced back to both the quantisation and especially the strict reduction in the message passing step of the model. As most of the relevant target clusters are comprised of more than 8 crystals, the message passing in the GravNet layer is not extensive enough to incorporate all crystals of the underlying target cluster. Making it harder for the model to aggregate the information of all associated crystals within one node.

To conclude, the evaluation on the single photon sample shows a comparable performance in cluster-finding metrics for the high-granularity models, as well as a drastically improved position resolution. The energy resolution is improved over the GNN-ETM and behaves more like the ICN-ETM. Both the tower and flat energy cut lead to similar performing models in cluster finding, position and energy prediction. With a purity loss of the Tower 5x5 model in the backward endcap. Regarding the full hardware implementation, the quantisation of the model and the reduced message passing hinder the Quantised model, especially with regards to its energy resolution. The cluster finding metrics and position resolution also exhibit slight reductions, but the model still performs reliably.

### 7.3. Bhabha Sample

The Bhabha sample  $e^+e^- \rightarrow e^+e^-(\gamma)$  is the most common event type the L1 trigger encounters. As Bhabha events have to be identified and filtered by the L1 trigger, the clustering performance is important on this sample. Furthermore, by analysing this Bhabha sample, the clustering performance up to the highest expected cluster energies is probed realistically. To focus on the performance of the high-energy clusters, the applied selection

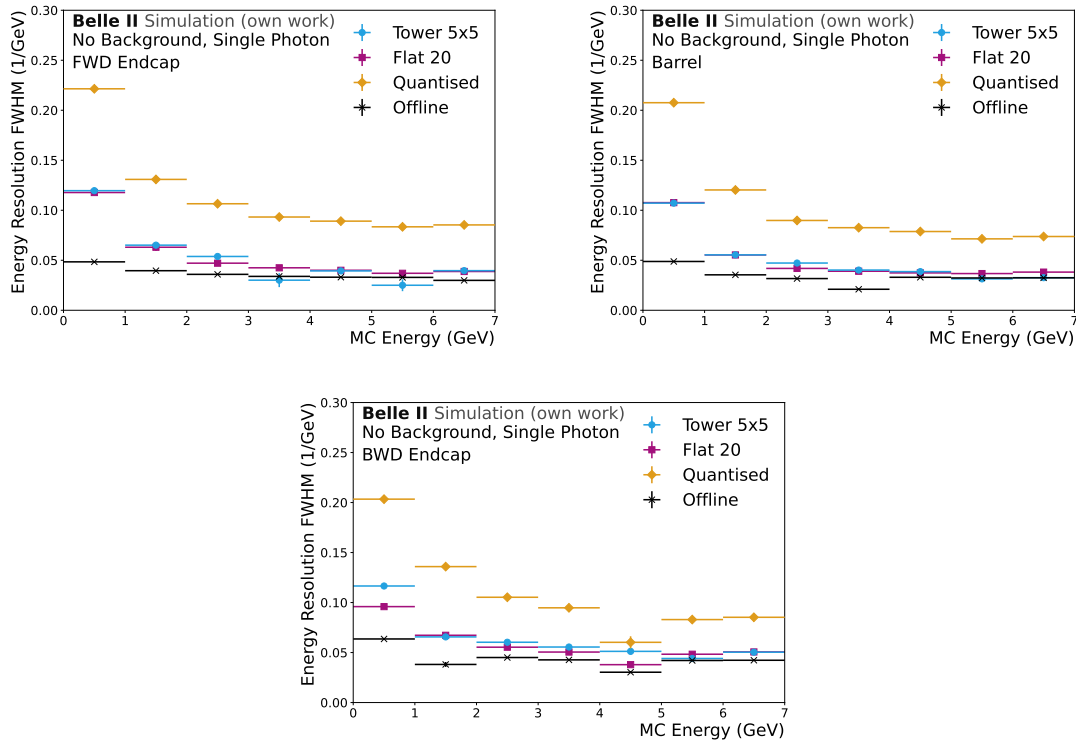


Figure 7.8.: Energy Resolution in the different detector regions. Determined for the two non-quantised Tower 5x5 and Flat 20 models and the Quantised model. A selection of only events with one target cluster is applied. Furthermore, the resolution is determined on the subset of MC-matched target clusters found by all models. As the target energy, the MC-information of the simulated photon is used. For comparison, the performance of the offline reconstruction on the same subset of clusters is shown in black.

removes a large part of the radiative Bhabha events. The most relevant clusters to be found are the positron clusters at around 4 GeV and the electron clusters at around 7 GeV. The sample uses the extrapolated Scenario 2 background, resulting in the low-energy regime, being dominated by background clusters, as can be seen in section 6.4.2. As this work focuses on the high-granularity models and the study of the input reduction methods, the performance of the ICN-ETM and GNN-ETM is not included. Especially as the lower granularity leads to an undesired subsampling of commonly found clusters for the resolution metrics. Introducing more artefacts and making the resulting distributions less significant.

In fig. 7.9, the comparison of the cluster-finding metrics of the high-granularity models is shown. Compared to the single photon sample, a much different behaviour of the metrics is visible. At low energies of up to 500 GeV, the target cluster distribution is background dominated. These background clusters are distributed across the whole detector and, therefore, isolated. However, the low energy makes them hard to distinguish from the other low-energy depositions that do not have an offline cluster attached. This leads to a reduced but robust performance metric. Additionally, there are also MC-matched target clusters with low energies following a distribution extending to much higher energies. These target clusters are caused by bremsstrahlung of the primary electron and positron, leading to comparably lower-energy clusters. As the distance between these clusters and the cluster caused by the primary lepton varies, these Bremsstrahlung clusters exhibit a varying amount of separation. These clusters are even harder to distinguish due to their lower energy compared to the overlaying primary cluster. This directly translates to an efficiency loss at energies of up to 4 GeV. Especially the Quantised model struggles to separate and identify these clusters. For the Quantised model, a major restriction is set by the reduced message passing, and hence reduced perceptive field. This inhibits the model from grasping the total two-cluster topology and separating them correctly. The energy range from 3.5 to 7 GeV is dominated by the electron and positron clusters; these clusters are high-energy clusters, making them better separated and easy to find. To support the above argument, target clusters out of the three different energy regimes are selected, and the distance to the closest other target clusters is shown in fig. 7.10. In all shown distributions, the impact of the radiative clusters is clearly visible with the peak at low distances. For low energies, the random positions of the background clusters lead to a pronounced tail extending to higher distances. For the highest energy bin, the distribution extends across the whole detector, and features a peak above 3 m at which the second particle of the Bhabha scattering is located.

Due to the event topology, no target clusters above 5 GeV are present in the backward endcap, leading to an undefined efficiency. While in the forward endcap, the radiative regime extends to 6 GeV. A further effect resulting in an efficiency reduction is caused by clusters close to the detector gaps between the endcaps and barrel, where only a part of the cluster is visible, and the Quantised model struggles to predict and reconstruct the cluster. However, this effect is numerically much less significant. The purity of the Flat 20 model exceeds that of the Tower 5x5 and Quantised model, which are reduced due to multiple predictions, leading to excessive predictions in the intermediate energy region. In combination with the efficiency loss, it is apparent that both models, the Tower 5x5 and the Quantised model, struggle to consistently separate and predict the correct number of clusters for these overlapping signatures.

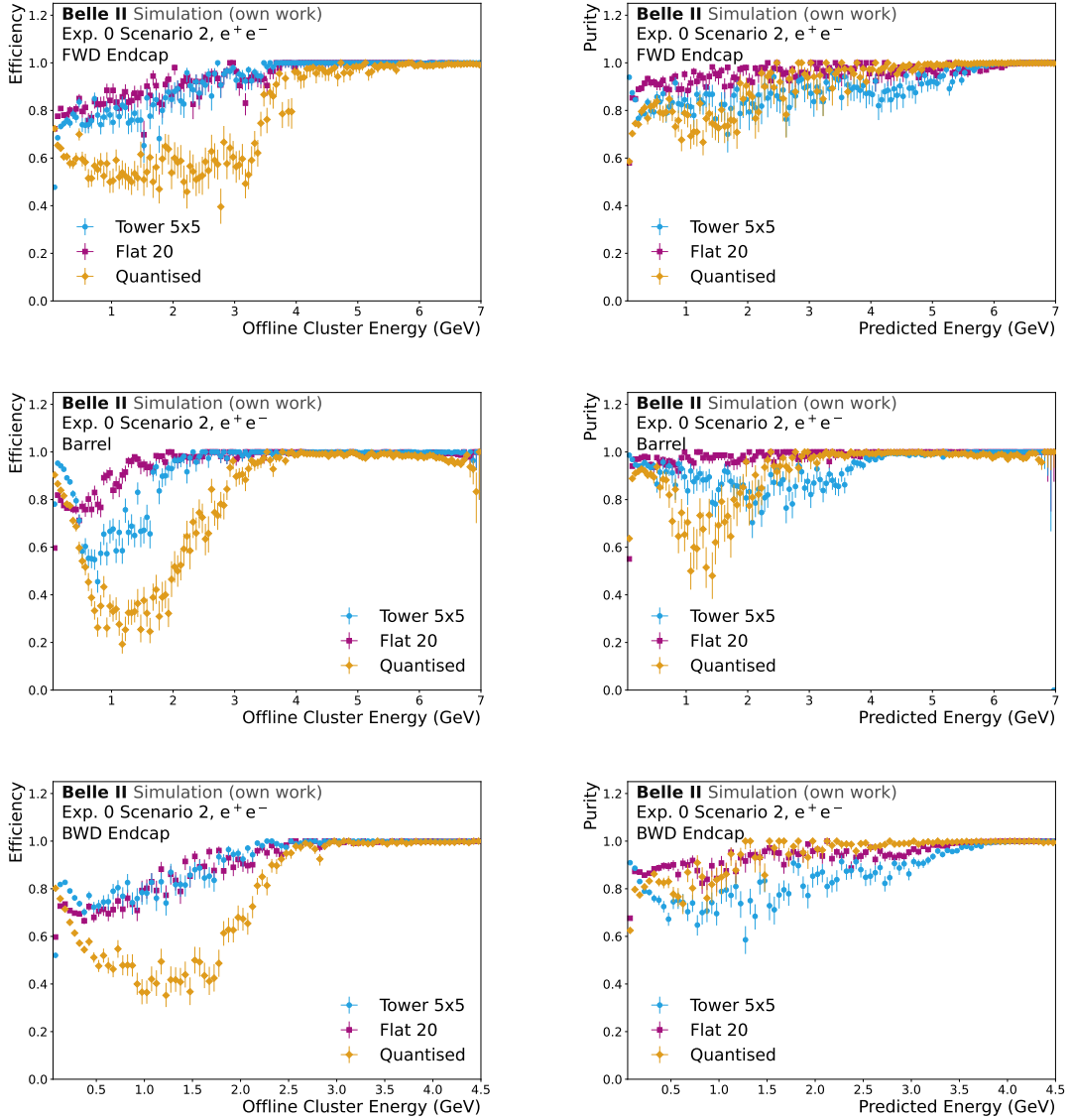


Figure 7.9.: Cluster finding efficiency and purity for the different detector regions in 50 MeV energy bins. Determined for the two non-quantised Tower 5x5 and Flat 20 models and the Quantised model. For each model, the respective targets in the Bhabha sample, with the in section 6.4.2 defined selections applied, are used.

