

---

# ELECTRON ENERGY REGRESSION IN THE CMS HIGH-GRANULARITY CALORIMETER PROTOTYPE

---

Roger Rusack<sup>1†</sup>, Bhargav Joshi<sup>1†</sup>, Alpana Alpana<sup>2†</sup>, Seema Sharma<sup>2†</sup>, Thomas Vadnais<sup>1†</sup>

<sup>1</sup> Department of Physics, University of Minnesota, Minneapolis, USA

<sup>2</sup> Department of Physics, IISER-Pune, India

<sup>†</sup> These authors contributed equally to this work

## ABSTRACT

We present a new publicly available dataset that contains simulated data of a novel calorimeter to be installed at the CERN Large Hadron Collider. This detector will have more than six-million channels with each channel capable of position, ionisation and precision time measurement. Reconstructing these events in an efficient way poses an immense challenge which is being addressed with the latest machine learning techniques. As part of this development a large prototype with 12,000 channels was built and a beam of high-energy electrons incident on it. Using machine learning methods we have reconstructed the energy of incident electrons from the energies of three-dimensional hits, which is known to some precision. By releasing this data publicly we hope to encourage experts in the application of machine learning to develop efficient and accurate image reconstruction of these electrons.

**Keywords** HGCAL · FAIR Data · Energy Regression · Machine Learning · DNN

## 1 Introduction

To measure the energy of particles produced in collisions at the large hadron collider (LHC) the Compact Muon Solenoid (CMS) experimental detector currently has in each of its two endcaps an electromagnetic calorimeter (ECAL), equipped with a preshower (ES) detector, and a hadronic calorimeter (HCAL). Between the interaction point (IP) where the collisions occur there is a silicon tracking detector to measure the momentum of charged particles as they move through the solenoidal magnetic field. Towards the end of this decade the LHC will be upgraded to the High-Luminosity LHC (HL-LHC) where the collision rate of the colliding beams will be increased by a factor of three or more. To cope with the high radiation levels from the particles produced in the collisions the calorimeters in the endcaps will be replaced with a new type of calorimeter, the high-granularity calorimeter (HGCAL), which tracks the progression of the loss of energy by high energy particles by sampling of the shower at different depths inside it. The HGCAL will be constructed from radiation hard silicon sensors, or plastic scintillator sensors, where the radiation levels are lower, that are sandwiched between passive layers of absorber material made of steel or lead. The location within the CMS detector and an outline of the design are shown in Fig. 1.

In the HGCAL there will be approximately three million detector channels in each of the two endcaps. The information of the energy deposited by particles and the time of their arrival in each channel is measured and digitized. This information is transmitted to off-detector electronics for processing and storage. How this information is used to reconstruct the energy of an incident electron, its impact on the calorimeter and its angle of incidence is a challenge that we discuss in this paper. In calorimetry the typical method to reconstruct electrons is with seeding and clustering methods. With the HGCAL design<sup>1</sup>, which has considerably more information available than in earlier examples of calorimeters, new algorithms based on modern machine learning (ML) methods can be developed to solve the reconstruction problem, which in a sense is like a three-dimensional image reconstruction problem. In this paper we discuss the problem of reconstructing high-energy electrons from the energy deposits in the sensors in the HGCAL.

For this we have generated a large volume of simulated data using the GEANT4[1] simulation package, which accurately simulates electromagnetic showers generated by electrons impacting the calorimeter. This data is available at Zenodo<sup>1</sup> and can be used to test new ML methods to address this problem. To accompany the data we provide exemplar software and metadata to permit non-specialist access to the data and development of novel solutions. The exemplar software describes how to access the data and provides a simple reconstruction example that is based on a Deep Neural Network (DNN). In this paper we describe the problem to be solved in more detail and the results that we have obtained with the DNN model.

