



# Software & Computing in ATLAS Trigger and Data- Acquisition

W.Vandelli  
CERN Experimental Physics Department/ADT

on behalf of  
ATLAS Collaboration

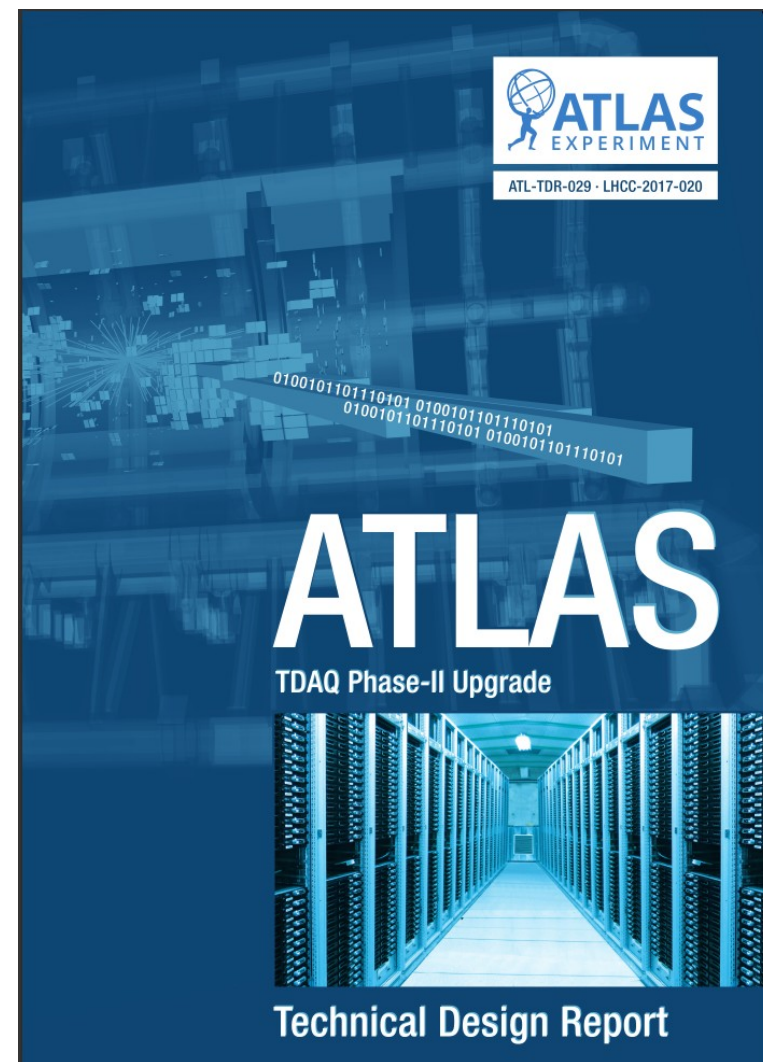
ATLAS TDAQ is a large mission-critical heterogeneous computing infrastructure

It operates in-house developed software and the DAQ team is responsible for all aspects from hardware procurement and installation to operation and upgrade

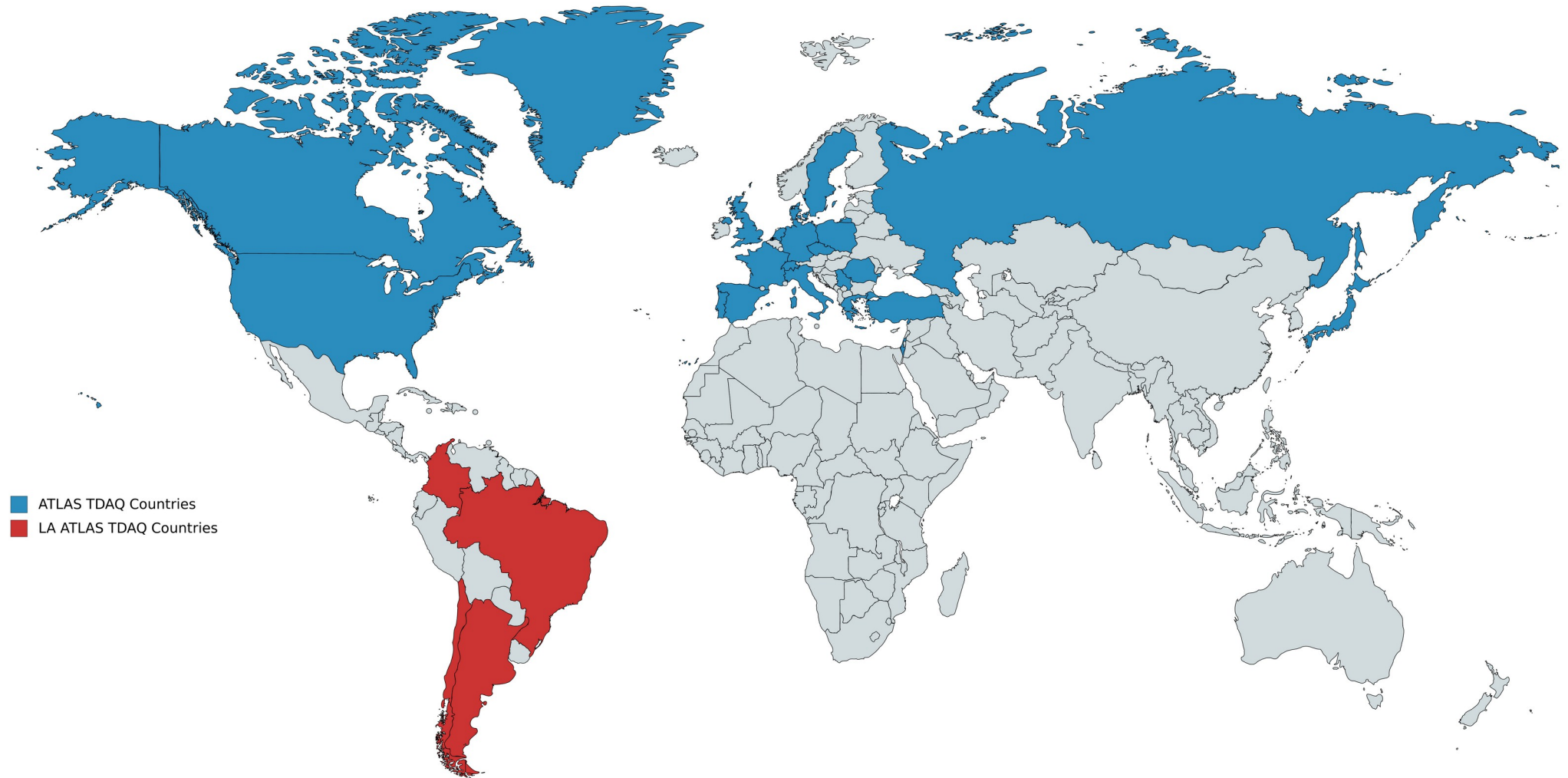
I will mainly focus on the challenges associated with the Phase-II/Run 4 upgrade (2024), even though Phase-I/Run 3 (2021) operation is not given

### Phase-I/Run 3 TDAQ

Processing Servers	~2000
Processing Appl.	~50000
Local Storage (TB)	~800
10 GbE Links	~500

