# LHCb: status and plans for the online system upgrade

Paolo Durante
*paolo.durante@cern.ch*
on behalf of the LHCb collaboration

# Outline

- New readout system
  - Slow & fast control
  - Optical links
  - Readout board
  - Event building
  - Data centre

- Online mile-stones

# Run3 upgrade

- Higher luminosity from LHC

- Sub-detectors will upgrade

- Removal of hardware trigger
  - Inefficient at high luminosity
  - New readout electronics

- Filter farm will need to handle:
  - Larger **event size** (~50KB to ~100KB)
  - Larger **event rate** (~1MHz to ~40MHz)

- New challenges for
  DAQ & High-Level Trigger

**Network – Projected Throughput [Tbit/s]**

| | Alice | Atlas | CMS | LHCb | LHCb today |
|---|---|---|---|---|---|
| 40.00 | | | ■ | ■ | |
| 30.00 | | | | | |
| 20.00 | | ■ | | | |
| 10.00 | ■ | | | | |
| 0.00 | | | | | ■ |

**2020**

Detector

~~Hardware trigger~~

Data acquisition

High Level Trigger farm

CERN long term storage

Offline physics analysis

# Run3 Online System

- Dimensioning the system:
  - **~10 000** optical links
    - detector to surface (**~350 m**)
    - up to **4.8 Gbps/link**
  - **~500** readout nodes (up to ~24 links/node)
  - **~40 MHz** event rate
  - **~100 KB** event size

- High bisection bandwidth in event builder network
  - ~32 Tb/s aggregate bandwidth
  - Leverage emerging 100G technologies

- Global configuration and control via ECS subsystem

- Global synchronization via TFC subsystem

Detector front-end electronics

Clock & fast commands

UX85B

x500 Event Builders (PC + readout board)

Clock & fast commands

TFC

throttle from PCIe40s

6 x 100 Gbit/s

Event Builder network

6 x 100 Gbit/s

subfarm switch

Online storage

subfarm switch

Eventfilter Farm ~ 80 subfarms

Point 8 surface

# ECS (Experiment Control System)

- Controls and monitors <u>all subsystems</u>
  - DAQ, TFC, HLT, farm…
- Continuity from current implementation
  - JCOP / DIM / WinCCOA / SMI++ / Recipes
- Already able to drive current readout board prototype, from input to output

- Frontends rely on GBT-SCA hardware by EP-ESE
- Low-level components are being implemented

Operation UI

HW Description UI

WinCC-OA

Control PC

Firmware

GbtServ
PCIe library ✓
PCIe driver ✓

DIMServer

FPGA

PCIe40

PCIe

Test UI

01/03/2016

Host PC

# TFC (Timing & Fast Control)

## Current status

- Already integrated in firmware
- Uses same readout board hardware as the DAQ (PCIe40)
- Can send fast commands to frontends
  - SciFi, UT, Muon ASICS already being tested
- Programmable internal throttle for bandwidth regulation

## Ongoing work

- PON (Passive Optical Network) technology integration (with EP-ESE)
- Clock phase tests on readout board (with CPPM)
- Continue feedback and compliance testing with frontend experts

# Long-distance optics


VTRx


MiniPOD™

- Counting room on surface
  - Power, cooling, space constraints in underground area
  - ~350 meter distance

- Based on EP-ESE technology
  - Rad-hard Versatile Link on frontends
  - Initially qualified for ~100m

- Fiber infrastructure by EN-EL
  - Pilot installation at end 2014

- Loopback tests in 2015
  - ~9 months, ~700 meters
  - Avago MiniPOD transceivers
  - Bit Error Rate $< 10^{-18}$
  - Full system equivalent: < 5 errors/day
  - ✓ LHCC milestone

- Continued tests in 2016
  - Versatile TX on frontend prototype
  - MiniPOD RX on readout board prototype

# Readout boards / Event builders

# Readout board hardware (PCIe40)

- **PCI Express add-in card**
  - Altera Arria10 FPGA
  - 100Gbps DMA engine
    to event-builder memory

- **High-density optical IO**
  - Up to 48 transceivers (Avago MiniPODs)
  - Reuse same HW for timing distribution system