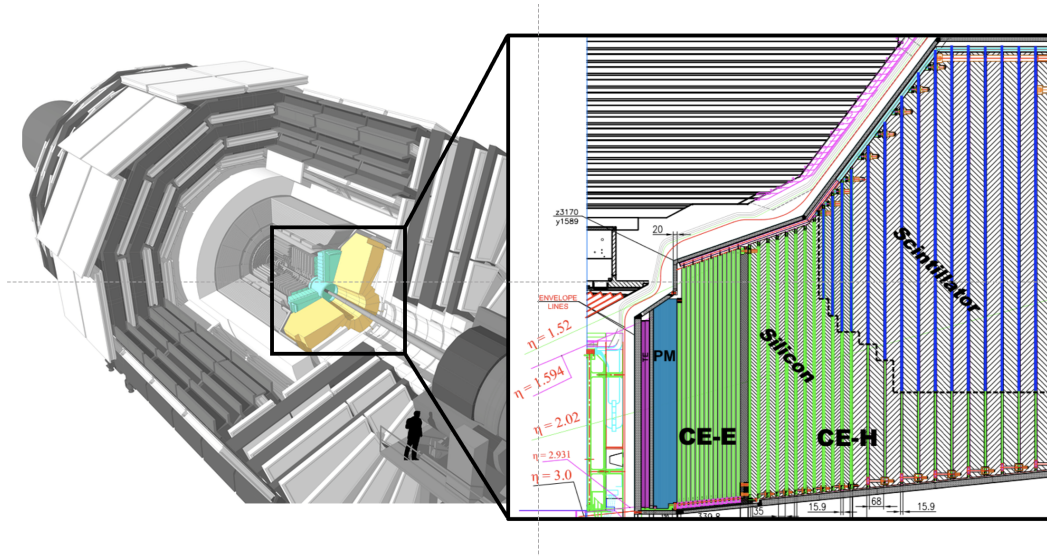


Figure 1: Current design of the CMS detector (left) to the human scale. The highlighted regions in blue and yellow color represent the ECAL and the HCAL detectors. These regions will be replaced by the newly designed calorimeter (right). It consists of three successive layers which combine the functionalities of both, the ECAL and the HCAL.

## 2 The High Granularity Calorimeter

The entire assembly of each of the two HGCAL calorimeters weighs approximately 230 T and will be used to measure the energies of particles produced at the IP with angles of approximately 10 to 30 degrees from the beam axis<sup>2</sup>. In the final detector the first 26 layers will form the electromagnetic (CE-E) [2] section which will have hexagonal silicon sensors of about 8" width divided into hexagonal cells with areas of 1.1 and 0.5 cm<sup>2</sup>. Behind the CE is the 21-layer hadronic section (CH). In this the first eight layers will consist of silicon sensors similar to the CE-E section, and the last 12 layers will have a mixture of silicon sensors and plastic scintillators.

### 2.1 The Prototype Setup

To evaluate the performance of the detector and to qualify many aspects of the design a large-scale prototype of the HGCAL was built and tested in the H2 beamline at CERN's Prévessin site (Figure 2). A beam of positrons is provided by Super Proton Synchrotron (SPS) accelerator. Since, the positron is an anti-particle of the electron differing only in electric charge, the response of the interaction of positron in the prototype is same as that of an electron without any external magnetic field. The prototype consisted of 3 sections, Electromagnetic (CE-E), Hadronic (CE-H) and a CALICE Analog Hadronic Calorimeter (AHCAL)[3, 4], arranged in series in that order. This is similar to the final configuration of the HGCAL. The CE-E [5] section consists of 28 sampling layers made using 14 double-sided mini-cassettes (Figure 3 right). Each cassette consist of a lead clad with stainless steel or Cu/CuW absorber sandwiched between two silicon sensor layers. The hexagonal silicon sensors are subdivided into 128 hexagonal silicon detector channels. Each channel is equipped with electronics to measure the energy and the time of the particle interactions in the sensor. The entire CE-E section corresponds to a total of 26 radiation lengths or 1.4 nuclear interaction lengths.

<sup>1</sup><https://zenodo.org/>

<sup>2</sup>The coverage is between 1.5 and 3.0 in pseudorapidity defined as  $\eta = -\ln|\tan \frac{\theta}{2}|$ , where  $\theta$  is the azimuthal angle relative to the beam axis.

