# Next Generation Triggers for CMS

**Work Package leaders and Task leaders:**
Andrea Bocci (CERN), Cristina Botta (CERN), Silvio Donato (INFN-Pisa), Jennifer Ngadiuba (Fermilab),
Felice Pantaleo (CERN), Giovanni Petrucciani (CERN), Marco Rovere (CERN),
Sioni Summers (CERN), Thiago Tomei (SPRACE)

**Steering committee inside CMS: L1T & HLT Coordination, WP leaders**
**Spokesperson and CERN Team Leader**
Cristina Botta (CERN), Silvia Goy Lopez (CIEMAT), Marino Missiroli (PSI/UZH), Marco Rovere (CERN)
Patty McBride (Fermilab), Luca Malgeri (CERN)

# NGT objectives

Enhance the triggers and the data collection and processing, and thus the scientific potential, of ATLAS and CMS in the HL-LHC phase beyond the currently projected scope

- Accelerate the evaluation and introduction of **novel computing, engineering and scientific ideas** already with demonstrators for Run3, but with main focus on HL-LHC

- Provide a major push to the work already ongoing in the experiments, by enabling **lines of research currently not feasible within existing financial, human and technology constraints**

- Provide **critical insight to develop data flows** for the even more ambitious objectives of a future collider, such as the Future Circular Collider (FCC) currently in its Feasibility Study phase

The EP, IT and TH Departments are also involved to **ensure that other current & future CERN experiments benefit from the results** in terms of computing frameworks and theoretical modelling.

All project results (IP) will belong to CERN and will be released under a valid open policy and IP generated will be released under appropriate open licenses **in compliance with the CERN Open Science Policy.**

# Conception of the "Next Generation Triggers" Proposal

**Sep 2022** — A group of private donors interested in supporting CERN scientific mission visits CERN to share ideas with CERN management and physics/computing experts

**Oct 2022** — Eric Schmidt contacts CERN impressed by its vision and contributions to the advancement of science and proposes that his foundation could **fund work on advanced Artificial Intelligence (AI) and computing techniques to improve ATLAS and CMS experiments data acquisition workflows**.

**Nov 2022**
**Dec 2022** — A task force composed of experts from EP, IT, TH, ATLAS, CMS and external experts works on a concept proposal. A 5-page concept is submitted to the **Eric and Wendy Fund for Strategic Innovation** for a **project of the value of 48M USD over 5 years**

**Feb 2023** — The Foundation informs CERN that they have **positively evaluated the proposal** and are ready to enter into a more detailed discussion of technical milestones, budget and legal/admin procedures

**June 2023** — The task force prepares a detailed proposal **which is validated from the administrative, legal, financial points of view, as well as international relation and reputation aspects, by the respective CERN competent bodies**

**Aug 2023**
**Sep 2023** — The proposal is **positively considered by CERN management**, legal negotiations with the Foundation lead to a draft Agreement. A proposal for approval to the CERN Council is prepared

**Oct 2023** — The proposal is approved by the Council, the final negotiations steps start

**Nov 2023** — The final Agreement is validated by the respective legal services, signatures finalised on Nov 22nd

# NGT Work Packages

**Two main experiment-specific work packages:**

- **Enhancing the ATLAS Trigger and Data Acquisition (WP2):**
  - Develop novel approaches to trigger event selection that will extend the ATLAS physics potential, in particular for search for events with exotic event topologies.
  - Develop state-of-the-art Machine Learning techniques for the online event selection at HL-LHC.
  - Enable the collection of richer data set of collision event, which will extend the reach of the ATLAS physics programme.
  - Extend the capabilities of the trigger and data acquisition system by efficient use of acceleration technologies.
- **Enhancing the CMS Real Time Data Processing (WP3):**
  - Redesign the data collection and scouting strategy to reduce the need to reject events in the Level-1 and High-Level CMS triggers aiming at complementing the current workflows
  - Replace the trigger filtering task with an event processing task similar to what happens with offline events stored on disk.
  - Leverage traditional physics-based algorithms and advanced AI solutions, remove the bottleneck currently implied by the real-time event selection and extend CMS discovery and precision measurement reach.

# NGT Work Packages

**Two general-purpose work packages:**

- **Infrastructure, Algorithms and Theory (WP1):**
  - Design and provision a dedicated hybrid computing infrastructure including a local cluster of low-latency/high-bandwidth interconnected GPUs, cloud resources and services and scalable AI workflows
  - *Develop software, tools, and methodologies to optimise Neural Networks on accelerated architectures and different performance/accuracy/energy consumption targets*
  - Develop novel event generators and analysis algorithms and common applications based on cutting edge AI and quantum-inspired technologies to support the design on novel detection strategies in the LHC experiments
  - *Implement a common set of tools and software engineering methodologies to efficiently develop, run and orchestrate optimised experiment and theory code on accelerated computing architectures*
- **Education Programmes and Outreach (WP4):**
  - Promote exchanges across computer scientists and physics researchers, academia and industry
  - Complement and extend existing education programmes to train data science and AI skills for the next generation of high-energy physicists
  - Leverage the rich ecosystem of education programmes at CERN, across HEP, and in computer sciences communities to provide targeted specialisation paths for new researchers

