# Allen: A High Level Trigger on GPUs for LHCb

**Dorothea vom Bruch**

LPNHE, CNRS

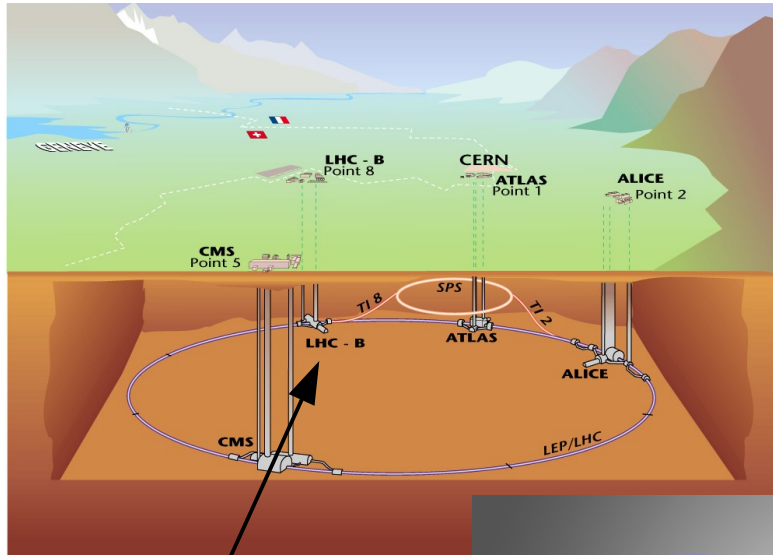Sorbonne University, Paris Diderot University
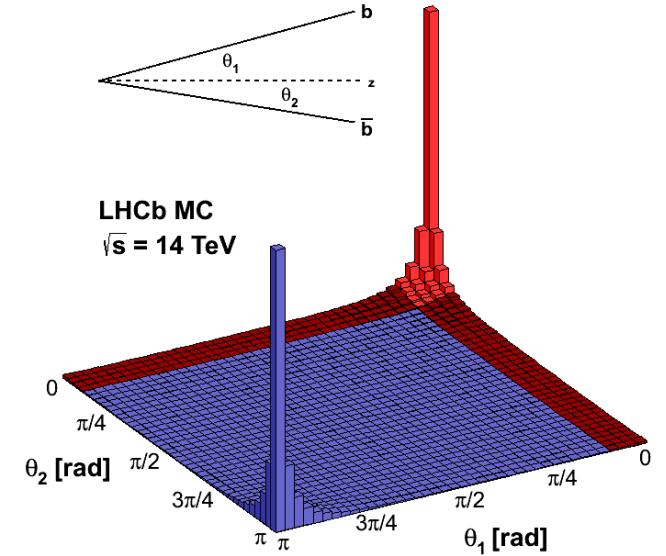
January 27th 2020

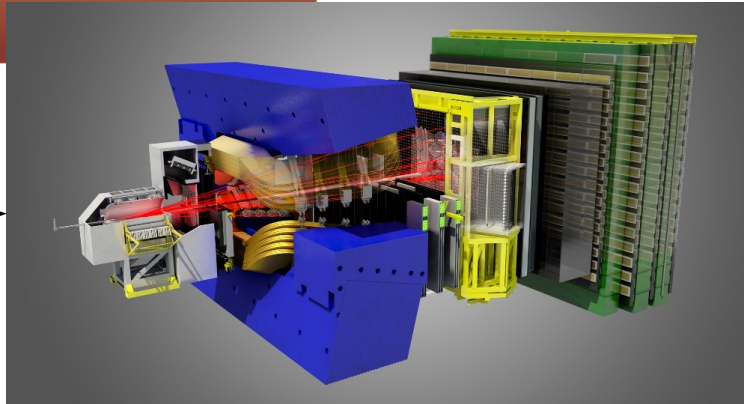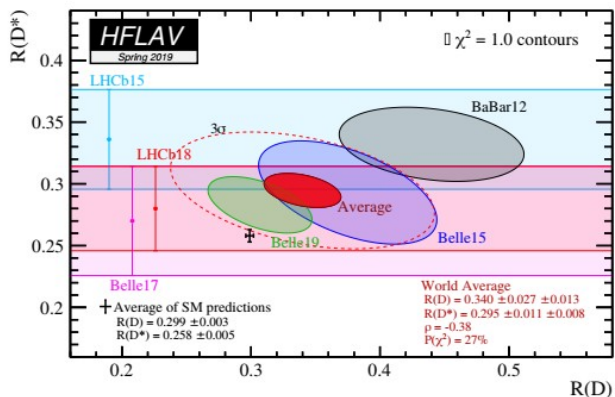LPNHE seminar

# LHCb

LHC @ CERN



General purpose detector in the forward region specialized in beauty and charm hadrons

# Highlights from Runs 1 & 2

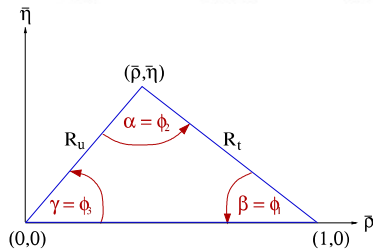## Lepton flavor universality

### b → clv



### b → sl⁺l⁻



T. Humair, Moriond 2019

## Constraining CKM angles



## CP violation in charm decays



Phys. Rev. Lett. 122, 211803 (2019)

## Pentaquarks



3

# Prospects for Run 3 (2021)

- Crucial to reduce uncertainties

  → manifestation of new physics?

- R(D*) is theoretically clean

  → reduction of statistical uncertainty necessary

- Runs 3 and beyond will shed light on the flavour

  anomalies currently observed

Projected absolute uncertainties on R(D$^*$)

Current sensitivities    Run 3    Run 4    Run 5

# LHCb upgrade for Run 3 (2021)



New Vertex locator
(Velo)
tracking detector

New UT
tracking detector

New scintillating fibre
(SciFi)
tracking detector

All readout
systems renewed

RICH 1 redesigned,
new photo detectors
for RICH1 & RICH2

Calorimeters & muons: redundant components
removed, new electronics, more shielding

5

# Reaching the MHz signal era



Run 3: Luminosity of $2 \times 10^{33}$ cm$^{-2}$s$^{-1}$, $\sqrt{s}$ = 14 TeV

# Reaching the MHz signal era



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

Run 3: Luminosity of $2 \times 10^{33}$ cm$^{-2}$s$^{-1}$, $\sqrt{s} = 14$ TeV

- General purpose LHC experiments: jets, electro-weak physics, Higgs physics
- Local characteristic signatures, e.g. high transverse energy
- Can trigger efficiently at ~100 kHz
  → hardware-level trigger possible

# Reaching the MHz signal era


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

Run 3: Luminosity of $2 \times 10^{33}$ cm$^{-2}$s$^{-1}$, $\sqrt{s}$ = 14 TeV


LHCb Simulation

- Too many interesting events
- No "simple" local criteria for selection
  → hardware-level trigger not an option

- General purpose LHC experiments: jets, electro-weak physics, Higgs physics
- Local characteristic signatures, e.g. high transverse energy
- Can trigger efficiently at ~100 kHz
  → hardware-level trigger possible