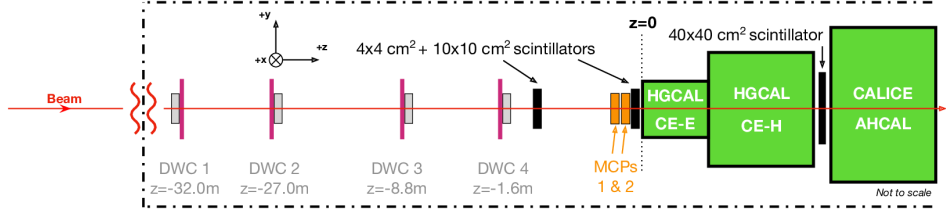


Figure 2: The test beam setup of the prototype along the H2 beam line. The four delay wire chambers (DWCs) track the position of the incoming positron. For triggering on signal events two plastic scintillators and fast multiplying tubes are used.

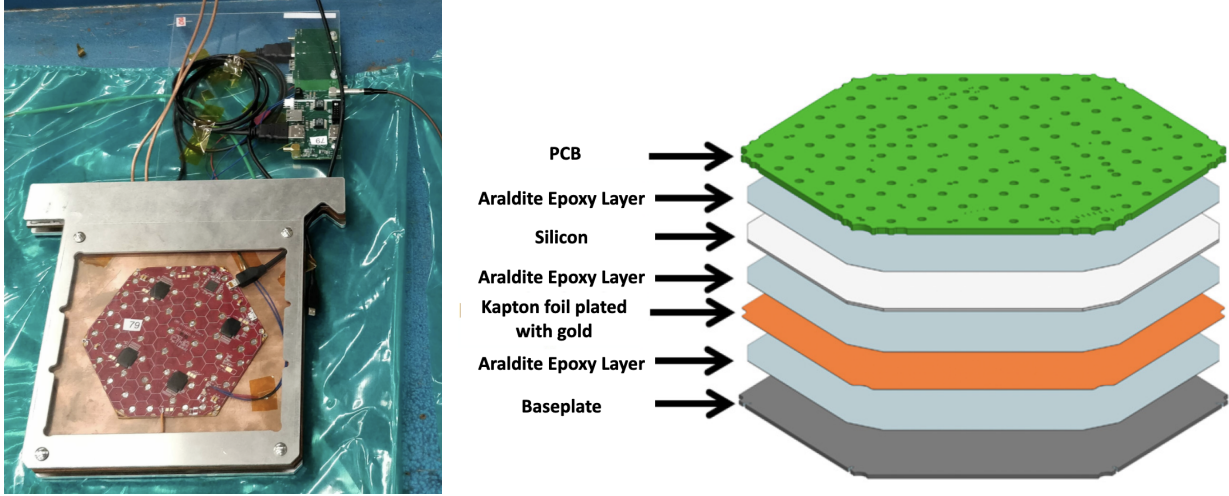


Figure 3: A front view of a prototype of the CE-E minicassette (left). It consists of two Hexagonal module mounted onto a Cu cooling plate on either side. The module is an assembly (right) of a baseplate made of copper or copper-tungsten, a  $100 \mu\text{m}$  thick gold-plated Kapton<sup>®</sup> sheet, a hexagonal silicon sensor, and a printed circuit board called 'hexaboard'. Araldite<sup>®</sup> is used as an epoxy to glue different components in the module.

In the prototype the CE-H [6] section was composed of 12 sampling layers each with seven Si modules arranged in daisy structure, each layer was sandwiched between a 40 mm thick steel plate. Due to the limited availability of silicon sensor modules, the last three layers of CE-H were equipped with only one sensor module placed at the center of the layer. The CE-H is followed by a 4.4 nuclear interaction length deep prototype of the AHCAL that was built with 39 sampling layers of SiPM-on-scintillator-tile active layers interspersed between steel absorbers.

