

Data Monte Carlo Preparation in CMS

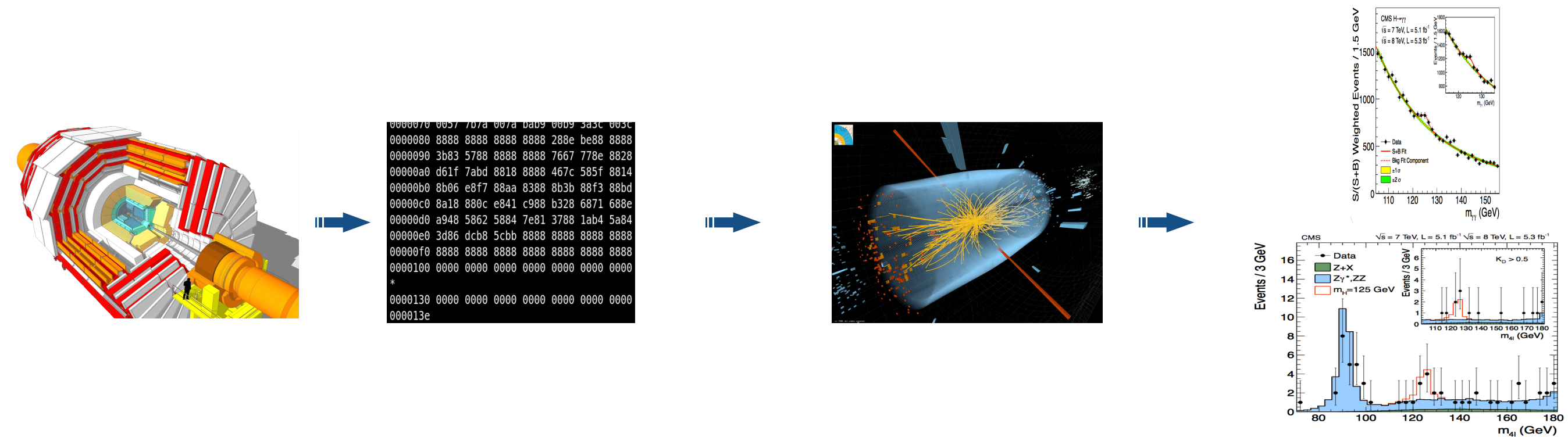


Gurpreet Singh Chahal (IPPP Durham University and Imperial College)

IRN Terascale, Durham, 7 Sept 2018

Outline

- Journey from data recording to data analysis
 - ▶ Data collection, and data quality monitoring
 - ▶ Data certification, and Monte Carlo production
 - ▶ CMS software development, and the CMS computing resources



Coordination of Data Preparation

Two main CMS coordination areas:

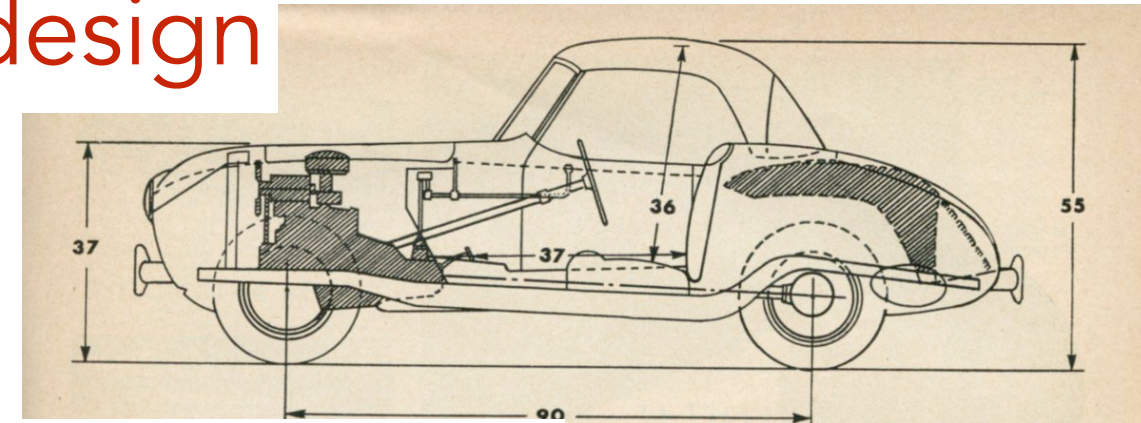
- **Offline & Computing**

- CMSSW software development, event reconstruction and simulation
- data processing and Simulated events generation, events storage and management

- **Physics Performance and Datasets (PPD)**

- Bridge between detector/physics groups and offline&computing
- data quality & certification
- alignment & calibrations
- software validation
- management and production of the Monte Carlo samples
- organisation and configuration of datasets and data processing

design



manufacturing



product industrialisation

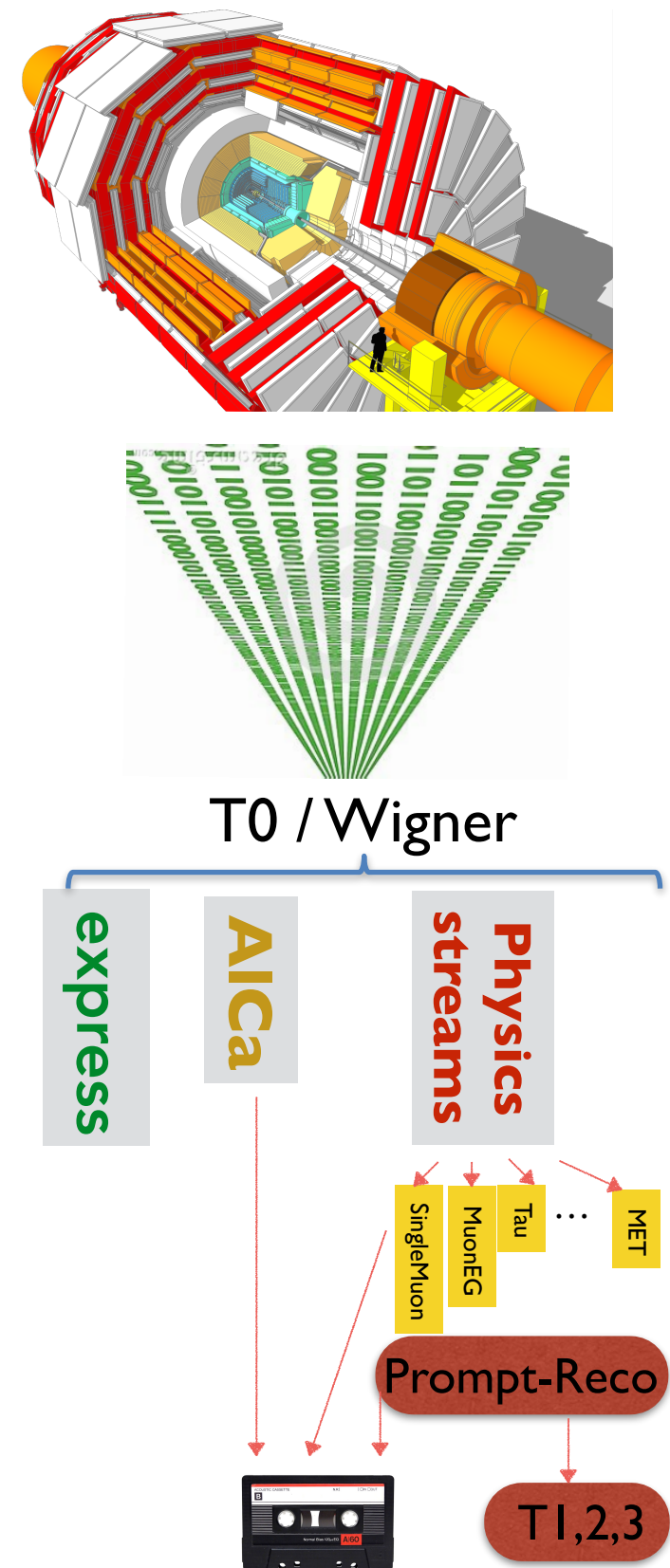




CMS Detector to Offline Platform

The collected events reach the Tier0 farm at CERN for tape archival, organisation and processing. Different **streams** and workflows for different **use-cases**

