# STRATEGIES AND FUTURE TRENDS FOR TRIGGER AND DAQ SYSTEMS IN LHC EXPERIMENTS
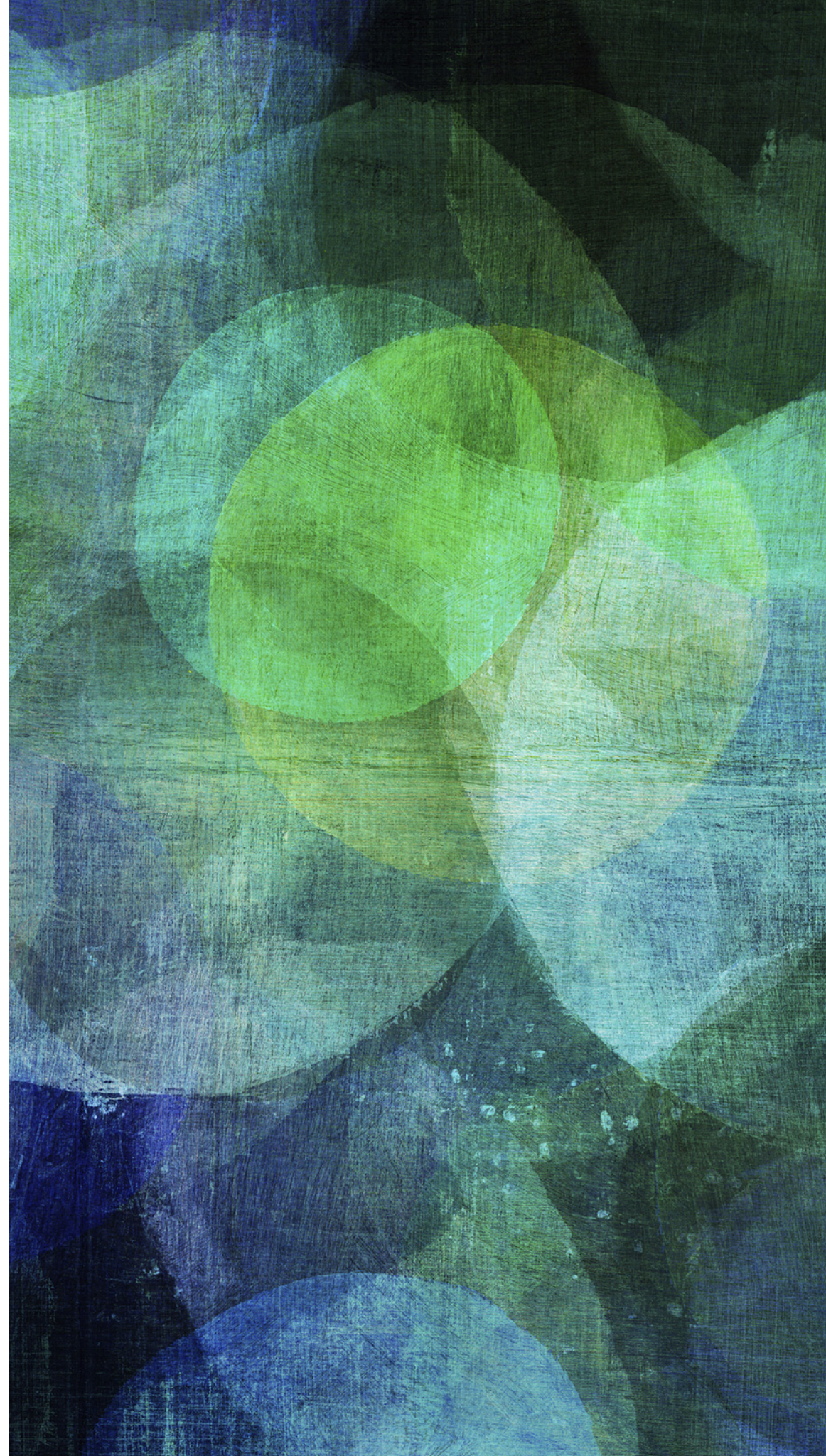
*F.Pastore (Royal Holloway Un. of London)*
*francesca.pastore@cern.ch*

# TRIGGERING AND TAKING DATA AT LHC
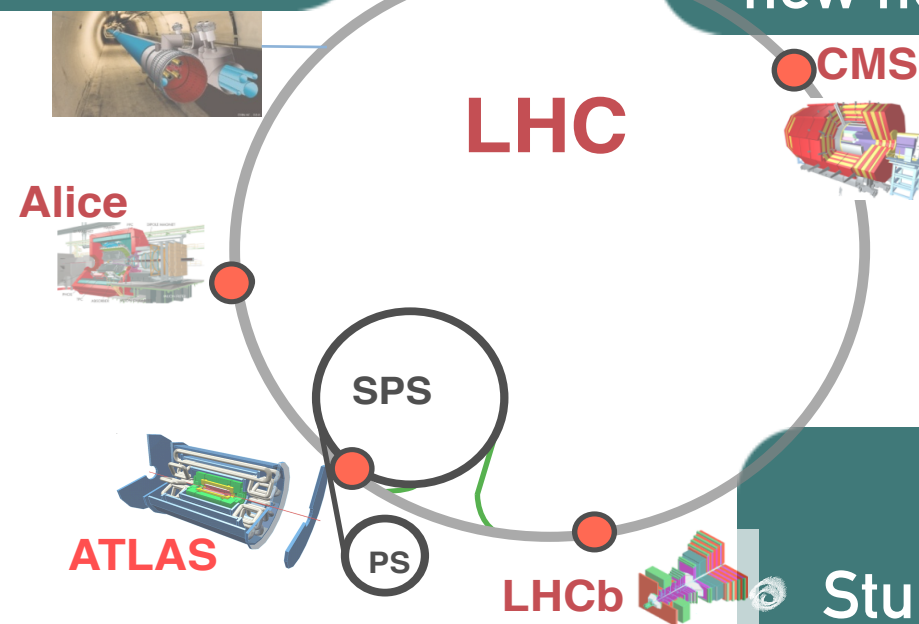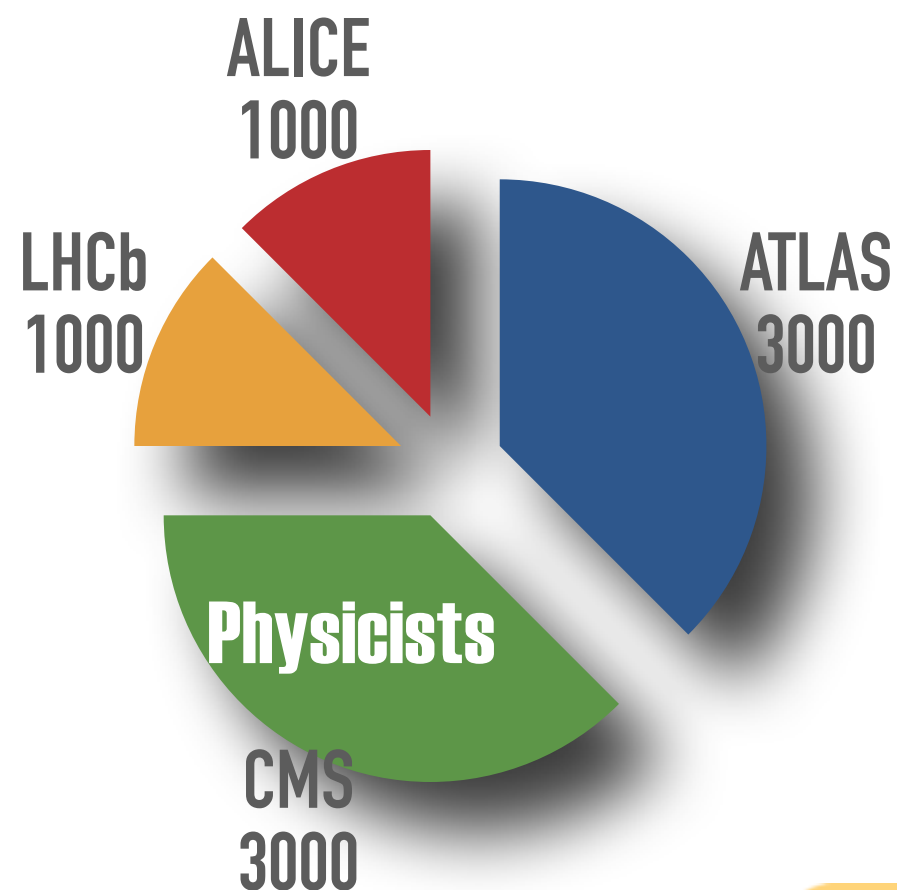
*TDAQ for large discovery experiments*

**Goal: explore TeV energy scale to find New Physics beyond Standard Model**

## ATLAS & CMS

- Completing the Standard Model and probing the Higgs sector
- Extending the reach for new physics beyond the Standard Model

## LHCb

- Study CP violation and rare decays in b- and c-quark sector
- Search for deviations of SM due to new heavy particles

## ALICE

Studying quark-gluon plasma, a complex system of strongly interacting matter produced by heavy ion collisions



ALICE
1000

LHCb
1000

ATLAS
3000

Physicists

CMS
3000

LHC

CMS

Alice

SPS

PS

ATLAS

LHCb

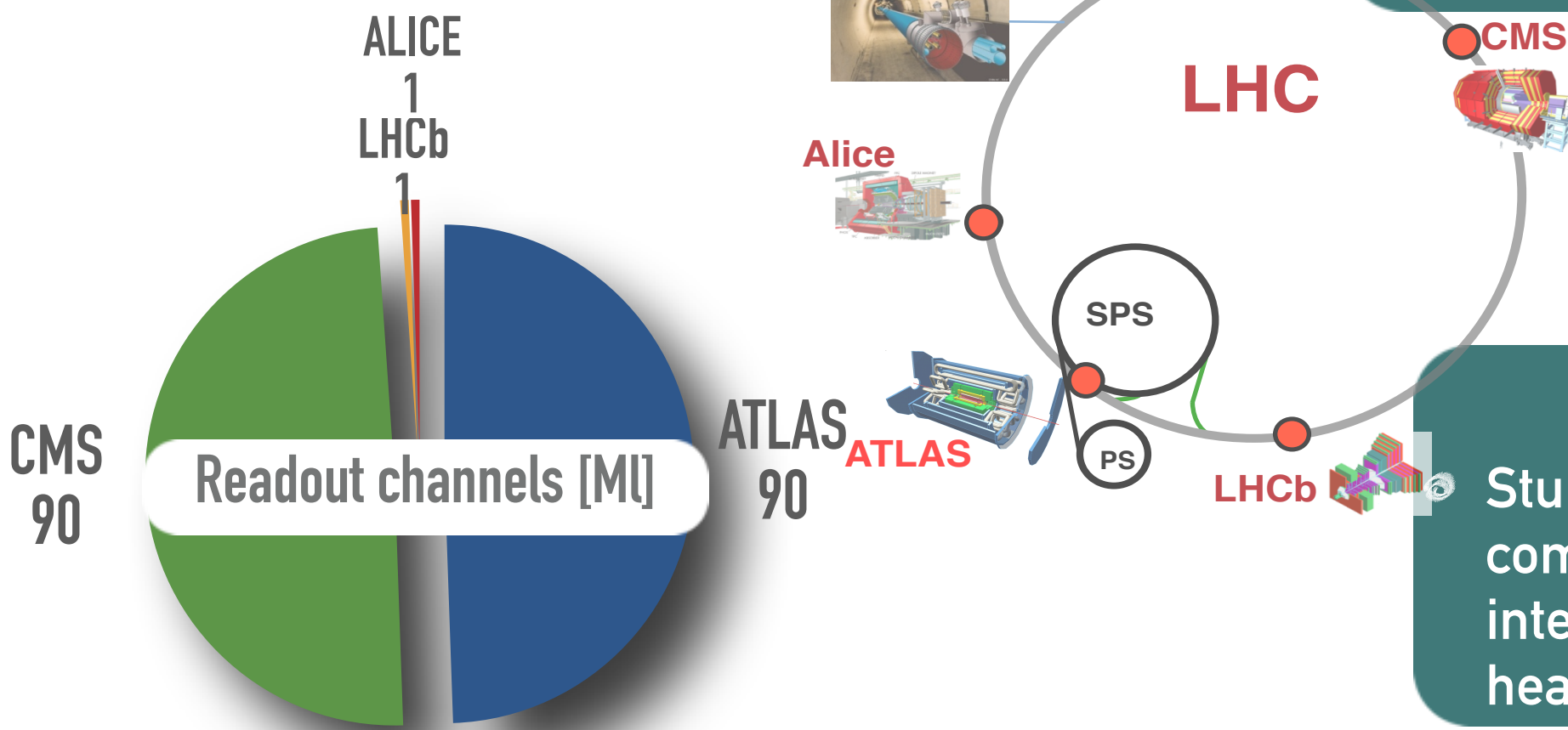**Proposed: 1992, Approved: 1996, Started: 2009**

**Goal: explore TeV energy scale to find New Physics beyond Standard Model**

## ATLAS & CMS

- Completing the Standard Model and probing the Higgs sector
- Extending the reach for new physics beyond the Standard Model

## LHCb

- Study CP violation and rare decays in b- and c-quark sector
- Search for deviations of SM due to new heavy particles

## ALICE

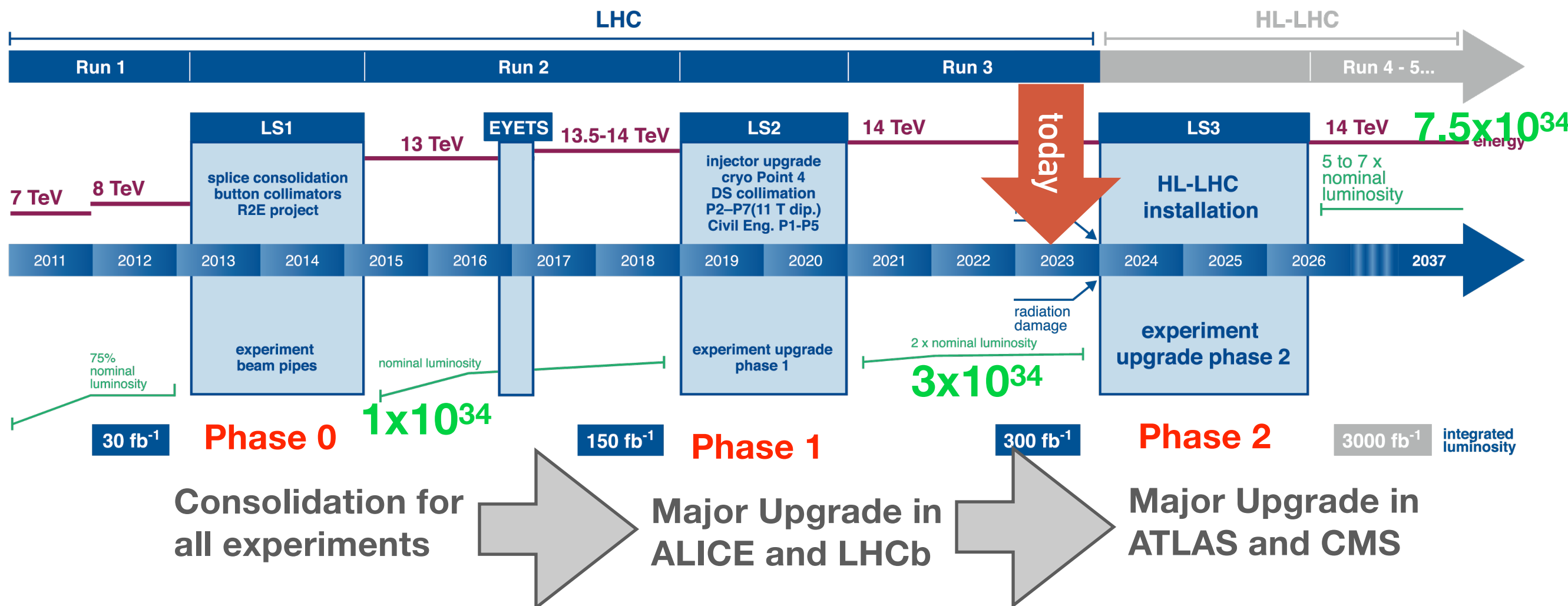Studying quark-gluon plasma, a complex system of strongly interacting matter produced by heavy ion collisions

ALICE
1
LHCb
1

CMS
90

ATLAS
90

**Readout channels [Ml]**

LHC

Alice

SPS

PS

ATLAS

CMS

LHCb

**Proposed: 1992, Approved: 1996, Started: 2009**

# LHC DATA DELUGE

p-p collisions
$E_{cms}$ = 13-14 TeV
$L = 10^{34}$ /cm$^2$ s
BC clock = 40 MHz

- High Luminosity with collisions close in time and space (1 collision/25ns)
  - fast electronics ➡ fast decisions
  - fine granularity detector ➡ high data volume
- Search for rare physics from hadronic collisions:
  - to store all the possibly relevant data is UNREALISTIC and often UNDESIRABLE
- Three approaches are possible:
  - Reduce the amount of data (packing and/or filtering)
  - Have faster data transmission and processing
  - Both!