### 3 Electromagnetic Showers

When energetic particles pass through a media, they typically lose energy through coulomb interactions with the electrons in the media. Energetic electrons ( $E \gg 1 \text{ GeV}$ ), on the other hand lose energy primarily via emission of *bremsstrahlung* radiations. When the electron passes through a dense media, it gets accelerated or decelerated quickly due to the strong electric fields of the nuclei which causes it to emit radiations or photons. Energetic photons, on the other hand, produce pairs of electrons and positrons as they interact with the nuclei of the atoms. This results in a cascade of secondary particles known as an Electromagnetic Shower<sup>4</sup> and the process continues until the energy of the decay products falls below a critical energy  $E_c$ .

These showers can be characterized by several parameters, which include the *radiation length* and *Molière radius*. The *radiation length* is defined as the distance over which the energetic electron loses  $1/e$  fraction of its energy. Thus, the "shower depth" can be written in terms of the *radiation length* as follows

$$X = X_0 \frac{\ln(E/E_c)}{\ln(2)} \quad (1)$$

where  $E_c$  is the critical energy<sup>3</sup> of electron in a given material.

As the electron dissipates energy, the size of the spread increases in directions orthogonal to its momentum. The *Molière* radius can be used to define the lateral spread of the shower till the critical energy reached as the electron traverses  $X_0$  through the medium. By definition, a cylinder of *Molière radius* contains about 90% of the total deposited energy.

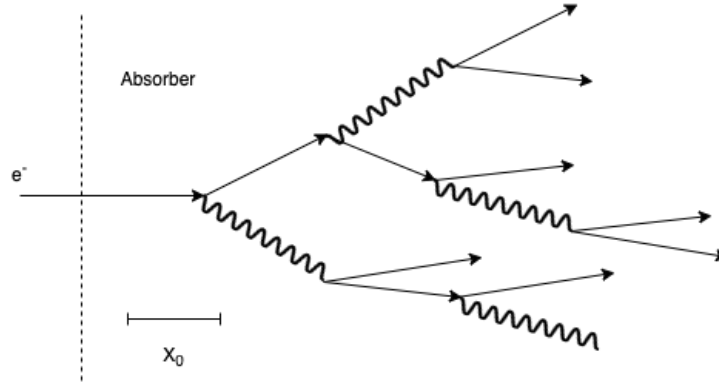


Figure 4: A schematic showing the development of an electromagnetic shower by an incoming electron in an absorber.

The electromagnetic calorimeters are designed to capture the highly energetic photons and electrons and measure their energies. They can also localise the position of the incoming particle in space and, in some cases, measure its direction. The part of the calorimeter that produces showers is known as the absorber material, whereas, the material that measures the energy is known as the active part. Ideally a calorimeter has a small  $X_0$  and *Molière* radius to contain the showers as effectively as possible. The electromagnetic calorimeters can either be of homogeneous type or of sampling type. Homogeneous calorimeters typically have one block of absorber, where the incoming particle dissipates energy and the active material surrounding it measures the energy. In a sampling calorimeter, there are alternating layers of absorbers and active materials, and the energy dissipated in one layer is measured using the energy deposited in the layers before and after the absorber. Finally, the sum of energies over all the layers gives the total energy deposited can be used to measure the energy of the incoming particle.

The energy resolution of a calorimeter gives its precision in measuring the energy. For an electromagnetic calorimeter, the energy resolution can be written as follows.

$$\frac{\sigma}{E} = \frac{S}{\sqrt{E}} \oplus \frac{N}{E} \oplus C, \quad (2)$$

where the first term on the right-hand side is the *stochastic* or *sampling* term, the middle term is the *noise* term and the last term is the *constant* term. The *stochastic* term arises from the fact that the number of primary and secondary particles produced in the interactions fluctuates. The *noise* term, on the other hand, comes from the noise in the detector electronics. Furthermore, this term receives contributions from other simultaneous interactions or collisions happening in the same event known as "pileup". Finally, the constant term is the measure of quality of the detector construction. It accounts for the imperfections in the geometry, non-uniformity in the response and energy losses that cannot be measured by its electronics.