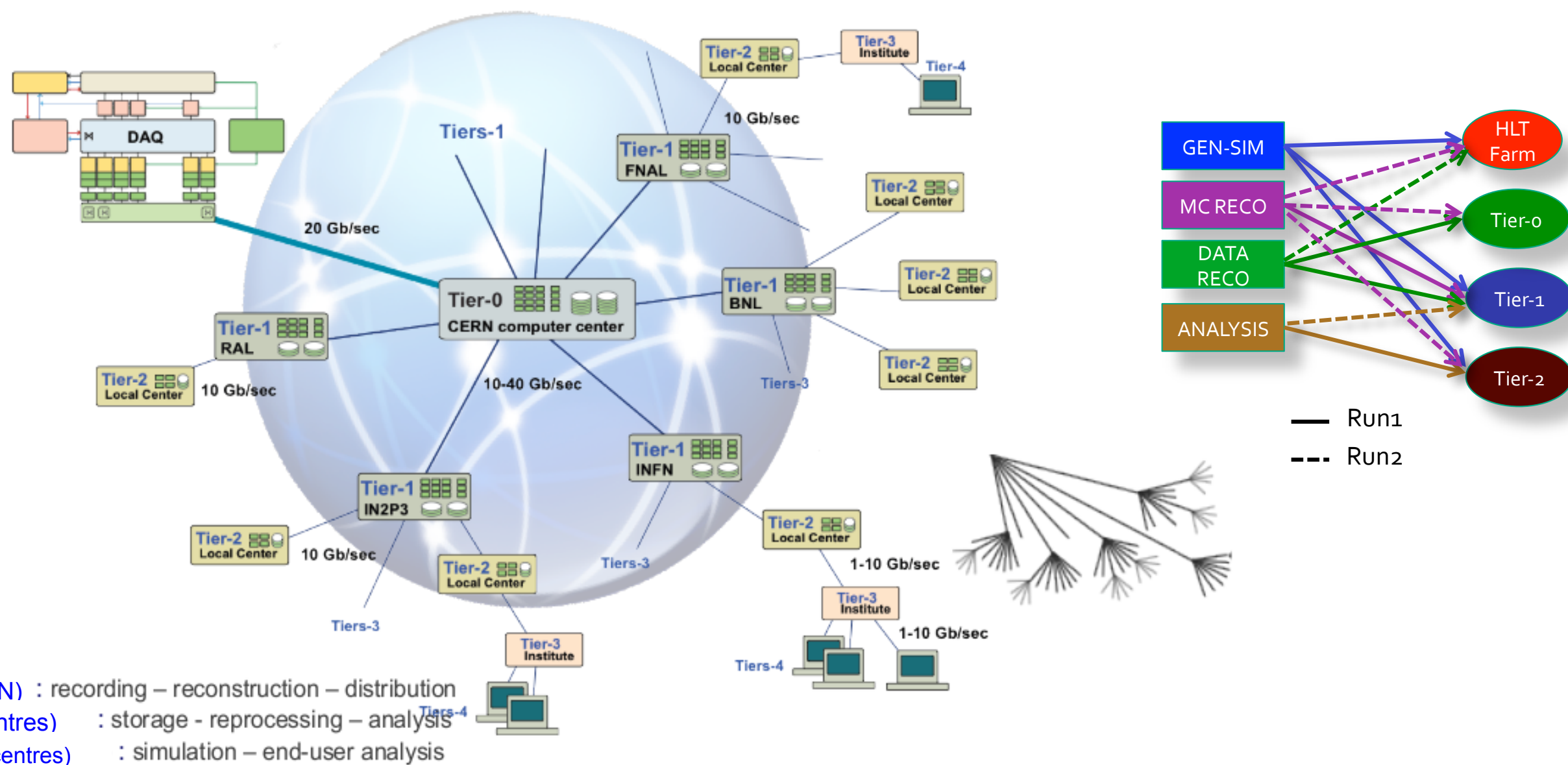
- **express**: available ~2h after data collection for prompt feedback & calibrations:
 - calibration (1/2) - detector (1/4) - physics (1/4) monitoring
- **Alignment & Calibration (AlCa)** streams
 - dedicated event selection & event content devised for calibration purposes
- **Physics Streams**: split into **primary datasets (PDs)** on the basis of the HLT results, and PDs are promptly reconstructed
 - **Store events with related topology** in the same PD to ease consumption;
 - **limit replication** of events (PDs overlap)
- Constraints from analysis:
 - definition centred on physics objects (e.g. SingleMuon, MuonEG...)
- Constraints from processing and handling:
 - approximately uniform average event rate across different PDs to ease distribution at the Tier2 centres
 - event rate > 10 Hz, to avoid small files & < 200 Hz
- On top of the PDs we can deploy "central skims" → customised event content + rate reduction using also RECO quantities
 - used for Detector Studies (DPG) or Physics Analysis Groups (POG-PAG)
- Other specialised streams (e.g "data parking", "data scouting"...)





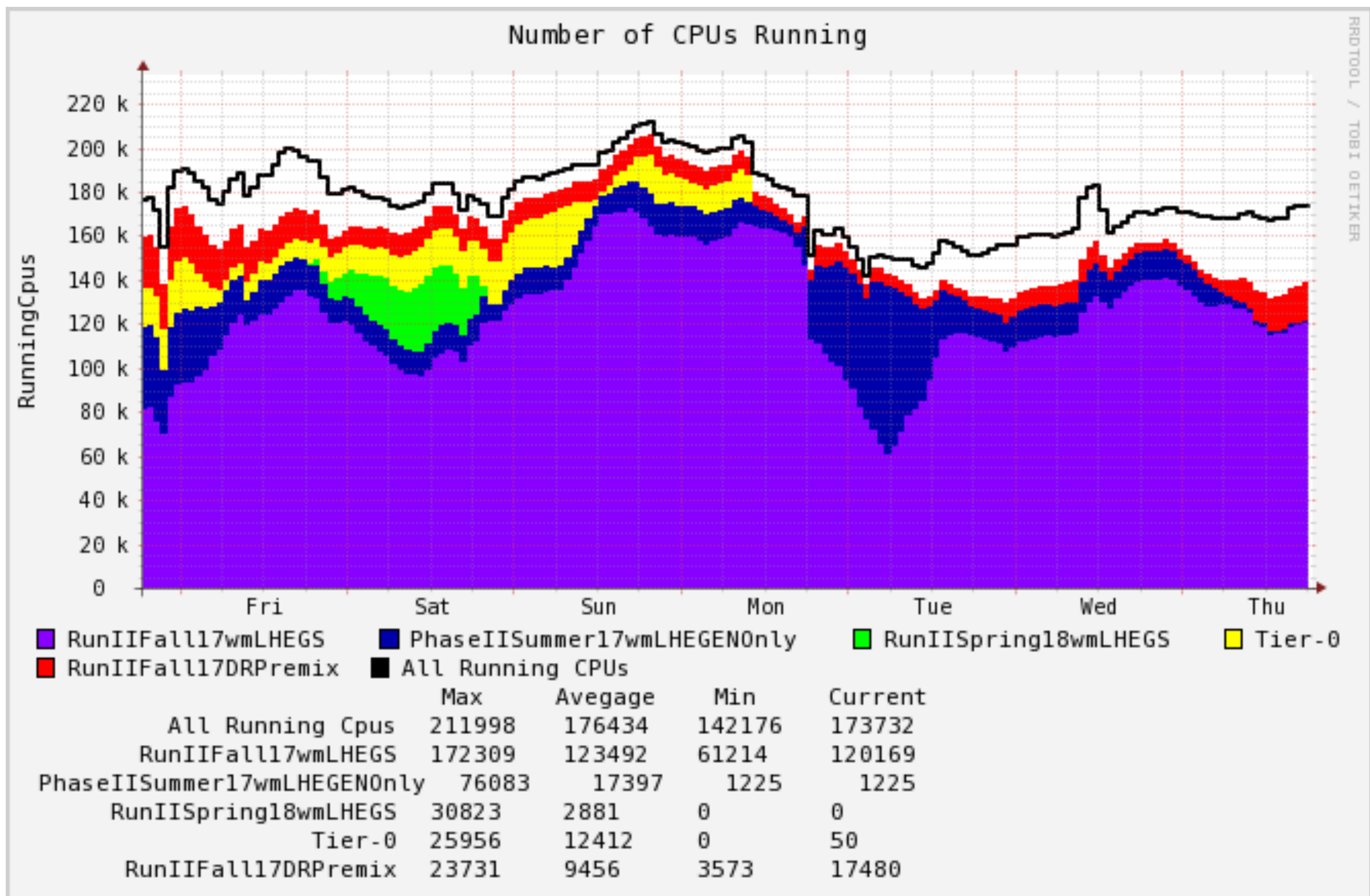
Worldwide LHC Computing Grid (WLCG)