**Buffering and parallelism**

**Maximum 1-2% deadtime**

**40 MHz COLLISION RATE**

Level-1

**DETECTOR CHANNELS**

Charge    Time    Pattern

**High speed electronics**

Energy    Tracks

**Level-1 triggers**

➡ Set max Readout rate

➡ Hardware, synchronous

➡ Readout parallelism

➡ Latency ~ μsec/event

Readout Buffers

**Readout links and buffering**

**Readout**

L1

Event building

**SWITCH NETWORK**

**Large data network with dedicated technology**

**DAQ**

HLT

Event filtering

**Dedicated PC farms**

Petabyte archive

**Computing Services**

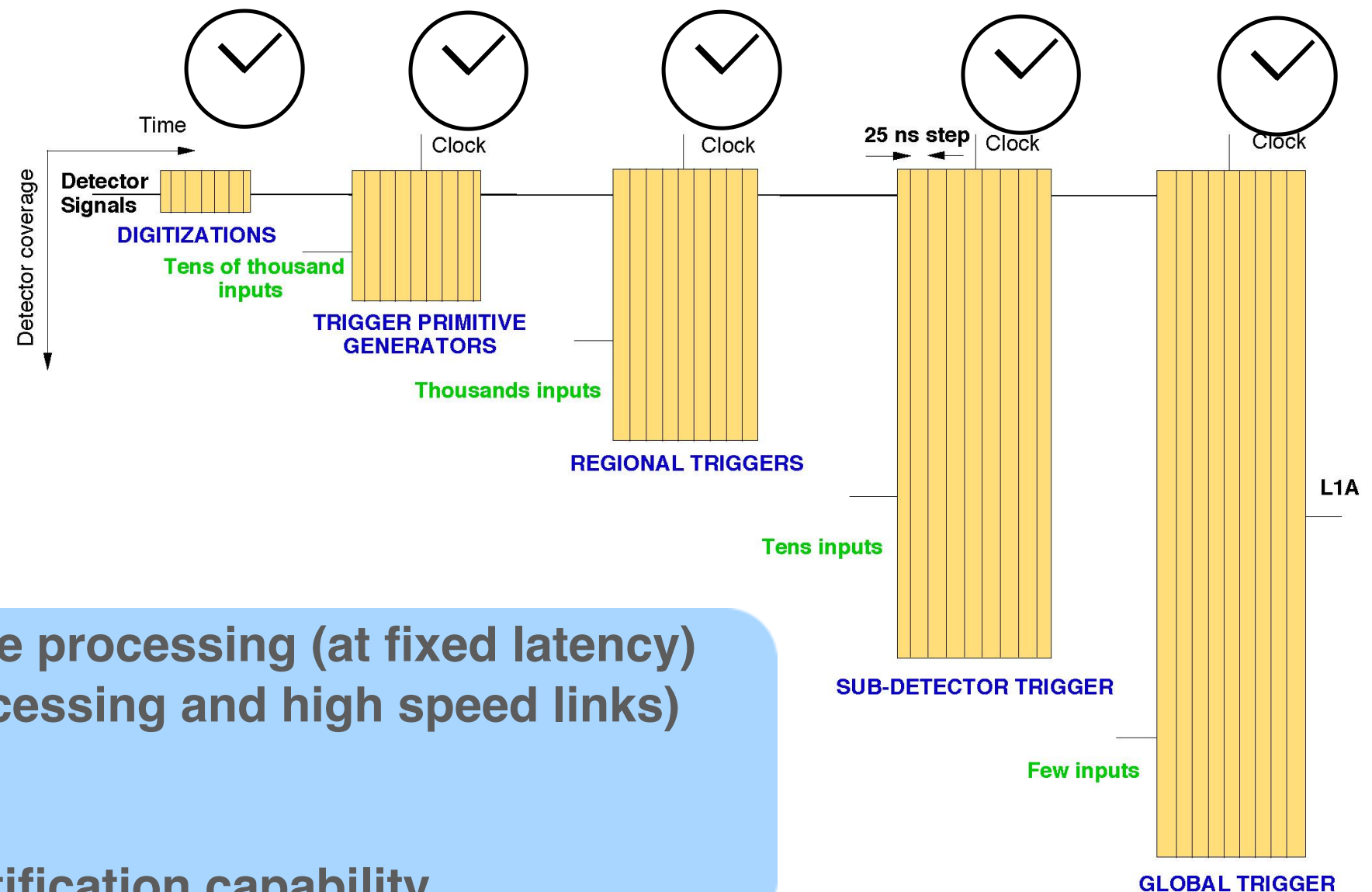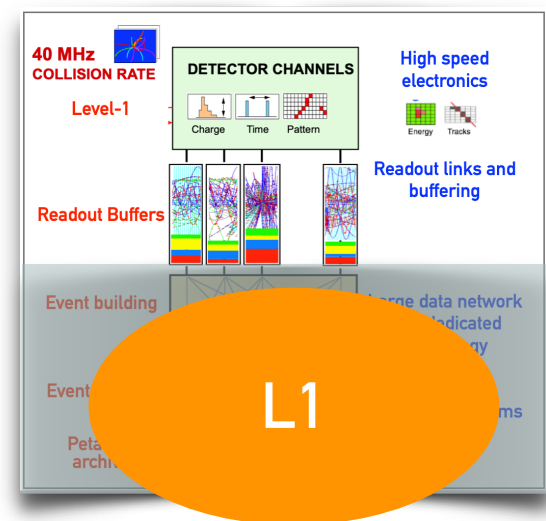**Higher level triggers**

➡ Set max storage rate

➡ Software, asynchronous

➡ Event parallelism

➡ Latency < 1 sec/event

# LEVEL-1 TRIGGER PRINCIPLES



➡ **Synchronous: pipeline processing (at fixed latency)**
➡ **Low latency (fast processing and high speed links)**
➡ **Scalable**
➡ **Massively parallel**
➡ **Bunch Crossing identification capability**

**Full synchronisation at 40 MHz (LHC clock)**
➤ large optical time distribution system

# Fast, robust electronics

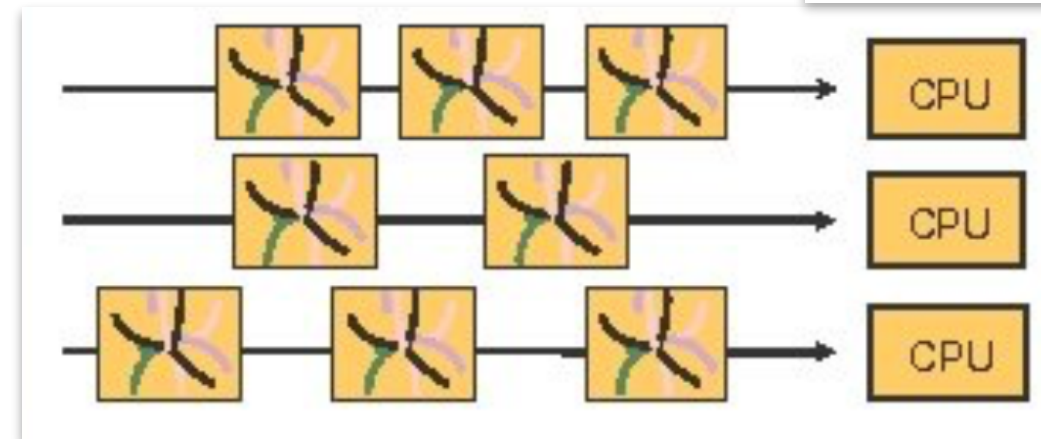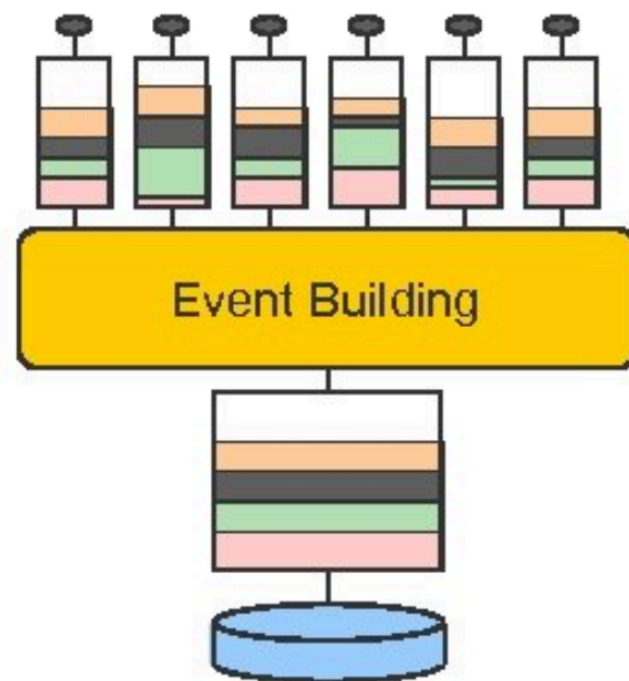| ALICE | No pipeline |
|-------|-------------|
| ATLAS | 2.5 µs |
| CMS | 3 µs |
| LHCb | 4 µs |

**Latency dominated by cable/transmission delay**

# HLT/DAQ REQUIREMENTS

**HLT**

➡ **Robustness and redundancy**
➡ **Scalability to adapt to Luminosity, detectors,…**
➡ **Flexibility (10-years experiments)**
➡ **Based on commercial products**
➡ **Limited cost**



## ATLAS/CMS Example

➡ **1 MB/event at 100 kHz for O(100ms) HLT latency**

   ➡ <u>Network</u>: 1 MB*100 kHz = **100 GB/s**

   ➡ <u>HLT farm</u>: 100 kHz*100 ms = **O($10^4$) CPU cores**

➡ Can add intermediate steps (level-2) to reduce resources, at cost of complexity (at ms scale)

*See S.Cittolin, DOI: 10.1098/rsta.2011.0464*

➤ Event Building and Filter farms on networks
   ➤ farm processing: one event per processor (larger latency, but scalable)
   ➤ additional networks regulates the CPU assignment

8

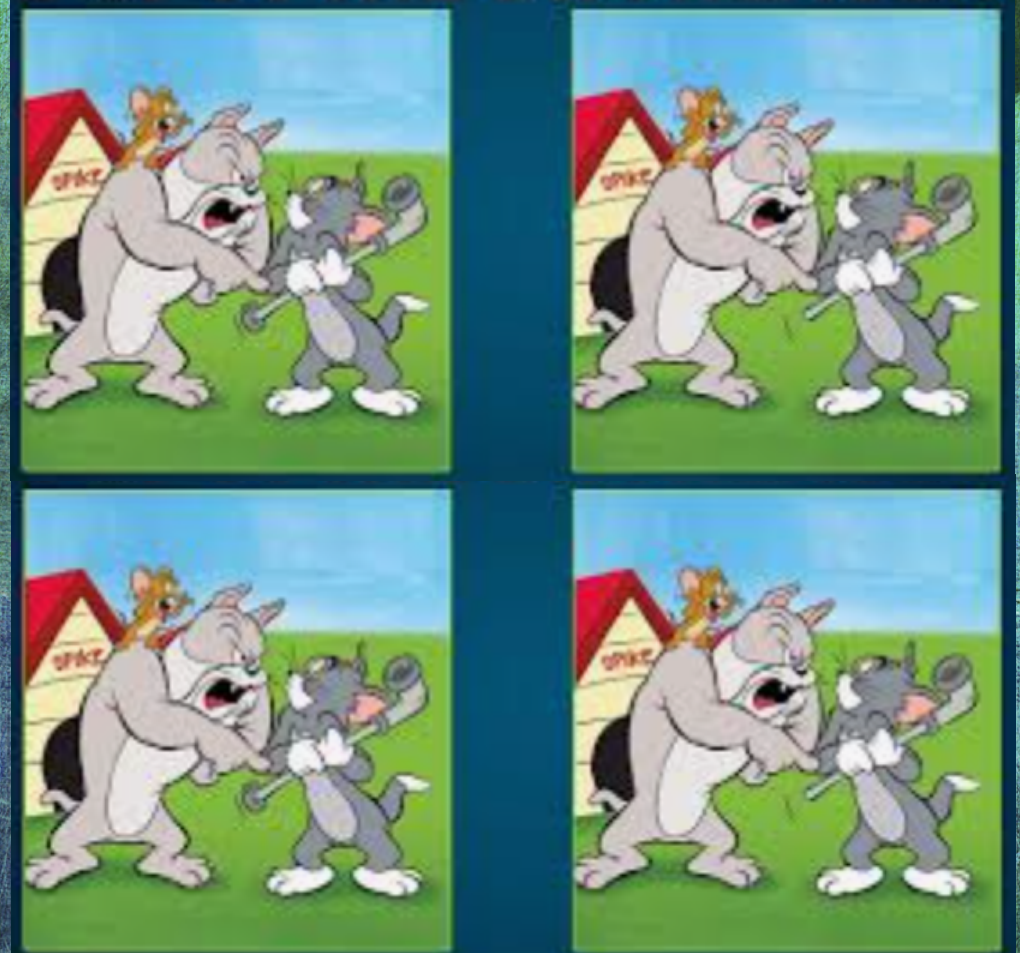# COMPARE 4 EXPERIMENTS

*How to maximise physics acceptance*

# DIFFERENT PHYSICS SEARCHES

…. and LHC operations

- **ATLAS/CMS: p-p collisions at full Luminosity**
  - search in high energy scale

