# Objectives of this talk

1. Motivate why it is interesting to perform variable latency charged particle reconstruction (tracking) in real-time, which at the LHC means at 30 MHz

2. Describe the challenges involved with delivering such a tracking for the LHCb experiment with reference to two specific architectures: x86 and GPU

3. Give my personal thoughts on what we have learned during this development process in LHCb, and thoughts about where this is going in the future

**Forward *spectrometer* optimized for precision physics**

# Why tracking @30 MHz?
# Why variable latency?

# Q : What is real-time?

DETECTOR → DATA IN → REAL TIME PROCESSING

SAVED DATA → STORAGE

DISCARDED DATA → DEV/NULL

**A : Any processing of data before it is permanently recorded**

# Why do we need to process data before recording it?

LHCb

CMS/ATLAS

Data volume
at detector
in Run 2

~30 Eb/
year

~1 Zb/
year

# Why do we need to process data before recording it?

LHCb

CMS/ATLAS

Data volume
at detector
in Run 2

~30 Eb/
year

Global internet
dataflow 2015

~1 Zb/
year

~640 Eb/
year

**Because HEP detectors produce too much data to store**

# Data volumes @ LHC after real-time processing

**LHCb**

**CMS/ATLAS**

Data volume
at detector
in Run 2

~30 Eb/
year

Global internet
dataflow 2015

~1 Zb/
year

~640 Eb/
year

Data volume
for analysts

~30 Pb/
year

~40 Pb/
year

**Real-time processing reduces data by 3-5 orders of magnitude**

Fixed latency

Data compression ← → Event selection

Variable latency

Fixed latency

Data compression ←——————→ Event selection

ATLAS/CMS/LHCb
first level calo &
muon triggers

Variable latency

**Distinguish fixed & variable latency, selection & compression**

Fixed latency

ATLAS/CMS/LHCb
first level calo &
muon triggers

Data compression → Event selection

ATLAS/CMS/LHCb
High Level Triggers

Variable latency

Fixed latency

ALICE upgrade
TPC processing

ATLAS/CMS/LHCb
first level calo &
muon triggers

Data compression → Event selection

ATLAS/CMS/LHCb
High Level Triggers

Variable latency

**Distinguish fixed & variable latency, selection & compression**

**Fixed latency**

ALICE upgrade
TPC processing

ATLAS/CMS/LHCb
first level calo &
muon triggers

Data compression ──────────► Event selection

ATLAS "trigger level analysis"
CMS   "data scouting"
LHCb  "real-time analysis"

ATLAS/CMS/LHCb
High Level Triggers

**Variable latency**

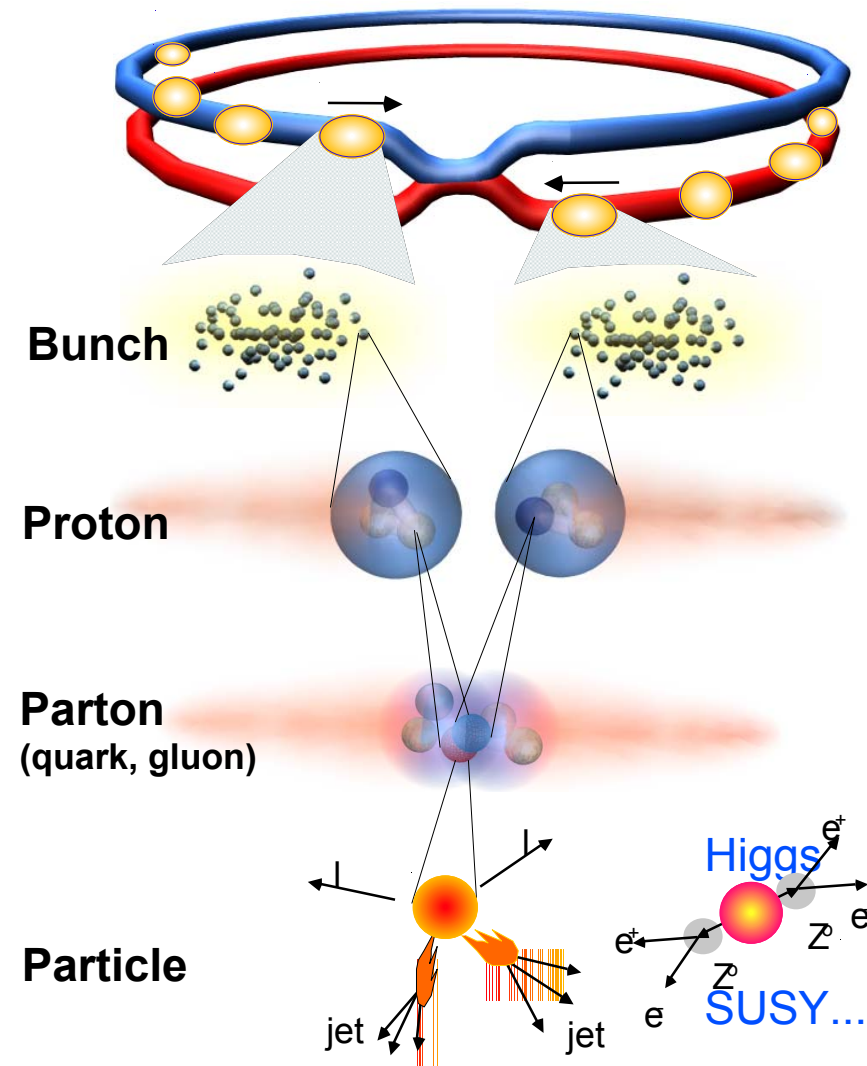**Distinguish fixed & variable latency, selection & compression**

# Traditional real-time processing, or "triggering"



CMS pp √s=7 TeV
∫ L dt = 4.98 fb⁻¹
Threshold : 15 GeV

L1 E/Gamma Trigger
Electrons from Z

L1_SingleEG15
● Barrel
□ Endcaps

Efficiency vs $E_T$ [GeV]

## Collisions at the LHC: summary



| | |
|---|---|
| Proton - Proton | 2804 bunch/beam |
| Protons/bunch | $10^{11}$ |
| Beam energy | 7 TeV ($7 \times 10^{12}$ eV) |
| Luminosity | $10^{34} cm^{-2} s^{-1}$ |
| Crossing rate | 40 MHz |
| Collision rate ≈ | $10^7$-$10^9$ |
| New physics rate ≈ .00001 Hz | |

Bunch

Proton

Parton
(quark, gluon)

Particle

Higgs

SUSY.....

jet    jet

**Event selection:
1 in 10,000,000,000,000**

P. Sphicas
Triggering

SSI 2006
July 2006

3

## Driven by fixed-latency selection, analysis on efficiency plateau

proton - (anti)proton cross sections

proton - (anti)proton cross sections

The plot which basically motivated the LHCb upgrade

**Fixed-latency CALO trigger only effective up to 4·10³²**

**We will have MHz of signals in our acceptance! Can only reject up to 1/60 efficiently with inclusive selections. Require real-time analysis beyond this.**

# From selection to compression : real-time analysis



Most physics measurements require only a signal candidate and information about the specific pp collision which produced it → the rest is pileup

The higher the luminosity, the larger the fraction of event data caused by pileup

Hence create more room for signal by compressing & removing pileup in real-time!

## So how do we carry out precise pileup suppression?

# We also need to align and calibrate our detector in real time

# So we did!



LHCb VELO
Preliminary

- x-translation
- y-translation

Empty markers = no update          23/04/2016 - 06/06/2016

LHCb Tracker
Preliminary

- IT1 ASide
- IT1 CSide
- IT1 Top
- IT1 Bottom

Empty markers = no update          06/05/2016 - 06/06/2016

L0 rates include readout dead time.

Few % control of calorimeter
response changes due to ageing!