# Change in trigger paradigm



**Access as much information about the collision as early as possible**

# Why no low level trigger?

Low level trigger on $E_T$ from the calorimeter

Low level trigger on muon $p_T$, $B \rightarrow K^*\mu\mu$



**Need track reconstruction at first trigger stage**

# Tracks in the LHCb detector



**Need information from many subdetectors → read out full detector**

# Run 2 versus Run 3 trigger

**LHCb Run 2 Trigger Diagram**

**40 MHz bunch crossing rate**

**L0 Hardware Trigger : 1 MHz readout, high $E_T/P_T$ signatures**

| 450 kHz $h^{\pm}$ | 400 kHz $\mu/\mu\mu$ | 150 kHz $e/\gamma$ |

**Software High Level Trigger**

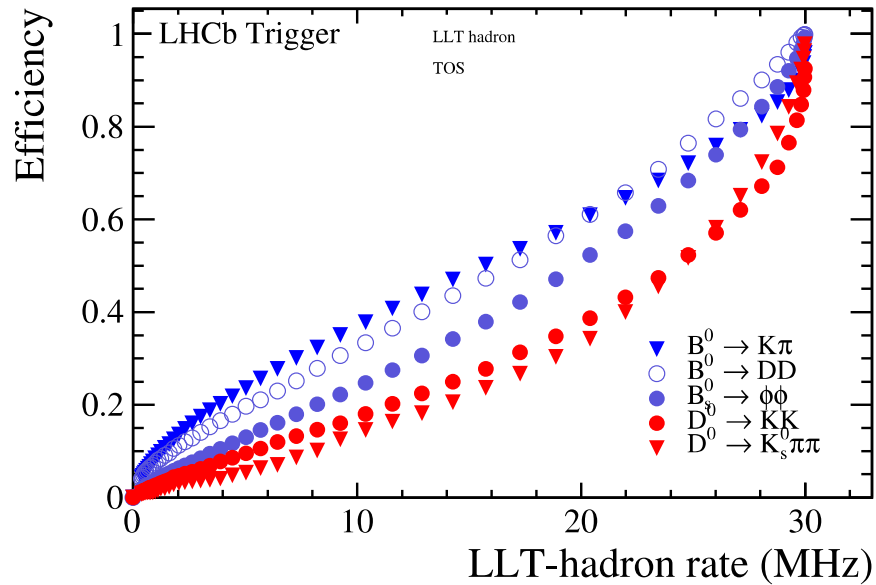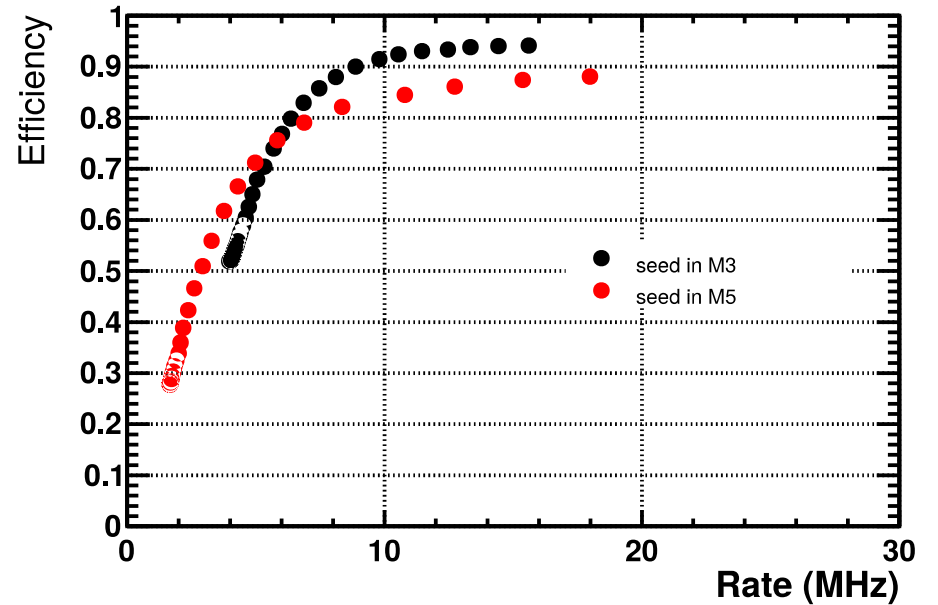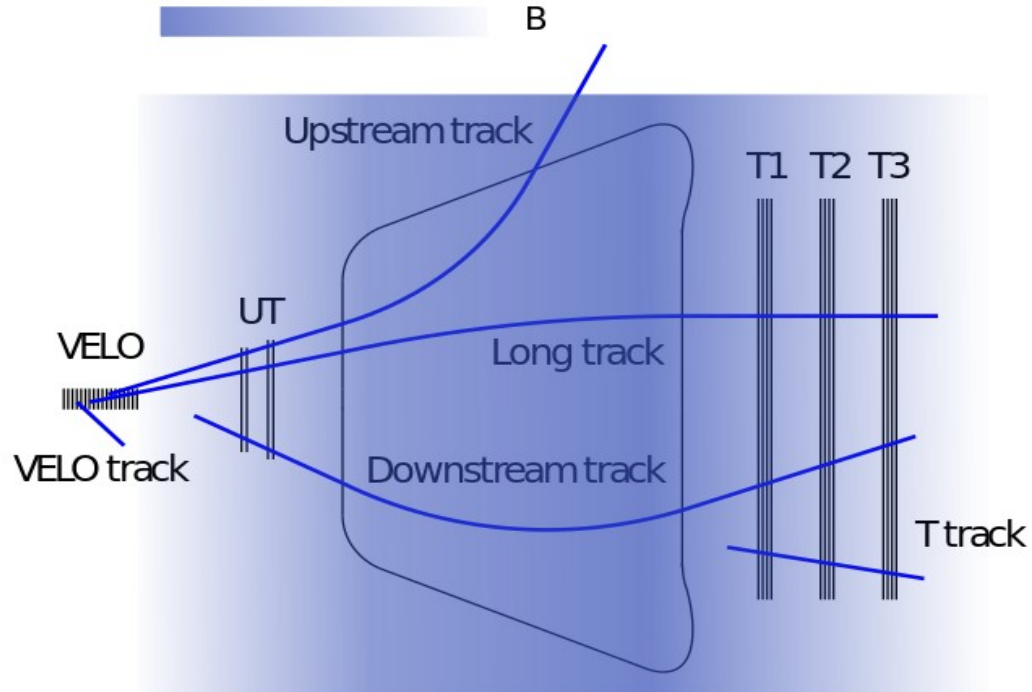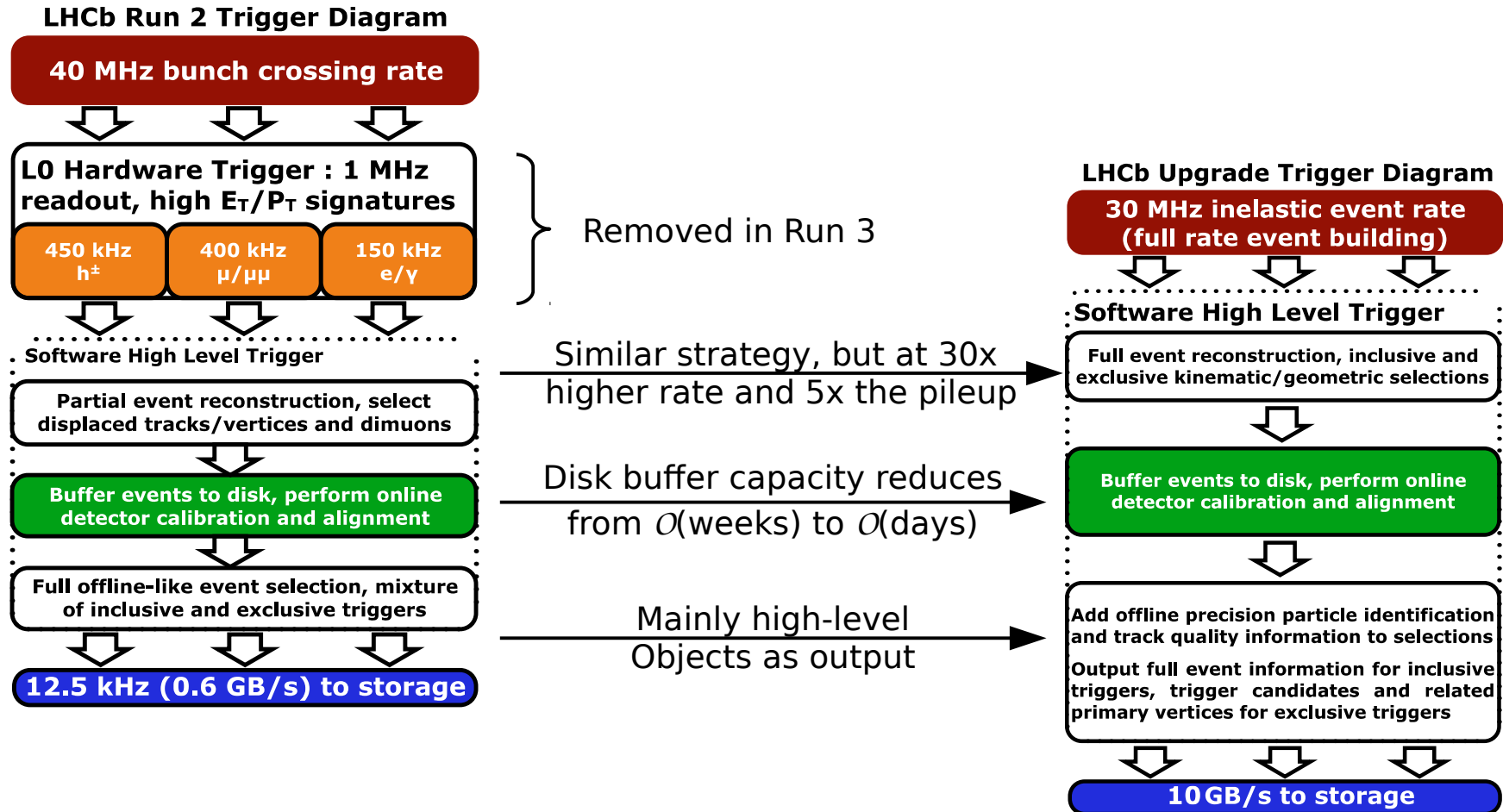**Partial event reconstruction, select displaced tracks/vertices and dimuons**