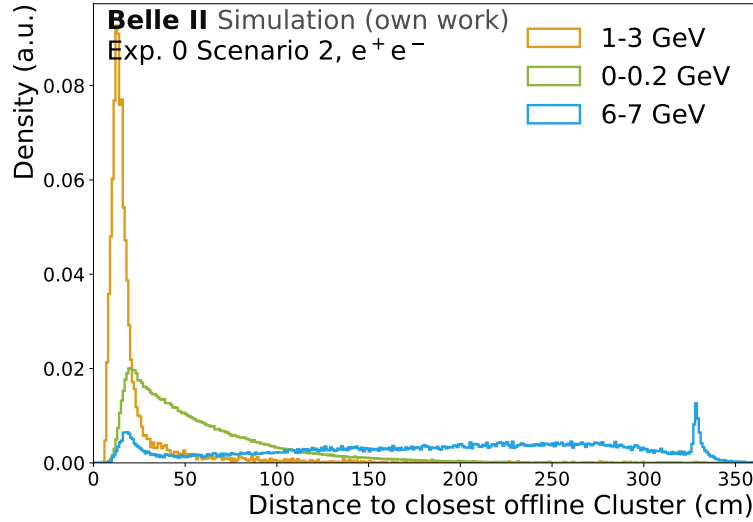


Figure 7.10.: Distance of selected offline clusters to the closest next offline cluster. Shown for low, medium and high cluster energies. This shows the difficulty in reconstructing the clusters with medium energies, as they are close to other offline clusters. Evaluated on the Bhabha sample with the, in section 6.4.2, defined selections applied.

Regarding the position resolution on MC matched target clusters, shown in fig. 7.11, all three models behave similarly, with the Flat 20 and Tower 5x5 model performing more precisely compared to the Quantised model. A worse resolution is observed in the endcaps compared to the barrel. Overall, the statistical fluctuations in the less populated energy regions are clearly visible. For the highly populated regimes, the models exhibit the same hierarchy as for the single photon sample.

Due to the sparse overall population at lower energies, the energy resolution is only determined for energies above 3.25 GeV. As the backward endcap is only populated by positrons and the forward endcap by electrons, the energy range for those detectors is adjusted to the respective peaks. The resulting plots are shown in fig. 7.12. The Flat 20 and Tower 5x5 models perform similarly and approach the performance of the offline resolution for higher energies. The Quantised model exhibits an up to twice as wide energy resolution as the other two models. This reduced energy resolution is partially caused by the inability of the model to separate overlapping clusters.

In conclusion, the Quantised model performs less reliably when tasked to separate the overlapping radiative clusters. Which can mainly be attributed to the restricted message passing and the large high energetic clusters. Nevertheless, all models perform well in finding the primary high-energy clusters, even under the increased background levels.

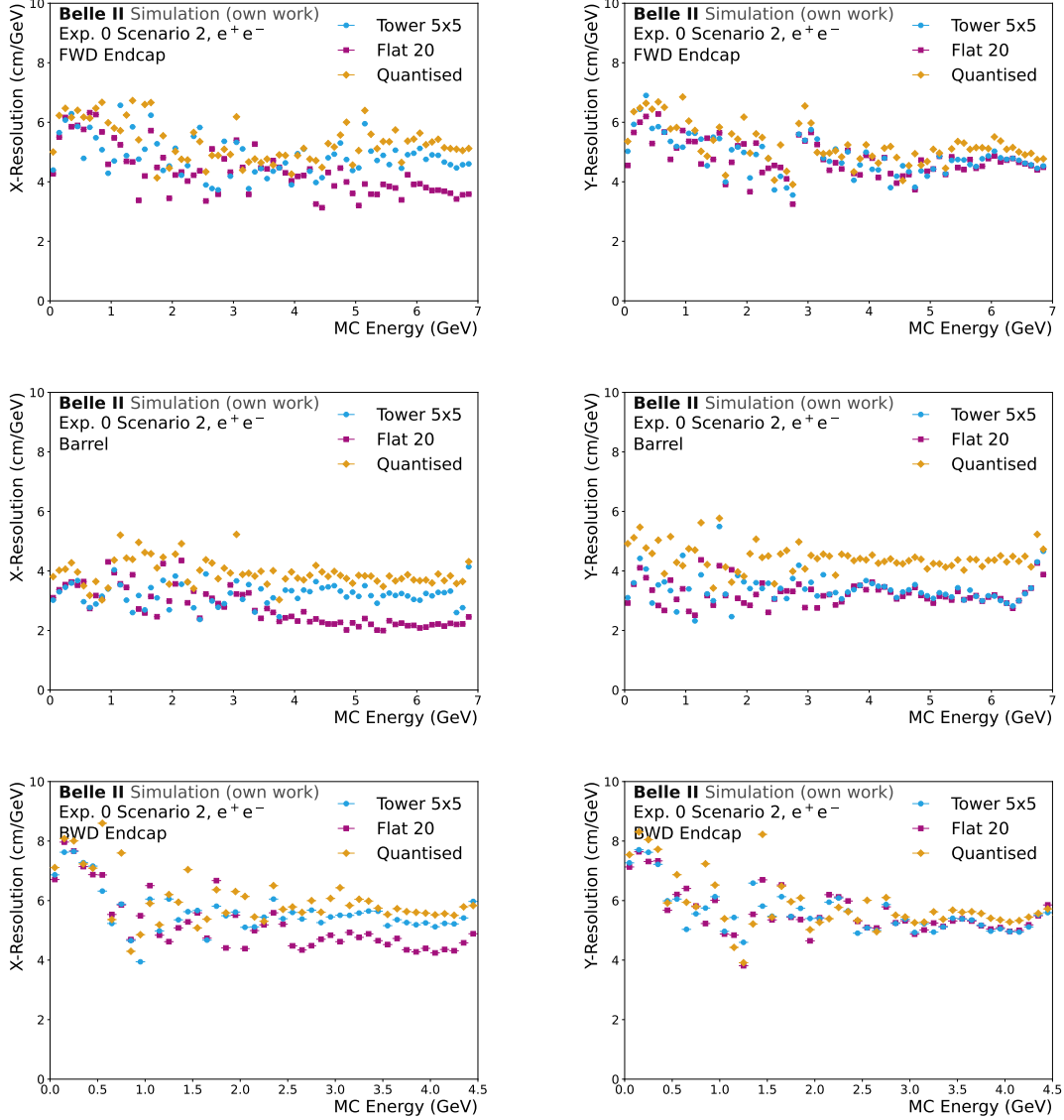


Figure 7.11.: Position Resolution in x and y direction for all detector regions. Determined for the two non-quantised Tower 5x5 and Flat 20 models and the Quantised model. The resolution is determined on the subset of MC-matched electrons, positrons and photons with associated clusters, found by all models. As the target position, the offline cluster information is used.

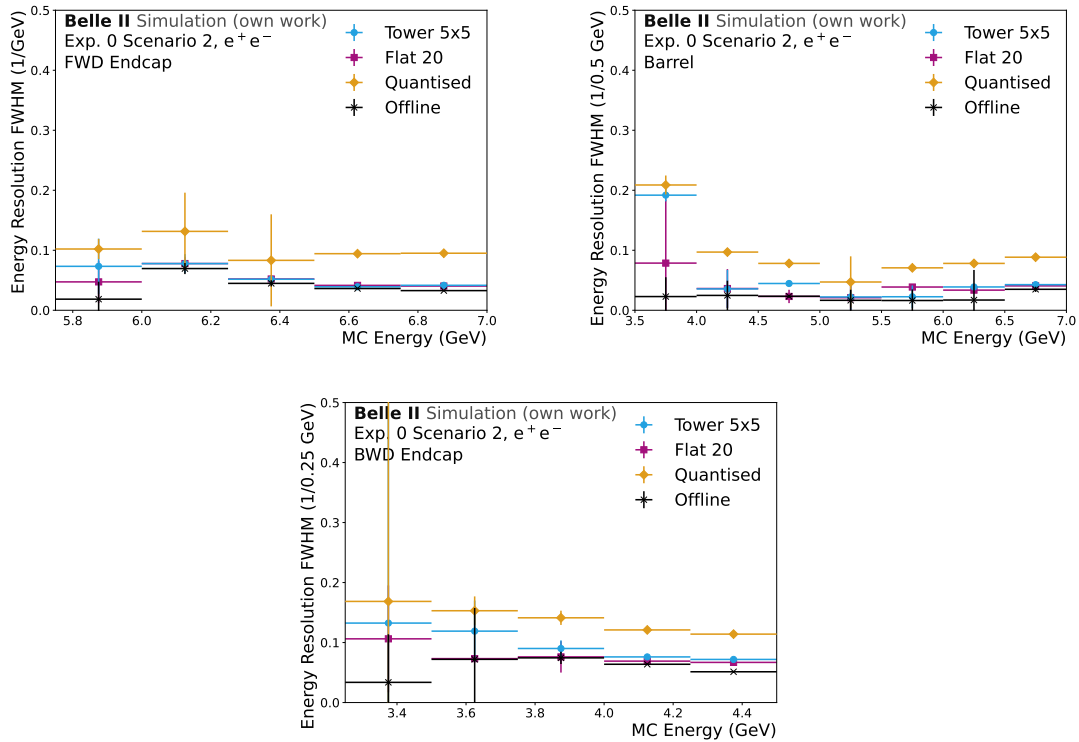


Figure 7.12.: Energy Resolution in the different detector regions. Determined for the two non-quantised Tower 5x5 and Flat 20 models and the Quantised model. The resolution is determined on the subset of MC-matched electrons, positrons and photons with associated clusters, found by all models. As the target energy, the MC-information of the simulated electrons and positrons is used. For comparison, the performance of the offline reconstruction on the same subset of clusters is shown in black.