- The distributed computing model in CMS is motivated by various factors:
 - The large quantity of data and computing requirement encouraged distributed resources
 - Ability to leverage resources at labs and university: Hardware, expertise, infrastructure.
 - Benefits of providing local control of some resources, and ability to secure local funding sources.
- ~30% of the resources are located at CERN, ~30% at T1s, and ~40% at T2s, (Total ~250 K cores)
 - Relies on the development of tools to make transparent access to the resources and efficient distributed
- Can only be successful with sufficient networking between facilities.
 - Availability of high performance networks has made the distributed model feasible.





Use of Computing Resources





CMS Software

- CMSSW: one release to rule them all
 - **GEN**erator, **SIM**ulation, **RECO**nstruction **ANALYSIS** workflows...
- C++ code and configuration handled via Python
 - “git” used for code versioning and integration
- Release schedule follows a “train model”:
dear developer: catch this train or wait for the next one
 - regular time-table of ~6 months (slightly tuned for major conferences or physics needs)
 - pre-releases are regularly produced while the release is under development
- Feature planning:
 - production releases: driven by physics/machine constraints & goals
 - e.g. Fall17 production is producing MC compatible to 2017 data

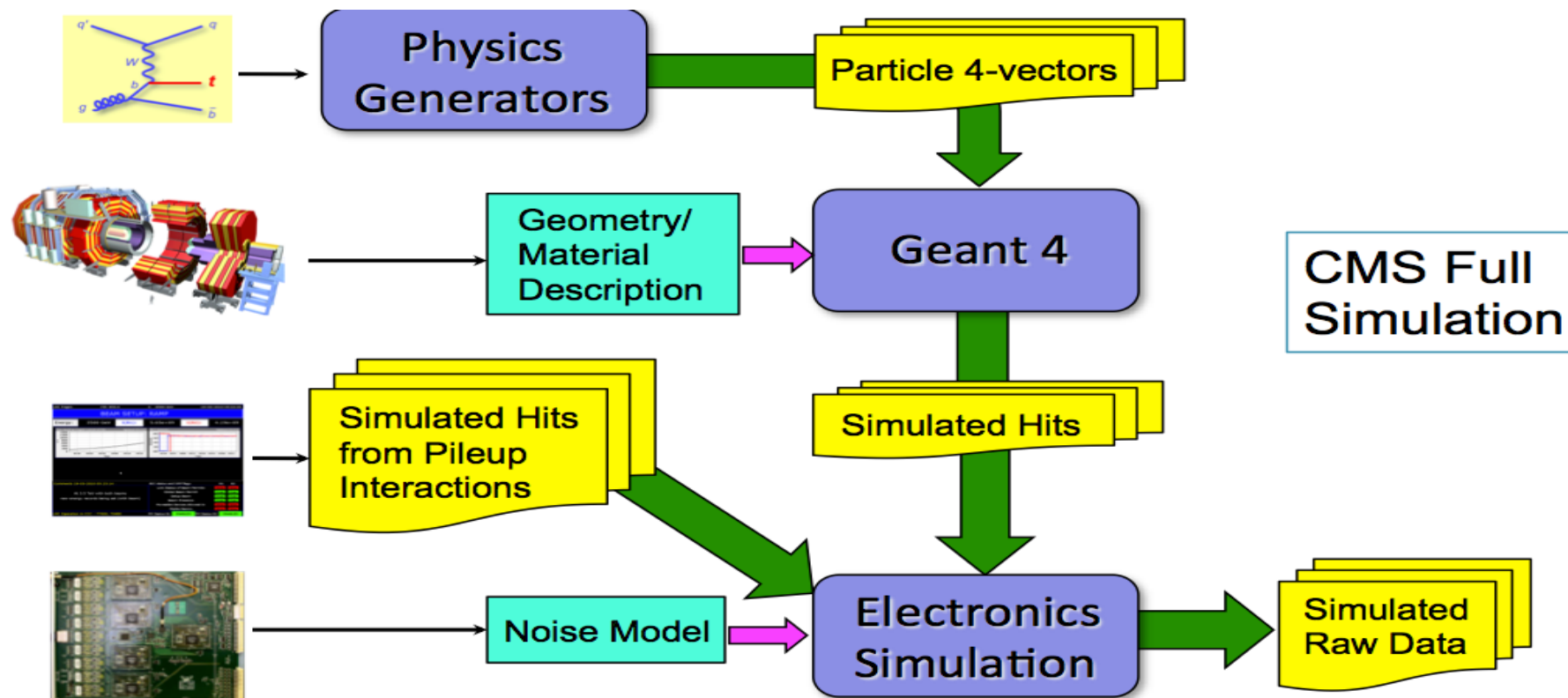
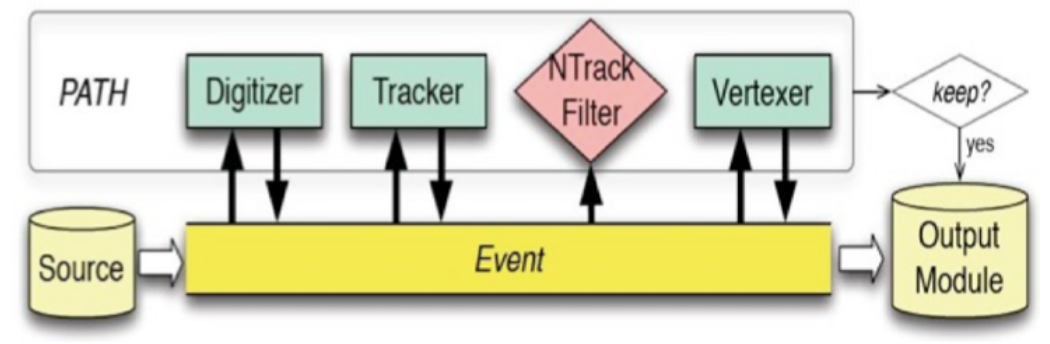




Simulation Flow

CMSSW simulation sequence:

- event simulation algorithms are implemented as “modules” communicating via the “Event”
- The Simulation sequence aims at producing MC truth + RAW data as if it came from CMS@P5



Alternatively: CMS “Fast Simulation” is a slightly less realistic but much faster simulation of low-level objects (hits, clusters)

Data tier

GEN, LHE

GEN-SIM

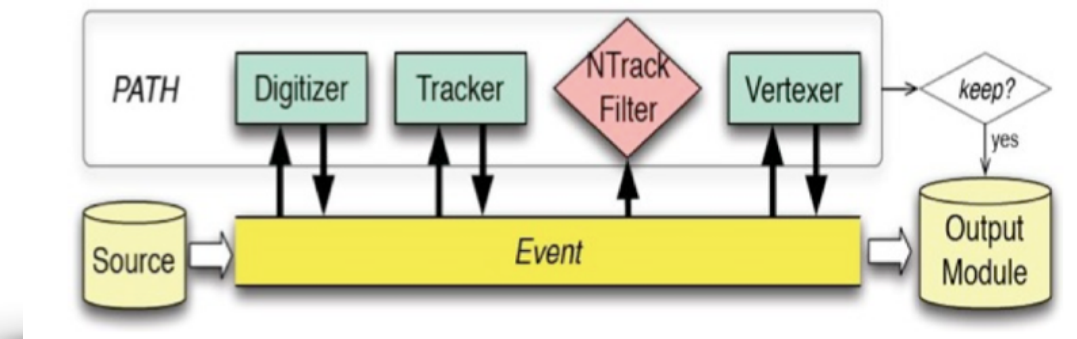
GEN-SIM-RAW



Reconstruction Flow

CMSSW reconstruction sequence:

- event reconstruction algorithms are implemented as “modules” communicating via the “Event”
- The Reconstruction sequence turns the binary output (RAW) from CMS/DIGI into **physically interpretable quantities** ready for data analysis
- Hits in the detector are aggregated in cluster and tracks, which in turn are matched to create **particle candidates** (PFAlgo): Tracks, muons, electrons, photons, jets ...

