

# *Heterogeneous techniques for rescaling energy deposits in the CMS Phase-2 endcap calorimeter*

**Bruno Alves**

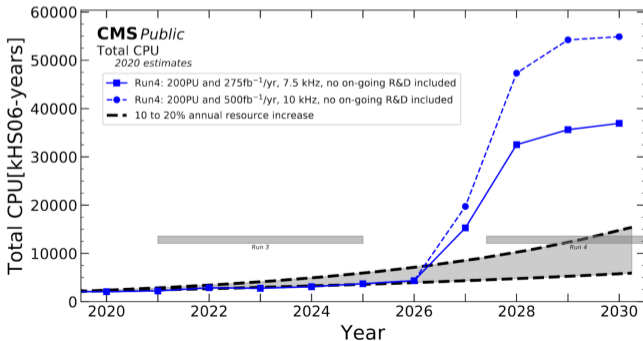
on behalf of the CMS Collaboration

**Tuesday 18 May, 2021**

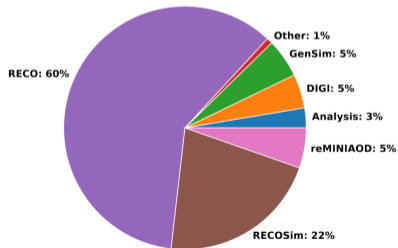


# Motivation

- High Luminosity LHC will start its operation in 2027:
  - **high luminosity** ( $\sim 5 \times 10^{34} \text{ cm}^2 \text{ s}^{-1}$ )
  - **large pileup** (up to 200 collisions per bunch crossing)
- CMS Projections: **significant gap** between future CPU needs and availability
- The biggest contributor to CPU usage is event reconstruction ( $\sim 6\%$  by HGCAL)



**CMS Public**  
Total CPU HL-LHC fractions  
2020 estimates



# High-Granularity CALorimeter (HGCal) @ CMS

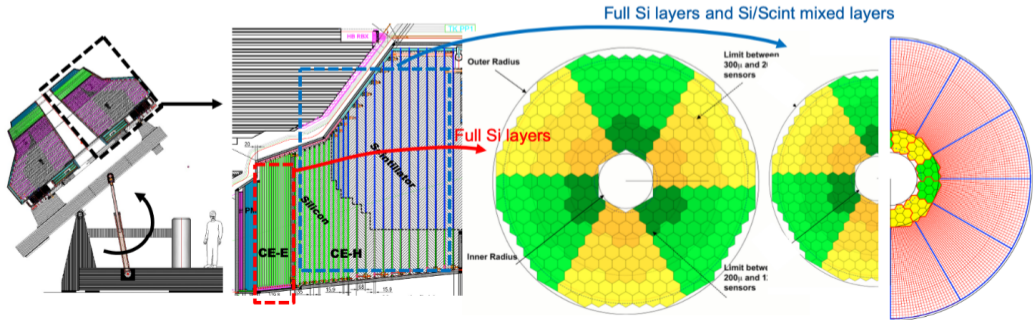
## ⇒ Silicon sensors

- measure electromagnetic and hadronic showers
- 28 (em) + 22 (had) layers
- ~30K wafers (~6M channels)
- area of 620 m<sup>2</sup>

## ⇒ Plastic scintillator tiles + SiPM

- measure hadronic showers
- 14 layers
- ~4K tiles (~240K channels)
- area of 400 m<sup>2</sup>