## 7.4. Dimuon Sample

The relevant cluster energies for the dimuon sample are limited to about 500 MeV, due to the minimal ionising properties of muons. At energies below 200 MeV, the Scenario 2 background dominates the target cluster distribution. On this sample, the performance on the muonic cluster shape is probed. As this is not resolved by the  $4 \times 4$  TCs, only the high-granularity models are studied. In fig. 7.13, the cluster-finding metrics evaluated on the dimuon sample are displayed. The Tower 5x5 model outperforms both other models at low and high energies, in terms of cluster-finding efficiency.

The main cause for the Flat 20 and Quantised models to lose efficiency at low energies is rooted in target clusters comprised of only a singular crystals. These single-crystal targets are very hard to distinguish from the background crystals that have no cluster associated with them, as only the time and energy of the crystal can be used. With respect to the number of target clusters in these regions, it is visible that the flat 20 MeV cut retains many more target clusters, which are caused by background. Therefore, the tower preprocessing provides a much cleaner input, making it easier for the model to predict the remaining clusters correctly.

At high energies, the efficiency reduction is caused by two effects, first of all, radiative muon events, where the emitted photon exhibits a low enough energy to pass the applied selection and overlaps with the muonic cluster. In these cases, it is much more difficult to identify the muonic signature on top of the larger electromagnetic cluster. The Tower 5x5 model still predicts these more difficult cluster topologies, even though it was not trained on these. The second effect is caused by clusters which are close to the edge of the timing window. A subset of the crystals associated with the cluster is retained by the applied timing cut, while the rest is not. Especially in the flat preprocessing, singular low-energy crystals are retained, leading to a background-like signature, with a higher-energy target cluster associated. This effect is reduced by the tower cut, compared to the flat cut.

With respect to the position resolution on MC matched target clusters, shown in fig. 7.14, all models exhibit similar performance, with the Quantised model performing overall a bit worse. In the energy resolution, computed with regard to the offline cluster energy, shown in fig. 7.15, the trend established in the single photon sample continues with the Quantised model performing much worse in predicting the energy of the cluster correctly. The Flat 20 and Tower 5x5 models are performing comparably, with a slight advantage of the Flat 20 model.

Concluding on the dimuon sample, the performance of the two non-quantised models is rather similar, except that the Tower 5x5 model exhibits a much improved and more stable efficiency across the whole detector and energy range, caused by the inherent background reduction of low-energy crystals. The Quantised model performs similarly to the Flat 20 model in the cluster finding metrics, though its resolution is reduced.

## 7.5. $B\bar{B}$ Sample

The  $B\bar{B}$  sample probes more densely populated event topologies. These also include a more varied selection of cluster topologies, as also hadronic showers are included. For the analysis

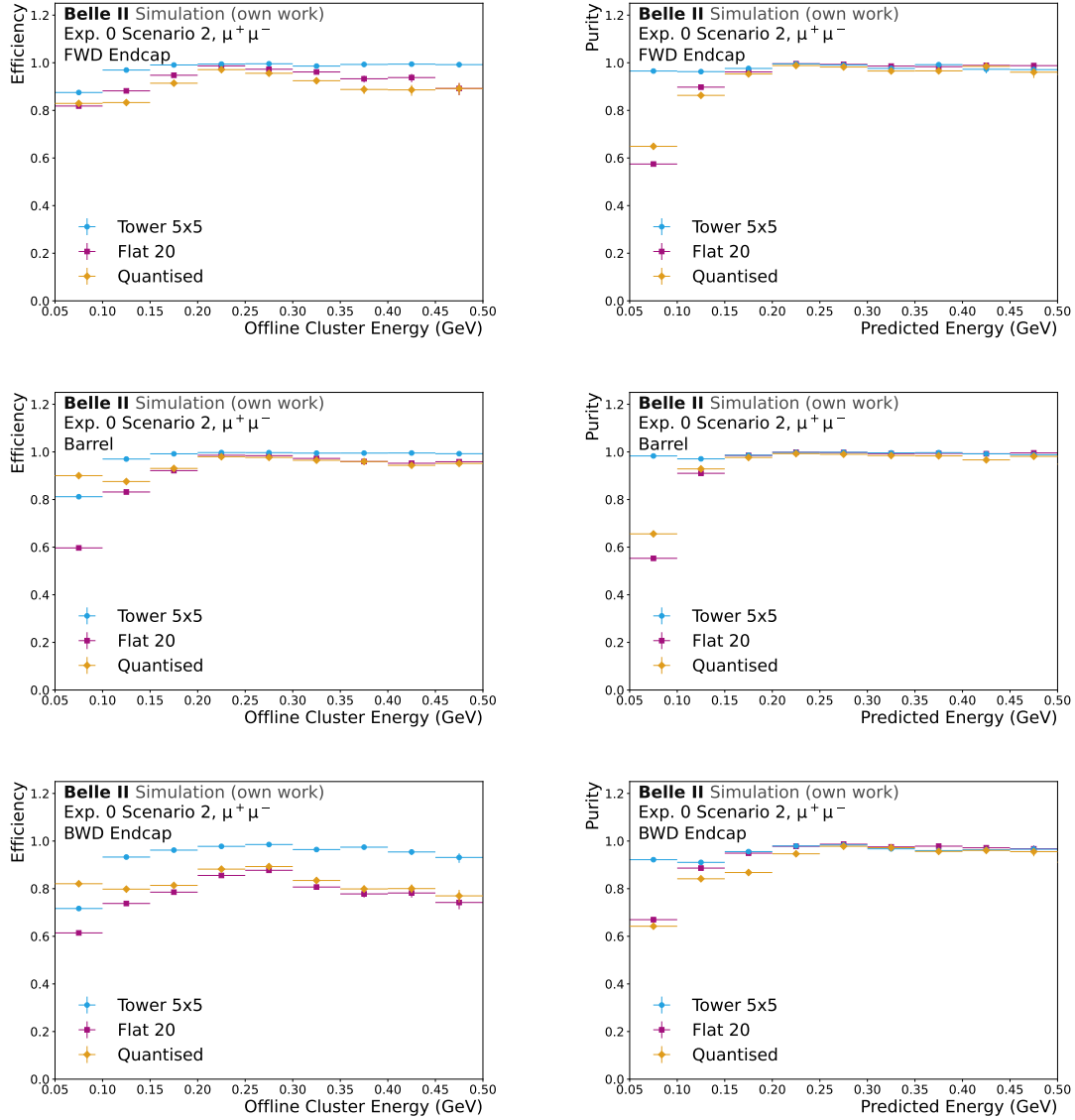


Figure 7.13.: Cluster finding efficiency and purity for the different detector regions in 50 MeV energy bins. Determined for the two non-quantised Tower 5x5 and Flat 20 models and the Quantised model. For each model, the respective targets in the dimuon sample, with the in section 6.4.3 defined selections applied, are used.

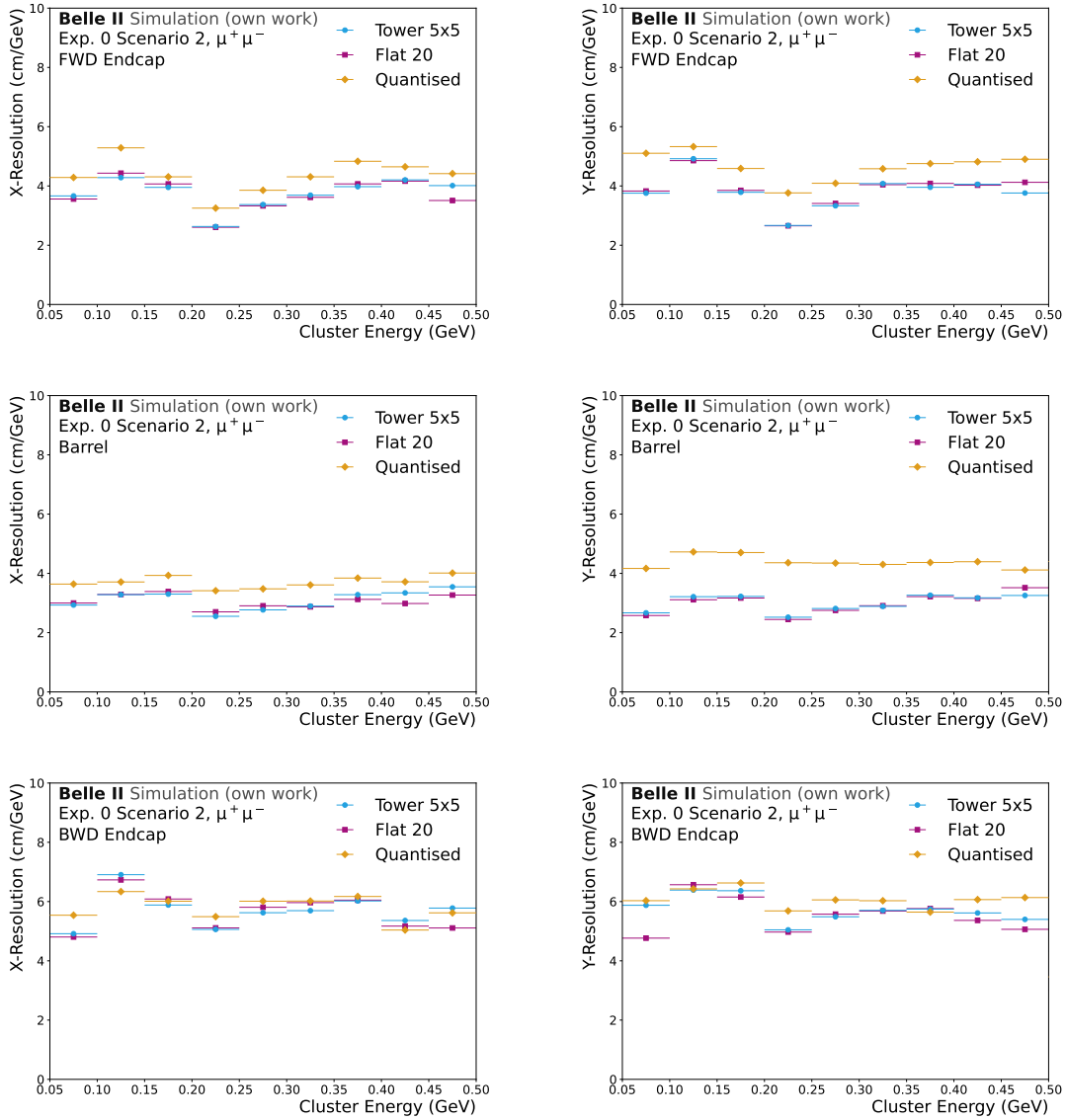


Figure 7.14.: Position Resolution in x and y direction for all detector regions. Determined for the two non-quantised Tower 5x5 and Flat 20 models and the Quantised model. The resolution is determined on the subset of MC-matched clusters, found by all models. As the target position, the offline cluster information is used.

