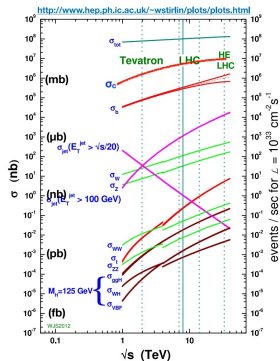# GPU-based software trigger for LHCb experiment
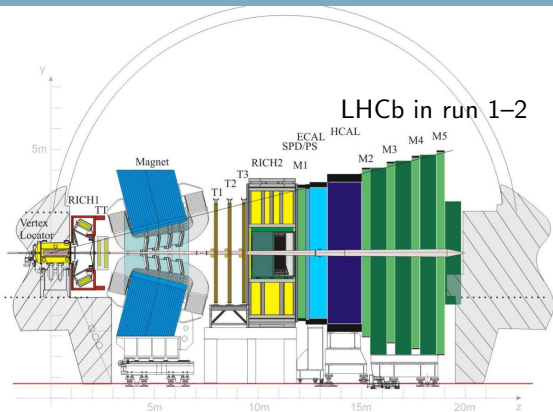
Anton Poluektov
on behalf of LHCb collaboration
RTA project

Aix Marseille Univ, CNRS/IN2P3, CPPM, Marseille, France
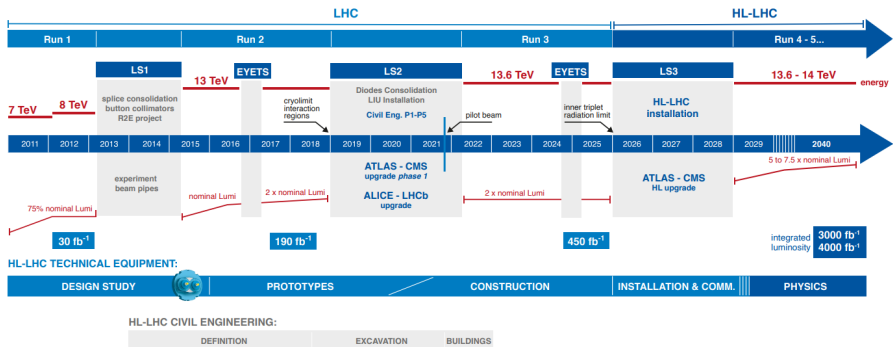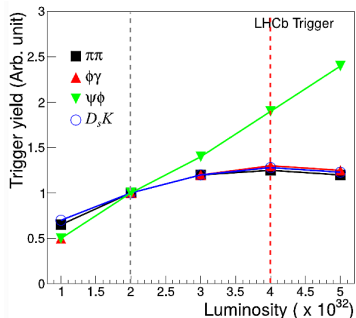
9 April 2024