20

# We also need to measure our efficiencies in real-time!

| Species | Low momentum | High momentum |
|---|---|---|
| $e^{\pm}$ | | $B^+ \to J/\psi K^+$ with $J/\psi \to e^+e^-$ |
| $\mu^{\pm}$ | $B^+ \to J/\psi K^+$ with $J/\psi \to \mu^+\mu^-$ | $J/\psi \to \mu^+\mu^-$ |
| $\pi^{\pm}$ | $K_S^0 \to \pi^+\pi^-$ | $D^{*+} \to D^0\pi^+$ with $D^0 \to K^-\pi^+$ |
| $K^{\pm}$ | $D_s^+ \to \phi\pi^+$ with $\phi \to K^+K^-$ | $D^{*+} \to D^0\pi^+$ with $D^0 \to K^-\pi^+$ |
| $p$ , $\bar{p}$ | $\Lambda^0 \to p\pi^-$ | $\Lambda^0 \to p\pi^-$ ; $\Lambda_c^+ \to pK^-\pi^+$ |



Unlike ATLAS and CMS, LHCb must maintain a data-driven permille level control of its efficiency across the kinematic and geometric acceptance of the detector. Requires collecting an extremely wide range of tag-and-probe samples in real time.

21

**Full flexibility to store "additional" detector information if required by some analyses**

# LHCb upgrade dataflow



40 Tb/s
30 MHz non-empty pp

1-2 Tb/s
0.5 - 1.5 MHz

80 Gb/s

LHC bunch crossing (30MHz)

DETECTOR READOUT

PARTIAL RECONSTRUCTION (HLT1)

Buffer

REAL-TIME ALIGNMENT & CALIBRATION

FULL RECONSTRUCTION (HLT2)

26% FULL

OFFLINE PROCESSING

68% TURBO & real-time analysis

6% CALIB

OFFLINE PROCESSING

Physics analysis

**RECONSTRUCTION STEP**    **OUTPUT OBJECTS**

EXECUTION ORDER

VELO tracking with simplified Kalman Filter → VELO tracks

Primary Vertex finding → Primary Vertices (PV)

VELO → UT tracking inital momentum estimate → upstream tracks

UT → T stations tracking with $p_T > 500$ MeV/c → long tracks

Kalman Filtering

Fake track rejection → fitted long tracks

**"Traditional" inclusive selections selecting bunch crossings. Must be based on tracks, so require 30 MHz tracking at $2 \cdot 10^{33}$!**

Because LHCb is a dipole spectrometer, tracking inherently requires non-local data from multiple subdetectors to be brought together.

You can build a fixed-latency track trigger but you will have to build the biggest part of the detector readout for it anyway — might as well just read everything out upfront and work in variable latency.

This is *not* an argument about e.g. not using FPGAs, just you first build events, then process them in whatever way is most cost-effective.

27

| CMS detector | LHC Run-2 | HL-LHC Phase-2 | |
|---|---|---|---|
| Peak $\langle PU \rangle$ | 60 | 140 | 200 |
| L1 accept rate (maximum) | 100 kHz | 500 kHz | 750 kHz |
| Event Size | 2.0 MB [a] | 5.7 MB [b] | 7.4 MB |
| Event Network throughput | 1.6 Tb/s | 23 Tb/s | 44 Tb/s |
| Event Network buffer (60 seconds) | 12 TB | 171 TB | 333 TB |
| HLT accept rate | 1 kHz | 5 kHz | 7.5 kHz |
| HLT computing power [c] | 0.5 MHS06 | 4.5 MHS06 | 9.2 MHS06 |
| Storage throughput | 2.5 GB/s | 31 GB/s | 61 GB/s |
| Storage capacity needed (1 day) | 0.2 PB | 2.7 PB | 5.3 PB |

**LHCb 2021 real-time tracking has to handle the *same data volume* as ATLAS/CMS HL-LHC upgrades! But earlier and for less money...** 28

# Challenges and solutions

RECONSTRUCTION STEP | OUTPUT OBJECTS

EXECUTION ORDER

VELO tracking with simplified Kalman Filter → VELO tracks

Primary Vertex finding → Primary Vertices (PV)

VELO → UT tracking inital momentum estimate → upstream tracks

UT → T stations tracking with $p_T > 500$ MeV/c → long tracks

Kalman Filtering

Fake track rejection → fitted long tracks

# Let's look at this sequence in more detail

| Inputs | | RECONSTRUCTION STEP | | OUTPUT OBJECTS |
|---|---|---|---|---|
| VELO raw data | | VELO tracking with simplified Kalman Filter | | VELO tracks |
| VELO tracks | | Primary Vertex finding | | Primary Vertices (PV) |
| VELO tracks and UT raw data | | VELO → UT tracking inital momentum estimate | | upstream tracks |
| Velo-UT tracks and SciFi raw data | | UT → T stations tracking with $p_T > 500$ MeV/c | | long tracks |
| Long tracks | | Kalman Filtering | | fitted long tracks |
| Fitted long tracks | | Fake track rejection | | |

EXECUTION ORDER

**To run at 30 MHz we need to get data off the detector, transform it to the global coordinate system, and do pattern recognition**

Timing fraction within the HLT1 sequence [%]

- PrPixelTrackingFast — 26.52 %
- VPCLustering — 16.39 %
- PatPV3D — 13.63 %
- PrStoreUTHit — 11.91 %
- PrStoreFTHit — 9.82 %
- FetchDataFromFile — 6.54 %
- createFTClusters — 4.77 %
- createUTLiteClusters — 3.15 %
- PrVeloUTFast — 2.89 %
- ForwardTrackingFast — 2.38 %
- DummyEventTime — 1.51 %
- other — 0.51 %

PATTERN RECOGNITION
DATA PREPARATION

**Early 2018 after about 3 years of work to make the framework thread-safe. Data preparation as important as pattern reco!**

# How to improve it?

1. Do what you can on the readout boards! Output the data in the most useful format possible, perform clustering in the readout if you can.

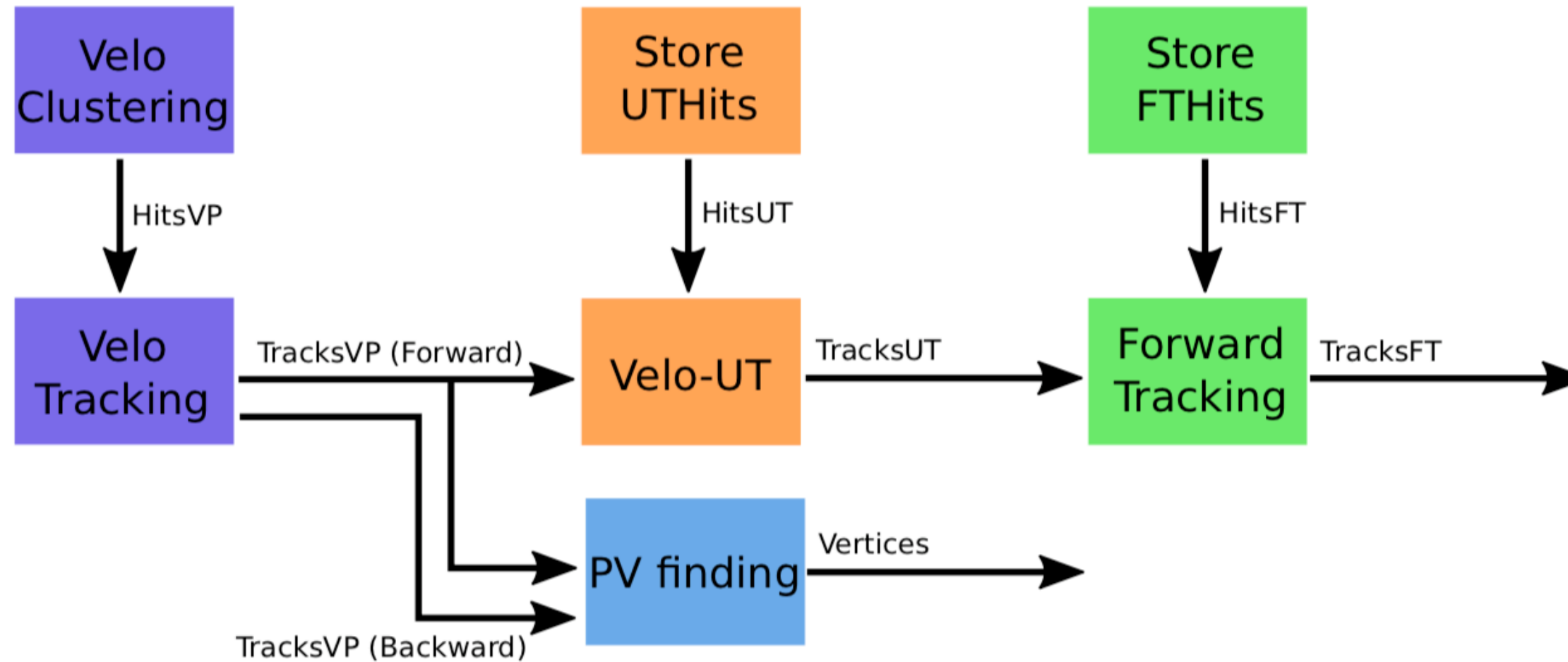2. Write custom throughput oriented data structures which only contain the absolute minimum needed by pattern recognition. "Plain old data".