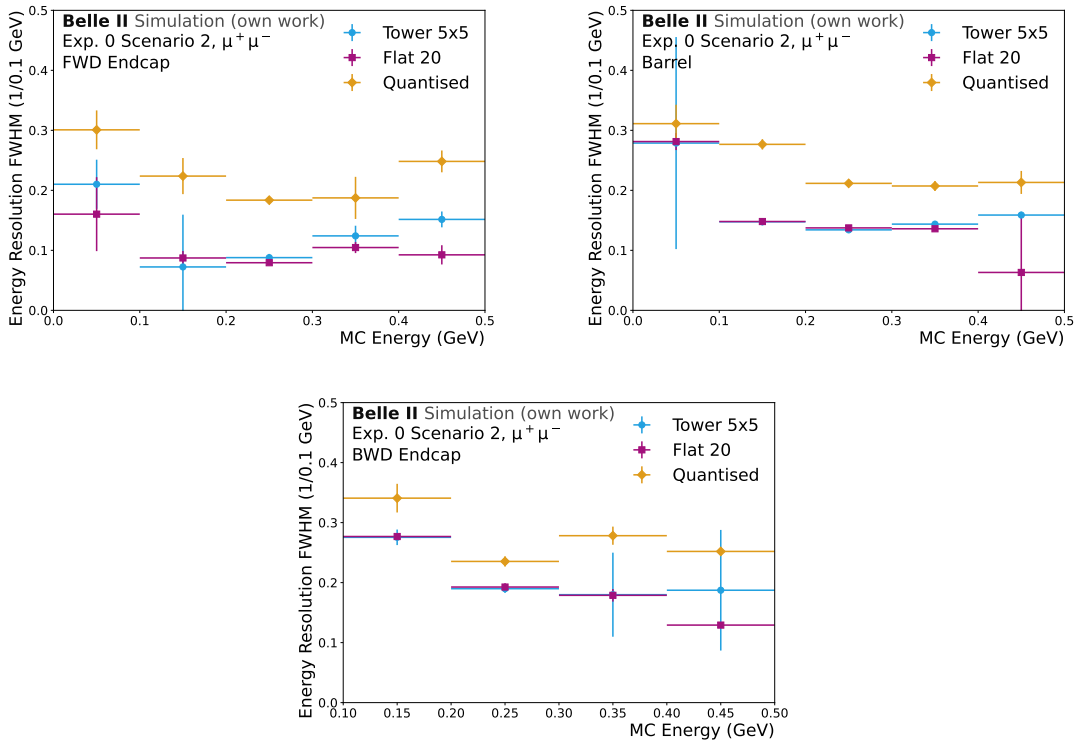


Figure 7.15.: Energy Resolution in the different detector regions. Determined for the two non-quantised Tower 5x5 and Flat 20 models and the Quantised model. The resolution is determined on the subset of MC-matched clusters, found by all models. As the target energy, the MC-information of the simulated particles, matched to the clusters, is used. For comparison, the performance of the offline reconstruction on the same subset of clusters is shown in black.

of the ICN-ETM, the artificial limit of only 6 clusters per event is removed, as this leads to an artificial loss in found clusters, and therefore reduces the statistics for this analysis.

### 7.5.1. Granularity Comparison

The cluster-finding metrics computed on all respective target clusters, without any selections applied, are shown in fig. 7.16. The GNN-based models exhibit a robust performance, with slight deviations from 1, especially within the lower energy regime. Both the efficiency and purity decrease at low energies are primarily caused by single-crystal contributions, which, apart from an energy-based classification, only exhibit timing information to classify as a cluster candidate or not. At higher energies, even larger clusters may not produce a condensation point and therefore cluster prediction. For the Flat 20 and Tower 5x5 model, this is often the case when the predicted cluster is close to a different higher-energy cluster. The Flat 20 model exhibits a notable efficiency loss for clusters contained in the backward endcap, compared to the other detector regions. This can be attributed to the much larger number of clusters in this region, which are retained by this input reduction method. However, as these additional clusters are comprised of a few crystals, the reconstruction of these clusters is challenging.

For the ICN-ETM, a dramatic performance loss is observed. The purity loss can be attributed to the larger amount of single TC depositions, which are caused solely by the beam background and do not have an offline cluster attached, as well as the double predictions, e.g. caused by the  $\square$ -pattern. Furthermore, the effect of extended clusters spanning across the detector gaps comes into play. The efficiency loss is caused by interconnected TC patterns, which are caused by multiple offline clusters, but the ICN-ETM only predicts one cluster, for example, in the case of the  $\square$  pattern.

For the position resolution, shown in fig. 7.17, only MC-matched photons are used as targets. As the MC-matching to the offline clusters struggles at low energies, following additional selection is applied to the clusters. Retaining clusters for which the relative ratio between the offline cluster energy  $E_{\text{offline}}$  and matched MC-energy  $E_{\text{MC}}$  fulfils

$$\frac{E_{\text{offline}} - E_{\text{MC}}}{E_{\text{MC}}} \in [-1, 1]. \quad (7.3)$$

Resulting in a similar model hierarchy as on the single photon sample. The backward endcap exhibits larger statistical fluctuations and an overall reduced position resolution, also by the offline reconstruction.

The improvement within the energy resolution for larger MC energies is also observed for the  $B\bar{B}$  sample in fig. 7.18. For the energy resolution, only the MC-matched photons are used. As the cluster finding performance of the ICN-ETM breaks down in this sample, the energy resolution is not determined for the ICN-ETM, as the underlying distributions no longer follow a double-sided crystal ball, but are composed of multiple peaks. By only fitting the peak closest to 0, the energy resolution of the ICN-ETM becomes less significant as only a small fraction of the found clusters is contained in this peak. Furthermore, the selection of only commonly found clusters further affected the distributions of the other models. Hence, the ICN-ETM is excluded for this part of the analysis.

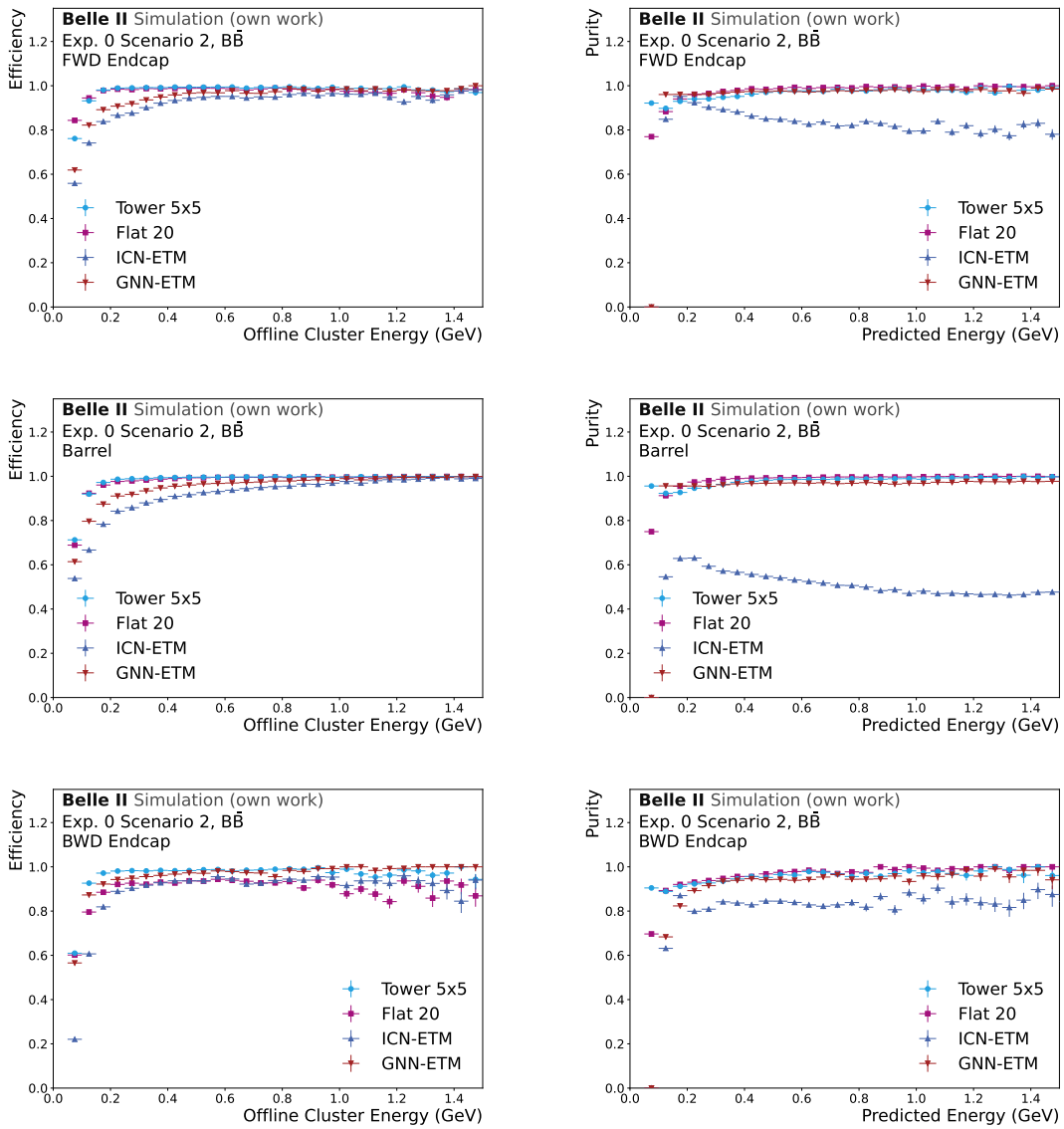


Figure 7.16.: Cluster finding efficiency and purity for the different detector regions in 50 MeV energy bins. Determined for the high-granularity Tower 5x5 and Flat 20 and the low granularity ICN-ETM and GNN-ETM. For each model, the respective targets in the  $B\bar{B}$  sample are used.