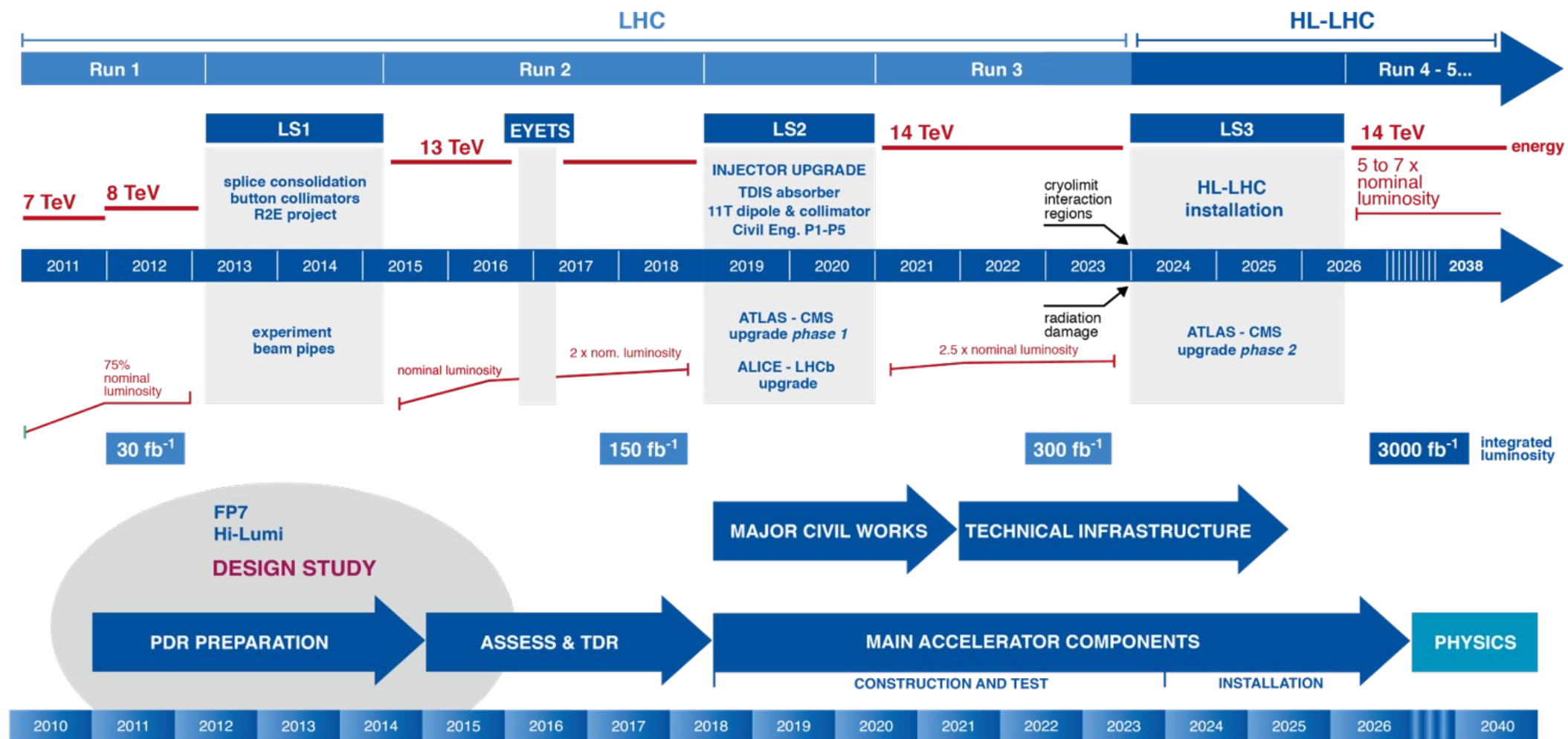


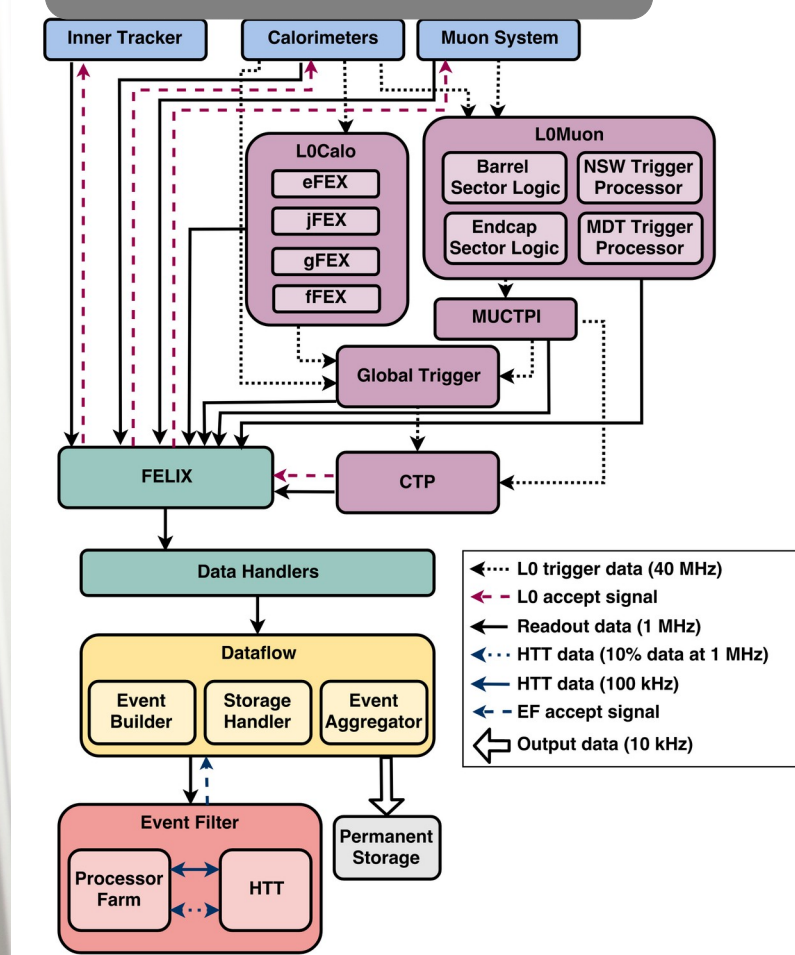
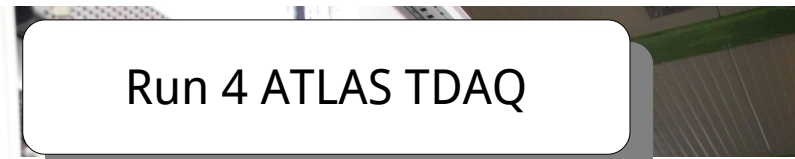
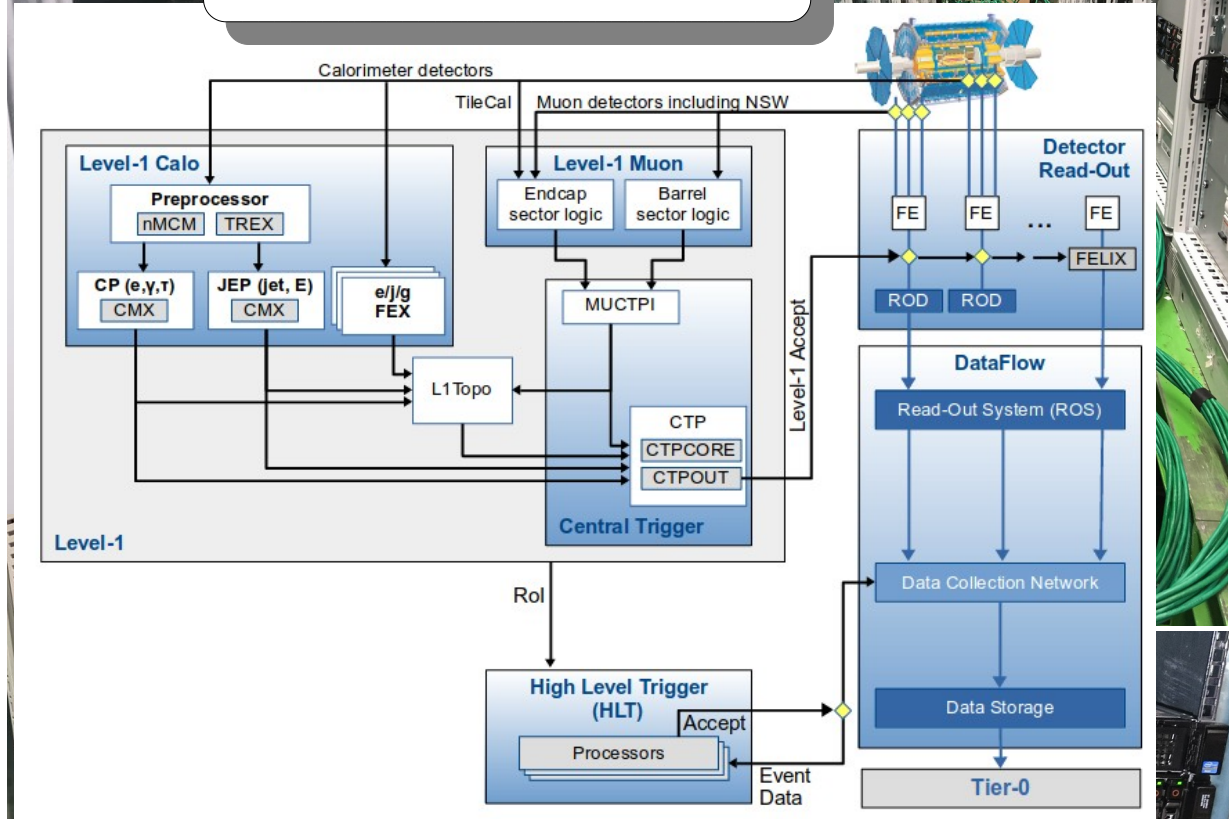
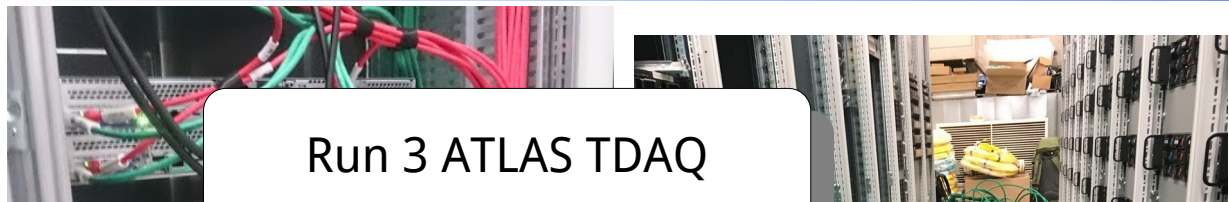
ATLAS Collaboration  
 Technical Design Report for the Phase-II Upgrade of  
 the ATLAS TDAQ System  
<https://cds.cern.ch/record/2285584>