**Buffer events to disk, perform online detector calibration and alignment**

**Full offline-like event selection, mixture of inclusive and exclusive triggers**

**12.5 kHz (0.6 GB/s) to storage**

Removed in Run 3

Similar strategy, but at 30x higher rate and 5x the pileup

Disk buffer capacity reduces from $\mathcal{O}$(weeks) to $\mathcal{O}$(days)

Mainly high-level Objects as output

**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**

12

# Trigger in Run 3 (2021)

**LHCb Upgrade Trigger Diagram**

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

40 Tbit/s
30 MHz

**Software High Level Trigger**

Full event reconstruction, inclusive and exclusive kinematic/geometric selections

1-2 Tbit/s
1 MHz

**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
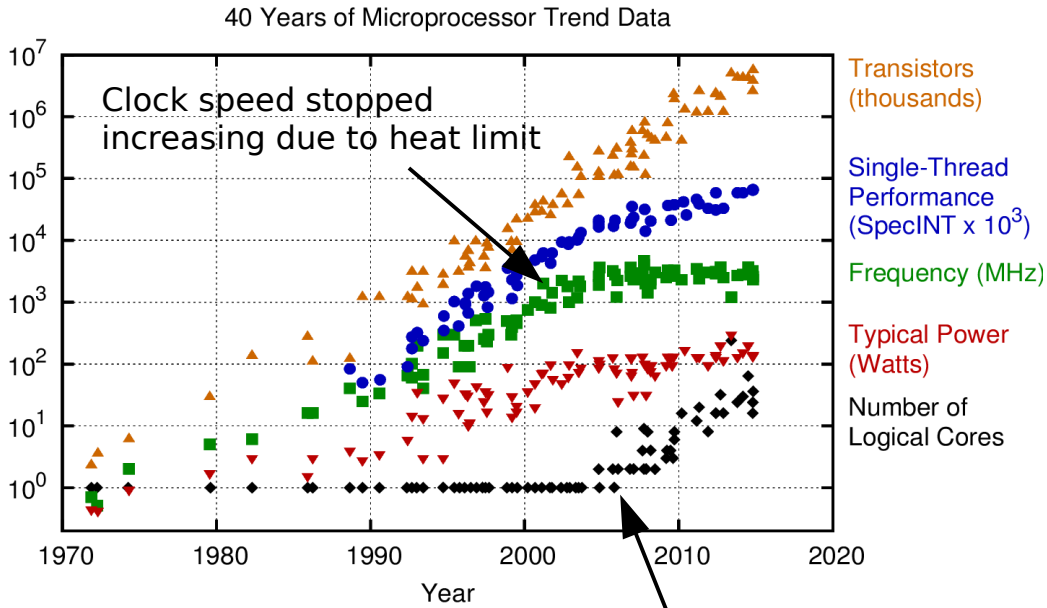
**10 GB/s to storage**

**High Level Trigger 1 (HLT1)**
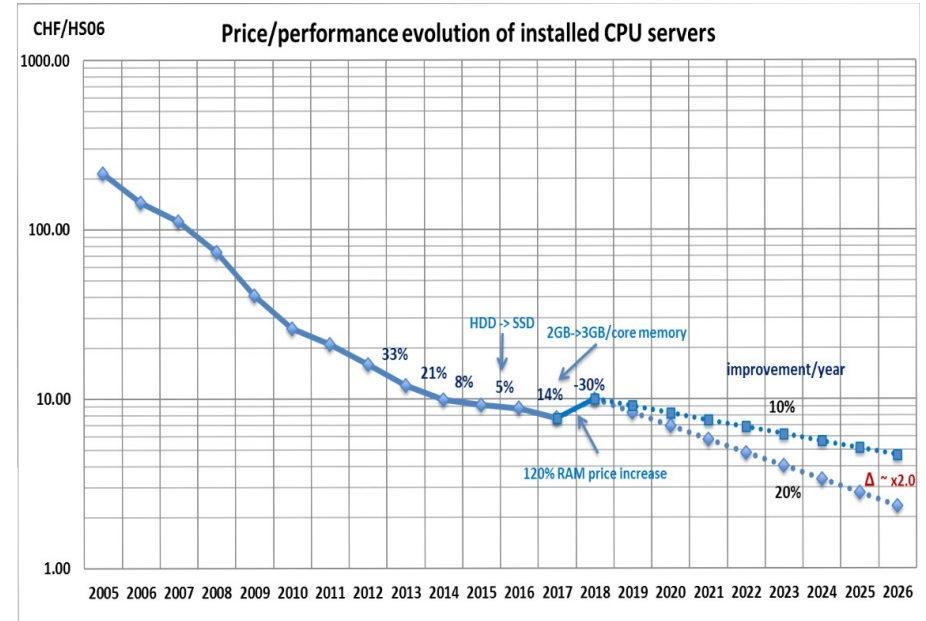- Full charged particle track reconstruction
- Few inclusive single or two-track selections
- Reduce event rate by roughly factor 30

- **High Level Trigger 2 (HLT2)**
- Aligned and calibrated detector
- Offline-quality track reconstruction
- Particle identification
- Full track fitting

# Trigger in Run 3 (2021)

**LHCb Upgrade Trigger Diagram**

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

40 Tbit/s
30 MHz

**Software High Level Trigger**

Full event reconstruction, inclusive and exclusive kinematic/geometric selections

1-2 Tbit/s
1 MHz

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**

**High Level Trigger 1 (HLT1)**

- Full charged particle track reconstruction
- Few inclusive single or two-track selections
- Reduce event rate by roughly factor 30

**Track reconstruction @ 30 MHz is a huge computing challenge!**

# Today's computing landscape



40 Years of Microprocessor Trend Data

Clock speed stopped increasing due to heat limit

Transistors (thousands)

Single-Thread Performance (SpecINT x $10^3$)

Frequency (MHz)

Typical Power (Watts)