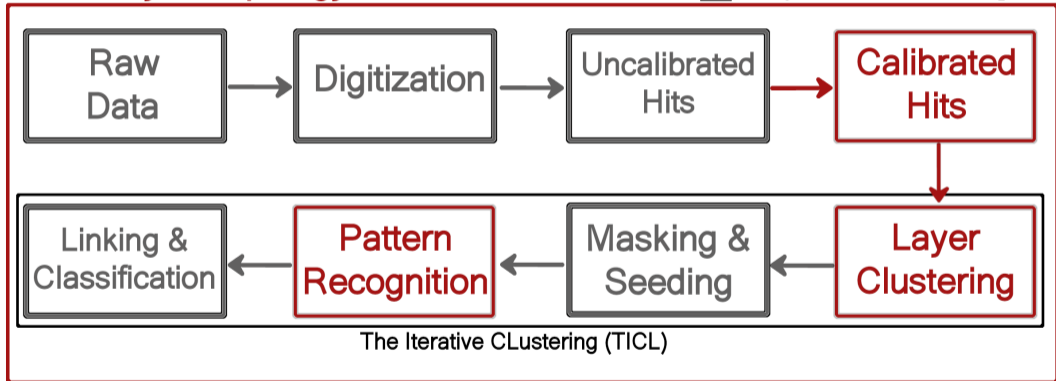
⇒ **Total weight:** ~ 220 × 2 tonnes



# HGCAL Reconstruction Chain

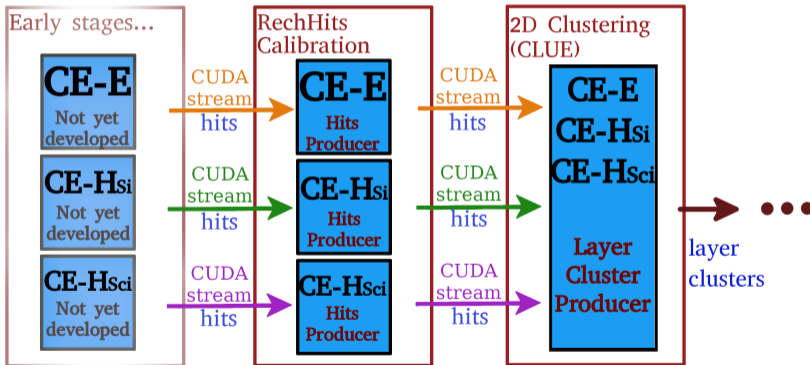
## Geometry & Topology

 Under GPU development  
 Not yet under GPU development



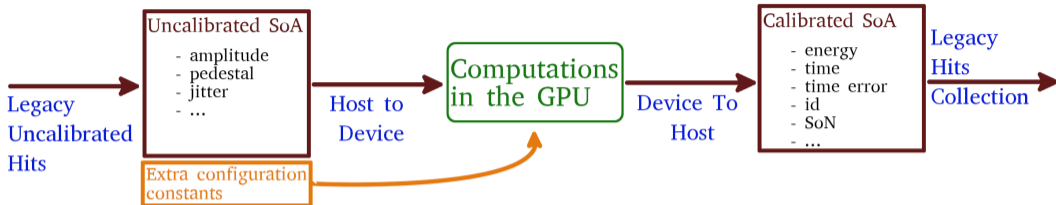
# Energy rescaling of uncalibrated energy deposits in the GPU

- Each CUDA thread maps to a single hit (no hit-communication required)
- The subdetectors allow the usage of CUDA streams for parallelization



- The advantages of **Structures of Arrays** (SoAs) are exploited

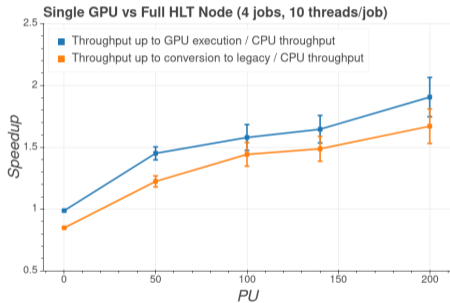
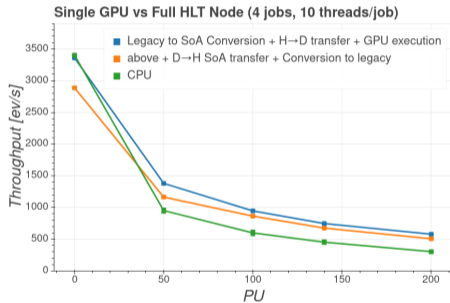
# Energy rescaling of uncalibrated energy deposits in the GPU



- The workflow is **split into four** for flexibility and validation purposes:
  1. Conversion of the CMSSW legacy uncalibrated hits into a SoA followed by its transfer to the GPU;
  2. Execution CUDA kernels in the GPU;
  3. Transfer the calibrated SoA back to the CPU (optional);
  4. Conversion of the SoA to legacy calibrated hits and storage (optional).
- Code **fully integrated** in the official CMS software.
- Validation shows **perfect agreement**.