*Text like this means there's a strong involvement from CMS members*

**More detailed information on WP1 and WP4 will be provided in future presentations.**

# NGT Proposal and Budget
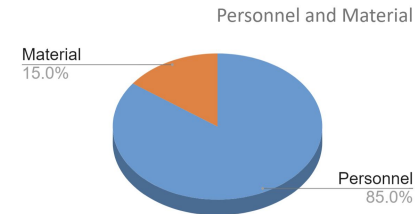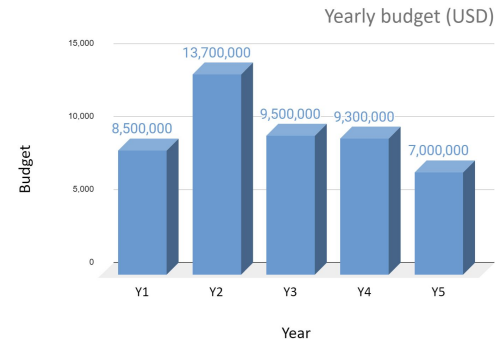
**The project funding is 48M USD distributed on 5 years:**
− 19M for common aspects (including funding for an HPC-like GPU cluster): **Management and Work Package 1 (WP1)**
− 13M for ATLAS, 13M for CMS (to be split 50%/50% for HLT & L1T): **Work Packages 2 and 3 (WP2, WP3)**
−  3M for education/outreach: **Work Package 4 (WP4)**

**The Project objectives are implemented through a series of Milestones:**
− every WP has a specific set of yearly milestones attached to the planned Tasks (14 milestones per each of the 5 years of the project across the 4 WPs, each milestone has a specific format and measurable deliverables e.g. report, software, demo, event, etc…)

**The full proposal with detailed Tasks, Milestones and deliverables for each of the WP can be found in agenda**

| WP | Personnel costs (USD) | Material costs (USD) | Total Cost (USD) |
|---|---|---|---|
| Management | 1.8M | 0.0M | 1.8M |
| WP1 | 11.2M | 6.0M | 17.2M |
| WP2 | 12.4M | 0.6M | 13.0M |
| WP3 | 12.4M | 0.6M | 13.0M |
| WP4 | 3.0M | 0.0M | 3.0M |
| Total | 40.8M | 7.2M | 48.0M |
| Percentage | 85% | 15% | 100% |



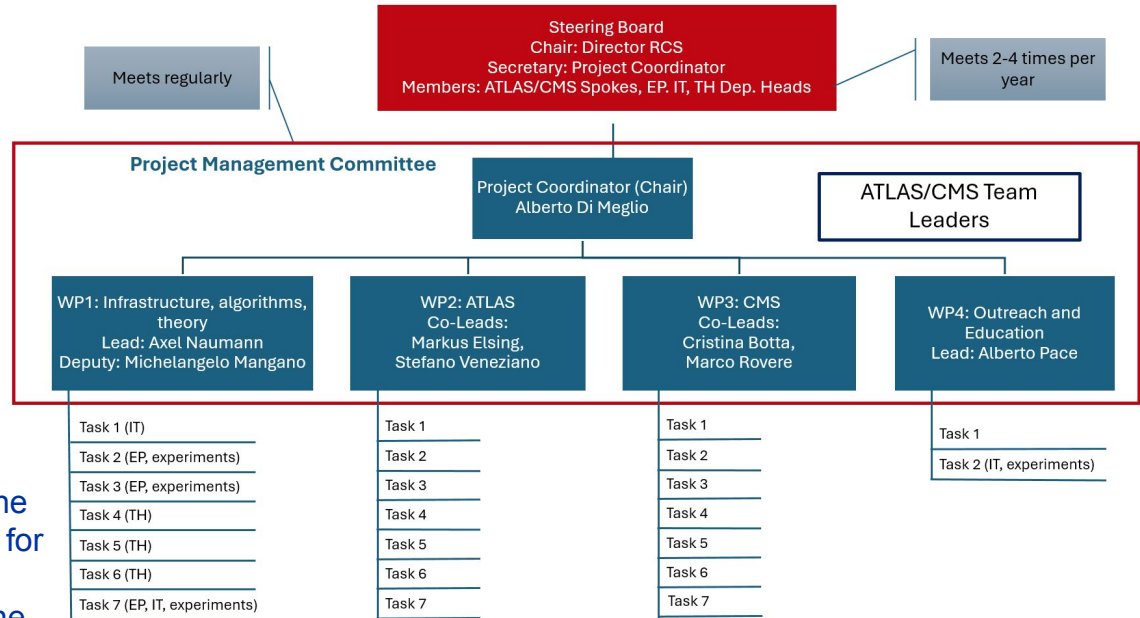Yearly budget (USD)



Personnel and Material

*The yearly donation from the fund is based on a preliminary budget estimate associated to each milestone*
***It will be paid based on the evidence of reaching the milestones for the previous year through the submission of the associated deliverables***

6

# NGT Organigram

**The Steering Board (SB) has just approved on Monday the following organigram:**

- A **Management Committee (PMC)** provides day-to-day executive governance of the project, it is led by a Project Coordinator and composed of the Work Package leads and ATLAS/CMS Team leaders