LHCb in run 1–2

Forward spectrometer, optimised for *b* and *c* decays. $2 < \eta < 5$

- Excellent vertex resolution (weak decays)
- High-precision tracking before and after the magnet
- PID in broad range of momenta $3 < p < 150\,\mathrm{GeV}$
- Efficient trigger, including fully-hadronic final states, $\sim$12 kHz output rate
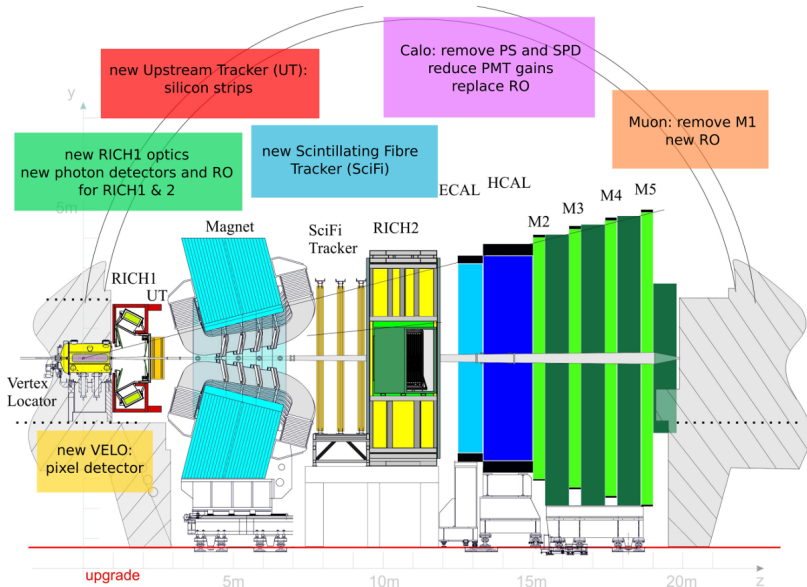
- LHC Run 2 finished in 2018
  - LHCb: $\int \mathcal{L}dt = 9\,\mathrm{fb}^{-1}$ collected in 2010-2018
- Long shutdown until 2022: upgrade of the machine and detectors
  - LHCb Upgrade I: major upgrade/replacement of the subsystems and readout
- Run 3 until 2026 $\to$ HL-LHC upgrade $\to$ Run 4 . . .
  - LHCb goal: $50\,\mathrm{fb}^{-1}$ by the end of Run 4 $\to$ Upgrade II
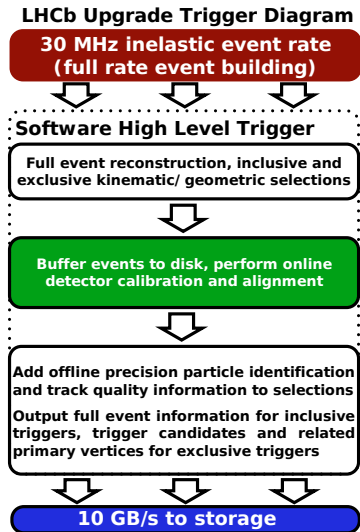
# LHCb upgrade case



- Instantaneous luminosity:
  $4 \times 10^{32}$ (Run 2) $\rightarrow 2 \times 10^{33}\,\mathrm{cm^{-2}\,s^{-1}}$
- Run 1–2 trigger:
  - First stage: hadrware L0 (40→1 MHz) using high $p_T/E_T$ signatures
  - 1 MHz limit saturates hadronic modes already in Run 2
    (higher rate ⇒ higher thresholds)
- The only solution: read full event at bunch-crossing rate and apply track reconstruction/IP selections.
- Upgrade/replace subsystems:
  - Cope with higher occupancy.
  - Faster/higher precision tracking
- Fully replace DAQ and trigger.

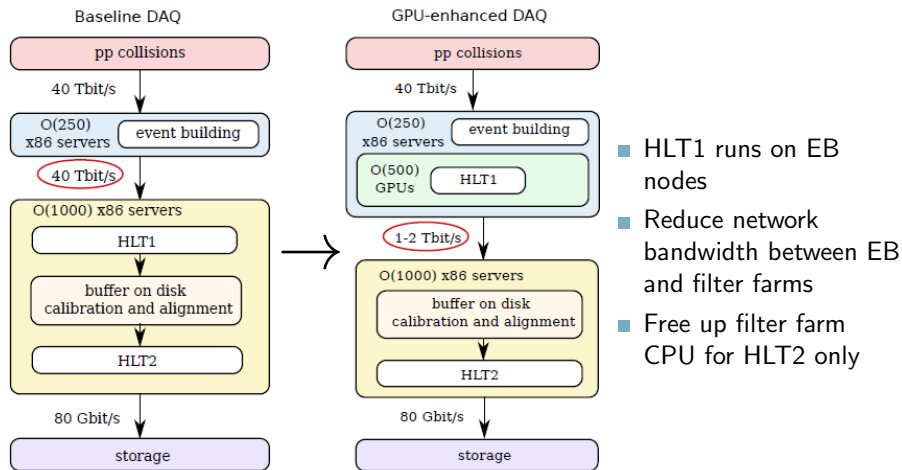Complete replacement of DAQ, fully software trigger (HLT1 + HLT2)