3. Work with SOA structures wherever possible to enable vectorization.

4. Minimize copying of information by breaking up large structures, for example tracks, into smaller pieces — for example track parameters and indices pointing to the track hits in one place, tracks states in another, fit results in a third. Prefer to join these later when needed.

# So what does the new sequence then concretely look like?

cross section at y=0

390 mrad

70 mrad

15 mrad

1 m

66 mm

interaction region showing
$2\times\sigma_{beam}$ = ~12.6 cm

Upstream track

T1  T2  T3

TT

VELO

Long track

VELO track

Downstream track

T track

1 m

cross section at y=0

390 mrad

70 mrad

x

z

15 mrad

66 mm

interaction region showing
$2 \times \sigma_{beam} = \sim 12.6$ cm

y

x

φ

**VELO tracks from the beamline traverse lines of constant φ**
**When extrapolating, looking for N nearest neighbours in φ is more**
**effective than searching for hits in a search window in φ**

36

**Required HLT1 throughput achieved in 2019! Gains from plain and local data and vectorization add non-linearly.**

## Proposal for an HLT1 implementation on GPUs for the LHCb experiment

R. Aaij[1], J. Albrecht[2], M. Belous[a,3], T. Boettcher[4], A. Brea Rodríguez[5], D. vom Bruch[6], D. H. Cámpora Pérez[b,7], A. Casais Vidal[5], P. Fernandez Declara[c,7], L. Funke[2], V. V. Gligorov[6], B. Jashal[9], N. Kazeev[a,3], D. Martínez Santos[5], F. Pisani[d,e,7], D. Pliushchenko[f,3], S. Popov[a,3], M. Rangel[10], F. Reiss[6], C. Sánchez Mayordomo[9], R. Schwemmer[7], M. Sokoloff[11], A. Ustyuzhanin[a,3], X. Vilasís-Cardona[8], M. Williams[4]

**Exploits flexibility of our Run 3 DAQ by implementing HLT1 directly in the servers receiving the data from the detector. Judged viable by external review, full cost-benefit analysis ongoing to decide if we will use this in Run 3.**

**Run 3: Baseline**

pp Collisions

5 TB/s

EB

5 TB/s

x86 CPU farm

HLT1

buffer on disk
calibration and alignment

HLT2

10 GB/s

Storage

**Run 3: GPU-enhanced**

pp Collisions

5 TB/s

EB on FPGAs
HLT1 on GPUs
(FPGAs & GPUs hosted
on same EB servers)

0.2 TB/s

x86 CPU farm

buffer on disk
calibration and alignment

HLT2

10 GB/s

Storage

32 Tb/s

10 Tb/s

100G IB

10GbE

25GbE

PCIe

~250 event building units

subfarm of 40 servers    subfarm of 40 servers    subfarm of 40 servers    subfarm of 40 servers

Up to 100 subfarms (up to 4000 event filtering servers)

