

# Heterogeneous computing for the local reconstruction algorithms of the CMS calorimeters

Andrea Massironi<sup>1,2</sup>, Viktor Khristenko<sup>2</sup> and Mariarosaria D'Alfonso<sup>3</sup>

<sup>1</sup>INFN Milano-Bicocca, Edificio U2, Piazza della Scienza 3, 20126 Milano, Italy

<sup>2</sup>CERN, Espl. des Particules 1, 1211 Meyrin, Switzerland

<sup>3</sup>MIT, 77 Massachusetts Ave, Cambridge, MA 02139, USA

E-mail: [andrea.massironi@cern.ch](mailto:andrea.massironi@cern.ch)

on behalf of the CMS collaboration

**Abstract.** The increasing LHC luminosity in Run III and, consequently, the increased number of simultaneous proton-proton collisions (pile-up) pose significant challenges for the CMS experiment. These challenges will affect not only the data taking conditions, but also the data processing environment of CMS, which requires an improvement in the online triggering system to match the required detector performance. In order to mitigate the increasing collision rates and complexity of a single event, various approaches are being investigated. Heterogeneous computing resources, recently becoming prominent and abundant, may be significantly better performing for certain types of workflows. In this work, we investigate implementations of common algorithms targeting heterogeneous platforms, such as GPUs and FPGAs. The local reconstruction algorithms of the CMS calorimeters, given their granularity and intrinsic parallelizability, are among the first candidates considered for implementation in such heterogeneous platforms. We will present the current development status and preliminary performance results. Challenges and various obstacles related to each platform, together with the integration into CMS experiments framework, will be further discussed.

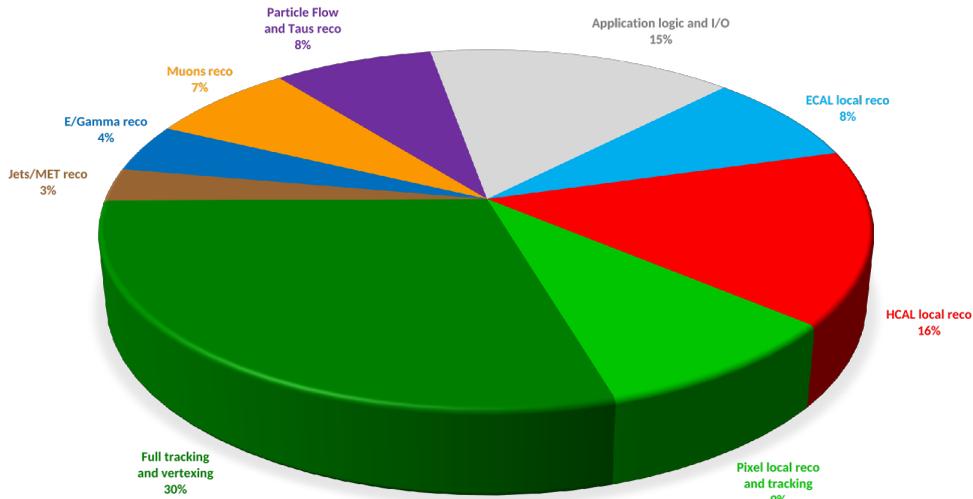
## 1. Introduction

The Compact Muon Solenoid (CMS) experiment [1] is a general purpose experiment at the Large Hadron Collider (LHC) at CERN, designed to search for the Standard Model (SM) Higgs boson and for new physics beyond the SM. A key ingredient of the detector is the calorimetry system, made of the electromagnetic calorimeter (ECAL) [2] and the hadronic calorimeter (HCAL) [3]. Although based on different detector and readout technologies, both parts have to face similar aspects concerning their local reconstruction of deposited energy in each cell. The ECAL is a homogeneous and hermetic calorimeter containing 61200 lead tungstate ( $\text{PbWO}_4$ ) scintillating crystals mounted in the barrel (EB), closed at each end by endcaps (EE) each containing 7324 crystals. The scintillation light is detected by avalanche photodiodes (APDs) in EB and by vacuum phototriodes (VPTs) in EE. The HCAL is composed of a brass and scintillator calorimeter with a central barrel (HB), and two endcaps (HE), steel and quartz fibre forward calorimeters, and an additional outer calorimeter in the barrel region. Both the HB and HE detectors consist of alternating layers of brass and plastic scintillator, whose scintillation light is collected by wavelength-shifting fibre and transmitted to the front-end readout electronics. The

total number of channels is 2592 in HB and 6768 in HE, with an increase for Run III to 9072 for HB.

The non-negligible number of channels and the complexity of the local reconstruction of the energy deposit per channel make the computational time consumed by the calorimeters non negligible, and a possible place of improvement of the High Level Trigger (HLT) time, see Figure 1.

In the following sections new prospects to improve this last point, keeping the current accuracy in the reconstructed energy deposit are reported.



**Figure 1.** Pie chart of time consumption by the HLT. The calorimeter local reconstruction consumes about 25% of the time.

