# Allen: A High Level Trigger on GPUs for LHCb

## Physics and throughput performance

**Dorothea vom Bruch**

on behalf of the LHCb collaboration

LPNHE, CNRS

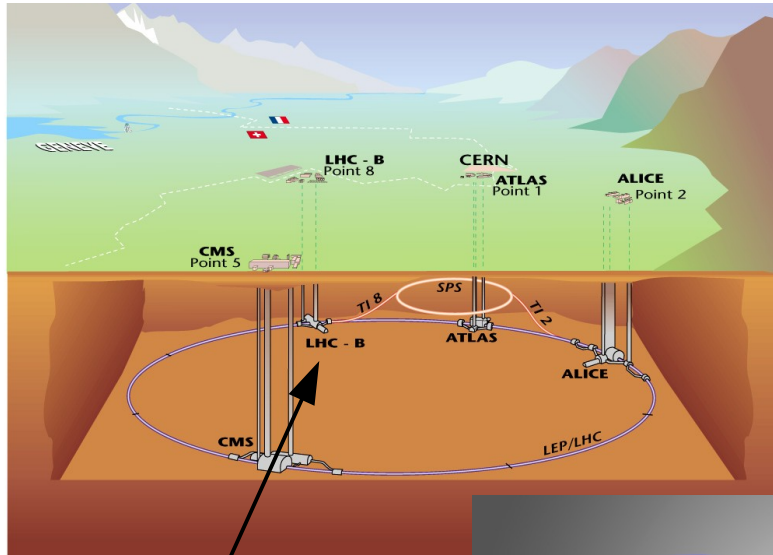Sorbonne University, Paris Diderot University

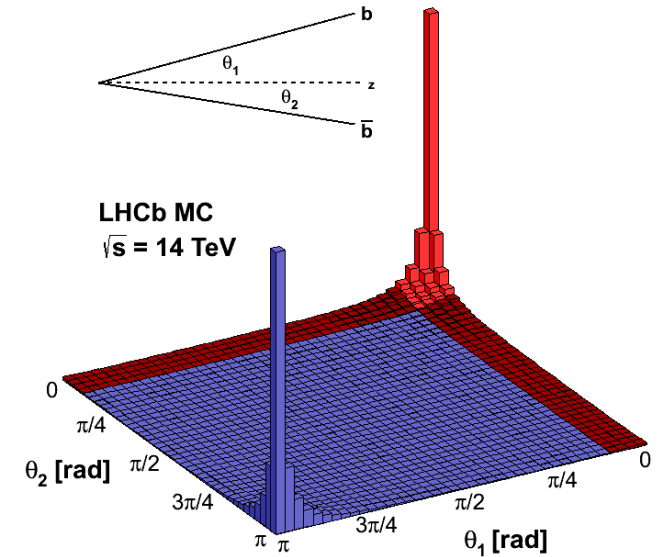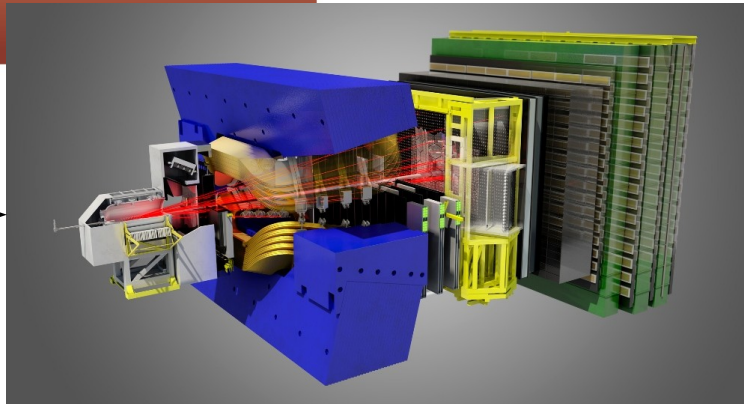November 6[th] 2019
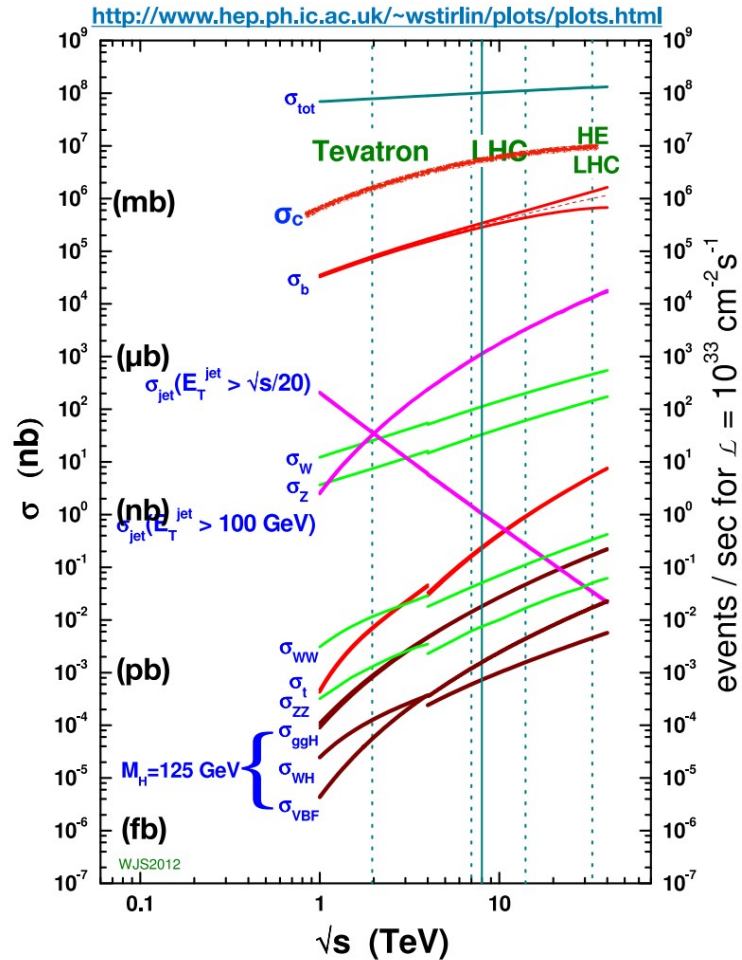
CHEP 2019, Adelaide