**Exploit empty slots on the event building servers — opportunistic but efficient**
**Each GPU eats 6 GB/s — first integration tests look fine for I/O, final ongoing**

```
                    ┌──────────────┐
                    │   Raw data   │
                    └──────┬───────┘
                           │
                   ◇───────────────◇
                   │     Global     │
                   │   Event Cut    │
                   ◇───────────────◇
                           │
        ┌──────────────┐      ┌──────────────┐      ┌──────────────┐
        │ Velo decoding│      │  UT decoding │      │ Muon decoding│
        │ and clustering│     └──────┬───────┘      └──────┬───────┘
        └──────┬───────┘             │                     │
               │              ┌──────────────┐      ┌──────────────┐
        ┌──────────────┐      │  UT tracking │      │   Muon ID    │
        │ Velo tracking│      └──────┬───────┘      └──────┬───────┘
        └──────┬───────┘             │                     │
               │              ┌──────────────┐      ┌──────────────┐
        ┌──────────────┐      │ SciFi decoding│     │  Find sec-   │
        │ Simple Kalman│      └──────┬───────┘      │ ondary vertices│
        │    filter    │             │              └──────┬───────┘
        └──────┬───────┘      ┌──────────────┐             │
               │              │ SciFi tracking│     ◇───────────────◇
        ┌──────────────┐      └──────┬───────┘      │ Select events  │
        │ Find primary │             │              ◇───────────────◇
        │   vertices   │      ┌──────────────┐             │
        └──────────────┘      │ Parameterized│      ┌──────────────┐
                              │ Kalman filter│      │Selected events│
                              └──────────────┘      └──────────────┘
```
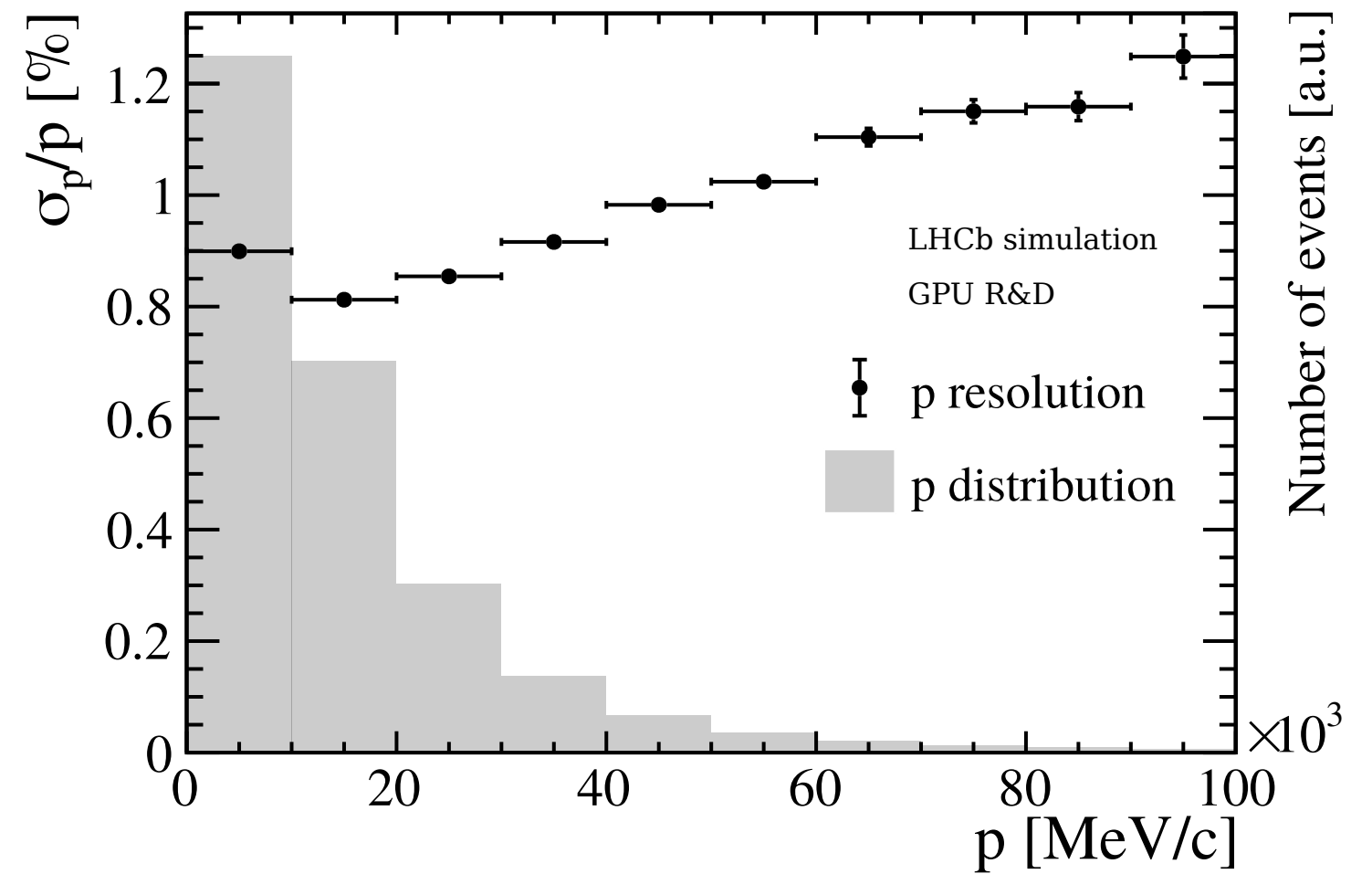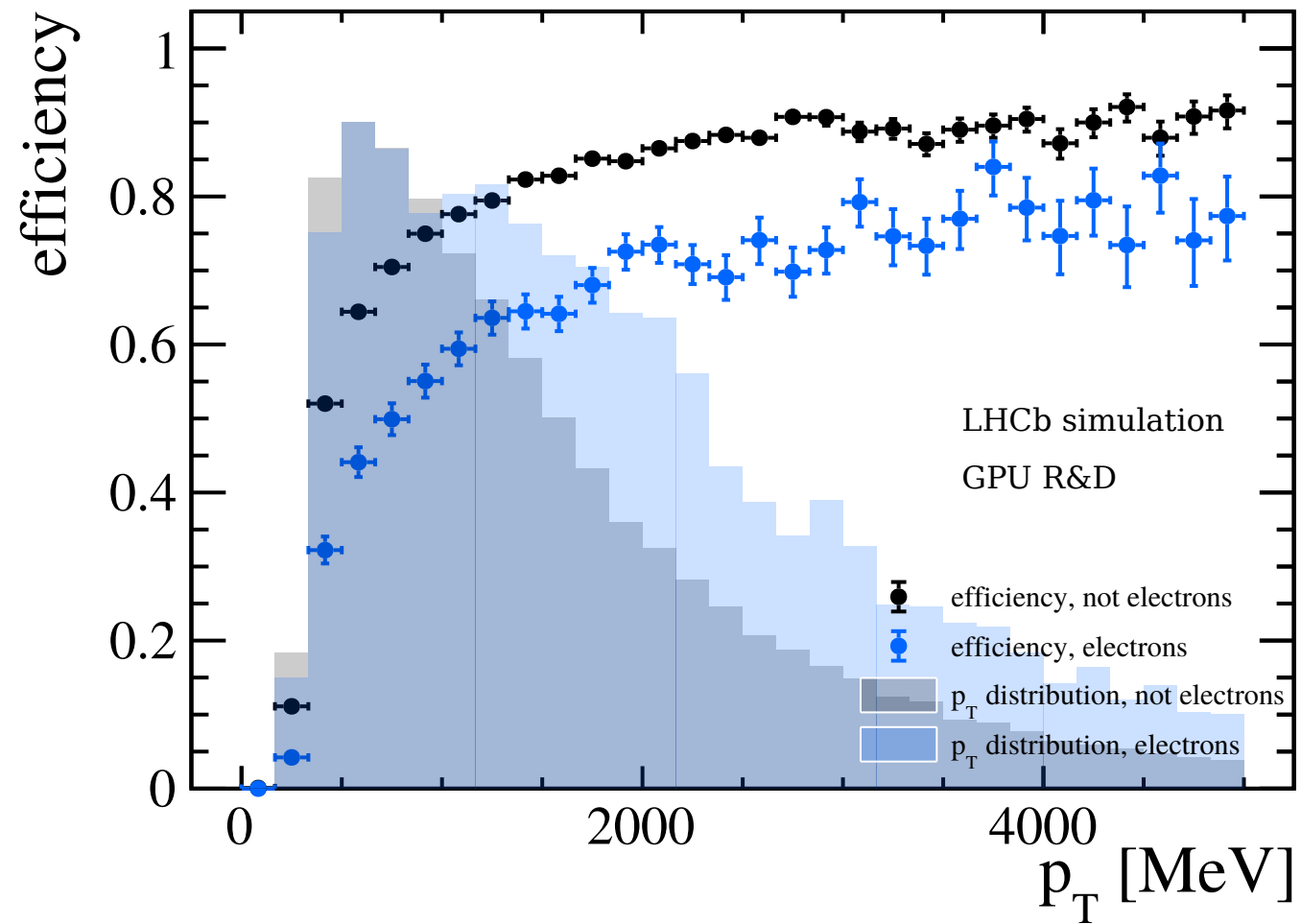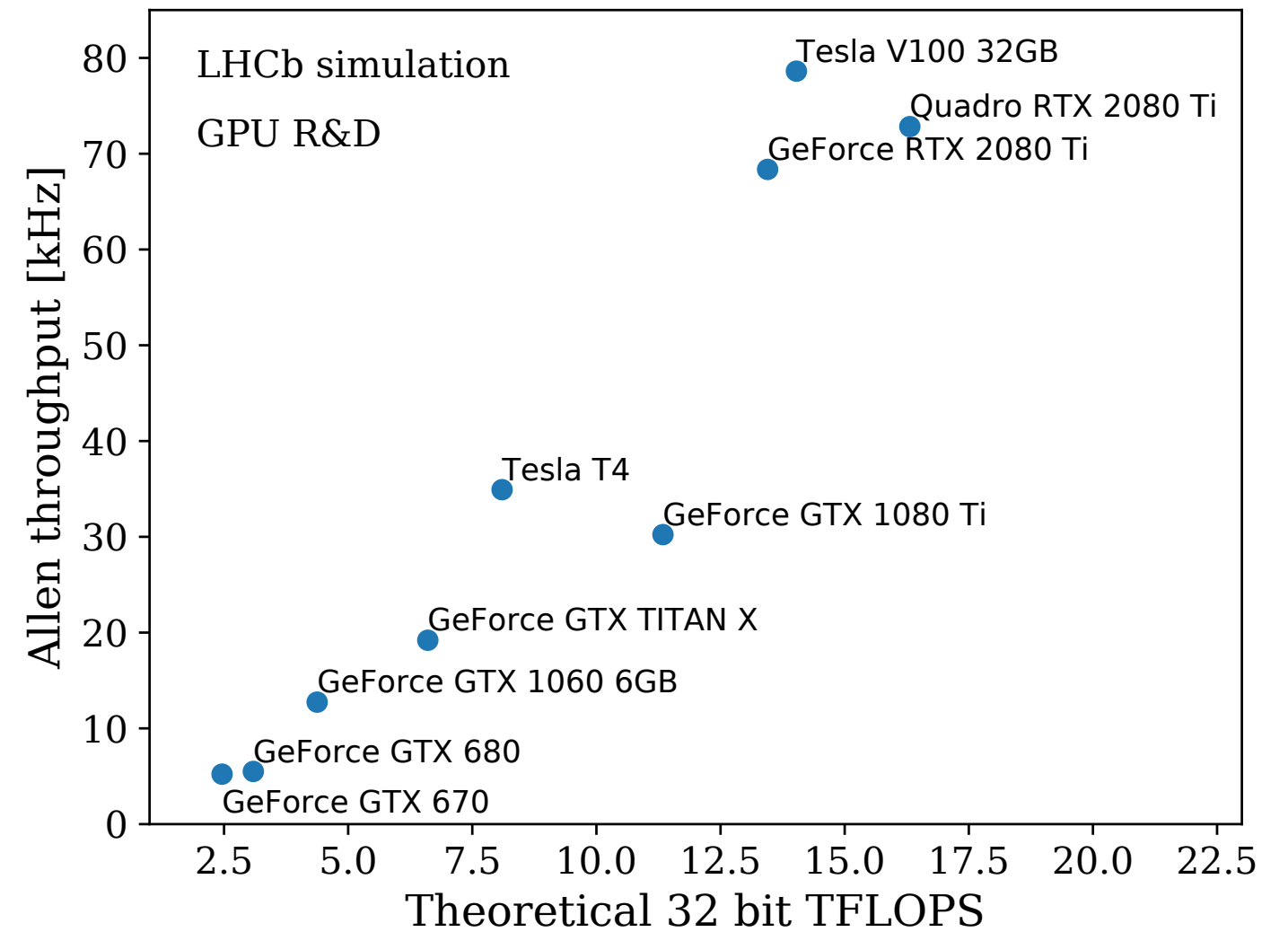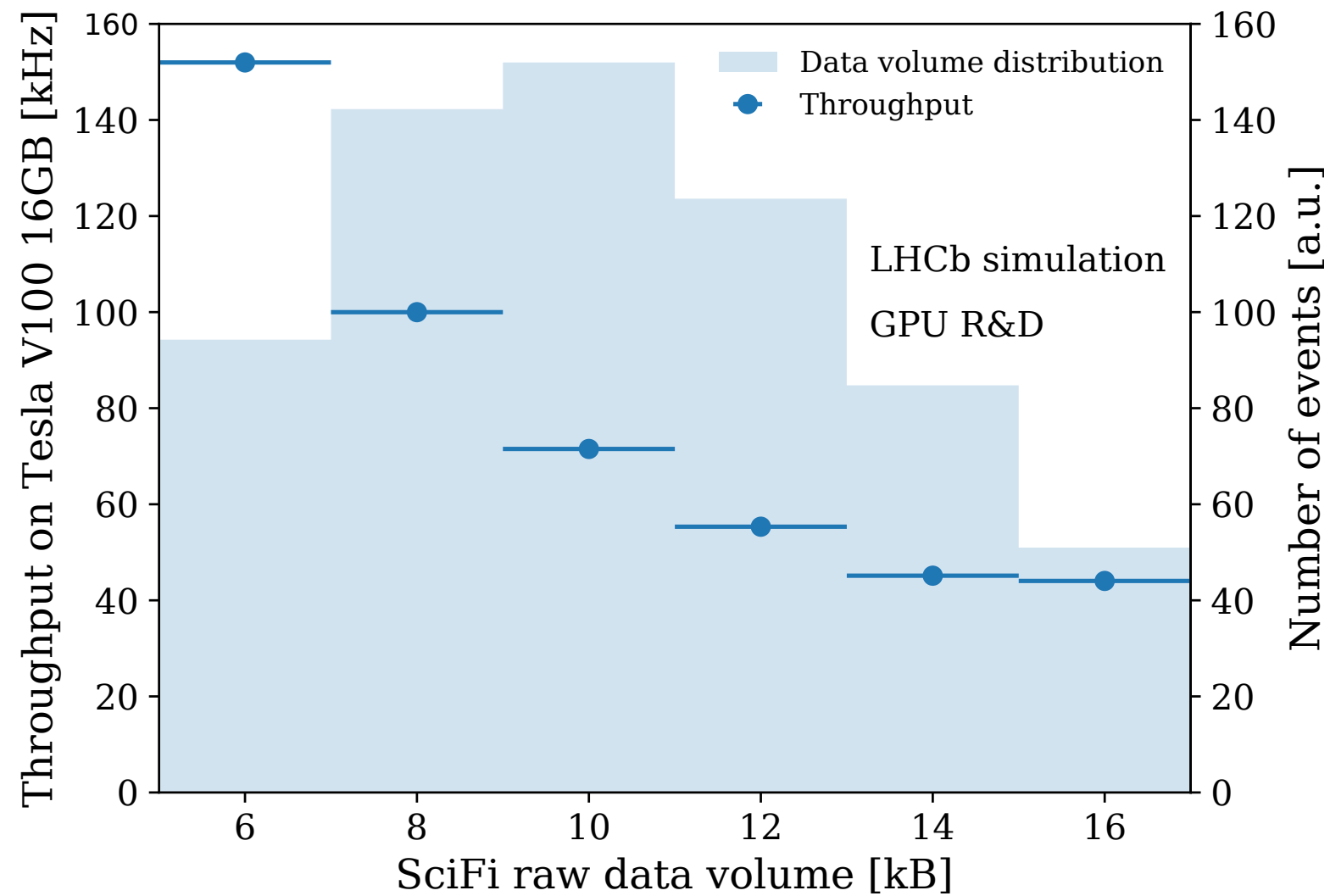
**Are really the same as multithreaded x86. Optimal degree of paralellism/branching is different, but plain local data is key!**