- **LHCb: p-p collisions at reduced Luminosity**
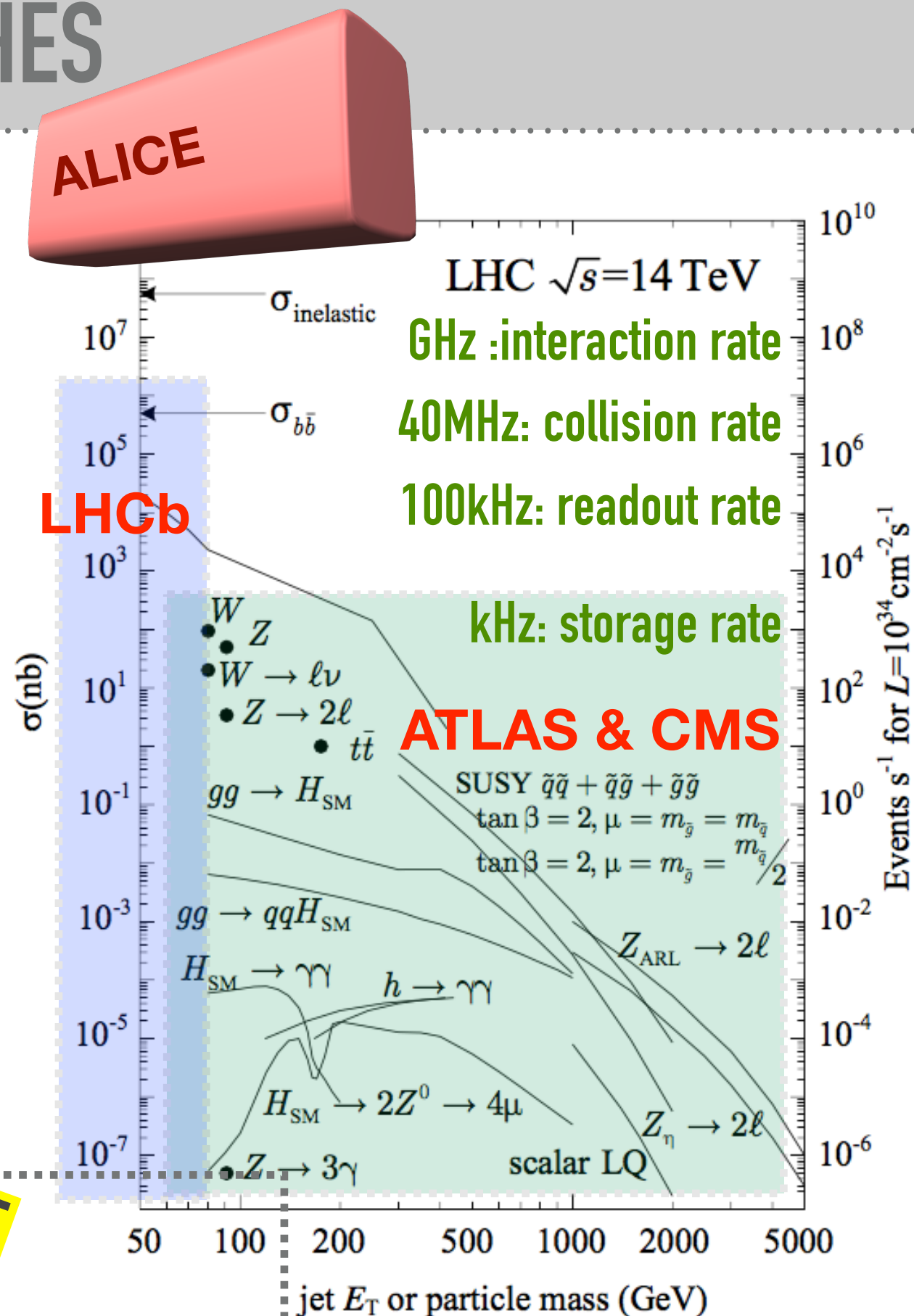  - search complex topologies of b-quark decays
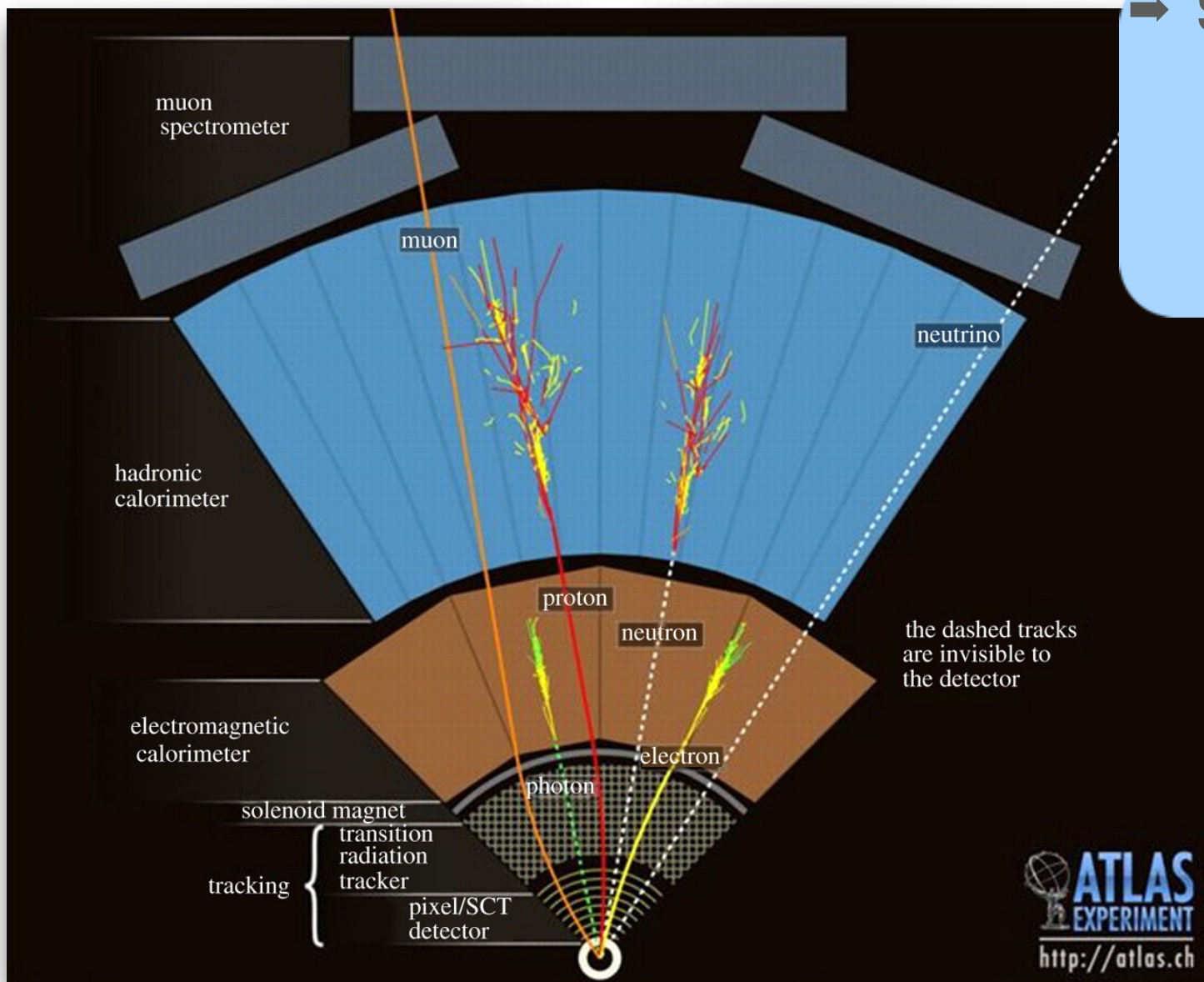
- **ALICE: heavy-ion collisions ~2000 mb**
  - search in high energy density

**DIFFERENT**

- ➡ **Expected rates and S/B ratio**
- ➡ **Signal topology and complexity**
- ➡ **Size of event (number of channels, particle multiplicity)**

**ALICE**

**LHCb**

**ATLAS & CMS**

LHC $\sqrt{s}=14\,\mathrm{TeV}$

GHz :interaction rate

40MHz: collision rate

100kHz: readout rate

kHz: storage rate

$\sigma_{\mathrm{inelastic}}$

$\sigma_{b\bar{b}}$

$W$

$Z$

$W \to \ell\nu$

$Z \to 2\ell$

$t\bar{t}$

$gg \to H_{\mathrm{SM}}$

SUSY $\tilde{q}\tilde{q} + \tilde{q}\tilde{g} + \tilde{g}\tilde{g}$
$\tan\beta = 2, \mu = m_{\tilde{g}} = m_{\tilde{q}}$
$\tan\beta = 2, \mu = m_{\tilde{g}} = m_{\tilde{q}}/2$

$gg \to qqH_{\mathrm{SM}}$

$H_{\mathrm{SM}} \to \gamma\gamma$

$h \to \gamma\gamma$

$Z_{\mathrm{ARL}} \to 2\ell$

$H_{\mathrm{SM}} \to 2Z^0 \to 4\mu$

$Z_\eta \to 2\ell$

$Z \to 3\gamma$

scalar LQ

$\sigma$(nb)

Events s$^{-1}$ for L=$10^{34}$cm$^{-2}$s$^{-1}$

jet $E_{\mathrm{T}}$ or particle mass (GeV)

➡ **Search in high-energy scale**

   ➡ Discover large mass particles through their <u>high-energy</u> products

   ➡ **Discovery** = inclusive selections

$$\frac{everything}{Higgs} = \frac{\sigma_{tot}}{\sigma_{H(500\,\mathrm{GeV})}} \approx \frac{100\,mb}{1\,pb} \approx 10^{11}$$
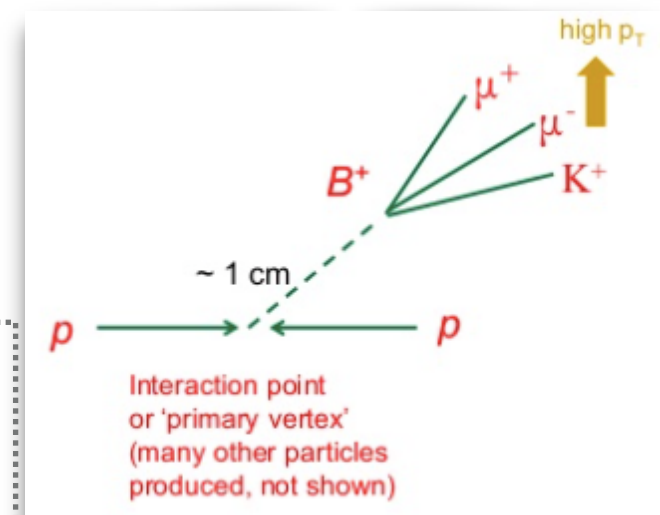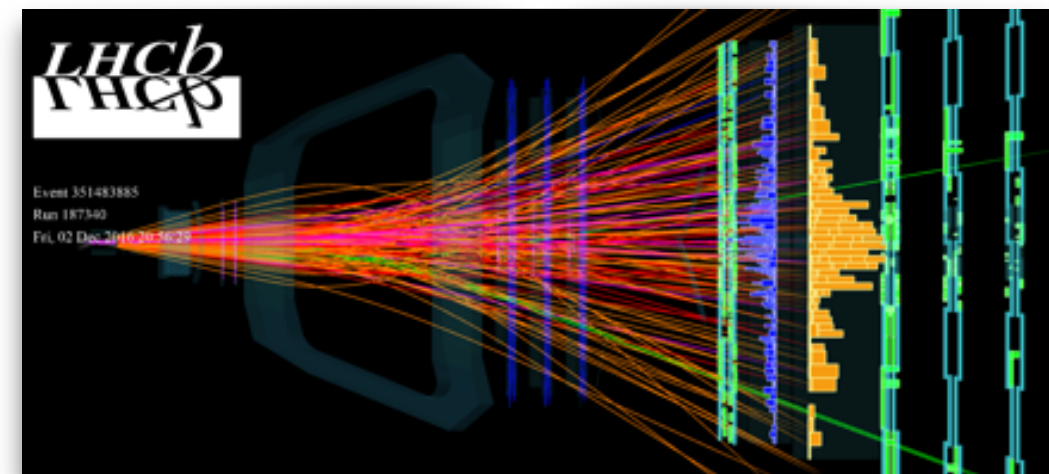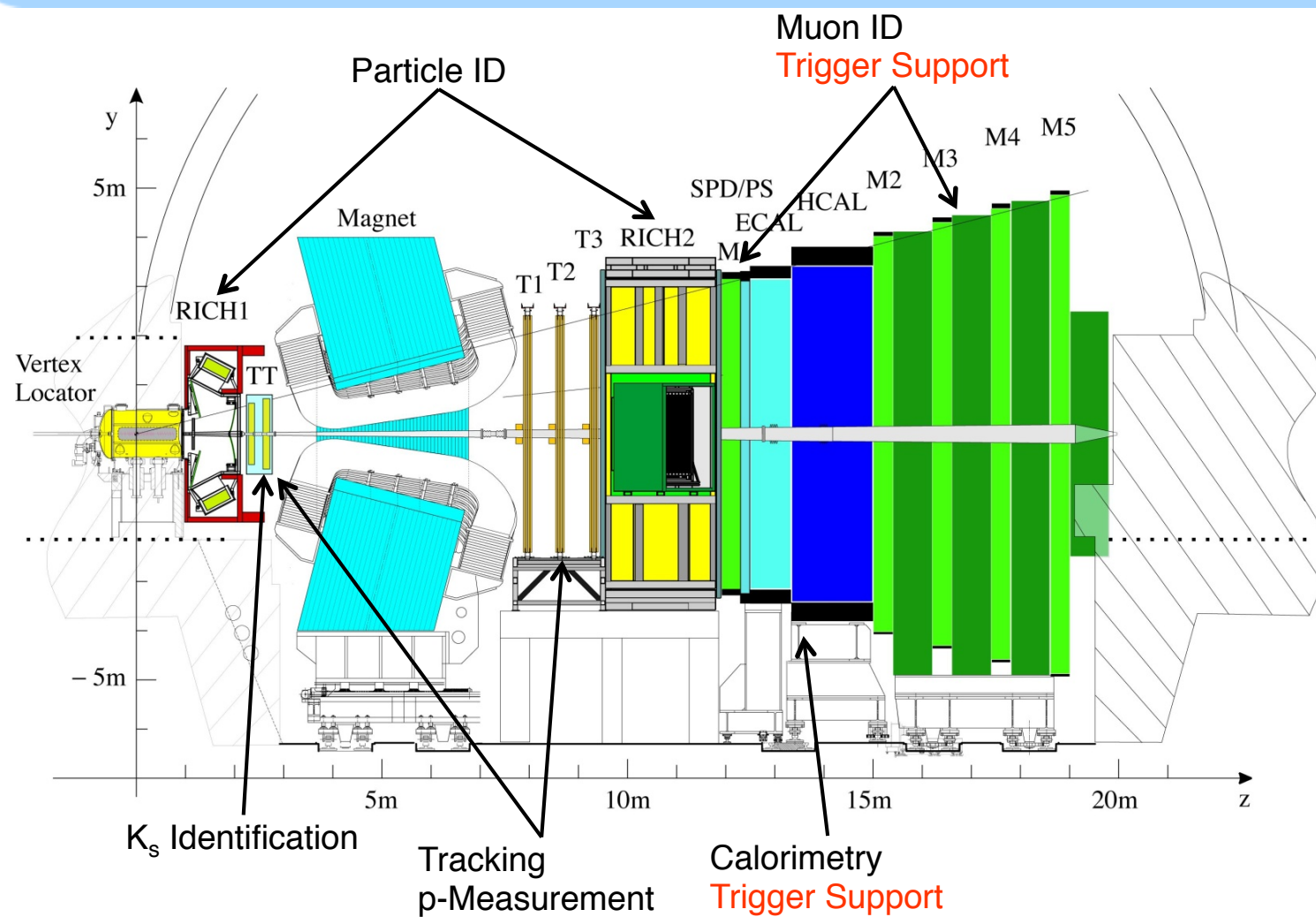
**approximately 10$^6$ rejection**

➡ **Easy selection of high-energy <u>leptons</u> over background ==> @L1**

   ➡ Against thousands of particles/collisions (typically low momentum jets)

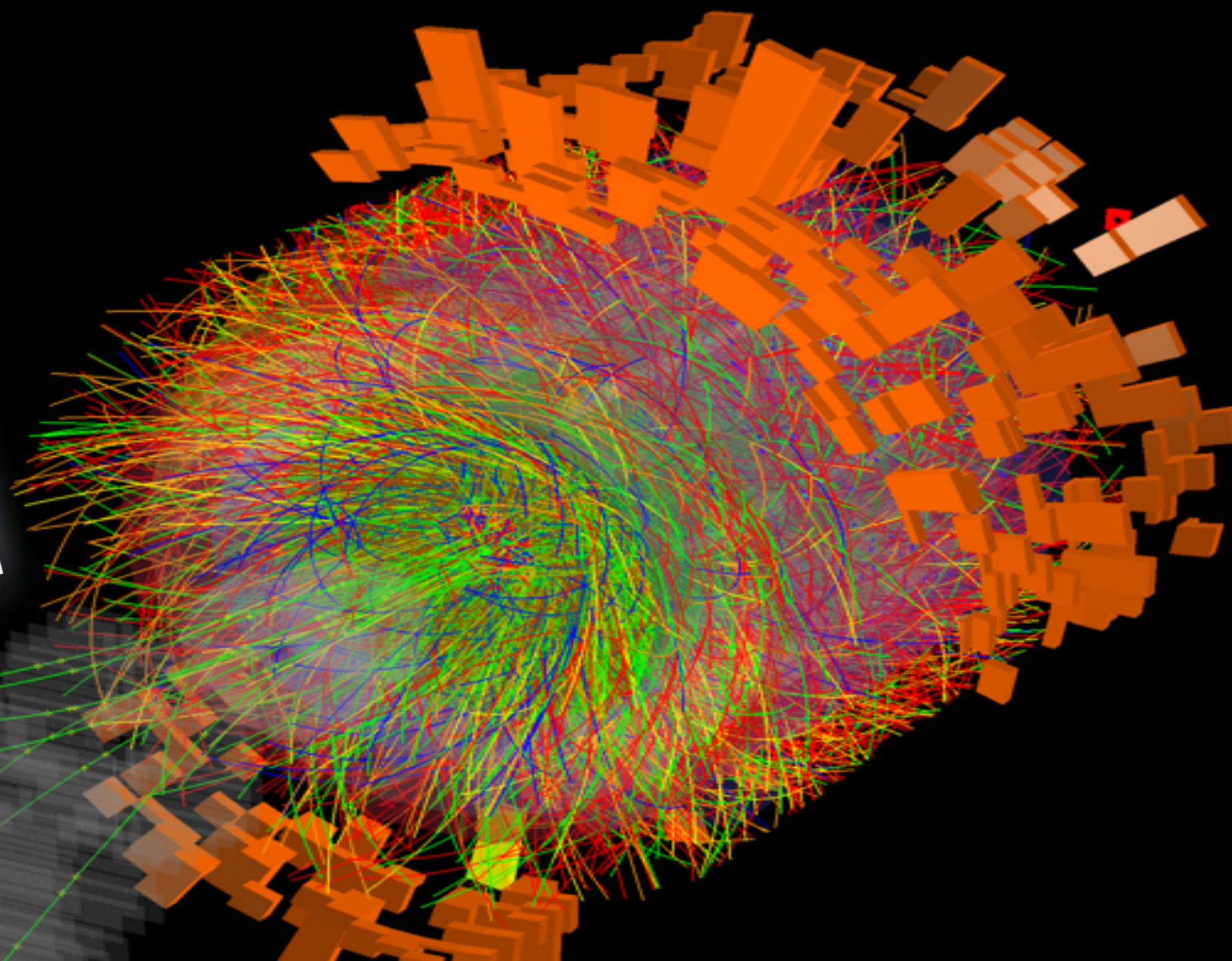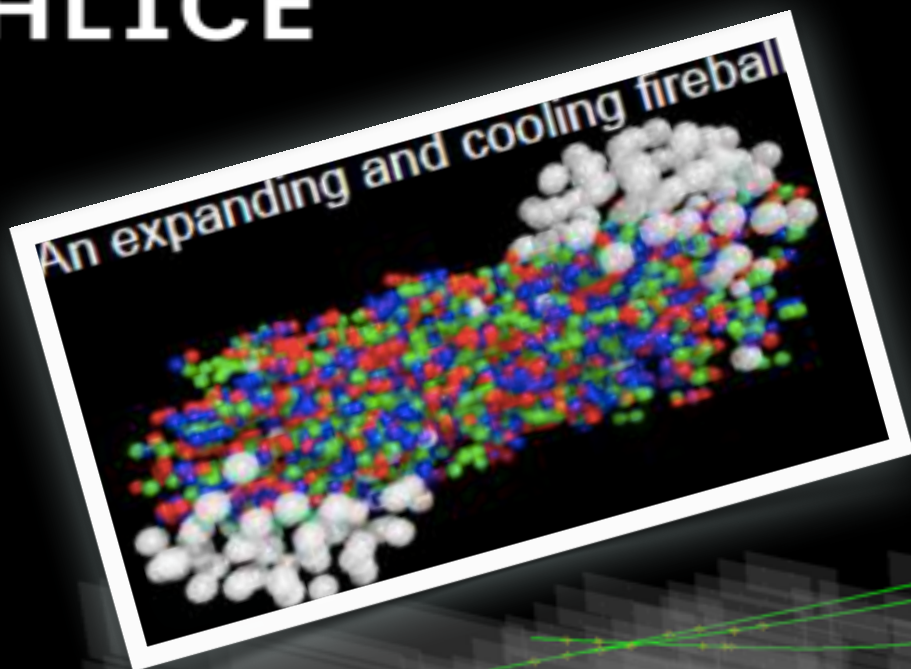➡ **Remember: 90M readout channels and full Luminosity ==> 1 MB/event**