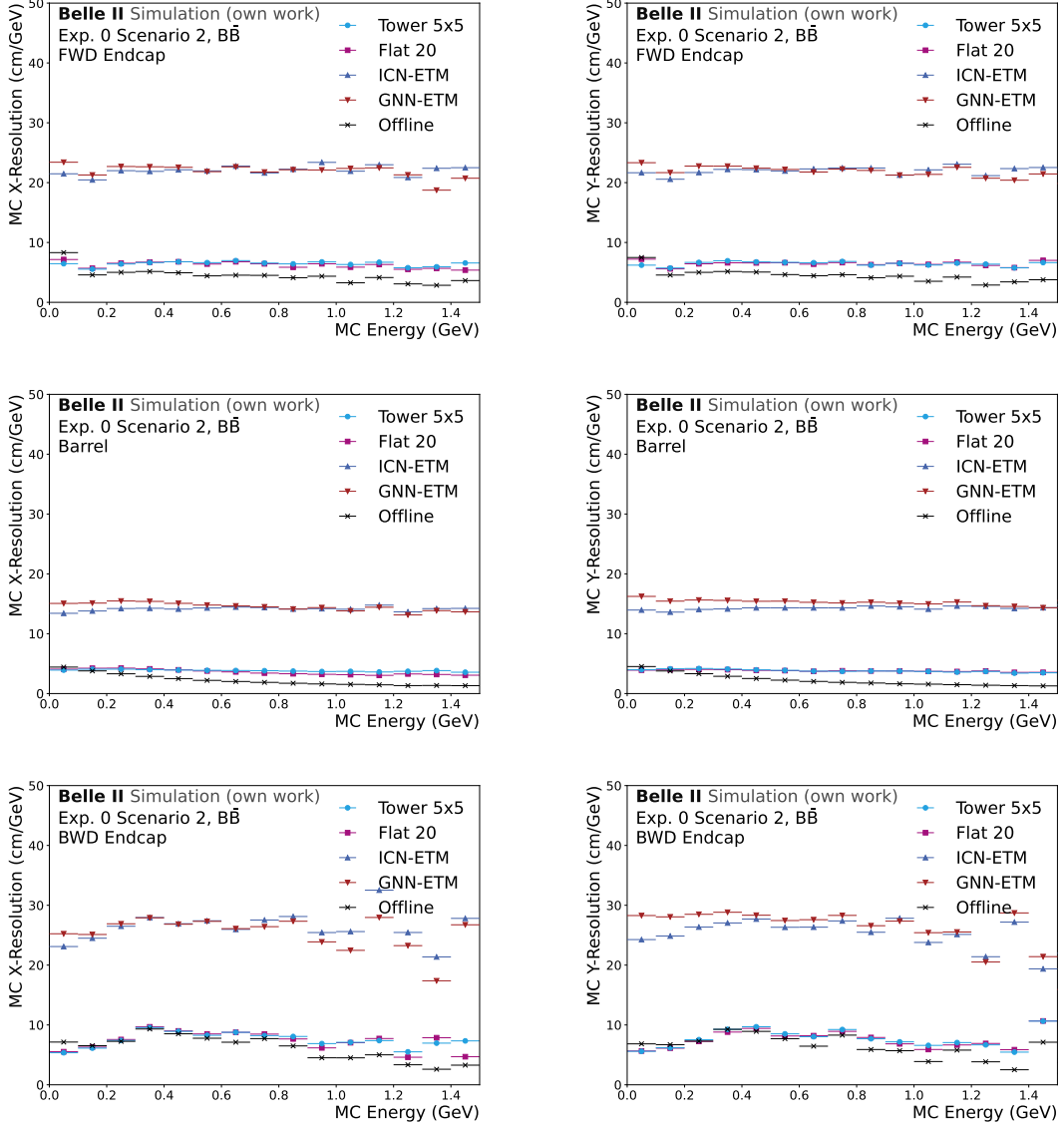


Figure 7.17.: Position Resolution in x and y direction for the different detector regions. Determined for the high-granularity Tower 5x5 and Flat 20 and the low granularity GNN-ETM on the  $BB$  sample. The resolution is determined on the subset of target clusters MC-matched to photons, found by all models. As the target position, the extrapolated MC-information of the underlying photon is used. For comparison, the performance of the offline reconstruction on the same subset of clusters is shown in black.

The resulting energy resolutions exhibit larger uncertainties, especially in the endcaps. This is driven by the lower statistics and the artefacts introduced by, e.g. double predictions, affecting the underlying distributions. Hence, increasing the uncertainties of the underlying fits. The Flat 20 model proves to be more robust than the other two models. While both the GNN-ETM and the Tower 5x5 model still perform reasonably well. Compared to the single photon sample, all models, as well as the offline reconstruction, exhibit a large decrease in overall performance.

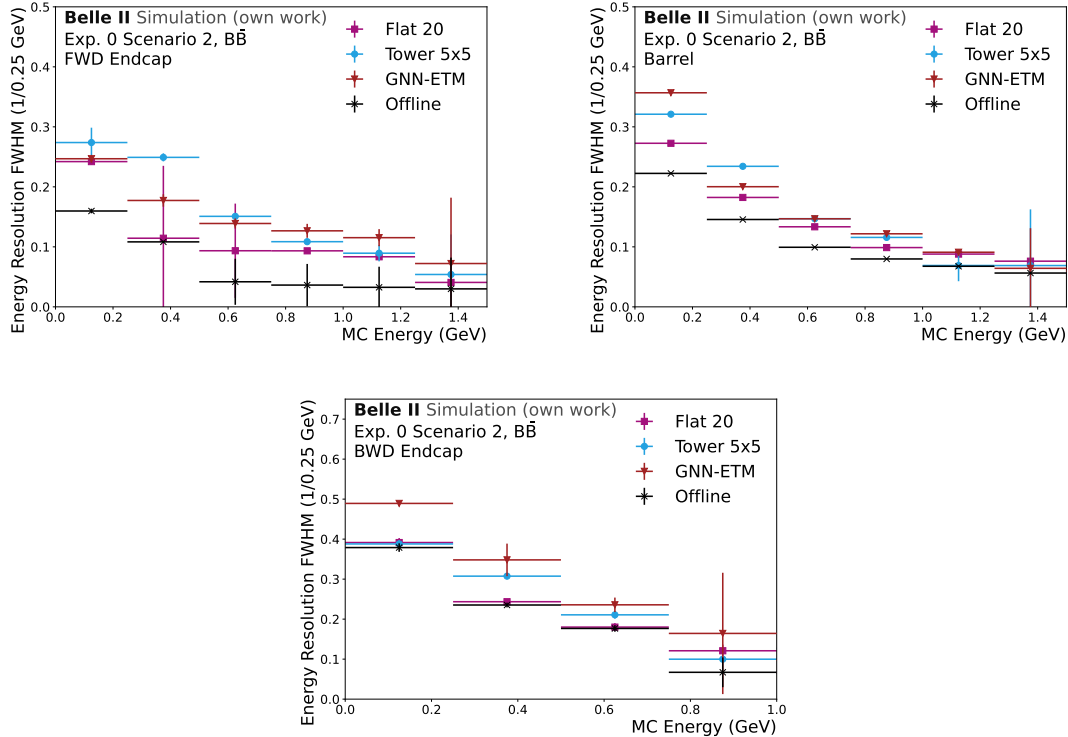


Figure 7.18.: Energy Resolution in the different detector regions. Determined for the high-granularity Tower 5x5 and Flat 20 and the low granularity ICN-ETM and GNN-ETM on the  $B\bar{B}$  sample. The resolution is determined on the subset of target clusters MC-matched to photons, found by all models. As the target energy, the MC-information of the simulated photon is used. For comparison, the performance of the offline reconstruction on the same subset of clusters is shown in black.

### 7.5.2. Quantisation Comparison

When comparing the cluster-finding metrics between the high-granularity models, shown in fig. 7.19, an almost constant efficiency loss can be observed for the Quantised model across all detector regions and energies. The reduced efficiency of the Quantised model is caused mainly by missing predictions. In the endcap, both the Quantised and the Flat 20 model suffer under the single crystals of clusters, which are predominantly outside of the trigger time window, especially at higher energies. This leads to predictions that are positioned

at these crystals, but the actual cluster position is far enough off that the prediction is not matched correctly. This is further accentuated by the reduced position resolution in the endcap under this high occupancy. Compared to the single photon sample, the purity drop of the Tower 5x5 and Quantised model is not visible; this can be explained by the fact that there are proportionally much more target clusters available in these events. Hence, a double prediction or a bad resolution leads more seldom to a reduced purity, as another target cluster can be matched with the prediction.

The position resolution on MC-matched photons, of all models, is reduced in the endcaps, seen in fig. 7.20. Especially in the backward endcap, the behaviour is much noisier. As this is also present in the offline reconstruction, this shows that for these energy regions and the highly occupied backward endcap, mainly caused by the reduced statistics.

The energy resolution is determined on MC-matched photons and shown in fig. 7.21. It exhibits an overall worse performance than for the single photon sample, while still following the same underlying shape. The Flat 20 model performs best, approaching the performance of the offline resolution, followed by the Tower 5x5 and Quantised models.

In the  $B\bar{B}$  sample, the higher-granularity models outperform the lower-granularity models, especially the ICN-ETM struggles to find the target clusters consistently. The high-granularity models still prove to be reliable even for those complicated events and cluster shapes. The Quantised model exhibits the same reduced energy and position resolution as in the other samples. Furthermore, a decrease across all energies and detector regions in the cluster-finding efficiency is observable for the Quantised model.

### 7.5.3. Signal/Background Classification Performance

Additionally to the energy and position prediction, the models predict a binary classifier for each cluster, indicating whether the cluster is caused by a MC-particle or not. As the ICN-ETM does not incorporate such a classifier, only the GNN-based clustering algorithms are compared. For the sake of completeness, the ROC curves for these models in the barrel region of the detector, evaluated for offline clusters with energies between 150-250 MeV, are shown in fig. 7.22. However, no deeper study on the signal classifier is conducted in this thesis; hence, the other detector parts and energy regimes are omitted.

The ROC curves of the different models behave comparably with only small deviations. Therefore, the classifier works as intended, and a reduction of background contributions can be made. Due to the different input reduction methods and model performances, these ROC curves have to be set into perspective with the absolute number of found signal and background clusters and can not be compared easily. For this, the number of found clusters, in the same detector part and energy range, is separated into signal and background clusters and listed in table 7.1. As the found clusters of the Quantised and Flat 20 model are more likely to be background clusters, the performance of the signal classifier has to outperform this inherently larger fraction of background clusters, to result in the overall same ratio of signal to background clusters. Therefore, the performance of the Tower 5x5 model is more promising. However, this requires a more detailed analysis.