40

# Physics performance



As good as x86 baseline (no approved plots for the baseline...)
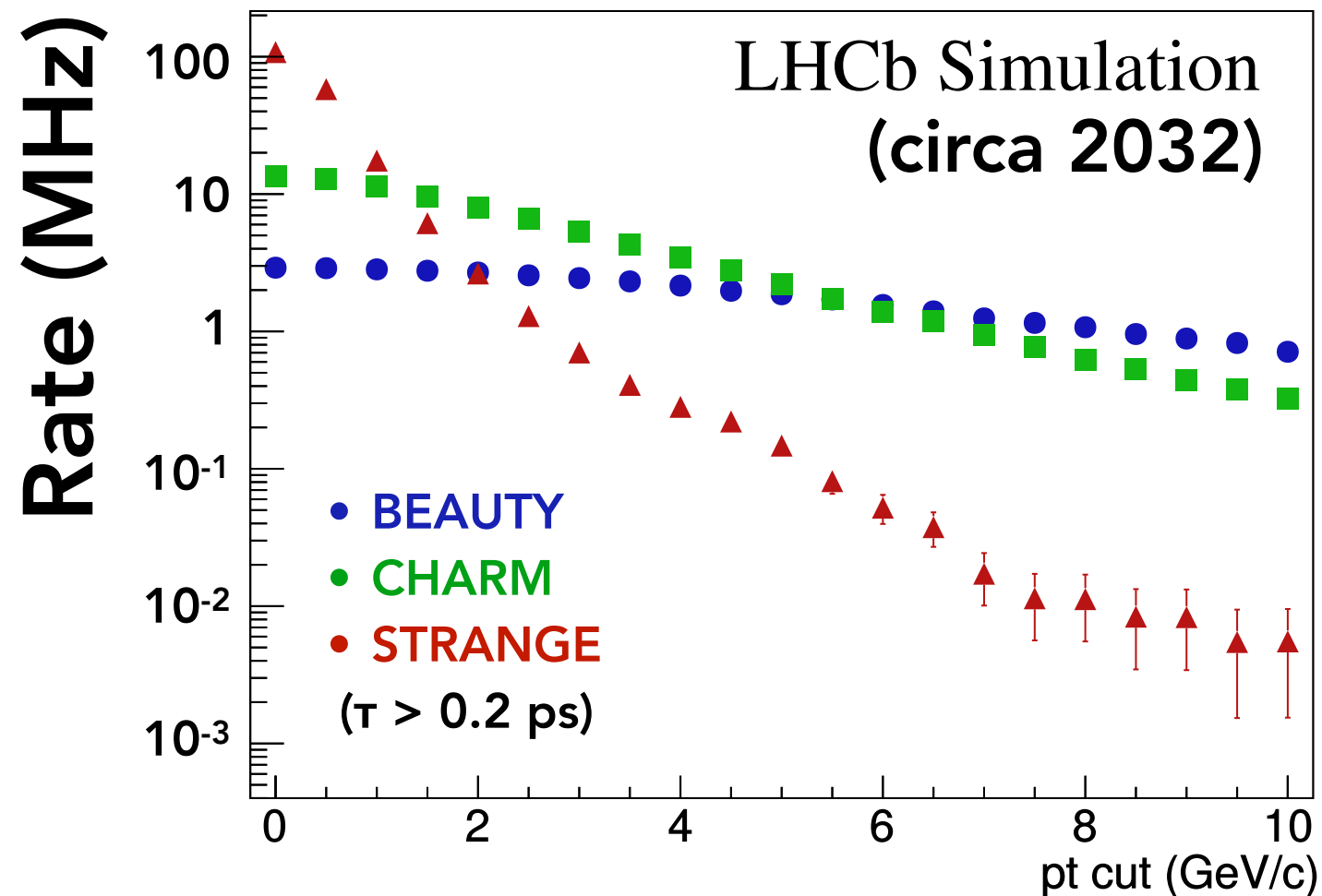
# GPU throughput scaling