➡ **Precision measurements and rare decays in the B system**
  ➡ Large production ($\sigma_{BB} \sim 500$ μb), but still $\sigma_{BB}/\sigma_{Tot} \sim 5 \times 10^{-3}$
  ➡ Interesting B decays are quite <u>rare</u> (BR $\sim 10^{-5}$ )



➡ **Single-arm spectrometer and low L ==> reduced event size**
➡ **Selection of B mesons ==> search for B-decay topologies**
  ➡ related to high mass and long lifetime of the b-quark

An expanding and cooling fireball

Run:244918
Timestamp:2015-11-25 11:25:36(UTC
System: Pb-Pb
Energy: 5.02 TeV
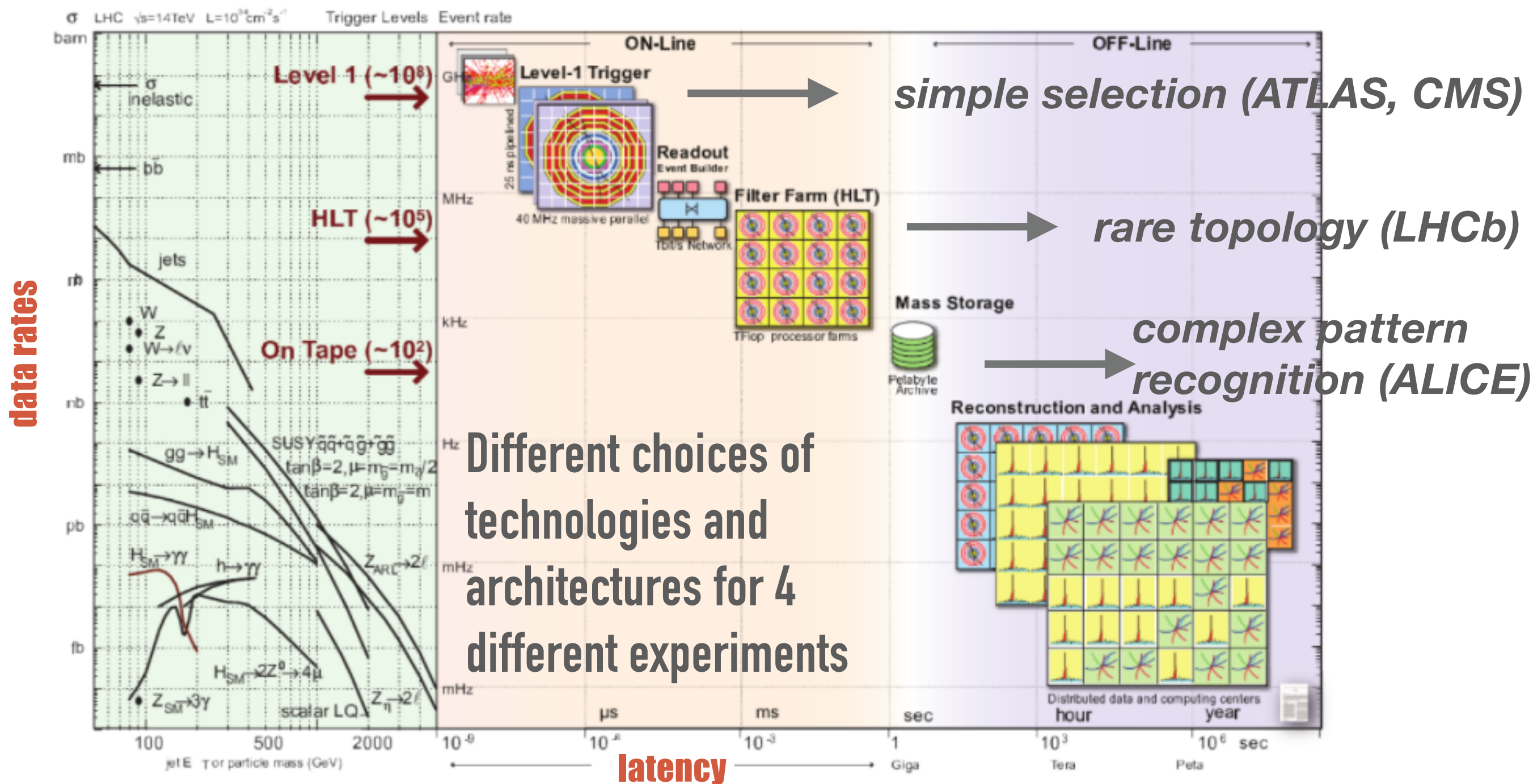
➡ **Physics of strongly interacting matters & quark-gluon plasma, with nucleus-nucleus interactions**
   - ➡ High particle multiplicities (~8000 particles/d$\eta$)
   - ➡ Identify heavy short-living particles
   - ➡ By selecting low-$p_T$ tracks (>100 MeV)

**Different choices of technologies and architectures for 4 different experiments**

*simple selection (ATLAS, CMS)*

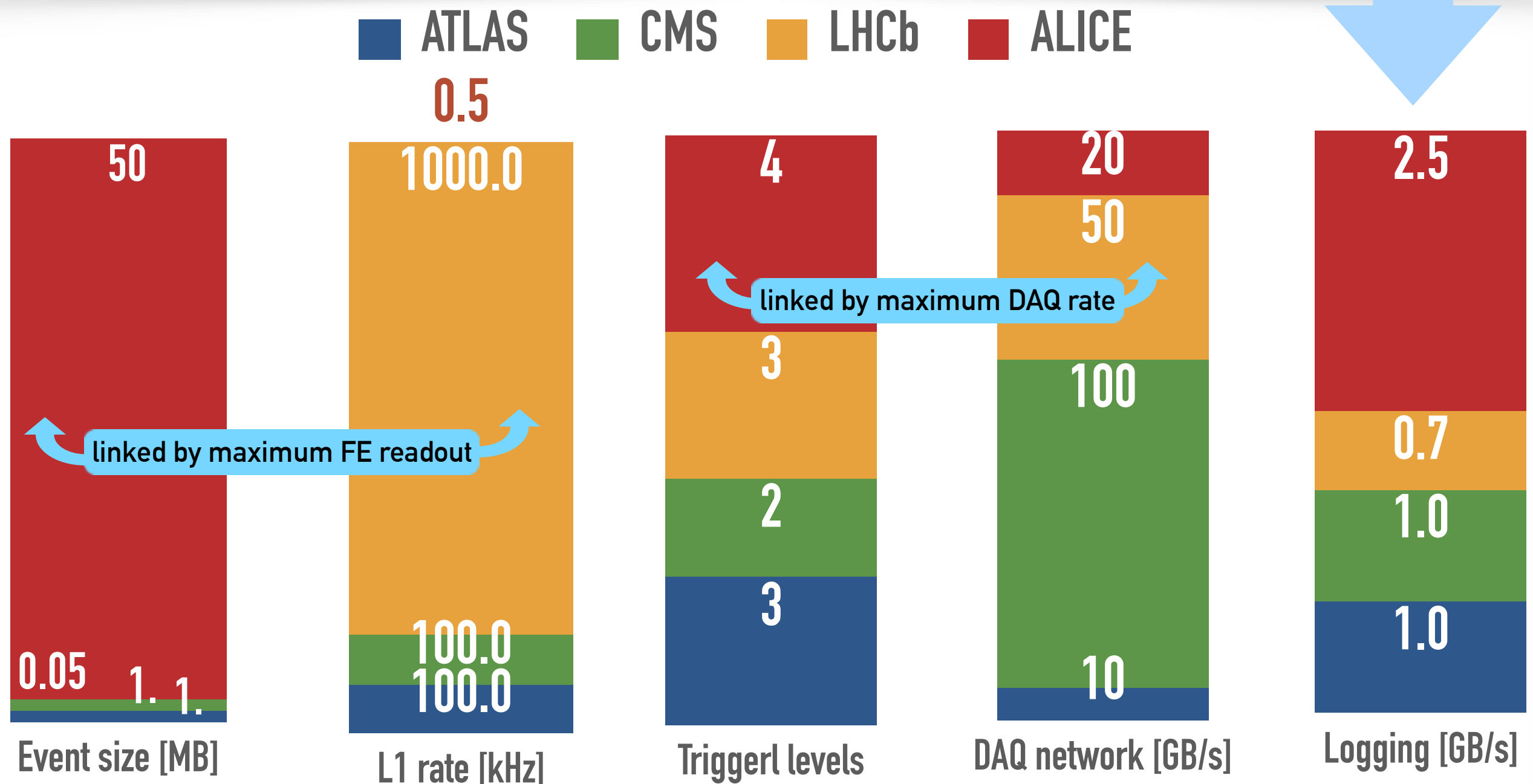*rare topology (LHCb)*

*complex pattern recognition (ALICE)*

➡ **ATLAS/CMS: Trigger power:** reducing the data-flow at the earliest stage

➡ **ALICE/LHCb: Large data-flow:** low trigger selectivity due to large irreducible background

# COMPARING BY NUMBERS

LHC experiments share the same CERN budget for computing resources, which is the constrain between trigger and DAQ power

Allowed storage and processing resources
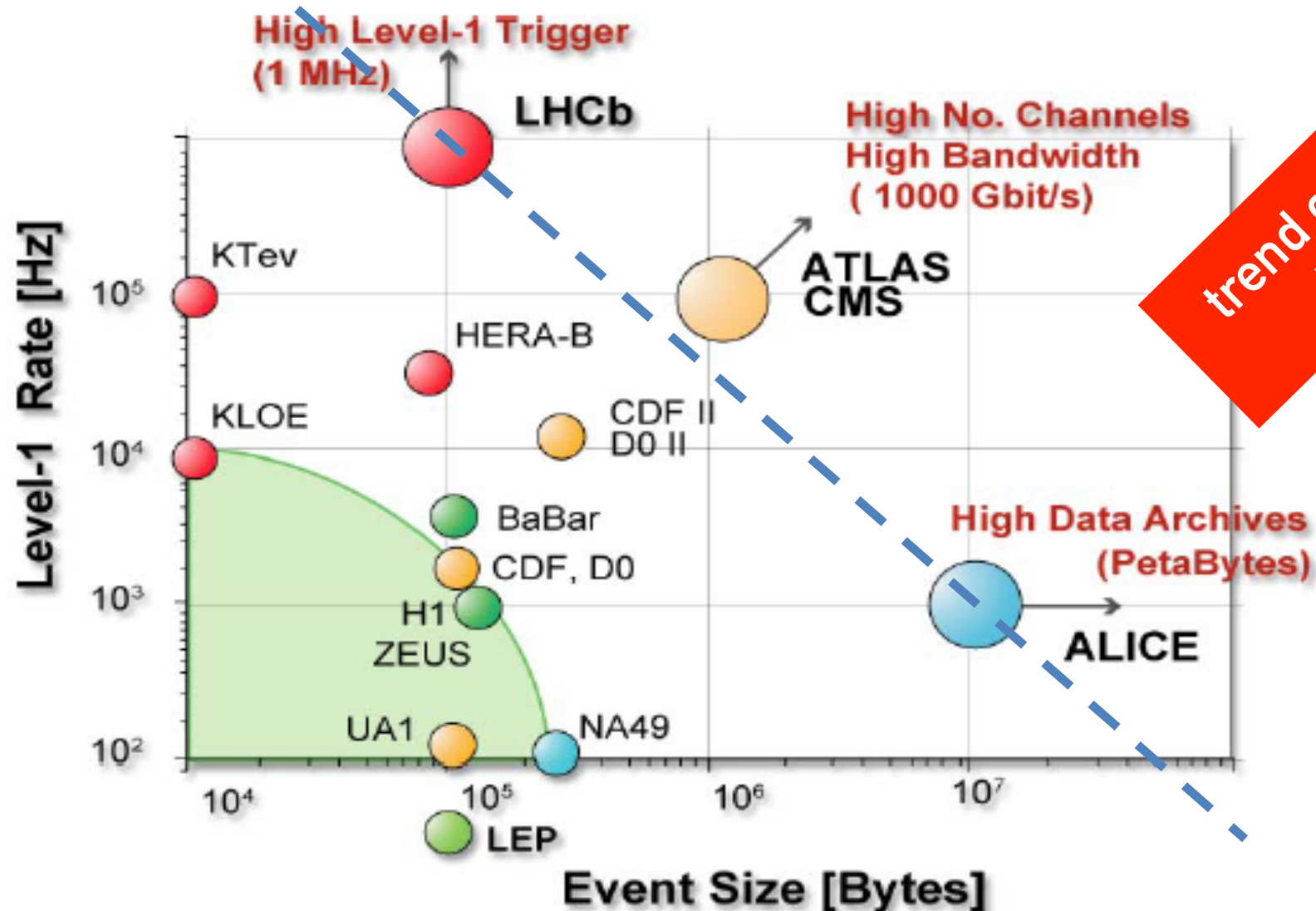
*Design values in 2009*



