# The LHCb trigger and data acquisition system in Run 3
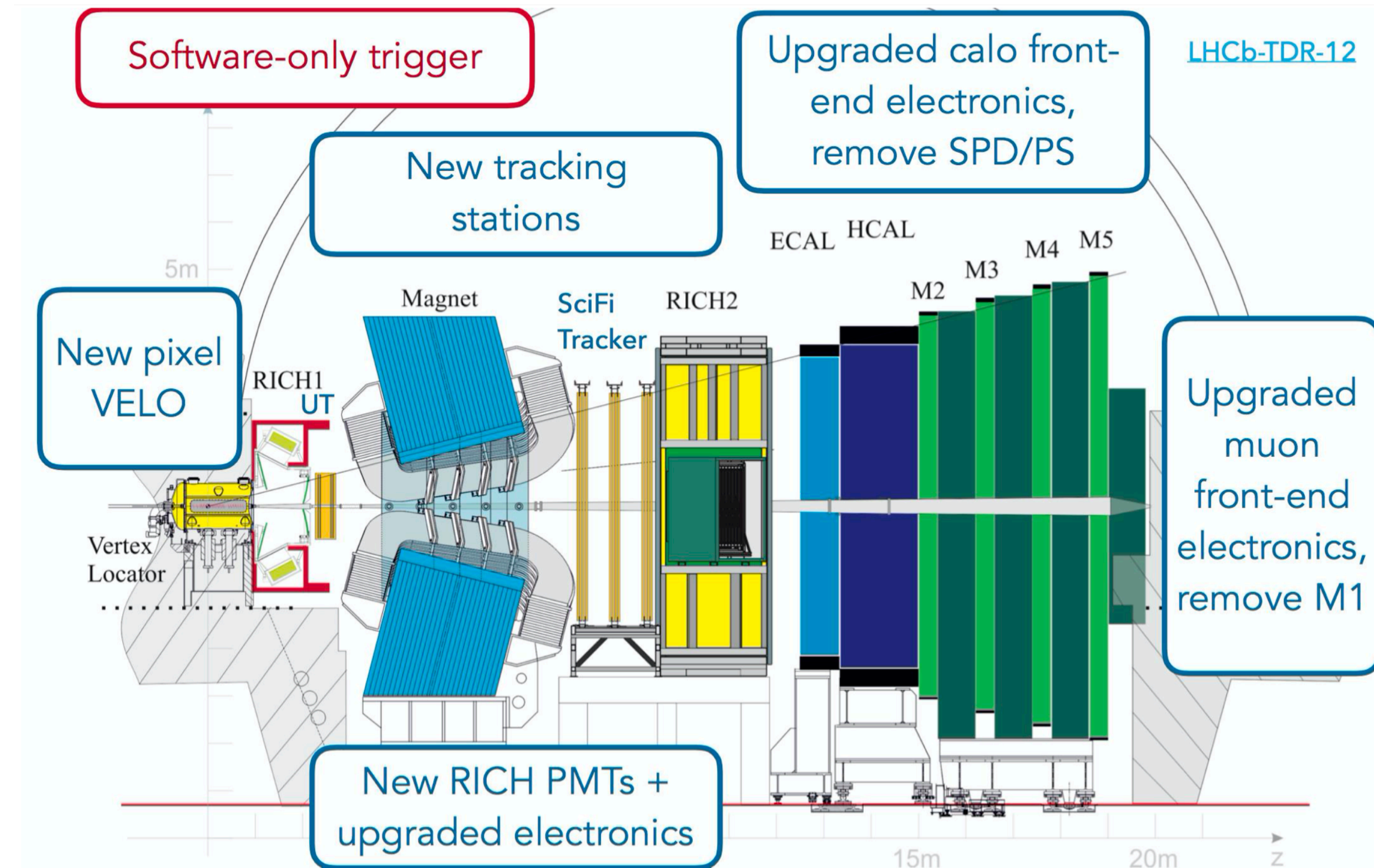
Daniel Hugo Cámpora Pérez

**Maastricht University**

**Department of Data Science and Knowledge Engineering**

**Daniel Cámpora**

# About the hardware

About parallel algorithms

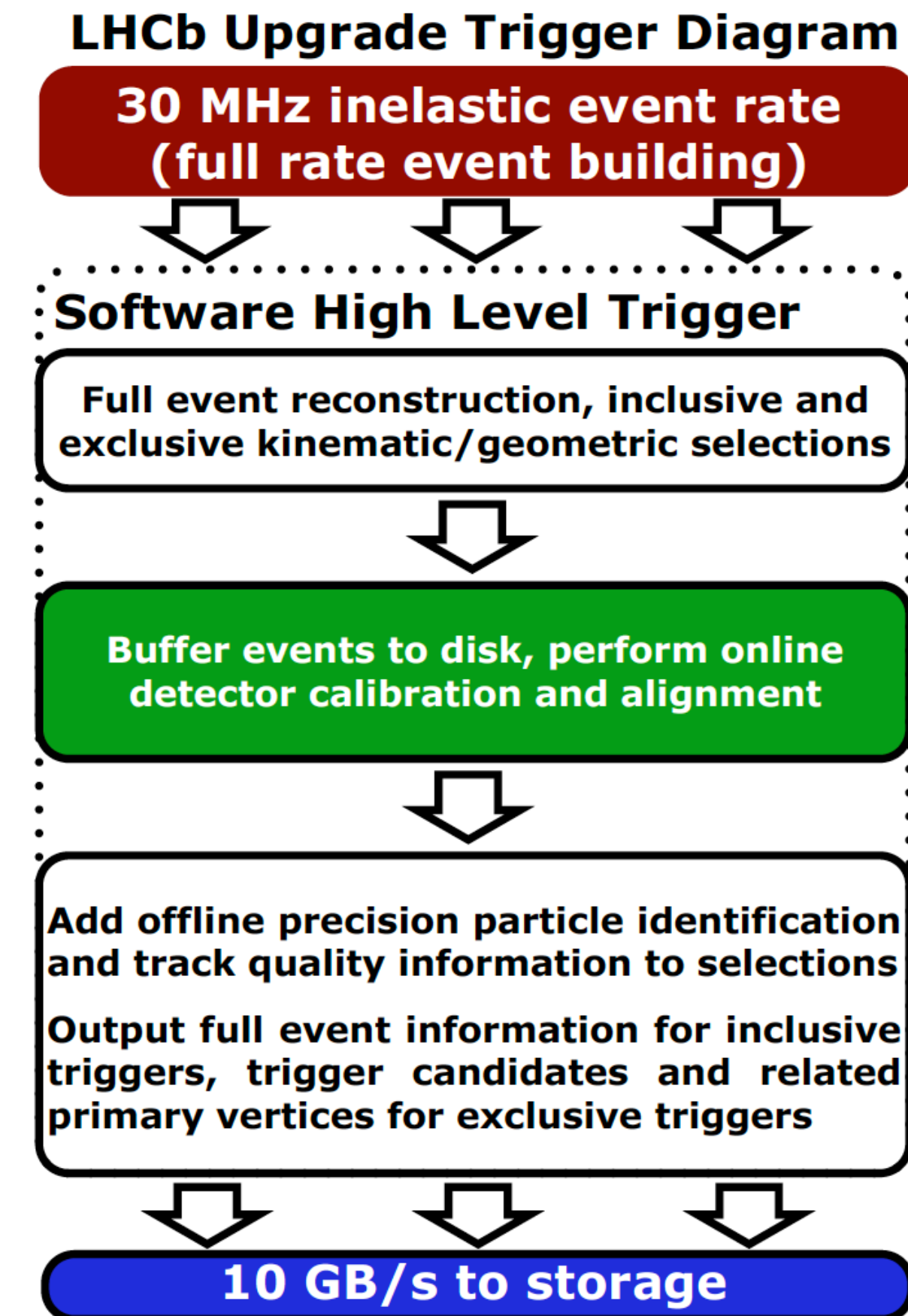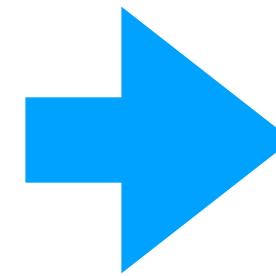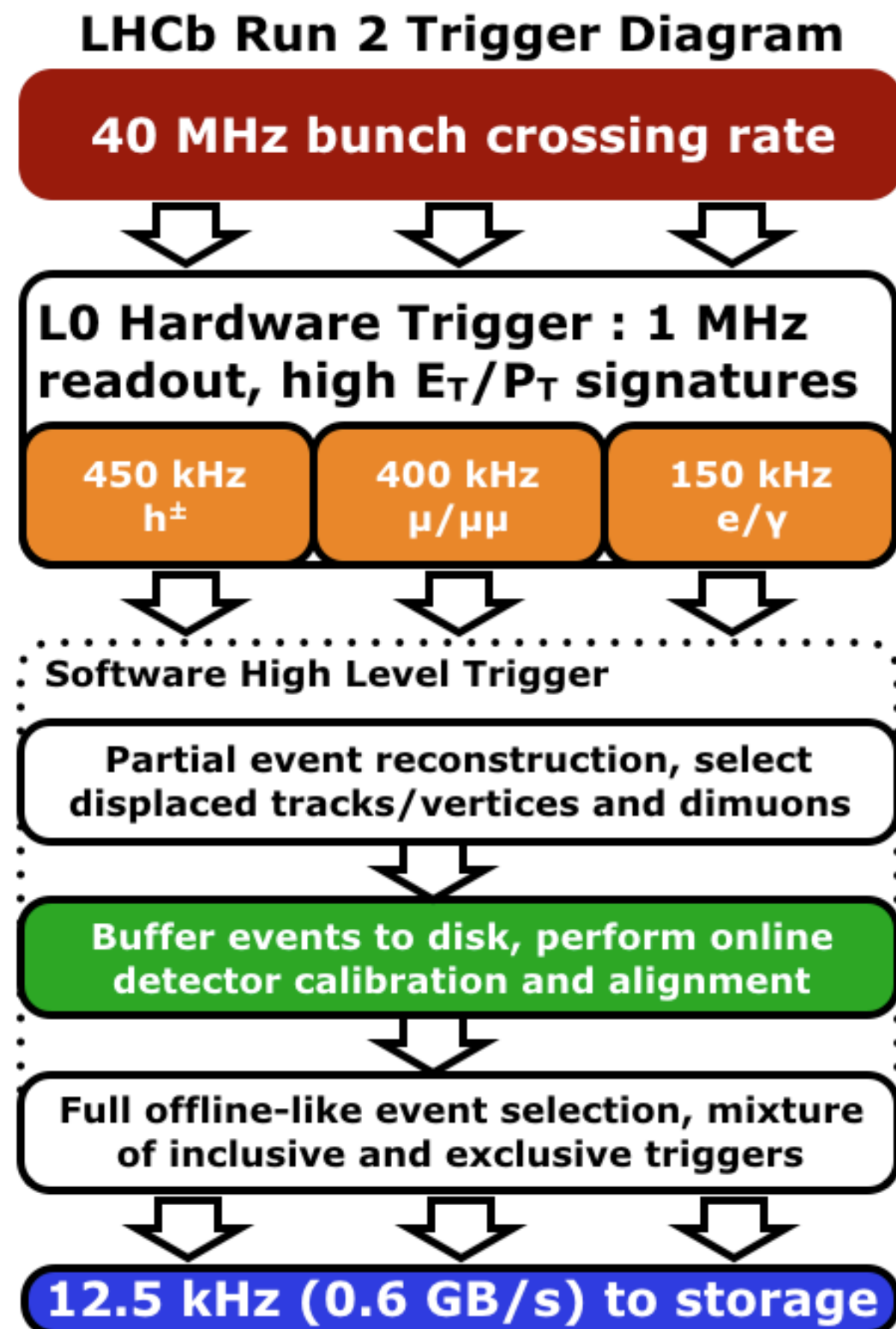About commissioning

# The LHCb U1 upgrade



The LHCb detector at CERN:

- Single-arm forward spectrometer for high-precision flavour physics
- High precision tracking and vertexing
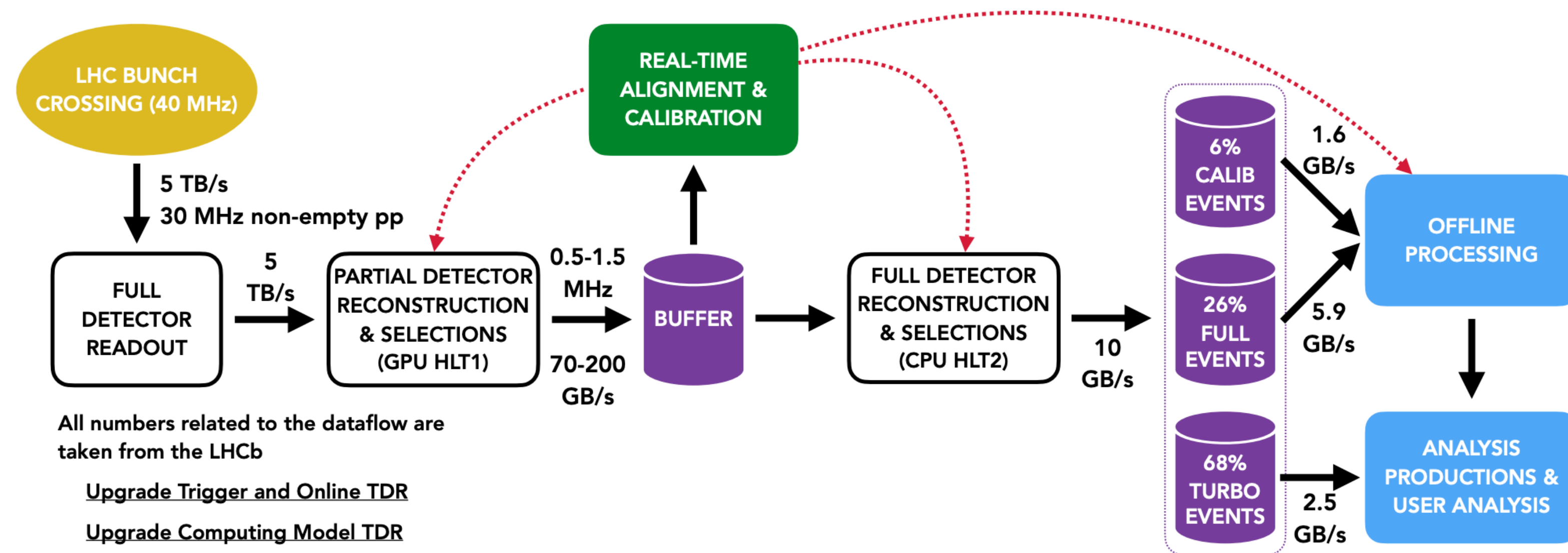- Complemented with excellent PID

The U1 upgrade

- Instantaneous luminosity will increase by x5
- Major upgrade in all sub-detectors to handle increased rates
- Software-only trigger!