**Linear scaling of throughput vs. occupancy, and throughput vs. the theoretical TFLOPS of each card. Optimal use of hardware!**
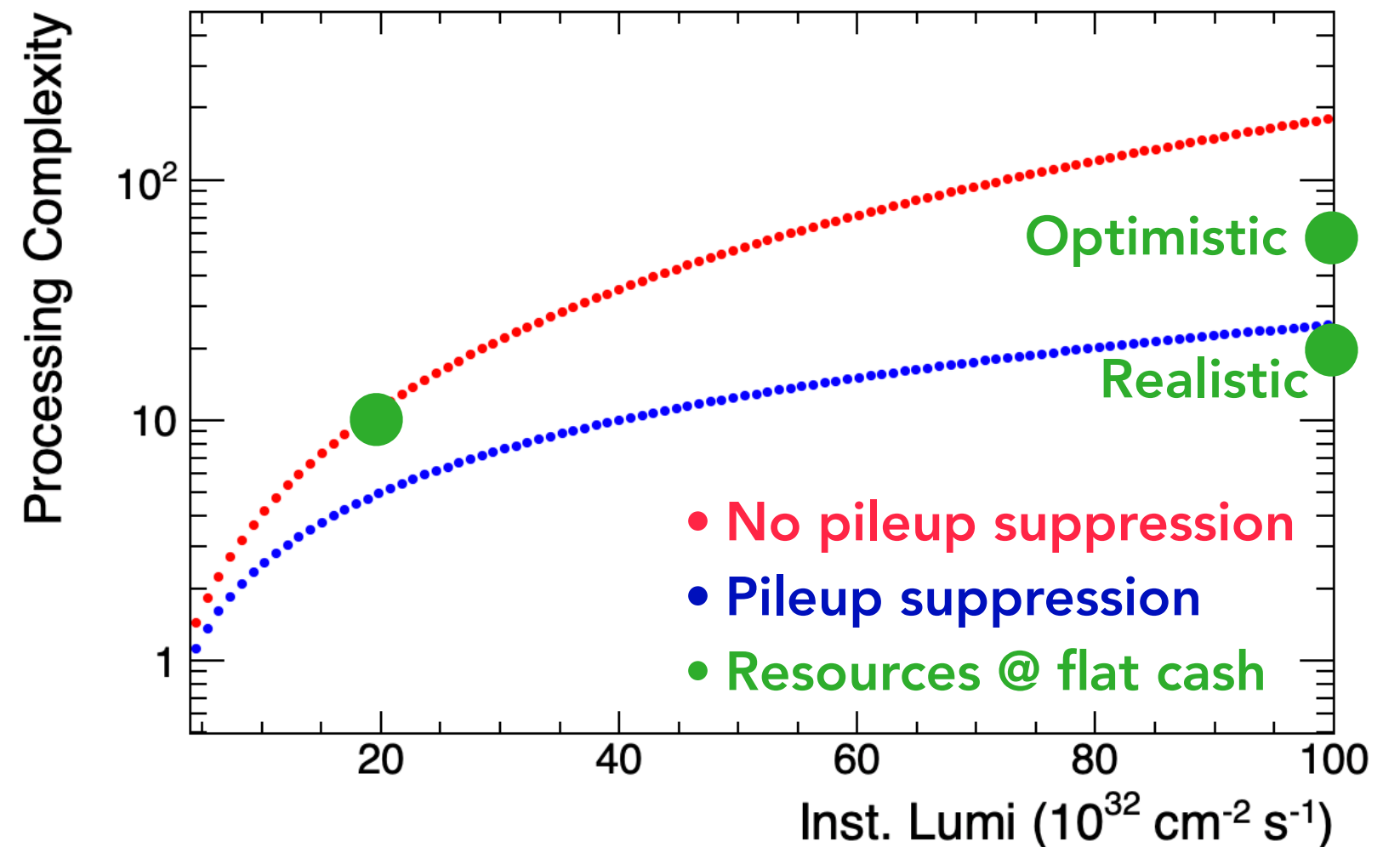
# Looking towards the future

# Looking beyond to a potential second LHCb upgrade

## Partially reconstructed signals



LHCb Simulation
(circa 2032)

Rate (MHz)

- BEAUTY
- CHARM
- STRANGE

(τ > 0.2 ps)

pt cut (GeV/c)

Fine print: this plot assumes that processing complexity goes linearly with detector occupancy, which is in itself an optimistic assumption before we even get to the pileup suppression part!



Processing Complexity

Optimistic

Realistic

- No pileup suppression
- Pileup suppression
- Resources @ flat cash

Inst. Lumi ($10^{32}$ cm$^{-2}$ s$^{-1}$)

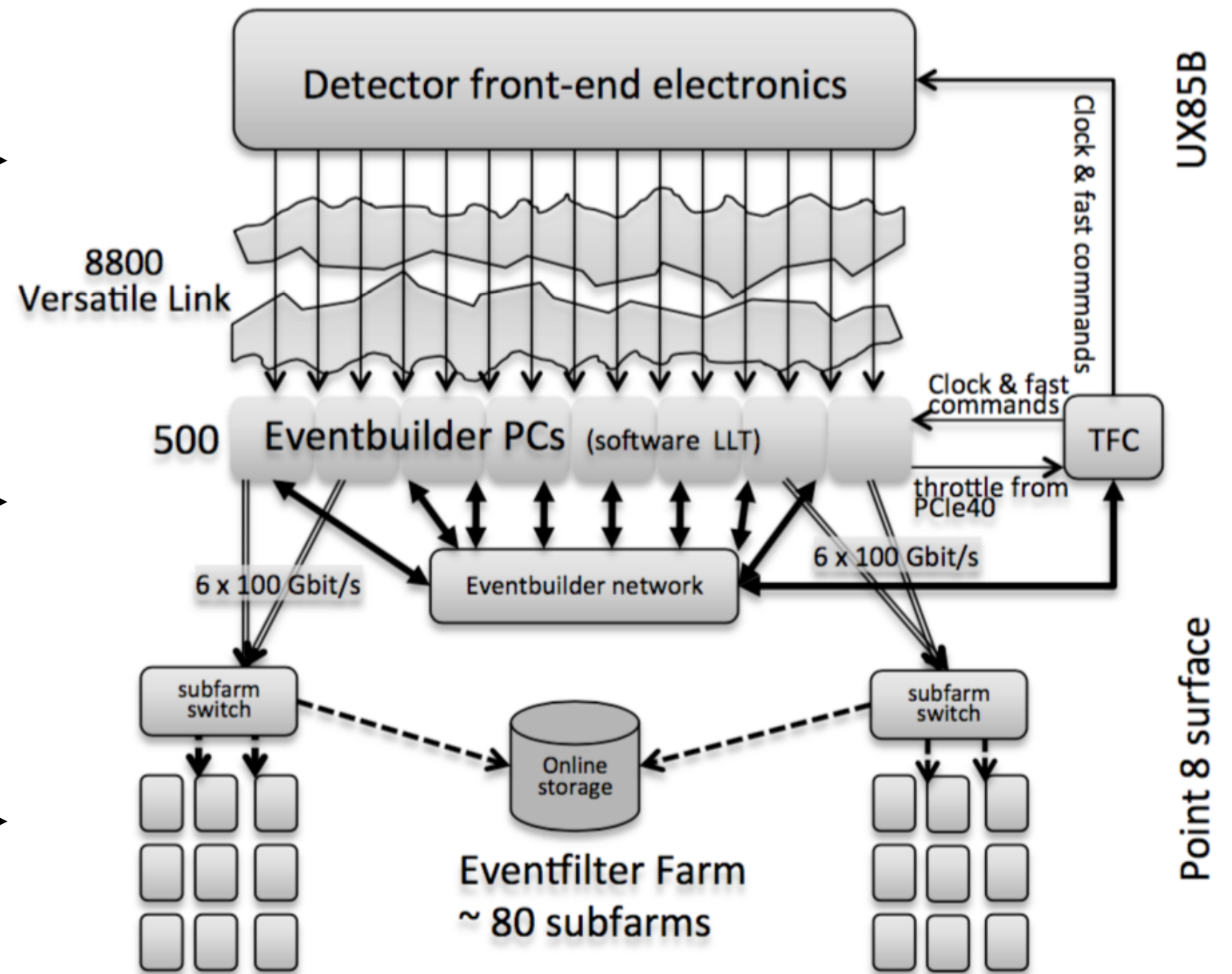**How to suppress pileup with O(60) pp collisions per bunch crossing?**