- A **Steering Board (SB)** led by the CERN Director for Research and Computing and composed of CERN stakeholders' representatives provides oversight and alignment to CERN, ATLAS and CMS objectives and interests.
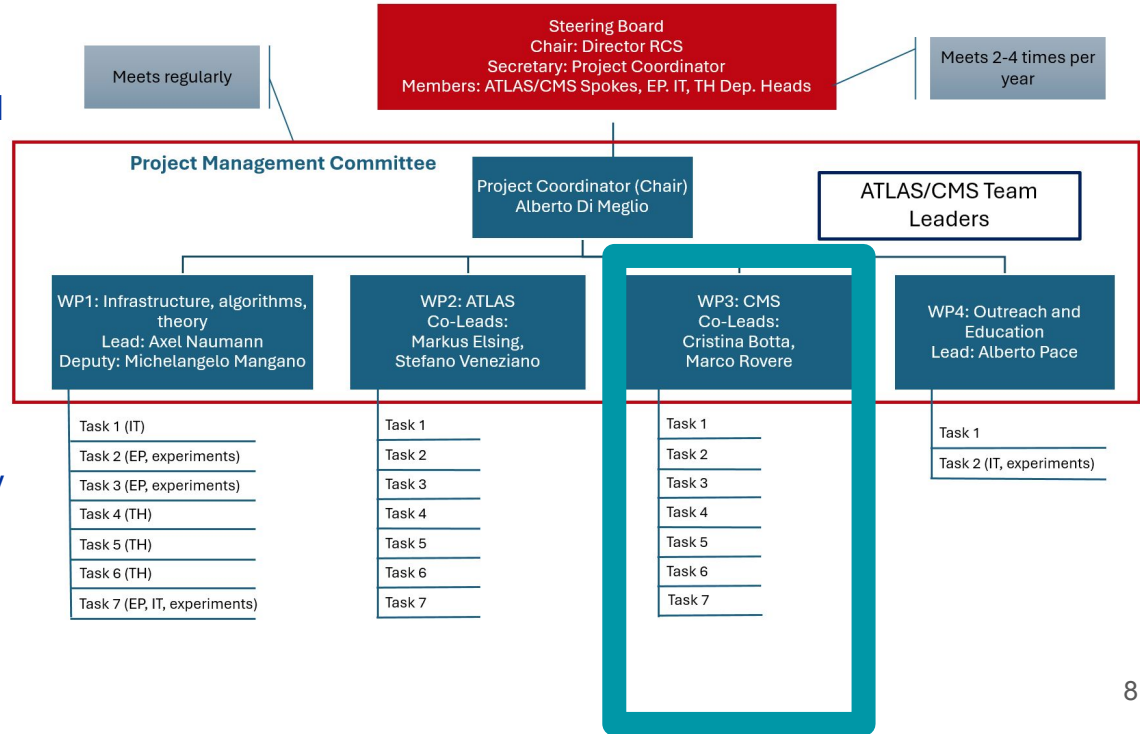
→ the SB has selected the **PC** and endorsed the names proposed by Experiments/Departments for the **WP leaders**

→ the PCM has to still endorse the names of the **Task leaders**, which are proposed by the WP leaders (jointly with the Experiments Steering Committees)



Steering Board
Chair: Director RCS
Secretary: Project Coordinator
Members: ATLAS/CMS Spokes, EP. IT, TH Dep. Heads

Meets regularly

Meets 2-4 times per year

Project Management Committee

Project Coordinator (Chair)
Alberto Di Meglio

ATLAS/CMS Team Leaders

WP1: Infrastructure, algorithms, theory
Lead: Axel Naumann
Deputy: Michelangelo Mangano

WP2: ATLAS
Co-Leads:
Markus Elsing,
Stefano Veneziano

WP3: CMS
Co-Leads:
Cristina Botta,
Marco Rovere

WP4: Outreach and Education
Lead: Alberto Pace

| WP1 | WP2 | WP3 | WP4 |
|---|---|---|---|
| Task 1 (IT) | Task 1 | Task 1 | Task 1 |
| Task 2 (EP, experiments) | Task 2 | Task 2 | Task 2 (IT, experiments) |
| Task 3 (EP, experiments) | Task 3 | Task 3 | |
| Task 4 (TH) | Task 4 | Task 4 | |
| Task 5 (TH) | Task 5 | Task 5 | |
| Task 6 (TH) | Task 6 | Task 6 | |
| Task 7 (EP, IT, experiments) | Task 7 | Task 7 | |

# NGT WP3: implementation in CMS

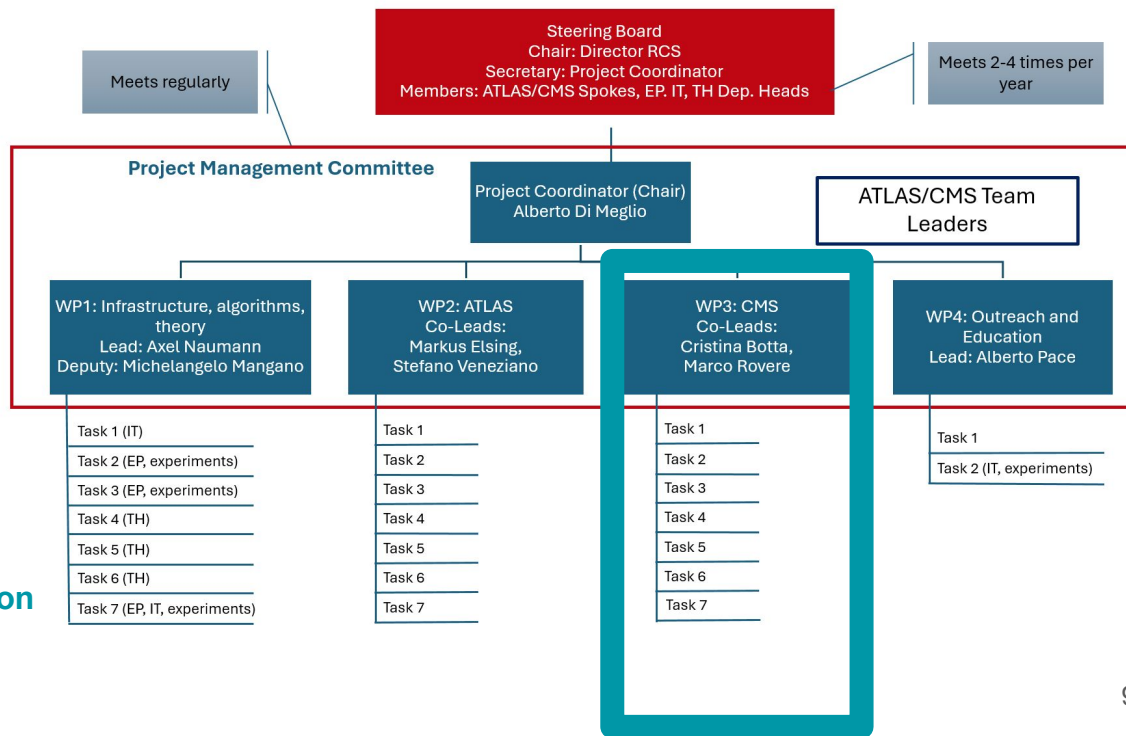**WP leaders, L1T & HLT Coordination, Spokesperson and CERN Team leader oversee the implementation in CMS**

