# Next Generation Triggers for CMS

**Work Package leaders and Task leaders: WP3**

Andrea Bocci (CERN), Cristina Botta (CERN), Silvio Donato (INFN-Pisa), Emilio Meschi (CERN)

Jennifer Ngadiuba (Fermilab), Felice Pantaleo (CERN), Giovanni Petrucciani (CERN),

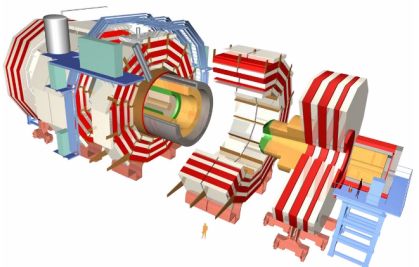Marco Rovere (CERN), Sioni Summers (CERN), Thiago Tomei (SPRACE)

**Task leaders and contacts from CMS in WP1, WP4**

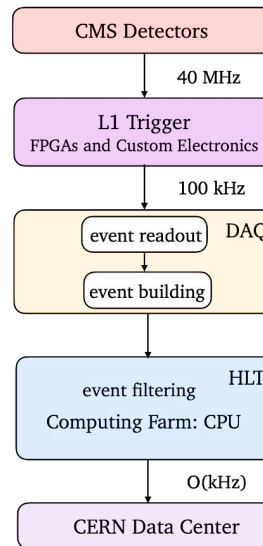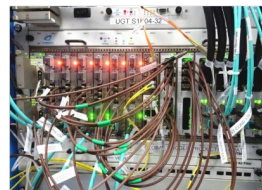Andrea Bocci (CERN), Maurizio Pierini (CERN), Felice Pantaleo (CERN)

# The CMS Trigger system (design)


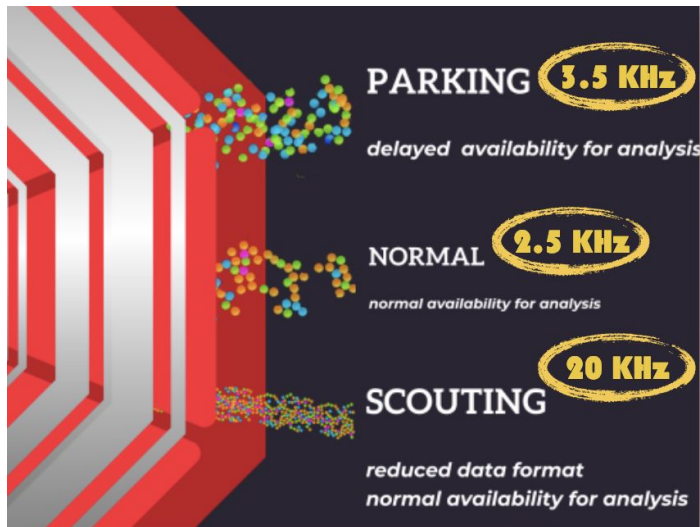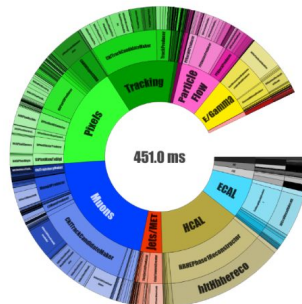CMS L1 Trigger


CMS High Level Trigger



- **Perfect compromise between technological limitations and theory motivated New Physics at O(100-1000)GeV scale with ewk-like cross sections**:
  - ○ **99.75%** events rejected at L1T, **98%** rejected at HLT

# The CMS Trigger system (now)



PARKING **3.5 KHz**
delayed availability for analysis

NORMAL **2.5 KHz**
normal availability for analysis

SCOUTING **20 KHz**
reduced data format
normal availability for analysis

- **After Higgs boson discovery and null results of NP searches in Run 1 and 2:**
  - In Run-3 pushed much further the potentialities of the CMS Data Acquisition system with: **Parking**, **HLT Scouting** and **L1T Scouting demonstrator**
  - Enlarged acceptance for **B-physics, VBF, Long-Lived particles, HH production**
  - Thanks to these new ideas & heterogeneous computing at HLT