# Maintaining the flexibility of our processing will be crucial

GBT link : 4.8 Gb/s Upgrade I
Assume evolution to 10 Gb/s for HL-LHC
using aggressive error handling : missing
factor 5 compared to data rate growth.

Event-building : current network is 500
servers with 100 Gb/s links. 200 Gb/s
readily available, keep an eye on price/
performance scaling beyond this?

Farm : carry out R&D in next years on
optimal use of hybrid architectures (GPU/
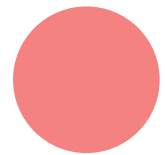CPU/FPGA), remain flexible



**We now have two viable HLT1 models, on x86 and on GPU, already for Run 3! Ability to exploit hybrid architectures crucial to maximize physics/Euro in the long term.**

45

# Personal observations on working in a hybrid world

1. The computing landscape is moving towards hybrid architectures. We are developing the skills to move with it!

2. If the basic principles of high throughput software are respected, a well designed software architecture will perform on x86, GPU, or FPGA systems. Functional design and uniform API helps to achieve this.

3. High-throughput software is far from what universities teach physics students no matter the architecture. Learning CUDA, HLS or C++17 is the same for them. Recognise the importance of new skills in the field.

4. A variable latency trigger is a home for API designers, physicists and selection authors, throughput experts, algorithm designers… it's a very diverse community and personal architecture preferences are real. It is more work to keep a diverse community coherent, but it's worth it.
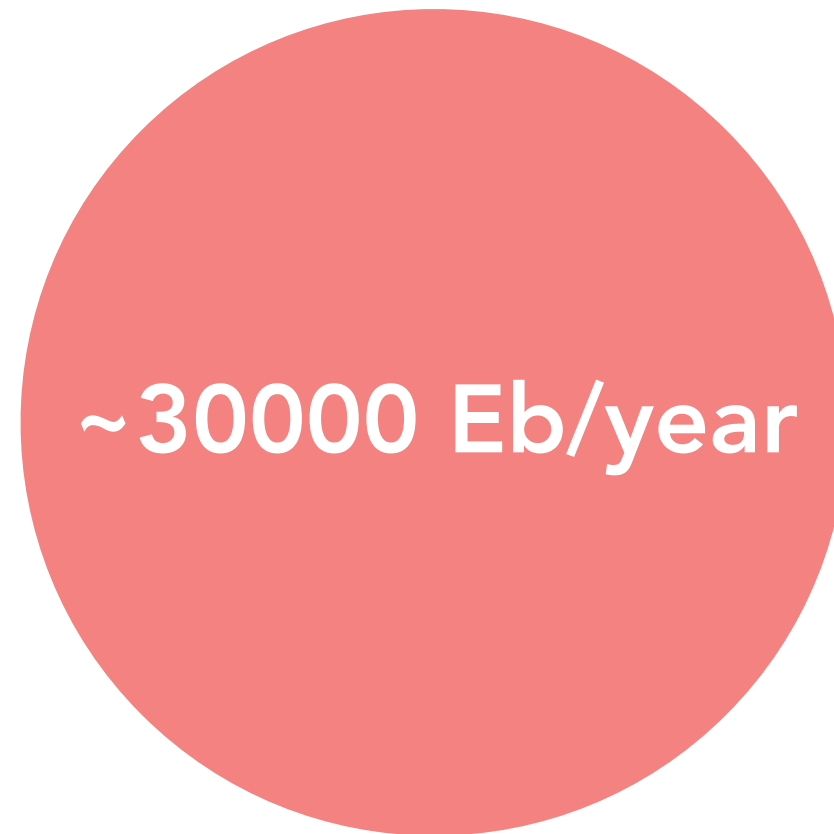
# Conclusions and final thoughts

**LHCb 2032**

**>1000 Eb/year**

**ATLAS+CMS 2027**

**260 Eb/year**

Square Kilometre Array (2030s)

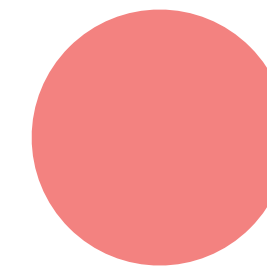**~30000 Eb/year**

Sequence genome of all humans on Earth

**8000 Eb**

Global internet dataflow 2021

**2800 Eb/year**

# Backup

# LHCb analysis methodology and role of calibration samples

## Trigger Efficiency
Tag-and-probe calibration method exists & widely used

## Tracking efficiency
Tag-and-probe
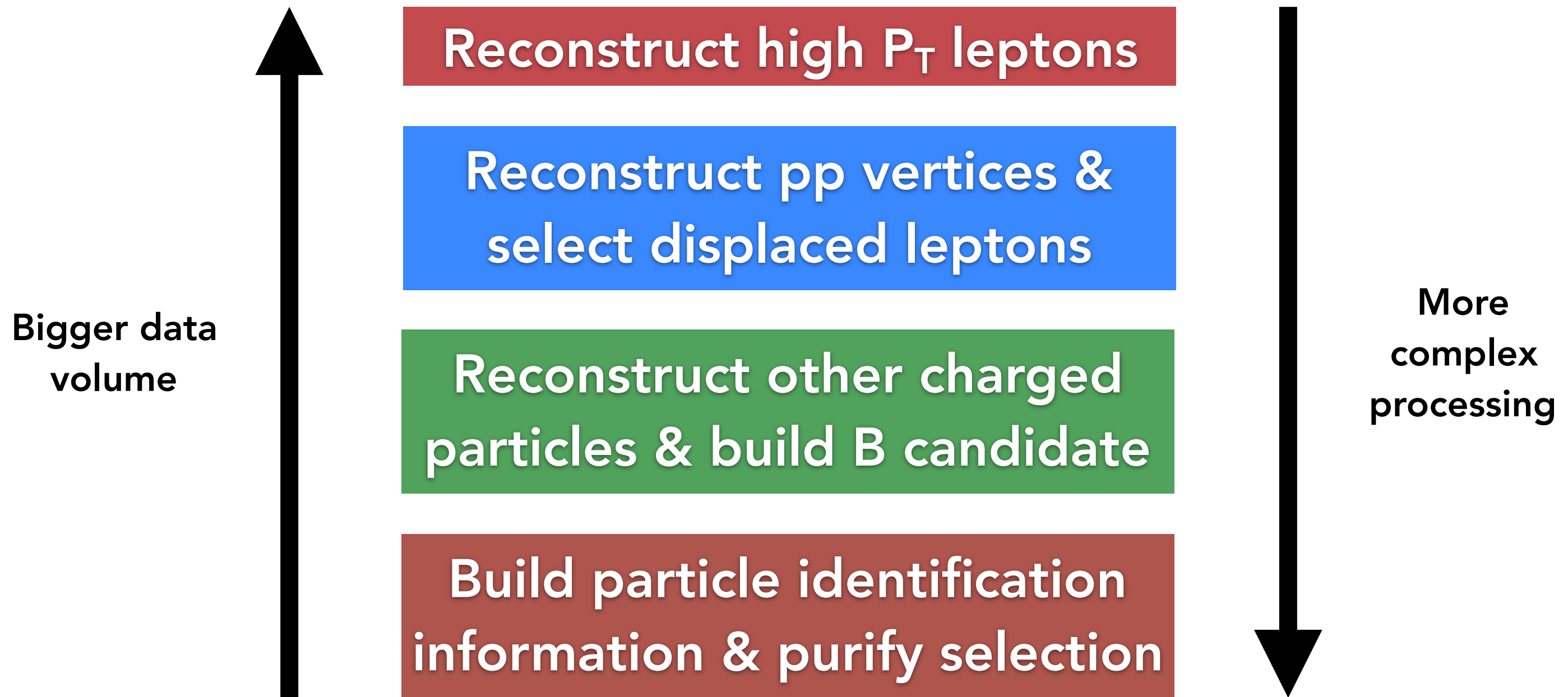
| Existing | Developing |
|----------|------------|
| μ | e,π,K,p |