− Hiring will start in Q1 of 2024
− CERN positions : Doctoral/Technical students, Graduates, staff → **we request all positions to be open to member states and not-member states (still under discussion at CERN management level)**
− Selection committee (when required by CERN policy/rules) will include a non-CERN member of CMS (to be defined)
− Allocate an important fraction (~15%) of the overall budget to host experts on site (project-associate) and to allow subsystency for supervisors of Doctoral students to stay at CERN for some time
− Overheads for each position: 3% of material/travels, 10% of typical Staff salary for supervision and M&O-A (when applicable)



Steering Board
Chair: Director RCS
Secretary: Project Coordinator
Members: ATLAS/CMS Spokes, EP. IT, TH Dep. Heads

Meets regularly

Meets 2-4 times per year

**Project Management Committee**

Project Coordinator (Chair)
Alberto Di Meglio

ATLAS/CMS Team Leaders

WP1: Infrastructure, algorithms, theory
Lead: Axel Naumann
Deputy: Michelangelo Mangano

WP2: ATLAS
Co-Leads:
Markus Elsing,
Stefano Veneziano

WP3: CMS
Co-Leads:
Cristina Botta,
Marco Rovere

WP4: Outreach and Education
Lead: Alberto Pace

Task 1 (IT)
Task 2 (EP, experiments)
Task 3 (EP, experiments)
Task 4 (TH)
Task 5 (TH)
Task 6 (TH)
Task 7 (EP, IT, experiments)

Task 1
Task 2
Task 3
Task 4
Task 5
Task 6
Task 7

Task 1
Task 2
Task 3
Task 4
Task 5
Task 6
Task 7

Task 1
Task 2 (IT, experiments)

# NGT WP3: implementation in CMS

**The WP leaders, L1T & HLT Coordination, Spokesperson and CERN Team leader oversee the implementation in CMS**

**Tasks and task leaders** (to be endorsed by PCM)

1.1 **$R^3$ Faster Reconstruction for HLT**
     **M. Rovere (CERN)**

1.2 **Optimized data structures for HLT**
     **F. Pantaleo (CERN)**

2. **Towards a distributed HLT architecture**
     **A. Bocci (CERN)**

3. **Reduction of the RAW data size for HLT**
     **S. Donato (INFN-Pisa)**

4. **Optimal calibrations for HLT**
     **T. Tomei (SPRACE)**

5. **Enhancing L1T Scouting for HL-LHC**
     **G. Petrucciani (CERN)**

6. **Practical real-time AI for L1T**
     **S.P. Summers (CERN)**

7. **L1T anomaly detection & data compression**
     **J. Ngadiuba (FERMILAB)**

# WP3: Enhancing the CMS Real Time Data Processing

- Experiments can only process a fraction of the LHC collisions, while the majority have to be discarded by the real-time selection system (the trigger system) due to limited storage resources.
- This work package aims to remove this limitation by **greatly enhancing the CMS trigger systems**.
- We aim to redesign the data collection strategy to reduce (and when possible remove) the need to reject events and to optimize the information stored per event.
- We aim to achieve this by leveraging traditional physics-based algorithms and advanced AI solutions.
- Two directions will be explored:
  - **extensions and alternative approaches for the Level-1 Trigger scouting,** to enhance the capability to perform physics analyses on all collisions events with only Level-1 Trigger information
  - next-generation High-Level Trigger system ("**Real-time Reconstruction Revolution**" or **R³**) capable of performing real-time full event reconstruction on all events selected by the L1T, without further filtering before data storage.
- These two major steps would allow the CMS collaboration to remove the bottleneck induced by the real-time event selection and extend our discovery and precision-measurement reach.
- Proof-of-concept R&D projects on these activities are ongoing. Some of them already established the validity of the ideas behind the proposed tasks.

# WP3 previous talks

- Offline and Computing Week: [link](link)
- Trigger Studies Group(TSG) meeting: [link](link)
- L1T workshop in Athens: [link](link)
- Management board meeting: [link](link)