## 4 Dataset

The dataset consists of simulations of reconstructed hits, known as "rechits", produced by positrons passing through the HGCal test beam prototype. For simulations, Monte Carlo method is used to produce positrons with energy ranging from 10 to 350 GeV. In the next step, GEANT4 [1] package is used to simulate their interactions with the detector material. The conditions used in generating positrons are fine tuned to account for real detector effects such as energy losses in the beam. The simulated hits are then digitized using the CMS software. The digitized information was then processed through the CMS software to reconstruct the signals as hits within the detector. The rechits along with their

<sup>3</sup>[https://pdg.lbl.gov/2022/AtomicNuclearProperties/critical\\_energy.html](https://pdg.lbl.gov/2022/AtomicNuclearProperties/critical_energy.html)

details pertaining to signal reconstruction was stored in root [7] format. These files were then skimmed using uproot [8] package to obtain the final dataset. A set of preselections is applied to ensure that the event selection is identical to the one used in performing the analysis [5] published by the CMS collaboration. The hits are chosen to have a minimum energy of  $0.5 \text{ MIP}^4$ , which is well above the HGICAL noise levels. Events with more than 50 hits in CE-H layers are rejected. The track of electron extrapolated using the hits from the DWC chambers is required to be within a  $2 \times 2 \text{ cm}^2$  window within the first layer. The final dataset is a set of 3.2 million events, each event containing position coordinates of rechits within the detector and their calibrated energies. HDF5 format is used to organize the data in hierarchical arrays. The file contains following the arrays:

- **nhits**: An integer array representing number of reconstructed hits (rechits) in each event.
- **rechit\_x**: A nested array of length equal to the number of events and sub-arrays of length of nhits. Each sub-array contains a floating value representing x-coordinate of the position of the rechits in units of centimeters.
- **rechit\_y**: A nested array with a structure and size same as rechit\_x. Each floating value represents the y-coordinate of the position of a rechit in units of centimeters.
- **rechit\_z**: A nested array with a structure and size same as rechit\_x. Each floating value represents the z-coordinate of the position of a rechit in units of centimeters.
- **rechit\_energy**: A nested array with a structure and size same as rechit\_x. Each floating value represents the calibrated energy of a rechit in units of MIPs.
- **target**: The true energy of the incoming positron in units of GeV.

To ensure the FAIR-ness of the publication of the dataset, it has been published [9] on Zenodo [10] platform, which was launched in May 2013 as part of the OpenAIRE project, in partnership with CERN. The dataset[9] consists of two files in *gzip* format. These can be uncompressed to obtain two files in HDF5 format. The smaller sample of 648,000 events with a label "0001" has a file size of 2.8 GB and the full dataset with a label "large" has a file size of 14.0 GB. The code to unpack and use the dataset has been made available on Github<sup>5</sup>. The metadata describing the contents of the file are available in JSON format under the same repository.

## 5 Summary

The purpose of the release of the dataset is to make it open for everyone for building models for estimating the resolution with better precision, develop visualization tools and benchmarking ML techniques such as Generative Adversarial Networks (GANs), which can be used for generating EM showers with reduced computational time. For the purpose of exploring the dataset, the source code of the simple DNN model that was developed in python for energy regression has been added to the aforementioned Github repository. The repository has been built using the "cookiecutter" template used by the FAIR4HEP group for ensuring Findability and reproducibility of the results. An example notebook in the repository also demonstrates a way to make event displays (Figure 5) of individual events in the dataset.

After training on the simulated dataset using a fully connected DNN, the performance of the network can be evaluated by computing the energy resolution in different bins of energies. To achieve this, the difference between measured and true energies from the simulations are plotted for energies ranging from 20 to 300 GeV in 14 bins of 25 GeV width. In each bin, the resulting distribution has a shape of a Gaussian distribution. This distribution is then fit using a  $\chi^2$  minimization technique to obtain the mean and the variance. The mean represents the bias in the estimation in each bin, whereas the ratio of the variance to the mean gives the estimate of the energy resolution. Without any contributions from pileup, the *noise* term in (Equation 1) is assumed to be zero. The squares of the resolutions obtained from the 14 energy bins can be fitted as the sum of quadratures of the *stochastic* term and the constant term. The slope and the intercept of the linear fit (Figure 6) provides an estimate for the *stochastic* term and the constant term respectively.