**ATLAS** **CMS** **LHCb** **ALICE**

linked by maximum DAQ rate

linked by maximum FE readout

| Event size [MB] | L1 rate [kHz] | Triggerl levels | DAQ network [GB/s] | Logging [GB/s] |
|---|---|---|---|---|
| 50 | 0.5 | 4 | 20 | 2.5 |
| | 1000.0 | 3 | 50 | 0.7 |
| | | 2 | 100 | 1.0 |
| 0.05  1.  1. | 100.0  100.0 | 3 | 10 | 1.0 |

$$R_{DAQ} = R_T^{max} \times S_E$$



*faster L1 electronics*

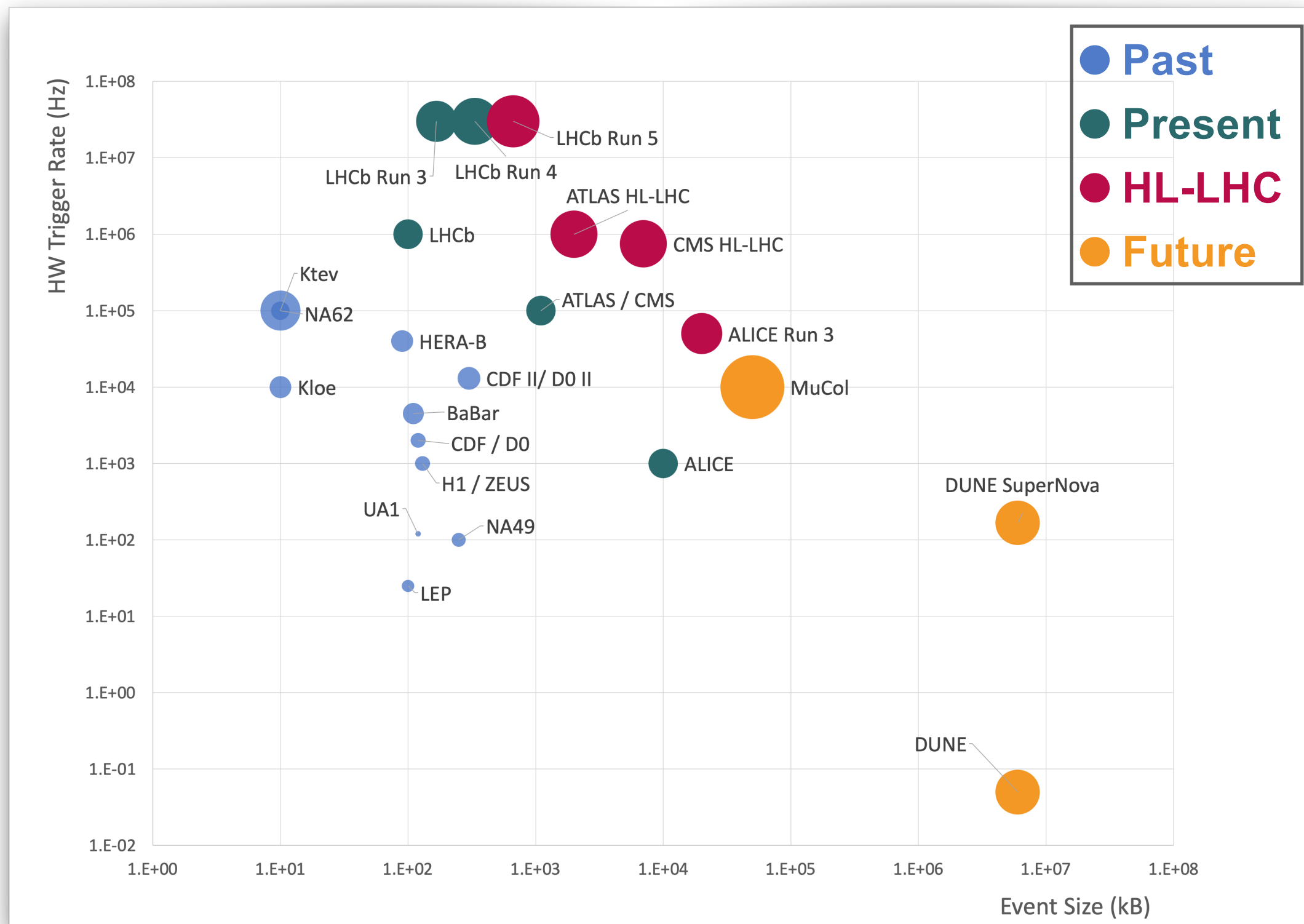*more channels, more complex events*

**ATLAS/CMS**

**Data to Process:**

100 kHz * 1 MB = 100 GB/s

**Data to Store:**

~ 1 PB / year /experiment

**As the data volumes and rates increase, new architectures need to be developed**

*Courtesy of A.Cerri*

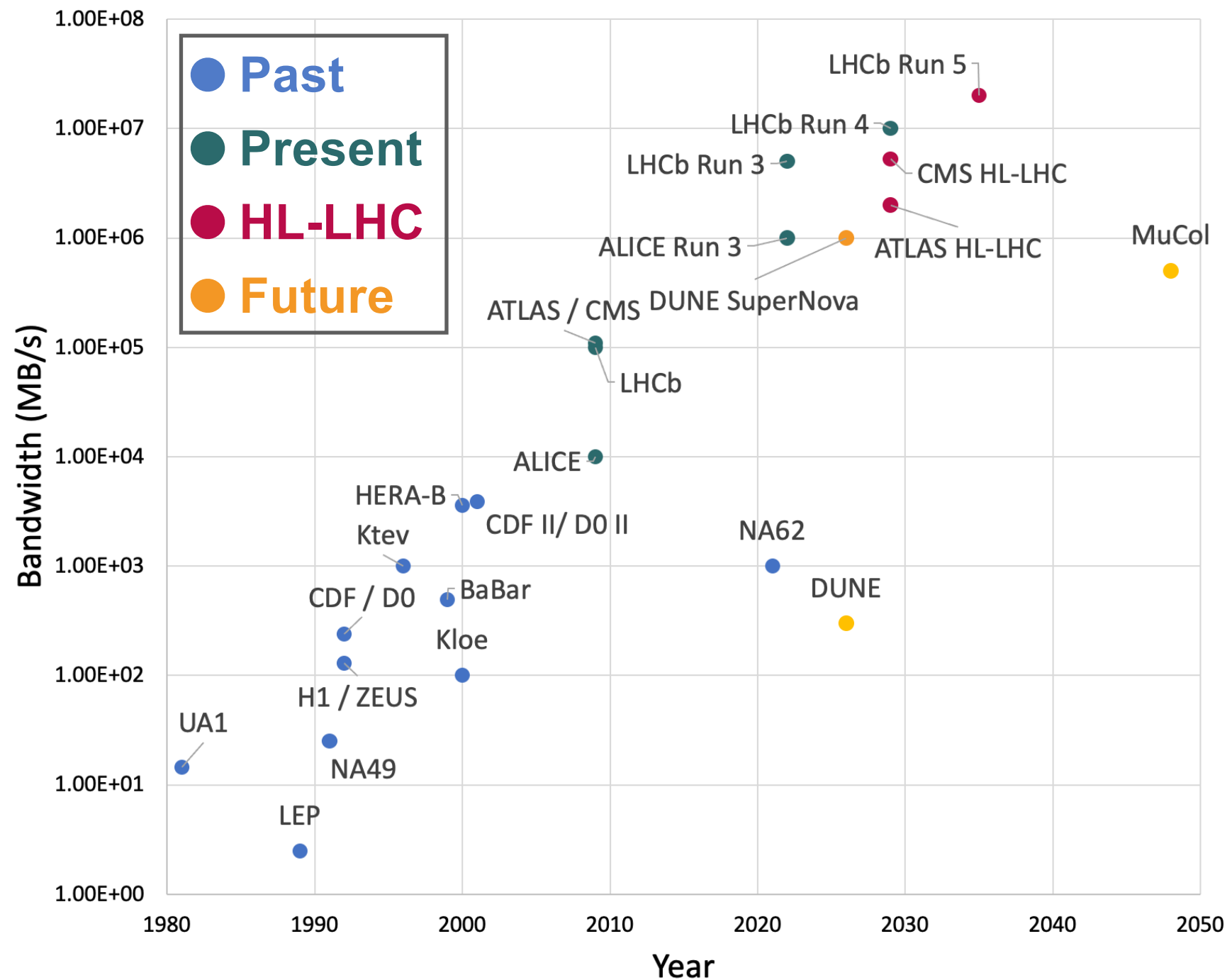*Courtesy of A.Cerri*

# FUTURE TRENDS FOR HIGH-LUMINOSITY

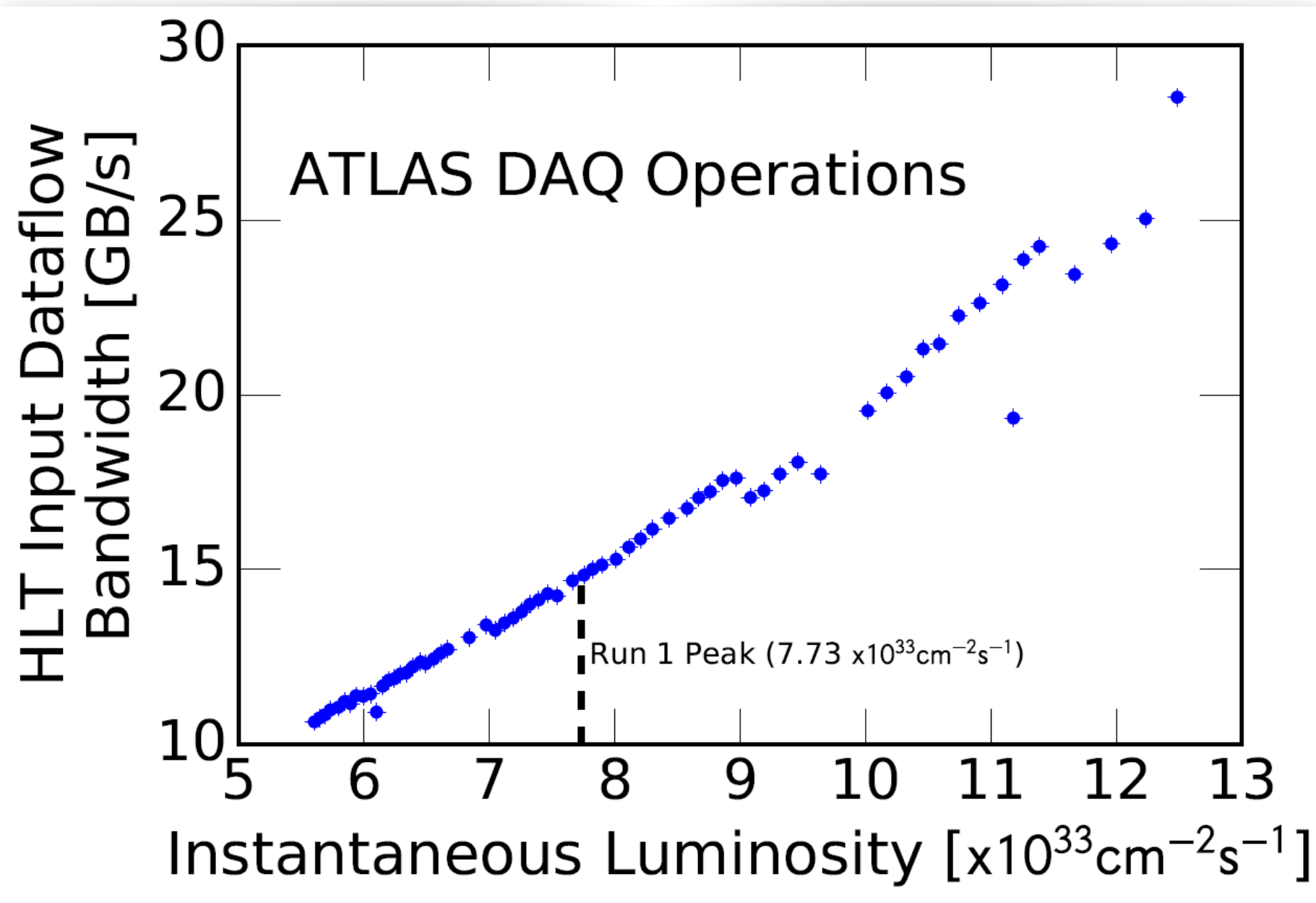*What about … tomorrow?*

# Design Luminosity x7.5

➡ **200 collisions per bunch crossing (any 25 ns)**

➡ **~ 10 000 particles per event**

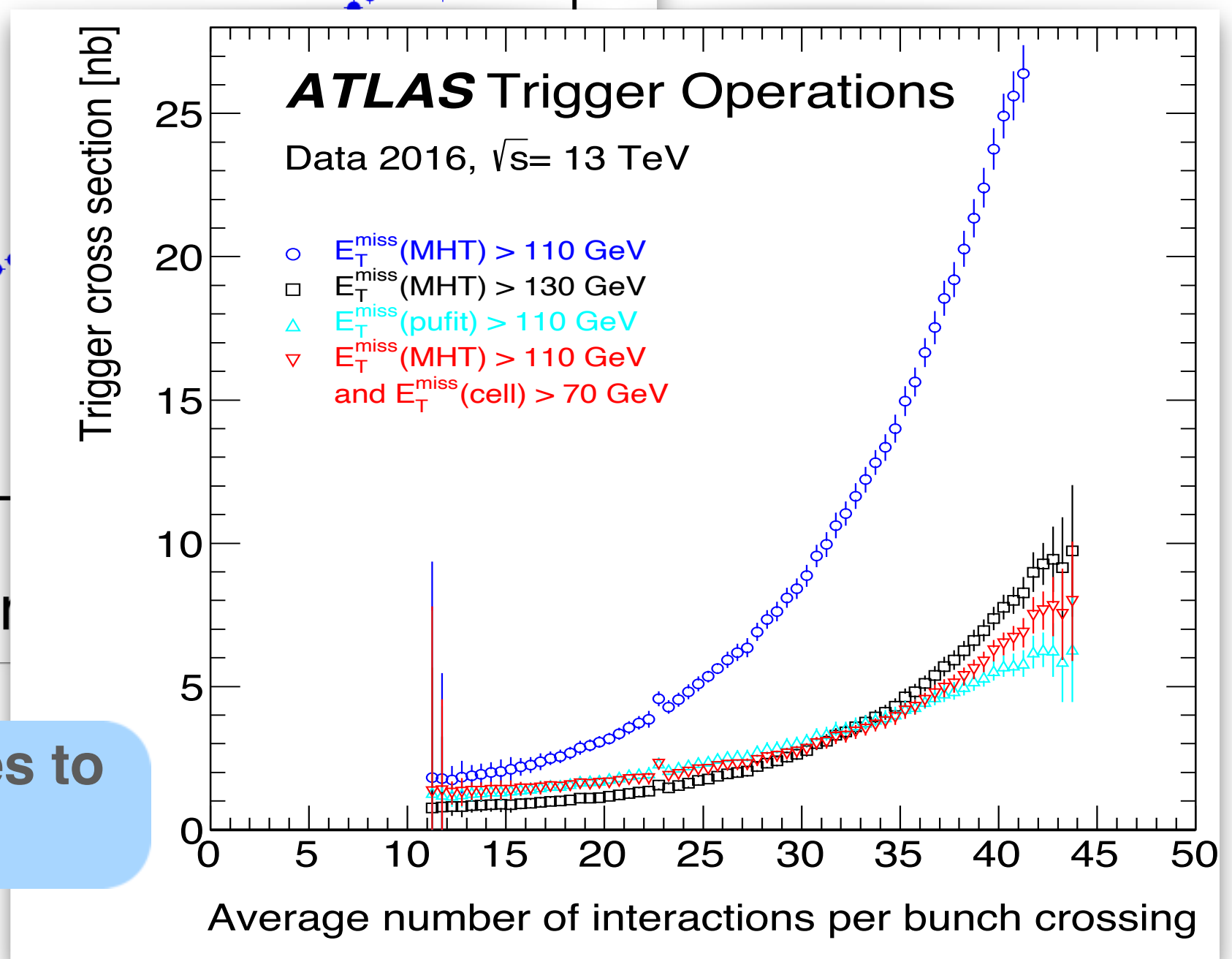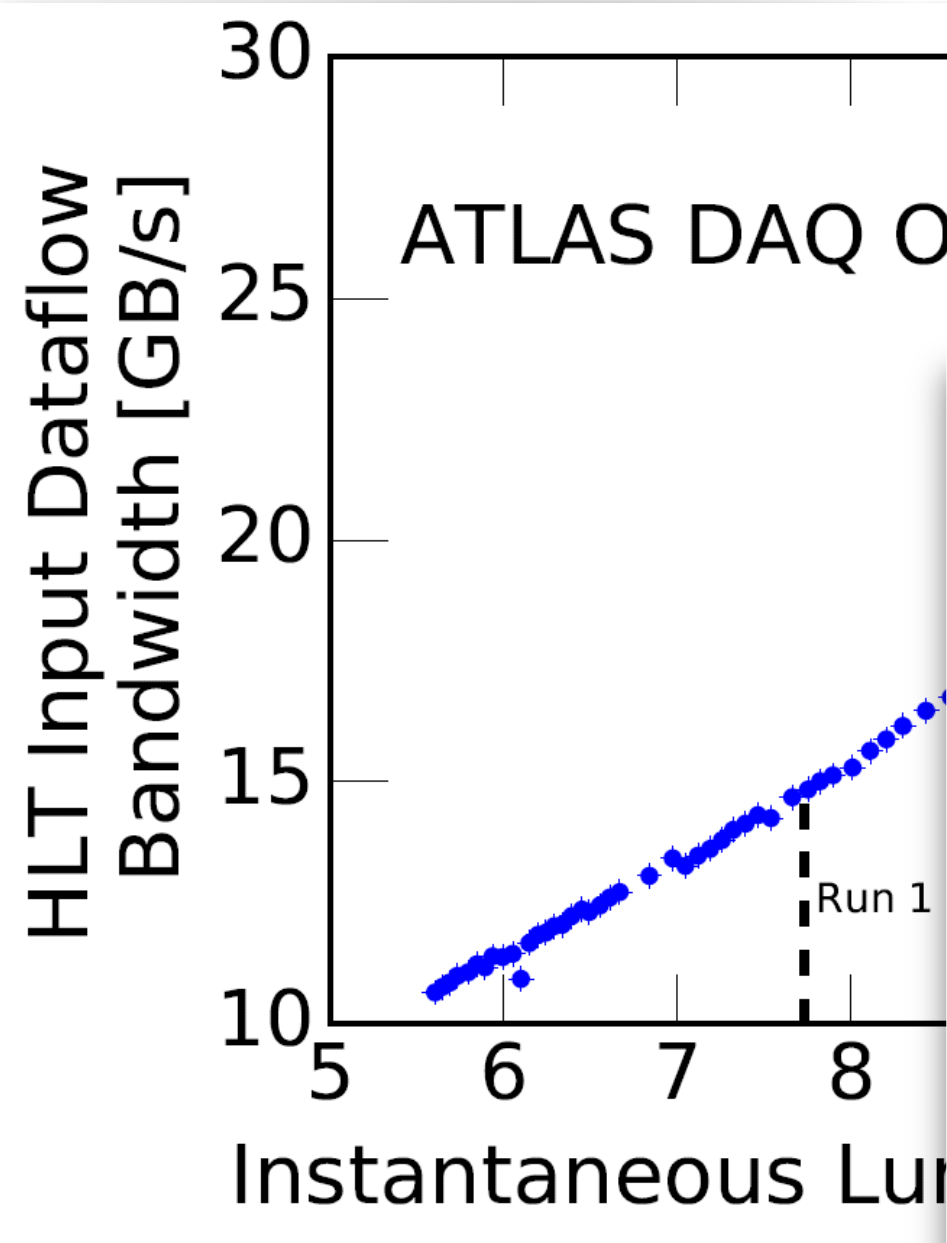➡ **Mostly low $p_T$ particles due to low transfer energy interactions**