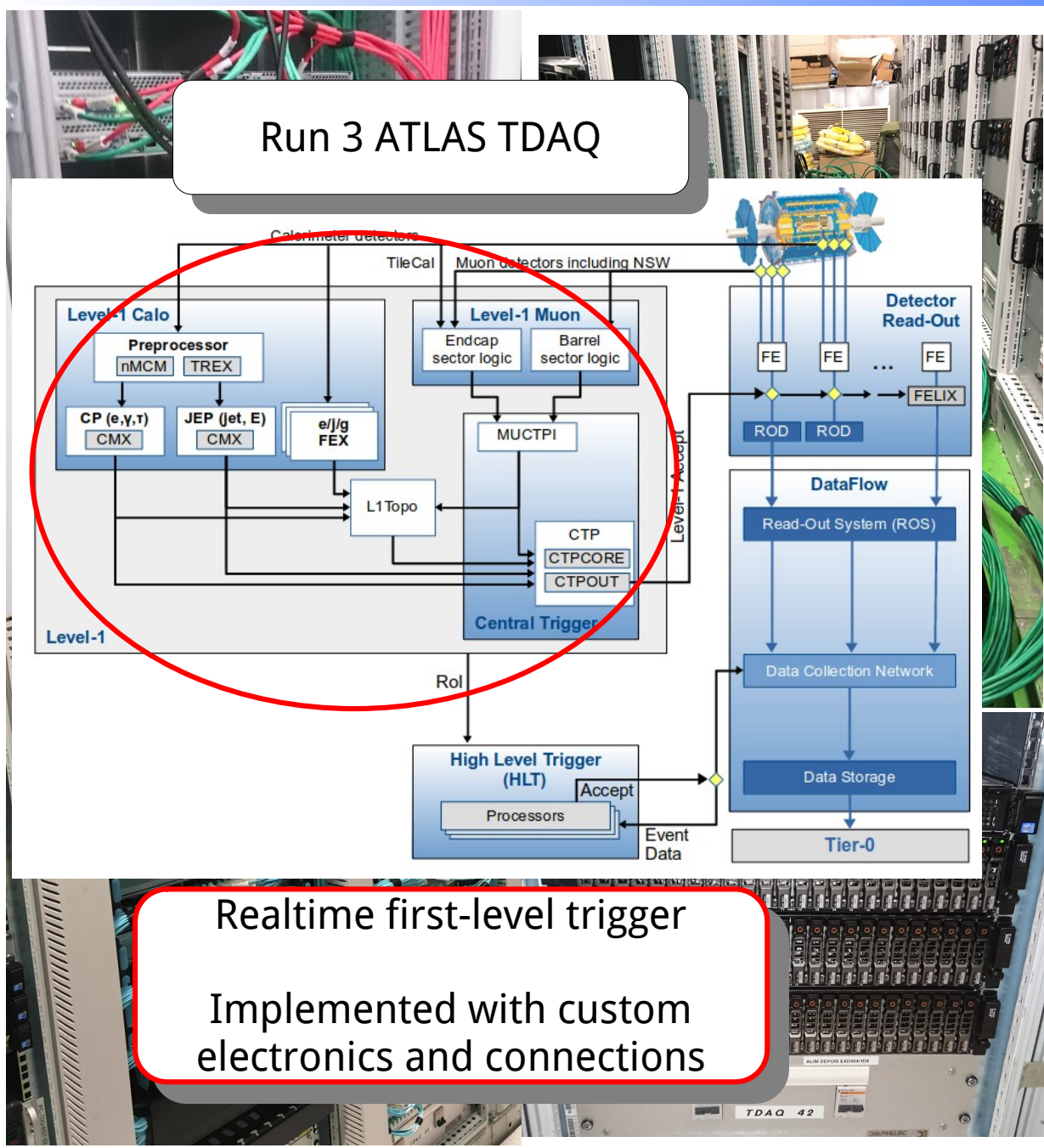


Created with mapchart.net ©

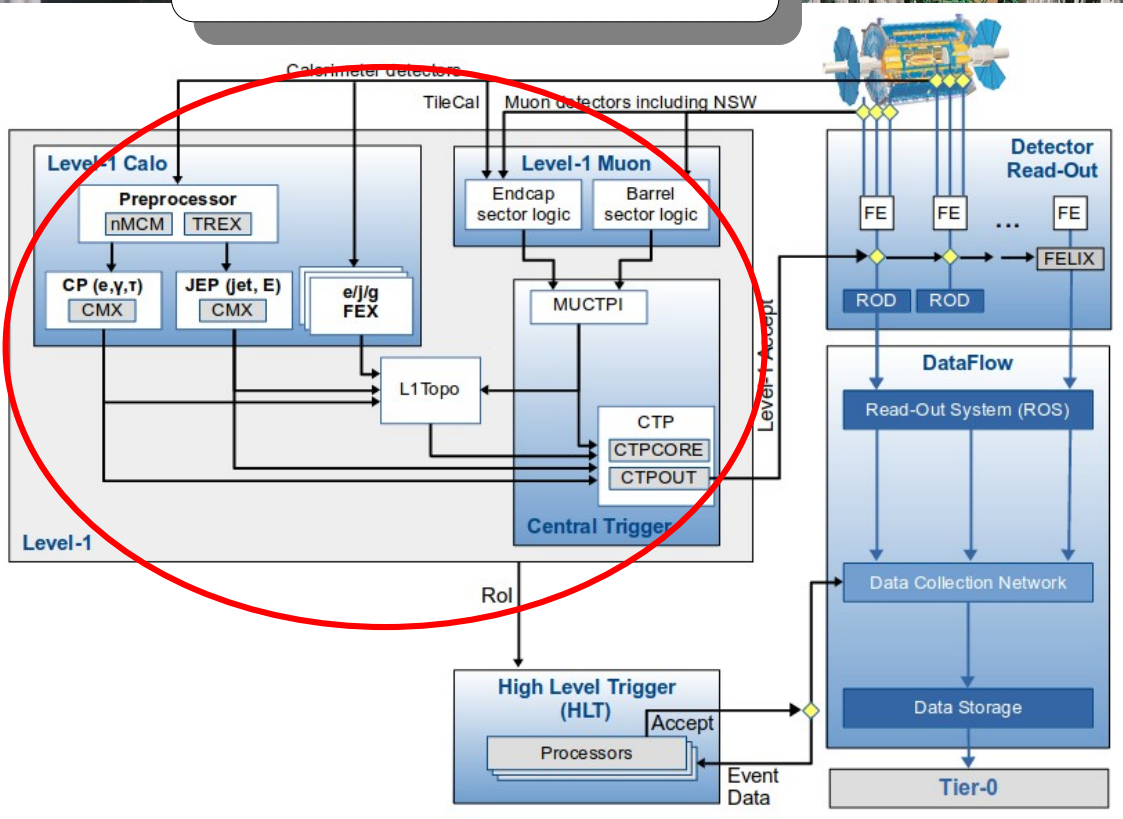
## LHC / HL-LHC Plan





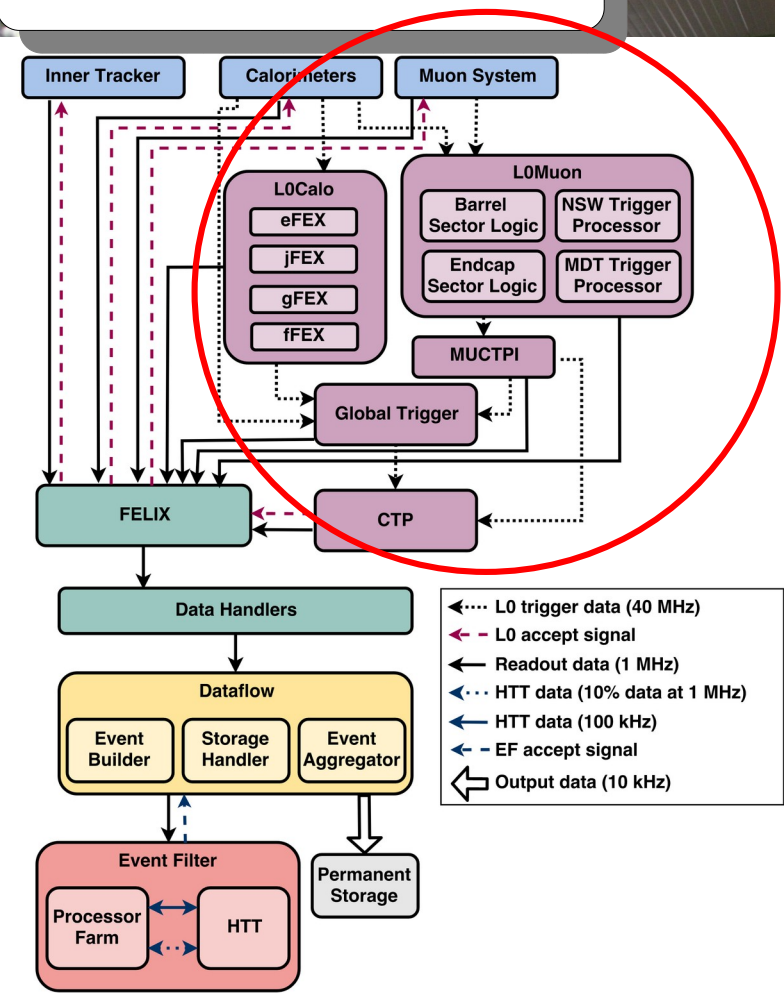


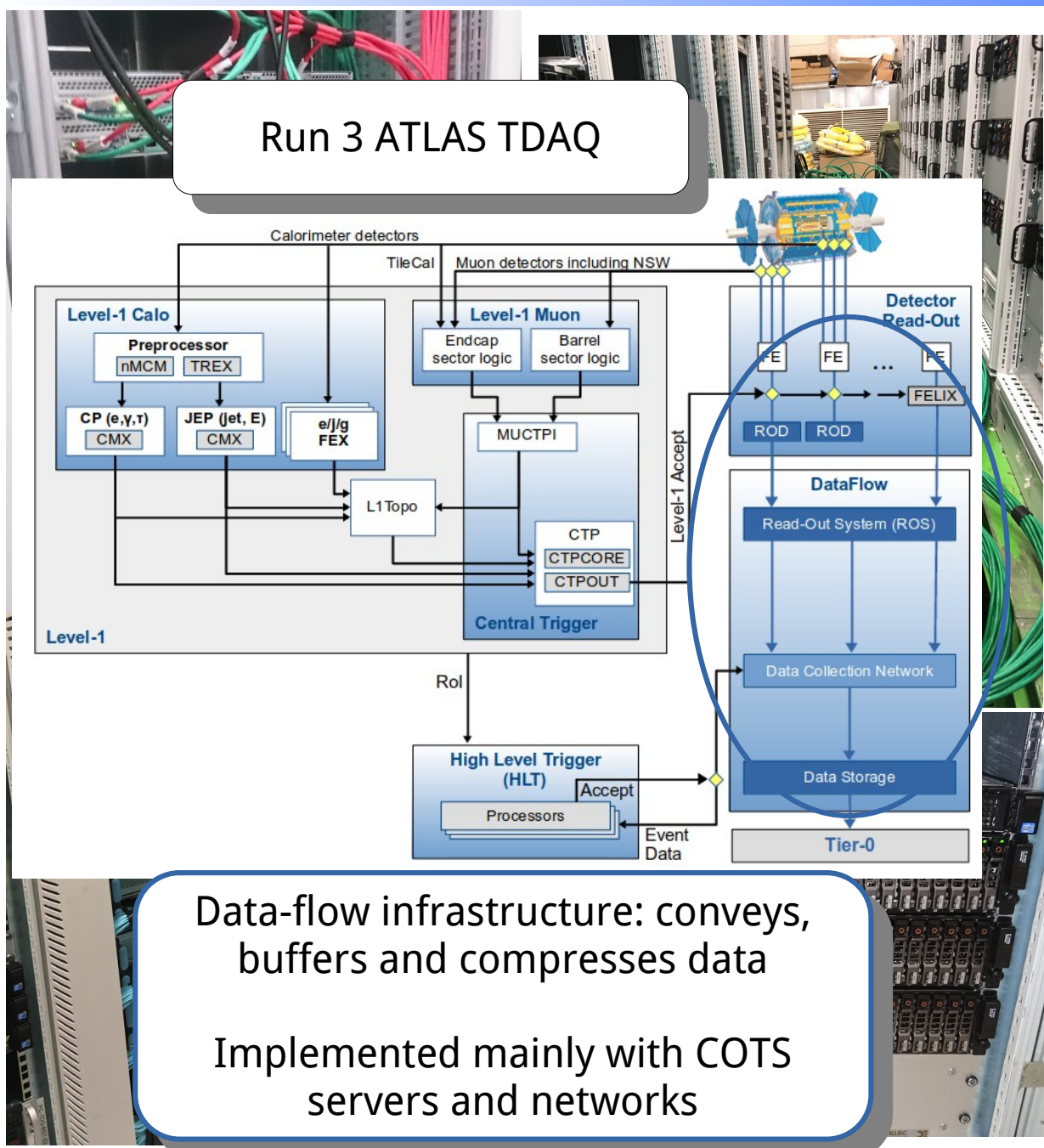
Run 3 ATLAS TDAQ



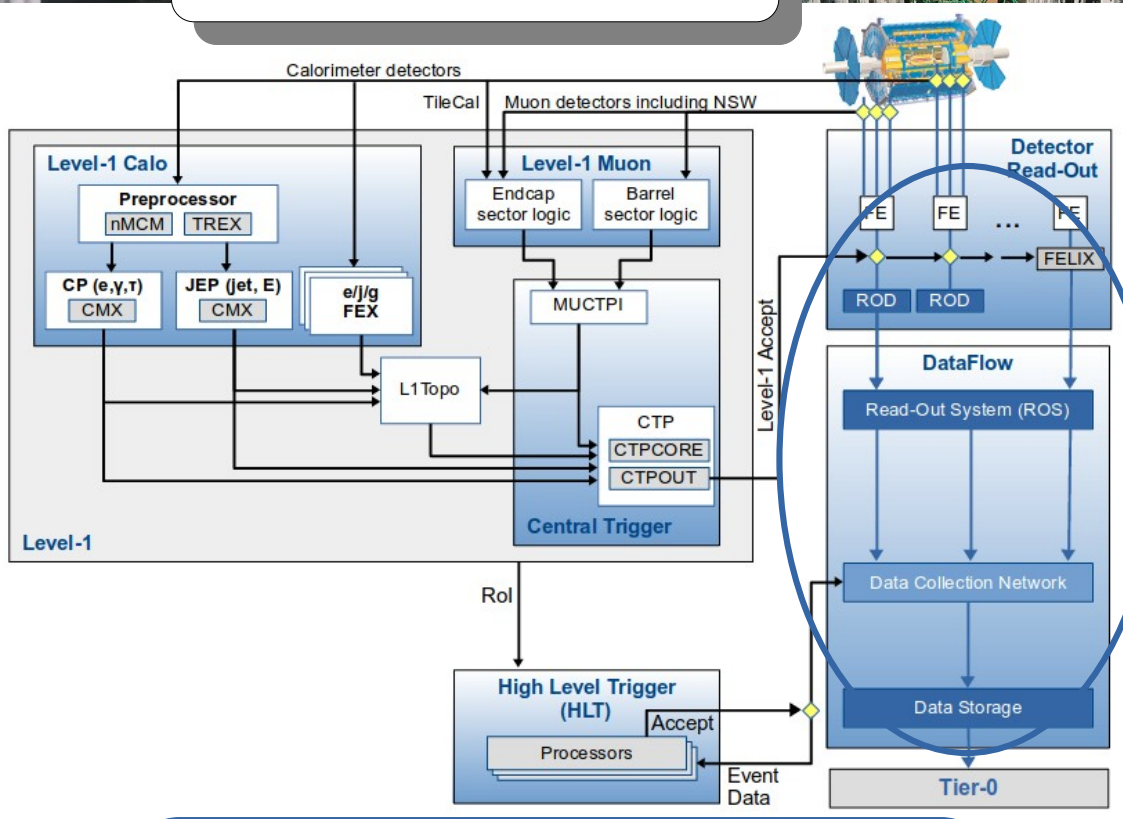
Realtime first-level trigger  
Implemented with custom electronics and connections