- **Decouple FPGA from network**
  - Maximum flexibility in network technology

- **Exploit commercial technologies**
  - PCI Express Gen3 interconnect
  - COTS servers designed for GPU acceleration

- **2nd generation readout board**
  - Developed at CPP Marseille

- **Pre-production launched**
  - Ready end of Q2

- **Market survey completed**
  - Tender in H2 2016

# Readout board firmware

- Common architecture for all subdetectors

- Joint effort of LAPP, CERN, CPPM

- Highly configurable through frontend-specific parameters

- Handles ECS, TFC & DAQ, subdetector-specific logic

N+M = 24

# Readout unit dataflow

- A single Readout unit must sustain ~400Gbps IO bandwidth

- Optimize memory bandwidth
  - Design for zero-copy operations and RDMA over the network
  - Organize dataflow according to topology and IO resources
  - Exploit full network bandwidth

# Event-building software (DAQPIPE)

- Recreate distributed event-building dataflow of LHCb Run3

- Modular architecture, "drivers" for each network technology under evaluation

- Leverage existing HPC sites to assess scalability

- Close collaboration with the industry through CERN OpenLab

- Already achieving ~86Gbps with Infiniband EDR
  - Meets our target
  - Reduced scale setup

- Off-site tests for Infiniband EDR and Intel Omni-Path currently ongoing…

| Data generator |
|---|

| Readout module | Dataflow manager | Builder unit |
|---|---|---|

| Event-builder core |
|---|

| Ethernet driver | InfiniBand driver | Omni-Path driver |
|---|---|---|

# 128node InfiniBand scalability (CINECA)



- Cores allocated: 8/16

- Size of buffer: ~128KB

- Average bw: 24.98 Gb/s

Performance limited by non-exclusive utilization of cluster

Matteo Manzali - INFN CNAF - Università degli Studi di Ferrara

# Future data centre at Point 8

| | Turnkey commercial solution (Requires minimal support from CERN engineering groups) | Leverage existing infrastructure at Point 8 |
|---|---|---|
| Building | Buy pre-fabricated containers from a commercial supplier | Accommodate the farm in an existing building (SX8 hall) |
| Cooling | Cooling solution depends on the vendor (e. g. free air cooling) | Passive rear-door heat exchangers using primary water from existing cooling towers<br>▪ Compatible with DCLC for hot spots<br>▪ Test setup at the pit to evaluate performance with on site "warm" cooling water |

A review to decide the most cost-effective solution will take place in April.

# Online mile-stones

✓ Q1 2016: long-distance fibres validation (successful)

✓ Q2 2016: event-building with 32-port switch (successful)

- Thanks to the Bologna, CNAF and Ferrara teams
→ would be good to repeat this in a lab-setup (better tuning etc…)

❑ Q2 2016: event-builder scalability to 600 nodes

✓Achieved with 128 nodes, larger scale test still this year
 need to get a slot from external HPC sites (France, US…)

❑ Q2 2017: data-centre solution defined

❑ Q4 2017: full DAQ network test

❑ Q3 2018: network technology decision

❑ Q3 2019: ready for commissioning

# Thank you for your time

# Setup



- Review of TDR suggested a test installation at P8 to ensure the viability of the long distance read-out
- Installed 3 Trunk cables with 144 fibers each from different providers at P8
  - 1 x 144 Fibers from Fibernet, pre-connectorized
  - 1 x 144 Fibers from CERN, pre-connectorized
  - 1 x 144 Fibers from CERN, spliced
- All 432 Fibers have been tested and have expected Attenuation
- Currently Measuring the spliced fibers
  - Worst case scenario due to additional attenuation of splice point
- Testing on loop with 2 x 350m
  - Higher attenuation than in final setup
  - AMC40 transmitter has more optical power though
  - Longer fiber just about compensates for the stronger transmitter

# Result so far

- Setup has been running for over 9 months now
- $10^{18}$ bits tested so far
- 0 Errors
- BER $< 10^{-18}$ with 95% confidence

$$CL = 1 - e^{-N \times BER_s} \times \sum_{k=0}^{E} \frac{(N \times BER_s)^k}{k!}$$