HL-LHC tt̄ event in ATLAS ITK
at <μ>=200

**Physics program for the future
is towards more rare processes
at the same energy scale**
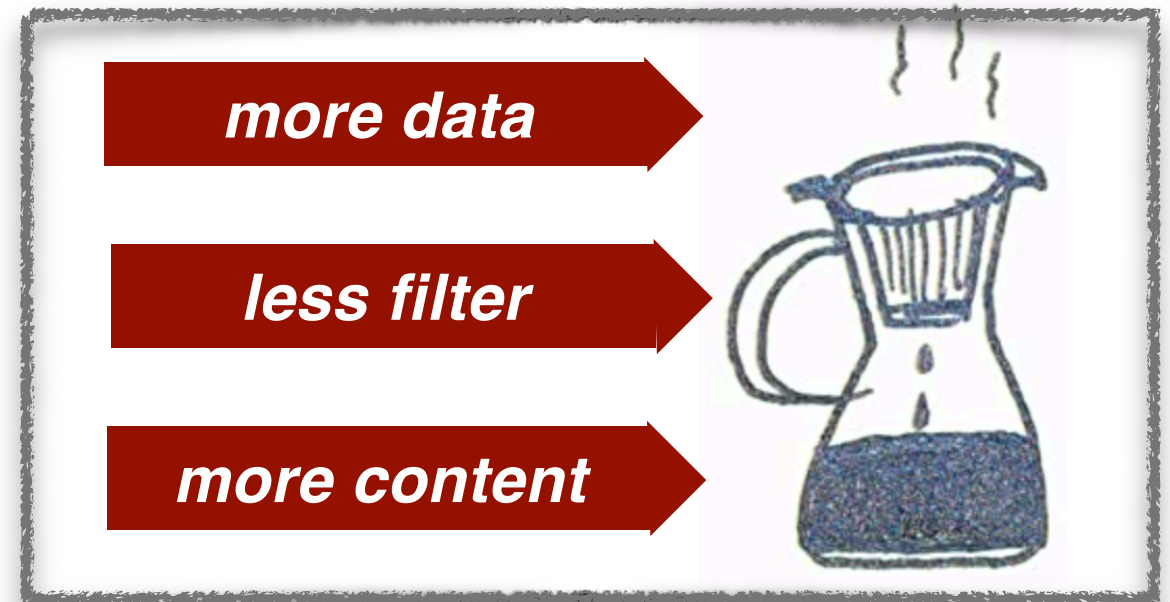
**Very large uncertainties to take into account!**

# ADDITIONAL COMPLICATION AT HL-LHC

## Luminosity x10, complexity x100: we cannot simply scale current approach
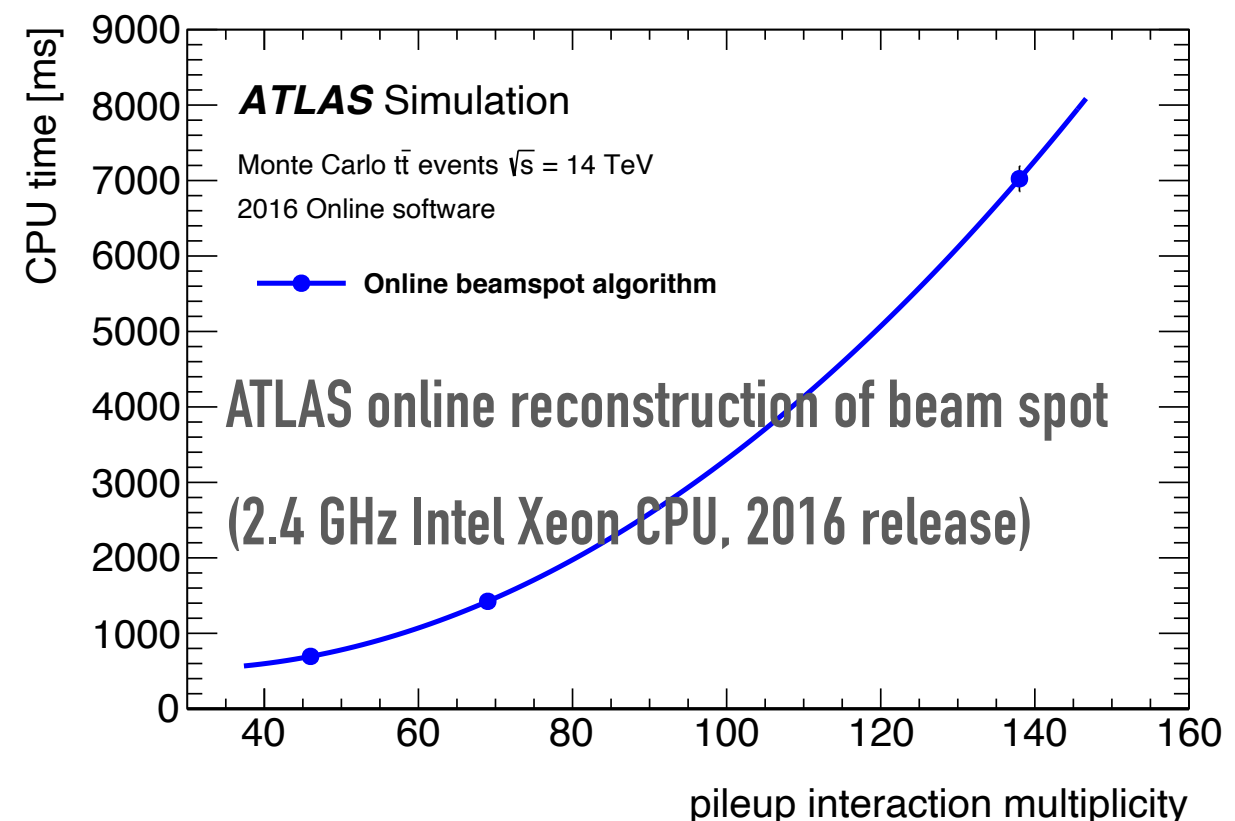
## x10 higher Luminosity means...

➡ **More interactions per BC (pile-up)**
  - ➡ Less rejection power (worse pattern recognition and resolution)
  - ➡ Larger event size

➡ **Larger data rates:**
  - ➡ FE readout rate @L1: 0.1 ➡ 1 MHz
  - ➡ DAQ throughput:  1 ➡ 50 Tbps

*ATLAS/CMS numbers*

**more data**

**less filter**

**more content**

## But cannot...

➡ **Increase trigger thresholds**
  - ➡ Need to maintain physics acceptance
➡ **Scale dataflow with Luminosity**
  - ➡ **H/W**: more parallelism ➡ more links ➡ more material and cost
  - ➡ **S/W**: processing time not linear ~ L

**ATLAS** Simulation

Monte Carlo $t\bar{t}$ events $\sqrt{s}$ = 14 TeV
2016 Online software

— Online beamspot algorithm

ATLAS online reconstruction of beam spot

(2.4 GHz Intel Xeon CPU, 2016 release)

CPU time [ms] — pileup interaction multiplicity

# THE REAL–TIME ADVENTURE

reduce latency

**Sequential Processing** — Single Core CPU — Single Core CPU Hyper-Threaded — Multi Core CPU — Graphics Processing Unit (GPU) → FPGAs → **Parallel Processing** — custom ASICs →

*Latency ranging from 100 to 2 μs*

LHCb
250 Eb/year

SKA
30000 Eb/year

Human Genome
8000 Eb/year

Exabytes ($10^{18}$ Bytes)!!

ATLAS/CMS
260 Eb/year

LHCb
1000 Eb/year

Global Internet
2800 Eb/year

| 2021 | 2023 | 2025 | 2027 | 2029 | 2031 |

*See Openlab workshop*

Trigger-less DAQ

Readout

Logic

Buffers

Tension between TDAQ architecture and FE complexity

High performance farms

Triggering detectors

**What we do?**



Trigger-less DAQ

high detector granularity

**Tension between TDAQ architecture and FE complexity**

Readout

Logic

Buffers

**High performance farms**

**Triggering detectors**

refine calibrations, as offline

complex ASIC logic

*LHCP-2022*

# BE SMARTER! INCREASE RESOLUTION FOR BETTER S/B

**What we do?**

**How?**

**Example**

**Tension between TDAQ architecture and FE complexity**

**Trigger-less DAQ**

Readout

Logic

Buffers

**High performance farms**

**Triggering detectors**

high detector granularity

high speed electronics/links

R&D on detectors Front-End

refine calibrations, as offline

large buffers, long latency

tight: offline=online (LHCb, ALICE)
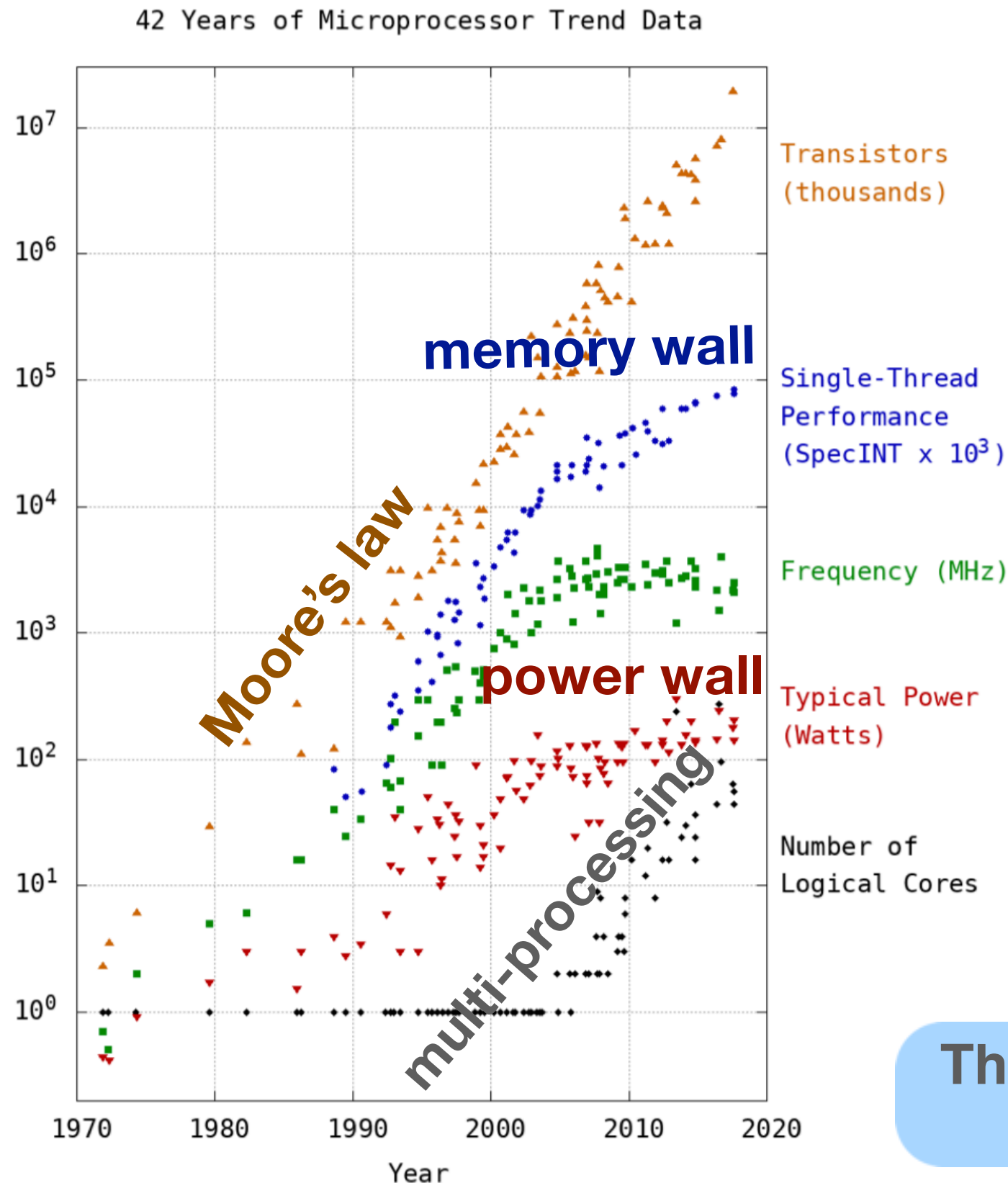soft: decouple trigger/DAQ (ATLAS, CMS)

complex ASIC logic

trigger-driven design

hardware track trigger (CMS)

42 Years of Microprocessor Trend Data

Data Source: https://github.com/karlrupp/microprocessor-trend-data

- ▸ **CPU frequencies are plateauing**
- ▸ **Local memory/core is decreasing**
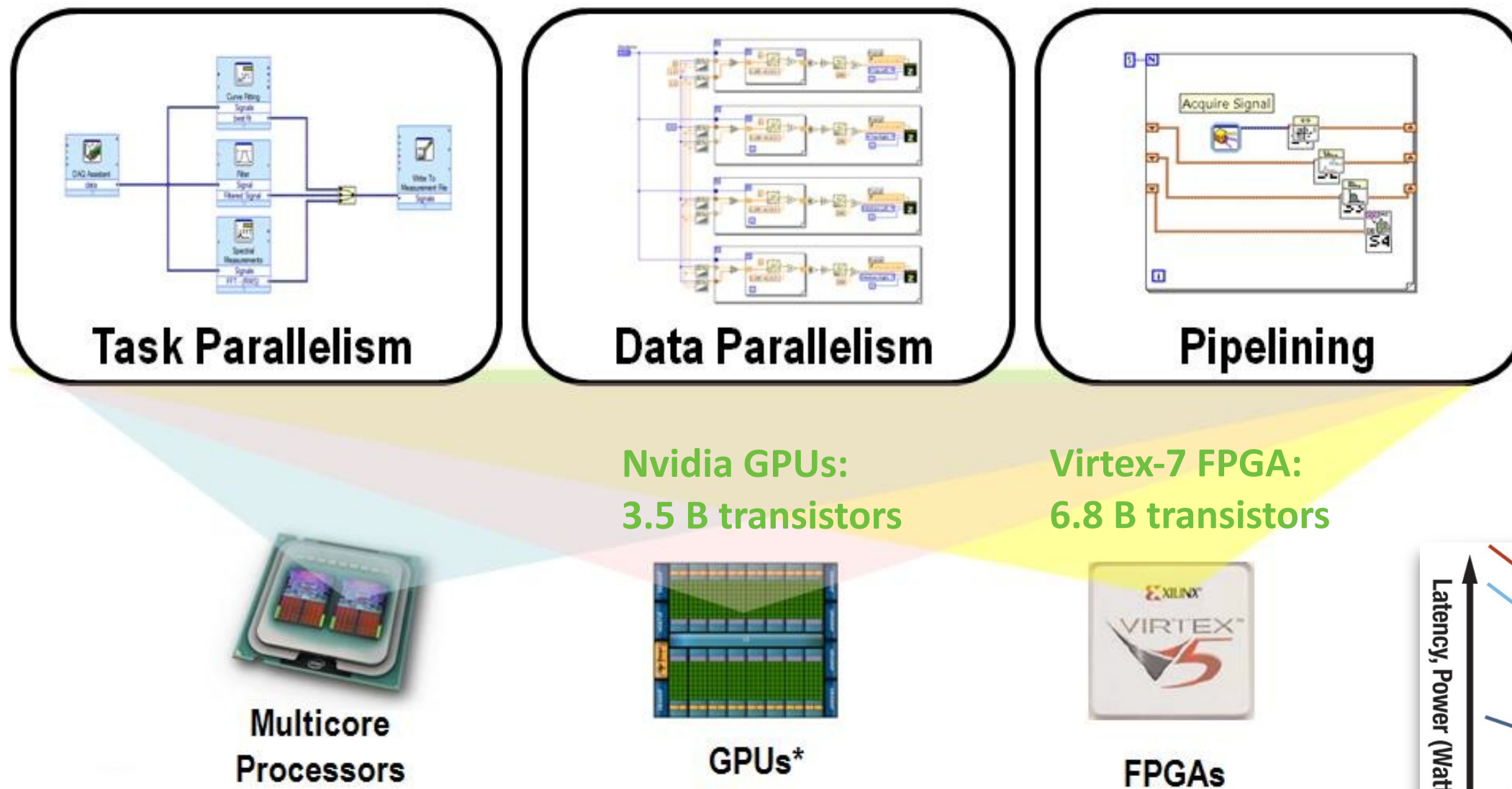- ▸ **Number of cores is increasing**