Run 4 ATLAS TDAQ

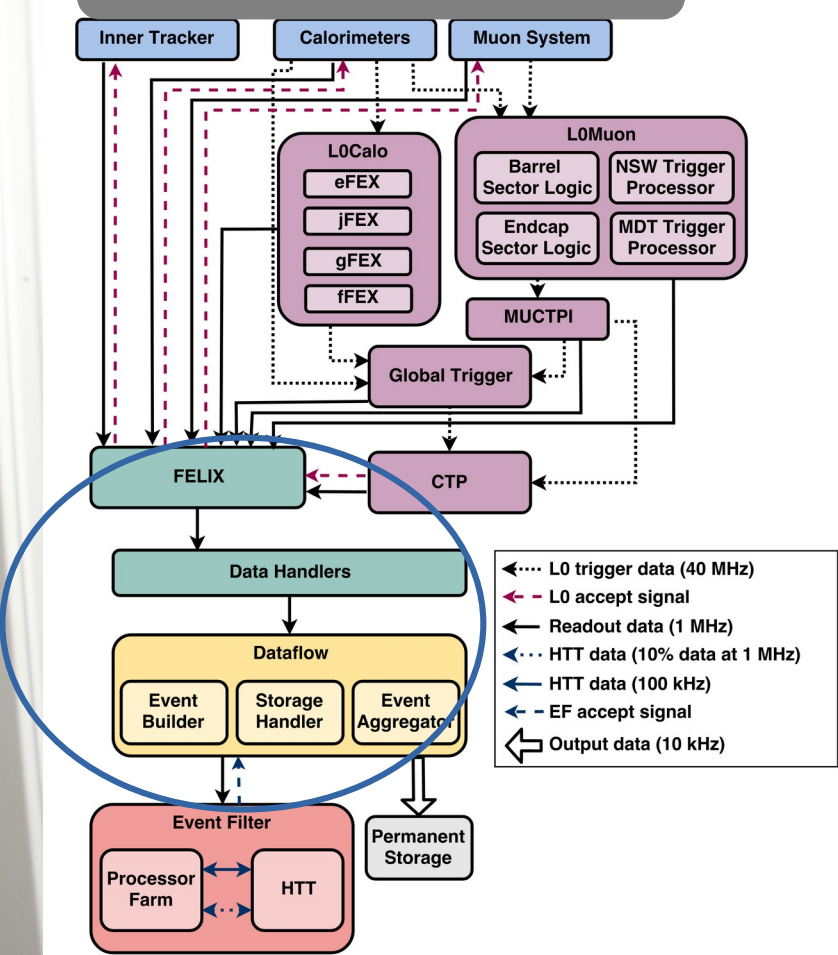




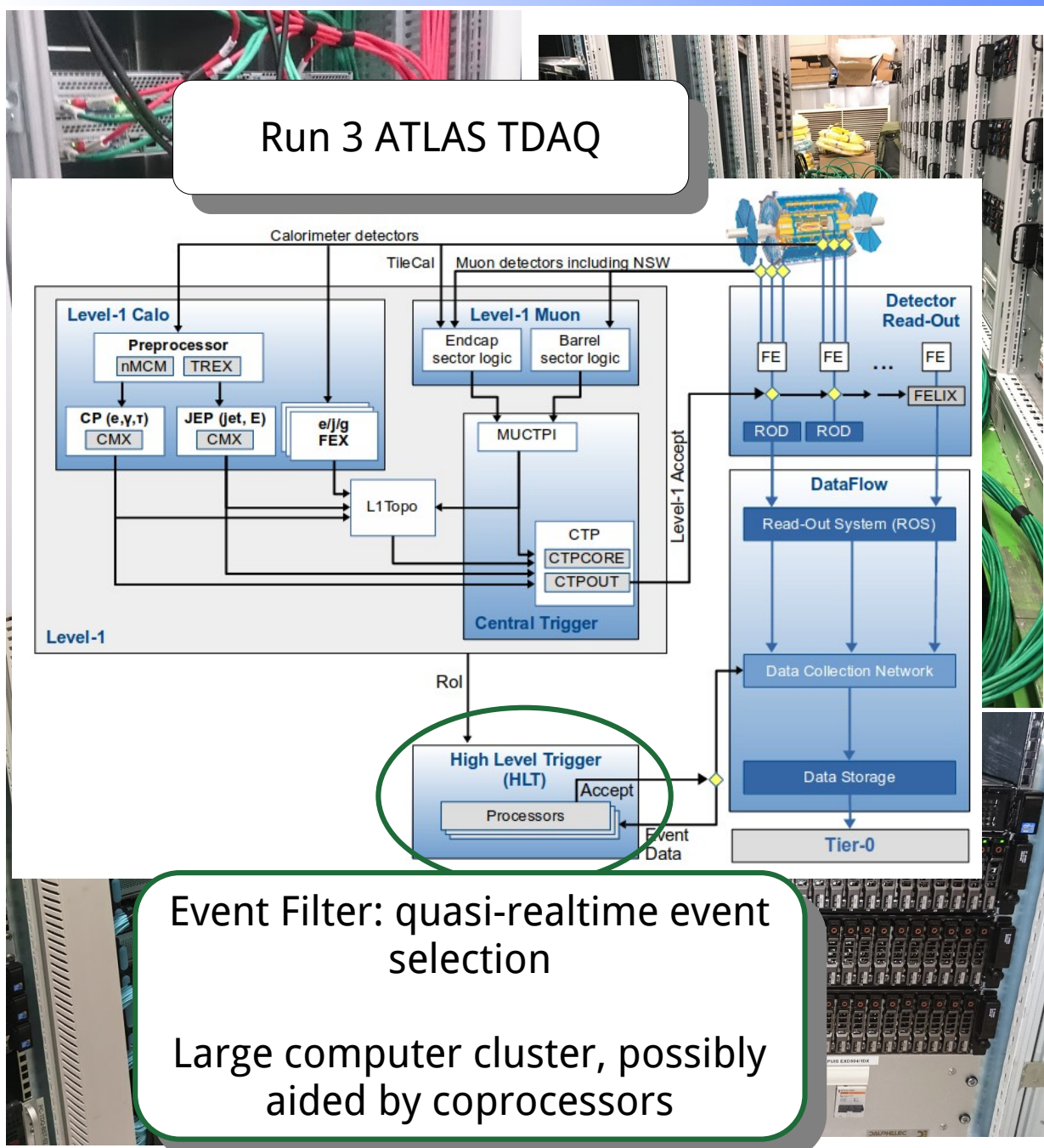
Run 3 ATLAS TDAQ



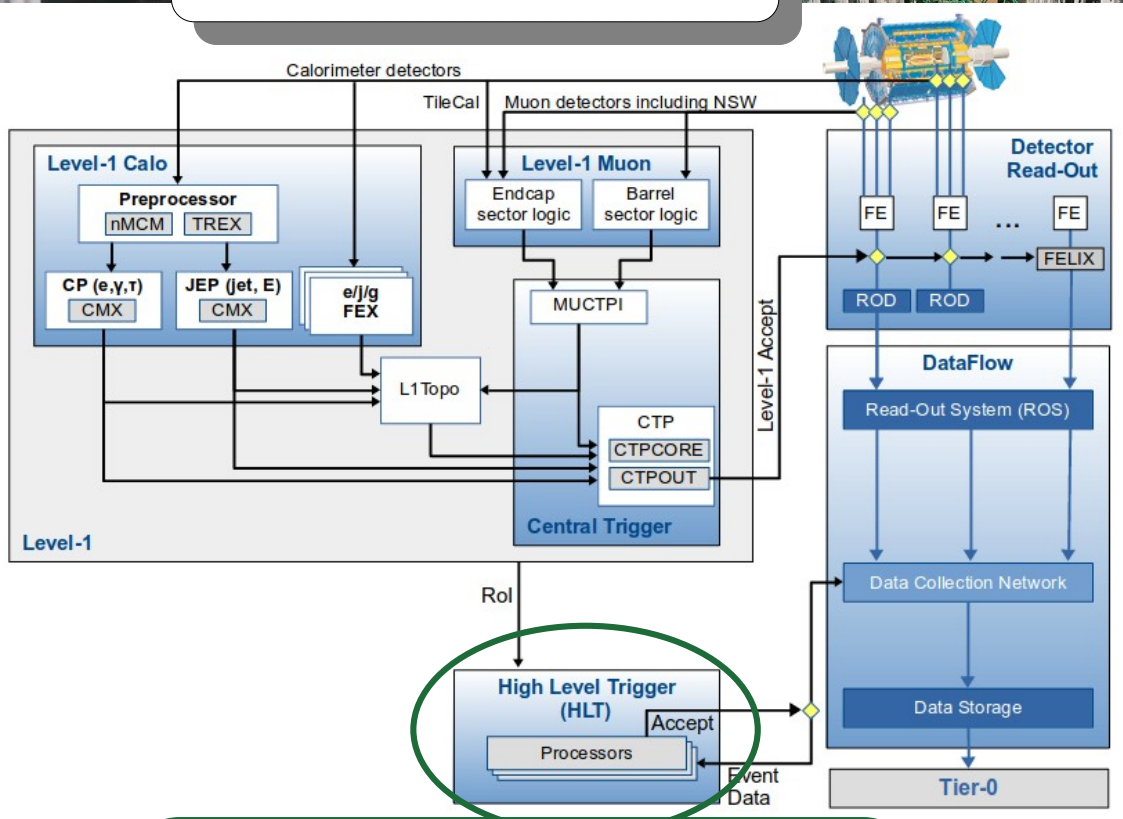
Run 4 ATLAS TDAQ



Data-flow infrastructure: conveys, buffers and compresses data  
 Implemented mainly with COTS servers and networks



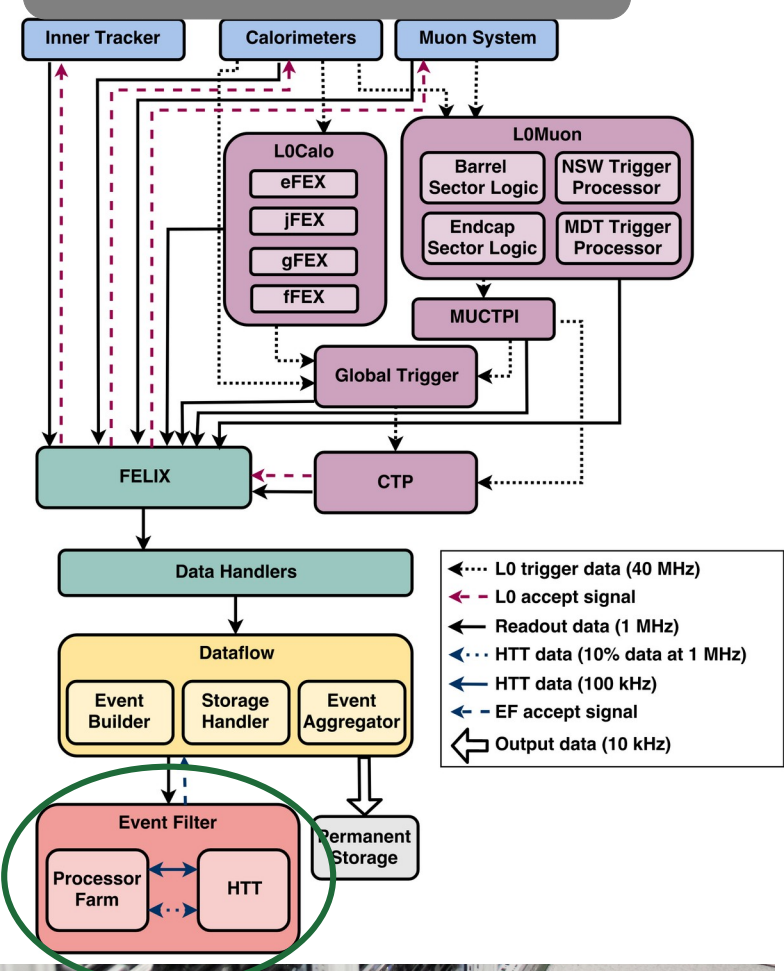
## Run 3 ATLAS TDAQ



Event Filter: quasi-realtime event selection

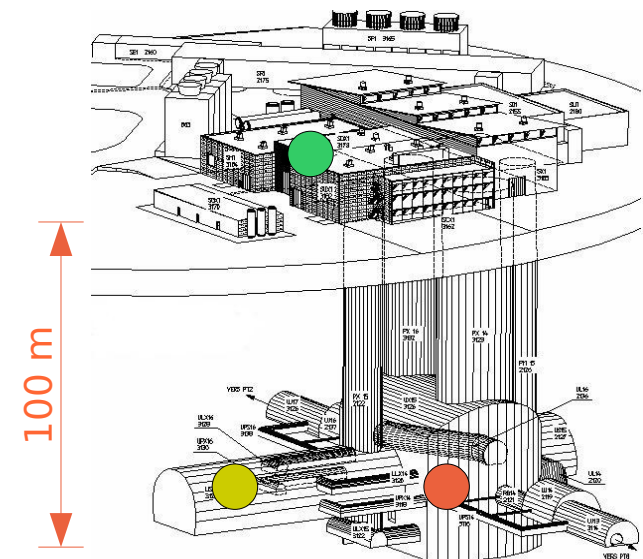
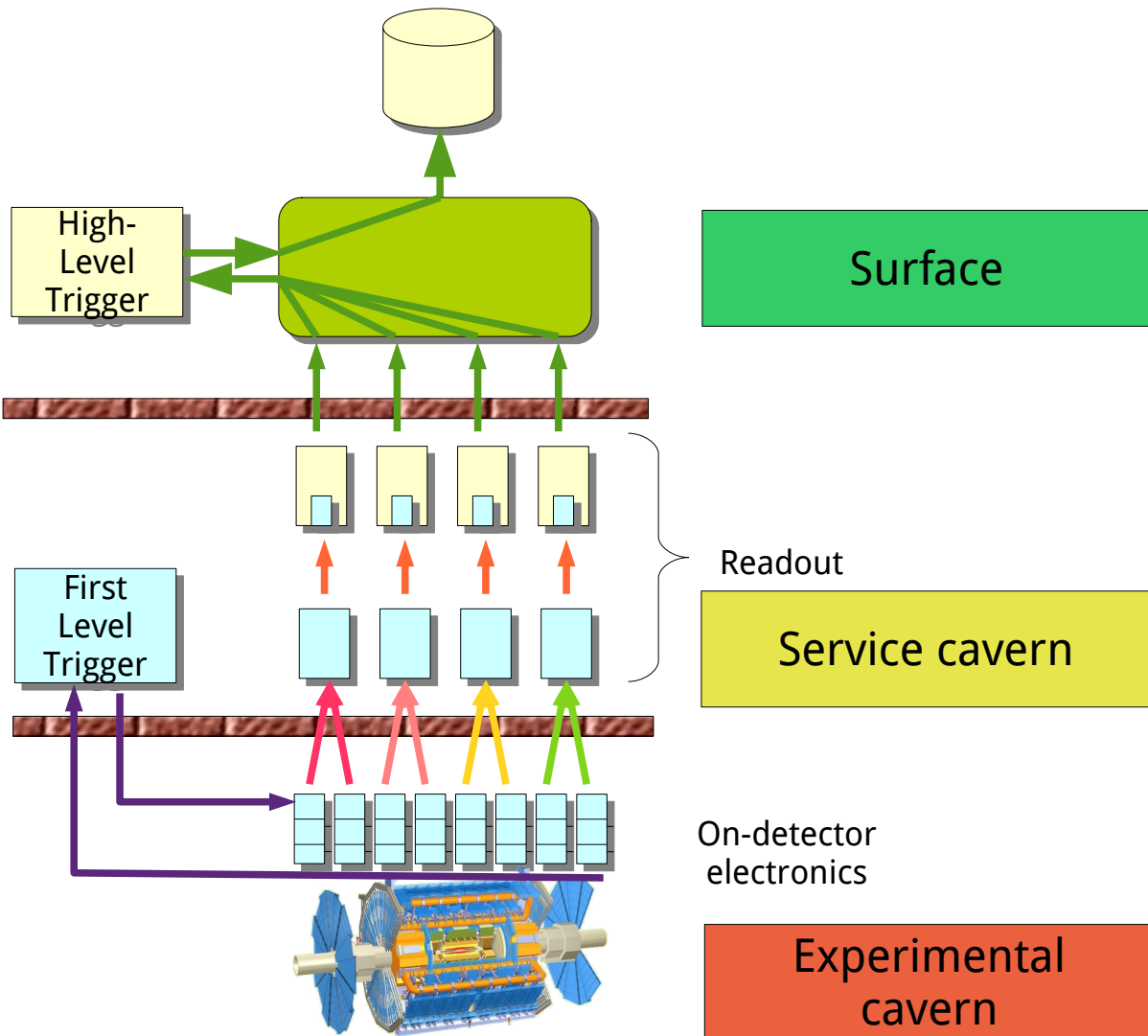
Large computer cluster, possibly aided by coprocessors