# WP3: Enhancing the CMS Real Time Data Processing

## HLT Related Tasks (3.1.1, 3.1.2, 3.2, 3.3, 3.4)

A. Bocci (CERN), S. Donato (INFN-Pisa), F. Pantaleo (CERN),
M. Rovere (CERN),  T. Tomei (SPRACE)
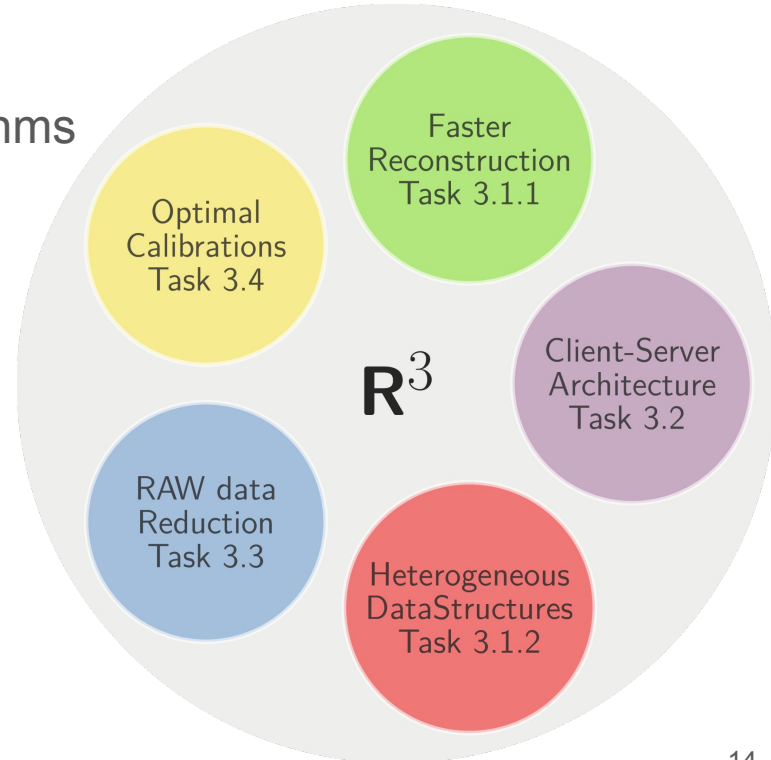
# The **R**eal-time **R**econstruction **R**evolution (**R³** - Rcube)

- **Overcome the two main limitations of the HLT**
  - the quality of the online reconstruction is limited by the processing capacity of the HLT farm, as complex algorithms can be run only on a fraction of the events
  - the HLT output rate is limited by the storage capacity and processing power of the offline computing infrastructure.
- **What if we could** …
  - have *offline-like quality* calibrations and reconstruction at the HLT?
  - store *all events* in nano-AOD format?

**The goal of the R³ project is to address these limitations through a comprehensive work program, consisting of five synergic tasks.**

# The **R**eal-time **R**econstruction **R**evolution (**R³** - Rcube)

The **R³** project seeks to

- develop next-generation reconstruction algorithms
  - greatly speed up reconstruction
  - improve online reconstruction quality to match full reconstruction
- leverage
  - heterogeneous compute resources
  - data-oriented programming
  - proven methodologies
  - advanced AI-based techniques



$R^3$

Faster Reconstruction Task 3.1.1

Optimal Calibrations Task 3.4

Client-Server Architecture Task 3.2

RAW data Reduction Task 3.3

Heterogeneous DataStructures Task 3.1.2

# Work Package 3: HLT Tasks

- **3.1.1**: Develop modern, heterogeneous-friendly, algorithms & data structures for reconstruction of all CMS physics objects
- **3.1.2**: Optimized data structures for heterogeneous architectures
- **3.2**: Distributed processing for HLT within CMSSW
  - support offloading of algorithms to remote accelerators
- **3.3**: Investigate & develop approaches for RAW data size reduction
  - *e.g.* lossless, lossy data compression or replacing raw with local reco data
  - both approaches already being investigated for Heavy Ion running
- **3.4**: Solutions to buffer the HLT input data for a few hours, run fast online calibrations and use them for final HLT processing
- *All activities will be consolidated under HLT Upgrade in the Trigger-HLT in CMS*

Total resources: 6.4M personnel + 100k material (over 5 years)

# PRELIMINARY allocation of resources

**Budget**



- RCube
- DataStructures
- ClientServer
- RawPrime
- OptimalCalibration

**FTEs**



- RCube
- DataStructures
- ClientServer
- RawPrime
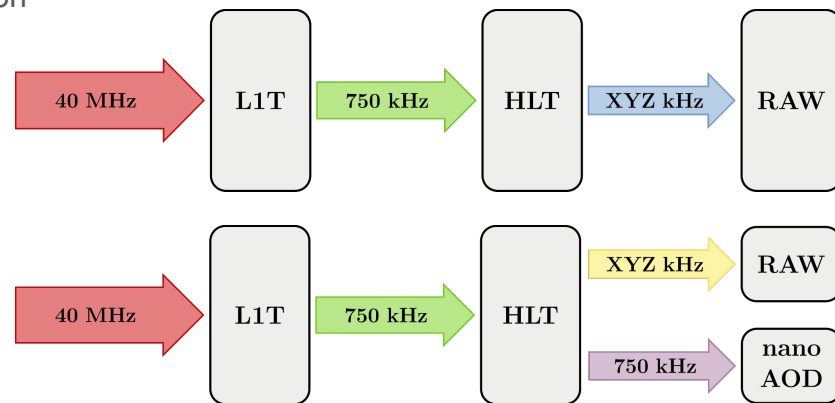- OptimalCalibration

**Personnel Costs by Category over 5 years**



- Overhead 14.6%
- Staff 14.6%
- Undergraduate 4.5%
- Graduate 37.9%
- Project Associates 12.2%
- Doct 16.2%

## Preliminary positions opening expected in 2024

- ○ 3.1.1: 1 staff, 1 fellow, 2 doctoral
- ○ 3.1.2: 1 fellow
- ○ 3.2:    1 doctoral
- ○ 3.3:    1 doctoral
- ○ 3.4:    1 fellow, 1 doctoral
- ○ 3 project-associates (PJAS) foreseen, when and where needed.