# LHCb

LHC @ CERN

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

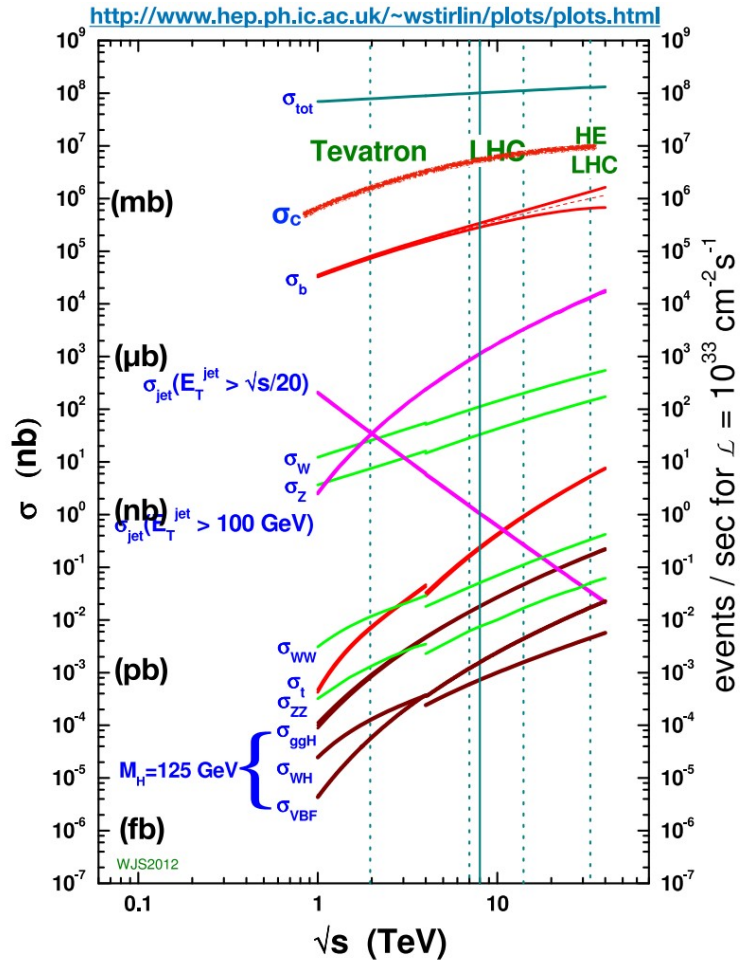# Reaching the MHz signal era



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

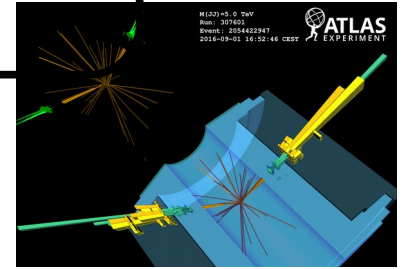# Reaching the MHz signal era

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



- General purpose LHC experiments
- Local characteristic signatures
- 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

- Too many interesting events
- No "simple" local criteria for selection
  → hardware-level trigger not an option
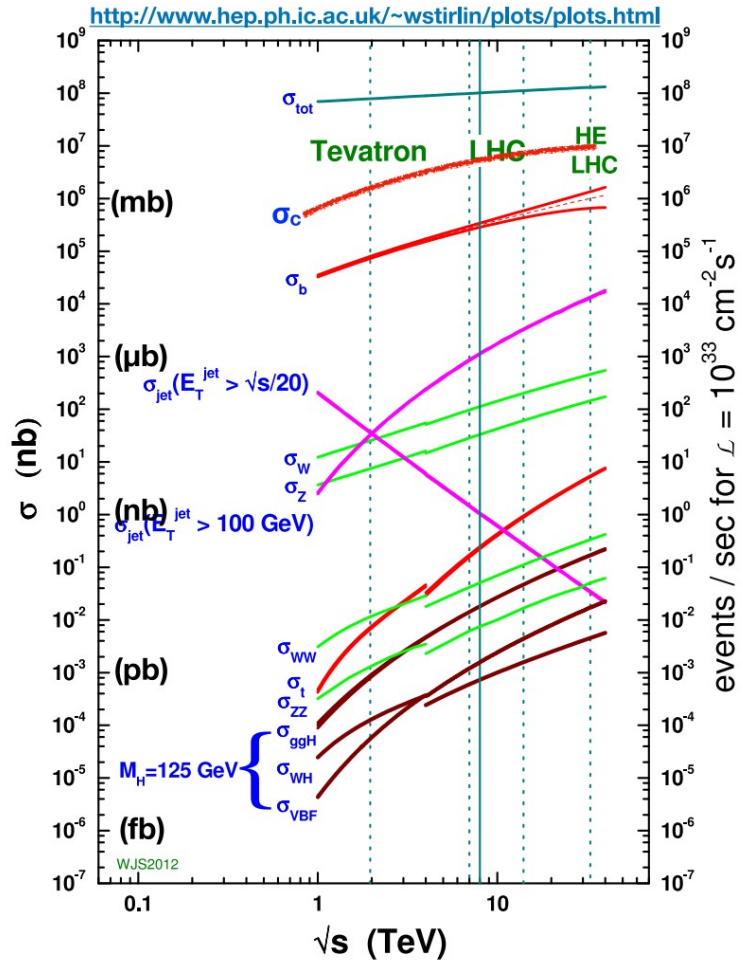
- General purpose LHC experiments
- Local characteristic signatures
- Can trigger efficiently at ~100 kHz
- Hardware-level trigger possible