## Run 4 ATLAS TDAQ

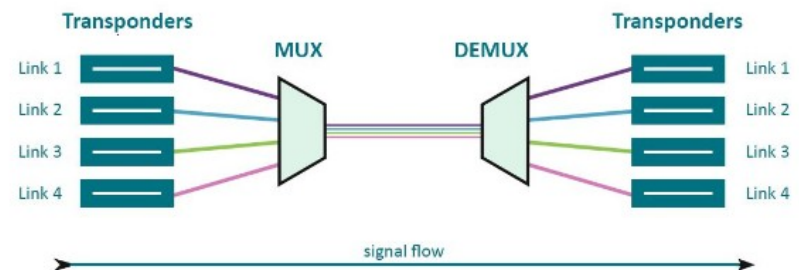


←···· L0 trigger data (40 MHz)  
 ←- - L0 accept signal  
 ← Readout data (1 MHz)  
 ←···· HTT data (10% data at 1 MHz)  
 ←···· HTT data (100 kHz)  
 ←- - EF accept signal  
 ← Output data (10 kHz)

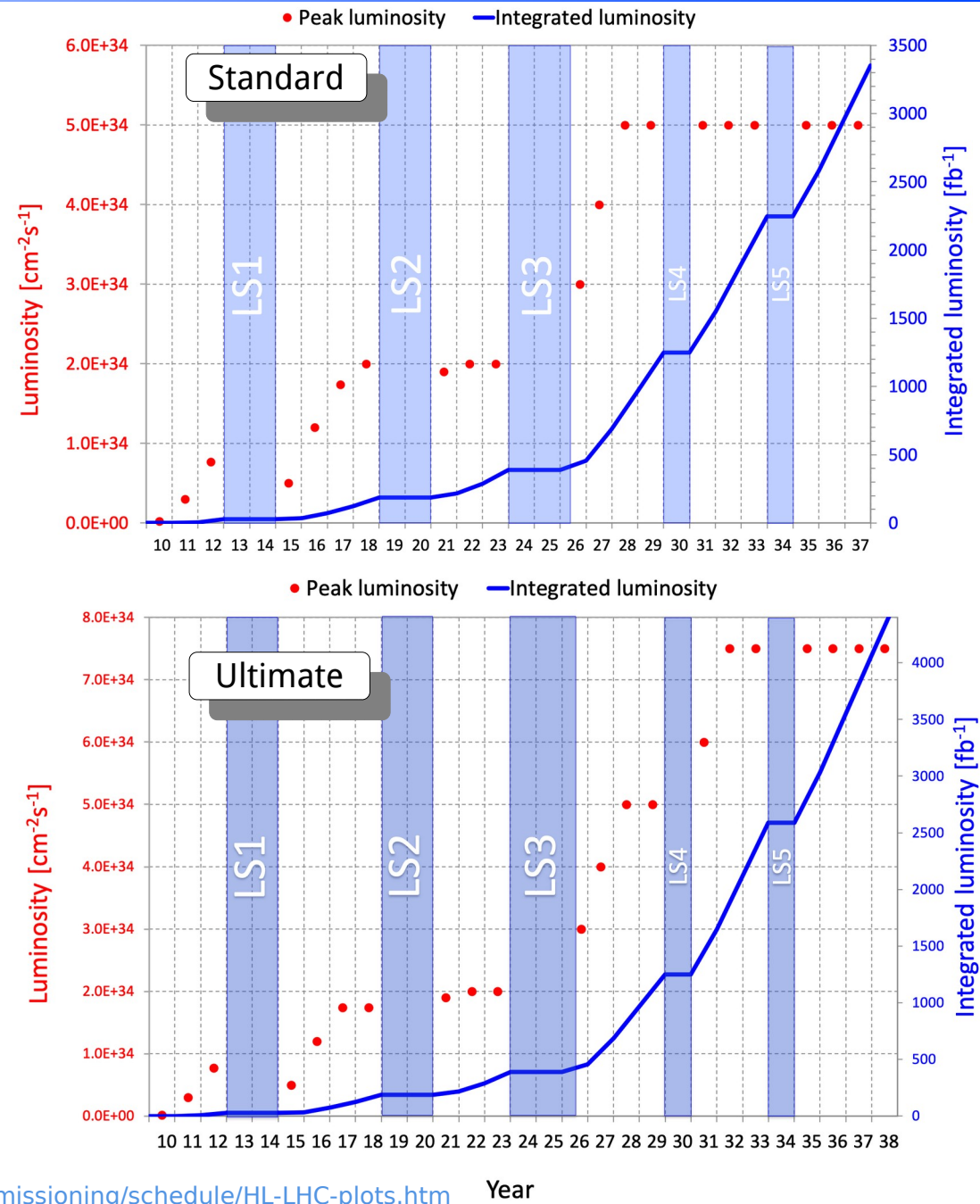




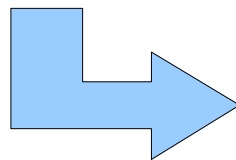
- e.g. considering WDM (wavelength division multiplexing) to reduce the number of long range fibres
  - trade-off like fibres vs transceivers cost



- HL-LHC programme aims at a total integrated luminosity of at least 3000 fb<sup>-1</sup>
  - ten-fold increase wrt Run 1/2/3 aggregate
- Corresponding increase in peak instantaneous luminosity
  - $L \sim 5 \cdot 10^{34}$  cm<sup>-2</sup>s<sup>-1</sup> (ultimate 7.5 · 10<sup>34</sup> cm<sup>-2</sup>s<sup>-1</sup>)
  - achieved mainly via pileup  $\langle \mu \rangle$ : 140 (ultimate 200)
- For reference Run 3 operation point:
  - $L \sim 2 \cdot 10^{34}$  cm<sup>-2</sup>s<sup>-1</sup> -  $\langle \mu \rangle \sim 50$



- The challenging and broad HL-LHC programme requires trigger thresholds comparable with the current ones, e.g.:
  - electroweak scale requires low  $p_T$  leptons
  - searches for new physics with low  $\Delta m$
  - HH measurements requires low  $p_T$  jets /b-jets
- At fixed threshold, trigger rates scale with peak luminosity
  - worsened by pileup environment

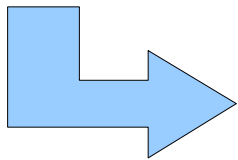


Major increase in readout and recording rates

Trigger Selection offline threshold (GeV)	Run 1	Run 2	HL-LHC
Isolated single e	25	27	<b>22</b>
Isolated single $\mu$	25	27	<b>20</b>
Di- $\gamma$	25, 25	25, 25	<b>25, 25</b>
Di- $\tau$	40, 30	40, 30	<b>40, 30</b>
Four-jet w/ b-jets	45	45	<b>65</b>
$H_T$	700	700	<b>375</b>
MET	150	200	<b>200</b>

	Run 3	Run 4
Readout rate (MHz)	0.1	1 (4)
Recording rate (kHz)	1.5	10

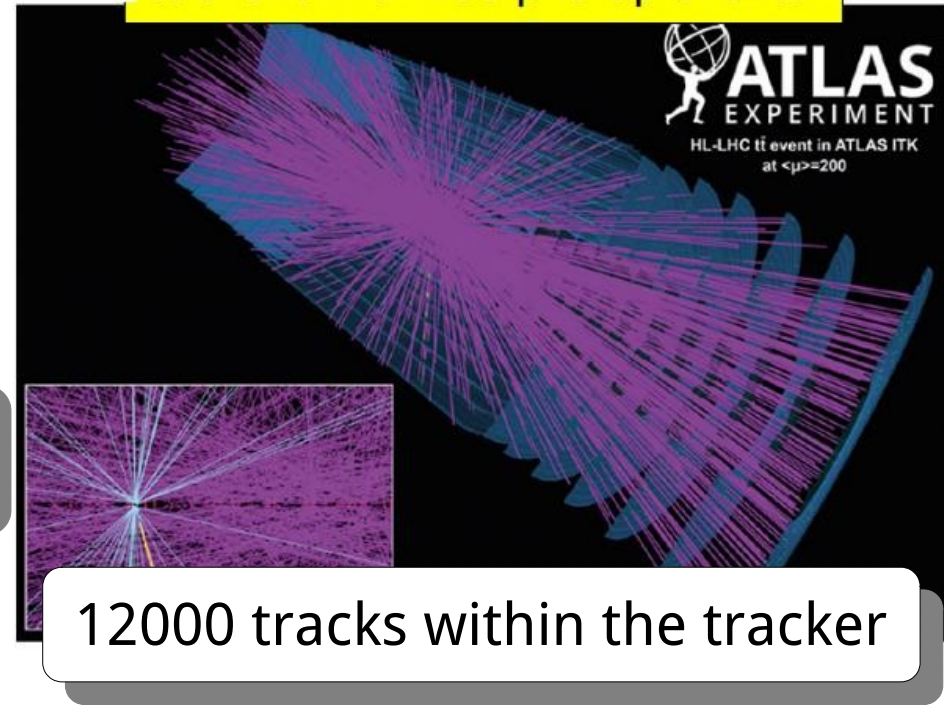
- High-granularity to cope with pileup
  - both for readout and trigger
  - complete replacement of inner detector → ITk



Larger event size

- Higher readout rate needs overhaul of detector front-end electronics
  - occasion to increase first level-trigger latency
    - *currently limited by on-detector buffer depths*
  - adopt unified readout link technology
    - *GBT/Versatile*

t $\bar{t}$  event with 200 pile-up events



12000 tracks within the tracker

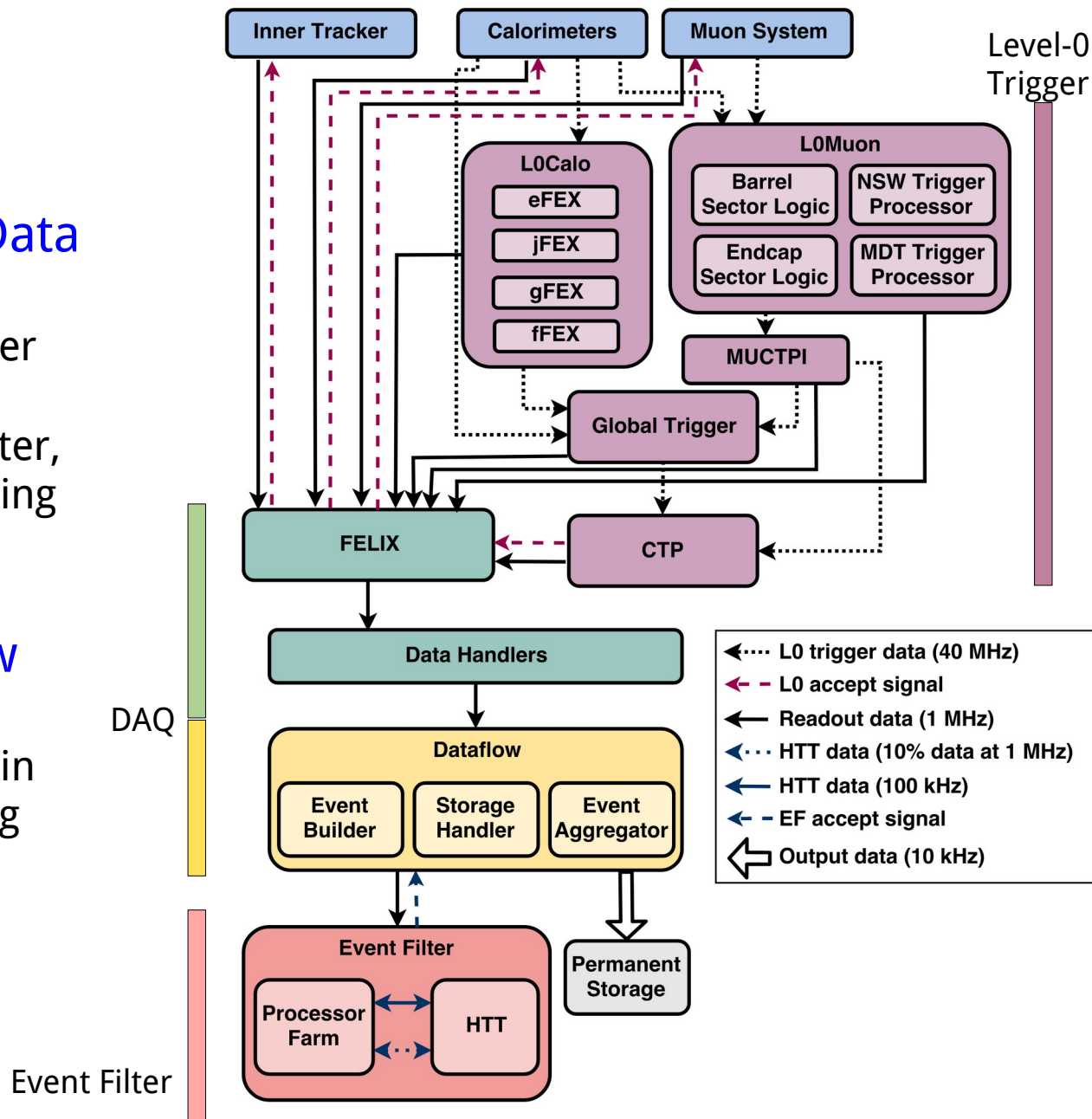
	Run 3	Run 4
First-level trigger latency ( $\mu$ s)	2.5	10
Event size (MB)	2.5	>5