# Upgraded DAQ+trigger: functional diagram

**LHCb Upgrade Trigger Diagram**



**30 MHz inelastic event rate (full rate event building)**

**Software High Level Trigger**

Full event reconstruction, inclusive and exclusive kinematic/ geometric selections

Buffer events to disk, perform online detector calibration and alignment

Add offline precision particle identification and track quality information to selections

Output full event information for inclusive triggers, trigger candidates and related primary vertices for exclusive triggers

**10 GB/s to storage**
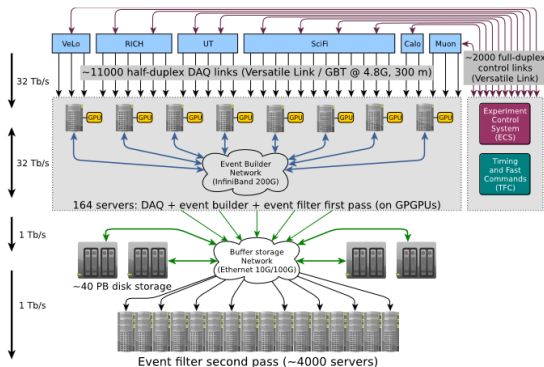
HLT1:    [LHCb upgrade computing TDR]

- Subdetector reconstruction:
  - VELO: clustering, tracking, vertex reconstruction
  - UT, SciFi: tracking
  - Muon: Hit-track matching
- Global event reconstruction:
  - Track fit (Kalman filter)
  - Reconstruction of secondary vertices
- Selections:    [LHCb-PUB-2019-013]
  - Single displaced tracks
  - Two-track displaced vertices
  - Single displaced muons
  - Low-mass displaced two-muon vertices
  - High-mass dimuons

Baseline CPU-based design was replaced by GPU-accelerated one



- HLT1 runs on EB nodes
- Reduce network bandwidth between EB and filter farms
- Free up filter farm CPU for HLT2 only

Warning: the exact numbers for BW, N(servers) have evolved since then

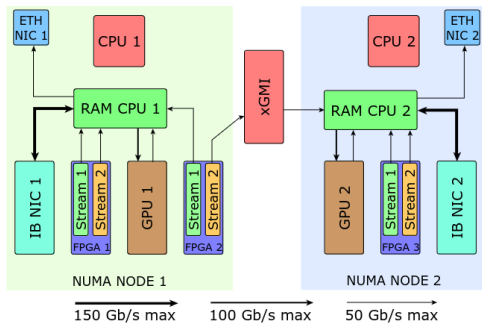# Upgraded LHCb DAQ: current implementation



- Event rate: $30\,\mathrm{MHz}$ non-empty bunch crossing
- Event size: $\sim 100$ kB
- Input bandwidth: $\sim 32$ Tbit/s

- New PCIe40 readout boards
  - 24 optical inputs, PCIe interface
- Event builder network using commercial technology
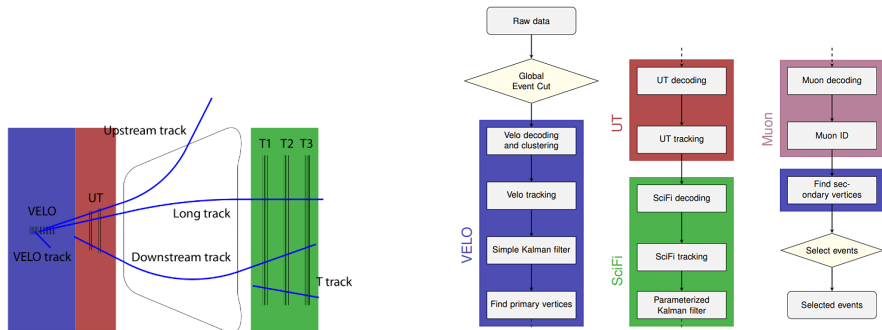  - 200 Gbit/s InfiniBand© network with remote direct memory access