GPUs

CMS Paper EXO



451.0 ms

(a) CPU-only

357.7 ms
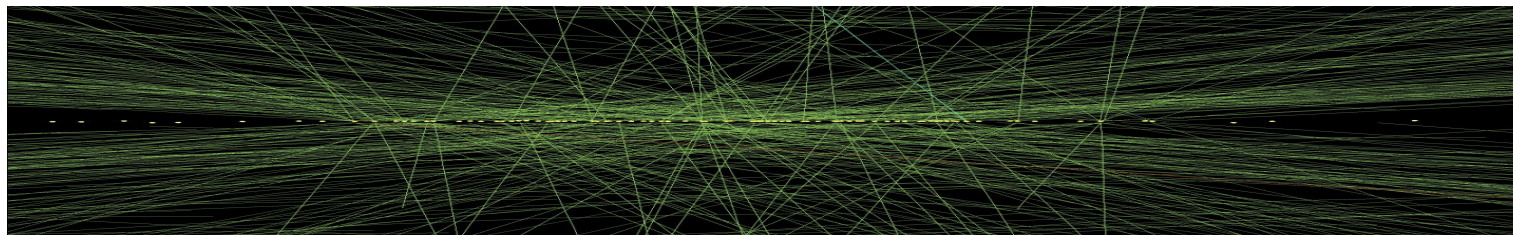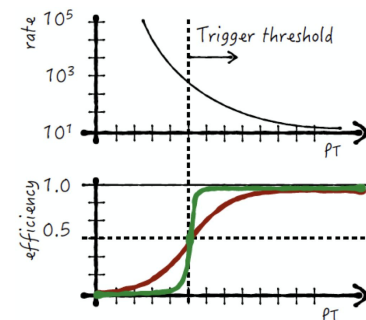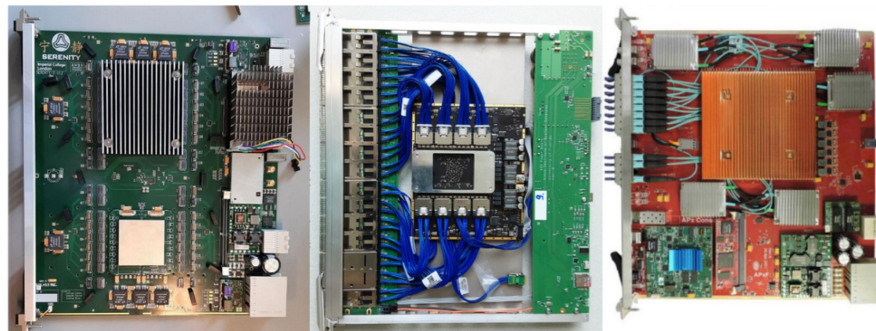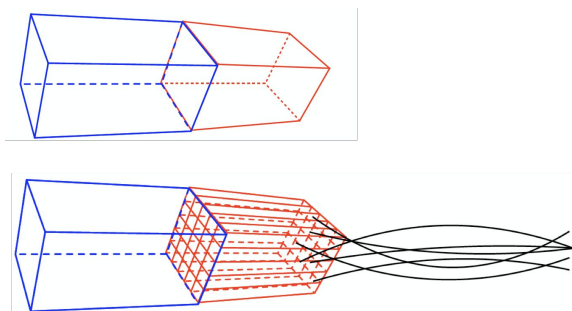
(b) CPU and GPU

# The CMS Trigger for HL-LHC

- At HL-LHC, up to 200 pile-up interactions: **CMS is upgrading the L1T and HLT to enable the same physics program we are doing now (at @60 PU)**



**New CMS L1 Trigger: input data from 2 Tb/s to 63 Tb/s (~1/10 of the internet traffic), 12.5µs to take decision**

# Rethinking the CMS Real-Time data processing (with NGT @ HL-LHC)

- **But what if New Physics doesn't look as expected,** and it's instead buried under the bulk of the background events we are throwing away due to the trigger selections?