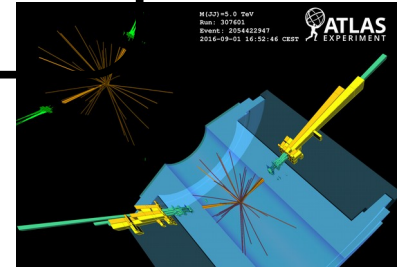# Change in trigger paradigm



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

# Tracks in the LHCb detector



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

# Trigger upgrade for Run 3 (2021)

**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**

# Trigger upgrade for Run 3 (2021)



**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

9

# 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!**

# Architecture for high level trigger?

40 Years of Microprocessor Trend Data



**Graphics Processing Units (GPUs) have thousands of cores**

# Amdahl's law

Speed-up factor vs N processors

Parallel fraction
- 90%
- 75%
- 50%

Speedup in latency = 1 / (S + P/N)

S: sequential part of program

P: parallel part of program

N: number of processors

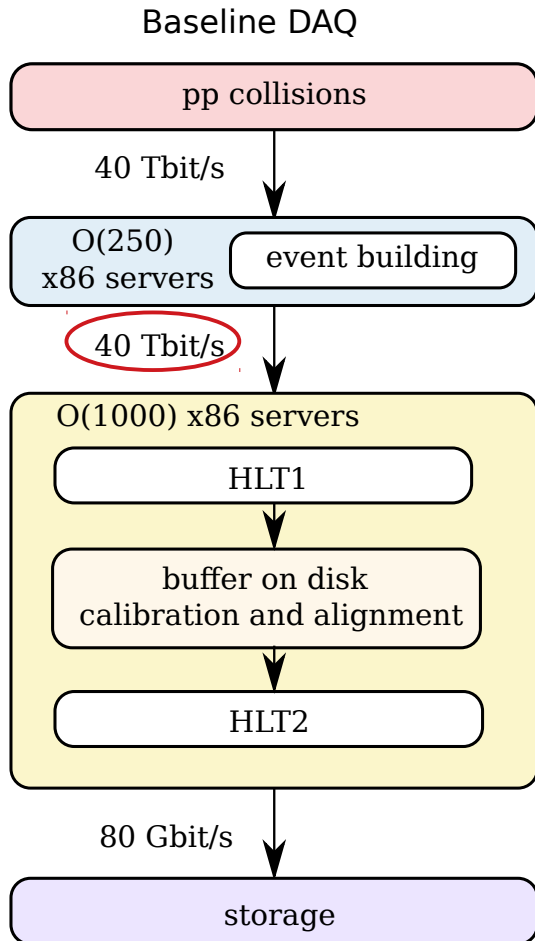**Can we use the FLOPS available on a GPU to run HLT1 @ 30 MHz?**

# Where to place the GPUs?

Baseline DAQ

# Where to place the GPUs?

# Where to place the GPUs?

# LHCb HLT1 elements



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

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

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

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

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

17

# 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<br>→ no stringent latency requirements | GPUs have higher latency than CPUs,<br>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 | GPUs have 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 | GPUs have higher latency than CPUs,<br>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 R&D project

- Fully standalone software project: https://gitlab.cern.ch/lhcb/Allen
- Only requirements:

  C++17 compliant compiler, CUDA v10, boost, ZeroMQ
- Built-in physics validation
- Configurable sequence, custom memory manager
- Cross-architecture compatibility

- Project started in February 2018
- After 15 months of development time:

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

- Talk on software challenges by D. Cámpora: Monday, Track 5

- 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)

- Process thousands of events in parallel
- Single precision only

Within one block:
intra-event parallelization

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

# Velo detector

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

# Velo detector: clustering

26 planes of silicon pixel detectors

Clustering with bit masks



1 m

cross section at y=0

390 mrad

70 mrad

15 mrad

66 mm

interaction region showing
$2 \times \sigma_{beam} = \sim 12.6$ cm
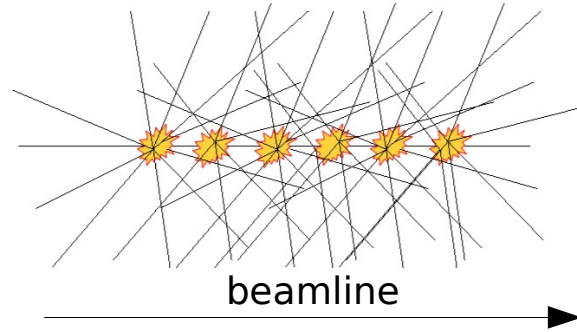
# Velo detector: track reconstruction
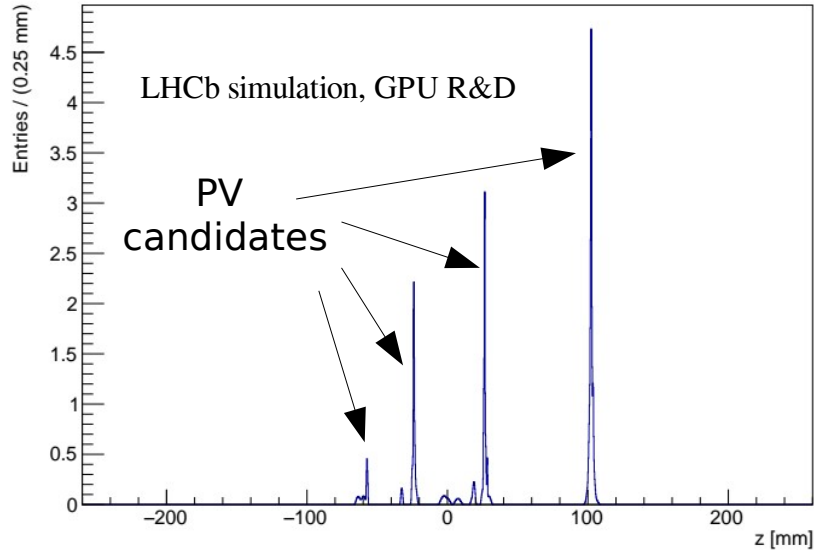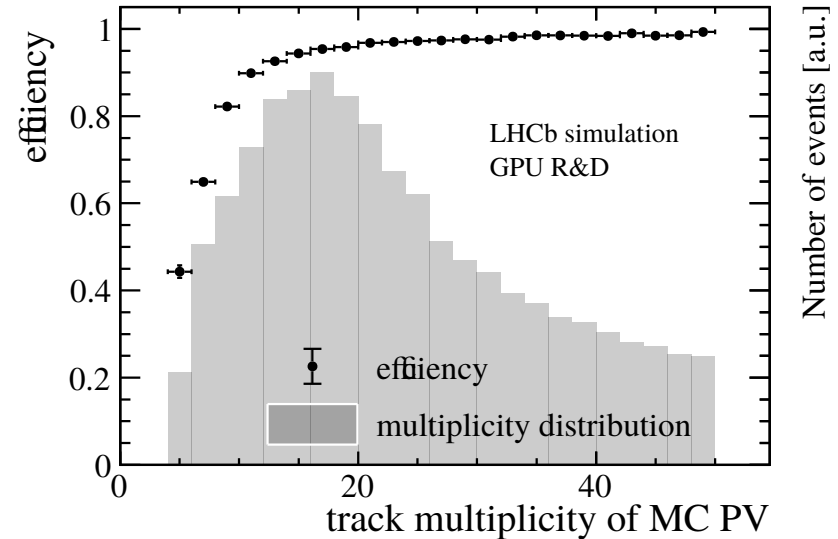
1) Sort hits by φ