## 6 Acknowledgements

This work has been supported by the Department of Energy, Office of Science, Office of Advanced Scientific Computing under award number DE-SC0021395. The authors would like to express their gratitude to the CMS Collaboration, and in particular to the CMS HGICAL community for making the providing the configurations files to generate simulated events. We would also like to thank our colleagues from the FAIR4HEP group for discussions and their invaluable inputs and suggestions for writing this paper.

<sup>4</sup>Minimum Ionizing Particle (MIP) is the unit used to count the energy of digitized hits.

<sup>5</sup><https://github.com/FAIR-UMN/FAIR-UMN-HGICAL>

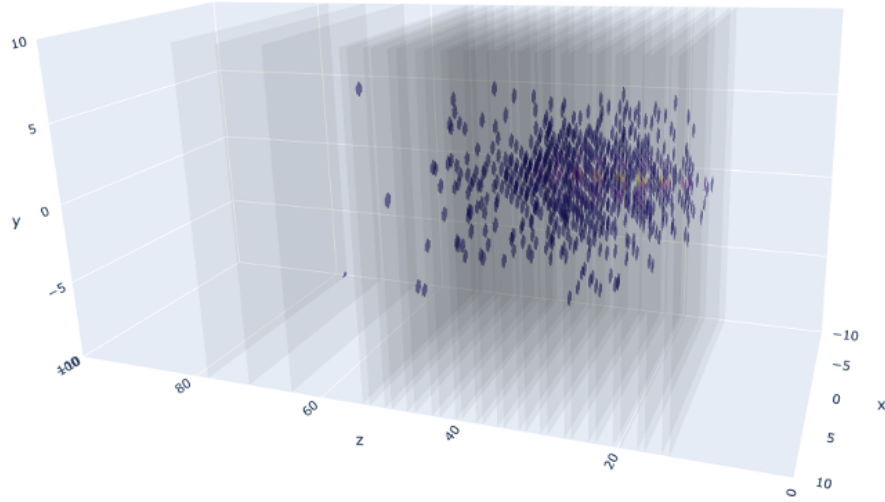


Figure 5: An event display of a simulated event of a 100 GeV positron passing through the prototype. The energy of reconstructed hits is measured in the units of Minimum Ionizing Particles (MIPs).