## Particle identification
Tag-and-probe

Tag-and-probe calibrations exist for all charged particle species and for $\pi^0/\gamma$, with new sources added over time to improve coverage
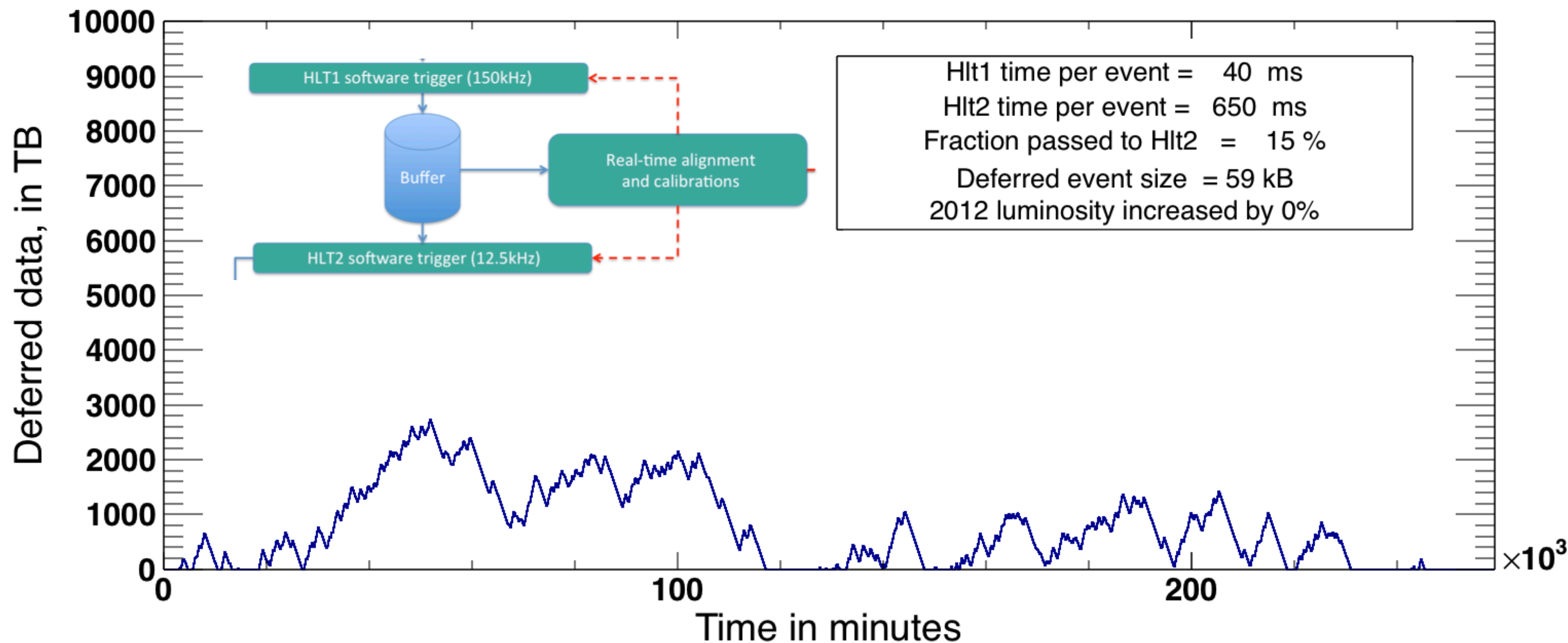
**Data driven efficiency calibration key to precision physics**
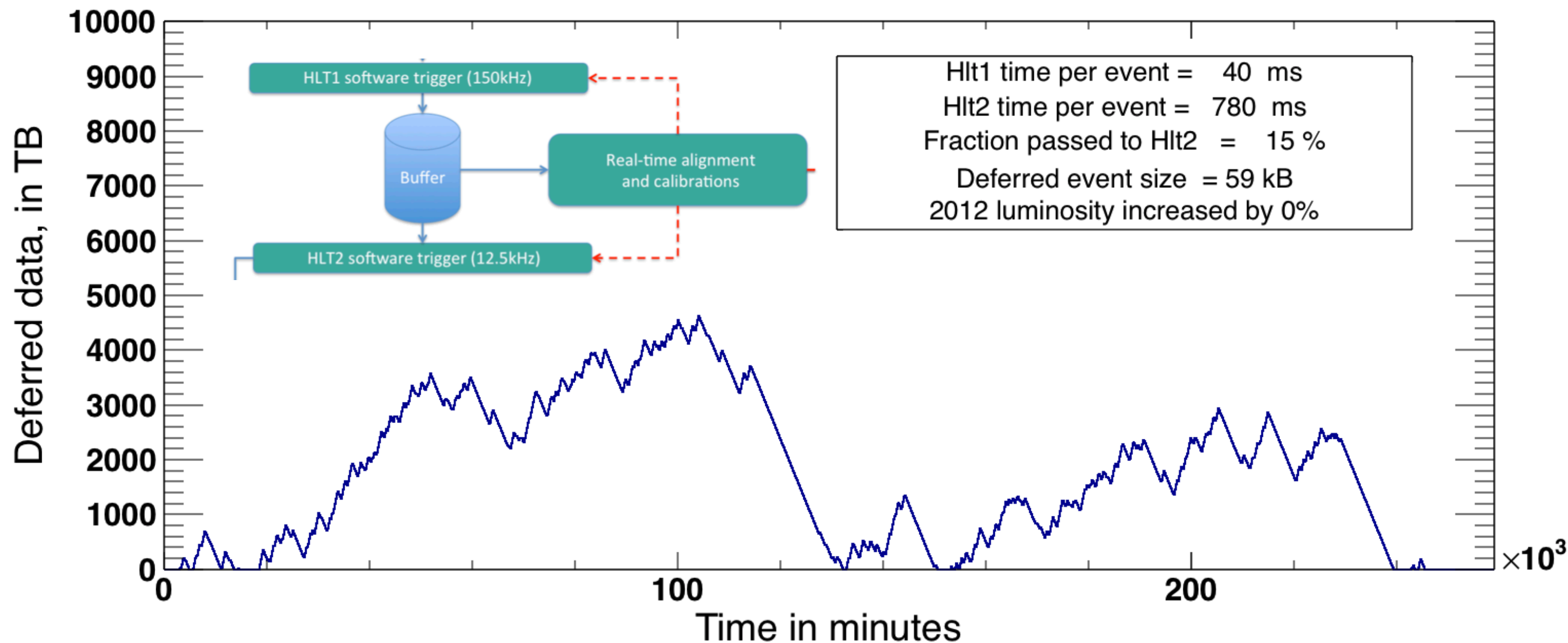
# What is a cascade buffer?

Reconstruct high $P_T$ leptons

Reconstruct pp vertices & select displaced leptons

Reconstruct other charged particles & build B candidate

Build particle identification information & purify selection

Bigger data volume

More complex processing

**A staged data reduction using increasingly complex algorithms**

**Use simulation to ensure robustness if timing estimates wrong**

# And what about data volumes?



**Data volume increases quadratically even with 0 background.
Select pp collisions, not bunch crossings, in real time!**