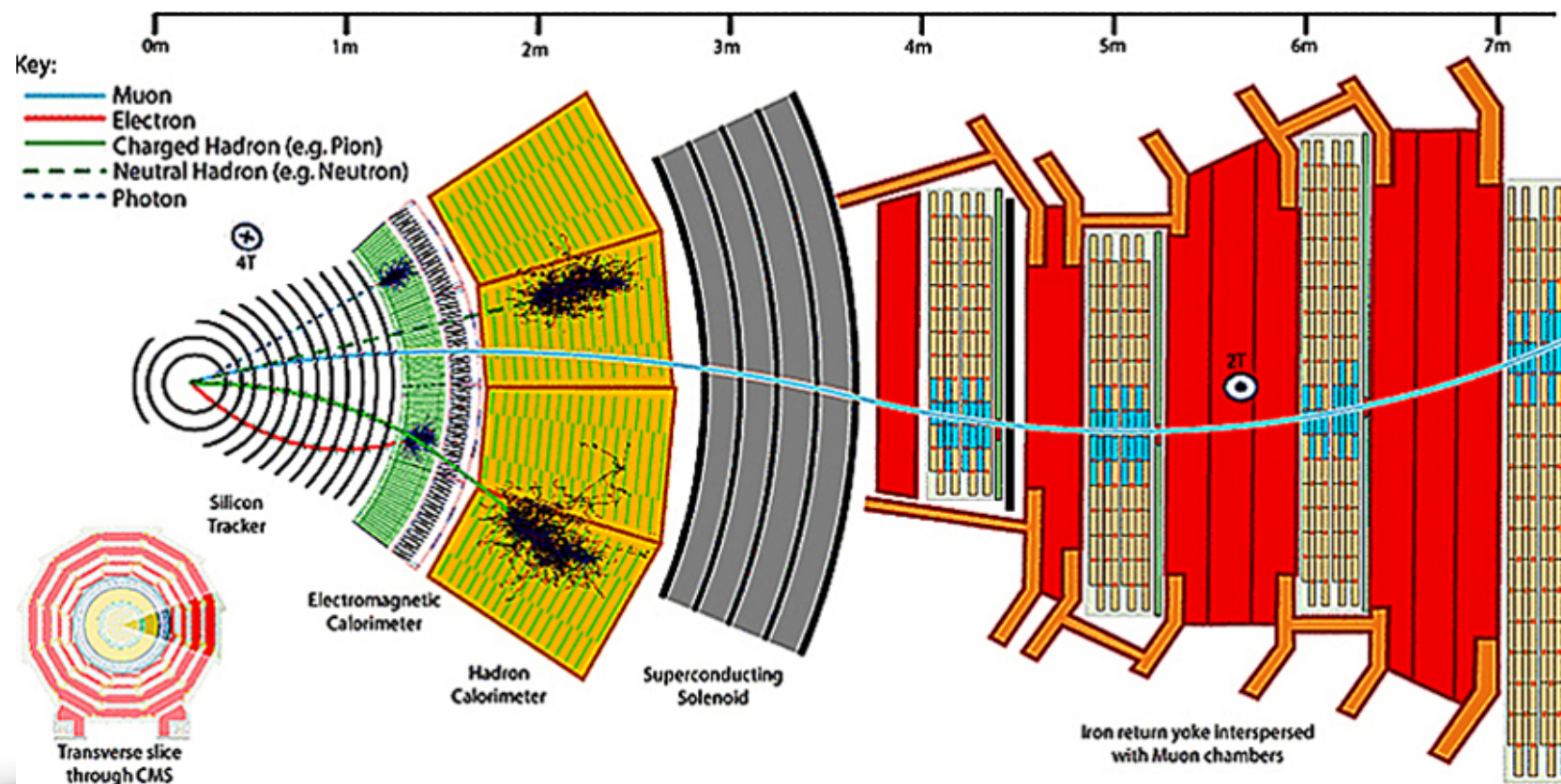


Data tier

RECO,AOD

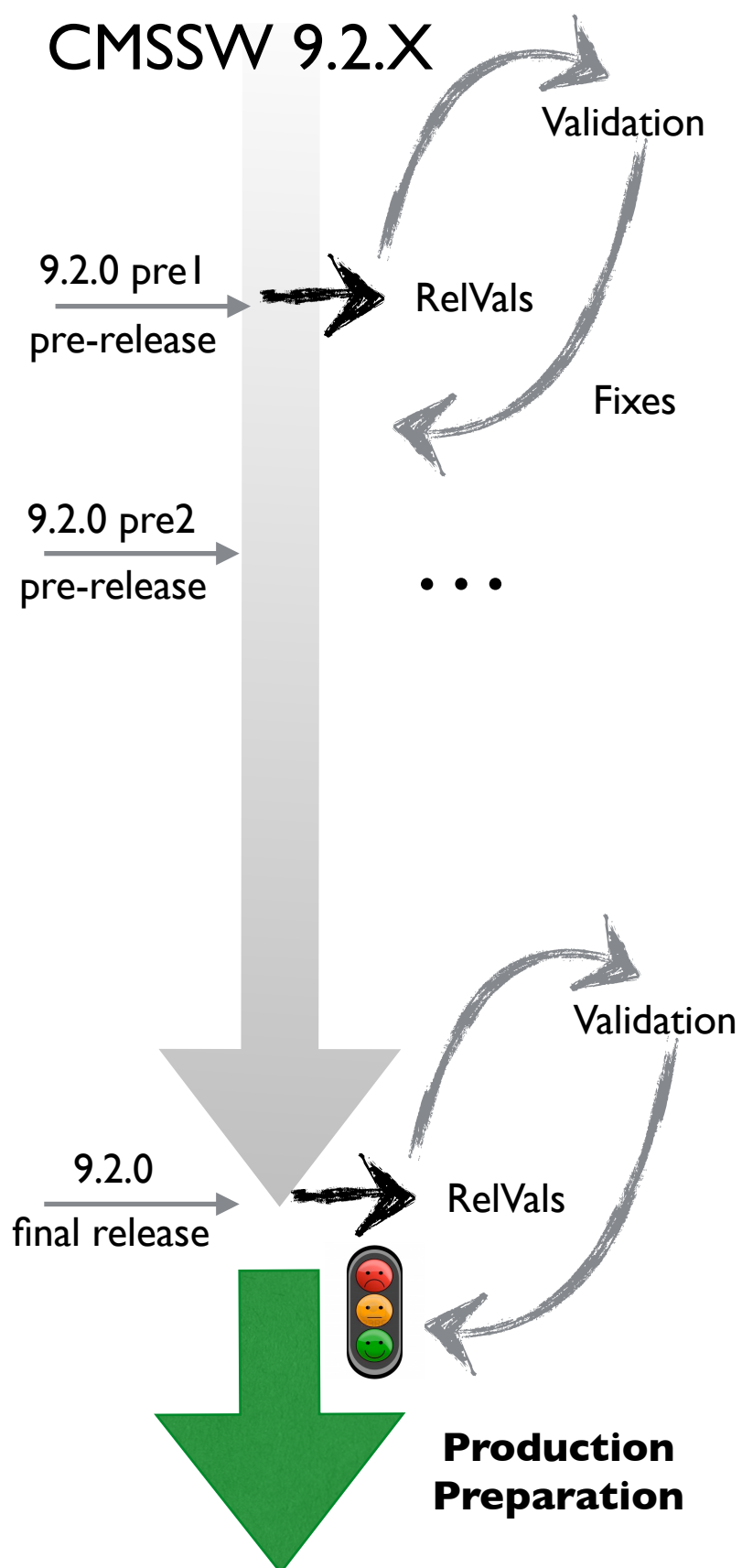
MINIAOD

MINIAODSIM





CMSSW release validation

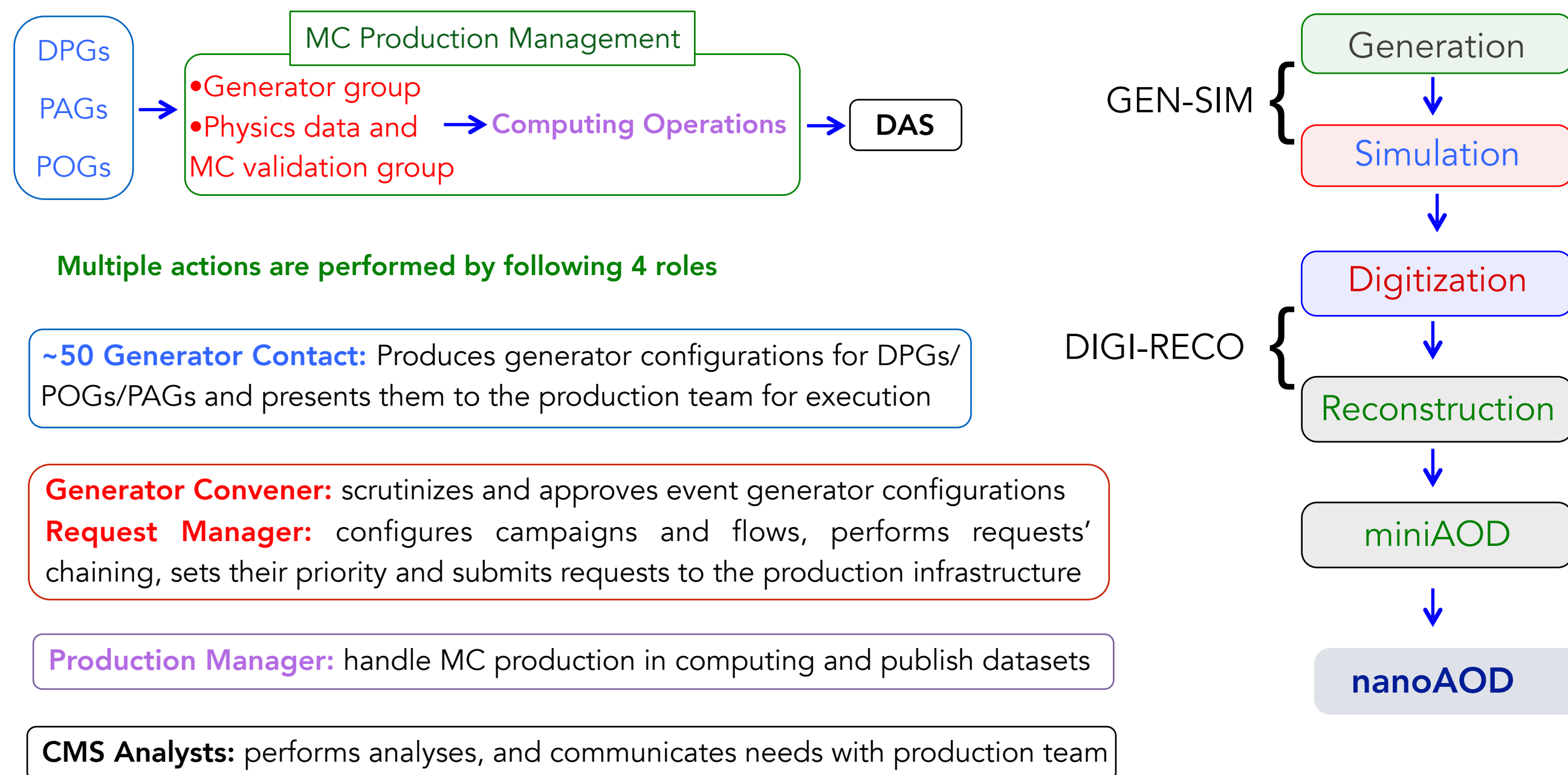


- Release integration bound to Quality Assurance Tests → Data Quality Monitoring (DQM)
 - unit tests & regression tests
 - small scale production tests: **Release Validation Test (RelVal)** producing DQM plots
- Once a major release (X.Y.Z) is green-lighted → start **preparation of the campaign** (re-reconstruction or MC production)
 - finalisation of the alignment and calibration conditions (and their validation)