2) Triplet seeding

3) Triplet forwarding

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

Point of closest approach of tracks to beamline

LHCb simulation, GPU R&D

PV candidates

PV reconstruction efficiency

LHCb simulation
GPU R&D

efftiency

multiplicity distribution

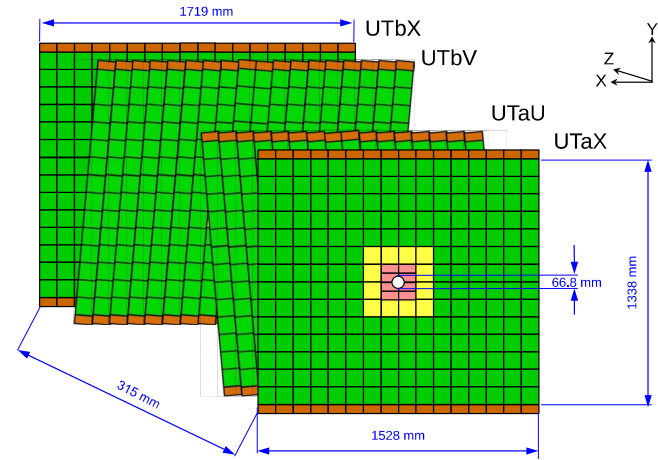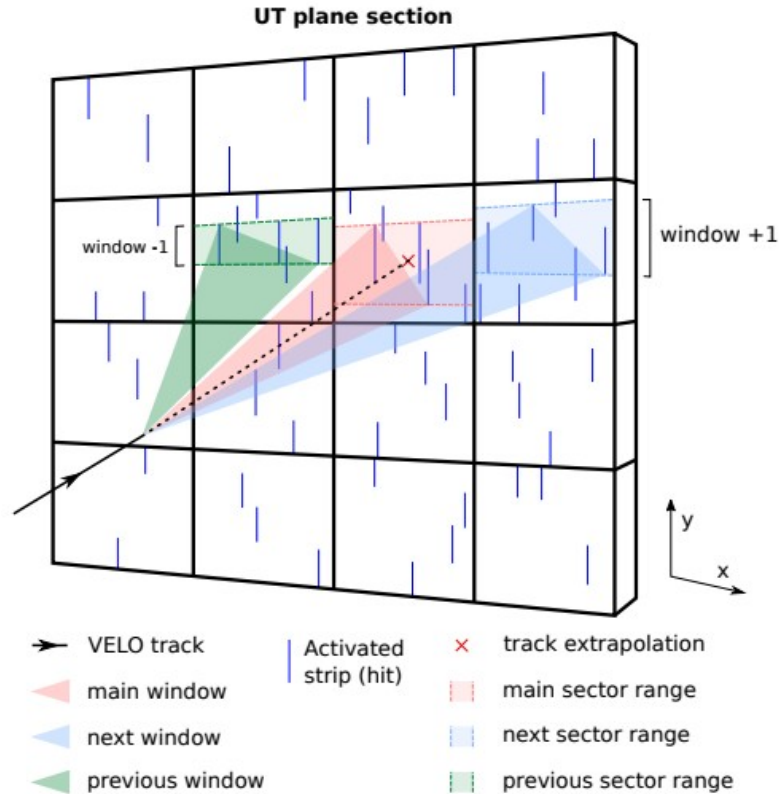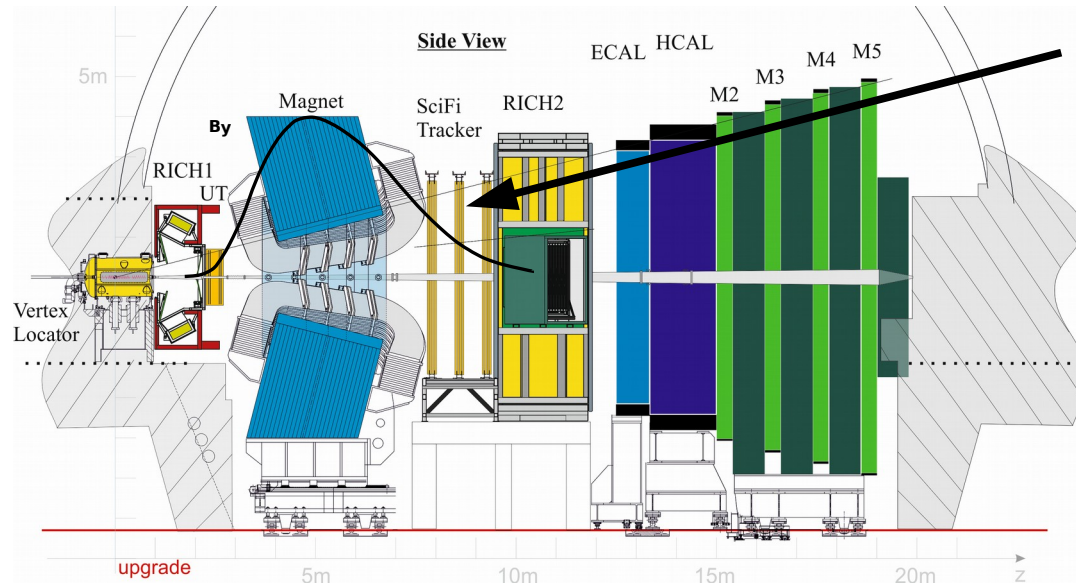effiency

