# Online data handling and storage at the CMS experiment
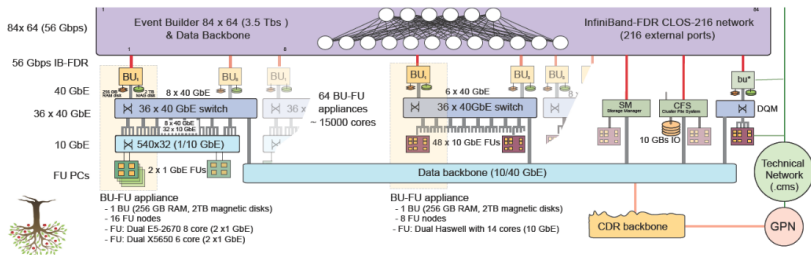
Lavinia Darlea, on behalf of CMS DAQ Group



MIT/DAQ CMS

April 10, 2015

# Introduction



Storage and Transfer System in the DAQ chain

- input: end of the DAQ chain, as described by Emilio* in the previous talk
- last part of the data flow: ensure safe storage and transfer to Tier0

*E. Meschi, File-based data flow in the CMS Filter Farm, CHEP 2015

# Storage and Transfer System Role

## 3 Implementation Stages

- merge the filter units output as to obtain 1 data and 1 metadata file/LS/stream
- buffer the data should the connection with Tier0 at CERN be lost
- copy the final files according to their intended destination:
    - Tier0 for the main data streams
    - various sub–detectors for online consumption: DQM, EventDisplay
    - store locally for local calibration of various sub–detectors
- ensure hand–shake with Tier0 for proper accountability

# Storage and Transfer System Requirements

## Merger System

- "merge" data at the BU level such as to obtain 1 file/BU/LS/Stream (mini–merger)
- centralize and merge all the BU outputs such as to obtain 1 file/LS/Stream (macro–merger)
- latency: a maximum of 2LS (1LS = 23s) delay in the macro–merger is considered acceptable
- provide input for the online monitoring system − 1 additional metadata file per data file (see next talk by Srecko*)
- not only "concatenate", but deal with special files, such as histograms and jsn files

*S. Morovic, A scalable monitoring for the CMS Filter Farm based on elasticsearch, CHEP 2015

## Storage and Transfer

- buffer a minimum of 3 days of continuous running (estimated 250TB)

- aggregated SM input from the 62 BUs is expected to reach a maximum of 2GB/s – mini merger write to LFS (Lustre FS)

- the macro–merger needs to consume this data online (2GB/s read the fragments, 2GB/s write the final merged file): 4GB/s(*)

- the transfer system is expected to transfer most of the data to Tier0 at 1GB/s

- overall: LFS needs to guarantee a total of sustained 7GB/s parallel r/w

## 2 available options

- "A"dditive
  - mini–mergers write a file/BU/LS/Stream, macro–merger merges them and makes them available for the TS
  - easy debugging, reliable, "standard" logic
- "C"opyless
  - mini–mergers write in parallel in the final file, macro–merger checks for completion and makes it available for the transfer system
  - reduce the required bandwidth with 4GB/s, fast due to parallel writing in the same file, more sensitive to corruption

Lustre FS architecture

- MDT: E2724, 16 drives of 1TB in one volume group plus 8 hot spares
- 2 OST controllers E5560 + expansion shelves DE6600, each 60 disks of 2TB
- servers: 6 DELL R720



Front OST



Disk shelves

# Storage

## High Availability

- all devices are dual powered (normal and UPS)
- all servers configured in active/passive failover mode
- volumes repartition to provide full shelf failure redundancy
- LFS availability: 40GE and InfiniBand (56Gb) data networks



Volumes configuration

Commissioning Acceptance

Proven steady 10GB/s rate in
r/w mode



Merger emulation

Proven steady 7.5GB/s rate

## Inheritance and Progress

- use old transfer system as a base
- new features have been added
  - identify and set the final destination of each stream per run
  - new logic in the bookkeeping and hand–shake protocol between the CMS site and Tier0

LFS bandwidth benchmarking

## LFS Validation

- usable space of 350TB
- tests done with different number of BU units
- obvious non–linear behaviour with the number of BUs
- saturation is expected around 8.5GB/s

# Conclusion



Mergers monitoring sample



Mergers delays sample

## Mergers Validation

- stable behaviour in 3 months of cosmic runs
- general latencies within the requirements
- proven reliability and availability

# Conclusion

## Transfer System

- successfully upgraded to transfer DAQ2 merged files
- work in progress: benchmarking and central management

# Questions?