## 2. Local reconstruction

Both ECAL and HCAL sample their signals at 40 MHz, meaning one sampled point every 25 ns. The LHC collision rate is the same, and there is the possibility for every bunch crossing (BX) to have an energy deposit in the calorimeters. The length of a pulse from an energy deposit extends over several samples and thus the contribution from several BXs (out-of-time pile-up) can influence the estimation of the nominal BX energy deposit. An out-of-time pile-up aware local reconstruction has been developed and it fits possible energy deposits both in the nominal BX and in the neighbour ones. The reconstruction algorithm relies on a Fast Non-Negative Least Square (FNNLS) minimization [4, 5], exploiting the non-negativity constraint of energy deposits. The algorithm is based on the minimization of the following  $\chi^2$ :

$$\chi^2 = \left( \sum_{j=1}^{N_{\text{pulse}}} A_j \vec{p}_j - \vec{S} \right)^T \mathbf{C}^{-1} \left( \sum_{j=1}^{N_{\text{pulse}}} A_j \vec{p}_j - \vec{S} \right) \quad (1)$$

where:

- $A_j$  is the target at the  $j$ -th BX, the reconstructed energy deposit for each BX,
- $\vec{p}_j$  is the nominal pulse shifted in time by  $j \times 25$  ns,
- $\vec{S}$  is the vector of the sampled signal from the detector (10 for ECAL and 8 for HCAL),

- $\mathbf{C}$  is the total covariance matrix, defined in Equation 2,

$$\mathbf{C} = \mathbf{C}_{\text{noise}} \oplus \sum_{j=1}^{N_{\text{pulse}}} A_j^2 \mathbf{C}_{\text{pulse}}^j \quad (2)$$

with  $\mathbf{C}_{\text{noise}}$  being the covariance of the noise, while  $\mathbf{C}_{\text{pulse}}$  is the covariance related to the uncertainty in the pulse shape itself,  $\vec{p}$ , mainly coming from time of flight variation due to beam spot width. The FNNLS is based on Cholesky decomposition given the symmetry of the matrix  $\mathbf{C}$ .

### 3. Expose parallelism in the code

In order to make the local reconstruction of the calorimeters faster, a massive parallelism approach has been followed, with two main guidelines during the development: accuracy and performance. Although major changes are needed to expose parallelism and then exploit new architectures, such as GPUs and FPGAs, the first requirement is to keep, within the required accuracy, the output of the algorithm the same. This basic requirement makes sure that all the workflow of the ECAL and HCAL reconstruction remains unchanged, such as clustering, corrections, regressions and object identifications and selections. The second step is the improvement in terms of throughput: faster algorithms mean more data being handled in the same amount of time. Although final results are not ready, first studies show encouraging results in terms of performance.

The code has been rewritten in order to expose as much as possible the parallelism. A first test was based on a single monolithic kernel, that, although achieving the accuracy requirement of the algorithm, it did not expose enough parallelism to obtain a big gain in terms of time performance. In addition to the simple parallelization based on several channels in ECAL and HCAL, further tests have been performed in order to parallelize the per-channel definition of the  $\chi^2$ , and the preparation of the required inputs. For example, the sampled signal  $\vec{S}$  needs to undergo a baseline subtraction process, that can be parallelized given its dimension, 10 for ECAL and 8 for HCAL. At the same time the matrix  $\mathbf{C}$  of Equation 1 can be precomputed exposing the different entries of the matrix to parallel computing.

Major re-writing of the code has been performed, while the most updated developments are still work in progress. The code has been written in order to be compiled with CUDA [6], and exploits the vast availability of third party libraries, such as Eigen [7], that has been adapted to be used with CUDA.

### 4. Summary

The CMS calorimeters local reconstruction will face new challenges due to the increased pile-up in LHC Run III. In order to optimize performance in terms of data throughput, a heterogeneous approach has been followed, exploiting new technologies, such as GPUs. Preservation of the accuracy in the re-written algorithm has been achieved and hints of improvements have been observed. Further work is needed to fully exploit the available parallelism. This work is part of a global effort of CMS in heterogeneous computing, addressing not only calorimeter local reconstruction, but also tracking and new detectors with a huge number of channels and complexity.

### References

- [1] S. Chatrchyan *et al.* CMS Collaboration, “The CMS experiment at the CERN LHC,” JINST 3, S08004 (2008).
- [2] CMS Collaboration, “CMS: The electromagnetic calorimeter. Technical design report,” CERN-LHCC-97-33, CMS-TDR-4.
- [3] CMS Collaboration, “The CMS hadron calorimeter project: Technical Design Report,” CERN-LHCC-97-031, CMS-TDR-2.

- [4] Jason Cantarella and Michael Piatek, “Tsmnls: A solver for large sparse least squares problems with non-negative variables,” <http://arxiv.org/abs/cs.MS/0408029> journal = CoRR,
- [5] Jay Lawhorn and CMS HCAL Collaboration, “New method of out-of-time energy subtraction for the CMS hadronic calorimeter”, doi 10.1088/1742-6596/1162/1/012036, IOP Publishing, 2019.jan.1162
- [6] <https://docs.nvidia.com/cuda/>
- [7] Eigen - a C++ template library for linear algebra: matrices, vectors, numerical solvers, and related algorithms. <http://eigen.tuxfamily.org>.