**Daniel Cámpora**

# Someone had to pull the 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

**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

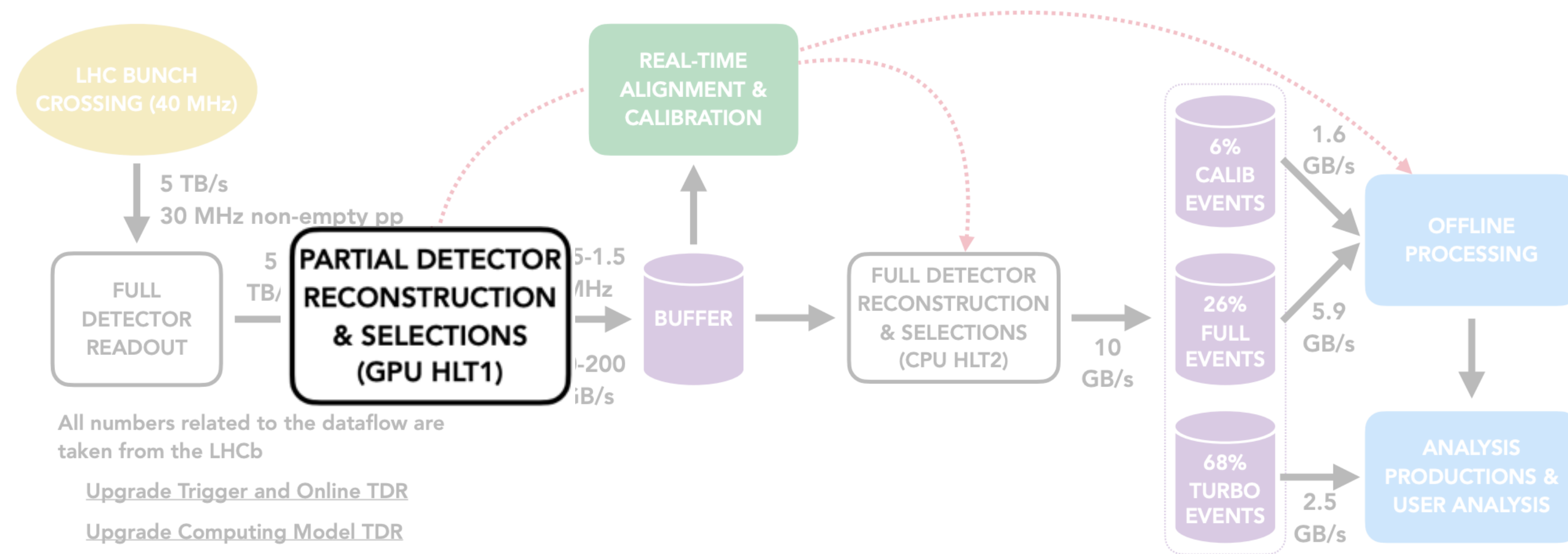10 GB/s to storage

4

**Daniel Cámpora**

# The LHCb data-flow



- Detector data received by O(500) FPGAs and built into events in the event building (EB) farm servers

- 2-stage software trigger, HLT1 & HLT2

- Real-time alignment & calibration

- After HLT2, 10 GB/s of data for offline processing
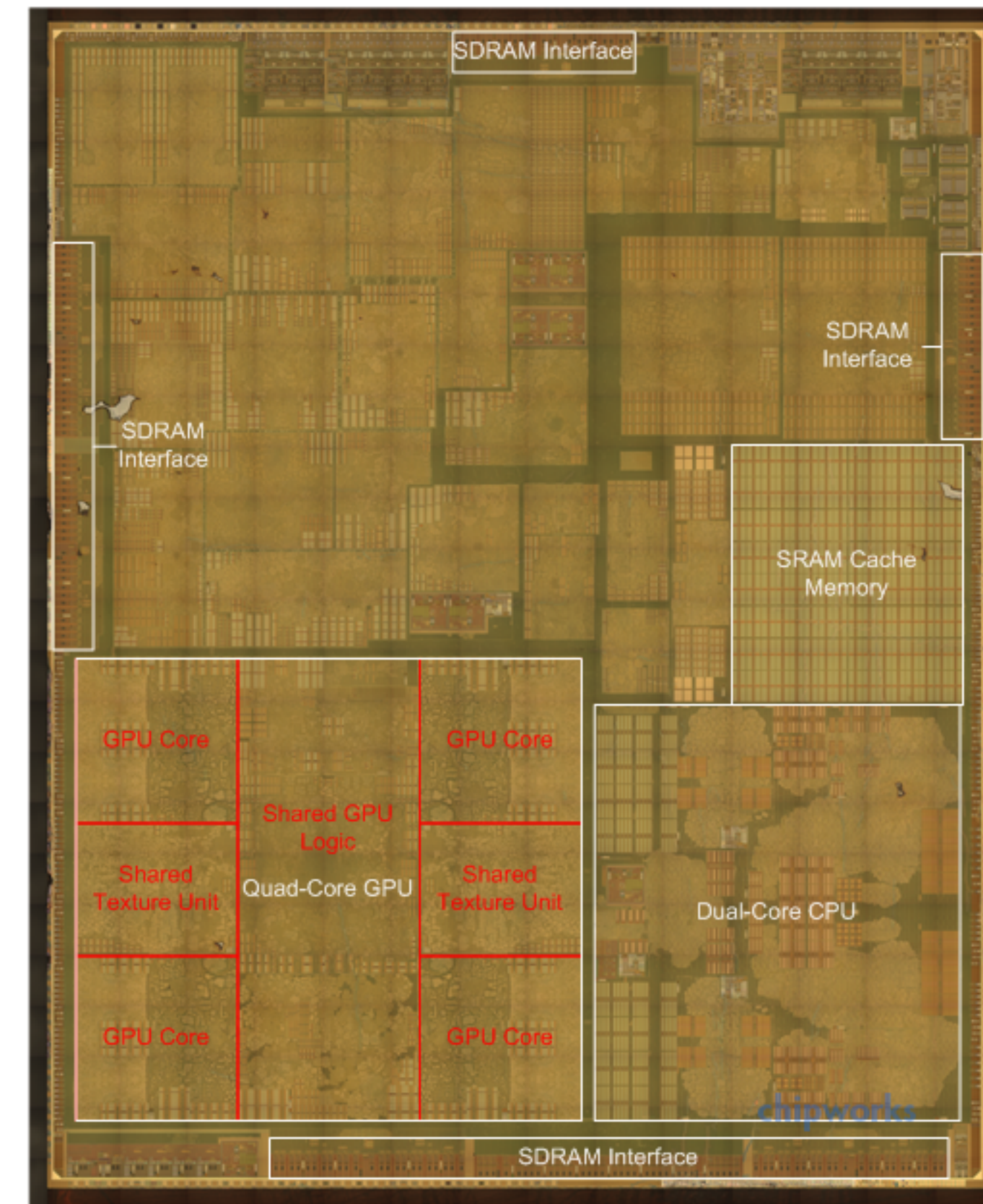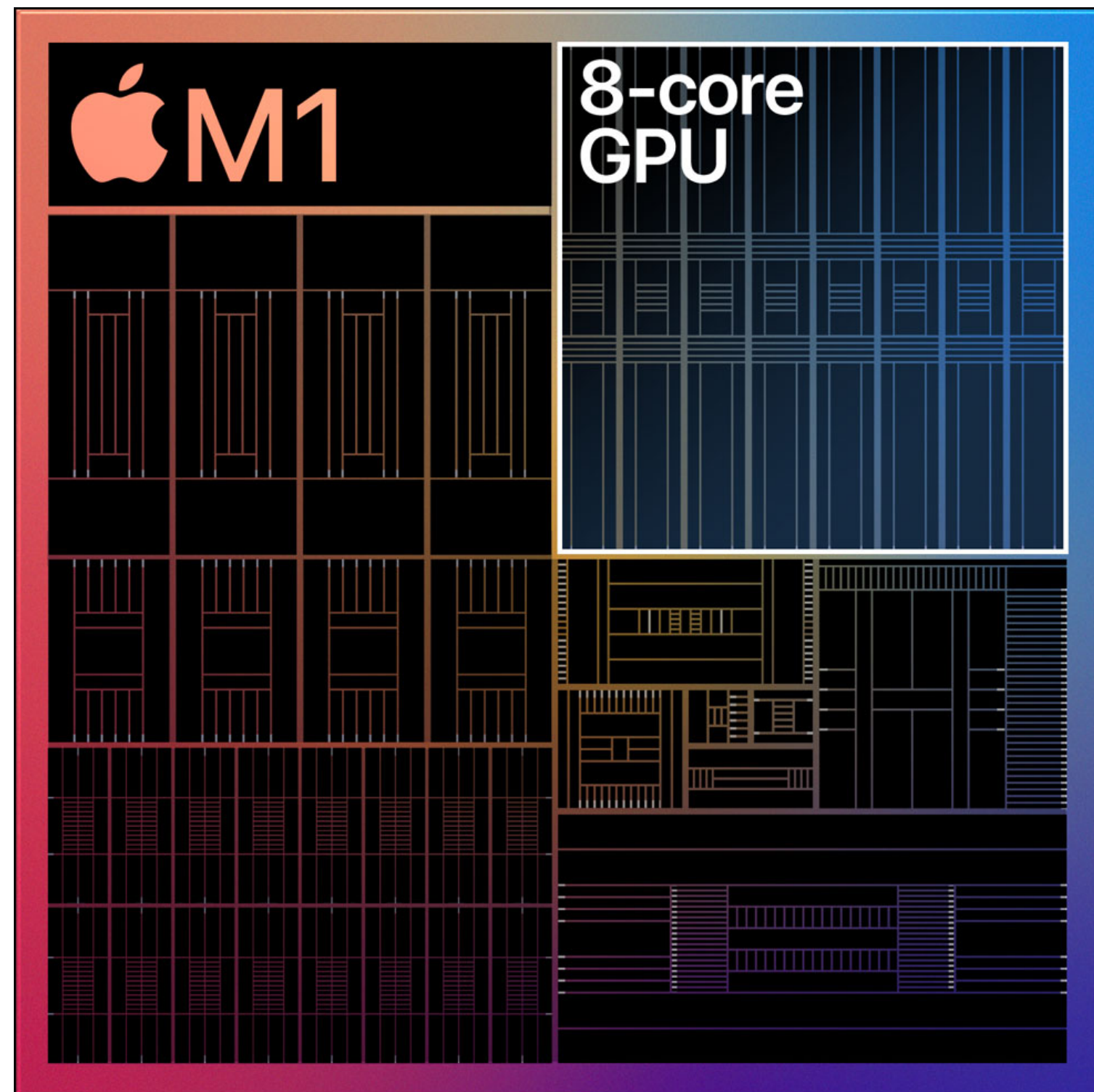
**Daniel Cámpora**

# The LHCb first level trigger



- **The goal of HLT1:**

  - Be able to intake the entirety of the LHCb raw data (5 TB/s) at 30 MHz

  - Perform partial event reconstruction & coarse selection of broad LHCb physics cases

  - Reduce the input rate by a factor of 30 (~ 1 MHz)

  - Store selected events in intermediate buffer for real-time alignment and calibration

  **But why GPUs?**

**Daniel Cámpora**

# CPU radiography

**Daniel Cámpora**

# GPUs: Parallel processors

- GPUs are processors specialized to perform graphic-oriented workloads

**Daniel Cámpora**

# How does a dedicated GPU card look like?

**Daniel Cámpora**

# How does a dedicated GPU card look like?



High-speed network to other GPUs

Stereo

Sync

Power connector

Memory

Processor

Memory

Display video / audio output

PCI-express connector to mainboard

**Daniel Cámpora**

# Strengths and limitations

| Memory | |
|---|---|
| Memory Size: | 24 GB |
| Memory Type: | GDDR6 |
| Memory Bus: | 384 bit |
| Bandwidth: | 768.0 GB/s |

| Render Config | |
|---|---|
| Shading Units: | 8192 |
| TMUs: | 256 |
| ROPs: | 96 |
| SM Count: | 64 |
| Tensor Cores: | 256 |
| RT Cores: | 64 |
| L1 Cache: | 128 KB (per SM) |
| L2 Cache: | 6 MB |

| Theoretical Performance | |
|---|---|
| Pixel Rate: | 162.7 GPixel/s |
| Texture Rate: | 433.9 GTexel/s |
| FP16 (half) performance: | 27.77 TFLOPS (1:1) |
| FP32 (float) performance: | 27.77 TFLOPS |
| FP64 (double) performance: | 867.8 GFLOPS (1:32) |

**https://www.techpowerup.com/gpu-specs/rtx-a5000.c3748**

- Excellent FP16 or FP32 performance. Avoid FP64.

- Room for growth: Tensor cores, RT cores.

- Limited memory: 24000 / 8192 = 2.9 MB per core.

**Daniel Cámpora**

# The main concepts behind the GPU HLT1

- Pipeline: raw-data in, selections out

- Scheduler: many events run concurrently

- Redesign algorithms:

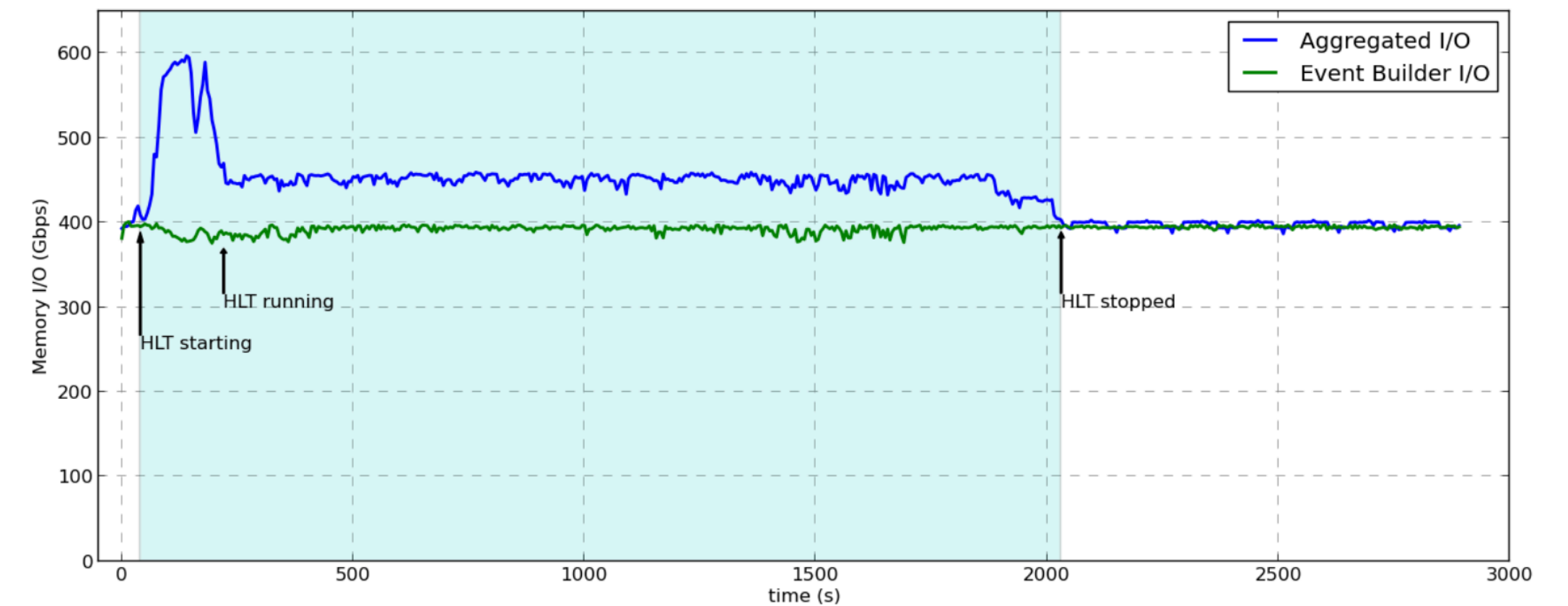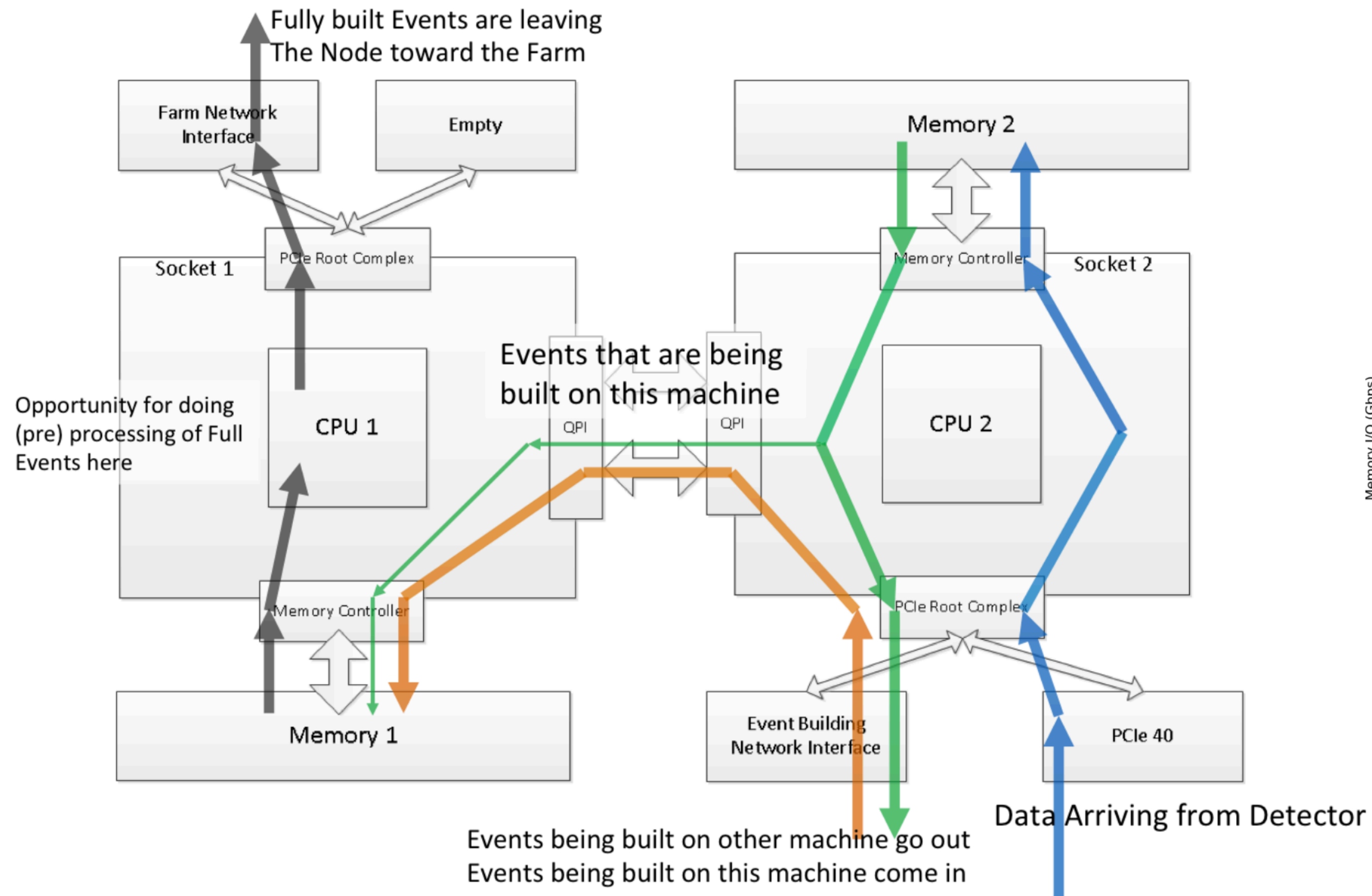  - Memory: very limiting O(MB) per core

  - Parallelism

**Daniel Cámpora**

# The LHCb DAQ

**Daniel Cámpora**

# Event Building



Event
Readout /
Building /
HLT1
filtering

32 Tb/s

1 Tb/s

Intermediate
storage

173 EB servers (with TELL40s)

17 storage servers

200G IB

100GbE

10GbE

- Event fragments are collected from the detector in readout cards (PCIe40)

- Data from these fragments is distributed to one destination at a time

- Events are fully built

- Events are processed by HLT1

- Finally, events are sent to the Event Filter Farm

**Daniel Cámpora**

# It works! (2014)



http://cds.cern.ch/record/1701361

**Daniel Cámpora**

# Putting all together

- Event builder farm equipped with 173 servers

- Each server has 3 free PCIe slots
  - Can be used to host GPUs
  - Sufficient cooling & power
  - Advantageous to have GPUs as self-contained processors
  - Sending data to GPU is like sending data to network card

- GPUs map well into LHCb DAQ architecture
- HLT1 tasks inherently parallelizable
- Reduced bandwidth network between EB & CPU HLT2
- Cheaper & more scalable than CPU alternative
- ➡**Was chosen as the baseline for the upgrade!**
- ➡**Is implemented with O(200) Nvidia RTX A5000 GPUs**

32 Tb/s

1 Tb/s

200G IB

100GbE

10GbE

173 EB servers (with TELL40s)

CPU+RAM1
RU  BU

CPU+RAM2
RU  BU

Readout  Readout  200G EB net  Accelerator  10G HLT net

Readout  200G EB net  Accelerator  Accelerator  10G HLT net

GPU-equipped event builder PC, with traffic of all three readout cards.