➡ **Exploiting CPU h/w, with more complicated programming**
  - ➡ Vectorisation, low-level memory…
➡ **Multithreading processing**
  - ➡ To reduce memory footprint
➡ **Use of co-processors:**
  - ➡ High Performance Computing (HPC) often employ GPU architecture to achieve record-breaking results!
➡ **Examples in LHC experiments:**
  - ➡ data reduction (ALICE & LHCb)
  - ➡ trigger selection (CMS & ATLAS)

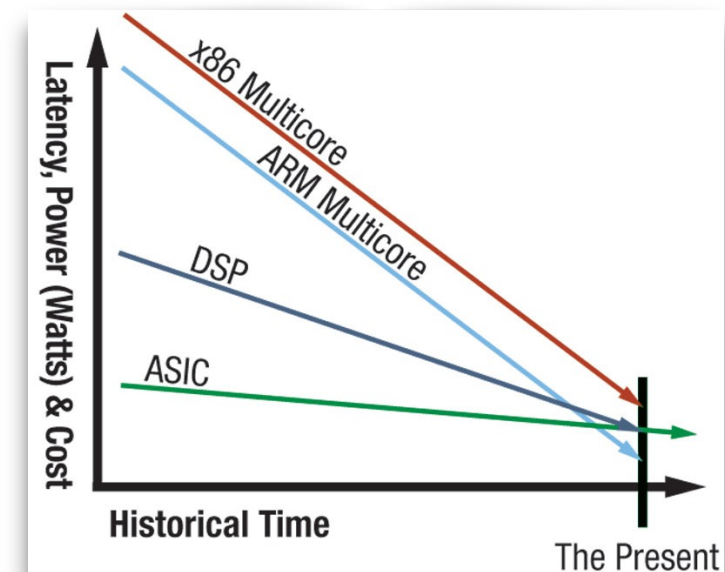**This requires fundamental re-write/ optimization of our software**

*Read: HPC computing*

**Task Parallelism**

**Data Parallelism**

**Pipelining**

Nvidia GPUs:
3.5 B transistors

Virtex-7 FPGA:
6.8 B transistors

**Multicore Processors**

**GPUs***

**FPGAs**

(*) Access to the nVIDIA® GPUs through the CUDA and CUBLAS toolkit/library using the NI LabVIEW GPU Computing framework.

Latency, Power (Watts) & Cost

x86 Multicore

ARM Multicore

DSP

ASIC

Historical Time

The Present

**The right choice can be combining the best of both worlds by analysing which strengths of FPGA, GPU and CPU best fit the different demands of the application**

➡ **Scientific computing is the third paradigm, complementing theory and experiment**
  ➡ Global scientific facilities (e.g., LIGO, LHC, Vera Rubin Observatory, the Square Kilometer Array)
➡ **Future trends in HPC focusing on:**
  ➡ Rise of massive scale commercial clouds (Google Kubernetes, serverless computing,….)
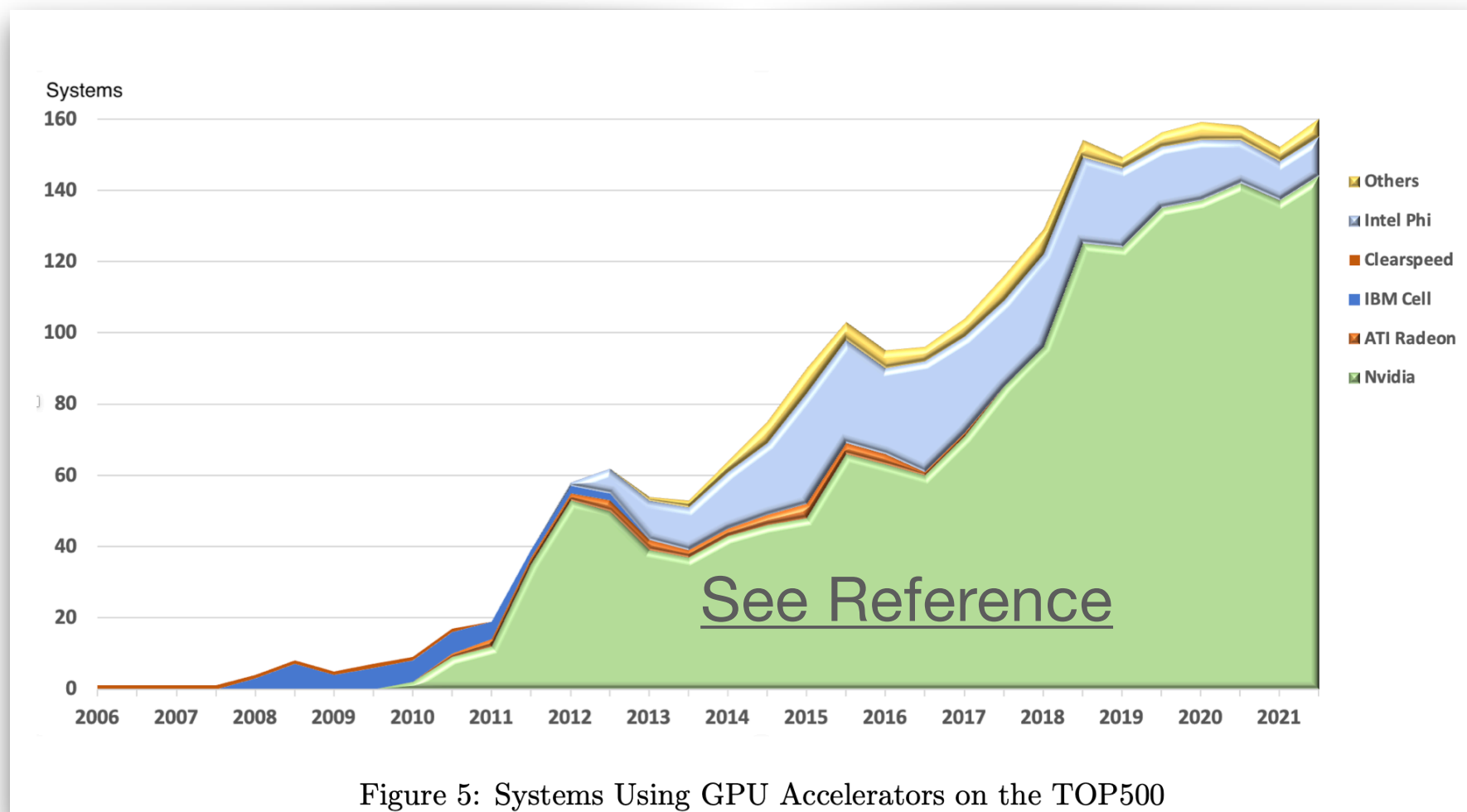  ➡ Evolution of semiconductor technology (chip size and packaging, see Amazon Graviton 3)
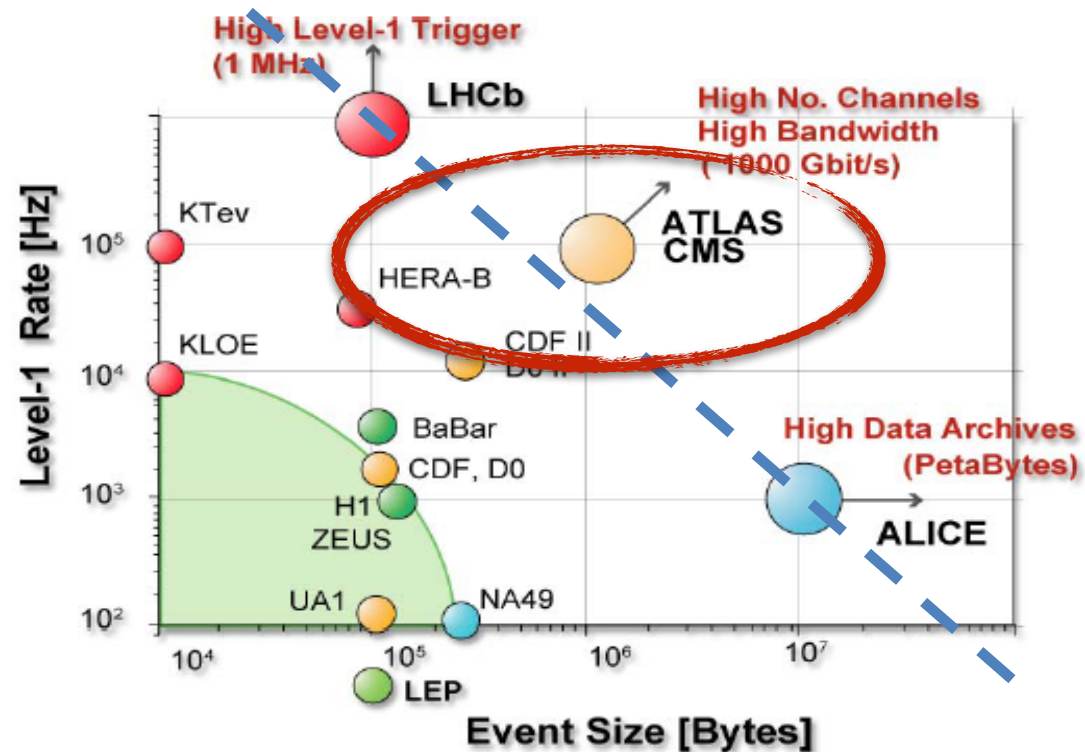


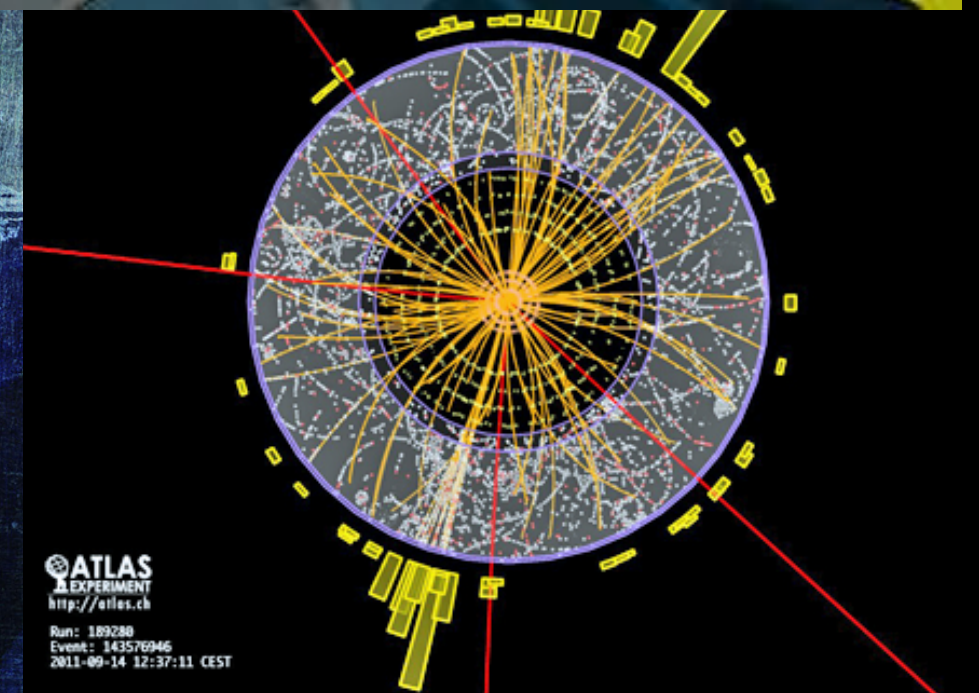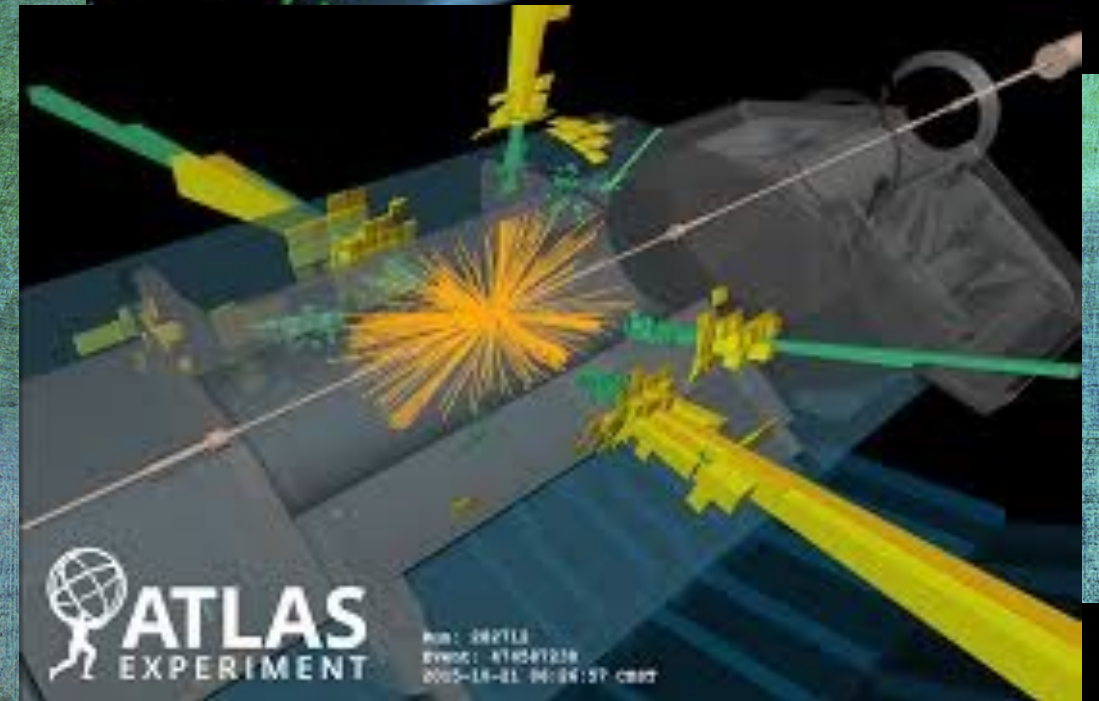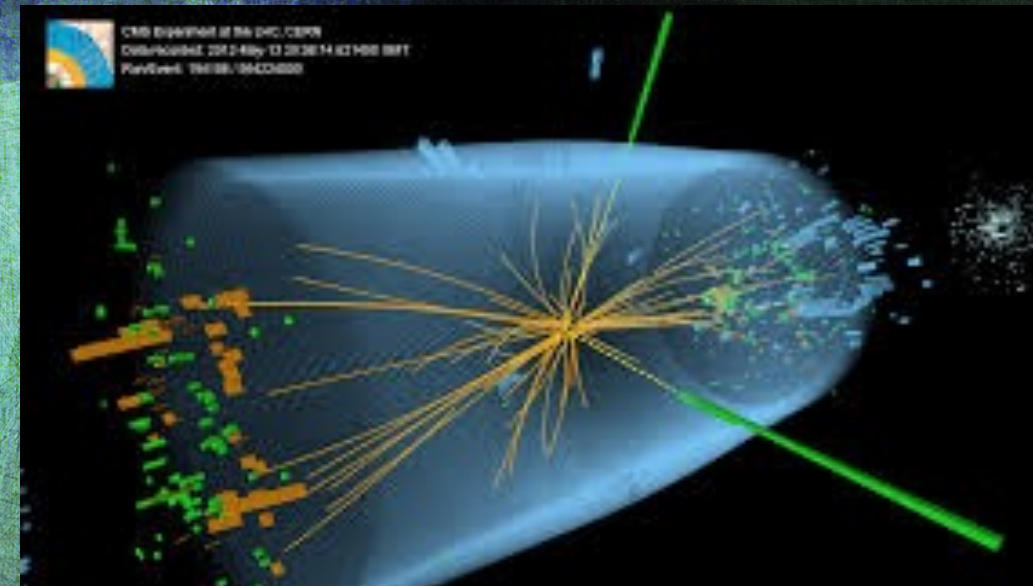Figure 5: Systems Using GPU Accelerators on the TOP500