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.
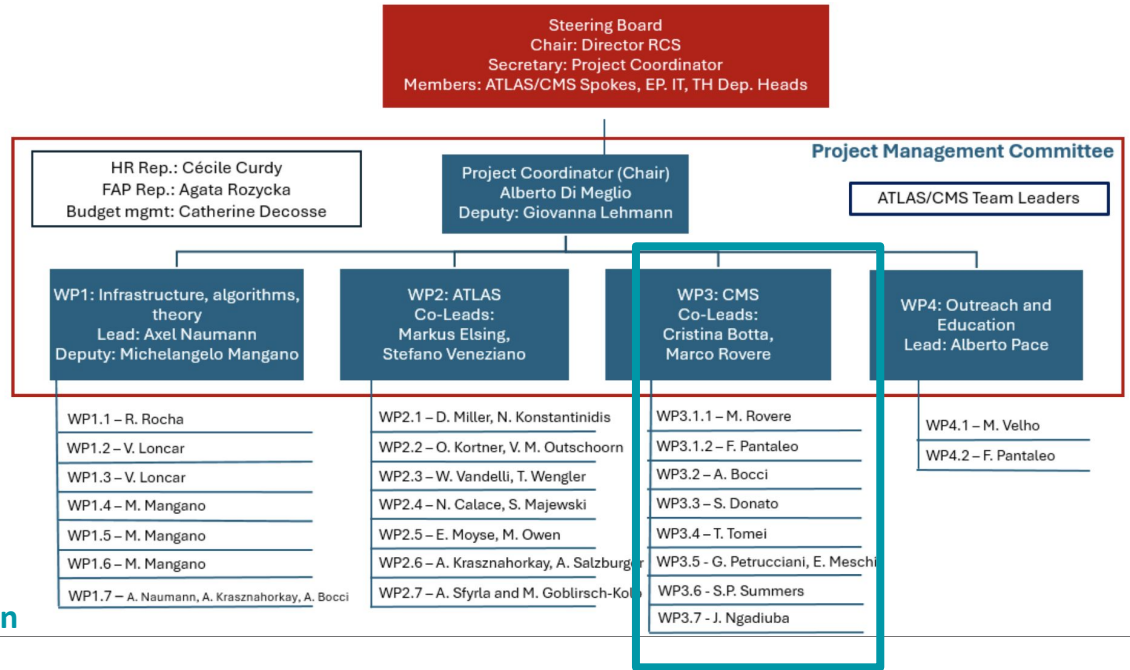
Leverage traditional physics-based algorithms and **advanced AI solutions**, remove the bottleneck currently implied by the real-time event selection and **extend CMS discovery and precision measurement reach**.

# NGT Implementation in CMS

**WP leaders, L1T/HLT/O&C/ML coordination, Spokesperson and CERN Team leader oversee the implementation in CMS (CMS Steering committee)**

**WP3 Tasks and Task Leaders**

- **1.1 R³ 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), E. Meschi (CERN)

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

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



**Participation of CMS also in WP1 and WP4 with M. Pierini (CERN) as CMS contact, A. Bocci (CERN), F. Pantaleo (CERN)**

# NGT @ L1T

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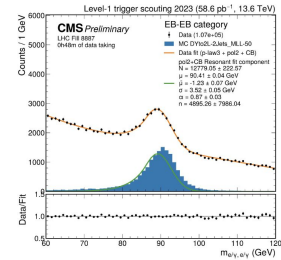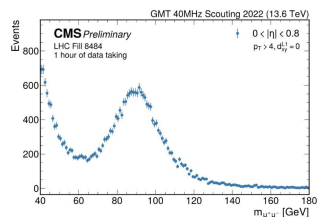Sioni Summers (CERN)

# 3.5: Enhancing L1T Scouting for HL-LHC

- **L1T Data Scouting: acquire and analyse the L1 Trigger information for all events (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 low and unknown mass
  - 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
- 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..
  - Demonstrated already in 2018 (muons), Run 3 (muons + calo)