Number of Logical Cores

Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2015 by K. Rupp

Multiple core processors emerge



Price/performance evolution of installed CPU servers

CHF/HS06

HDD -> SSD    2GB->3GB/core memory

33%  21%  8%  5%  14%  30%

improvement/year

10%

120% RAM price increase

20%

Δ ~ x2.0

# Amdahl's law

Speed-up factor vs N processors

**Parallel fraction**
- 90%
- 75%
- 50%

$$\text{Speedup in latency} = 1 / (S + P/N)$$

S: sequential part of program

P: parallel part of program

N: number of processors

**Can we use the FLOPS available on a GPU to run HLT1 @ 30 MHz?**

# Where to place the GPUs?

Baseline DAQ

```
┌─────────────────────────────┐
│        pp collisions        │
└─────────────────────────────┘
              │
   40 Tbit/s  ▼
┌─────────────────────────────┐
│ O(250)    ┌──────────────┐  │
│ x86 servers│ event building│ │
│           └──────────────┘  │
└─────────────────────────────┘
              │
   40 Tbit/s  ▼
┌─────────────────────────────┐
│ O(1000) x86 servers         │
│  ┌───────────────────────┐  │
│  │         HLT1          │  │
│  └───────────────────────┘  │
│             │               │
│             ▼               │
│  ┌───────────────────────┐  │
│  │    buffer on disk     │  │
│  │ calibration and alignment│ │
│  └───────────────────────┘  │
│             │               │
│             ▼               │
│  ┌───────────────────────┐  │
│  │         HLT2          │  │
│  └───────────────────────┘  │
└─────────────────────────────┘
              │
   80 Gbit/s  ▼
┌─────────────────────────────┐
│           storage           │
└─────────────────────────────┘
```

17

# Where to place the GPUs?



Baseline DAQ

pp collisions

40 Tbit/s

O(250) x86 servers — event building

40 Tbit/s

O(1000) x86 servers
- HLT1
- buffer on disk calibration and alignment
- HLT2

80 Gbit/s

storage

GPU-enhanced DAQ

pp collisions

40 Tbit/s

O(250) x86 servers — event building

O(500) GPUs — HLT1

1-2 Tbit/s

O(1000) x86 servers
- buffer on disk calibration and alignment
- HLT2

80 Gbit/s

storage

# Where to place the GPUs?



Baseline DAQ

pp collisions

40 Tbit/s

O(250) x86 servers — event building

40 Tbit/s

O(1000) x86 servers
- HLT1
- buffer on disk calibration and alignment
- HLT2

80 Gbit/s

storage

GPU-enhanced DAQ

pp collisions

40 Tbit/s

O(250) x86 servers — event building

O(500) GPUs — HLT1

1-2 Tbit/s

O(1000) x86 servers
- buffer on disk calibration and alignment
- HLT2

80 Gbit/s

storage

GPUs naturally integrate into LHCb's DAQ

**If HLT1 can run on 500 GPUs → Save money on network → Buy GPUs instead**

19

# LHCb HLT1 elements



**Velo**
- Decode raw data
- Clustering of measurements
- Track reconstruction
- Primary vertex reconstruction

**SciFi**
- Decode raw data
- Track reconstruction

**Muons**
- Decode raw data
- Match hits to tracks

**UT**
- Decode raw data
- Track reconstruction

Track fit: Kalman filter

Find secondary vertices

**Selections**
- 1-track selection
- 2-track selection
- Based on p, $p_t$, displacement, vertex criteria and muon identification

# How does HLT1 map to GPUs?

| Characteristics of LHCb HLT1 | Characteristics of GPUs |
|---|---|
| Intrinsically parallel problem:<br>  - Run events in parallel<br>  - Reconstruct tracks in parallel | Good for<br>  - Data-intensive parallelizable applications<br>  - High throughput applications |
|  |  |
|  |  |
|  |  |

# How does HLT1 map to GPUs?

| Characteristics of LHCb HLT1 | Characteristics of GPUs |
| --- | --- |
| Intrinsically parallel problem:<br>  - Run events in parallel<br>  - Reconstruct tracks in parallel | Good for<br>  - Data-intensive parallelizable applications<br>  - High throughput applications |
| Huge compute load | Many TFLOPS |
| | |
| | |
| | |

# How does HLT1 map to GPUs?

| Characteristics of LHCb HLT1 | Characteristics of GPUs |
|---|---|
| Intrinsically parallel problem:<br>  - Run events in parallel<br>  - Reconstruct tracks in parallel | Good for<br>  - Data-intensive parallelizable applications<br>  - High throughput applications |
| Huge compute load | Many TFLOPS |
| Full data stream from all detectors is read out → no stringent latency requirements | Higher latency than CPUs, not as predictable as FPGAs |
|  |  |
|  |  |

# How does HLT1 map to GPUs?

| Characteristics of LHCb HLT1 | Characteristics of GPUs |
|---|---|
| Intrinsically parallel problem:<br>  - Run events in parallel<br>  - Reconstruct tracks in parallel | Good for<br>  - Data-intensive parallelizable applications<br>  - High throughput applications |
| Huge compute load | Many TFLOPS |
| Full data stream from all detectors is read out → no stringent latency requirements | Higher latency than CPUs, not as predictable as FPGAs |
| Small raw event data (~100 kB) | Connection via PCIe → limited I/O bandwidth |
| | |

# How does HLT1 map to GPUs?

| Characteristics of LHCb HLT1 | Characteristics of GPUs |
|---|---|
| Intrinsically parallel problem:<br>  - Run events in parallel<br>  - Reconstruct tracks in parallel | Good for<br>  - Data-intensive parallelizable applications<br>  - High throughput applications |
| Huge compute load | Many TFLOPS |
| Full data stream from all detectors is read out → no stringent latency requirements | Higher latency than CPUs, not as predictable as FPGAs |
| Small raw event data (~100 kB) | Connection via PCIe → limited I/O bandwidth |
| Small event raw data (~100 kB) | Thousands of events fit into O(10) GB of memory |

**Perfect fit!**

# The Allen project

- Fully standalone software project: https://gitlab.cern.ch/lhcb/Allen
- Only requirements: a C++17 compliant compiler & CUDA v10.1
- Built-in physics validation
- Configurable sequence, custom memory manager
- Cross-architecture compatibility

<br>

- Project started in February 2018
- Roughly 14 part-time developers, mostly students
- 2 almost full-time
- After 15 months of development time:

  project reviewed as viable solution for Run 3 (starting in 2021)

<br>

- Named after Frances E. Allen

# HLT1 on GPUs

Raw data

Selection decisions

Individual events

Block (0,0)    Block (0,1)    ...    Block (0,n)

Block (1,0)    Block (1,1)    ...    Block (1,n)

Block (m,0)    Block (m,1)    ...    Block (m,n)

copy    stream 1    stream 1

execute    stream 1    stream 1

copy    stream 1    stream 2

execute    stream 1    stream 2

time

Thread (0,0)   Thread (0,1)   ...   Thread (0,N)

Thread (M,0)   Thread (M,1)   ...   Thread (M,N)

Within one block:
intra-event parallelization

# Software framework

## Static scheduler
### For sequence of algorithms

```
SEQUENCE_T(
    velo_estimate_input_size_t,
    prefix_sum_velo_clusters_t,
    velo_masked_clustering_t,
    velo_calculate_phi_and_sort_t,
    velo_fill_candidates_t,
    velo_search_by_triplet_t,
    velo_weak_tracks_adder_t)
```

## Memory manager for
### GPU memory

**Bachelor, Master and PhD students contribute
in only a few months time**

# Recurrent tasks of HLT1

**Raw data decoding**

- Transform binary payload from subdetector raw banks into collections of hits (x,y,z) in LHCb coordinate system
- Parallelizable over events, all subdetectors and readout units

**Track reconstruction**

- Consists of two steps:
  - Pattern recognition: Which hits belong to which track? → Huge combinatorics
  - Track fitting: Done for every track
- Parallelizable over events, combinations of hits, and tracks

$$f(x) = ... +/- ...$$

**Vertex finding**

- Where did proton-proton collisions take place?
- Where did particles decay within the detector volume?
- Parallelizable over events, combinations of tracks

[x 0.2mm]