- **Two-Level Trigger and Data Acquisition System**

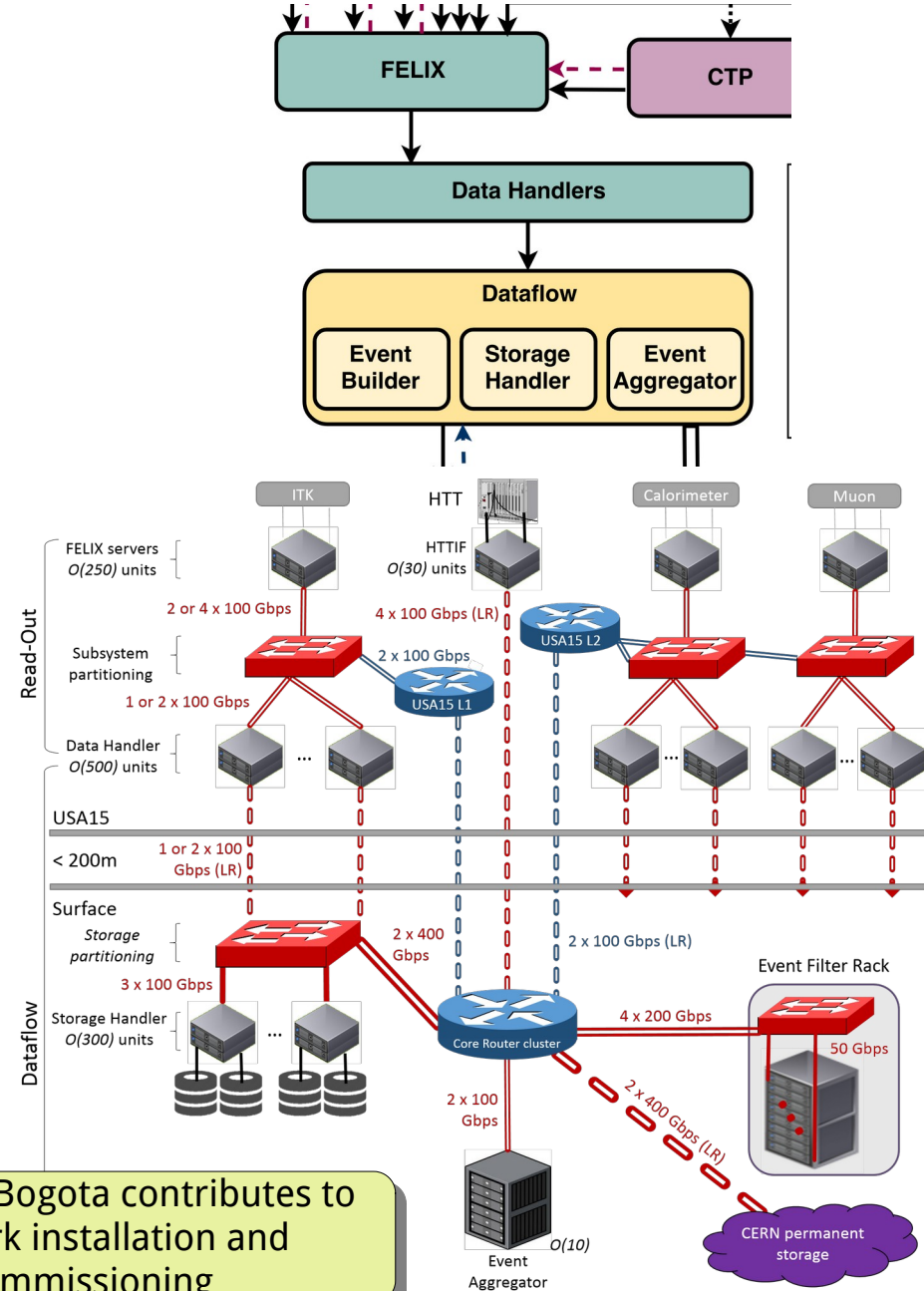
- hardware-based L0 trigger system
- software-based Event Filter, aided by dedicated tracking accelerator

- **Storage-based data-flow infrastructure**

- decouple realtime domain from software processing
- enable advanced data processing strategies



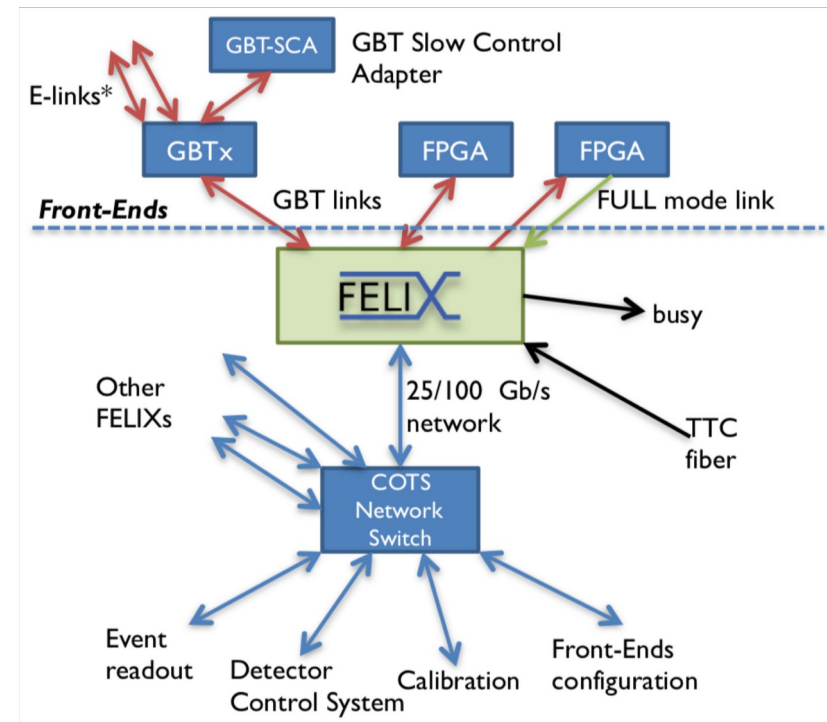
- **DAQ infrastructure responsible for**
  - interfacing the detector readout links to a commercial network domain
  - buffering the data and serving them to the Event Filter processors
  - discarding rejected events and formatting selected data for offline transfer
- **Largely implemented with commodity off the shelf hardware**
- **Backbone is a multi-layered sliced network**
  - baseline design based on Ethernet, do not exclude HPC technologies
  - ~2500x 100 Gbps
  - ~200x 200 Gbps
  - ~70x 400 Gbps



Network Design, Installation, Operation, Simulation

Colombia/Bogota contributes to network installation and commissioning

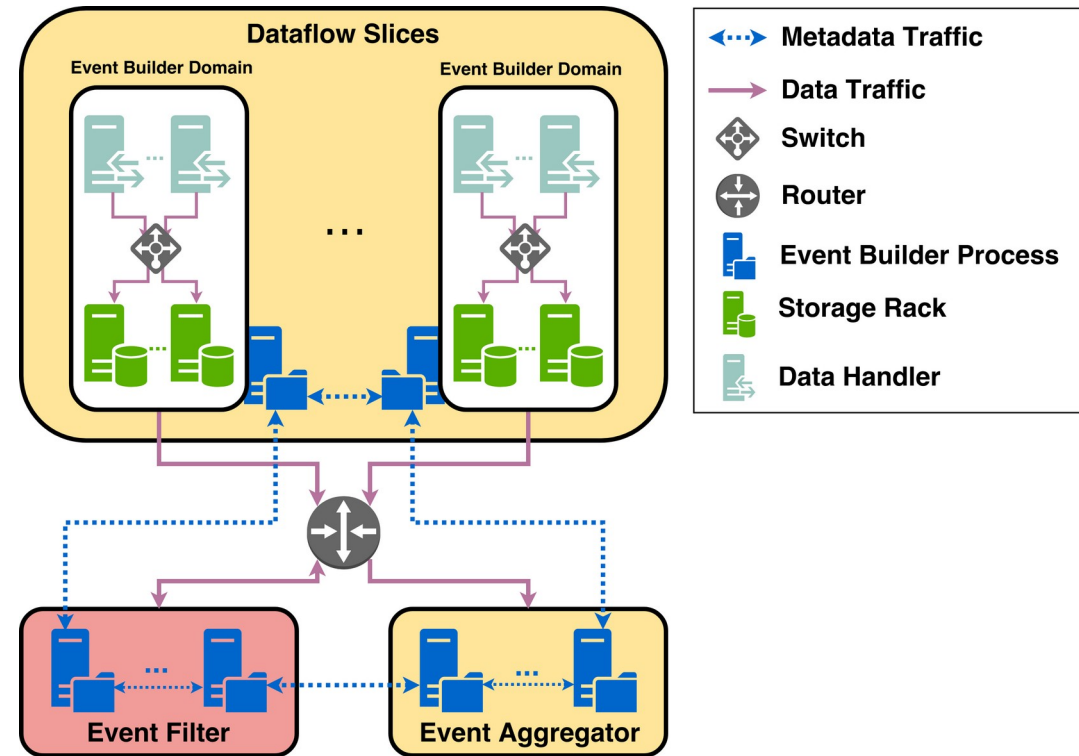
- Detector interfacing relies on a concept being deployed for Run 3
  - extended to the whole ATLAS
- Front-end Link Exchange (FELIX) acts a heterogeneous router
  - translates between network and serial links
  - distributes timing and trigger signals
  - as detector-agnostic as possible
    - *still provision for detector specific functions*
- Implementation based on commercial servers equipped with custom FPGA-based PCIe interfaces
  - plan for 48 10Gbps links per card
  - ~550 cards serving almost 20000 links
- In Run 3 a single FELIX server will handle 40 MHz of data frame rate at 15 Gbps
  - challenging software operating close to the hardware
    - *6 core 3 GHz CPU → ~500 clock cycles per frame*
  - expect at least factor 10 increase in Phase-II



Phase-I FELIX Interface Card

High-Performance I/O Software, Linux Drivers

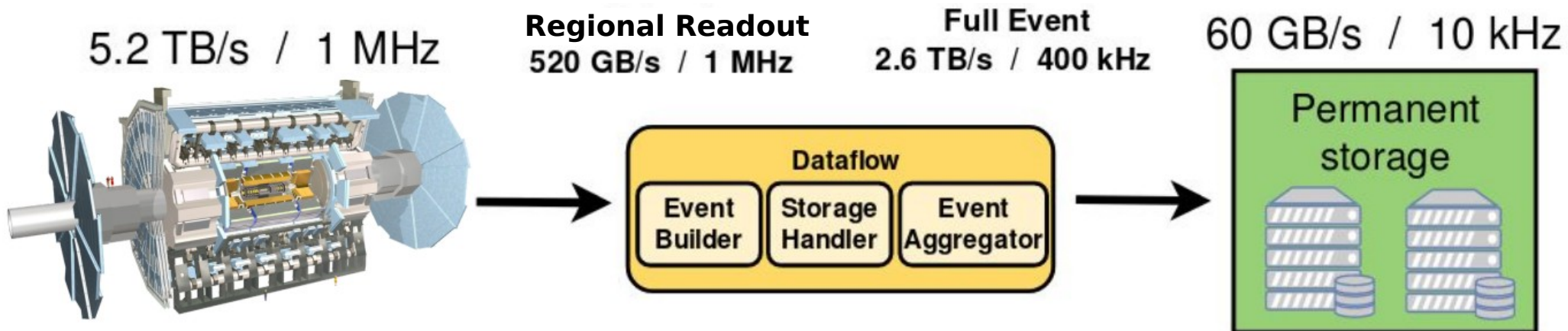
- Extend the DAQ buffering capabilities using a large storage infrastructure
  - decouple realtime domain (Level-0) and software domain (Event Filter)
  - enable delayed processing or fail-over scenarios
- Event Filter computer farm may be operated similarly to a batch system
  - quasi-realtime data stream required for online physics and detector monitoring