 - finalisation of the parameters for the Pile-Up overlay (PU scenario)
 - preparation of the injection machinery for the central processing by computing



Overview of MC production

- ~20 groups: physics analyses (PAGs), detector performance (DPGs), physics object studies (POGs)
- 100s of physics analyses; 1000s of MC samples needed; Billions of events produced (Run-I+ Run-II+Phase-II upgrade)
 - Over 20 Billion simulated events produced in 1 year in 2016-17
- Strong and efficient production infrastructure required: bookkeeping and interface to computing resources

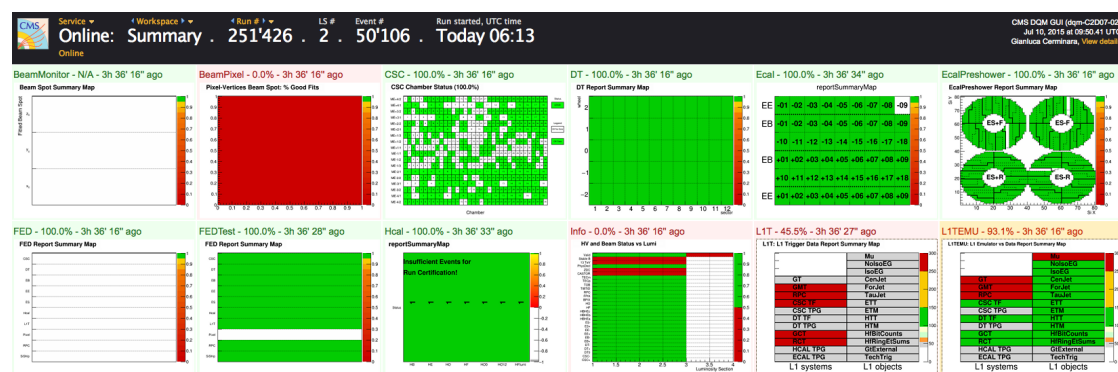
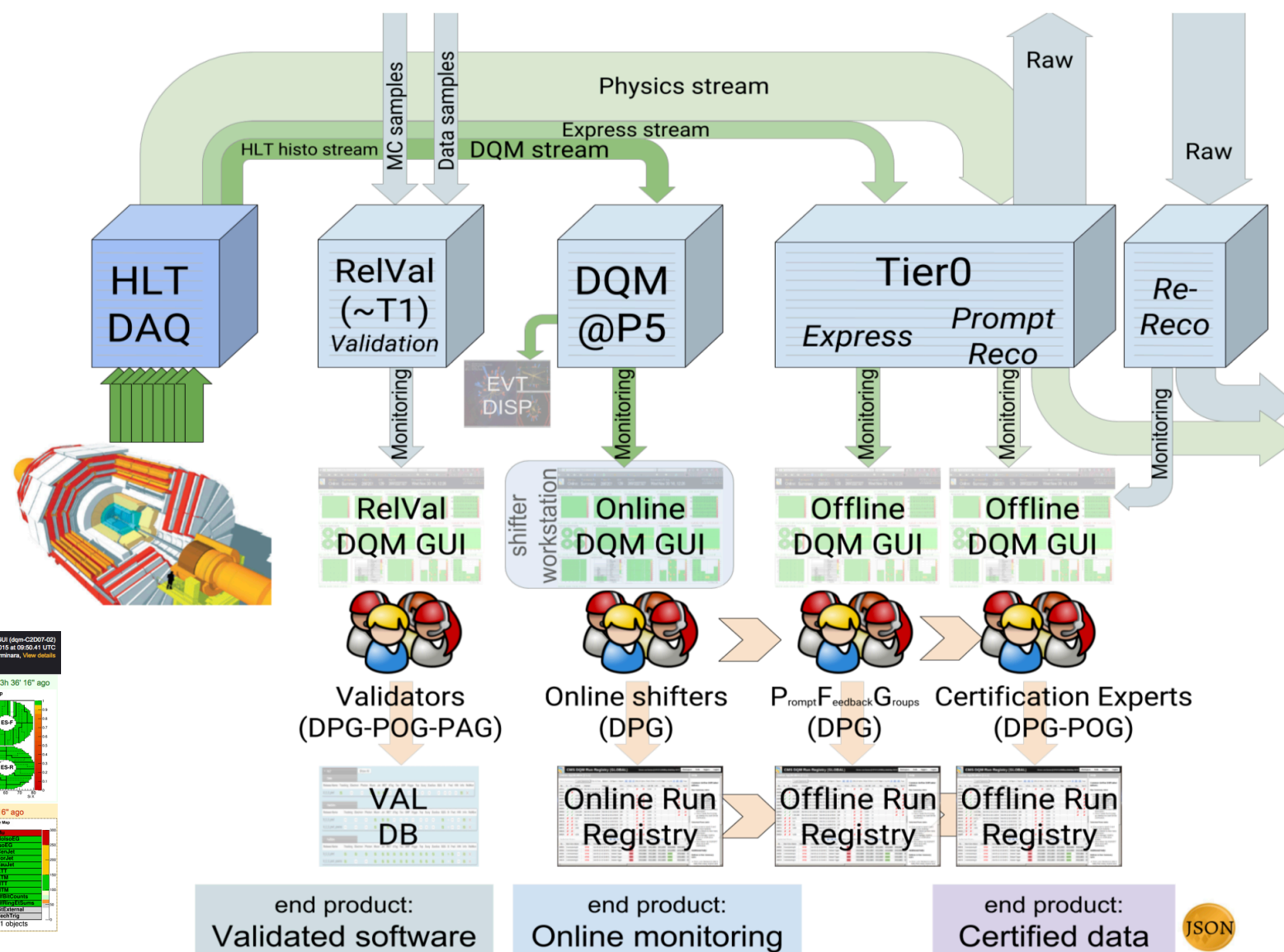