Current configuration: 164 2-CPU server nodes



2-CPU server node hardware diagram

- **CPU:** ×2 AMD EPYC 7502, 32 cores
- **GPU:** ×2 NVIDIA RTX A5000
- **RAM:** 512 GB DDR4

- **Network:** ×2 NVIDIA ConnectX-6 HDR
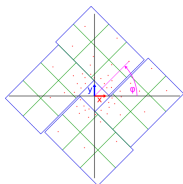- **Readout:** ×3 PCIe40 FPGA boards

- Framework for GPU-based execution of an algorithm sequence
  [GitLab repo], [Documentation]
- Cross-architecture compatibility:
  Runs on CPU, NVidia GPU (`CUDA`), AMD GPU (`HIP`)
- Algorithm sequences defined in `python`, generated at runtime
- Three levels of parallelism:
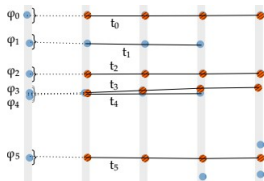  Intra-collision (tracks, clusters), collisions, collision batches

# Allen project: parallel algorithms

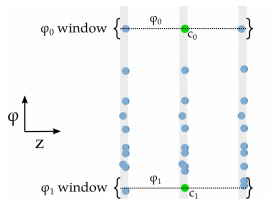Fast parallel algorithms developed for tracking, vertexing etc.
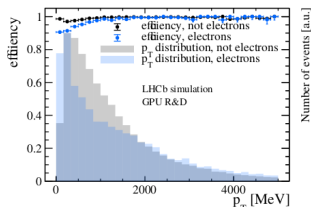
E.g. reconstruction of tracks in VELO:

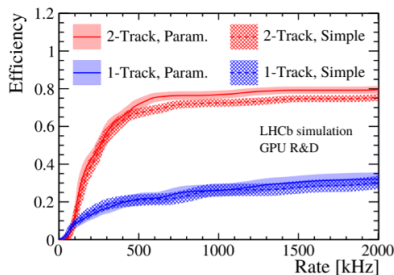

1. Sorting hits in $\phi$



2. Triplet seeding



3. Triplet forwarding



VELO tracking performance

[D. Campora, N. Neufeld, A. Riscos Nez, IPDPSW 2019]

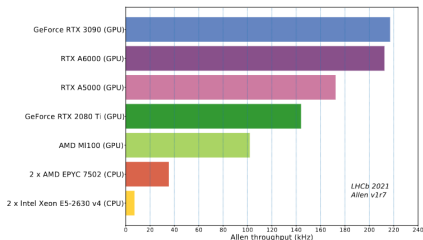| Trigger | Rate [kHz] |
|---|---|
| 1-Track | $215 \pm 18$ |
| 2-Track | $659 \pm 31$ |
| High-$p_T$ muon | $5 \pm 3$ |
| Displaced dimuon | $74 \pm 10$ |
| High-mass dimuon | $134 \pm 14$ |
| Total | $999 \pm 38$ |