**ICHEP's talk by Emilio**
Run3 L1T Scouting demonstrator

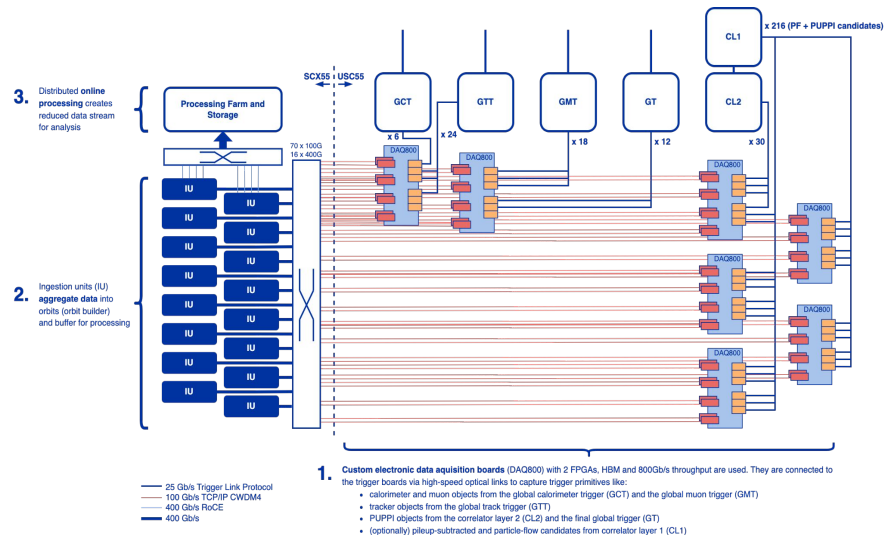# 3.5: Enhancing L1T Scouting for HL-LHC

- **Goal of this project:** explore the physics opportunities and technical feasibility using **different L1 inputs with respect to TDR baseline and R&D to investigate different implementations**



**ICHEP's talk by Emilio**
L1T Scouting
baseline for HL-LHC

- How: investigate **different prototype analyses of increasing complexity** in terms of data processing and bandwidth:
  - from dilepton resonances to soft hadronic final states with NN, GNN taggers, multi-BX signatures
  - ML algorithms for soft objects (e/γ/jet) reconstruction & tagging
- 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 from DAQ board: 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 ( Large Versal HBM chips? ) to have more links on the same FPGA and more on-board resources to allow more data aggregation and pre-processing for Scouting

9

# 3.5: Enhancing L1T Scouting for HL-LHC

- **On going:** some analyses defined and exploration started for some
  - Exclusive rare decays: $W \rightarrow 3\pi$, $W \rightarrow Ds\ \gamma \rightarrow KK\pi\ \gamma$, $W \rightarrow \pi\ \gamma$, $H \rightarrow \phi\phi \rightarrow 4K$, $Bs \rightarrow \tau\tau$. **Multiple soft objects** inspired from arXiv:1902.05535 (new scalars with long decay chains). **Resonances in $\mu\mu$** (in progress), **ee, $\gamma\gamma$, $\tau\tau$** (students identified)

- **On going:** procurement of hardware in good shape and work on firmware & software for demonstrator system in in B40 DAQ Lab started
  - Representative algorithms are being tested on different architectured and benchmarked on the available hardware
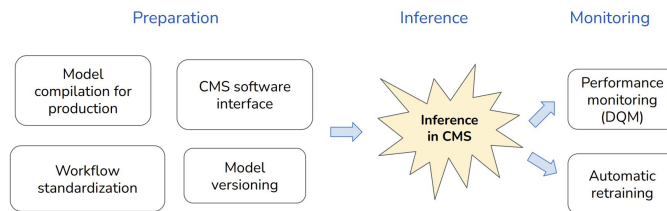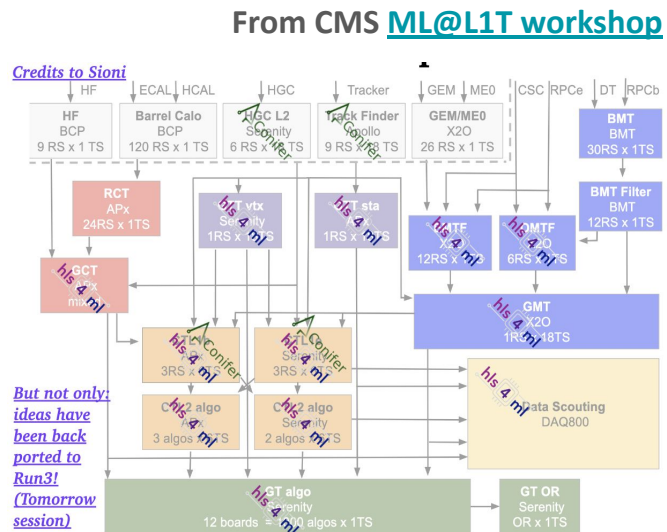