Data Quality Monitoring (DQM)

- DQM is the tool to produce plots while running RECO (or any CMSSW workflow)
- 2 main applications:
 - **online:** samples events after HLT and plots quantities with very low latency
 - live monitoring of detector performance during data taking
 - **offline:** reads all events while they are reconstructed
 - data certification
 - release validation
- **DQM GUI** → front-end to browse histograms for a given dataset/run

DQM Workflow

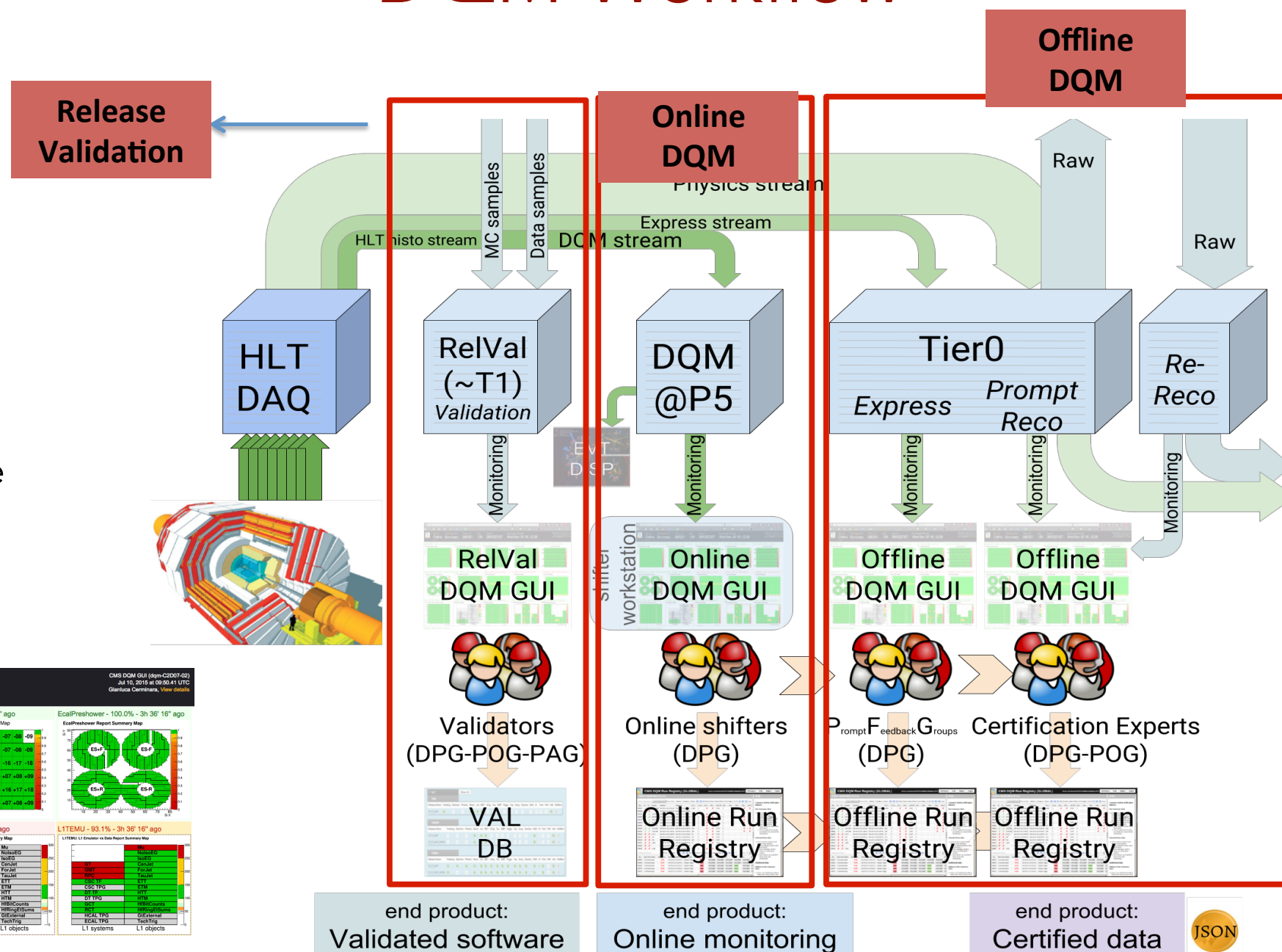




Data Quality Monitoring (DQM)