Rates of HLT1 lines on minimum bias events

Efficiency of 1-Track and 2-Track selections with $B_s^0 \to \phi\phi$ MC

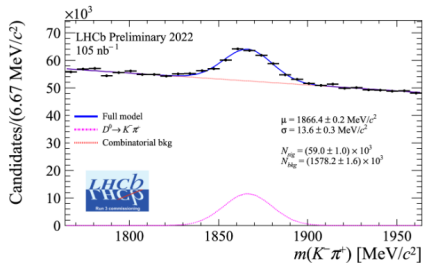| Signal | GEC | TIS -OR- TOS | TOS | GEC $\times$ TOS |
|---|---|---|---|---|
| $B^0 \to K^{*0}\mu^+\mu^-$ | $89 \pm 2$ | $91 \pm 2$ | $89 \pm 2$ | $79 \pm 3$ |
| $B^0 \to K^{*0}e^+e^-$ | $84 \pm 3$ | $69 \pm 4$ | $62 \pm 4$ | $52 \pm 4$ |
| $B_s^0 \to \phi\phi$ | $83 \pm 3$ | $76 \pm 3$ | $69 \pm 3$ | $57 \pm 3$ |
| $D_s^+ \to K^+K^-\pi^+$ | $82 \pm 4$ | $59 \pm 5$ | $43 \pm 5$ | $35 \pm 4$ |
| $Z \to \mu^+\mu^-$ | $78 \pm 1$ | $99 \pm 0$ | $99 \pm 0$ | $77 \pm 1$ |

Efficiencies of HLT1 selection for benchmark signals

HLT1 throughput for various GPU cards.
[LHCB-FIGURE-2020-014]



$D^0 \to K^- \pi^+$ peak directly from HLT1 (2022)
[LHCB-FIGURE-2023-009]

# HLT2 signal rates



http://www.hep.ph.ic.ac.uk/~wstirlin/plots/plots.html

- Signal rates at $\mathcal{L} = 2 \times 10^{33}\,\mathrm{cm}^{-2}\,\mathrm{s}^{-1}$:
  - $O(10)\,\mathrm{MHz}$ charm
  - $O(1)\,\mathrm{MHz}$ beauty

- Output bandwidth limited to 10 GB/s. Up to $100\,\mathrm{kHz}$ with full event size of 100 kB.

- Need to reduce the event size for higher rate

# Persistency model

Selective persistency: write out only the "interesting" part of the event.



- Turbo stream:
  - Minimum output: only HLT2 signal candidates

  Limitations: cannot refit tracks and PVs offline, rerun flavour tagging etc.
  Advantage: Event size $O(10)$ smaller than RAW

Selective persistency: write out only the "interesting" part of the event.



- Turbo stream:
  - Minimum output: only HLT2 signal candidates
  - Optionally: (parts of) *pp* vertex (*e.g.* "cone" around candidate for spectroscopy searches)

  Limitations: cannot refit tracks and PVs offline, rerun flavour tagging etc.
  Advantage: Event size $O(10)$ smaller than RAW

# Persistency model

Selective persistency: write out only the "interesting" part of the event.



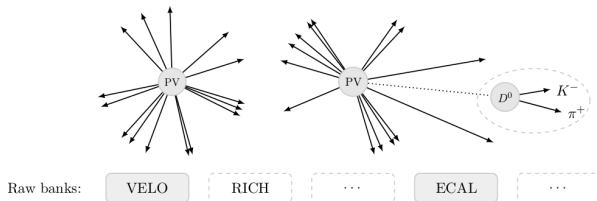Raw banks: VELO | RICH | ⋯ | ECAL | ⋯

- Turbo stream:
    - Minimum output: only HLT2 signal candidates
    - Optionally: (parts of) *pp* vertex (*e.g.* "cone" around candidate for spectroscopy searches)

    Limitations: cannot refit tracks and PVs offline, rerun flavour tagging etc.

    Advantage: Event size $O(10)$ smaller than RAW
- FULL stream: all reconstructed objects in the event
    - + selected RAW banks

Selective persistency: write out only the "interesting" part of the event.