Tracks from primary vertex

Primary vertex

B⁺

μ⁺

J/Ψ

K⁺

μ⁻

B decay vertex

# Velo detector: clustering

26 planes of silicon pixel detectors

Clustering with bit masks

# Velo detector: track reconstruction

1) Sort hits by φ

2) Triplet seeding

3) Triplet forwarding

Track reconstruction efficiency for tracks originating from B decays



D. Campora, N. Neufeld, A. Riscos Núñez: "A fast local algorithm for track reconstruction on parallel architectures", IPDPSW 2019

# Velo detector: primary vertex reconstruction



beamline

**Florian Reiss (LPNHE)**

Point of closest approach of tracks to beamline



LHCb simulation, GPU R&D

PV candidates

PV reconstruction efficiency



LHCb simulation
GPU R&D

efficiency

multiplicity distribution

track multiplicity of MC PV

# UT detector: track reconstruction

## 4 planes of silicon strip detectors



**UT plane section**

window -1 [
window +1 ]

- ➤ VELO track
- | Activated strip (hit)
- ✕ track extrapolation
- ◄ main window
- ◄ next window
- ◄ previous window
- main sector range
- next sector range
- previous sector range



UTbX
UTbV
UTaU
UTaX
1719 mm
1338 mm
66.8 mm
315 mm
1528 mm

Track reconstruction efficiency for tracks originating from B decays



efficiency
Number of events [a.u.]

- efficiency
- $p_T$ distribution

LHCb simulation
GPU R&D

$p_T$ [MeV]

33

# SciFi detector

12 layers of scintillating fibres
Efficiency of fibres ~ 98-99%

UT track

x u v x    x u v x    x u v x

T1    T2    T3

x u v x

x 3

53cm
SiPM

5.0°

fibres

mirror

fibres

SiPM

Stereoangles

X: 0°    U: -5°    V: +5°

y

z

-5° +5°

x

# SciFi detector: track reconstruction

Track reconstruction efficiency for tracks
originating from B decays

Momentum resolution



**Renato Quagliani (LPNHE)**

# Kalman filter

Improved track description → better impact parameter resolution



- Simple: Simplified Kalman filter with constant momentum assumption
- Param.: Parameterized Kalman filter with momentum estimate from SciFi track reconstruction

# Muon identification

Four multi-wire proportional chambers
Interleaved with iron walls

SciFi track

Muon identification efficiency

# Ingredients for selections

**Primary vertices**



**Secondary vertices**



**Momentum**



**Impact parameter**



**Selections**
- 1-track selection
- 2-track selection
- Based on $p$, $p_t$, displacement, vertex criteria and muon identification

**Tracks**



**Muon identification**



38

# Event selection

| 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 | $999 \pm 38$ |

**Event rate reduced from 30 MHz to 1 MHz**

# Selection efficiencies

## Allen

Selection efficiencies, values given in %

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

TIS: Trigger independent from signal
TOS: Trigger on signal

**Consistent physics performance with TDR,
which assumed running on x86 architecture**

## Technical Design Report



$B^0 \to K^*[K^+\pi^-]\mu^+\mu^-$

LHCb
Simulation

$B_s^0 \to \phi[K^+K^-]\phi[K^+K^-]$

LHCb
Simulation

# Full HLT1 running on GPUs

Physics performance matches HLT1 requirements

What about the throughput performance?

# Throughput on various GPUs



Full HLT1 sequence

**The system can run on 500 GPUs
→ network cost saving → no additional cost from using GPUs**

# Allen publication

- First publication submitted: arXiv:1912.09161

# Integration test with event building server

Impact on event building when running Allen?



Monitoring temperatures, memory bandwidths, processing rate, …

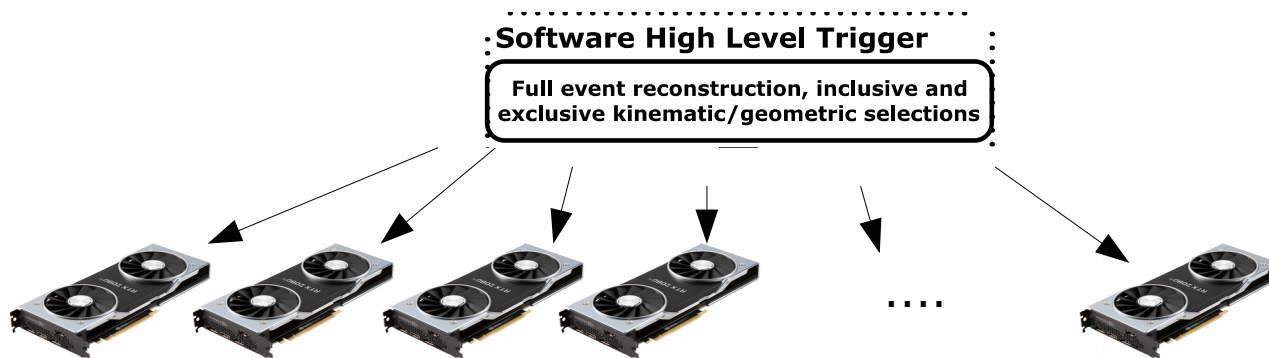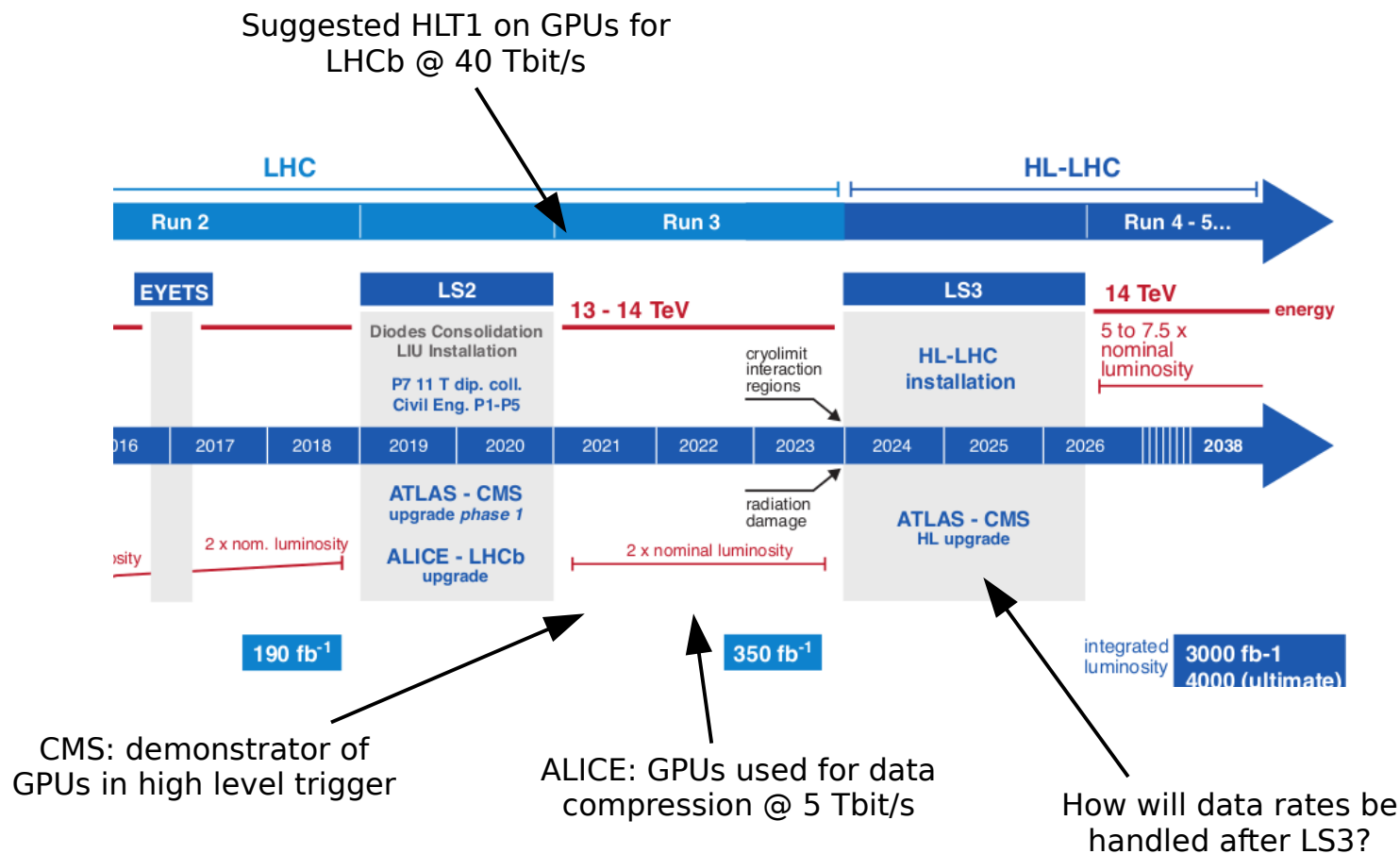# The Allen team