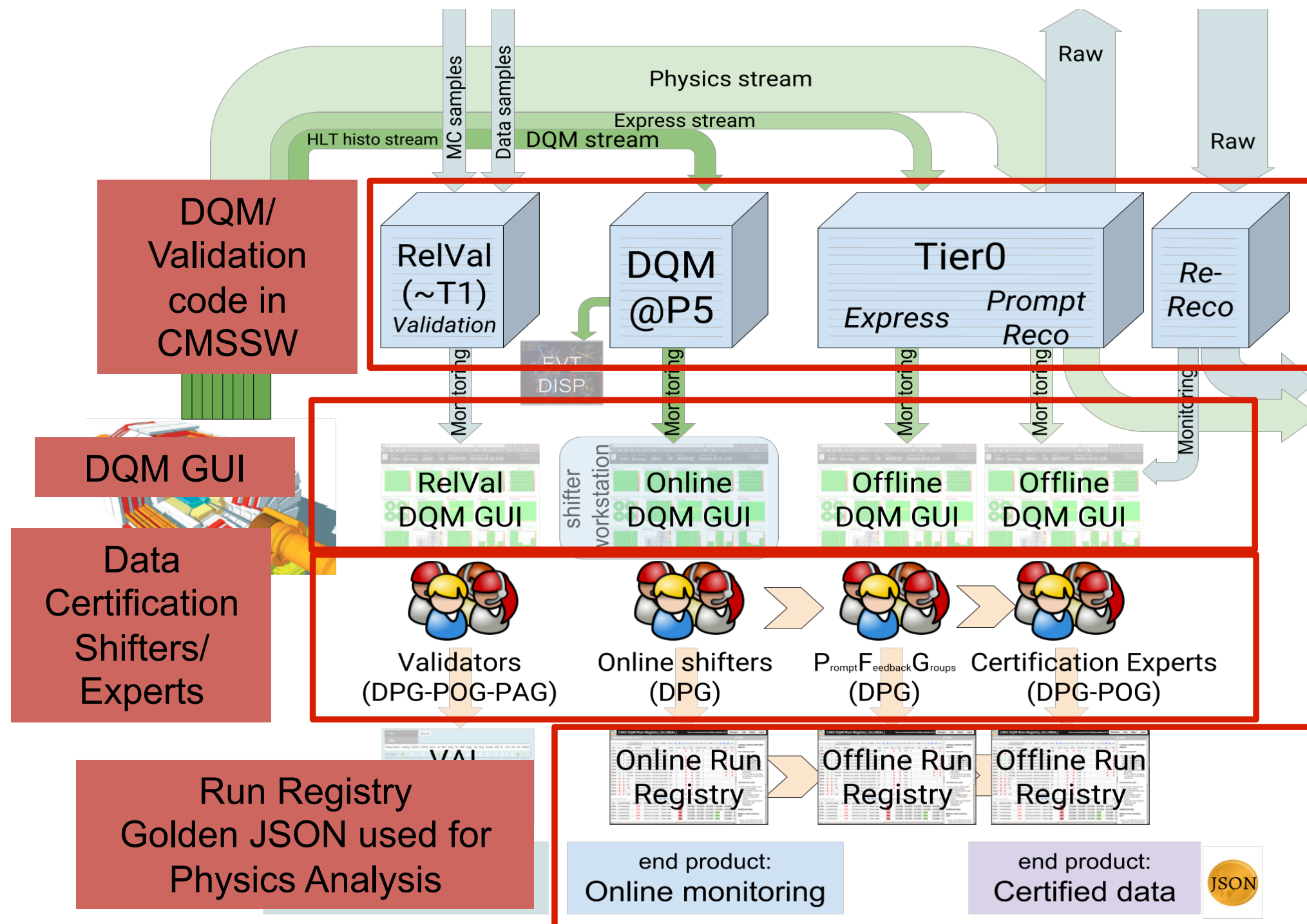
- DQM is the tool to produce plots while running RECO (or any CMSSW workflow)
- 2 main applications:
 - **online:** samples events after HLT and plots quantities with very low latency
 - live monitoring of detector performance during data taking
 - **offline:** reads all events while they are reconstructed
 - data certification
 - release validation
- **DQM GUI** → front-end to browse histograms for a given dataset/run

DQM Workflow





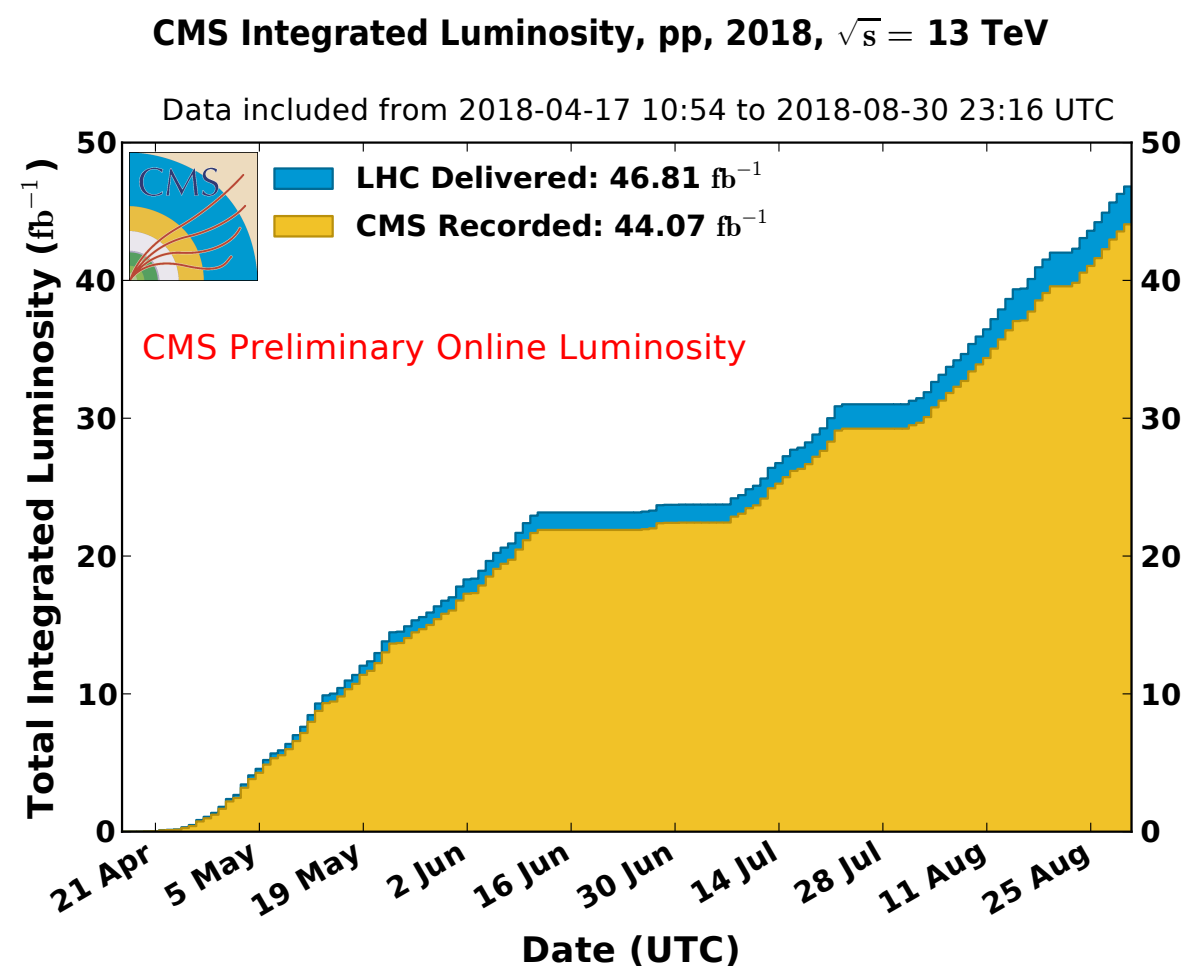
DQM Workflow and Data Certification





Certification & Analysis

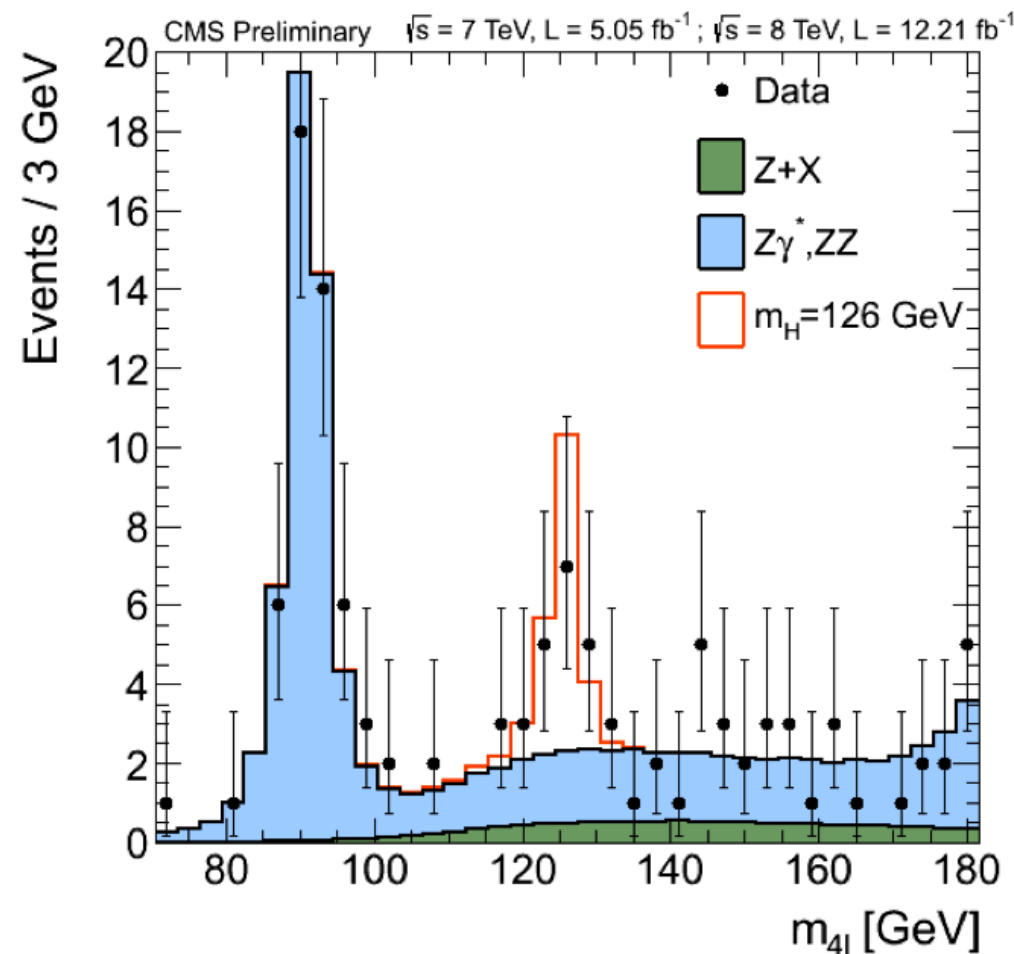
- Selection of LumiSections (LS) (≈ 23 s of run) considered GOOD for physics
 - distributed in JSON format
 - weekly for PromptReco
 - after each major re-reco pass
- Several “flavors” of the JSON file, like:
 - **golden** → requires all sub-detectors/POGs to be “GOOD”
 - **muon-only** → no requirements on calorimeters
- How do I use the JSON file:
 - to be used in CRAB to run only on CERTIFIED LSs of your dataset