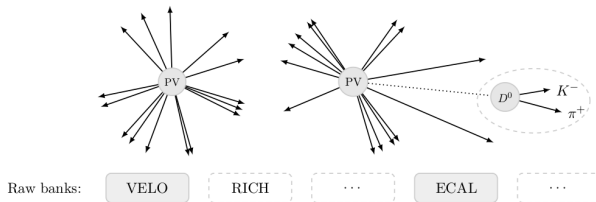Raw banks: VELO RICH ⋯ ECAL ⋯

- `Turbo` stream:
    - Minimum output: only HLT2 signal candidates
    - Optionally: (parts of) *pp* vertex (*e.g.* "cone" around candidate for spectroscopy searches)

    Limitations: cannot refit tracks and PVs offline, rerun flavour tagging etc.
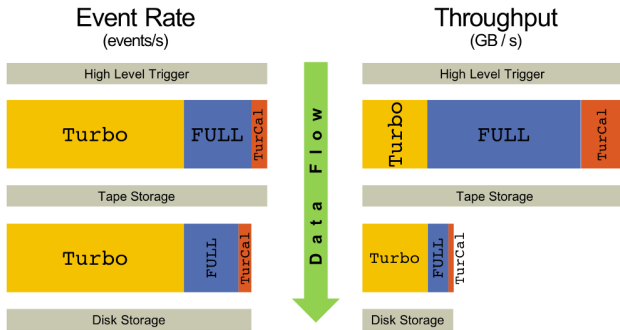    Advantage: Event size $O(10)$ smaller than RAW
- `FULL` stream: all reconstructed objects in the event
    - + selected RAW banks
- `TurCal` stream: HLT2 candidates and selected RAW banks
    Used for offline calibration and performance measurement
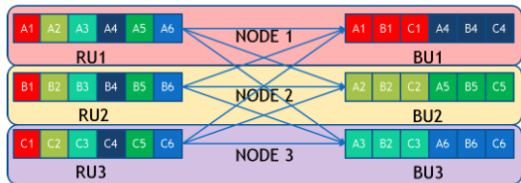
## Rate and bandwidth to tape

| stream | rate fraction | throughput (GB/s) | bandwidth fraction |
|--------|---------------|-------------------|--------------------|
| FULL   | 26%           | 5.9               | 59%                |
| Turbo  | 68%           | 2.5               | 25%                |
| TurCal | 6%            | 1.6               | 16%                |
| total  | 100%          | 10.0              | 100%               |

## Disk bandwidth

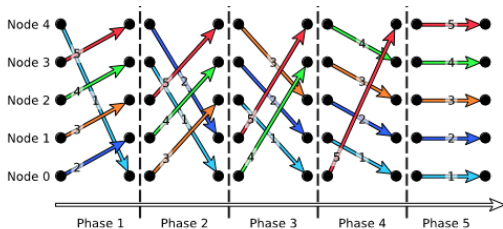| stream | throughput (GB/s) | bandwidth fraction |
|--------|-------------------|--------------------|
| FULL   | 0.8               | 22%                |
| Turbo  | 2.5               | 72%                |
| TurCal | 0.2               | 6%                 |
| total  | 3.5               | 100%               |

# Summary

- LHCb started taking data after upgrade in 2022 (Run 3)
    - Commissioning with LHC in 2022
    - 2023: run with open VELO after incident; UT commissioning
    - Plan to run with maximum performance in 2024–2025
- Aim to increase instantaneous luminosity to $2 \times 10^{33}\,\mathrm{cm}^{-2}\,\mathrm{s}^{-1}$ (5 times pre-upgrade).
- Major redesign of readout and trigger compared to Run 2
- Remove hardware L0 stage, read out full detector at $30\,\mathrm{MHz}$ non-empty bunch crossing rate
    - Need to cope with 32 Tbit/s input bandwidth:
      *highest in any physics experiment to date*
- HLT filtering farm:
    - Architecture of split trigger with disk buffer, alignment and calibration $\Rightarrow$ offline-quality output.
    - GPU-based HLT1 stage in the event builder farm
    - CPU-based HLT2
    - Increase physics output by moving most of signal rate to `Turbo` stream (reduced size, no RAW information).
    - 10 GB/s output bandwidth to tape for further analysis

# Backup

# Event builder



Event building process



Linear shifting scheduling