**Versal Premium Series VPK180 Evaluation Kit**
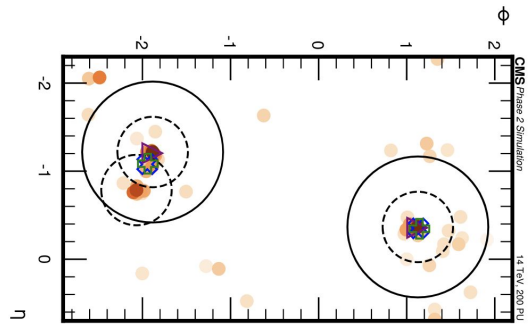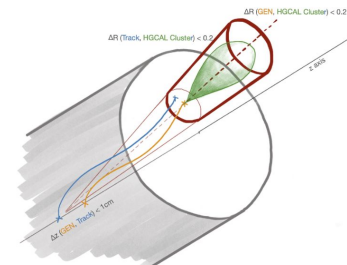
**Versal HBM Series VHK158 Evaluation Kit**

- **NGT team so far:** Cristina Botta, Valentina Camagni, Gianluca Cerminara, Emilio Meschi, Łukasz Michalski, Matteo Migliorini, Giovanni Petrucciani, Leah-Louisa Sieder

# 3.6: Practical real-time AI for L1T/Scouting

- Develop and operate advanced and fast **ML algorithms for improving the L1T reconstruction (used by standard triggers, and L1 Scouting)**
  - Into FPGAs of **Global Trigger, Correlator and Global Track Trigger subsystems**
  - This is already happening, but the project will **boost it** profiting from experience on advanced ML methods
- Develop practices for training & deployment of ML algorithms in L1T **(MLOps)**
  - updating & redeploying algorithms for changes in detector conditions
  - tracking, archiving & retrieving ML models (e.g. for consistent emulation)
  - gain 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 (implemented in different subsystems, only limited amount of information is propagated between subsystems)

**From CMS ML@L1T workshop**



Nice **talk from Dylan** on Edge ML at on-going SMARTHEP School **11**

# 3.6: Practical real-time AI for L1T/Scouting

- **On going:** improving baseline algorithms to expand baseline acceptance
  - Multiclass discriminator to improve and harmonize IDs for e/γ, pions, pile-up separation for PF reconstruction: **goal increase efficiency at low pT for e/γ and missing energy reconstruction**
  - Multi-class jet taggers **to improve current b and tau tagging**

- **On going:** develop prototype MLOps training and model tracking for correlator ML algorithms
  - Gaining experience with Run-3 (Anomaly Detection models recently deployed)

- **NGT team so far: Cristina Botta, Chris Brown, Gianluca Cerminara, Duc Hoang, Leon Joel Kerner, Kyungmin Park, Stella Scheafer, Sioni Summers**