16

**Daniel Cámpora**

About the hardware

# About parallel algorithms

About commissioning

# How to make a good parallel algorithm

- What degrees of parallelism does your problem have?

- How could you map this parallelism onto your hardware architecture?

- What memory patterns can you identify?

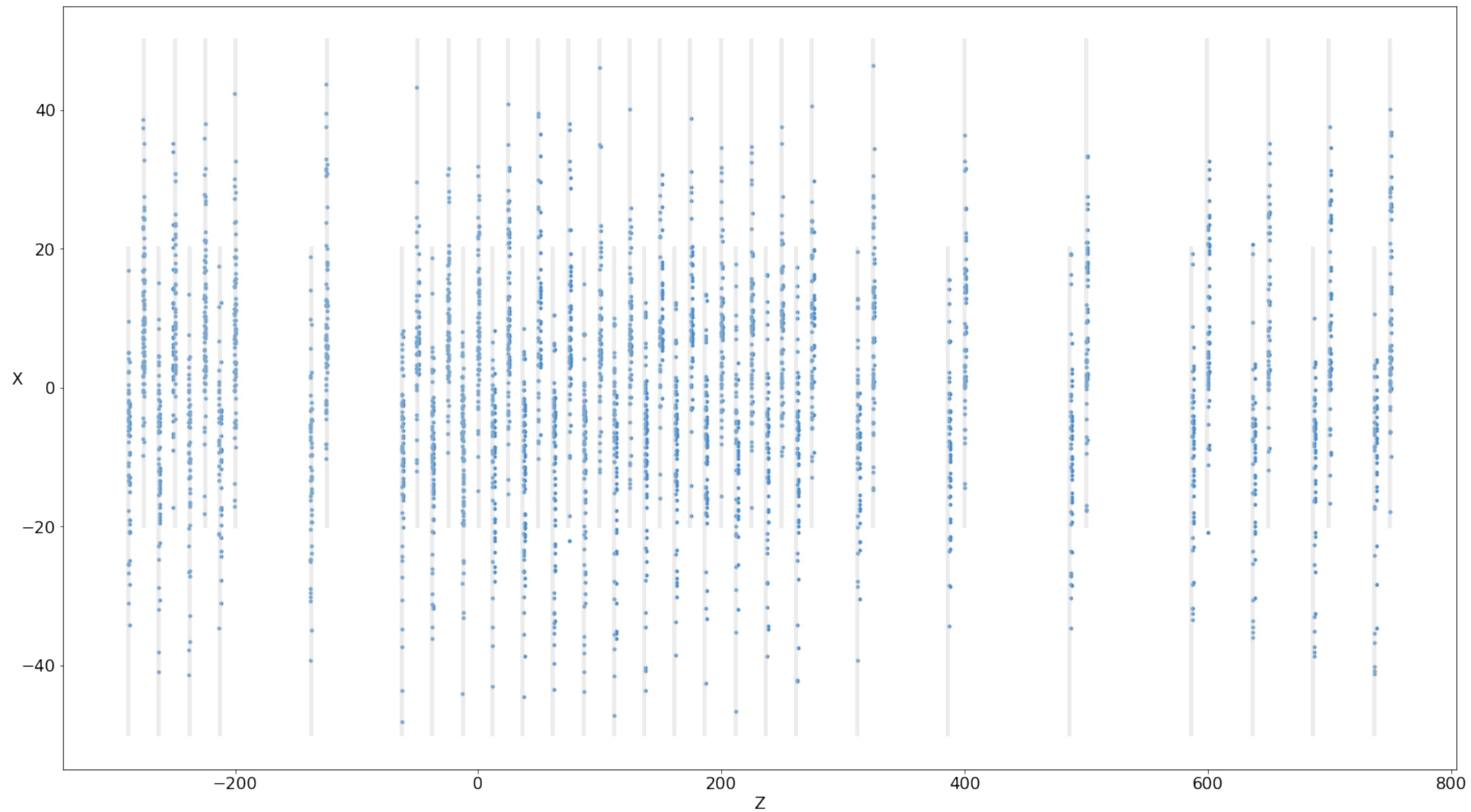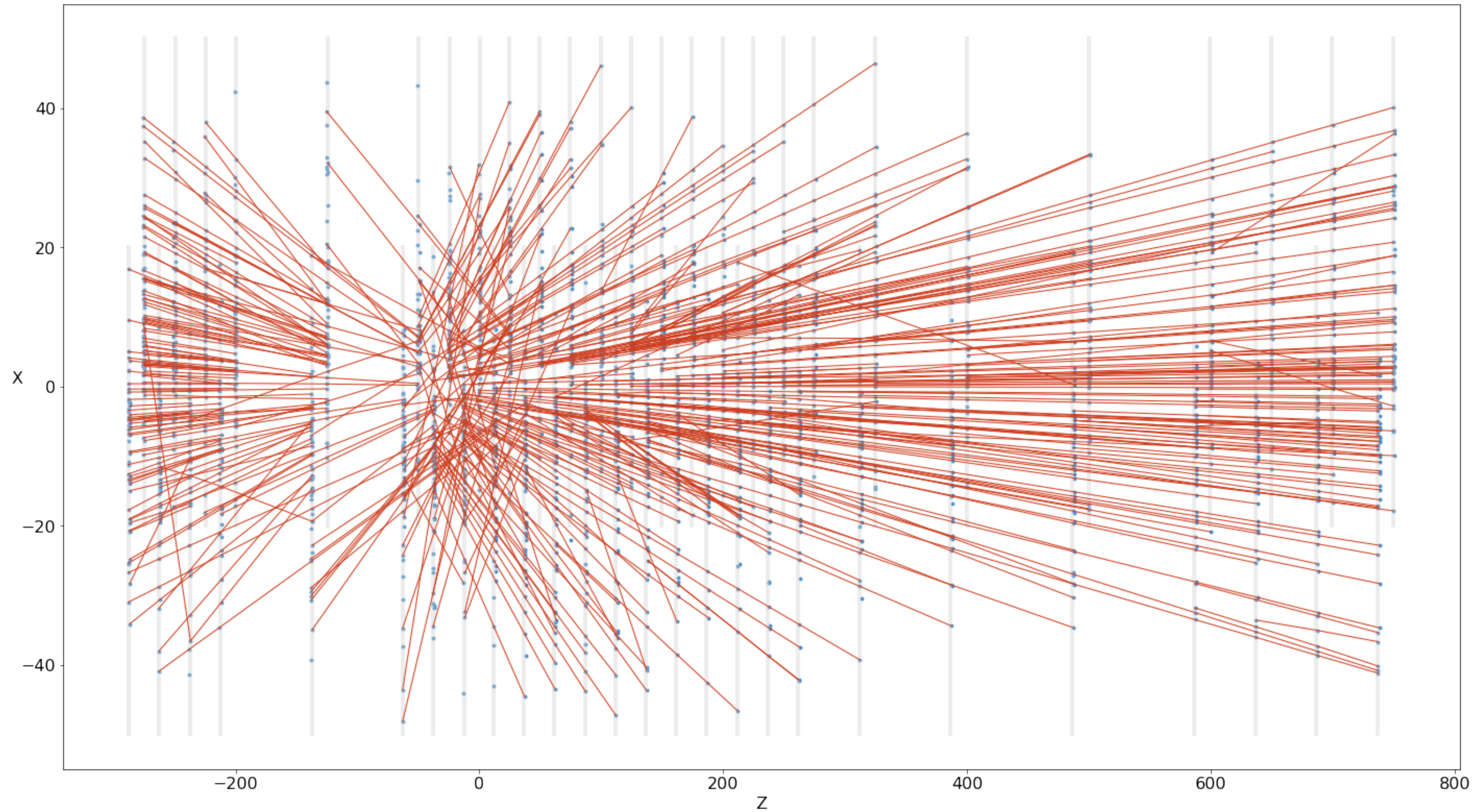- How can you map these patterns onto your hardware architecture?

**Daniel Cámpora**

# Track finding at LHCb





**VELO subdetector**

- The tracking system is composed of detectors: VELO, UT, SciFi (T stations).

- A magnetic field bends charged particles, we find out their **momentum**, **charge,** and the **collision / decay vertices** of the event

**Maastricht University**

**Daniel Cámpora**

# Going from this

**Daniel Cámpora**

# into this

**Daniel Cámpora**

# (at a high rate)



30,000,000 per second

**Daniel Cámpora**

# Taking VELO tracking as an example

- Rich literature of tracking methods

- Each event is physically independent

- Each track is independent

- Tracks come from a collision or secondary vertex

- VELO tracks are straight lines

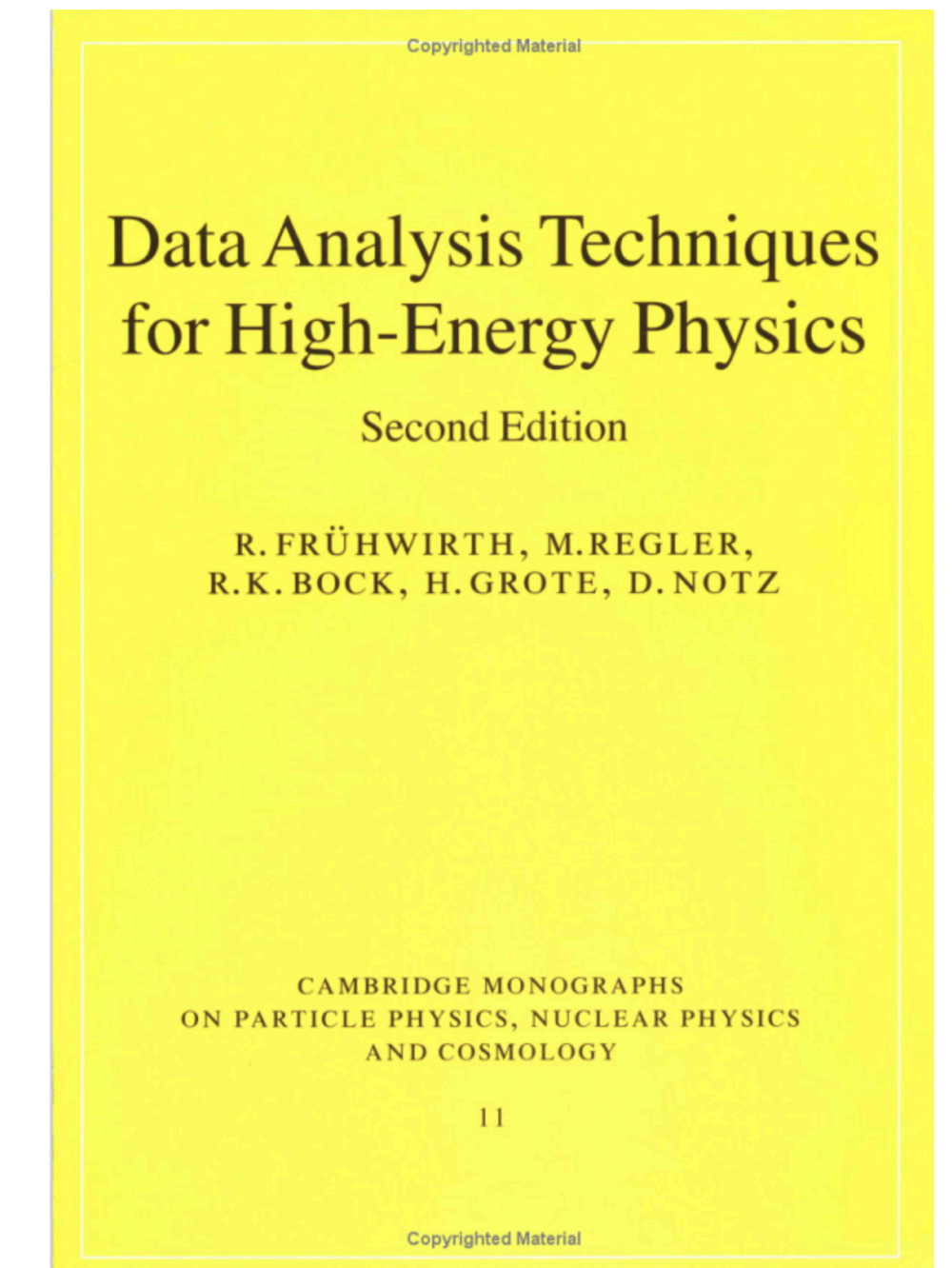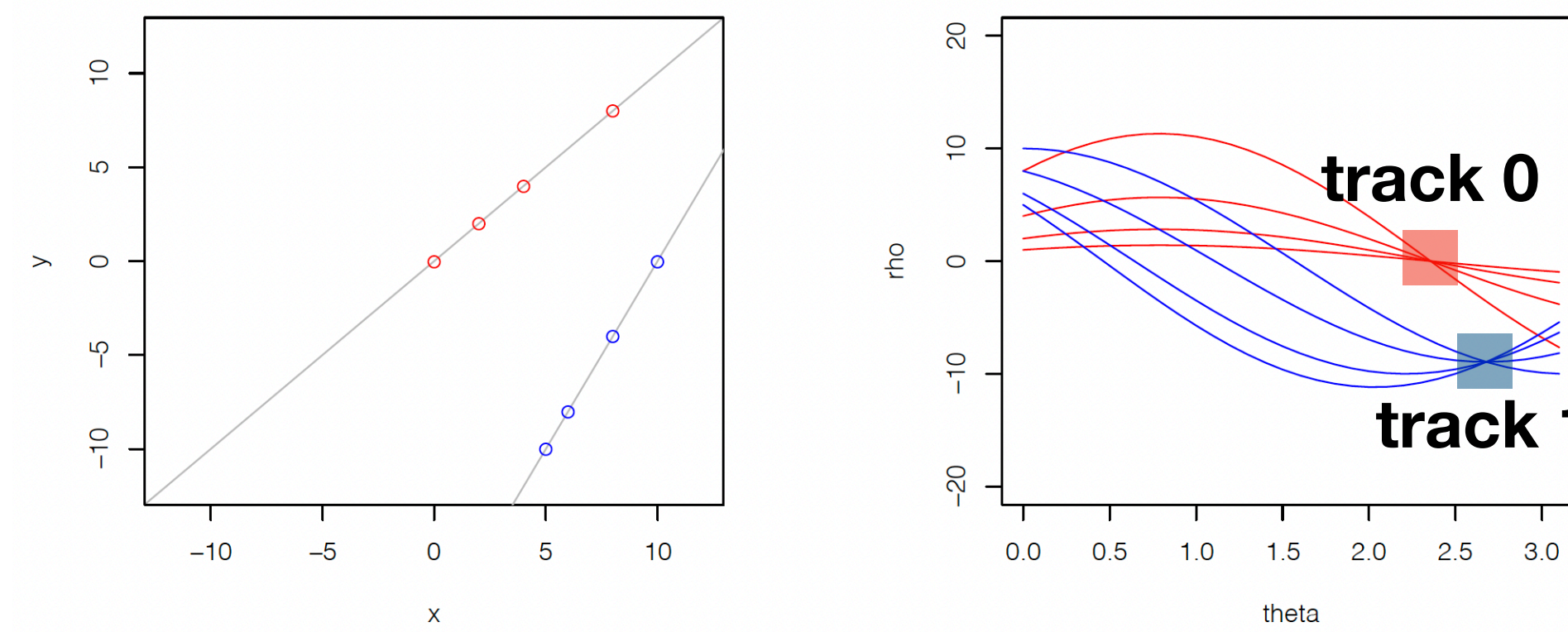- There is a geometrical distribution of hits across modules

**Daniel Cámpora**

# *Rich literature of tracking methods*

Local methods iteratively build a track and *follow* it:



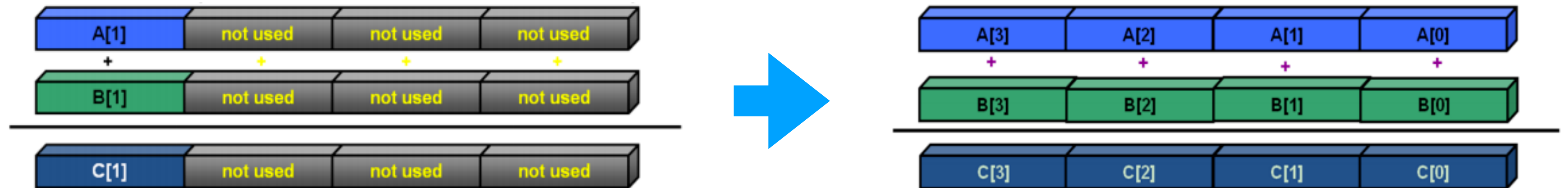Global methods map the problem to other equivalent formulations:

**Daniel Cámpora**

# *Each event is physically independent*

- We can therefore process them in separate CPU threads (as Gaudi does).