To conclude the evaluation, the potential of the single crystal input for high beam background conditions is undeniable. It proves to be more reliable in high background conditions and

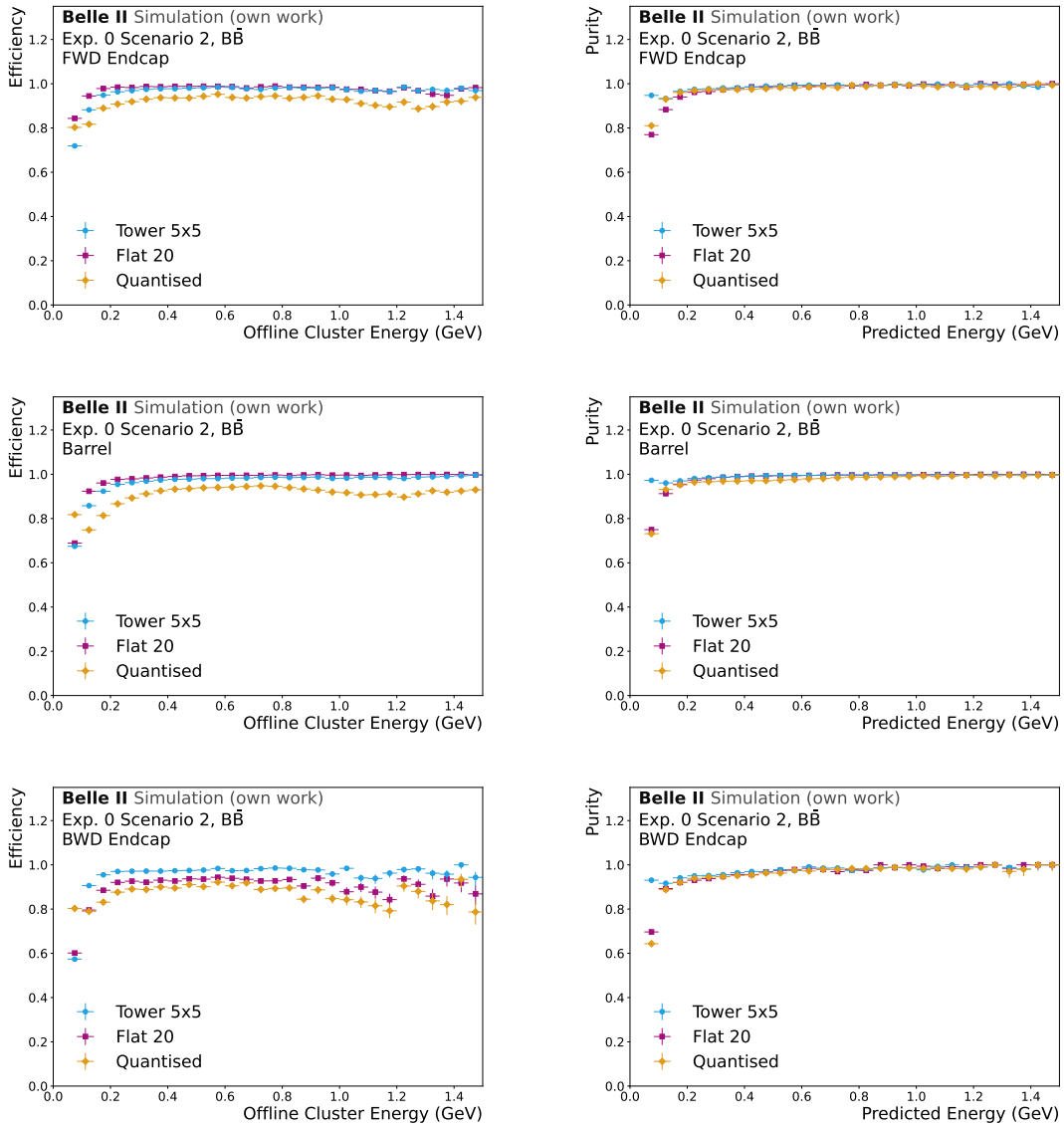


Figure 7.19.: Cluster finding efficiency and purity for the different detector regions in 50 MeV energy bins. Determined for the two non-quantised Tower 5x5 and Flat 20 models and the Quantised model. For each model, the respective targets in the  $B\bar{B}$  sample are used.

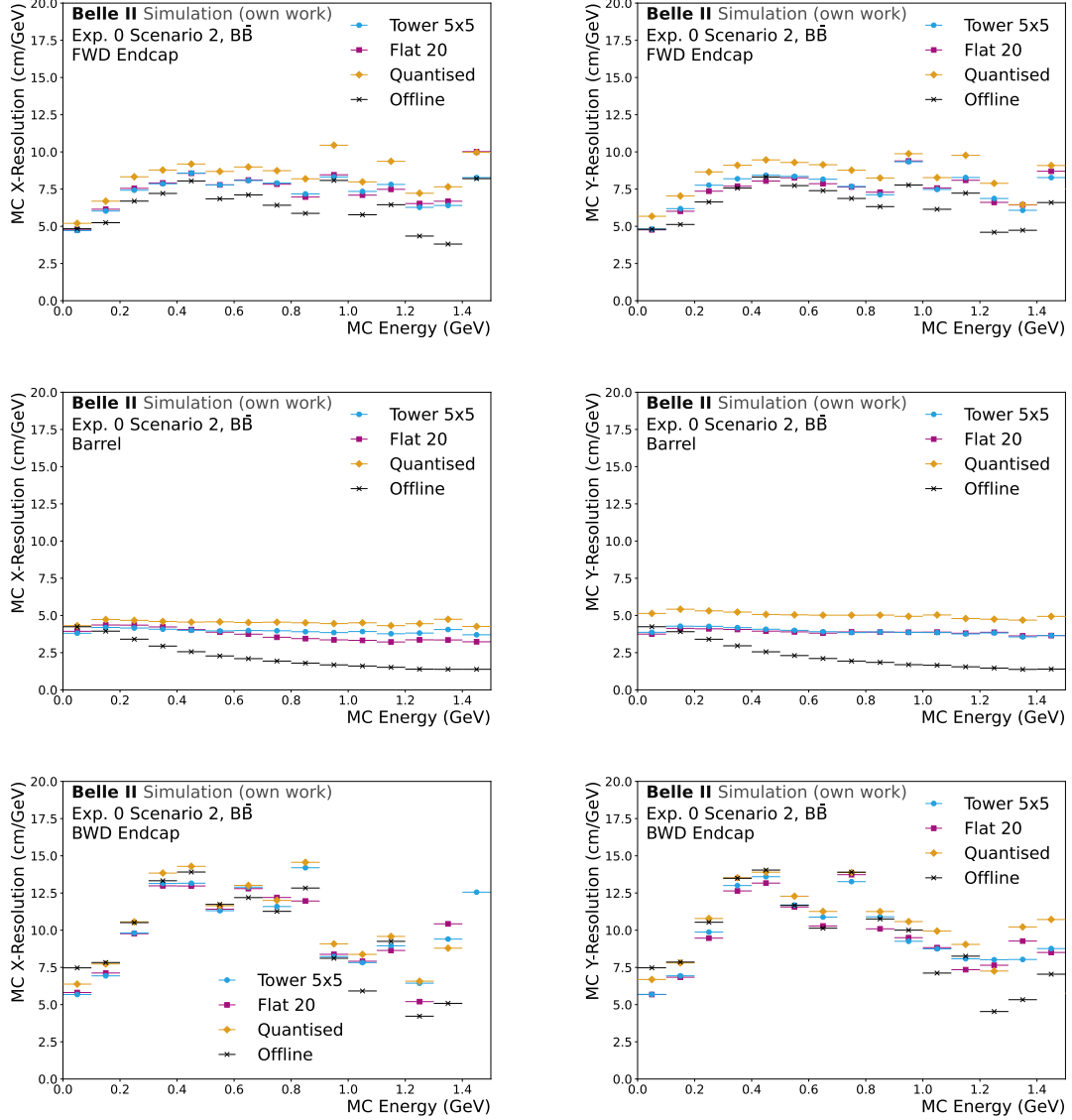


Figure 7.20.: Position Resolution in x and y direction for the backward endcap region. Determined for the two non-quantised Tower 5x5 and Flat 20 models and the Quantised model. The resolution is determined on the subset of target clusters MC-matched to photons, found by all models. As the target position, the extrapolated MC-information of the simulated photon is used. For comparison, the performance of the offline reconstruction on the same subset of clusters is shown in black.

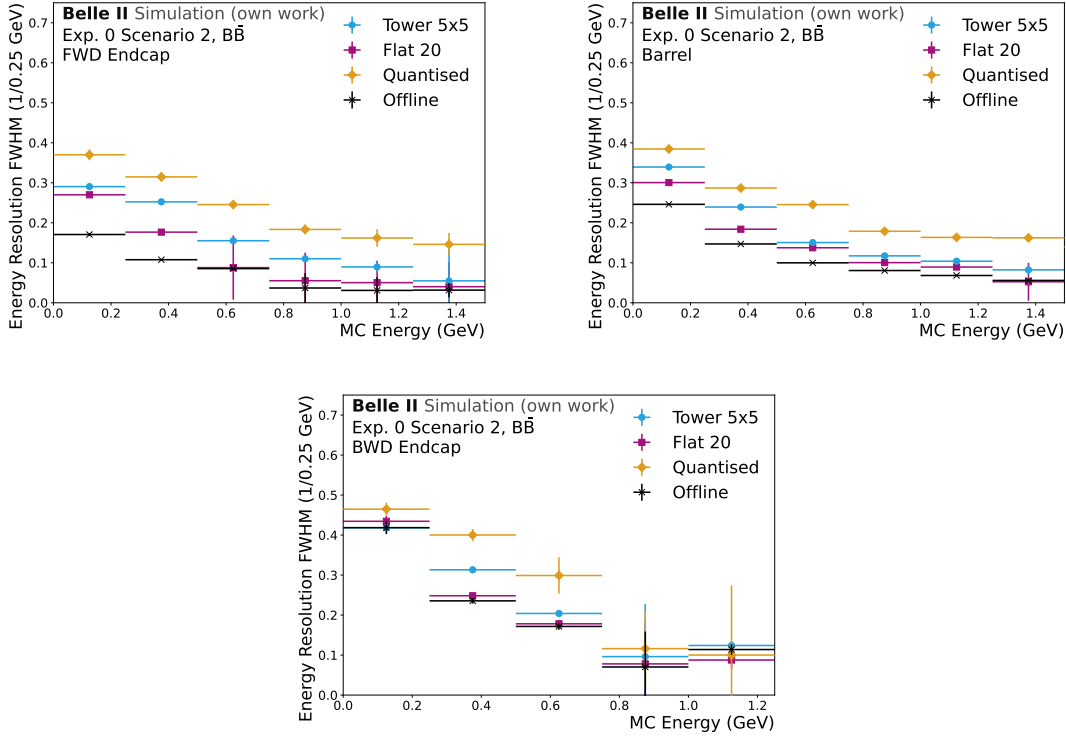


Figure 7.21.: Energy Resolution in the different detector regions. Determined for the two non-quantised Tower 5x5 and Flat 20 models and the Quantised model on the  $B\bar{B}$  sample. The resolution is determined on the subset of target clusters MC-matched to photons, found by all models. As the target energy, the MC-information of the simulated photon is used. For comparison, the performance of the offline reconstruction on the same subset of clusters is shown in black.