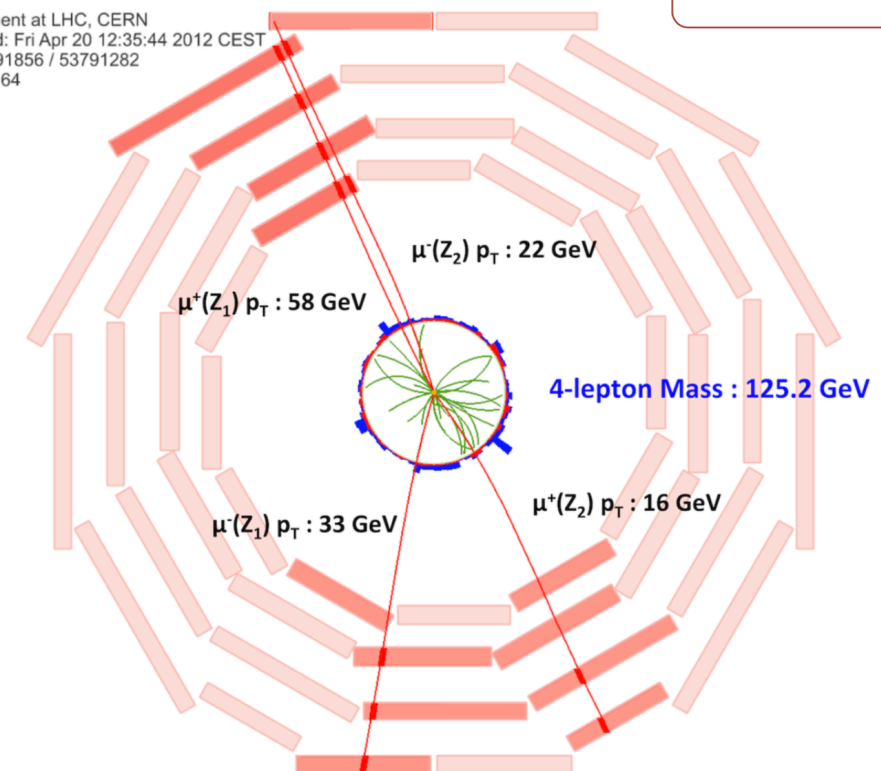


Running on the Data/MC

- Once you discovered/moved all the datasets you need to run and you have information about release, global tags, good data, etc. You can do analysis → CRAB (**C**ms **R**emote **A**nalysis **B**uilder)



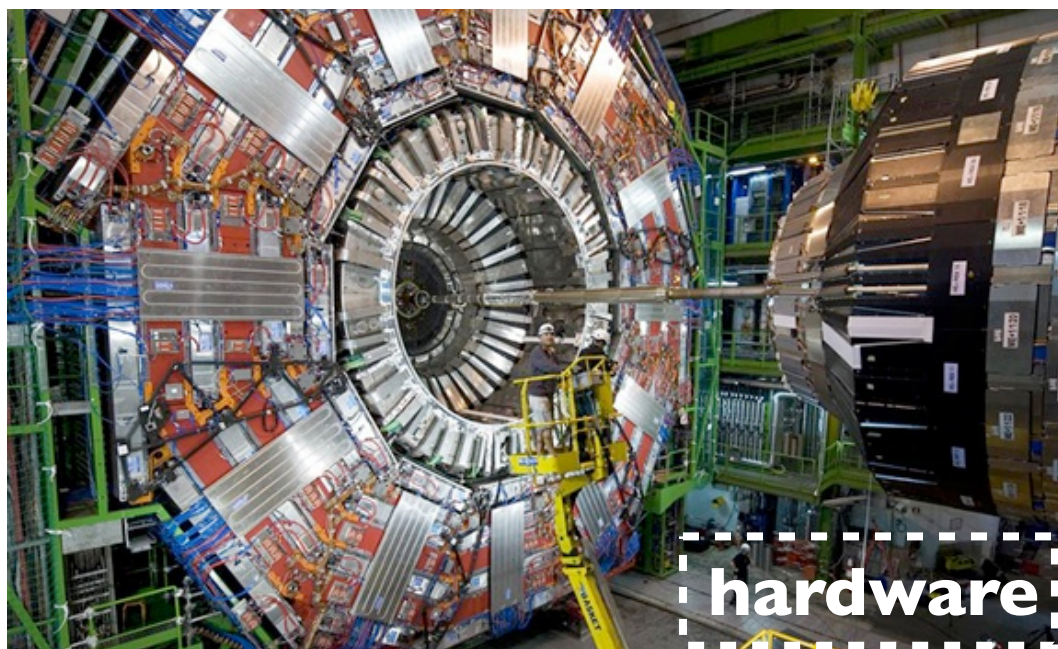
CMS Experiment at LHC, CERN
Data recorded: Fri Apr 20 12:35:44 2012 CEST
Run/Event: 191856 / 53791282
Lumi section: 64



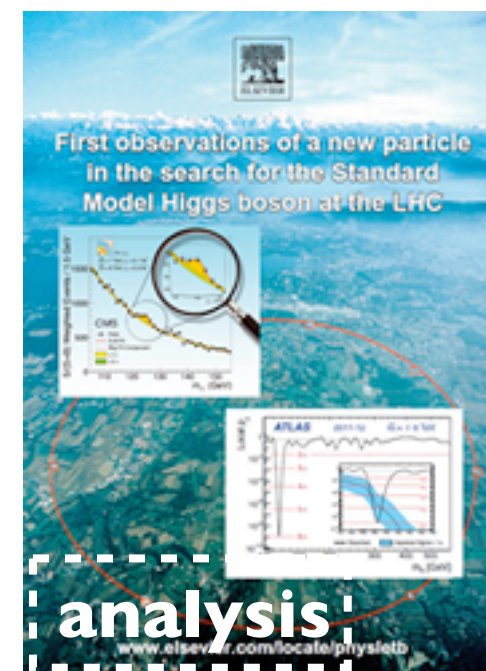


Outlook

- (some) Physicists tend to think that CMS is just



&



- This is **almost true!** However to exploit the hardware at its best and publish top quality analysis you need more:
 - software, data preparation and computing reached in CMS an unprecedented level of complexity
 - In ideal conditions, 4 billion MC events/month
- The main goal is to get high quality data, MC and software
 - Few hundred physicists are always working to achieve these goals
 - CMS is collaborating with industry (IBM...) to use the machine learning and deep learning techniques to improve the data quality monitoring and data certification



Thank you!