arXiv:2402.01876

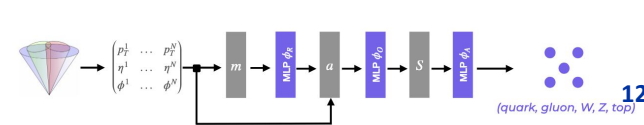a) Multilayer Perceptron MLP

*(quark, gluon, W, Z, top)*

b) Deep Sets DS

*(quark, gluon, W, Z, top)*

c) Interaction Network IN

*(quark, gluon, W, Z, top)*

# 3.7: L1T Anomaly Detection and Data Compression

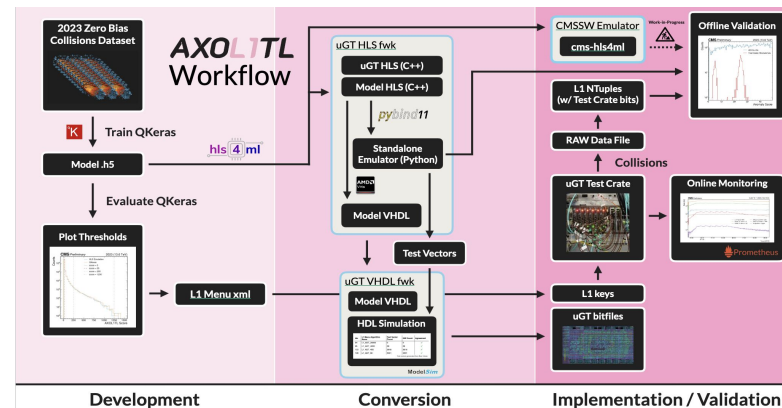- **Another way:** trigger on "anomalous" events using ML unsupervised models (sub-μs autoencoders!)
  - CMS recently established a **new trigger paradigm with sub-μs autoencoders** for anomaly detection (Abhijith's talk at ICHEP): AXOL1TL deployed in L1T Global Trigger



- **Next: Commissioning and validation of the Run 3 unsupervised AD algorithms**
  - Learn end-to-end algo integration and system operations
  - Optimize the model to be more robust against conditions (e.g. pileup)
- **Next: 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
- **Next: Develop new AD model suited for HL-LHC L1 trigger system using TPs/particle-based info**
  - Design an innovative GNN-based autoencoder for low latency and resources which will serve as baseline AD at HL-LHC
- **Next: Develop intelligent data compression and/or reduction for L1 Scouting system (to allow long-term storage of more scouting data)**

# 3.7: L1T Anomaly Detection and Data Compression

- **On going:** full training-to-deployment workflow has large number of steps
  - Model versioning, rate stability monitoring …
  - **CMSSW-hls4ml emulator workflow established to be followed by every future NN!**
- **On going:** improving baseline: more robust against PU and other system changes to make it usable for physics searches



- **NGT team so far: Diptarko Choudhury, Sabrina Giorgetti, Maciej Glowacki, Eric Moreno, Jennifer Ngadiuba, Sioni Summers**

# NGT @ HLT

Andrea Bocci (CERN), Silvio Donato (INFN-Pisa), Felice Pantaleo (CERN), Marco Rovere (CERN), Thiago Tomei (SPRACE)

# R³ Has ambitious goals

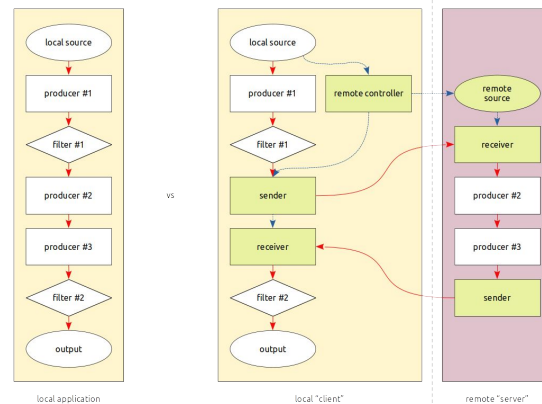- R³ aims to transform the HLT event reconstruction by developing a suite of algorithms that rethink the process entirely, rather than just speeding up existing ones. Depending on the level of speed-up required, innovative approaches will be applied as needed to meet live physics analysis requirements. Key efforts include optimizing data structures for accelerators like GPUs, redesigning CMSSW as a distributed application with minimal code impact, and leveraging high-speed interconnects to reduce latency.

- R³ will also reduce disk usage by compressing or simplifying data, and compute necessary conditions at HLT to match offline reconstruction quality, ensuring high physics performance with minimal disk space.

# 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**

40 MHz → L1T → 750 kHz → HLT → XYZ kHz → RAW

40 MHz → L1T → 750 kHz → HLT → XYZ kHz → RAW
HLT → 750 kHz → nano AOD

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

- **Ongoing**:
  - Start a systematic investigation of the reconstruction performance in different scenarios:
    i. HLT Phase2 reconstruction
    ii. Offline Phase2 reconstruction
    iii. Run3 Offline reconstruction
  - Critically compare the 3 scenarios, understand the differences between them, highlight the bottlenecks
  - The outcomes of this study will be the backbone of the EOY report
  - All measurements and considerations are available at this link
  - Start investigating the current status of the Monte Carlo truth information at large in CMS
    i. This is in view of the work scheduled for next year, 2025.