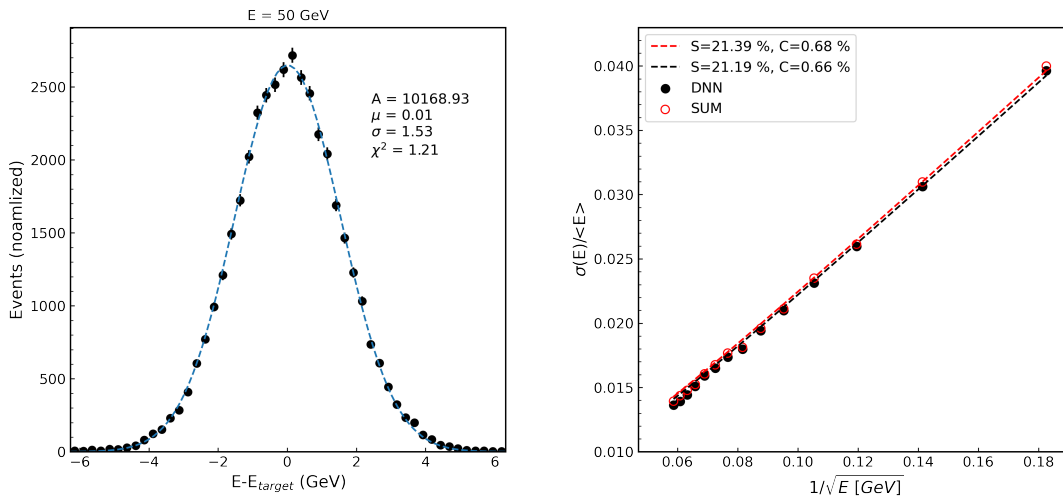


Figure 6: Predicted energies of positrons between a range of [40, 60] GeV particle as predicted by the DNN relative to the energy of the incoming particle (left). The data points are fitted with a Gaussian distribution using a minimum  $\chi^2$  fit. Resolution plotted as a function of the inverse of square root of the energy of the simulated particle (right). The resolutions obtained through summing the energies of rechits (black line) and those obtained through the output of DNN (red line) are comparable.

## References

- [1] S. Agostinelli et al. “Geant4—a simulation toolkit”. *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* 506.3 (2003), pp. 250–303. ISSN: 0168-9002. DOI: [https://doi.org/10.1016/S0168-9002\(03\)01368-8](https://doi.org/10.1016/S0168-9002(03)01368-8). URL: <https://www.sciencedirect.com/science/article/pii/S0168900203013688>.
- [2] The CMS HGCAL collaboration and et al. “Construction and commissioning of CMS CE prototype silicon modules”. *Journal of Instrumentation* 16.04 (Apr. 2021), T04002. DOI: 10.1088/1748-0221/16/04/T04002. URL: <https://dx.doi.org/10.1088/1748-0221/16/04/T04002>.

- [3] CALICE Collaboration. *Design, Construction and Commissioning of a Technological Prototype of a Highly Granular SiPM-on-tile Scintillator-Steel Hadronic Calorimeter*. 2022. arXiv: 2209.15327 [physics.ins-det].
- [4] The CALICE collaboration et al. “Construction and commissioning of the CALICE analog hadron calorimeter prototype”. *Journal of Instrumentation* 5.05 (May 2010), P05004. DOI: 10.1088/1748-0221/5/05/P05004. URL: <https://dx.doi.org/10.1088/1748-0221/5/05/P05004>.
- [5] CMS HGCAL collaboration and et al. “Response of a CMS HGCAL silicon-pad electromagnetic calorimeter prototype to 20–300 GeV positrons”. *Journal of Instrumentation* 17.05 (May 2022), P05022. DOI: 10.1088/1748-0221/17/05/P05022. URL: <https://dx.doi.org/10.1088/1748-0221/17/05/P05022>.
- [6] CMS Collaboration and CALICE Collaboration. “Performance of the CMS High Granularity Calorimeter prototype to charged pion beams of 20–300 GeV/c” (2022). DOI: 10.48550/ARXIV.2211.04740. URL: <https://arxiv.org/abs/2211.04740>.
- [7] I. Antcheva et al. “ROOT — A C++ framework for petabyte data storage, statistical analysis and visualization”. *Computer Physics Communications* 180.12 (2009). 40 YEARS OF CPC: A celebratory issue focused on quality software for high performance, grid and novel computing architectures, pp. 2499–2512. ISSN: 0010-4655. DOI: <https://doi.org/10.1016/j.cpc.2009.08.005>. URL: <https://www.sciencedirect.com/science/article/pii/S0010465509002550>.
- [8] Jim Pivarski et al. *scikit-hep/uproot3: 3.14.4*. Version 3.14.4. Feb. 2021. DOI: 10.5281/zenodo.4537826. URL: <https://doi.org/10.5281/zenodo.4537826>.
- [9] Bhargav Joshi and Alpina Sirohi. *Electron Energy Regression in High-Granularity Calorimeter Prototype*. Version v1. Zenodo, Jan. 2023. DOI: 10.5281/zenodo.7504164. URL: <https://doi.org/10.5281/zenodo.7504164>.
- [10] European Organization For Nuclear Research and OpenAIRE. *Zenodo*. en. 2013. DOI: 10.25495/7GXX-RD71. URL: <https://www.zenodo.org/>.