# Summary

- Allen is the first complete high throughput trigger implementation on GPUs

- Baseline HLT1 can run on GPUs

- Efficient selections enable full exploitation of statistics in Run 3

  → crucial to explore new physics scenarios

- Scaling of GPU performance → maximize physics discovery potential of LHCb
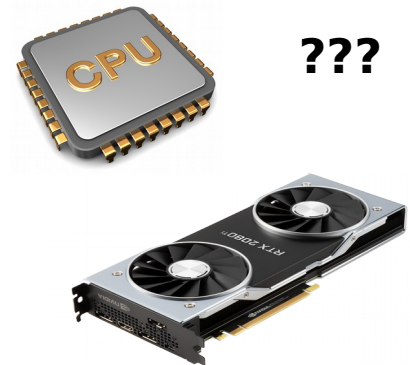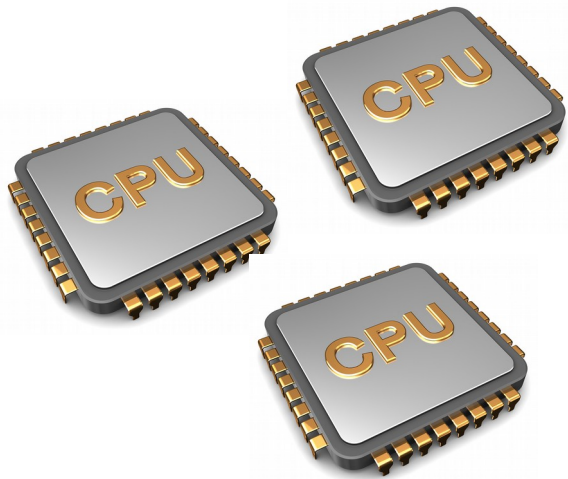
# LHC Schedule



Suggested HLT1 on GPUs for LHCb @ 40 Tbit/s

CMS: demonstrator of GPUs in high level trigger

ALICE: GPUs used for data compression @ 5 Tbit/s

How will data rates be handled after LS3?

# Outlook

- Features of Allen are not tied to HLT1 for LHCb
- High-throughput applications can profit from a similar setup
- Ideas for parallelization of algorithms can be useful for other applications

- Higher luminosity not only challenges real-time event selection
- Simulation production also requires major computing resources
- Need common effort to make best use of heterogeneous computing within high energy physics

# Backup

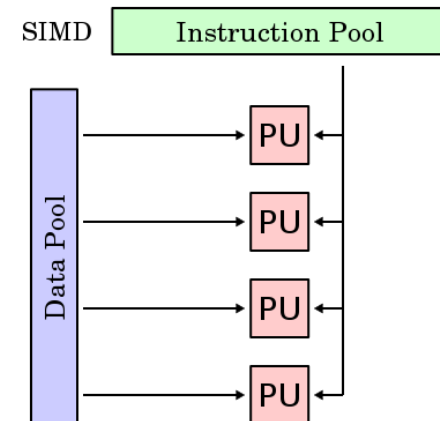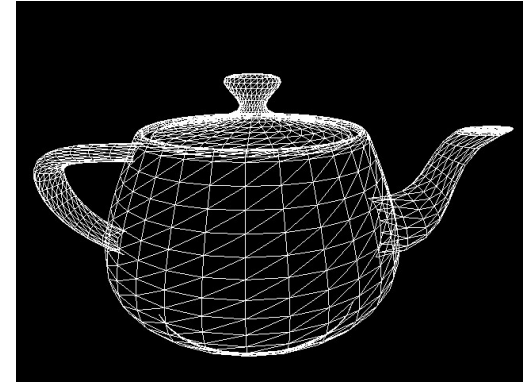# Graphics requirements

**Graphics pipeline**

- Huge amount of arithmetic on independent data:
  - Transforming positions
  - Generating pixel colors
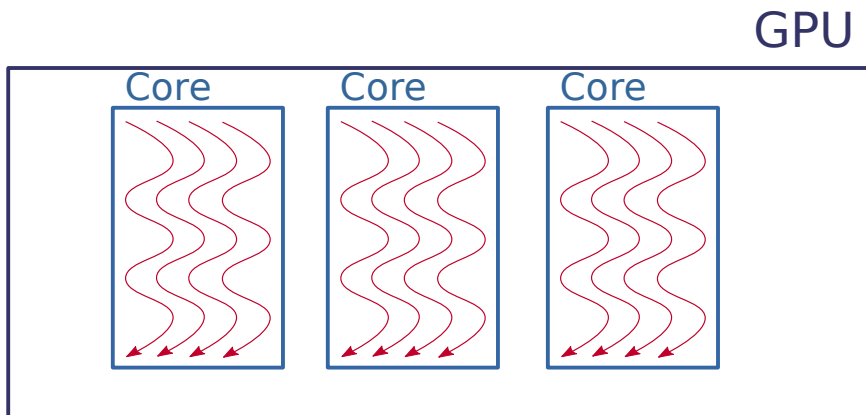  - Applying material properties and light situation to every pixel



**Hardware needs**

- Access memory simultaneously and contiguously
- Bandwidth more important than latency
- Floating point and fixed-function logic
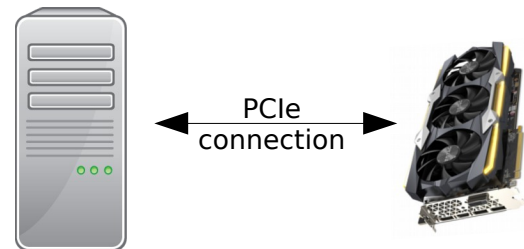
  → Single instruction applied to multiple data: SIMT

# GPU in a nutshell

- Core: multiple SIMT threads grouped together
- GPU: many cores grouped together