- On GPU, each event would be processed by a different *block.*

    - *Each block instantiates a different program, which is executed and managed by its own scheduler.*
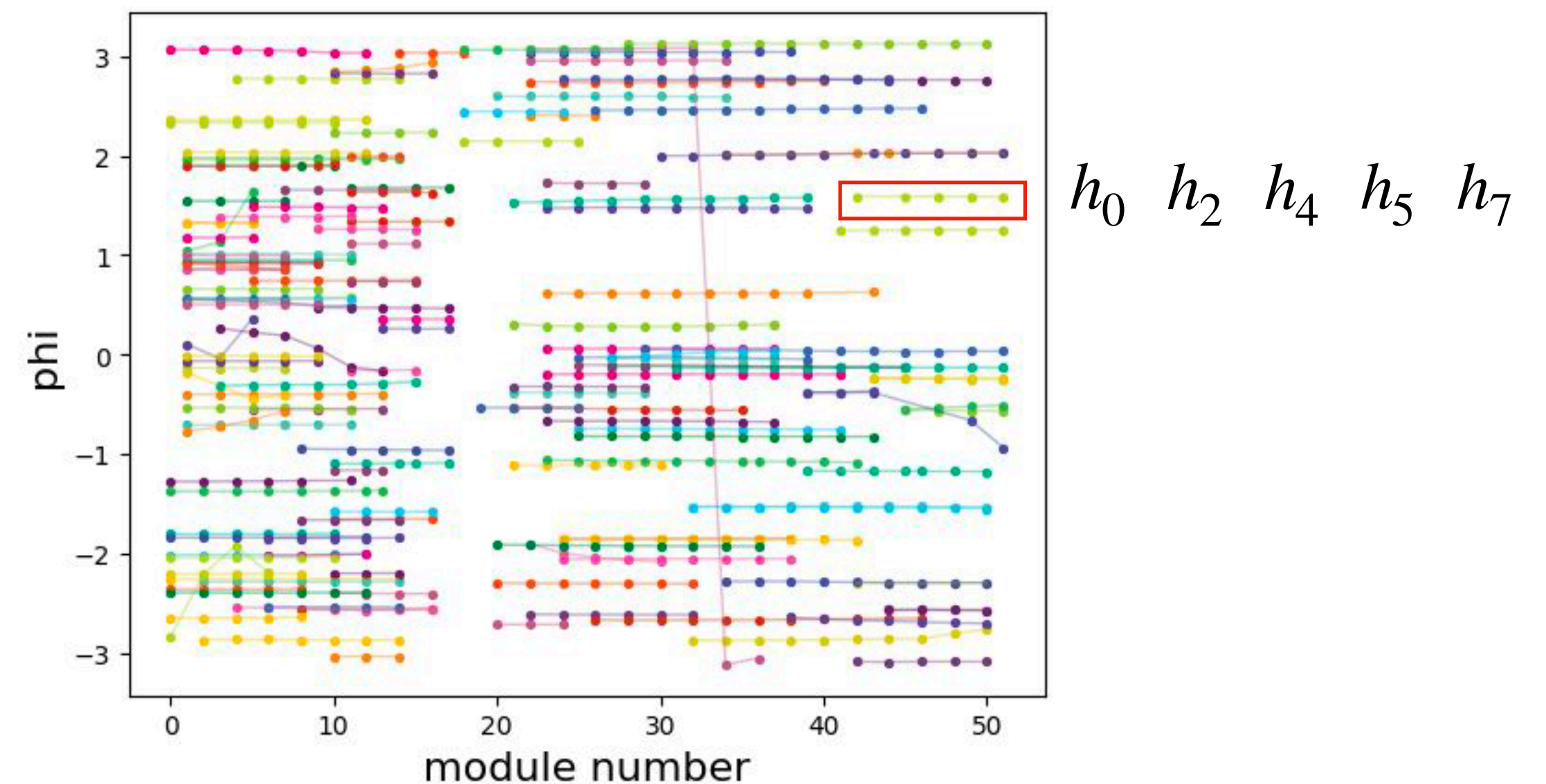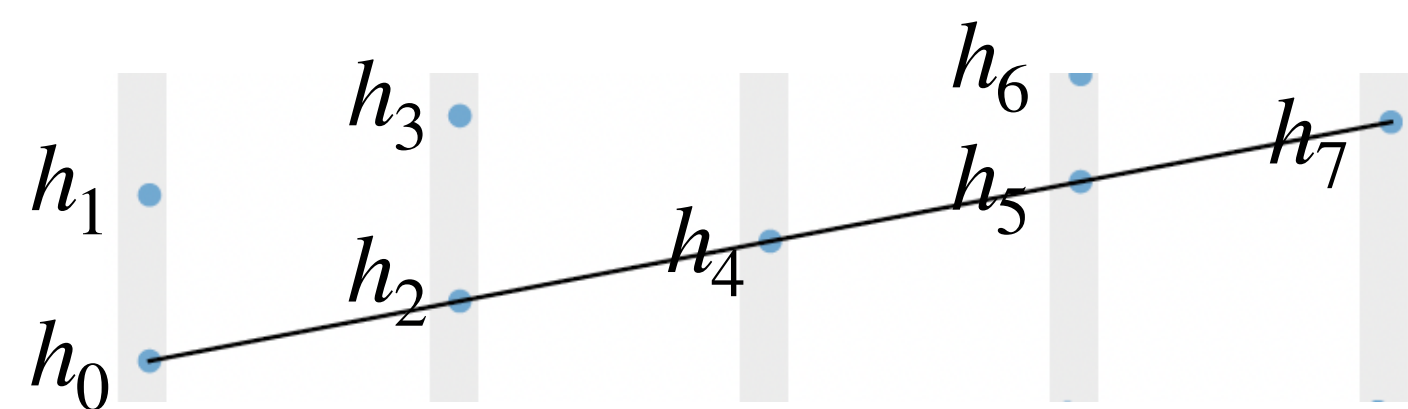
**Daniel Cámpora**

# *Each track is independent*

- Each thread can perform *Single-Instruction Multiple-Data* (SIMD) operations, processing several tracks in one go.

  - Modern CPUs have *SSE* and *AVX* extensions, leading to 4 or 8 FP32 simultaneous operations

  - GPUs use *warps* to control 32 threads at the same time

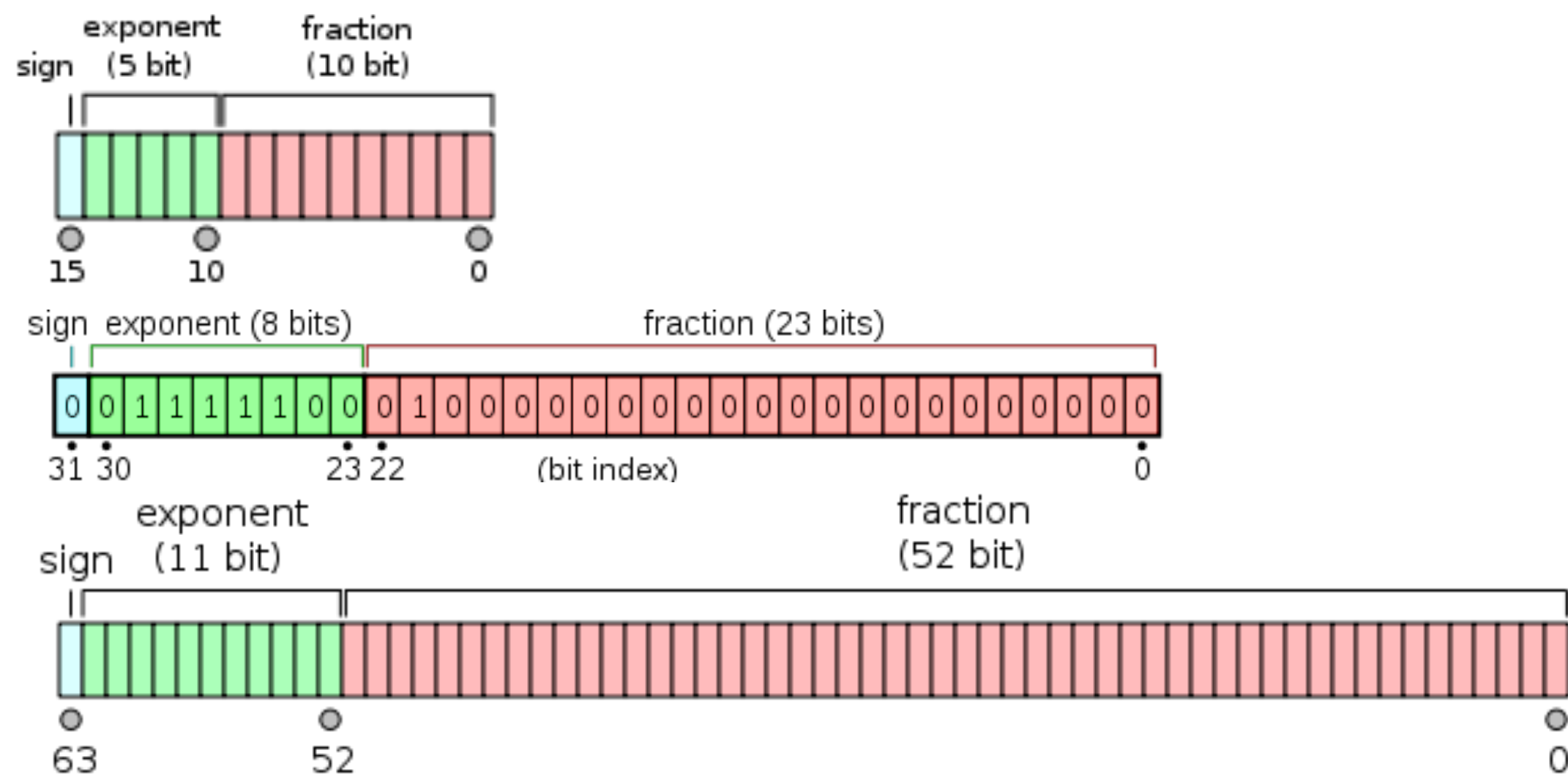  - Operations must be homogeneous to take advantage

**Daniel Cámpora**

# *Tracks come from vertices*

- We can exploit geometrical properties from the tracks to prepare an efficient data-structure.

  - In this case, sorting by phi is a good idea

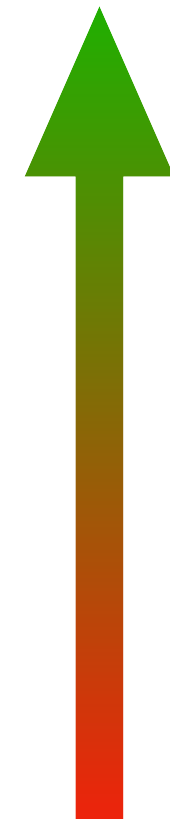  - Multi-dimensional structures (think 4D tracking) also exist

**Daniel Cámpora**

# *VELO tracks are straight lines*

- The model for these tracks is not complex

- Use the necessary precision for your problem. The lower the better.

  - Distinguish between **arithmetic** and **storage**



| Arithmetic | Storage |
|---|---|
| Double precision | Double precision |
| Double precision | Single precision |
| Single precision | Single precision |
| Single precision | Half precision |
| Half precision | Half precision |

Precision ↑

Speed ↓

**Daniel Cámpora**

# Little exercise:
# What's wrong with this piece of code?

```
1  __global__ void shared_memory_example(float* dev_array) {
2    for (int i = threadIdx.x; i < 256; i += blockDim.x) {
3      dev_array[i] = 1 / std::sqrt(2. + dev_array[i]);
4    }
5  }
```

**Daniel Cámpora**

# Little exercise (2)

```
1  __global__ void shared_memory_example(float* dev_array) {
2    for (int i = threadIdx.x; i < 256; i += blockDim.x) {
3      dev_array[i] = 1 / std::sqrt(2. + dev_array[i]);
4    }
5  }
```
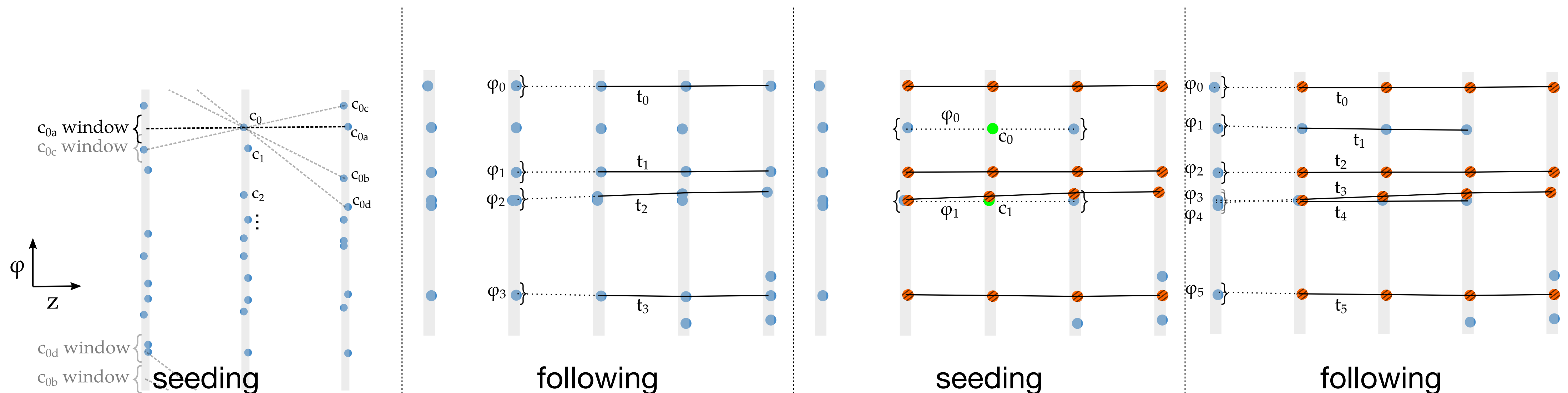
*Use compiler flag -Wdouble-promotion to avoid surprises!*

**Daniel Cámpora**

# Little exercise (3)

```
1  __global__ void shared_memory_example(float* dev_array) {
2    for (int i = threadIdx.x; i < 256; i += blockDim.x) {
3      dev_array[i] = 1 / std::sqrt(2. + dev_array[i]);
4    }
5  }
```

*Use compiler flag -Wdouble-promotion to avoid surprises!*

*... and come to the thematic CERN School of Computing to learn more!*
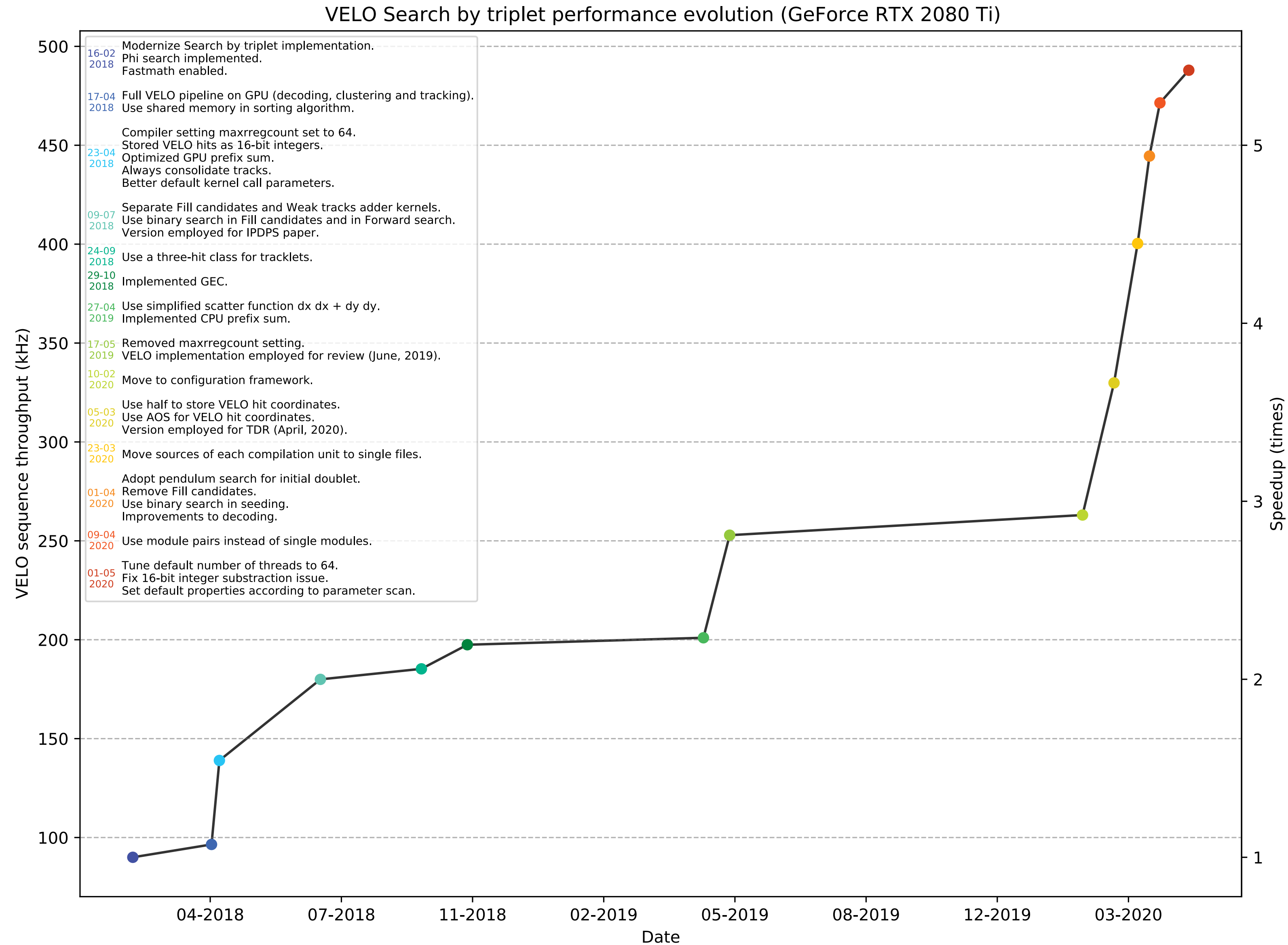
**Daniel Cámpora**