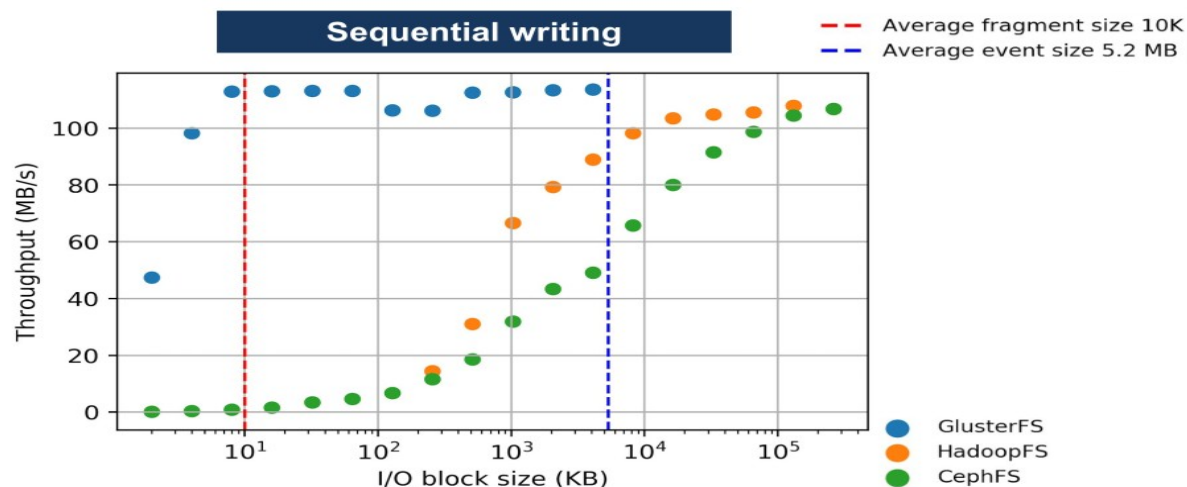
Component Connection	Traffic
Detector Front-ends to <a href="#">FELIX</a>	5.2 TB/s
<a href="#">FELIX</a> to Data Handlers	5.2 TB/s
Data Handlers to Event Builder/Storage Handler	5.2TB/s
Storage Handler to Event Filter	2.6 TB/s
Event Filter to <a href="#">HTTIF</a>	Event Filter to <a href="#">rHTT</a> Event Filter to <a href="#">gHTT</a>
Event Filter to Event Aggregator and Permanent Storage	60 GB/s

Storage, Databases, Simulation





- Understand feasibility of commodity hardware and software for the storage infrastructure
  - hardware infrastructure: novel solid state technologies, hierarchical storage, ...
  - storage software: distributed file system, distributed hash-tables, ...
  - operation model and interfaces: trade-off between compute, networking, storage

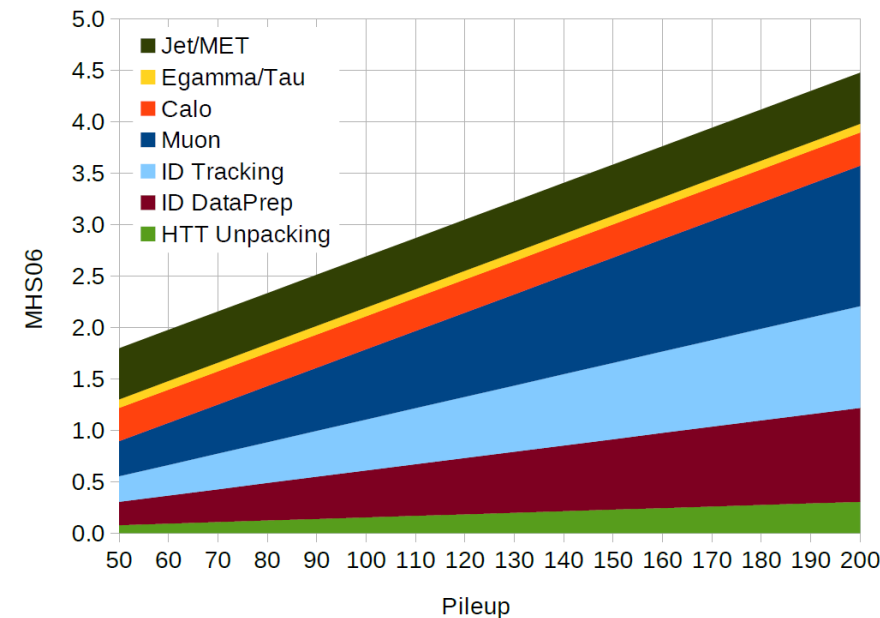


Argentina/UBA simulations of data movements

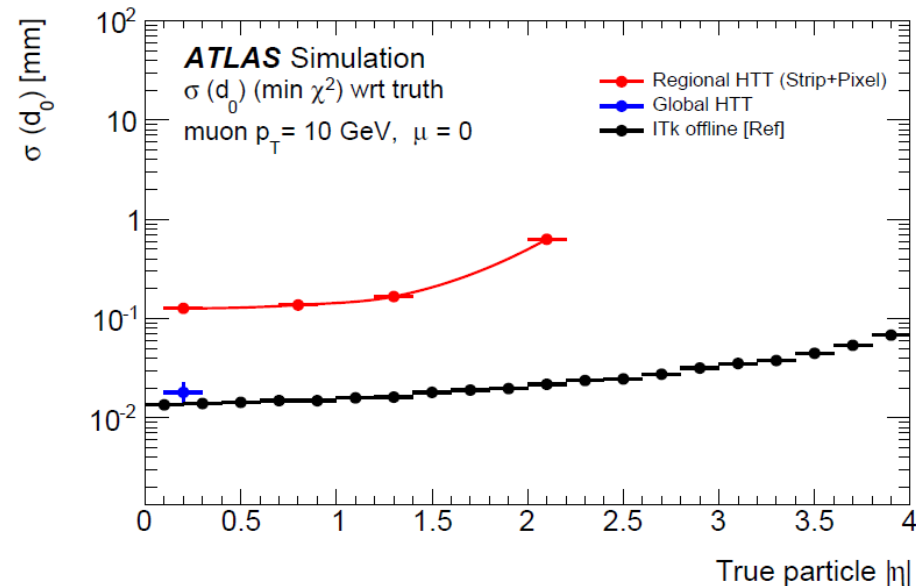
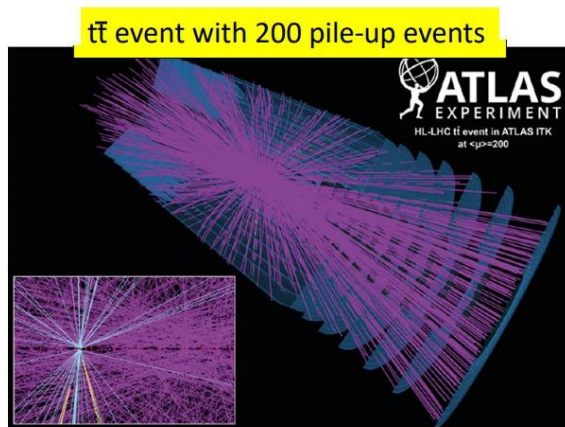
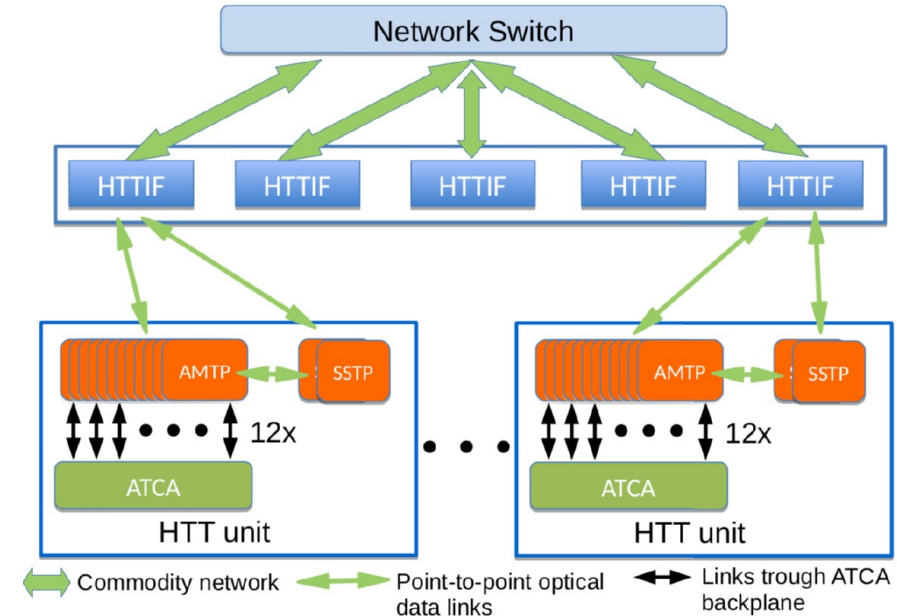
- Similar to Run 3 → large computer farm
  - aided by a dedicated tracking system
  - performs the last level of selection from 1 MHz to 10 kHz
- In high pileup environment tracking is key to recover algorithms performance and maintain low thresholds
  - separation of electrons and background jets
  - calculation of global event quantities like  $E_{T}^{\text{miss}}$
  - jet energy resolution
- Event Filter baseline implementation is based on CPUs
  - in parallel investigations of accelerators (GPGPU & FPGA) and associated dedicated algorithms

High-Performance Software, AI/ML, GPU/FPGA Programming

Brazil AI-based e/ $\gamma$  identification



- ITk tracking software → 10 times larger computer farm would be required
  - based on current tracking software
  - ongoing software optimisations potential to significantly reduce this estimate
- HTT (Hardware Track Trigger) massively parallel device
  - custom electronics using Associative Memories (AM ASICs) for pattern recognition and FPGAs for fitting
  - driven by the Event Filter requests

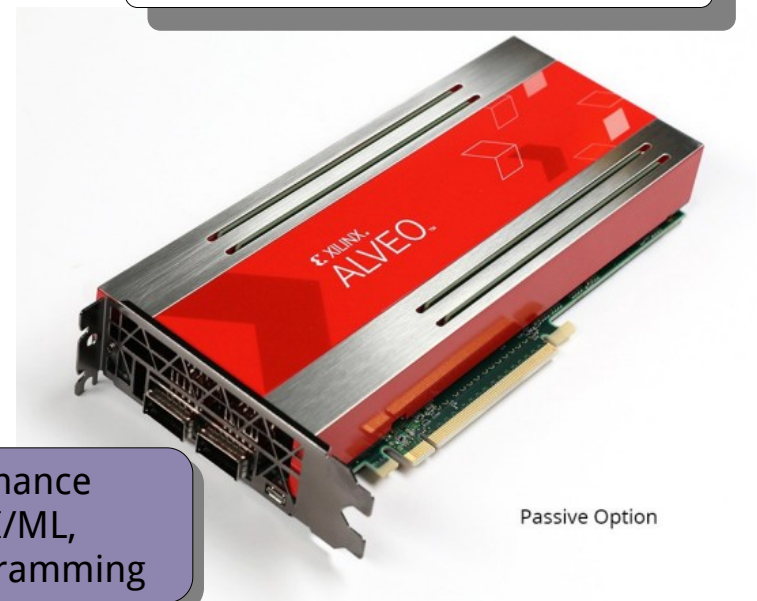


- Active investigations in COTS alternative to tracking in custom hardware
- Based on
  - CPU-based software implementation
  - coprocessor and specialised algorithms
- Studying both algorithmic and implementation improvements