**Data transfer to a GPU**



GPU



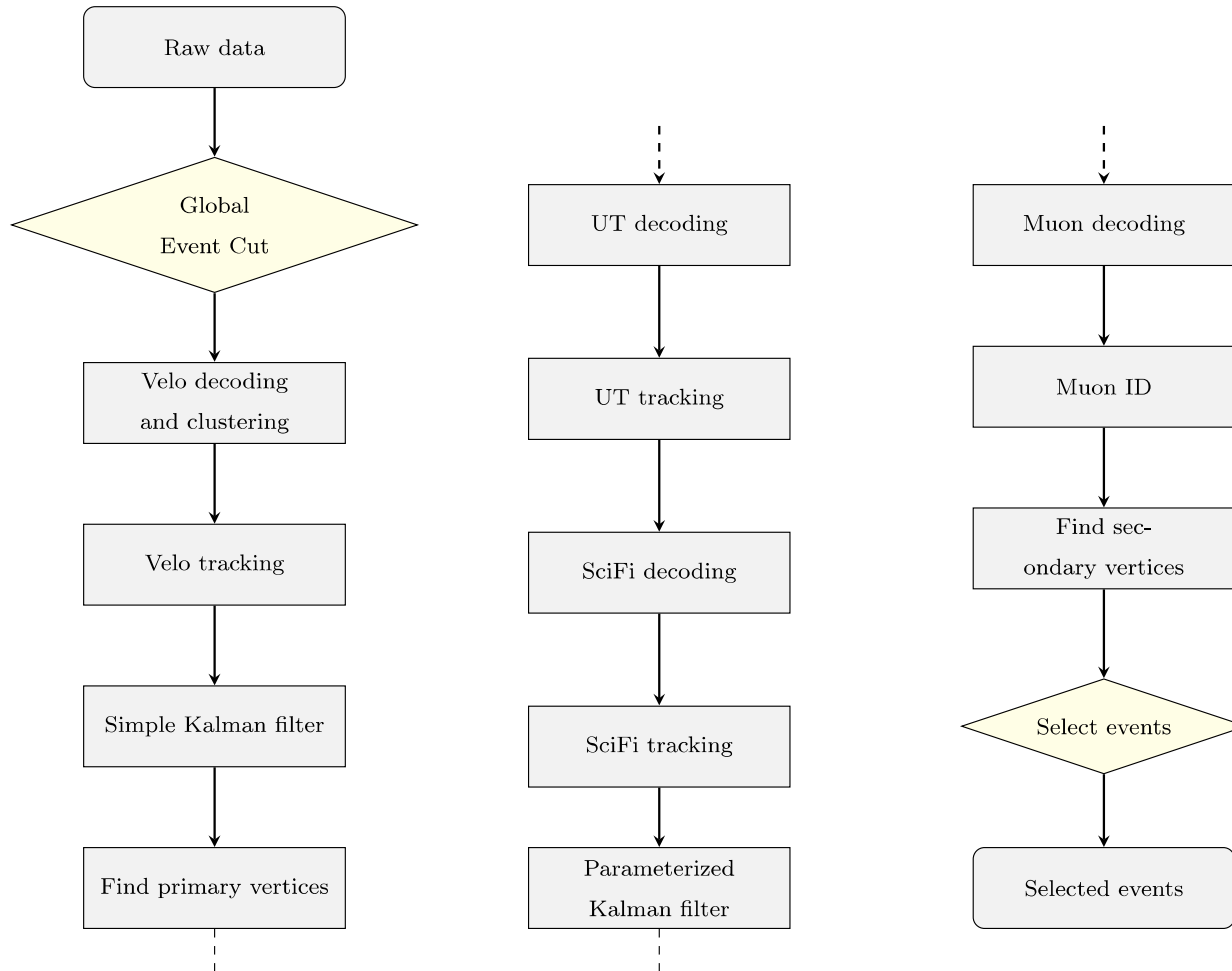| PCIe generation | 16 lanes | Year |
|:---:|:---:|:---:|
| 3.0 | 15.75 GB/s | 2010 |
| 4.0 | 31.5 GB/s | 2017 |

# Selections

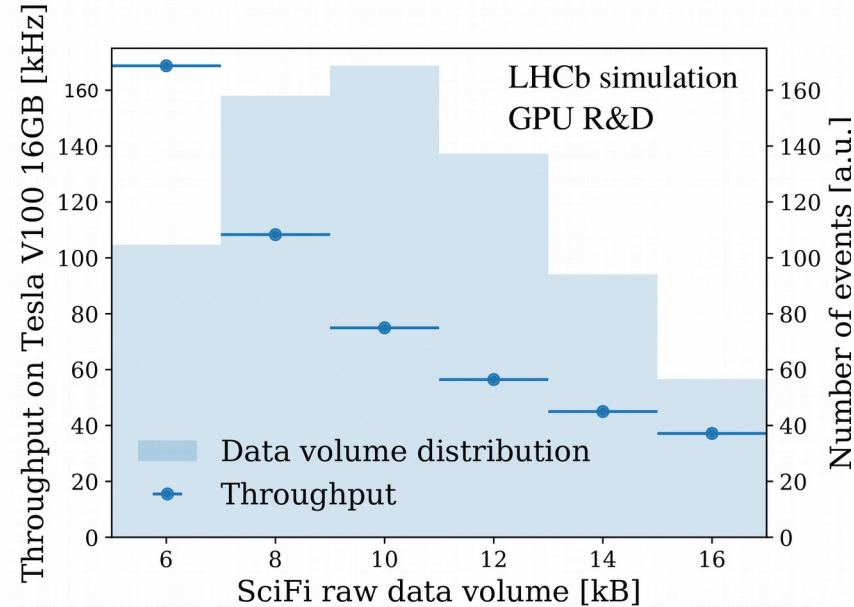| Selection name | Criteria |
| --- | --- |
| 1-Track | Single displaced track with high $p_T$ |
| 2-Track | Two-track vertex with significant displacement and $p_T$ |
| High-$p_T$ muon | Single muon with high $p_T$ |
| Displaced diumuon | Displaced di-muon vertex |
| High-mass dimuon | Di-muon vertex with mass near or larger than the J/Ψ |

Criteria applied to signal decays in efficiency calculations

| $b$ and $c$ hadrons | $p_T > 2$ GeV |
| --- | --- |
| | $\tau > 0.2$ ps |
| $b$ and $c$ hadron children | $p_T > 200$ MeV |
| | $2 < \eta < 5$ |
| | reconstructible in the Velo and SciFi detector (long track) |
| $Z$ children | $p_T > 20$ GeV |
| | $2 < \eta < 5$ |
| | reconstructible in the Velo and SciFi detector (long track) |

# HLT1 algorithms in Allen

# Throughput versus occupancy



- Data volume proportional to occupancy
- Low performance decrease at high occupancy

  → will be able to handle real data (likely higher in occupancy than simulation)
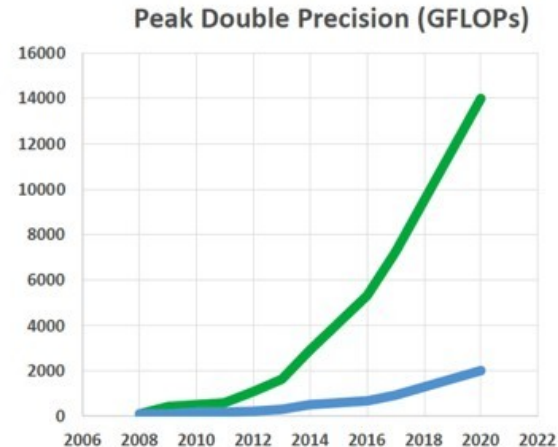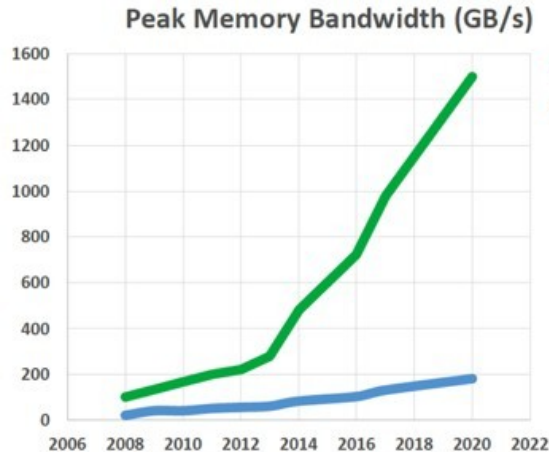
# GPUs for throughput measurement

CUDA streams

| Allen settings | Threads (-t) | Memory (-m) | Number of events (-n) | Repetitions (-r) |
|---|---|---|---|---|
| High | 12 | 700 | 1000 | 100 |
| Low | 2 | 700 | 1000 | 100 |