track multiplicity of MC PV

# UT detector



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

# UT detector: track reconstruction

## 4 planes of silicon strip detectors

P. Fernandez Declara, D. Campora Perez, J. Garcia-Blas, D. vom Bruch, J. Daniel Garca, N. Neufeld , IEEE Access 7 (2019)
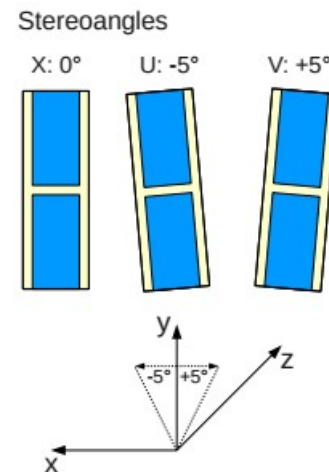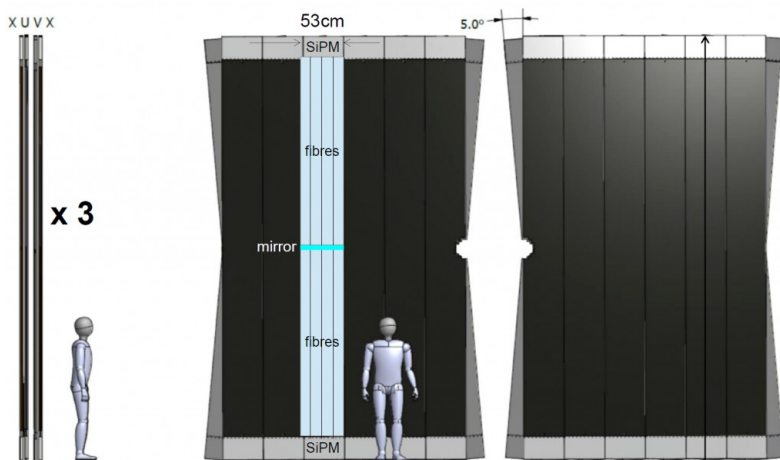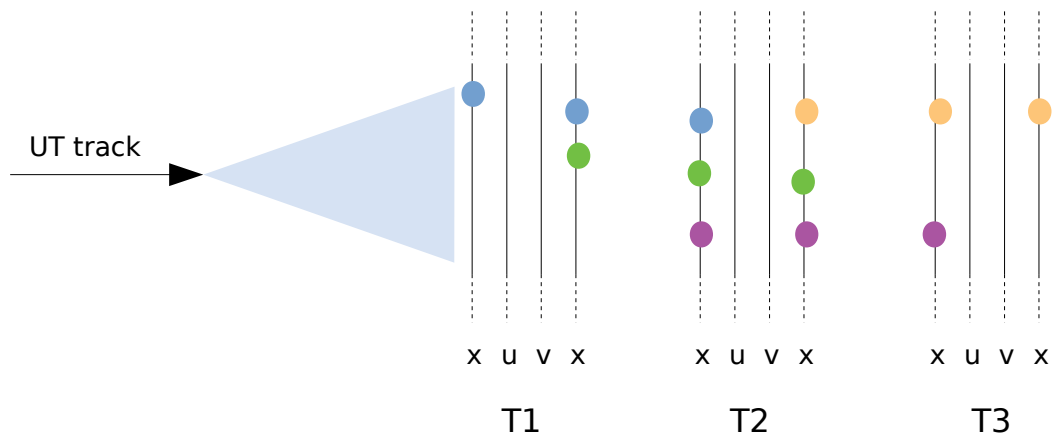
# SciFi detector
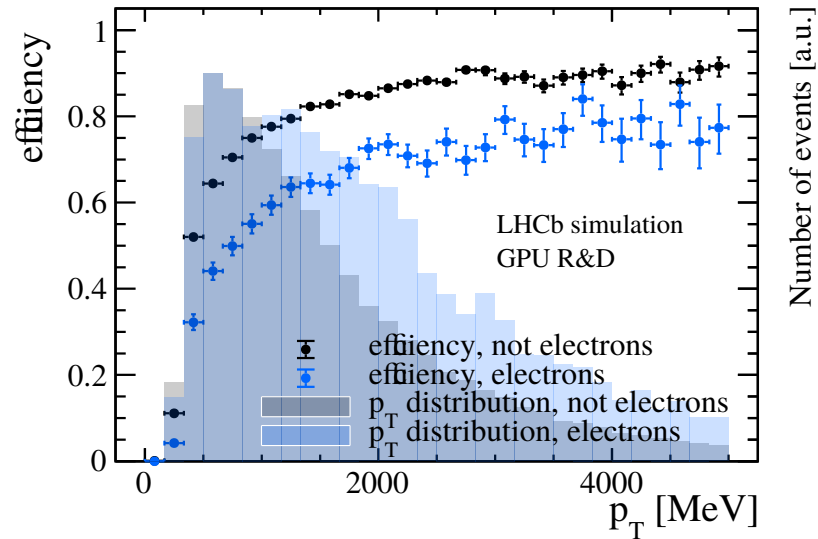