CL = Confidence Level
N = Number of bits tested
E = Errors
$BER_s$ = Specific Error Rate

| Confidence | 63% | 74% | 95% |
|---|---|---|---|
| Aggregated BER | $< 3.1 \times 10^{-18}$ | $< 4 \times 10^{-18}$ | $< 9 \times 10^{-18}$ |
| Error rate for 12k fibers | $< 16$ / day | $< 20$ / day | $< 45$ / day |

# Timing distribution at LHCb



- Unidirectional
- 1:32 max split ratio
- Low bandwidth
- Obsolete components
- Very low payload for Bunch synchronous data

# Upgrade proposal 1:TTC over GBT

From TFC point of view, we ensure constant:
- ✓ LATENCY: Alignment with BXID
- ✓ FINE PHASE: Alignment with best sampling point

LHC Clocks

S-ODIN

TFC+ECSInterface

TELL40s

■ = Receiver
■ = Transmitter

Some resynchronization mechanisms envisaged:
→ Within TFC boards
→ With GBT
  ✓ No impact on FE itself

GBT for TFC+ECS

GBT for DATA

Loopback mechanism:
→ re-transmit TFC word back
→ allows for latency measurement + monitoring of TFC commands and synchronization

FE ASIC

Front-Ends

**Courtesy: Federico Alessio (CERN)**

# Upgrade proposal 2:TTC over PON

- Major topology of growing Access Network Market known as:
  - Fiber To The X (FTTx) (http://en.wikipedia.org/wiki/Fiber_to_the_x)
- Point-to-MultiPoint (P2M)
- One single fibre in charge of both downstream and upstream transmissions (using wavelength multiplexing technique)



1490nm

1310nm

FTTH

FTTC

FTTB

**Courtesy S. Baron et al. (CERN)**
https://indico.cern.ch/event/400951/

# TTC-PON

- Standard: XG-PON1
- Fixed latency downstream
- Line rate
  - Down: 9.6 Gb/s
  - Up: 2.4 Gb/s
- Wavelength:
  - Down: 1577nm
  - Up: 1270nm
- Bidirectional
- Range: 1-1000m
- Split ratio: 1:128
- Error detection

OLT: Optical Line Terminal
ONU: Optical Network Unit

# AMC40 readout board

- Philosophy: one common board for DAQ and timing/control

- 1st generation readout board
  - Developed at CPP Marseille

- Mezzanine for ATCA crate

- Altera Stratix V FPGA

- ≤36 bidirectional optical links
  - Avago MiniPODs
  - GBT / GWT / 10GbE protocols

- Produced and made available to the collaboration (MiniDAQ1)

- Firmware designed for migration to 2nd generation readout board

# PCIe40 status

**Current status**

- Low-Level Interface available
  - Slow control through PCIe
  - 100Gbps DMA
  - On-board communication (I2C, SPI)
  - Filtered PLL clock
  - CvP (programming via PCIe)
  - Temperature and current monitoring
- New custom heat-sink design
- Procurement ongoing after tender

**Next plans**

- x24 GBT integration
- PON tests for timing distribution
- Optical links characterization
  - Eye aperture
  - BER
- Cooling characterization
  - *f( firmware, airflow, heat-sink )*
- Monitoring via IPMI

# PCIe Gen3 DMA performance

## 1 interface (optimized)

**DMA rate sampling (12h run)**



## 2 interfaces (unoptimized)

# Firmware status

## AMC40

- Readout over GBT / WideBus
- Fixed-size frontend protocol
- 6 optical links
- Front-end configuration via GBT-SCA
- Data monitoring inside FPGA
- UDP output protocol over 10GBASE-R
- Embedded FE data generator
- Data verification outside FPGA
- Configurable data rate throttle
- ☐ Variable-size FE protocol

## PCIe40

… all of the previous, plus:

- HW-specific low level interface
- PCI Express Gen3 DMA interface & driver

- ☐ Increase number of optical links
- ☐ Integrate LLI & DMA into firmware
- ☐ Implement remaining frontend protocols

# Event-building software performance

- DAQPIPE now supports following APIs:

- **Libfabric**

- **MPI**

- **TCP**

- Implemented automatic benchmarking and reporting

- (with web-based interface)

- On-the-fly monitoring (htopml)

- Benchmarks of

- Intel **Omni-Path**

- Infiniband **EDR**

      currently in progress!

- Already observing **86 Gb/s** on **EDR**