# *Geometrical properties*

- The geometry of the detector should be used to find good access patterns

- Principle of locality:

  - **Spatial locality** - Prefer to access neighbouring data in memory

  - **Temporal locality** - Prefer reusing accessed data soon after first access
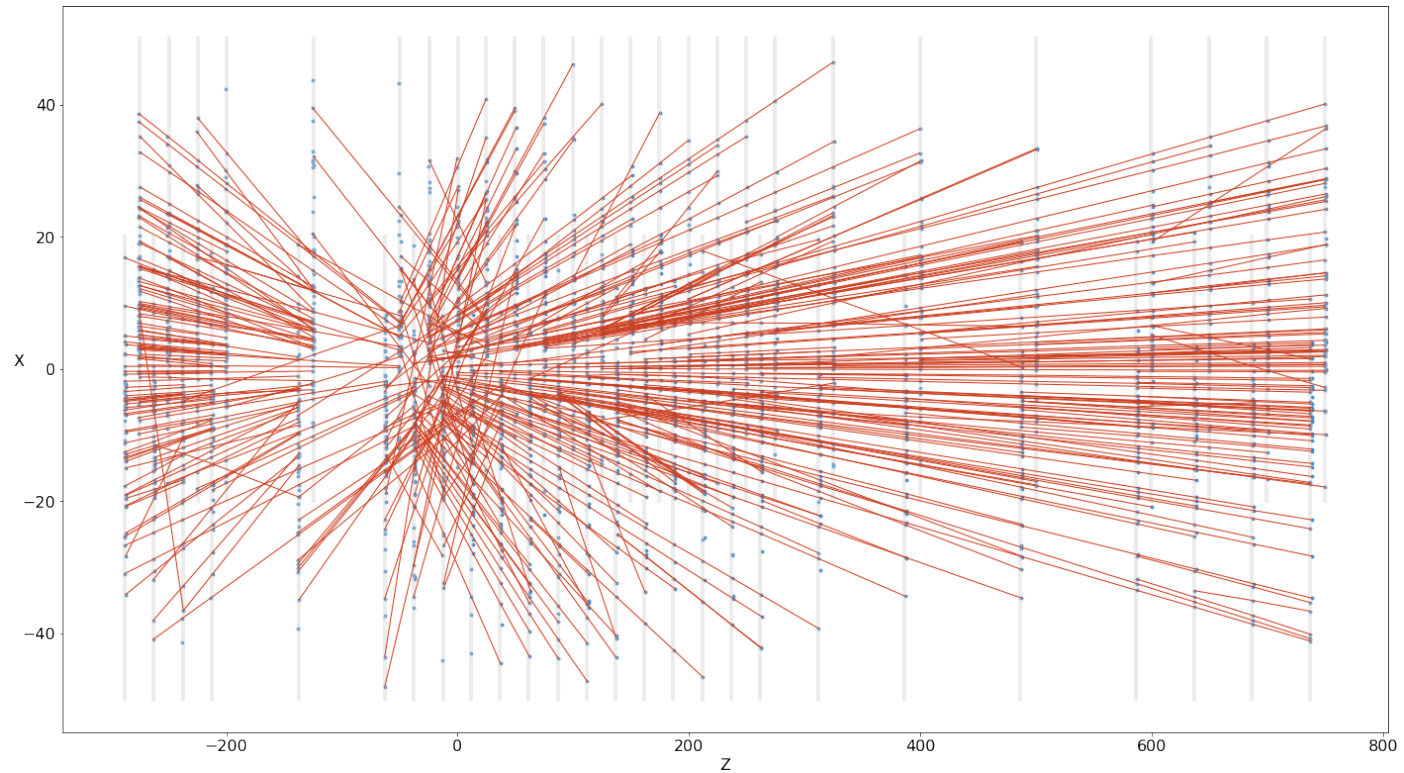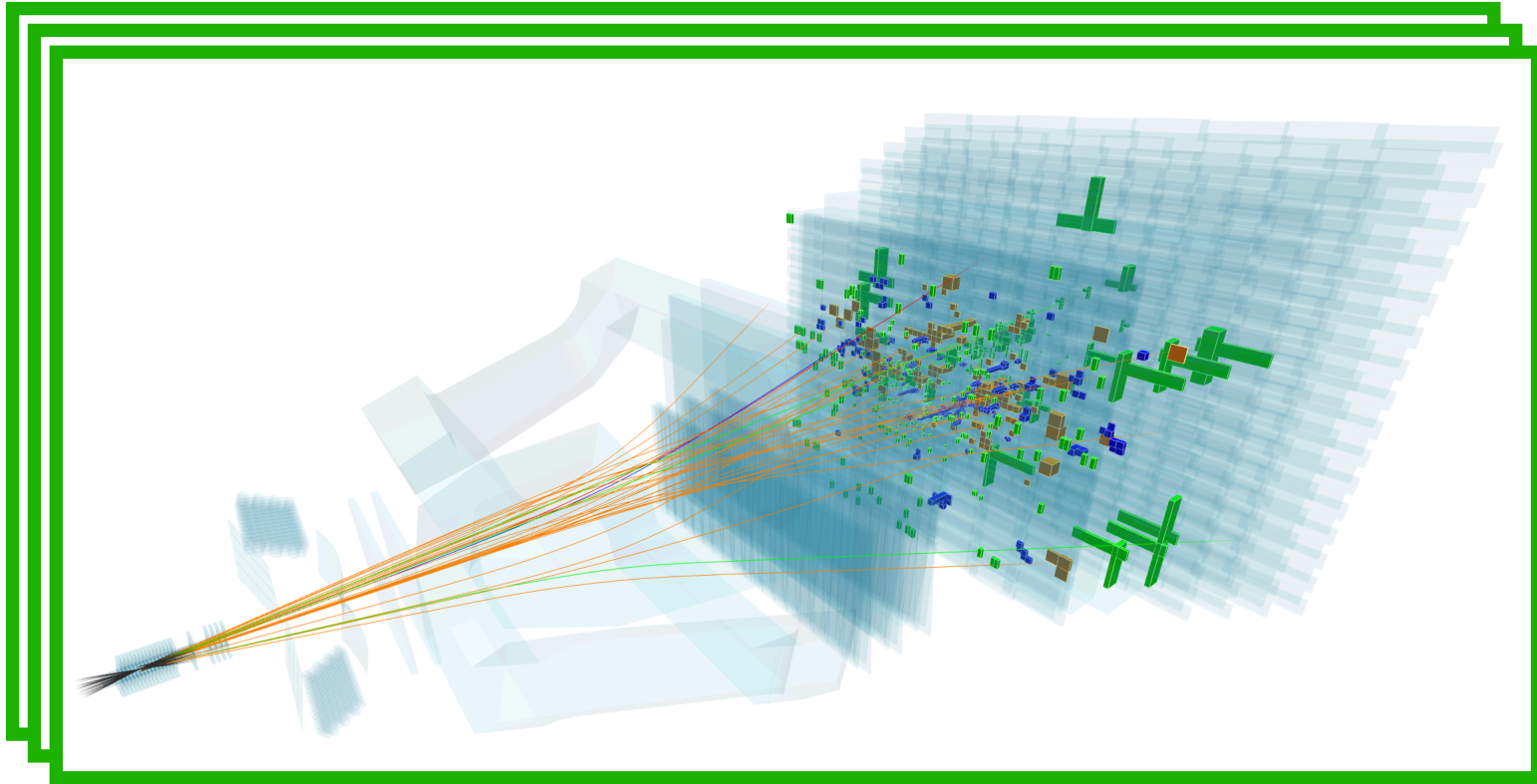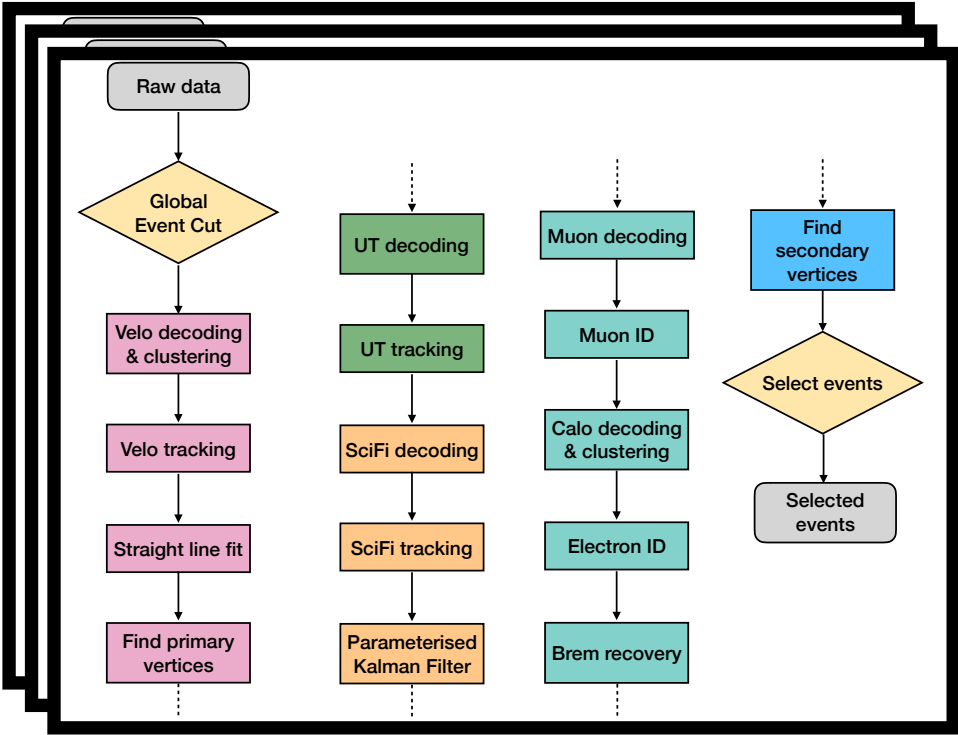
**Daniel Cámpora**

# Applying all these principles



VELO Search by triplet performance evolution (GeForce RTX 2080 Ti)

33

**Daniel Cámpora**

# Parallelization

We parallelize at three levels:

- **Sequences**
- **Events**
- **Intra-Algorithms**



| | Sequences | Events | Intra-Algorithms |
|---|---|---|---|
| **CPU** | | Threads | Vectorisation |
| **GPU** | Streams | Blocks | Threads |

**Daniel Cámpora**

# Track reconstruction at HLT1

**Velo tracking:**           *Journal of Computational Science, vol. 54, 2021*

- 52 silicon pixel modules with $\sigma_{x,y} \sim 5\ \mu m$

- Parallel local tracking algorithm: *Search by Triplet*

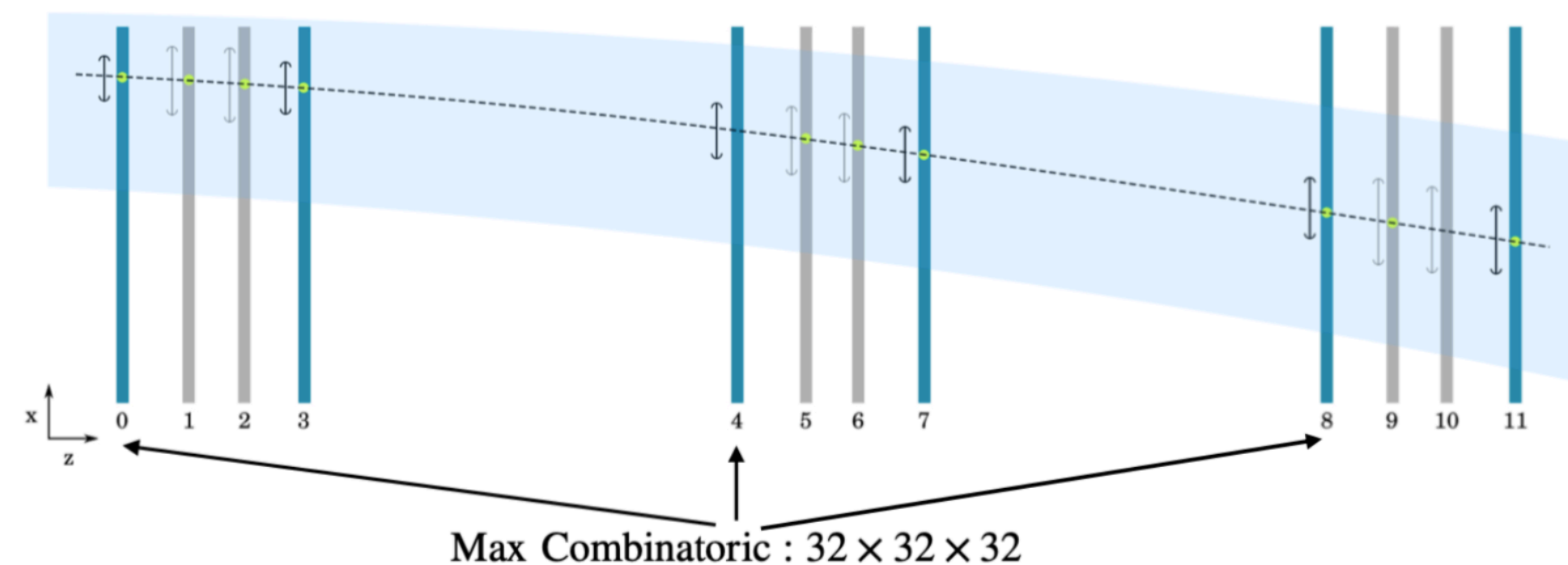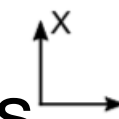- Tracks fitted with simple Kalman filter assuming straight line model