**SciFi**
- Decode raw data
- Track reconstruction

# SciFi detector

- 12 layers of scintillating fibres
- Efficiency of fibres ~98-99%
- Describe trajectories in magnetic field with parameterizations
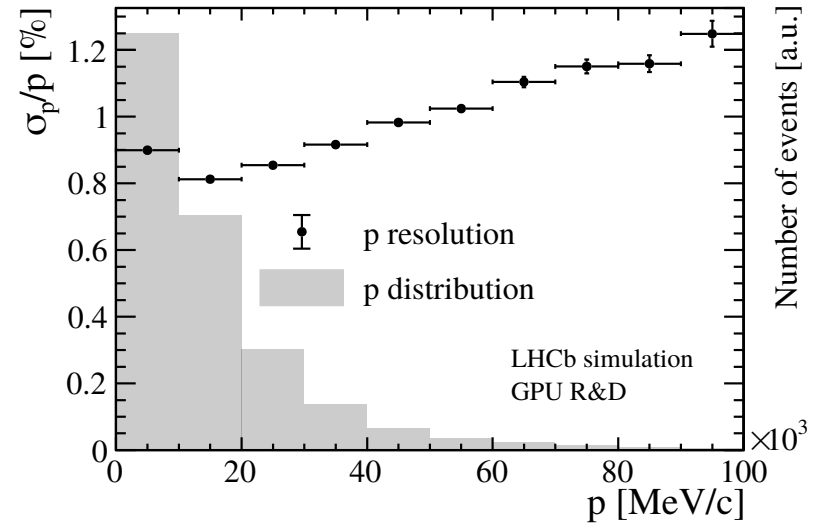  → no need to load large field map into GPU memory

UT track

x u v x
x u v x
x u v x
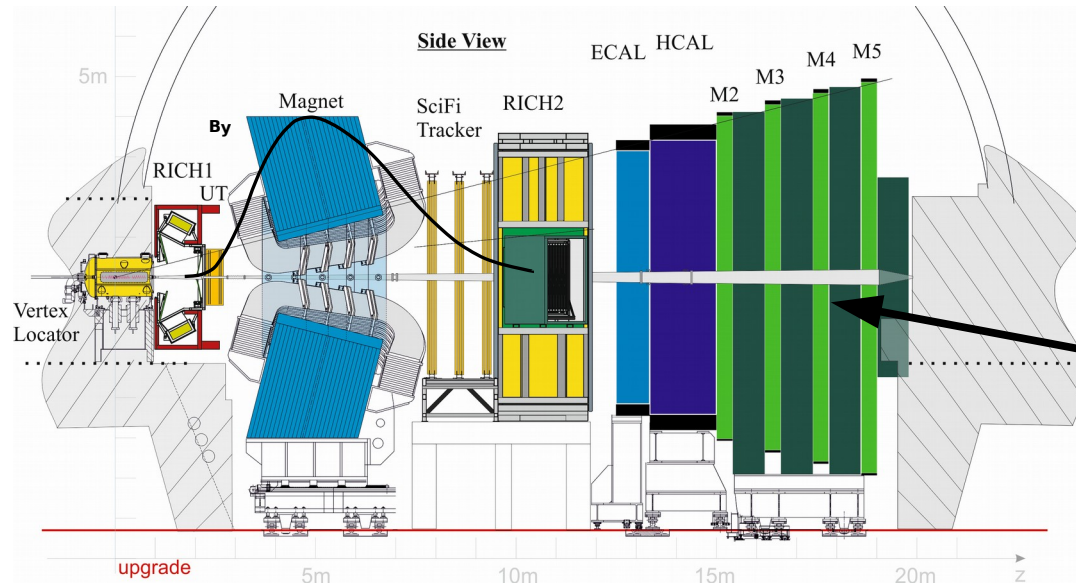
T1
T2
T3

XUVX

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



LHCb simulation
GPU R&D

efficiency, not electrons
efficiency, electrons
$p_T$ distribution, not electrons
$p_T$ distribution, electrons