# Task 3.1.1: R³ Faster Reconstruction (*task leader M. Rovere*)

- The successful Patatrack experience in CMS has shown that it is possible to improve the physics quality and reconstruction throughput of selected physics objects (pixel tracks) by leveraging heterogeneous architectures
- This required ~4 years of development to:
  - Study the performance of the current algorithm and identify bottlenecks
  - Rethink the algorithms and data structures targeting heterogeneous architectures
  - Develop, integrate and validate the results in CMSSW
  - Propagate the new objects to the rest of the reconstruction
- The R³ project will use a similar approach to redesign the most important physics objects:
  - Muons
  - Electrons and photons
  - Taus
  - Jets, MET and Particle Flow Global Event interpretation
- Perform offline-like full event reconstruction, in addition to the traditional event selection
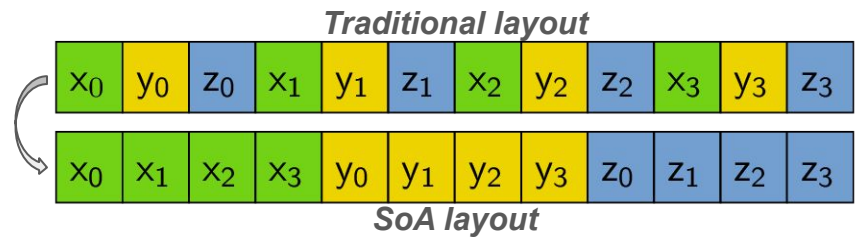
# Task 3.1.2: **R³** Optimized data structures for heterogeneous platforms *(task leader F. Pantaleo)*



**Traditional layout**

**SoA layout**

- The development of **data-oriented structures** ("Structure of Arrays", SoA for short) **will be fundamental for R³ to reach its goal**.
  - achieve better memory bandwidth and vectorization performance
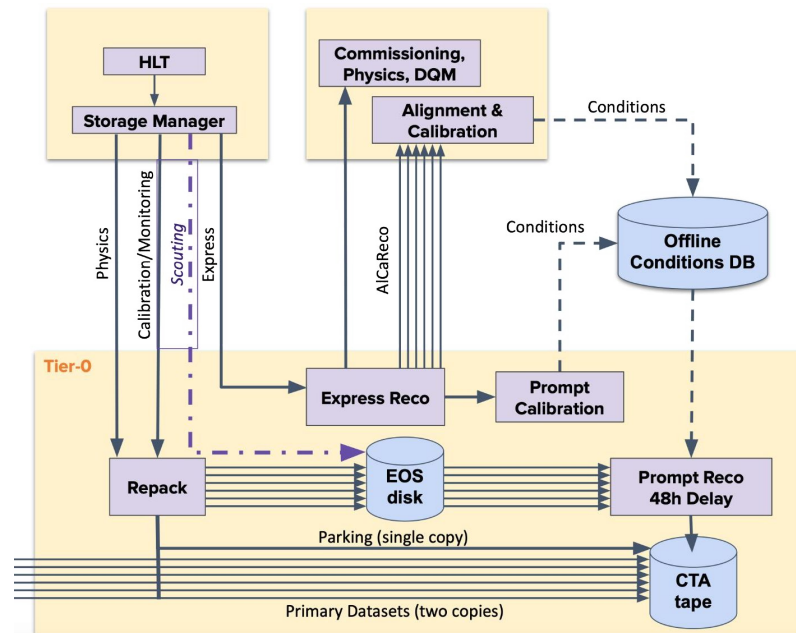  - provide a seamless interface to AI algorithms

- Their adoption in the HEP software stack requires the development of a user-friendly, **generic SoA implementation**.

- To achieve the best performance running real-time reconstruction, the I/O subsystem of the CMS framework should be extended to leverage direct data transfers between the network and storage subsystems on one side, and the accelerators on the other, bypassing the host CPU.

# Task 3.4: R³ Optimal calibrations for HLT *(task leader T. Tomei)*

- Design accelerated calibration workflows to achieve at HLT the same accuracy as the offline reconstruction
  - optimize the calibration process for the CMS detectors
  - introduce data buffering online
  - exploit predictive AI techniques
- Synergy with Run-3 operations
  - deploy a prototype applied to the HLT Scouting workflow during the last year of Run-3
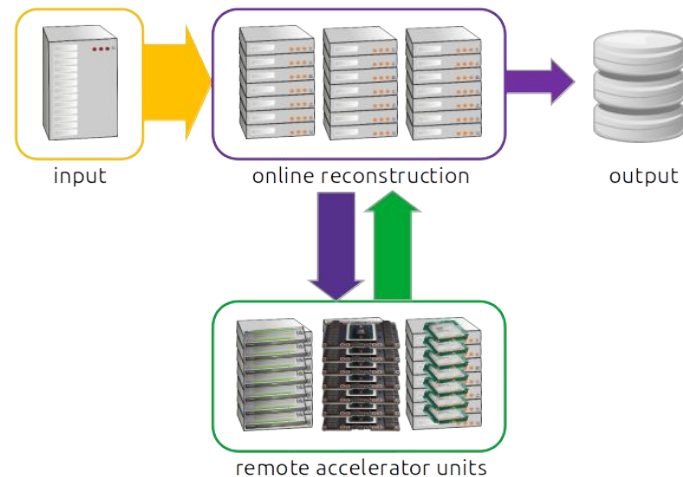  - rethink the hardware and software infrastructure for the calibration workflow



*the Prompt Calibration Loop as it is today*

# Task 3.2: R³ Towards a distributed HLT architecture
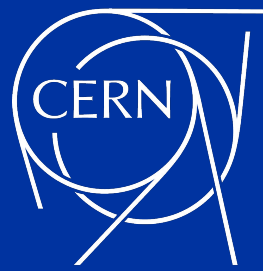## *(task leader A. Bocci)*

- Extend the **CMS data processing framework** to span multiple nodes, **adapt to different network topologies** and **leverage remote accelerators**, with little or no modification to the physics reconstruction code.