# Performance

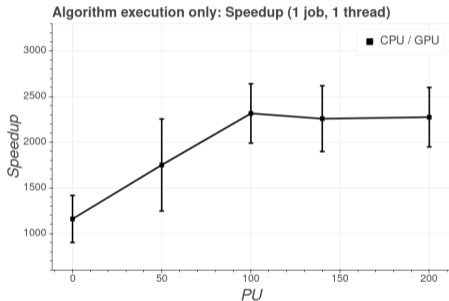
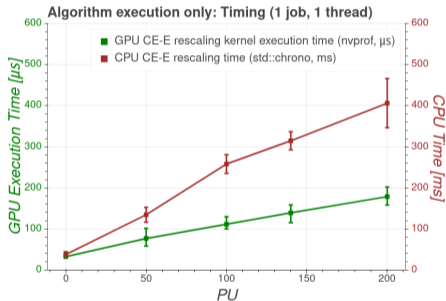
- **Full Node** (Intel(R) Xeon(R) Silver 4114 with 40 logical cores) vs. **Single GPU** (T4 Nvidia)



- **~50 MiB per CUDA stream** of GPU peak memory (**~550K hits/event** at PU200)

# Performance

- Timing of the **rescaling algorithm only**
- `nvprof` for kernel execution and `std::chrono` for CPU algorithm

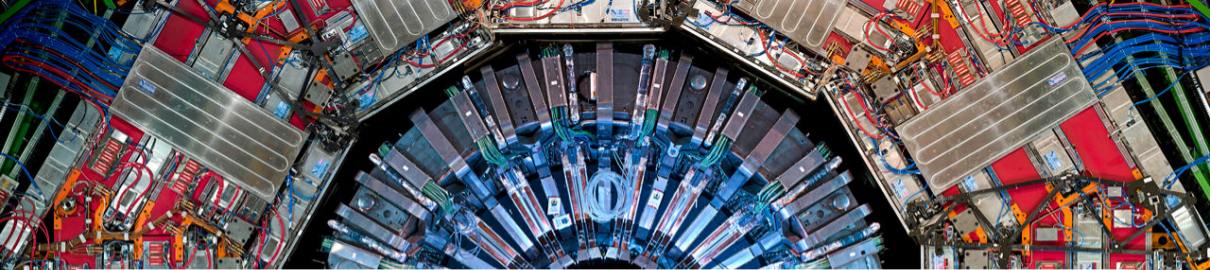


- Speedup of around **three orders of magnitude** (for CE-E with  $\sim 500$ K hits/event)

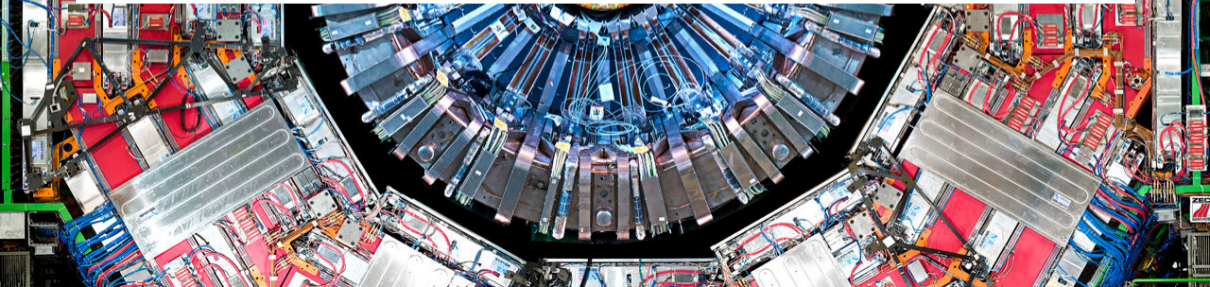


# Conclusions

- The **rescaling of HGCAL energy deposits** has been **ported to GPUs** and was **included in the CMS software** (first time for HGCAL)
- **Validation** and **performance measurements** were performed
- The results show **speedup benefits**
  - ⇒ **Cost benefits** are expected too
- This work paves the way for future HGCAL accelerator developments
  - ⇒ Link current GPU modules and test their performance
  - ⇒ Use abstraction libraries (eg. [alpaka](#))
  - ⇒ Port the **full HGCAL reconstruction** into GPUs



***Back-up***



# Acknowledgments

Bruno Alves thanks the following institutions for the funding provided:

- Fundação para a Ciência e Tecnologia (FCT): SFRH/BEST/150186/2019
- Organisation européenne pour la recherche nucléaire (CERN)