Momentum resolution



p resolution

p distribution
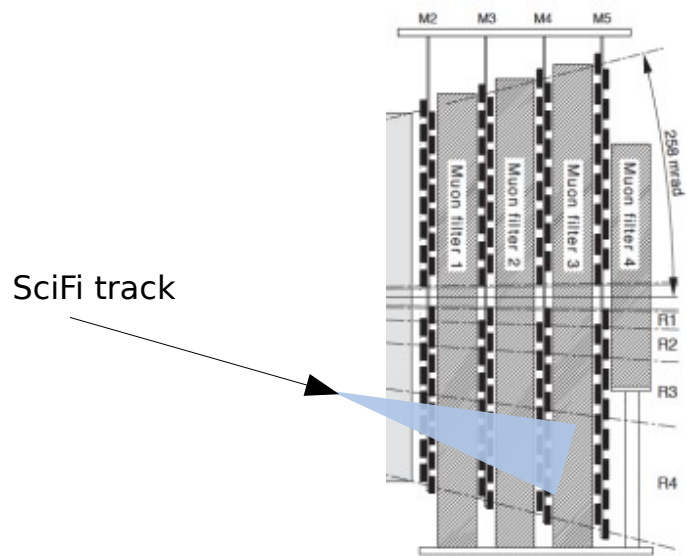
LHCb simulation
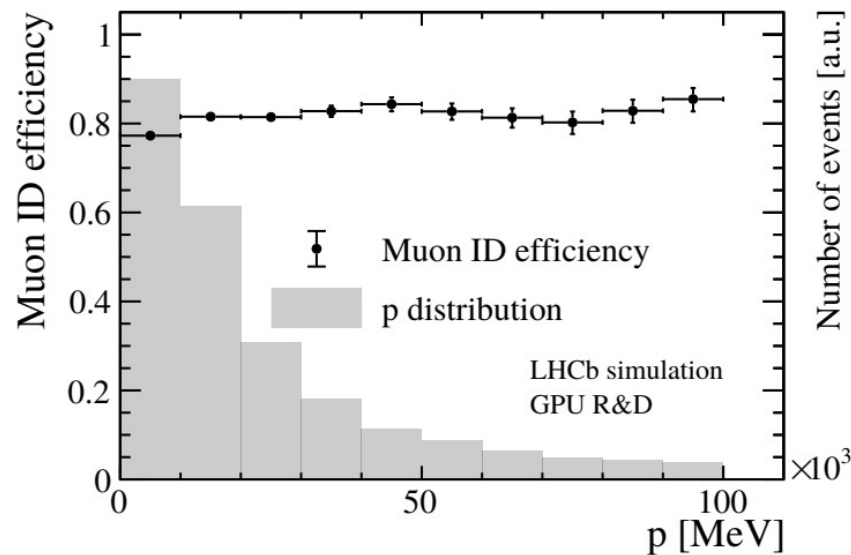GPU R&D

# Muon chambers



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

# Muon identification

Four multi-wire proportional chambers
Interleaved with iron walls

Muon identification efficiency

SciFi track

# Ingredients for selections
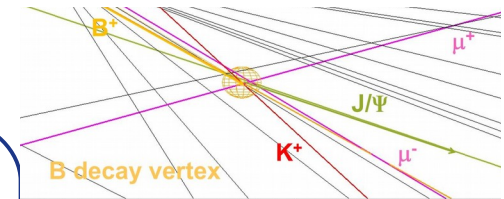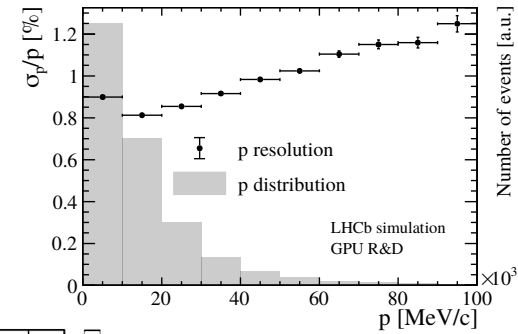
Primary vertices



Secondary vertices



Momentum



Impact parameter



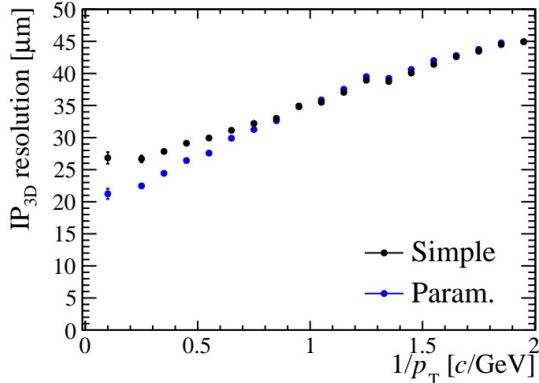**Selections**

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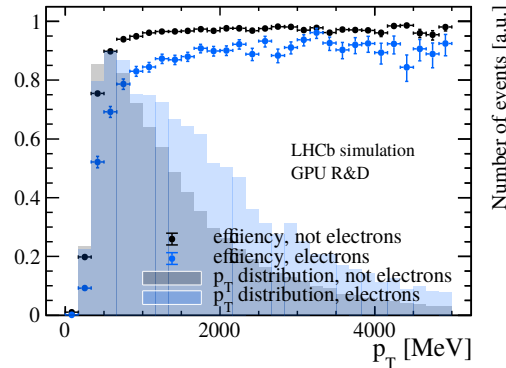- Based on p, $p_t$, displacement, vertex criteria and muon identification

Tracks
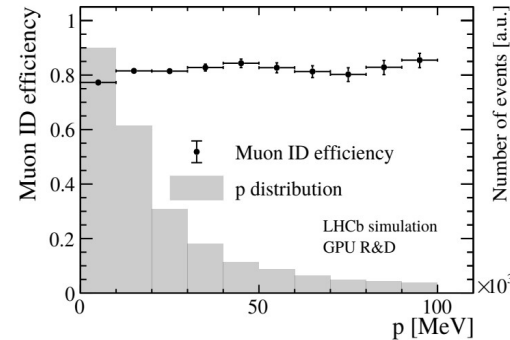


Muon identification

# Event selection

| Trigger | Rate [kHz] |
|---|---|
| 1-Track | 249 $\pm$18 |
| 2-Track | 663 $\pm$30 |
| High-$p_T$ muon | 1 $\pm$1 |
| Displaced dimuon | 50 $\pm$8 |
| High-mass dimuon | 101 $\pm$12 |
| Total | 971 $\pm$36 |

Selection efficiencies, values given in %

| Signal | GEC | TIS -OR- TOS | TOS | GEC $\times$ TOS |
|---|---|---|---|---|
| $B^0 \to K^{*0}\mu^+\mu^-$ | 89 $\pm$2 | 85 $\pm$2 | 78 $\pm$3 | 69 $\pm$3 |
| $B^0 \to K^{*0}e^+e^-$ | 84 $\pm$3 | 69 $\pm$4 | 62 $\pm$4 | 53 $\pm$3 |
| $B_s^0 \to \phi\phi$ | 83 $\pm$3 | 70 $\pm$3 | 65 $\pm$4 | 54 $\pm$3 |
| $D_s^+ \to K^+K^-\pi^+$ | 82 $\pm$4 | 62 $\pm$5 | 38 $\pm$5 | 32 $\pm$4 |
| $Z \to \mu^+\mu^-$ | 78 $\pm$1 | 97 $\pm$1 | 97 $\pm$1 | 75 $\pm$1 |

GEC: Global event cut
TIS: Trigger independent from signal
TOS: Trigger on signal

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

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

# Full HLT1 running on GPUs

Physics performance matches HLT1 requirements

What about the throughput performance?