- **Team so far**: Marco Rovere, Jan Gerrit Schulz, Luca Ferragina, Marco Musich, Davide Valsecchi

# Task 3.1.2: Optimized data structures (*task leader F. Pantaleo*)

- 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

**Traditional layout**

| $x_0$ | $y_0$ | $z_0$ | $x_1$ | $y_1$ | $z_1$ | $x_2$ | $y_2$ | $z_2$ | $x_3$ | $y_3$ | $z_3$ |
|---|---|---|---|---|---|---|---|---|---|---|---|

| $x_0$ | $x_1$ | $x_2$ | $x_3$ | $y_0$ | $y_1$ | $y_2$ | $y_3$ | $z_0$ | $z_1$ | $z_2$ | $z_3$ |
|---|---|---|---|---|---|---|---|---|---|---|---|

**SoA layout**

- 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.

- **Ongoing**:
    - Discussion with CMSSW FW Core team ongoing
    - More details discussed during a dedicated meeting
    - Establishing a small test bench for demonstrating different prototypes

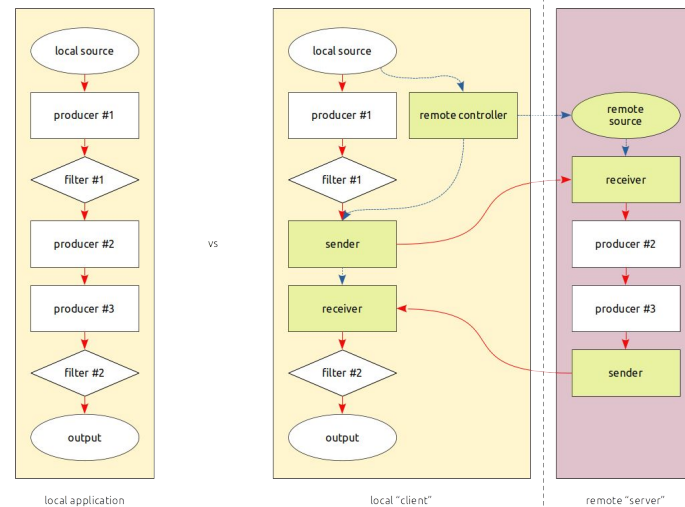- **Team so far**: Felice Pantaleo, Leonardo Beltrame

# Task 3.2: Evolving the CMSSW into a distributed application (*task leader A. Bocci*)

- The main goal is to **achieve more flexibility in the design of the HLT farm for Phase-2**:
  - independently *scale the amount of CPU and GPU* processing power;
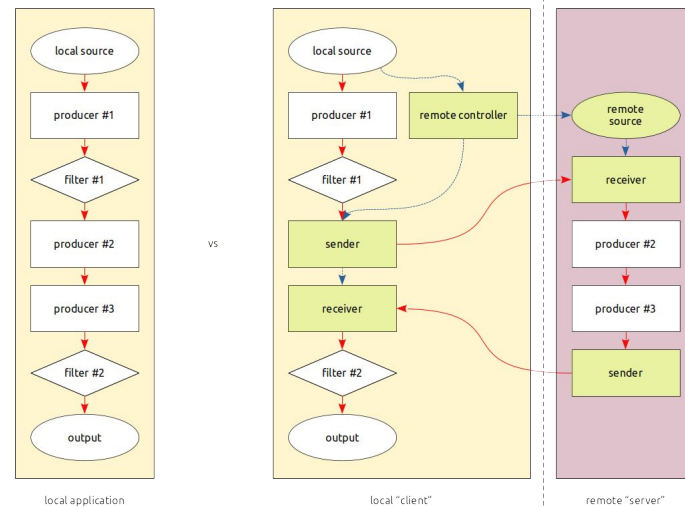  - support *different kind of accelerators*, like GPUs and FPGAs.