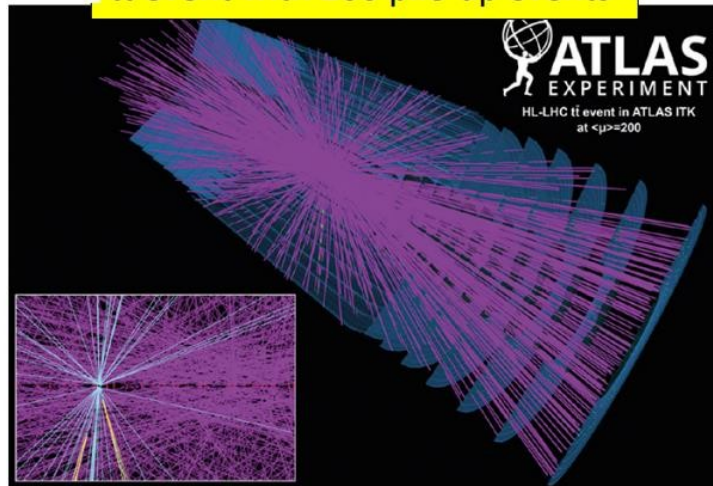
GPGPU



FPGA



$t\bar{t}$  event with 200 pile-up events



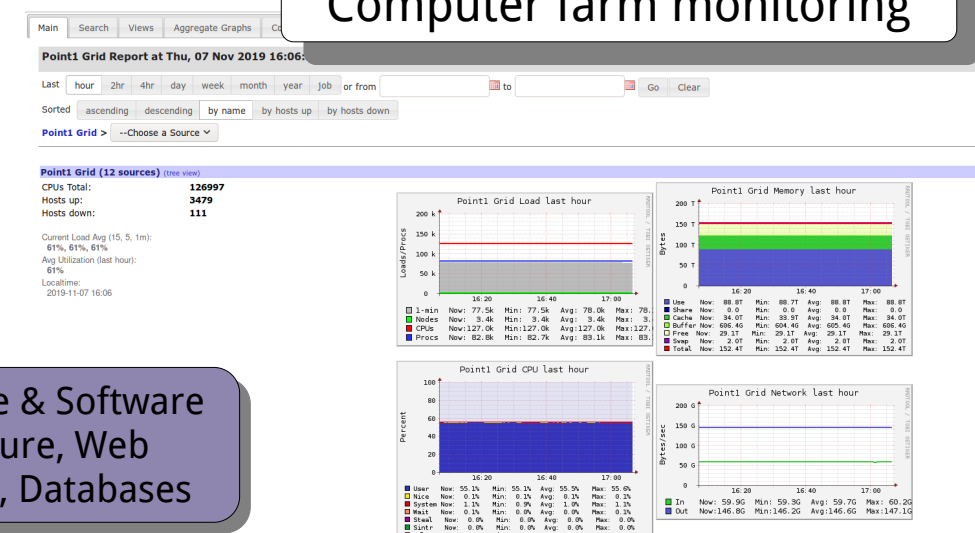
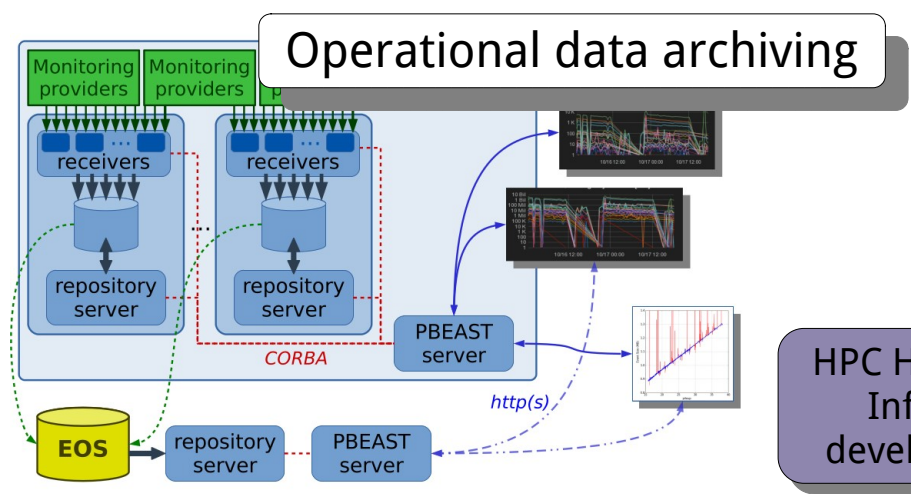
High-Performance Software, AI/ML, GPU/FPGA Programming

Passive Option

- **Managing, monitoring and configuring 50000+ applications**
  - on a heterogeneous cluster
  - on a mission-critical duty
- **Storage system enables staged filtering**
  - mechanism to control applications in different domains (realtime) and with different lifetimes
- **Orchestration systems current answer in cloud environment**
  - applicability to data-acquisition to be evaluated
- **More in general, require to scale all aspects of monitoring**
  - operational and hardware
  - detector
  - trigger and physics



## Computer farm monitoring



HPC Hardware & Software Infrastructure, Web development, Databases

ATLAS TDAQ sits at the centre of the action, uniting detector, physics, off-line computing

Phase-II upgrade major scale-up of the TDAQ computing infrastructure

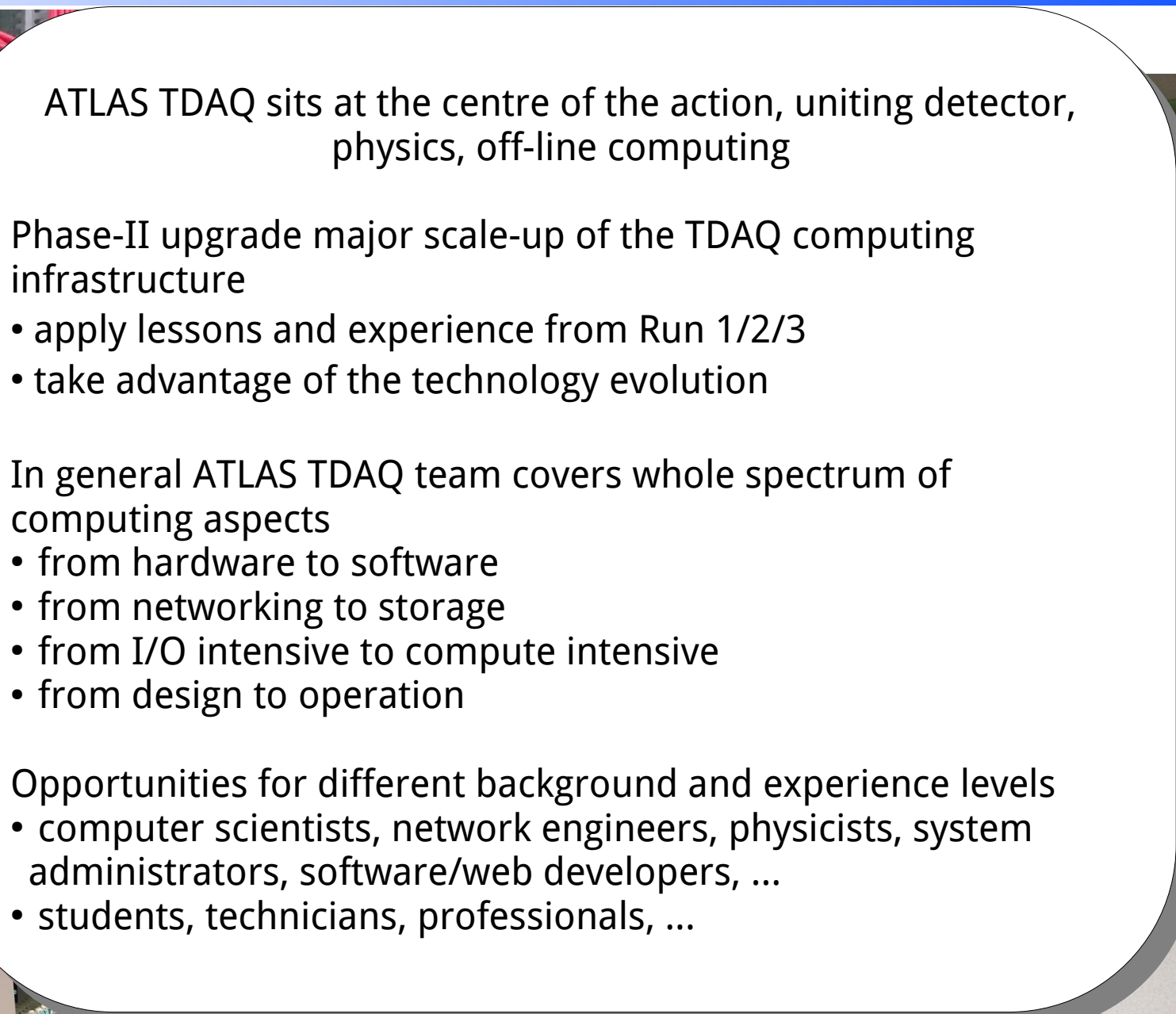
- apply lessons and experience from Run 1/2/3
- take advantage of the technology evolution

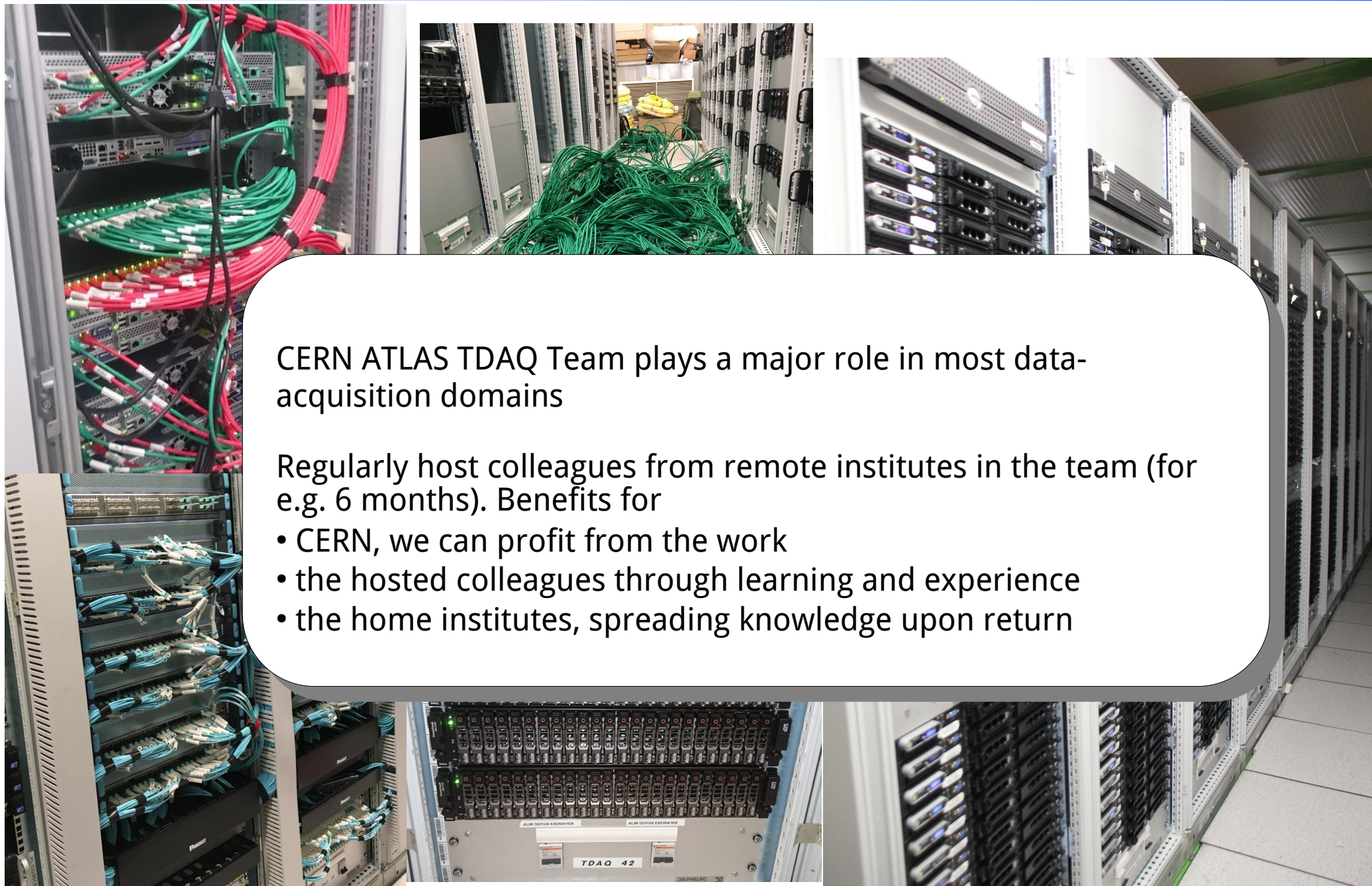
In general ATLAS TDAQ team covers whole spectrum of computing aspects

- from hardware to software
- from networking to storage
- from I/O intensive to compute intensive
- from design to operation

Opportunities for different background and experience levels

- computer scientists, network engineers, physicists, system administrators, software/web developers, ...
- students, technicians, professionals, ...





CERN ATLAS TDAQ Team plays a major role in most data-acquisition domains

Regularly host colleagues from remote institutes in the team (for e.g. 6 months). Benefits for

- CERN, we can profit from the work
- the hosted colleagues through learning and experience
- the home institutes, spreading knowledge upon return



**Bonus**



- Evolution path to a two-level hardware trigger included in the design
  - L0 – 4 MHz
  - L1 – 1 MHz
  - Event Filter – 10 kHz
- Possible transition from baseline to evolution driven by physics requirements
  - hadronic trigger rates
  - occupancy of inner layers of ITk
- Avoid the baseline TDAQ implementation restricting the trigger menu at the ultimate HL-LHC operating conditions
- Level-1 Trigger combines L0 objects with track information from a dedicated subsystem to discriminate against pileup in the calorimeter