# Throughput on various GPUs

Throughput of the full HLT1 sequence



| GPU | Throughput |
|---|---|
| Tesla V100-PCIE-32GB | 78.62 kHz |
| Quadro RTX 6000 | 72.83 kHz |
| GeForce RTX 2080 Ti | 68.36 kHz |
| Tesla T4 | 34.92 kHz |
| GeForce GTX 1080 Ti | 30.23 kHz |
| GeForce GTX TITAN X | 19.20 kHz |
| GeForce GTX 1060 6GB | 12.74 kHz |
| GeForce GTX 680 | 5.50 kHz |
| GeForce GTX 670 | 5.22 kHz |

LHCb simulation

GPU R&D

**HLT1 can run on 500 GPUs
→ Buy GPUs instead of expensive network**

# Allen scalability with GPU model

# The Allen team

# Summary

- Allen is the first complete high throughput trigger implementation on GPUs
- Developed a compact, modular and scalable framework
- Baseline HLT1 can run on GPUs
- Scaling of GPU performance → maximize physics discovery potential of LHCb
- Integration tests ongoing (see talk by D. Cámpora, Monday Track 5)
- HLT1 on GPUs is being considered as alternative to the baseline x86 architecture

**Software High Level Trigger**

**Full event reconstruction, inclusive and exclusive kinematic/geometric selections**

**500 GPUs**

# Backup

# 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?

# Graphics requirements

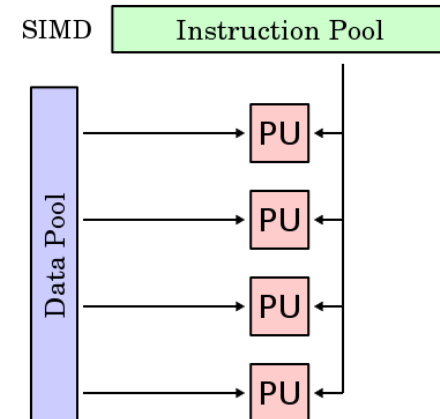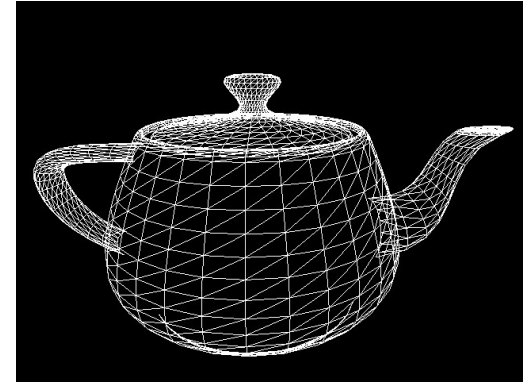**Graphics pipeline**

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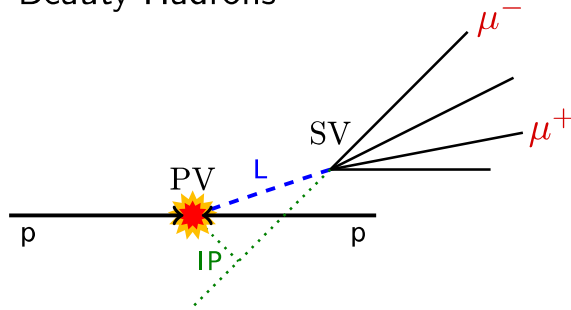  - Applying material properties and light situation to every pixel

**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

# Beauty and charm decays

Beauty Hadrons

Charm Hadrons

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

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

- $\tau$ ~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)

- $\tau$ ~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

46

# Why no low level trigger?

Low level trigger on $E_T$ from the calorimeter

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



**Need track reconstruction at first trigger stage**

# 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

# GPU in a nutshell

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

GPU



Data transfer to a GPU



PCIe connection

| 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)

# Algorithm breakdown



| Algorithm | % |
|---|---|
| search_by_triplet | 11.97 % |
| lf_triplet_seeding | 9.50 % |
| pv_beamline_multi_fitter | 5.35 % |
| muon_add_coords_crossing_maps | 4.67 % |
| lf_collect_candidates | 4.64 % |
| pv_beamline_peak | 4.50 % |
| scifi_direct_decoder_v4 | 4.33 % |
| lf_quality_filter_x | 3.88 % |
| lf_triplet_keep_best | 3.79 % |
| estimate_input_size | 3.63 % |
| compass_ut | 3.36 % |
| masked_velo_clustering | 3.08 % |
| lf_extend_tracks_x | 2.81 % |
| ut_search_windows | 2.69 % |
| calculate_phi_and_sort | 2.63 % |
| lf_fit | 2.05 % |
| lf_search_initial_windows | 2.00 % |

Showing only algorithms contributing ≥ 2%

# 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 |

# Throughput of x86 HLT1



LHCb Upgrade simulation

Scalar event model, maximal SciFi reconstruction

Scalar event model, fast SciFi reconstruction with tighter track tolerance criteria

Scalar event model, vectorizable SciFi reconstruction with entirely reworked algorithm logic

Fully SIMD-POD friendly event model, vectorizable SciFi and vectorized vertex detector and PV reconstruction, I/O improvements

20 physical cores per node