Table 7.1.: Total number of found signal and background clusters by the respective model. Determined on all respectively found offline clusters between 150-250 MeV on the  $B\bar{B}$  sample in the barrel region.

	Found Signal	Found Background
GNN-ETM	349136	7328
Tower 5x5	419425	11004
Quantised	355413	11649
Flat 20	423480	14562

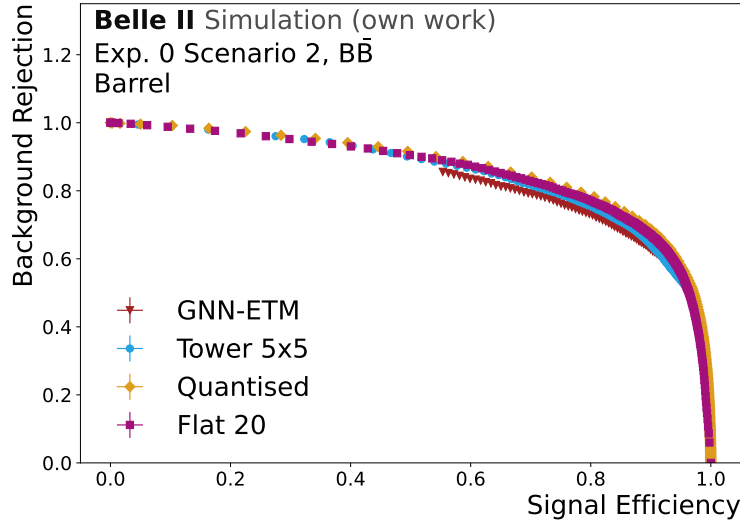


Figure 7.22.: ROC curve for the GNN-based models on the  $B\bar{B}$  sample in the barrel region, determined on all respectively found offline clusters between 150-250 MeV.

increases the position resolution significantly. The quantised implementation exhibits the improved position resolution, but in difficult cases, like overlapping clusters, the reduced message-passing hinders the model from separating the overlapping clusters reliably and reaching the baseline of the non-quantised models. The metric with the largest room for improvement is the energy resolution. Regardless of that, the incorporated signal classifier proves the capability of the GNN-based approach to suppress background contributions, showing the suitability of this approach to be used in the higher background conditions expected after the LS2 upgrade.

## 8. Conclusion

In this thesis, I present a fully implemented adaptation of the current ECL L1 trigger pipeline at Belle II, modified specifically for the LS2 upgrade. This proposal is based on the CaloClusterNet architecture and exploits the input resolution increase of the upgraded readout electronics.

I present two different reduction strategies to limit the input size per trigger window to 128 active crystals: a flat energy cut across the whole detector and a two-stage approach, comprised of determining high-energy trigger towers and subsequently their surrounding crystals as regions of interest. For both strategies, I train and compare the proposed non-quantised model with the currently deployed online, L1 trigger, and subsequent offline clustering algorithms. These comparisons show that the cluster-finding of my approach works for clean as well as highly occupied event signatures with high background conditions. While the current L1 trigger algorithm, the ICN-ETM, breaks down in the extrapolated high beam background conditions, with an up to 40% loss in purity. Both models outperform the currently deployed online algorithms, regarding their position as well as energy resolution. Even so, the performance increase in the energy resolution is less consistent, while still reaching up to 30%, the position resolution is persistently improved by 60 – 70%. The proposed tower cut provides an inherent reduction of background-induced crystals compared to the flat energy cut. The model trained on this tower cut sample has a slightly reduced cluster finding performance compared to the model trained and evaluated on the flat energy cut. Except for the dimuon and  $B\bar{B}$  sample, where the tower model proves to be more robust. The model trained on the flat energy cut exhibits overall an increased x-position and energy resolution. On some samples, reaching an increase in energy resolution of up to 30% over the tower model.

As proof of concept, the proposed quantised model is also implemented on an AMD VCK190 test board, ensuring the feasibility of this study. This board incorporates integrated AI Engines (AIEs), which can efficiently compute large parts of the model. Even so, the quantisation and reduced message passing limit the performance of the final model; both the potential of the higher input granularity as well as the practicality of this approach are proven. The quantised model performs well on finding and predicting isolated clusters. The performance, regarding the position resolution, is slightly decreased, compared to the non-quantised models, but still significantly better than the current ICN-ETM and GNN-ETM.

The overall performance reduction compared to the non-quantised models, in the energy

resolution and the separation ability of overlapping clusters, is mainly caused by the severe bottleneck in information exchange within the GNN. For future developments, the prioritisation of the message passing step, by removing computationally costly steps in the model, could optimise the overall performance and counteract these limitations. A major reduction in hardware complexity can be achieved by shifting to a static graph approach and removing the dynamic graph building, as these parts cannot be accelerated by the AIEs.

This implementation shows that the deployment of large graph neural networks is possible and feasible for real-time trigger environments by exploiting novel hardware designs. With the incorporated background classification, this proves to be a promising way to reduce the impact of the higher beam background conditions on the performance of the L1 trigger. Hence, this is an important development for the future operation of Belle II as well as other high-background collider experiments.

# Bibliography

- [1] Francesco Forti. *Snowmass Whitepaper: The Belle II Detector Upgrade Program*. arXiv:2203.11349 [hep-ex]. Mar. 2022. DOI: 10.48550/arXiv.2203.11349. URL: <http://arxiv.org/abs/2203.11349> (visited on 12/28/2025).
- [2] Youwen Xue, Taichiro Koga, and Junhao Yin. “2 Trigger rate extrapolation using the data on April 16th 2024”. en. In: ().
- [3] ByungGu Cheon. “Design of a Electromagnetic Calorimeter Trigger System for the Belle II Experiment”. en. In: *Journal of the Korean Physical Society* 57.6 (Dec. 2010), pp. 1369–1375. DOI: 10.3938/jkps.57.1369. URL: <https://www.jkps.or.kr/journal/view.html?doi=10.3938/jkps.57.1369> (visited on 02/26/2026).
- [4] Isabel Haide. *A Real-Time Graph Neural Network Trigger Algorithm for the Belle II Electromagnetic Calorimeter*. de. ISBN: 9781000184921. 2025. DOI: 10.5445/IR/1000184927. URL: <https://publikationen.bibliothek.kit.edu/1000184927> (visited on 12/28/2025).
- [5] I. Haide et al. *Real-time graph neural networks on FPGAs for the Belle II electromagnetic calorimeter*. arXiv:2602.15118 [physics]. Feb. 2026. DOI: 10.48550/arXiv.2602.15118. URL: <http://arxiv.org/abs/2602.15118> (visited on 02/21/2026).
- [6] Jan Kieseler. “Object condensation: one-stage grid-free multi-object reconstruction in physics detectors, graph and image data”. In: *The European Physical Journal C* 80.9 (Sept. 2020). arXiv:2002.03605 [physics], p. 886. ISSN: 1434-6044, 1434-6052. DOI: 10.1140/epjc/s10052-020-08461-2. URL: <http://arxiv.org/abs/2002.03605> (visited on 12/29/2025).
- [7] Shah Rukh Qasim et al. “Learning representations of irregular particle-detector geometry with distance-weighted graph networks”. In: *The European Physical Journal C* 79.7 (July 2019). arXiv:1902.07987 [physics], p. 608. ISSN: 1434-6044, 1434-6052. DOI: 10.1140/epjc/s10052-019-7113-9. URL: <http://arxiv.org/abs/1902.07987> (visited on 12/29/2025).
- [8] KEK. *SuperKEKB Photo Gallery*. ja. URL: <https://www.kek.jp/ja/about/pr/image/12623> (visited on 02/21/2026).
- [9] P. Ahlburg et al. “The new and complete Belle II DEPFET pixel detector: Commissioning and previous operational experience”. In: *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* 1068 (Nov. 2024), p. 169763. ISSN: 0168-9002. DOI: 10.1016/j.nima.

- 2024.169763. URL: <https://www.sciencedirect.com/science/article/pii/S0168900224006892> (visited on 02/25/2026).
- [10] K. Ravindran et al. *Operational experience and performance of the Silicon Vertex Detector after the first long shutdown of Belle II*. arXiv:2504.17715 [physics]. Apr. 2025. DOI: 10.48550/arXiv.2504.17715. URL: <http://arxiv.org/abs/2504.17715> (visited on 02/25/2026).
- [11] N. Taniguchi. “Central Drift Chamber for Belle-II”. en. In: *Journal of Instrumentation* 12.06 (June 2017), p. C06014. ISSN: 1748-0221. DOI: 10.1088/1748-0221/12/06/C06014. URL: <https://doi.org/10.1088/1748-0221/12/06/C06014> (visited on 02/25/2026).
- [12] Hulya Atmacan et al. “The Imaging Time-of-Propagation Detector at Belle II”. In: *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* 1080 (Nov. 2025). arXiv:2504.19090 [hep-ex], p. 170627. ISSN: 01689002. DOI: 10.1016/j.nima.2025.170627. URL: <http://arxiv.org/abs/2504.19090> (visited on 02/25/2026).
- [13] Kenta Uno. “Operation and performance of the Belle II Aerogel RICH detector”. In: *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* 1056 (Nov. 2023), p. 168635. ISSN: 0168-9002. DOI: 10.1016/j.nima.2023.168635. URL: <https://www.sciencedirect.com/science/article/pii/S0168900223006253> (visited on 02/25/2026).
- [14] B. Shwartz and on behalf of BELLE II calorimeter group on behalf of. “Electromagnetic calorimeter of the Belle II detector”. en. In: *Journal of Physics: Conference Series* 928.1 (Nov. 2017), p. 012021. ISSN: 1742-6596. DOI: 10.1088/1742-6596/928/1/012021. URL: <https://doi.org/10.1088/1742-6596/928/1/012021> (visited on 02/25/2026).
- [15] C. Ketter et al. *Design and Commissioning of Readout Electronics for a  $4\pi$  and  $4\pi$  Detector at the Belle II Experiment*. arXiv:2502.02724 [hep-ex]. Feb. 2025. DOI: 10.48550/arXiv.2502.02724. URL: <http://arxiv.org/abs/2502.02724> (visited on 02/25/2026).
- [16] T. Abe et al. *Belle II Technical Design Report*. arXiv:1011.0352 [physics]. Nov. 2010. DOI: 10.48550/arXiv.1011.0352. URL: <http://arxiv.org/abs/1011.0352> (visited on 12/28/2025).
- [17] E Kou et al. “The Belle II Physics Book”. In: *Progress of Theoretical and Experimental Physics* 2020.2 (Feb. 2020), p. 029201. ISSN: 2050-3911. DOI: 10.1093/ptep/ptaa008. URL: <https://doi.org/10.1093/ptep/ptaa008> (visited on 01/05/2026).
- [18] *Belle II Analysis Software Framework (basf2)*. URL: <https://zenodo.org/records/14710811> (visited on 01/26/2026).
- [19] T. Kuhr et al. “The Belle II Core Software”. en. In: *Computing and Software for Big Science* 3.1 (Nov. 2018), p. 1. ISSN: 2510-2044. DOI: 10.1007/s41781-018-0017-9. URL: <https://doi.org/10.1007/s41781-018-0017-9> (visited on 01/26/2026).