- Leverage the *performance portability* approach and the task-based framework used in CMSSW

- Avoid strict requirements in terms of hardware design of the hosting computing centre

- Anticipate the use of novel accelerator technologies: FPGAs, TPUs, quantum devices…



input → online reconstruction → output

remote accelerator units

# Task 3.3: R³ Reduction of the RAW data size for HLT
## *(task leader S. Donato)*

- A **limiting factor** in the amount of data that the HLT can select for offline storage and further processing is the **size of the RAW events**.

- Characterize multiple approaches to the **compression of RAW data**, with different trade-offs between the compression factor, latency, available hardware and impact on the final physics result.
    - lossless compression on accelerators
    - lossy compression algorithms
        - build upon the work done by the Heavy Ions group
    - physics-driven compression, replacing basic information with higher-level quantities
        - leverage prompt-reconstruction-level calibrations to reduce the physics impact

- reducing the RAW size:
    - increase the rate of full RAW data collected by the HLT
    - leave more space for the scouting data
    - reuse offline the high-level quantities reconstructed online

# WP3: Enhancing the CMS Real Time Data Processing
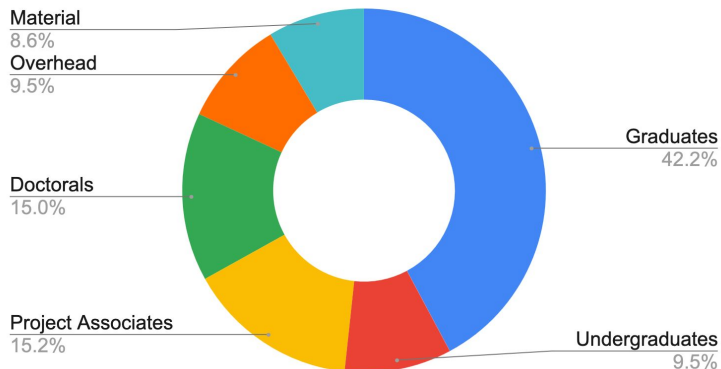
## L1T Related Tasks (3.5, 3.6, 3.7)

*Cristina Botta (CERN),* Jennifer Ngadiuba(Fermilab),
Giovanni Petrucciani (CERN), Sioni Summers(CERN)

# WP3: L1T tasks and preliminary resources

- **3.5:** R&D for for L1T data scouting
  - 3.5M (500k material)
- **3.6:** Development of ML algorithms and deployment practices for L1T
  - 1.5M
- **3.7:** Anomaly detection & auto-encoder based data compression
  - 1.5M

**Total resources: 6M personnel + 500k material (total on 5 years)**
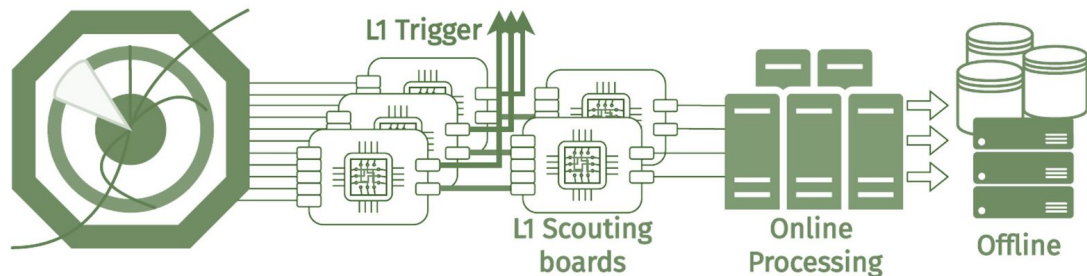
TOTAL over 5 years



Material 8.6%
Overhead 9.5%
Doctorals 15.0%
Project Associates 15.2%
Graduates 42.2%
Undergraduates 9.5%

Overhead = Supervision/Material-Travel

**Preliminary positions opening expected by the end of 2024**
- 3.5 : 1 fellow, 1 doctoral, 1 technical student, 1 project-associate
- 3.6 : 1 fellow, 1 doctoral
- 3.7 : 1 fellow, 1 doctoral

*All activities will be consolidated under the standard L1TDAQ meetings in CMS*
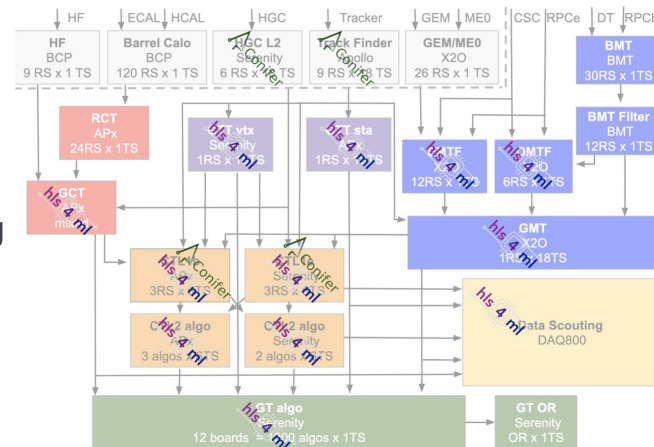
# L1T Scouting for HL-LHC