**Velo-UT tracking:**           *IEEE Access, vol. 7, pp. 91612-91626, 2019*

- 4 layers of silicon strips

- Velo tracks extrapolated to UT taking into account fringe B-field

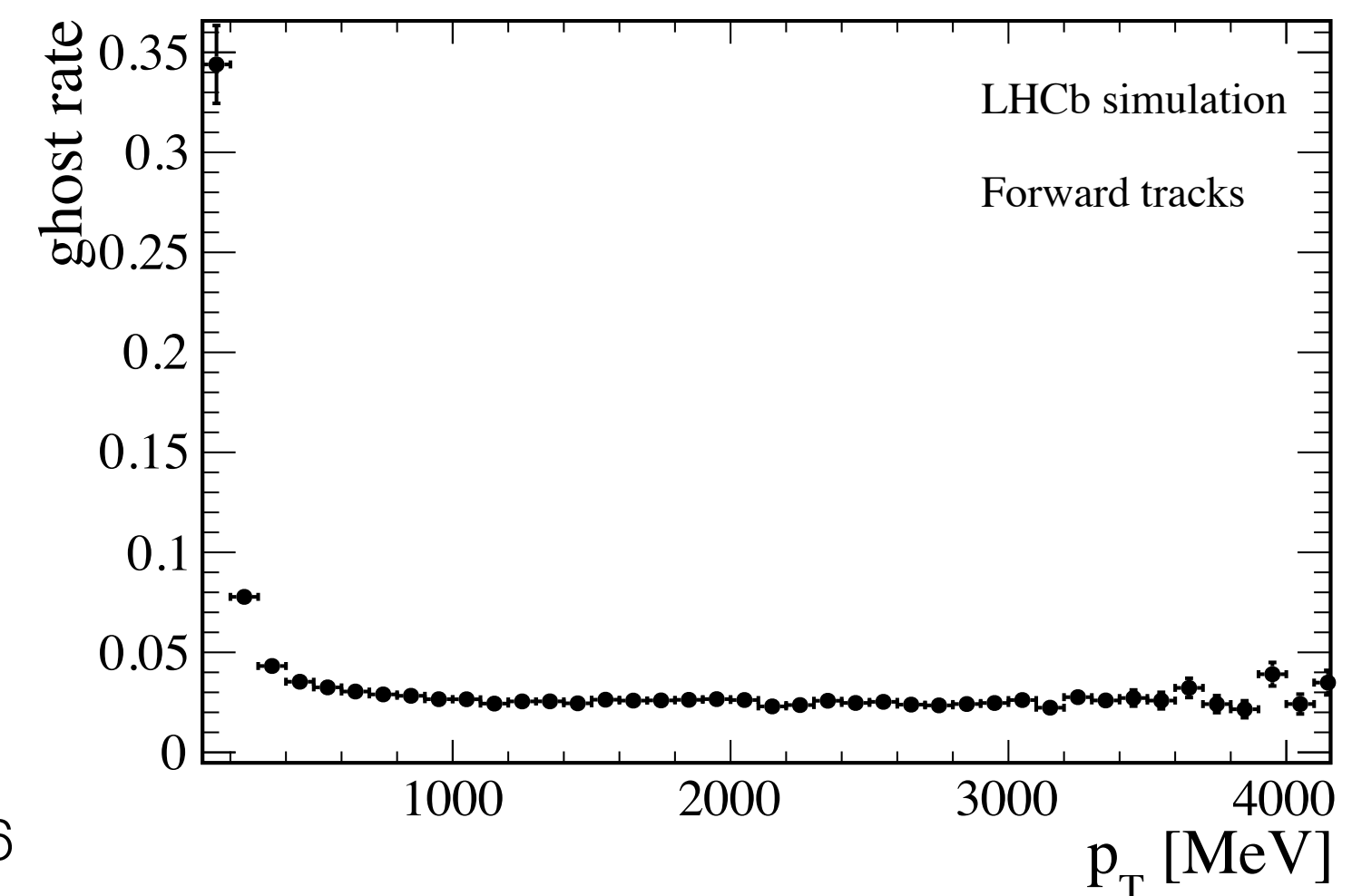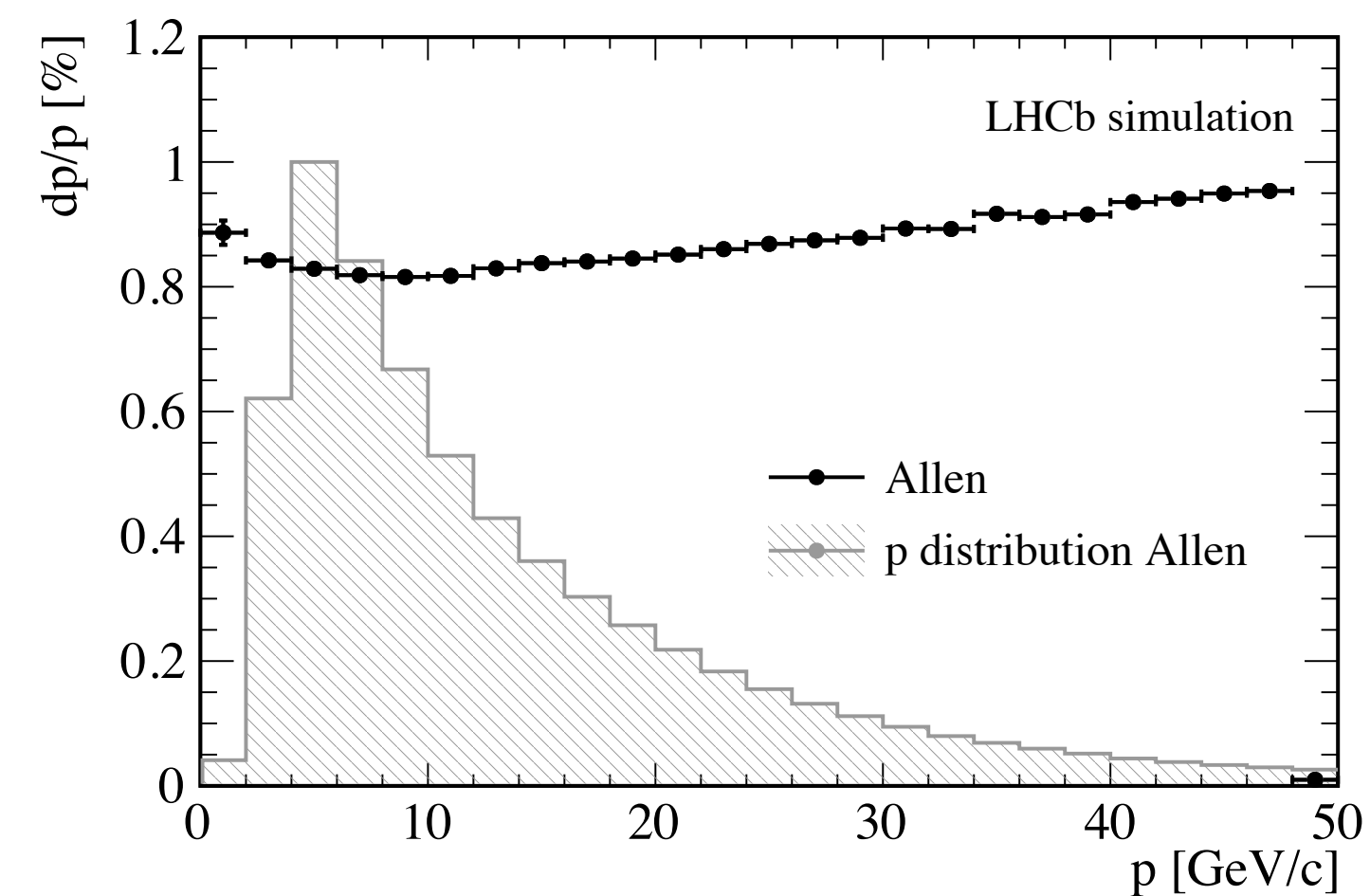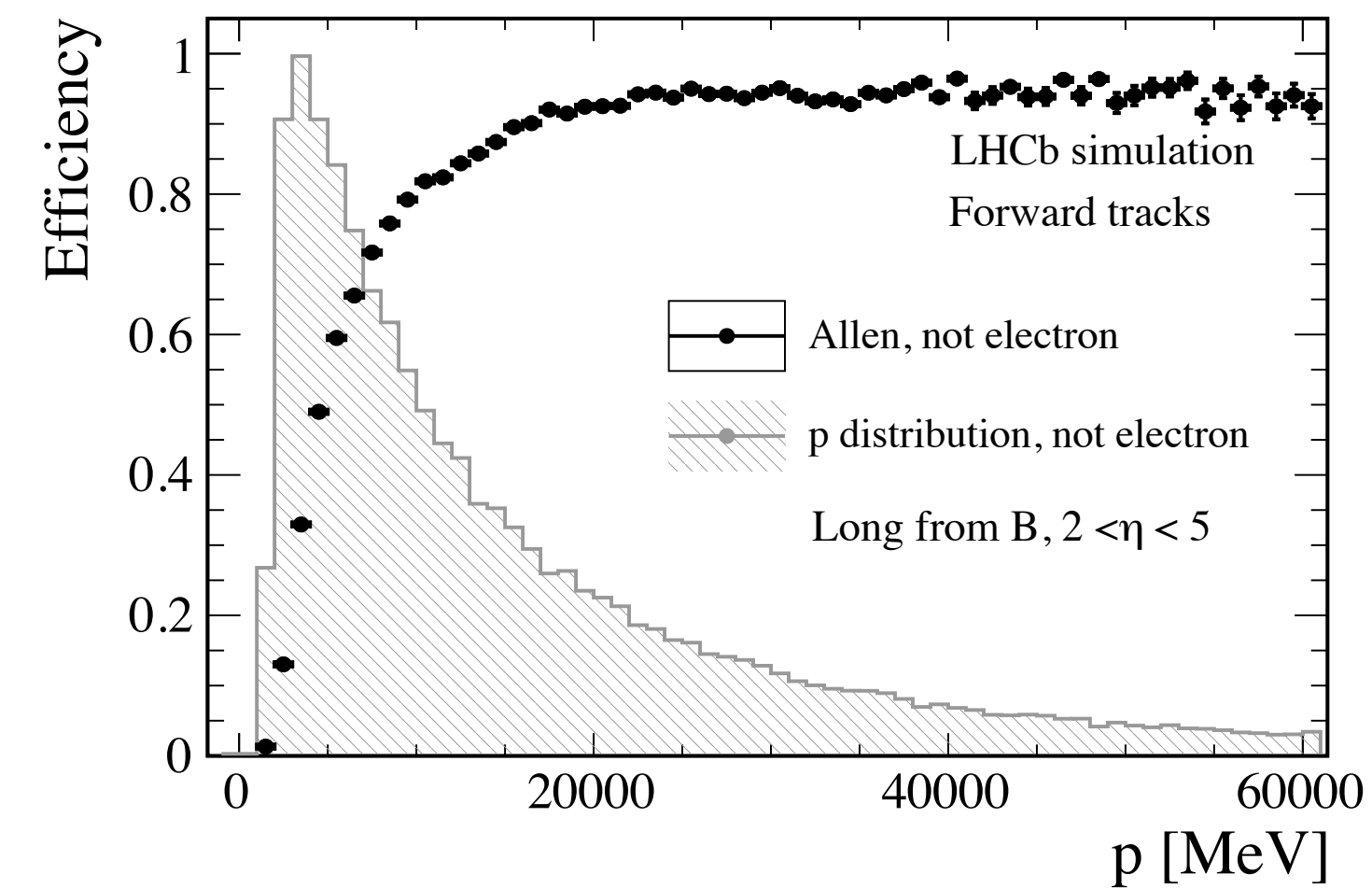- Parallelized tracklet finding inside search windows requiring at least 3 hits

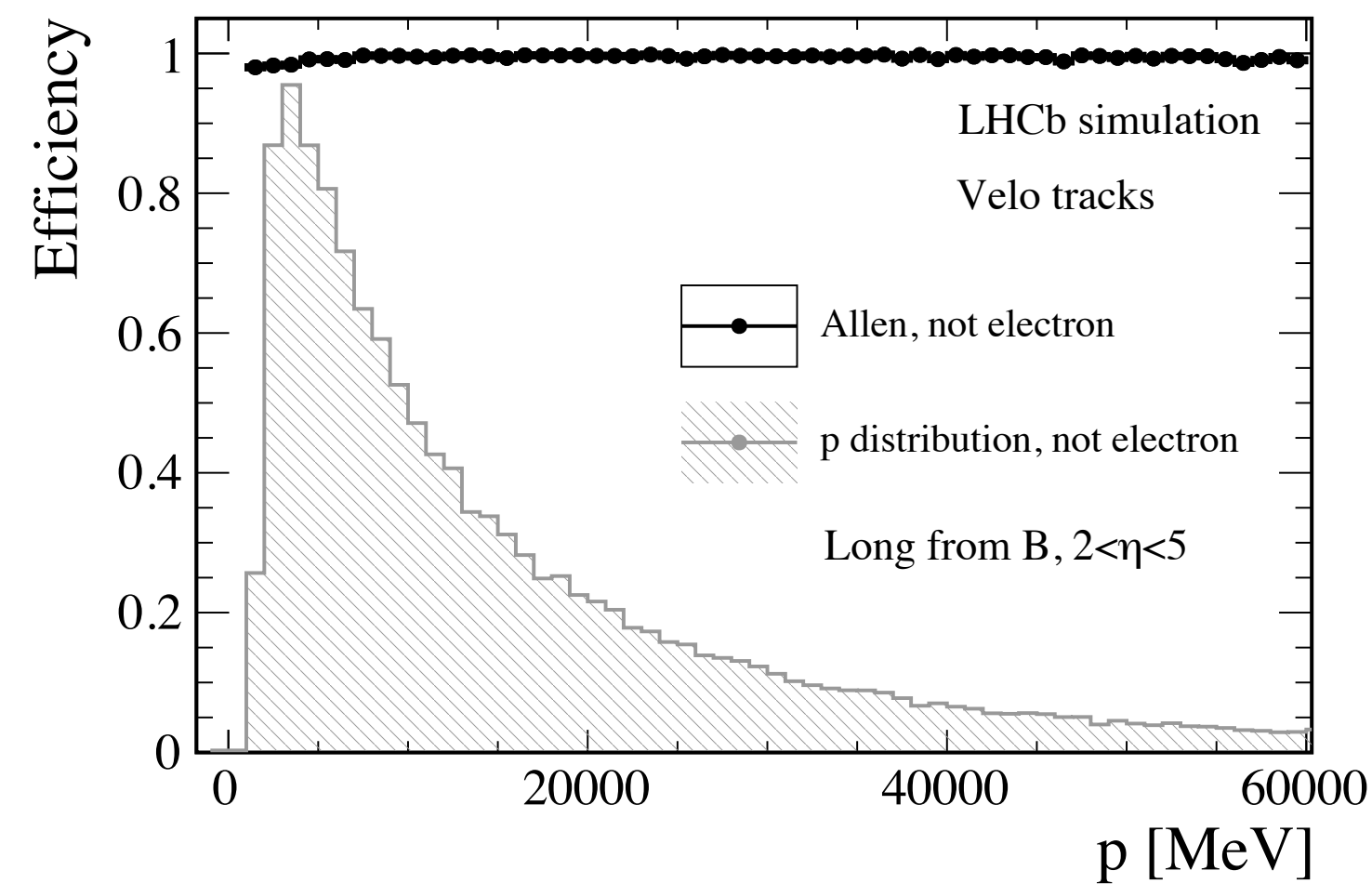**Forward tracking:**           *Comput Softw Big Sci 4, 7 (2020)*

- 3 stations with 4 layers of Scintillating Fibres

- Velo-UT tracks extrapolated using parametrization

- Parallelized *Forward algorithm* to reconstruct **long tracks:**

    - Search windows from on Velo-UT momentum estimate

    - Form triplets and extend to remaining layers

35
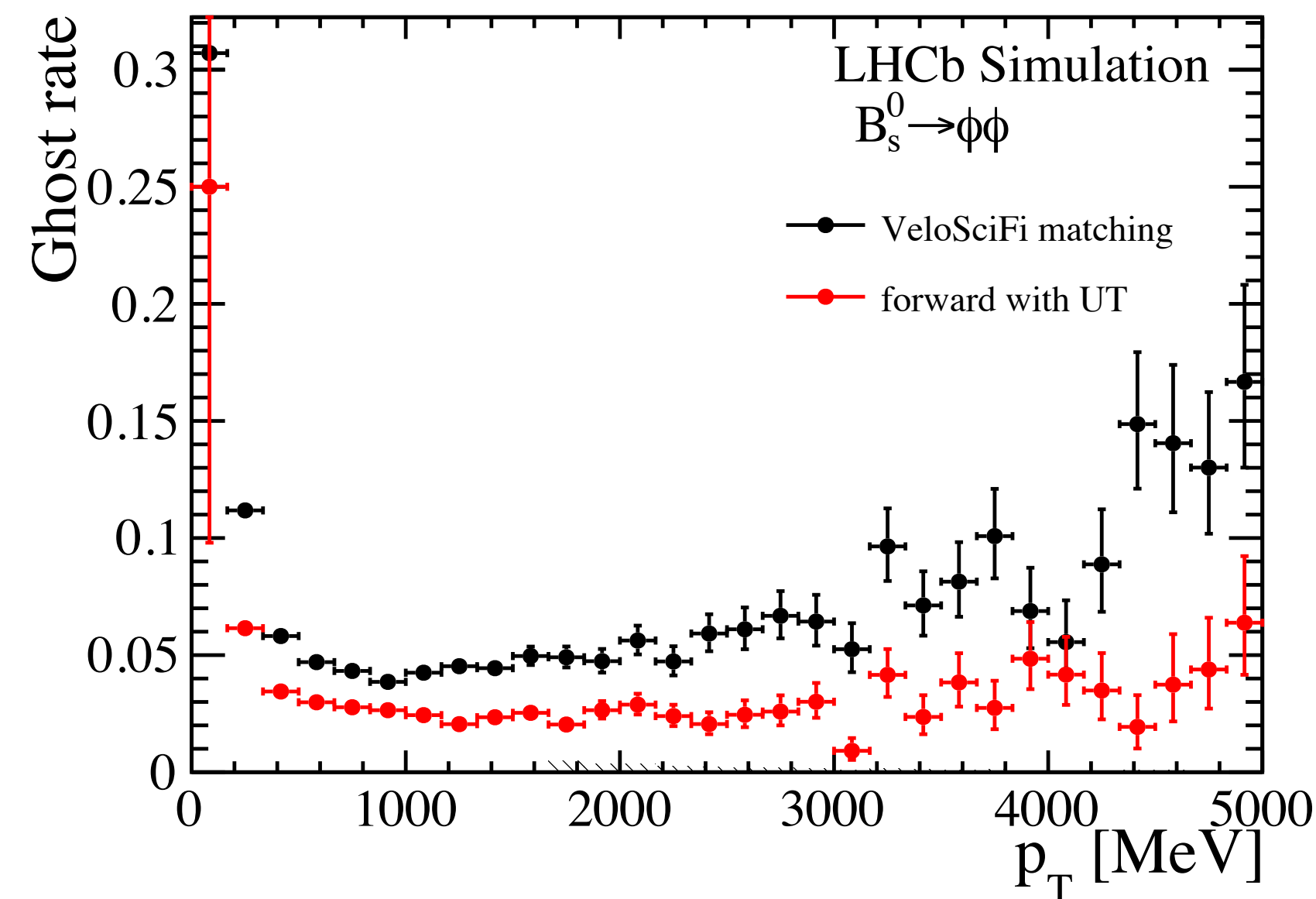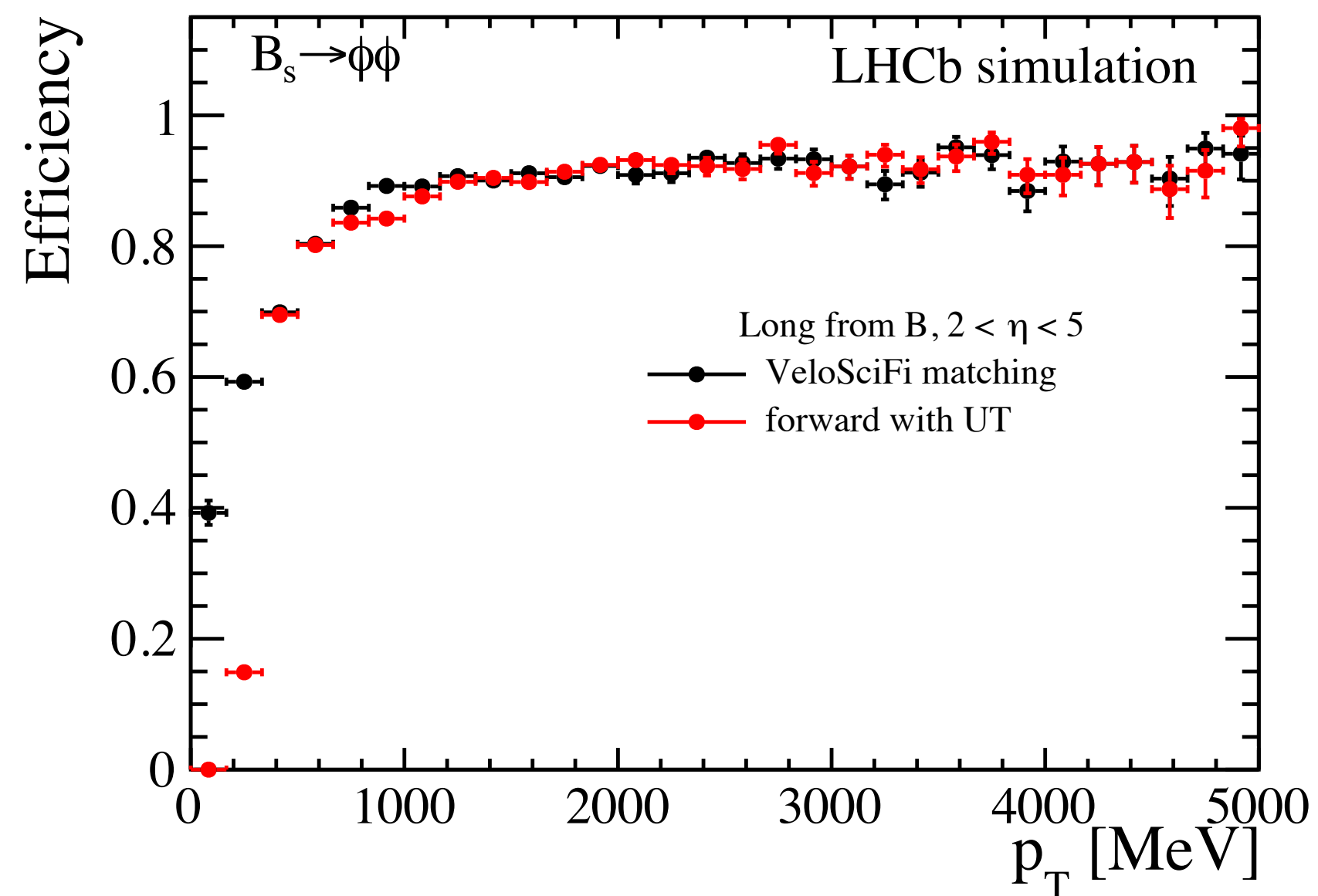
**Daniel Cámpora**