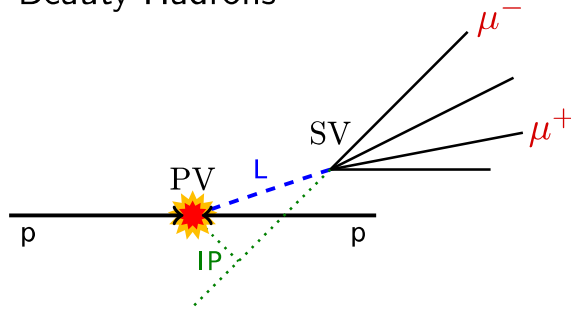| Card | # cores | Max freq. (GHz) | Cache (MiB, L2) | DRAM (GiB) | DRAM type | CUDA cap. | Allen settings |
|---|---|---|---|---|---|---|---|
| Geforce GTX 670 | 1344 | 1.06 | 0.5 | 1.95 | GDDR5 | 3.0 | Low |
| Geforce GTX 680 | 1536 | 1.14 | 0.5 | 1.95 | GDDR5 | 3.0 | Low |
| Geforce GTX 780 Ti | 2880 | 0.93 | 1.5 | 2.95 | GDDR5 | 3.5 | Low |
| Geforce GTX 980 | 2048 | 1.29 | 2 | 2.01 | GDDR5 | 5.2 | Low |
| Geforce GTX TITAN X | 3072 | 1.08 | 3 | 11.92 | GDDR5 | 5.2 | High |
| Geforce GTX 1060 6G | 1280 | 1.81 | 1.5 | 5.94 | GDDR5 | 6.1 | Low |
| Geforce GTX 1080 Ti | 3584 | 1.67 | 2.75 | 10.92 | GDDR5 | 6.1 | High |
| Geforce RTX 2080 Ti | 4352 | 1.545 | 6 | 10.92 | GDDR5 | 7.5 | High |
| Tesla T4 | 2560 | 1.59 | 4 | 15.72 | GDDR6 | 7.5 | High |
| Tesla V100 32GB | 5120 | 1.37 | 6 | 32 | HBM2 | 7.0 | High |

# Computing costs and prospects

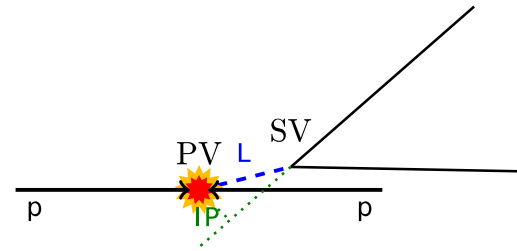| Architecture | GPU | CPU |
|---|---|---|
| Performance | > 60 kHz | 35 kHz on 20 cores |
| Amount | 500 | 1500 (12C / 24T each) |
| Type | RTX 2080 Ti GPUs | Intel Xeon Silver 4116 2.1 G + RAM |
| Price | 0.5M $ | 2M $ |
| Compactness (servers) | 250 (< 100 with PCIe4) | 750 (dual-socket) |

Peak Memory Bandwidth (GB/s)

Peak Double Precision (GFLOPs)

# Beauty and charm decays

Beauty Hadrons

$\mu^-$

SV

$\mu^+$

PV

L

p

p

IP

Charm Hadrons

SV

PV  L

p

p

IP

- B$^{\pm/0}$ mass ~5.3 GeV

  → Daughter p$_T$ $\mathcal{O}$(1 GeV)

- τ ~1.6 ps → flight distance ~1cm

- Detached muons from B→J/ΨX, J/Ψ → μ⁺μ⁻

- Displaced tracks with high p$_T$

- D$^{\pm/0}$ mass ~1.9 GeV

  → Daughter p$_T$ $\mathcal{O}$(700 MeV)

- τ ~0.4 ps → flight distance ~4mm

- Also produced from B decays

PV: Primary vertex
SV: Secondary vertex
IP: Impact parameter: distance between point
of closest approach of a track and a PV
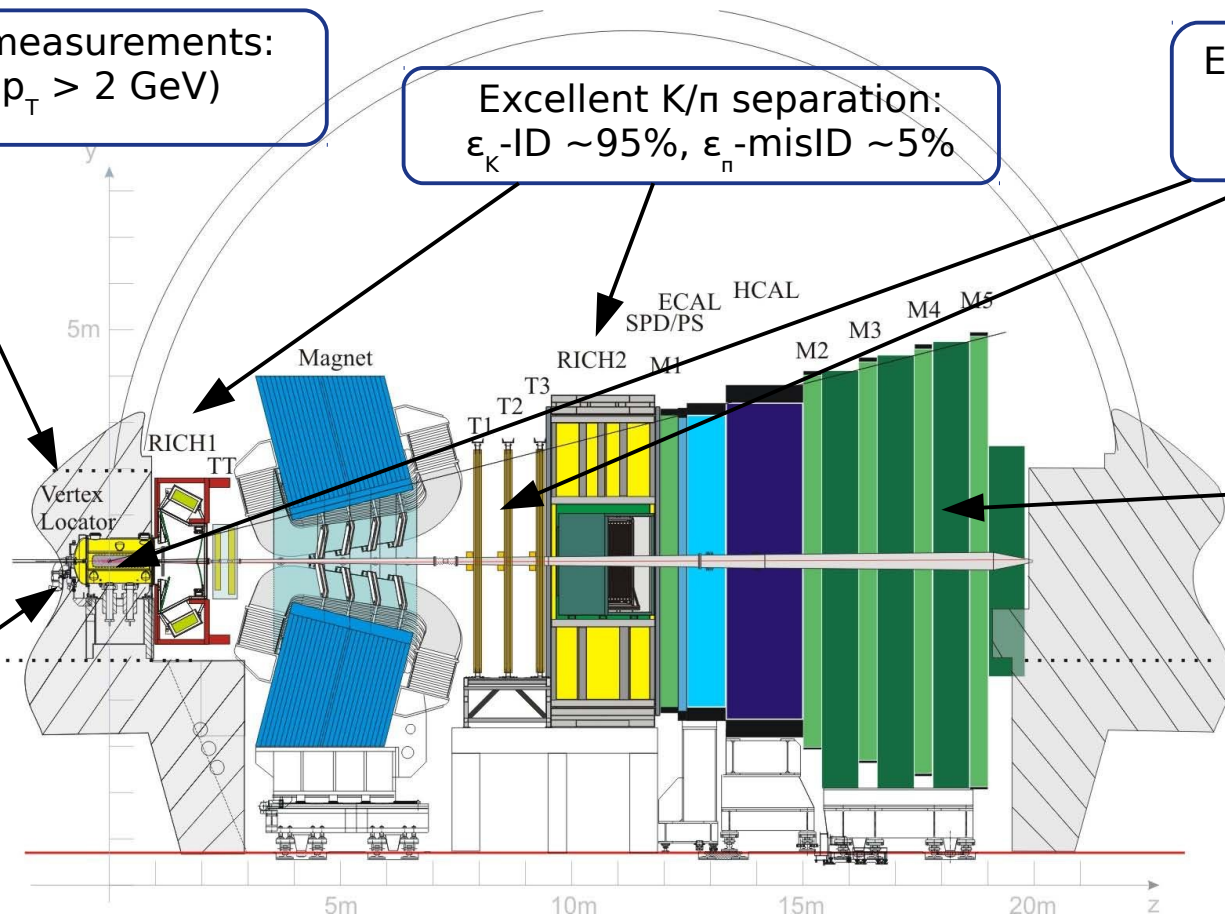
# LHCb detector, 2011 - 2018

Precise vertex measurements:
$\sigma_{IP} = 20$ μm ($p_T > 2$ GeV)

Excellent K/π separation:
$\varepsilon_K$-ID ~95%, $\varepsilon_\pi$-misID ~5%

Excellent momentum resolution:
$\Delta p/p$ ~0.5-1%

Excellent muon Identification:
$\varepsilon_\mu$-ID ~97%

Excellent decay time resolution:
$\sigma_\tau \sim 45$ fs
for b hadrons

ECAL HCAL
SPD/PS
RICH2 M1 M2 M3 M4 M5
T3
T2
T1
Magnet
RICH1
TT
Vertex
Locator

5m

5m 10m 15m 20m z

# GPU performance over time



Theoretical Peak Performance, Single Precision

# Throughput of x86 HLT1