- L1T Data Scouting: **acquire and analyse the L1 Trigger information for all events (*i.e.* at 40 MHz)**
- Look for physics signatures identifiable with **just L1 information** but that would evade the L1T → HLT → Offline chain, e.g.:
  - Too large "irreducible" backgrounds, e.g. narrow resonances of unknown mass (already done using HLT scouting for muons & hadronic resonances)
  - Signal identification requires an algorithm that can't fit the L1 fixed latency and resource budget, e.g. has too complex combinatorics on some events (soft multiple hadronic objects)
  - Signal identification requires time-correlation across several BXs, beyond what allowed by GT, e.g. very slow or long-lived BSM
- FPGA-equipped boards that receive L1 data via optical links and transfer it to PCs and the software world via TCP/IP or PCI express
- at HL-LHC: can profit from much improved L1T object reconstruction quality
  - L1 Tracker Tracks, Particle-flow linking, Pile-up per particle identification, etc.
  - But being demonstrated already in 2018 (muons), Run 3 (muons + calo)

24

# Task 3.5: Enhancing L1T Scouting for HL-LHC *(Task leader G. Petrucciani)*

- **Goal of this project:** explore the physics opportunities and technical feasibility using different L1 inputs, R&D to investigate different implementations
  - baseline as in L1T TDR foresee streaming through scouting all the high-level objects and PUPPI candidates: can this be extended to PF candidates, all L1 Tracks and others TPs?
- How: investigate **different prototype analyses of increasing complexity** in terms of data processing and bandwidth: mixing traditional algorithms with AI solutions
  - from classical dilepton resonances to soft hadronic final states with NN, GNN taggers
  - jet reconstruction & tagging with more complex algorithm
- How: incrementally build a **test system of increasing bandwidth & processing power** to run the benchmark analyses
  - investigate running algorithms on CPU / GPU / FPGA / AI engines
  - investigate different approaches for data acquisition: protocols (e.g. TCP/IP vs RoCE vs direct fpga-fpga links), networking (e.g. NICs vs accelerator cards or converged NIC+GPU), workflows (HLT-like, analysis-like, kafka,..)
  - prototype a Scouting DAQ board with newer technology wrt DAQ-800
    - Option: Large Versal HBM chips vs Virtex Ultrascale+ HBM to have more links per FPGA to allow data aggregation, zero suppression, demultiplexing and redistribution of larger data streams ?

# Task 3.6: Practical real-time AI for L1T *(Task leader S.P.Summers)*

- **Develop ML algorithms for improving the L1T reconstruction (used by standard triggers, and L1 Scouting)**
  - For Global Trigger (also for Run 3), Correlator and Global Track Trigger subsystems
  - This is already happening, but the project would boost it profiting from experience on advanced ML methods
- **Develop practices for training & deployment of ML algorithms in L1T**
  - updating & redeploying algorithms for changes in detector conditions
  - tracking, archiving & retrieving ML models (e.g. for consistent emulation)
  - gain operational experience from Run3 Global Trigger and develop for HL-LHC
- **Develop setup for complex trainings of multiple algorithms**
  - e.g. simultaneously optimize algorithms for particle ID, object reconstruction, and event selection which are implemented in different subsystems (only limited amount of information is propagated between subsystems)



26

# Task 3.7: L1T anomaly detection & data compression
*(Task leader J. Ngadiuba)*



Input collision    Regularize latent space to avoid overfitting    Sampled latent representation    Reco collision

$$\text{loss} = \| x - \hat{x} \|^2 + KL[\ N(\mu_x, \sigma_x), N(0, I)\ ]$$

- **Commissioning and validation of the Run 3 unsupervised anomaly detection (AD) algorithms**
  - Learn end-to-end algo integration and system operations
  - Optimize the model to be more robust against conditions (e.g. pileup)
- **Demonstrate end-to-end physics analysis using unsupervised anomaly detection on Run 3 data**
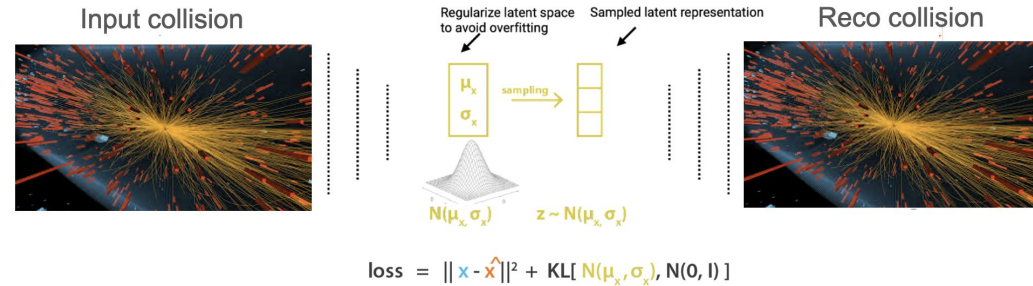  - Develop innovative methods for the analysis of anomalous data: a) seed/enhance supervised searches with the new AD algo and b) use new methods (e.g. active learning or clustering) to inform what to search next
- **Develop new AD model suited for HL-LHC L1 trigger system using TPs/particle-based info**
  - Design and optimize an innovative GNN-based autoencoder for low latency and resources which will serve as baseline AD at HL-LHC
- **Develop intelligent data compression and/or reduction for L1 Scouting system (to allow long-term storage of more scouting data)**
  - Develop an efficient event-level AI embedding method that allows a compression level able to store all data in an unbiased way while maintaining capabilities for interesting analysis-level downstream tasks (e.g. by taking inspiration from foundation models)

# Conclusions

# Conclusions

- We invite all the interested parties to contact the WP/Task leaders or directly get in touch with the L1T Scouting/Upgrade and HLT Upgrade groups.
- We remind that all the activities related to developments in WP3 (in terms of work/meetings/communications) will go through existing structures, i.e. the L1T Scouting and Upgrade Subsystems/DPG and the HLT Upgrade under Trigger-HLT.