# HLT1 tracking efficiency

- **Run 2 efficiency maintained at x5 instantaneous luminosity**

- Excellent track reconstruction efficiency (> 99% for VELO, 95% for high-p forward tracks)

- Good momentum resolution and fake rejection

36

**Daniel Cámpora**

# Tracking without the UT

- In 2022, the UT detector will unfortunately not be available for data-taking

- Tracking performance and throughput maintained, at the cost of larger fake rate

- Opportunity to commission 2 options, which **both maintain the current throughput:**
  - **Forward without UT**
  - **Seeding+Matching:**
    - Standalone SciFi reconstruction & matching to VELO seeds
    - Highly efficient for low momenta
    - Opens the door to additional physics cases in HLT1 (downstream and SciFi tracks)

**Daniel Cámpora**

# Vertex reconstruction

- Primary vertices found from **clusters** in the closest approach of tracks to the beamline
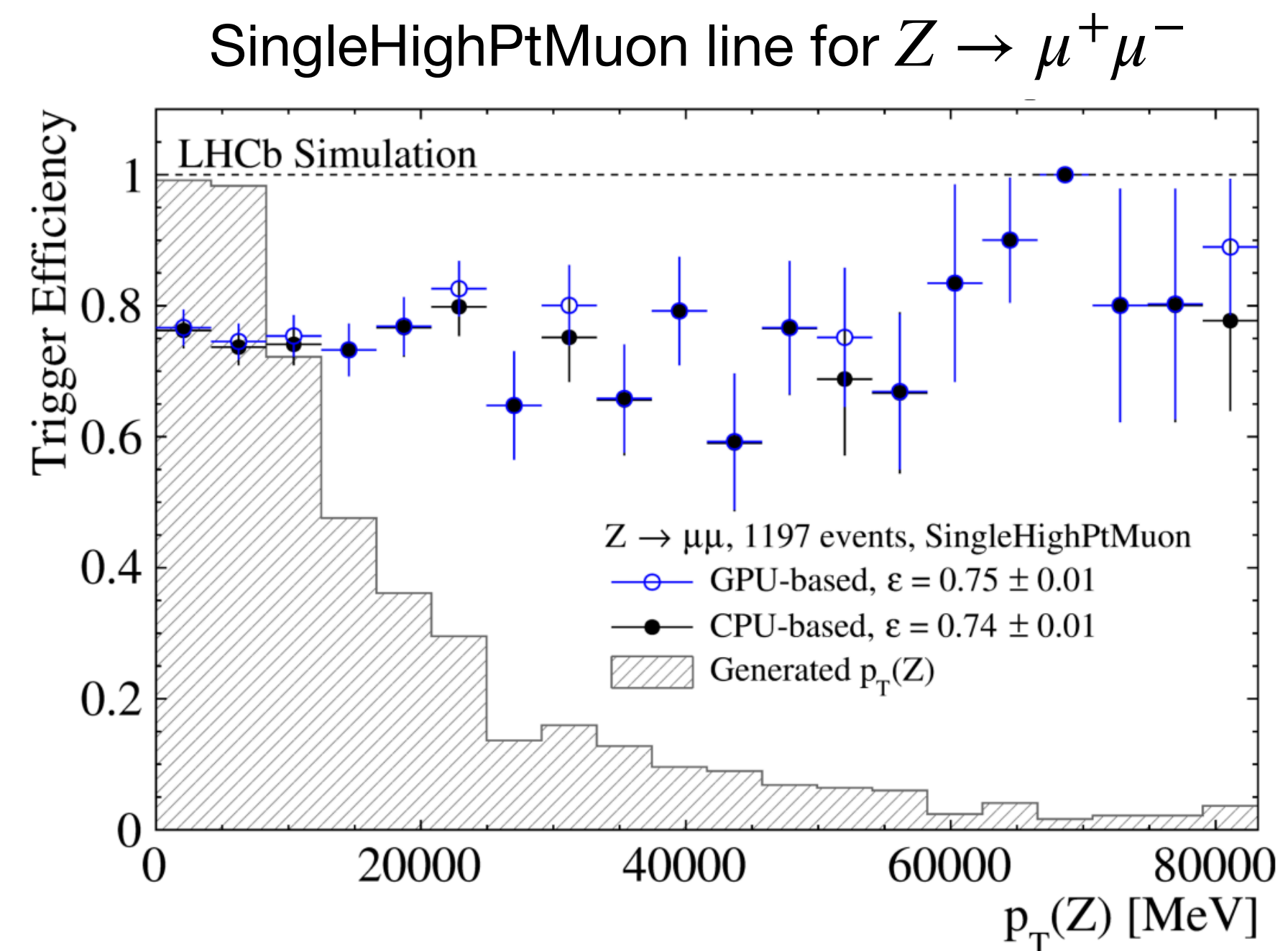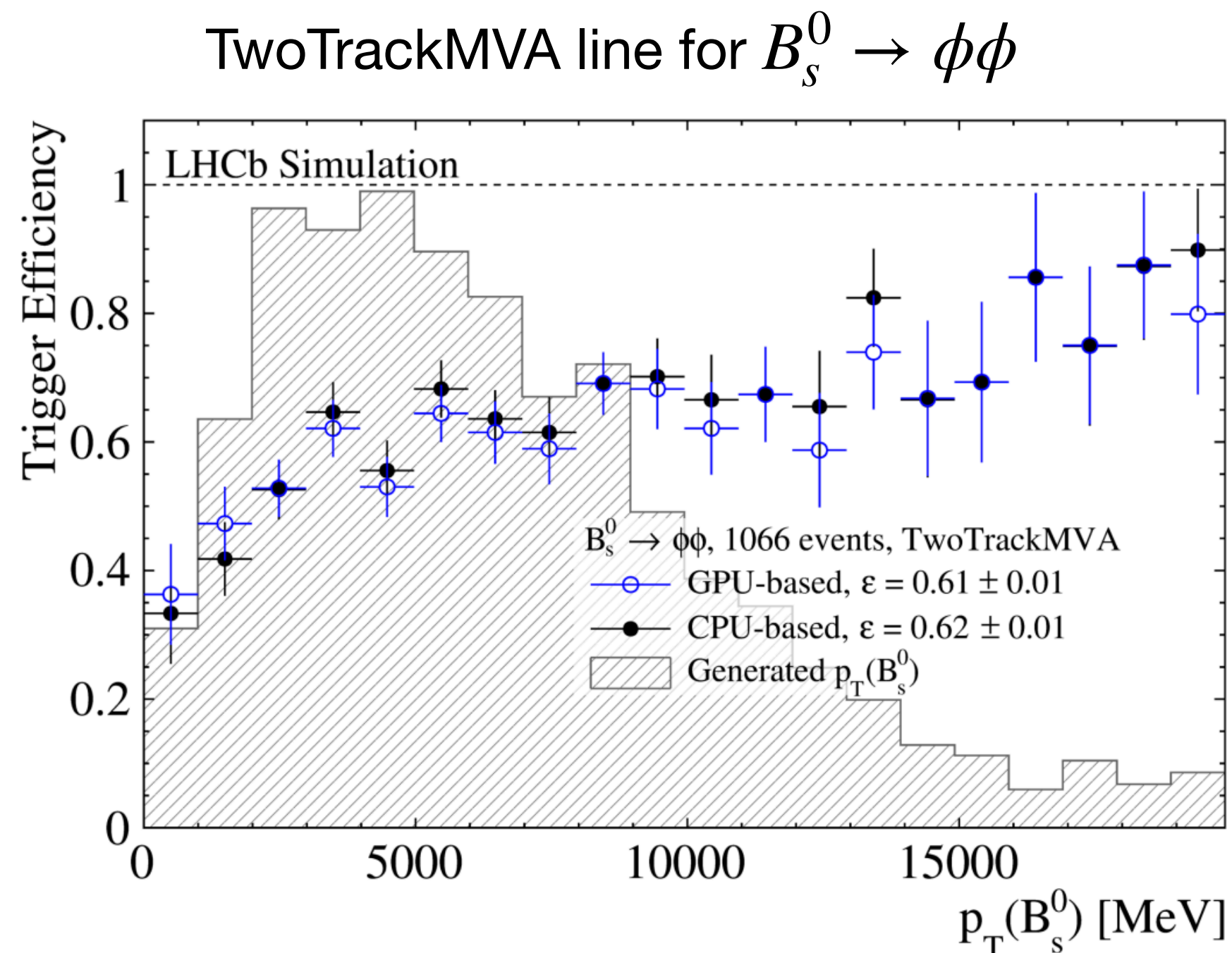
- 1-1 mapping between tracks and vertices requires **serialization**

  ‣ Instead, every track assigned to every vertex based on **weight**

- **Efficiency > 90%** for vertices with N. tracks > 10



*Comput. Softw. Big Sci. 6 (2022) no.1, 1*

**Daniel Cámpora**

# HLT1 selection performance

- Inclusive rate for the main HLT1 lines ~ 1 MHz

- O(30) lines implemented so far:
  - Cover majority of LHCb physics program (B, D decays, semileptonic, EW physics)
  - Special lines for monitoring, alignment and calibration
  - Additional trigger lines under development
  - Compatible performance between CPU and GPU

TwoTrackMVA line for $B_s^0 \to \phi\phi$

SingleHighPtMuon line for $Z \to \mu^+\mu^-$

**Daniel Cámpora**

# Throughput

- 30 MHz goal can be achieved with O(200) GPUs (maximum the Event Builder server can host is 500)

- Throughput scales well with theoretical TFLOPS of GPU card

- Additional functionalities are being explored



LHCb-FIGURE-2020-014

**Daniel Cámpora**

# About commissioning

# Previously...

LHCb October beam-test was major milestone for the Allen commissioning

Steps for the future:

- As more sub-detectors get installed, commission more parts of the **decoding**, **reconstruction** and **selection** chain

- Commission the **full chain** (EB → HLT1 → HLT2 → storage & offline )

- **Monitoring**

- Continue the *installation of GPU cards* in the LHCb Data center

- **Throughput**, **memory**, **cooling** and **stability** tests with larger-scale system

- *Data taking* with stable beams expected to start in spring 2022!



Data center

PCIe40 board

PCIe40 server

**Daniel Cámpora**

# Over the last months

The software is anything but frozen

Use the right tx and ty definition in TrackCheckerHistos.cpp
!839 · created 2 months ago by Lorenzo Pica  MC checking  RTA  bug fix  ci-te

Update Muon
!822 · created

Fix RICH line
!842 · created

Fix bug in Ma
!845 · created

Prepare Aller
!841 · created

Beam Gas Lir
!819 · created

Follow strear
!818 · created

Create and u
!832 · created

Use rocm-sm
!833 · created

make XZ fit independent of hit caching
!831 · created 3 months ago by Louis Henry  ⅄ matching_velosc

CI: whitelisted mi100-full: [run_throughput, RelWithDebInfo, hlt
!830 · created 3 months ago by Ryunosuke O'Neil  RTA  only Git

Explicitly enforce all Allen datatypes be trivially copyable
!821 · created 3 months ago by Daniel Hugo Campora Perez  RTA

Prevent full run jobs from starting when they shouldn't
!829 · created 3 months ago by Ryunosuke O'Neil  RTA  hlt1-thro

Removed duplication of channelIDs between LHCb and Detector
!825 · created 3 months ago by Sebastien Ponce  RTA

Remove PatPV
!820 · created 3 months ago by Daniel Hugo Campora Perez  RTA

Update build.rst
!828 · created 3 months ago by Daniel Hugo Campora Perez

Use CMAKE_TOOLCHAIN_FILE, various cmake changes
!797 · created 3 months ago by Daniel Hugo Campora Perez  Buil

"AVOID_HIP" hack for full run matrix to avoid running certain tests
!826 · created 3 months ago by Ryunosuke O'Neil  ⅄ dcampora_remove_un

adapt CI to CMAKE_TOOLCHAIN configuration; pesky HIP tests are disable
!868 · created 1 mo

Use proper bank si
!814 · created 3 mo

Update Refe
!901 · create

Update geo
!879 · create

Add algorith
!835 · create

Fix SMOG2 unstab
!872 · created 1 mo

Fix use of yaml-cp
!871 · created 1 mo

Fix incl   Decode R
!870 · c   !805 · crea

Add DD   Detect MI
!855 · c   !859 · crea

SMOG2   RetinaClu
!602 · c   !748 · crea

Resize   Fix docun
!772 · c   !851 · crea

Make A   Update Re
!863 · c   !852 · crea

Fix for f   Add the r
!858 · c   !849 · crea

Delayer   Better err
!846 · c   !840 · crea

Update refs
!904 · create

fix veloSP_validation input file
!886 · created 1 month ago by Giovanni Ba

Early Measurement HLT1 TwoTrack line f
!751 · created 5 months ago by Andre Gun

Refactor Allen geometries
!913 · created 2 weeks ago by Da

Remov
!869 · c

Update
!878 · c

tuned onetrack and twotrack
!909 · created 2 weeks ago by Ni

Update References for: Allen!90
!914 · created 2 weeks ago by So

Dielectron l
!774 · create

Run ch
!882 · c

Removed unused options cpu_o
!908 · created 2 weeks ago by Da

Draft: this is
!898 · create

Remov
!881 · c

Set logger::debug as minimum v
!912 · created 2 weeks ago by Da

Really fix #3
!899 · create

Remov
!875 · c

Update References for: LHCb!3
!910 · created 2 weeks ago by So

Fix check-e
!897 · create

Fixed in
!880 · c

Refactor store, create Allen::bu
!902 · created 3 weeks ago by Da

Cleanup line
!885 · create

Fix war
!877 · c

Unified validators for long and k
!862 · created 1 month ago by Ch

Convert DumpMuon
!884 · created 1 mont

Fix start->stop->start and allow
!900 · created 3 weeks ago by Roel Aaij  RTA  bug fix  ci-test-triggered  hlt1-throughput-decreased  new feature

Add reference bot fo
!890 · created 4 week

Move para
!812 · created 3 months ago by Dorothea Vom Bruch  RTA  ci-test-triggered  new feature

Add SingleHighPtNoMuIDMuon line to HLT1
!815 · created 3 months ago by Ross John Hunter  RTA  ci-test-triggered  new feature  selections

Improved type-ful store
!923 · created 1 week ago by Daniel Hugo Campora Perez  RTA  ci-test-triggered

New calo decoding
!691 · created 8 months ago by Jean-Francois Marchand  Calo  RTA  ci-test-triggered

Make check for banks with 5 most-significant bits set more robust
!920 · created 1 week ago by Roel Aaij  RTA  ci-test-triggered

Install sequences under AllenSequences
!903 · created 2 weeks ago by Patrick Spradlin  RTA  ci-test-triggered

Update selections documentation to reflect new event model
!915 · created 2 weeks ago by Thomas Boettcher  RTA

Fix multiple raw event locations in TransposeRawBanks
!850 · created 2 months ago by Roel Aaij  RTA  bug fix  ci-test-triggered

Allow empty selreport in OutputHandler.
!919 · created 1 week ago by Kate Abigail Richardson  RTA

Remove host buffers used for monitoring.
!918 · created 1 week ago by Daniel Hugo Campora Perez  RTA  ci-test-triggered

Update Allen to follow tuning of TwoTrackMVA model
!917 · created 1 week ago by Vladimir Gligorov  RTA

HLT1 Lumi Line
!743 · created 5 months ago by Shu Xian  Luminosity  RTA  ci-test-triggered

Test only relevant throughput tests
!916 · created 2 weeks ago by Daniel Hugo Campora Perez  RTA  only GitLab CI

**Daniel Cámpora**

# Commissioning is a changing objective

- Decoding status: ecal, muons, VELO, SciFi, plume (from now on dev for HLT1 and HLT2 simultaneous)

- Reconstruction status: Forward tracking, VELO-SciFi matching

- Selections: 30 lines, scalable (~10% throughput hit scaling to 100 lines)

- Monitoring: progressing well

- Throughput / stability tests: Up to 30 MHz without any dropped packets (*full chain*, calo clustering or passthrough 1/25). One can push up to 40 MHz, *Allen is not the bottleneck*

- All GPUs are installed, operational and tested (also all FPGAs, EB nodes, etc.)