- [20] S. Longo et al. “CsI(Tl) pulse shape discrimination with the Belle II electromagnetic calorimeter as a novel method to improve particle identification at electron–positron colliders”. In: *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* 982 (Dec. 2020), p. 164562. ISSN: 0168-9002. DOI: 10.1016/j.nima.2020.164562. URL: <https://www.sciencedirect.com/science/article/pii/S0168900220309591> (visited on 01/06/2026).
- [21] A. Natochii et al. “Measured and projected beam backgrounds in the Belle II experiment at the SuperKEKB collider”. In: *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* 1055 (Oct. 2023). arXiv:2302.01566 [hep-ex], p. 168550. ISSN: 01689002. DOI: 10.1016/j.nima.2023.168550. URL: <http://arxiv.org/abs/2302.01566> (visited on 12/28/2025).
- [22] SungHyun Kim et al. “Status of the electromagnetic calorimeter trigger system at Belle II.” en. In: *Journal of Physics: Conference Series* 928.1 (Nov. 2017). Publisher: IOP Publishing, p. 012022. ISSN: 1742-6596. DOI: 10.1088/1742-6596/928/1/012022. URL: <https://doi.org/10.1088/1742-6596/928/1/012022> (visited on 12/30/2025).
- [23] Yoshihito Iwasaki et al. “Level 1 trigger system for the Belle II experiment”. In: *2010 17th IEEE-NPSS Real Time Conference*. May 2010, pp. 1–9. DOI: 10.1109/RTC.2010.5750454. URL: <https://ieeexplore.ieee.org/document/5750454> (visited on 01/24/2026).
- [24] *AMD Versal AI Core Series Adaptive SoCs*. en. URL: <https://www.amd.com/en/products/adaptive-socs-and-fpgas/versal/ai-core-series.html> (visited on 02/28/2026).
- [25] H. Aihara et al. *The Belle II Detector Upgrades Framework Conceptual Design Report*. arXiv:2406.19421 [hep-ex]. July 2024. DOI: 10.48550/arXiv.2406.19421. URL: <http://arxiv.org/abs/2406.19421> (visited on 01/06/2026).
- [26] Claudionor N. Coelho et al. “Automatic heterogeneous quantization of deep neural networks for low-latency inference on the edge for particle detectors”. en. In: *Nature Machine Intelligence* 3.8 (Aug. 2021), pp. 675–686. ISSN: 2522-5839. DOI: 10.1038/s42256-021-00356-5. URL: <https://www.nature.com/articles/s42256-021-00356-5> (visited on 01/29/2026).
- [27] Lea Reuter et al. “End-to-End Multi-Track Reconstruction using Graph Neural Networks at Belle II”. In: *Computing and Software for Big Science* 9.1 (Dec. 2025). arXiv:2411.13596 [physics], p. 6. ISSN: 2510-2036, 2510-2044. DOI: 10.1007/s41781-025-00135-6. URL: <http://arxiv.org/abs/2411.13596> (visited on 02/25/2026).
- [28] Astropy Collaboration et al. “The Astropy Project: Sustaining and Growing a Community-oriented Open-source Project and the Latest Major Release (v5.0) of the Core Package”. In: 935.2, 167 (Aug. 2022), p. 167. DOI: 10.3847/1538-4357/ac7c74. arXiv: 2206.14220 [astro-ph.IM].

- [29] Jonas Eschle et al. “zfit: Scalable pythonic fitting”. In: *SoftwareX* 11 (Jan. 2020), p. 100508. ISSN: 2352-7110. DOI: 10.1016/j.softx.2020.100508. URL: <https://www.sciencedirect.com/science/article/pii/S2352711019303851> (visited on 02/25/2026).
- [30] Michael Zhu and Suyog Gupta. *To prune, or not to prune: exploring the efficacy of pruning for model compression*. arXiv:1710.01878 [stat]. Nov. 2017. DOI: 10.48550/arXiv.1710.01878. URL: <http://arxiv.org/abs/1710.01878> (visited on 02/26/2026).

# A. Appendix

## A.1. Single Photon Sample

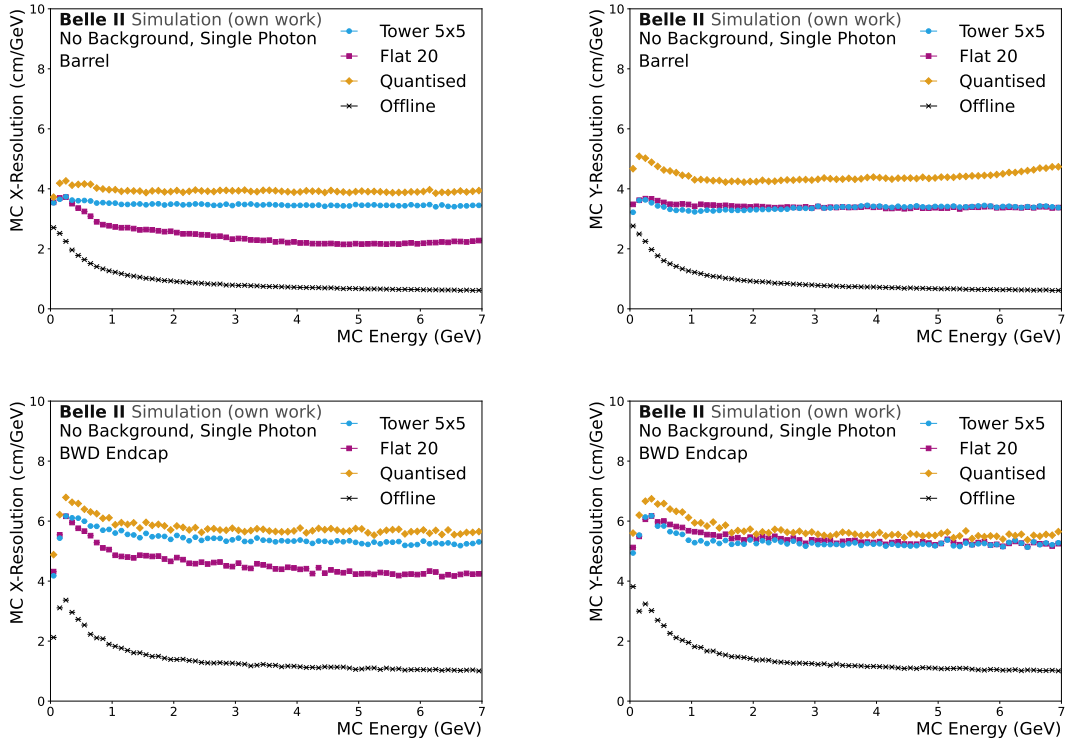


Figure A.1.: Position Resolution in x and y direction for the barrel and backward endcap region. Determined for the two non-quantised Tower 5x5 and Flat 20 models and the Quantised model. A selection of only events with one target cluster is applied. Furthermore, the resolution is determined on the subset of MC-matched target clusters found by all models. As the target position, the extrapolated MC-information of the simulated photon is used. For comparison, the performance of the offline reconstruction on the same subset of clusters is shown in black.



# List of Acronyms

**AIE** AI Engine. 10, 11, 25, 30, 32, 84

**ALP** Axion like particle. 14, 17

**ARICH** Aerogel Ring-Imaging Cherenkov. 4

**basf2** Belle II Analysis Framework. 4, 5, 9, 13, 15–17, 20, 37

**CDC** Central Drift Chamber. 3, 9, 21, 22

**CM** centre of mass. 44, 47

**CR** Connected Region. 5

**ECL** Electromagnetic Calorimeter. 1–4, 8, 9, 11, 12, 21–23, 25, 33, 34, 37, 44, 53, 83

**FPGA** Field Programmable Gate Array. 1, 9, 10, 25, 29, 30, 32, 53, VII

**FTL** Fast Timing Layer. 22

**FWHM** full width at half maximum. 41, 42

**GDL** Global Decision Logic. 9, 12

**GNN** Graph Neural Network. 2, 25, 26, 28, 57, 74, 78, 82, 84

**GRL** Global Reconstruction Logic. 9, 12

**HLT** High Level Trigger. 1, 7, 9, 21

**ICN-ETM** Isolated Cluster Number ECL Trigger Module. 11–14, 25, 28, 38, 40, 53, 57–61, 65, 74, 75, 77, 78, 83

**IR** interaction region. 19, 20

**ITT** Inner Tracking and Timing detector. 22

**KLM**  $K_L^0$  and  $\mu$  detector. 4, 9, 21, 22

- L1 trigger** Level 1 Trigger. 1, 2, 9–12, 14, 19–23, 25, 28, 30, 33, 42–44, 47, 53, 63, 83, 84
- LM** Local Maximum. 5, 7
- LS1** Long Shutdown 1. 19
- LS2** Long Shutdown 2. 1, 2, 13, 19, 20, 22, 25, 30, 33, 34, 82, 83
- LUT** Look-up table. 10
- MC** Monte Carlo simulation. 2, 4, 12–17, 19, 20, 35, 37, 39, 40, 43, 44, 46, 47, 49, 51, 54, 59–61, 63–65, 67–70, 72–74, 76–78, 80, 81, 89
- NoC** Network on Chip. 10
- PID** particle identification. 4
- PXD** Pixel Detector. 3, 22
- ShaperDSP** shaper-digitizer. 22, 23, 25, 30, 53
- SVD** Silicon Vertex Detector. 3, 22
- TC** Trigger Cell. 1, 11–17, 23–28, 33, 34, 38, 40, 43, 47, 57, 70, 74
- TOP** Time-Of-Propagation. 4, 9, 22
- UT** Universal Trigger Board. 10, 25, 28, 30
- VTX** Vertex Detector. 22
- VXD** Vertex Detector. 19

# Glossary

**Belle II** Multi-purpose detector, B factory and the successor of the Belle experiment . 1–3, 5, 9, 19, 20, 22, 23, 47, 83, 84

**GNN-ETM** Graph Neural Network (GNN) based ECL Trigger Master Module (ETM) performing online clustering and reconstruction as described in [4, 5]. 25, 26, 28–30, 38, 53, 54, 57–61, 63, 65, 75–77, 83, VII

**SuperKEKB** Electron-positron collider and accelerator at which the Belle II experiment is located. Successor of KEKB accelerator and collider. . 1, 19, 20