- The plan is to *extend CMSSW to a fully distributed application*, based on a loose set of requirements:
  - to support *arbitrary CMSSW configurations*
  - to *minimise the impact* and maintenance on the framework and reconstruction code:
    - → leverage the portability libraries and microarchitecture features
    - → do not rewrite the reconstruction code, or move it out of CMSSW
  - to *minimise inter-process traffic*:
    - → schedule entire sequences and tasks on a remote node
  - to *minimise latency* and overhead:
    - → leverage high speed, efficient interconnects
    - → RDMA protocols like InfiniBand and RoCE for direct network-to-memory and network-to-GPU copies

# Task 3.2: Evolving the CMSSW into a distributed application (*task leader A. Bocci*)
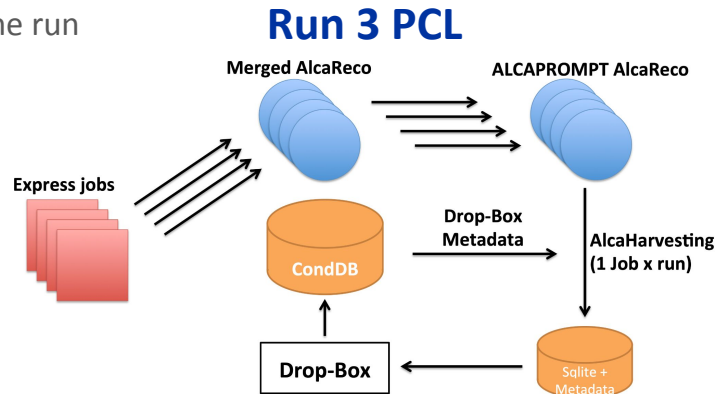
- The approach being considered to minimise the impact and maintenance burden is to leverage the Event-based interface of CMSSW:
  - wrap all communications in a small number of modules (a controller, a source, a sender, a receiver);
  - use the rest of CMSSW unchanged;
  - automatically support all kind of modules (legacy, alpaka-based, ML-based, …) that can run in CMSSW.
- The implementation of a client-server prototype is **in progress**.

- **Personpower**:
  - hiring of a doctoral student is in progress
  - continue the collaboration with non-CERN experts

- **Team so far**: Andrea Bocci, Fawaz Alazemi

# Task 3.3: 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
    - physics-driven compression, replacing basic information with higher-level quantities
        - build upon the work done by the Heavy Ions group
        - 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

- **Ongoing**: evaluate the impact of RAW data compression and of their replacement with low-level reconstructed quantities (RAW') in view of the EOY report
- **Team so far**: Silvio Donato, Simone Rossi Tisbeni

# Task 3.4: Optimal Calibration at HLT (*task leader T. Tomei*)

- Current status in Run 3:
  - Prompt offline reconstruction done with **Prompt Calibration Loop** (PCL) calibrations
  - Calibrations for HLT reconstructions run on different schedules instead (weekly, per-fill, …)
- Reminder of **Prompt Calibration Loop** (PCL):
  - Express stream (100 Hz) → calibration → prompt RECO of the run
  - Time window: 48h
  - Workflows:
    - Beamspot
    - SiPixel: bad components, alignment, Lorentz Angle
    - SiStrip: bad components, efficiency, gains
    - ECAL: pedestals
    - PPS: alignment, offset, timing
- **Small-scale prototype at the HLT**, explore tradeoffs:
  - Buffer RAW data for *N* hours → Latency vs. how much data to hold
  - Derive improved calibrations → Express rate vs. latency vs. accuracy
  - Inject and use them for online reconstruction of HLT scouting



**Run 3 PCL**

23

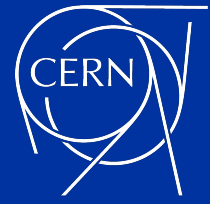# Task 3.4: Optimal Calibration at HLT (*task leader T. Tomei*)

- **Ongoing**:
  - documenting all the alignment and calibration workflows in CMS and trying to ascertain the more important ones, to focus on them.
    - ECAL and Tracker already addressed/ongoing
    - HCAL,  Muons, and all Physics Objects will follow

- **Team so far**: <u>Thiago Tomei</u>, Mateuzs Zaruki, Jessica Prendi, Marco Musich

# Welcome to the CMS NGT team!

cern-cms-ngt-l1t@cern.ch, cern-cms-ngt-hlt@cern.ch