- In general, more monitoring is needed to identify issues more quickly

- Advances every week: **excellent work atmosphere, good pace, responsive team**

**Daniel Cámpora**

# Data taking



- LHCb has been exercising its DAQ in parallel to the LHC commissioning

- Sub-set of detectors (Calorimeters, Muon stations, PLUME) already in the global partition of the Experiment Control System (ECS)

- System running 24/7 in parallel to sub-detector commissioning activities

**Daniel Cámpora**

# Data taking



- ~200 GPUs are installed in the EB
- HLT1 is already included in the global partition

46

**Daniel Cámpora**

# Data taking



- ~200 GPUs are installed in the EB
- HLT1 is already included in the global partition
- Here shown: Triggering on calorimeter clusters @ 20 MHz

**Daniel Cámpora**

# Greedy because we can

- For the moment, we are running either of

  `hlt1_pp_matching_no_gec`

  `hlt1_pp_no_gec_no_ut`

```
NVIDIA GeForce RTX 3090    ███████████████████████████    159.13 kHz (1.00x)
NVIDIA RTX A5000           ██████████████████████         124.70 kHz (1.00x)
NVIDIA GeForce RTX 2080 Ti ██████████████                 78.61 kHz (1.00x)
AMD EPYC 7502 32-Core      █                              8.62 kHz (1.01x)
                           +--+--+--+--+--+--+--+--+--+
                           0  20 40 60 80 100 120 140 160
```
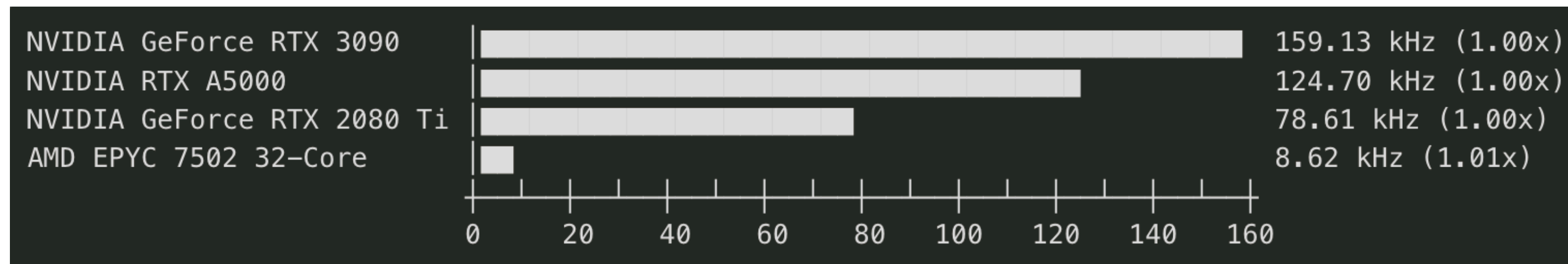
- *No UT* means drop in physics efficiency, but not in throughput.

- *No GEC* costs about 20% in throughput.

- A5000 is running at 125 kHz for `hlt1_pp_no_gec_no_ut`, so one expects drops if we run at full rate and luminosity

**Daniel Cámpora**

# So are we good?

- Rate of current runs is about 20 MHz, held back by various issues

  - Recent retina update, now it's RICH...

- In terms of pileup, this year we will run at mu=1. Next year, mu=7.

  - Nu, mu and pileup definitions: https://twiki.cern.ch/twiki/bin/view/LHCb/NuMuPileUp

  - Even when we run at 30 MHz, Allen will handle it this year.

- For the upcoming period we will need to do something to keep using the current sequences and scale up.

**Daniel Cámpora**

# Milestones - First trigger on real data

- Test the High Level Trigger 1 (HLT1):

  - Copy data to / from GPUs, perform decoding and reconstruction algorithms, select based on physics cuts

  - For now only based on calorimeter data: Decode 3x3 ecal clusters, trigger on > 400 MeV

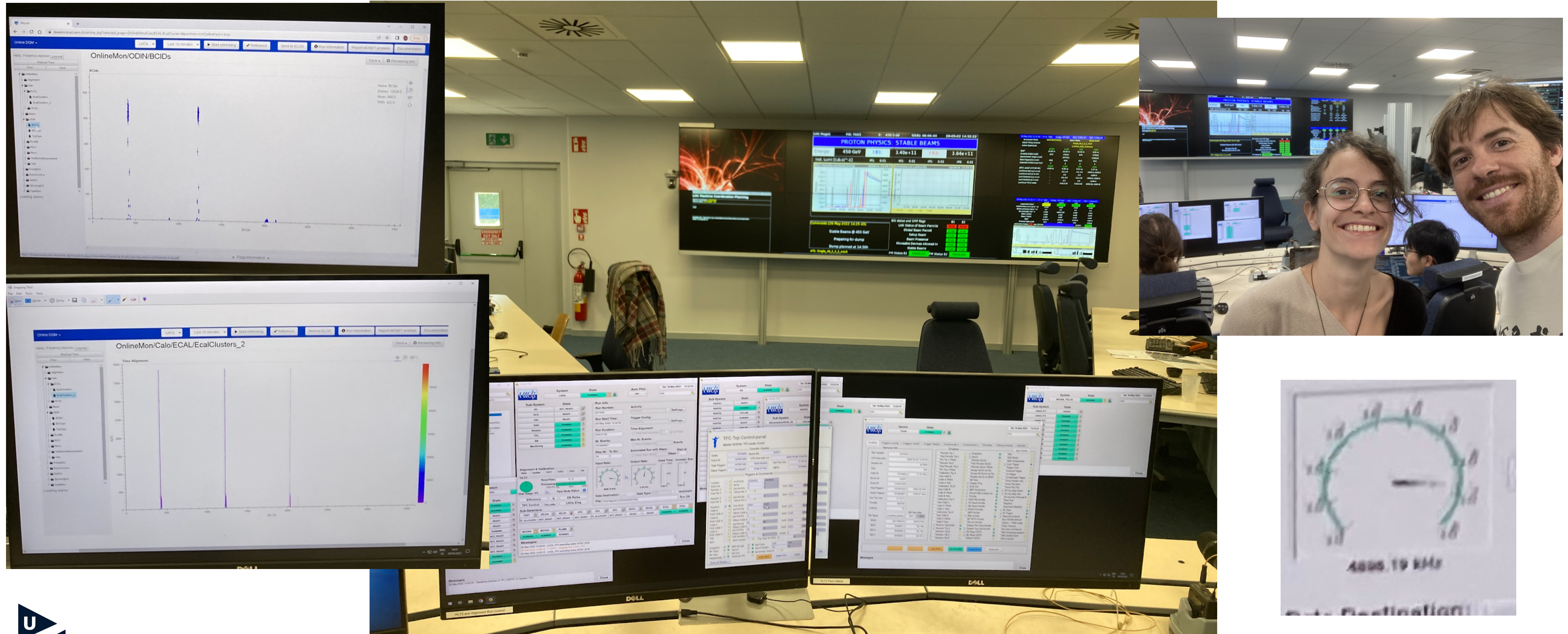    - What actually happened:

Saturday 28/05:
- Since 9:12: `CaloActivity`, triggering on ADC, no decoding fix
- 12:55: Calo experts suggest to switch to cluster trigger line, 400 to 2000 MeV
- **14:26:** Decoding fix implemented, first proper trigger from this point on. **Beam dump at 14:30.**

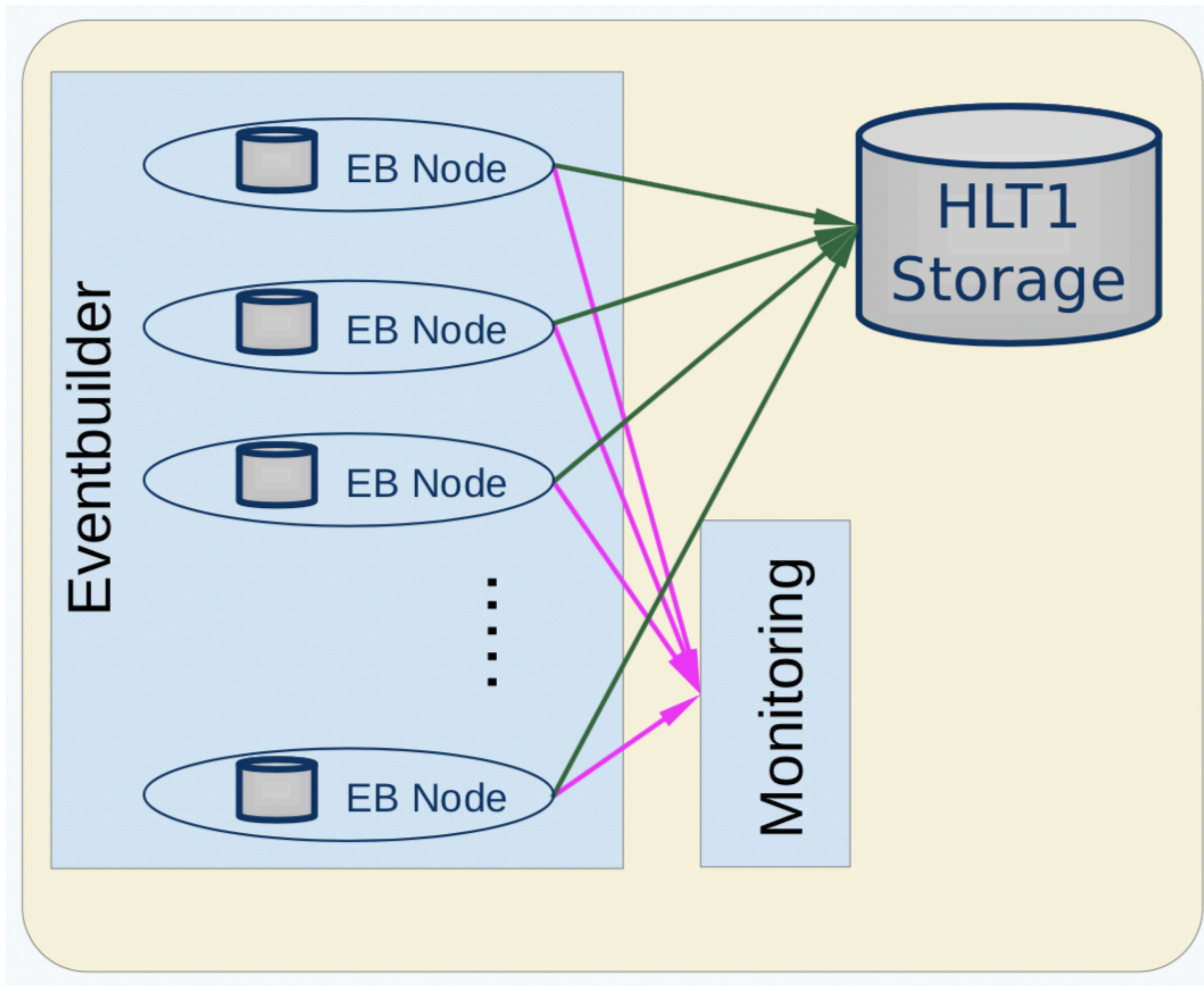Sunday 29/05:
- HLT1 ran all day with no major hiccups

**Daniel Cámpora**

# Milestones - First trigger on real data

- First trigger on real-data in LHCb Run 3 (28/05/22)

**Daniel Cámpora**

# Milestones - DAQ

- Large-scale Data Acquisition test



**Data-flow:**

- Up to 40 Tbps of data from the subdetectors if received by FPGA cards (PCIe40) hosted in the Event Building (EB) servers

- The EB network protocol bring the subdetector data fragments produced in a single event together and then group these events together into Multi-Event Packets (MEPs)

- Allen receives MEPs from the EB and processes them, producing output files in Markus' Data-Format (MDF)

**Daniel Cámpora**
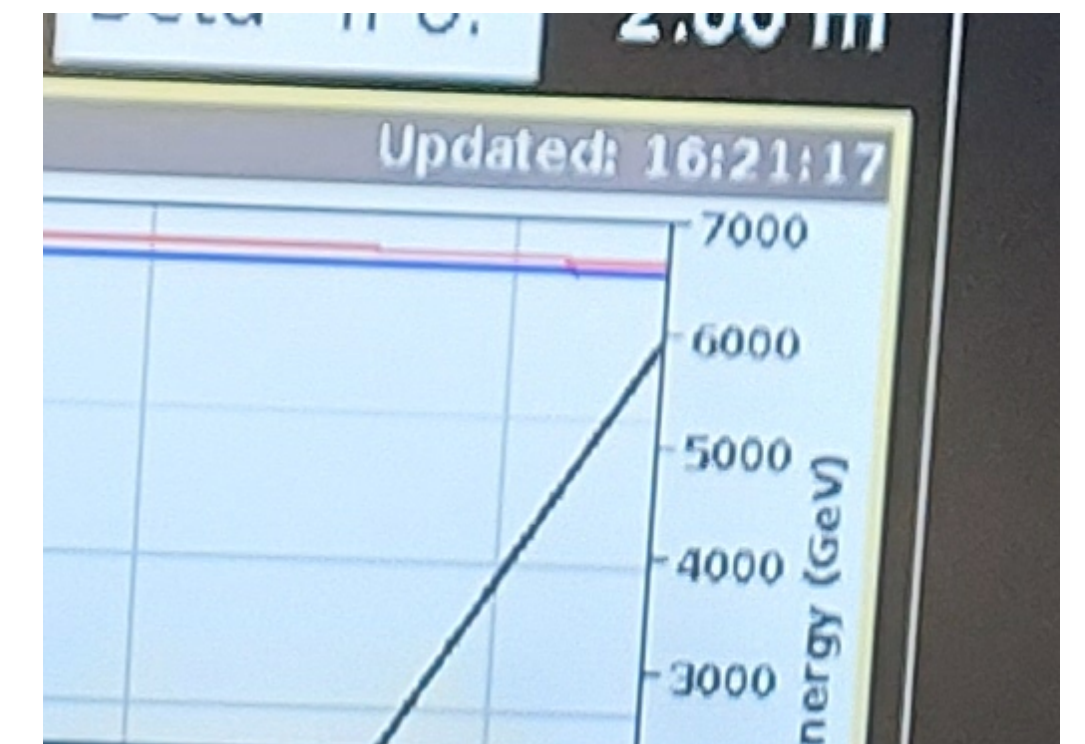
# Milestones - DAQ



2/06/22, 10 PM: First successful test at 10 MHz

- Take all TELL40s from all subsystems and try running the full DAQ system and push the limits of the TELL40, EB and storage

- Tested 1-2-5-10-15 MHz with no major issues encountered

- Latest tests pushing the rates to 20-30 MHz achieved in the next days after solving an issue

**Daniel Cámpora**

# Milestones - Run at design energy

- Highest energy in a particle physics accelerator (13.6 TeV, 05/07/22)

**Daniel Cámpora**

# Conclusions

- The R&D of many years (at least since 2012!) has finally been realized

  - LHCb has a new detector, a new DAQ and a new trigger

  - The specification has changed quite a lot, only possible with a dedicated collaboration and lots of hard work (it's also fun)

- Software has changed along the way, very different from 10 years ago

  - New architectures available, exciting time to develop new software

- Commissioning is a rocky road, but extremely rewarding

  - We are preparing the physics of tomorrow

  - You can make a great impact!

**Daniel Cámpora**

# Future work

Copa Mundial de la FIFA 2022™ · Hoy, 16:00

España      contra      Costa Rica

Fase de grupos · Grupo E

**Daniel Cámpora**

# Thank you for your attention!

Daniel Cámpora