**TOP500 today largely examples of a commodity monoculture: nodes with server-class microprocessors + GPUs**

# ATLAS AND CMS

*Studying the Standard Model at the high energy frontier*

## Same physics plans, different competitive approaches for detectors and DAQ

➡ **Same trigger strategy and data rates**

**1 MB * 100 kHz= 100 GB/s readout network**



*inclusive trigger selections*



ATLAS

CMS

➡ **Different DAQ architectures**
- ➡ **ATLAS**: minimise data flow bandwidth with multiple levels and regional readout
- ➡ **CMS**: large bandwidth, invest on commercial technologies for processing and communication

## Cannot do Event Building at 100 kHz

**CMS DAQ-1**

100 GB/s readout network in 2 steps

100 kHz Event Building factorised x8

*2 EB networks in blu*

*Filter network in green*



Myrinet (data concentrator)

1GB/s Ethernet (event builder)

➡ **Bet on exponential growth of technologies (networking/processing)**

➡ **Scalable and modular**

  ➡ Independent development of two network technologies

**Run-1 (as from TDR, 2002)**

➡ Myrinet + 1GBEthernet

➡ 1-stage building: 1200 cores (2C)

➡ HLT: ~13,000 cores

➡ 18 TB memory @100kHz: ~90ms/event

## Run 1: 100 GB/s network

**Myrinet widely used when DAQ-1 was designed**

➤ high throughput, low overhead
➤ direct access to OS
➤ flow control included
➤ new generation supporting 10GBE

## Run 2: 200 GB/s network

➤ Increased event size to 2MB
➤ Technology allows single EB network (56 Gbps FDR Infiniband)
➤ Myrinet —>10/40 Gbps Ethernet

Top500.org share by interconnect family

Myrinet

Custom

1 Gb/s Ethernet

10 Gb/s Ethernet

Share (%)

Infiniband

2002    2014    2018

## Choose best prize/bitps!

**Event size up to 1MB**

**Event size up to 2MB**

**100 kHz L1 rate**

**Myrinet**

**1 Gb/s Ethernet**

**100 GB/s 8 slices**

**CMS DAQ 1**

13000 core, 1260 host filter farm

max. 1.2 GB/s to storage

**100 kHz L1 rate**

**10/40 Gb/s Ethernet**

**56 Gb/s Infiniband**

**~200 GB/s**

**1 slice**

16000+ core, 900 host filter farm

**CMS DAQ 2**

~ 3-6 GB/s to storage

**HLT selections based on <u>regional readout and reconstruction</u>, seeded by L1 trigger objects (RoI)**



RoI=Region of Interest

➡ **Total amount of RoI data is minimal: a few % of the Level-1 throughput**

➡ one order of magnitude smaller readout network …

➡ … at the cost of a higher control traffic and reduced scalability

**Overall network bandwidth: ~10 GB/s    (x10 reduced by regional readout)**

Run 3



**complex data router to forward different parts of the detector data, based on the trigger type**

**Silicon tracking systems provide incredibly high resolution, crucial for controlling rates**



Vertex Reconstruction

Jet Tagging

Pile-up Removal

Missing Energy Reconstruction

## *Tracking challenges*

- ➤ Readout ~800M channels, ~50 Tbps
- ➤ Combinatorics ($10^4$ hits/BC)

**combinatorics scales like $L^N$**

L=luminosity,    N=number of layers

**Tracking reconstruction not feasible @40MHz, nor in few microseconds**

| | ATLAS [1] | CMS [2] |
|---|---|---|
| *data reduction @40MHz* | regions from L1 (RoIs) | stubs from hw coincidences |
| *track finding @1MHz* | Studying best algorithms to run in FPGAs and/or in GPUs | |
| *track fit @1MHz* | | |
| *precision tracking @100kHz* | optimized offline | optimized offline |



**stubs in CMS PT modules**

# LHCb, THE B-MESON OBSERVATORY

.................................................

*The lightest experiment to study the heavy b-quark*

http://lhcb-public.web.cern.ch/lhcb-public/

# LHCB TRIGGER STRATEGY

**LHCb 2012 Trigger Diagram**

**40 MHz bunch crossing rate**

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

| 450 kHz h± | 400 kHz μ/μμ | 150 kHz e/γ |
|---|---|---|

**Software High Level Trigger**

29000 Logical CPU cores

Offline reconstruction tuned to trigger time constraints

Mixture of exclusive and inclusive selection algorithms

**5 kHz (0.3 GB/s) to storage**

| 2 kHz Inclusive Topological | 2 kHz Inclusive/ Exclusive Charm | 1 kHz Muon and DiMuon |
|---|---|---|

**Input rate**

**L0 trigger**

**High Level**

**Low input rate and occupancy**

- ✦ Limited acceptance: 10 MHz
- ✦ Limited Luminosity = $2 \times 10^{32} cm^{-2}s^{-1}$

- ✦ Select Bs in hadronic triggers
- ✦ Reject complex/busy events

**60kB * 1MHz = 60 GB/s readout network**

- ✦ Multitude of exclusive selections

## LHCb 2015 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

**150 kHz**

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

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

**12.5 kHz Rate to storage**

**Can increase efficiency on B-hadrons? YES, use more precision!!**

**Real-time calibration and alignments**

**Synchronous with DAQ**

HLT-1

✦ Use tracks for selections on B-decay vertices (in 35ms)

**Split with a large buffer (4PB)!**

HLT-2

**Deferred Processing**

✦ Reconstruct with offline-like calibrations (in 350ms), becoming real-time physics analysis

prompt charm production cross-sections from LHCb turbo stream in Run2

## Can we get rid of FrontEnd raw data?

➡ **Event size/10 -> x10 rate, for free**

➡ **Tested on dedicated data streams in many experiments:**

    ➡ Full online reconstruction (**LHCb**)

    ➡ Data scouting (**ATLAS/CMS**)

        ➡ for some high rate signatures, save only reduced information

➡**Main data stream for LHCb & ALICE upgrade**

    ➡ **and be a guidance for all other experiments**



di-jet mass spectrum from CMS data-scouting in Run2

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

*40Tbit/s*

**30 MHz**

**Software High Level Trigger**

Full event reconstruction, inclusive and exclusive kinematic/geometric selections

HLT-1

*1-2 Tbit/s*

**1MHz**

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

HLT-2

**2-5 GB/s to storage**

**50 kHz**

*80 Gbit/s*

**FE readout & Event Building at 30 MHz (~40 Tbit/s)**

**Key strategy: reduce data size at FE and suppress pileup with tracking**

### Tracking at ~30 MHz ?

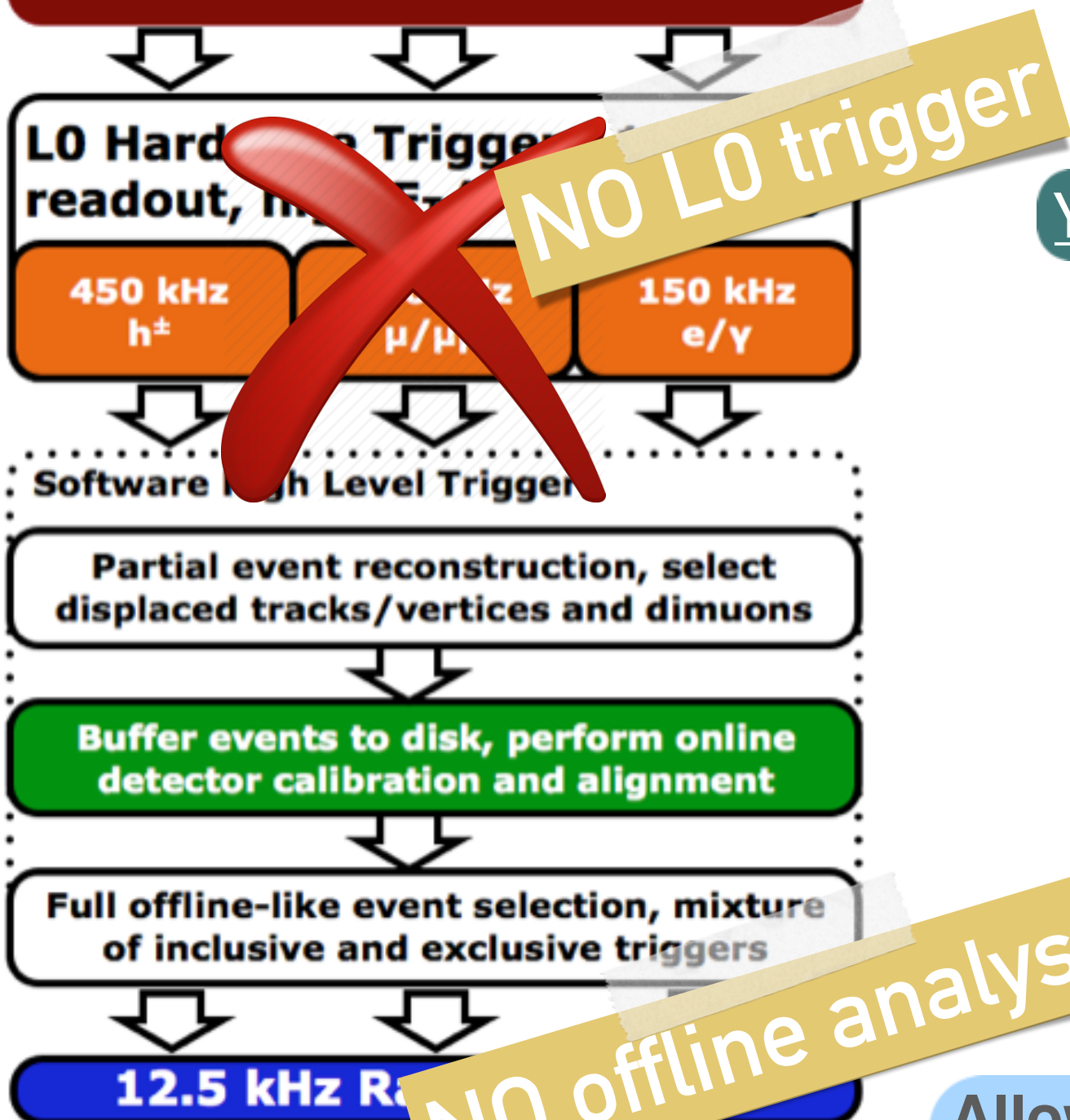✦ Run2: ~ 100k cores < 6 ms

✦ Run3: modern CPU & co-processors (FPGA/GPU)

Scintillating Fibre Tracker

VELO

Upstream Tracker

**Online Tracking**

Velo tracking

Velo-UT tracking
$p_T > 200$ MeV, $\delta p/p \sim 15\%$

Forward tracking
$p_T > 500$ MeV, $\delta p/p \sim 0.5\%$

PV finding

Rate reducing cuts
Output < 1 MHz

Muon Identification

Simplified Kalman fit

Particle Identification

arXiv:2105.04031

$$150kB \times 30MHz = 40Tbs$$

**Readout @ 30 MHz**
**Event size ~ 150kB**

➡ **Data reduction:**
  ➡ Custom FPGA-card (PCIe40) also used in ALICE
  ➡ Data-packing for sub-detectors (zero-suppression, clustering)

➡ **Massive link usage:**
  ➡ ~10,000 GBT (4.8 Gb/s, rad-hard)

**DAQ network < 40 Tbit/s**
**Record rate: <100 kHz**

**PCIe-gen3:** simple protocol, large bandwidth
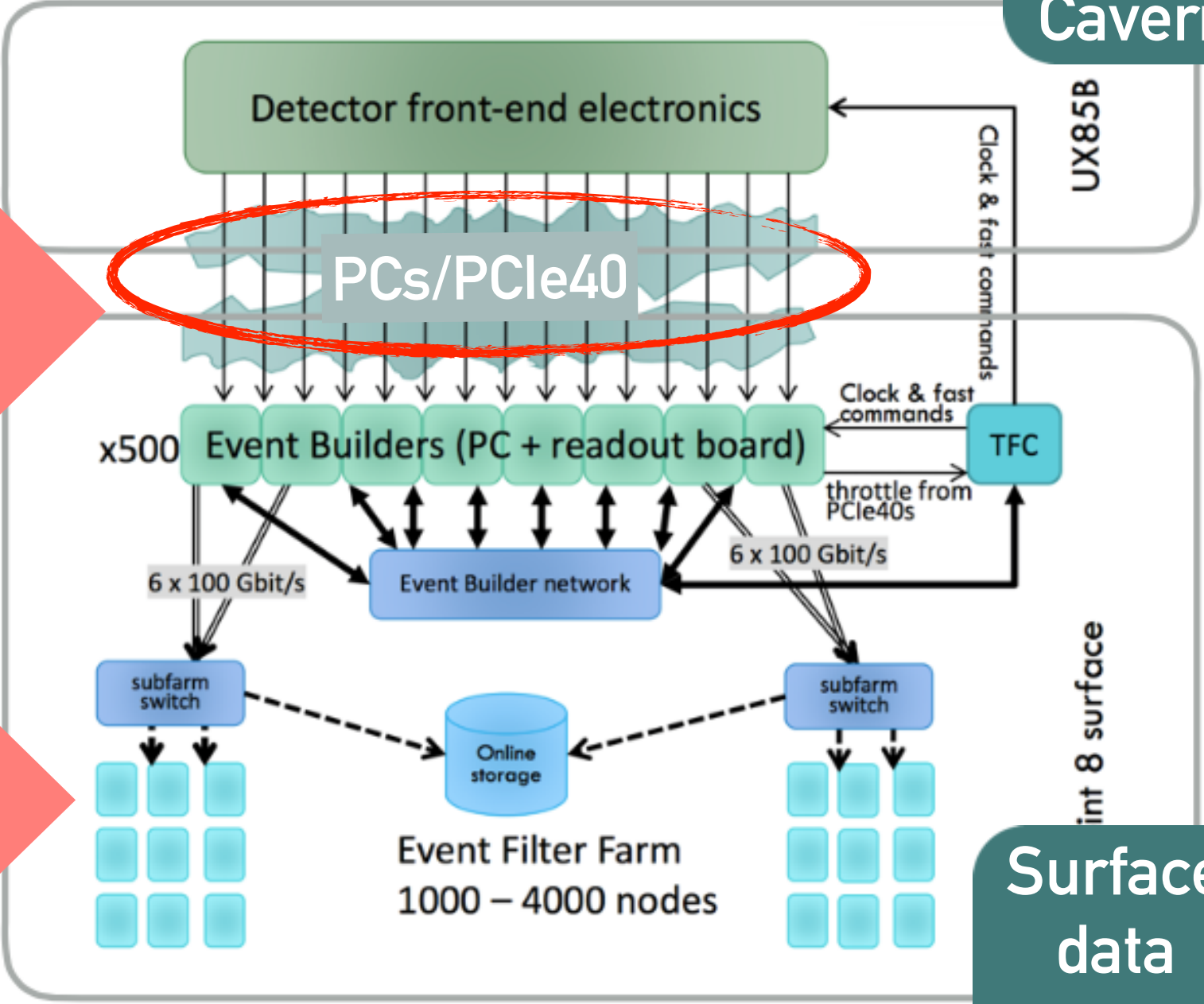**PCIe:** maximum flexibility in later networking choice



Inside Cavern

Surface data centre

UX85B

Clock & fast commands

Detector front-end electronics

PCs/PCIe40

x500 Event Builders (PC + readout board)

Clock & fast commands

TFC

throttle from PCIe40s

6 x 100 Gbit/s

Event Builder network

6 x 100 Gbit/s

subfarm switch

Online storage

subfarm switch

Event Filter Farm
1000 – 4000 nodes

*Ref for PCIe40*

**Large farm of equal nodes with 8 PCIe40 boards, specialised by firmware**



➡ **EB network is oversized: able to manage 64Tb/s (320 network cards x 200Gb/s)**

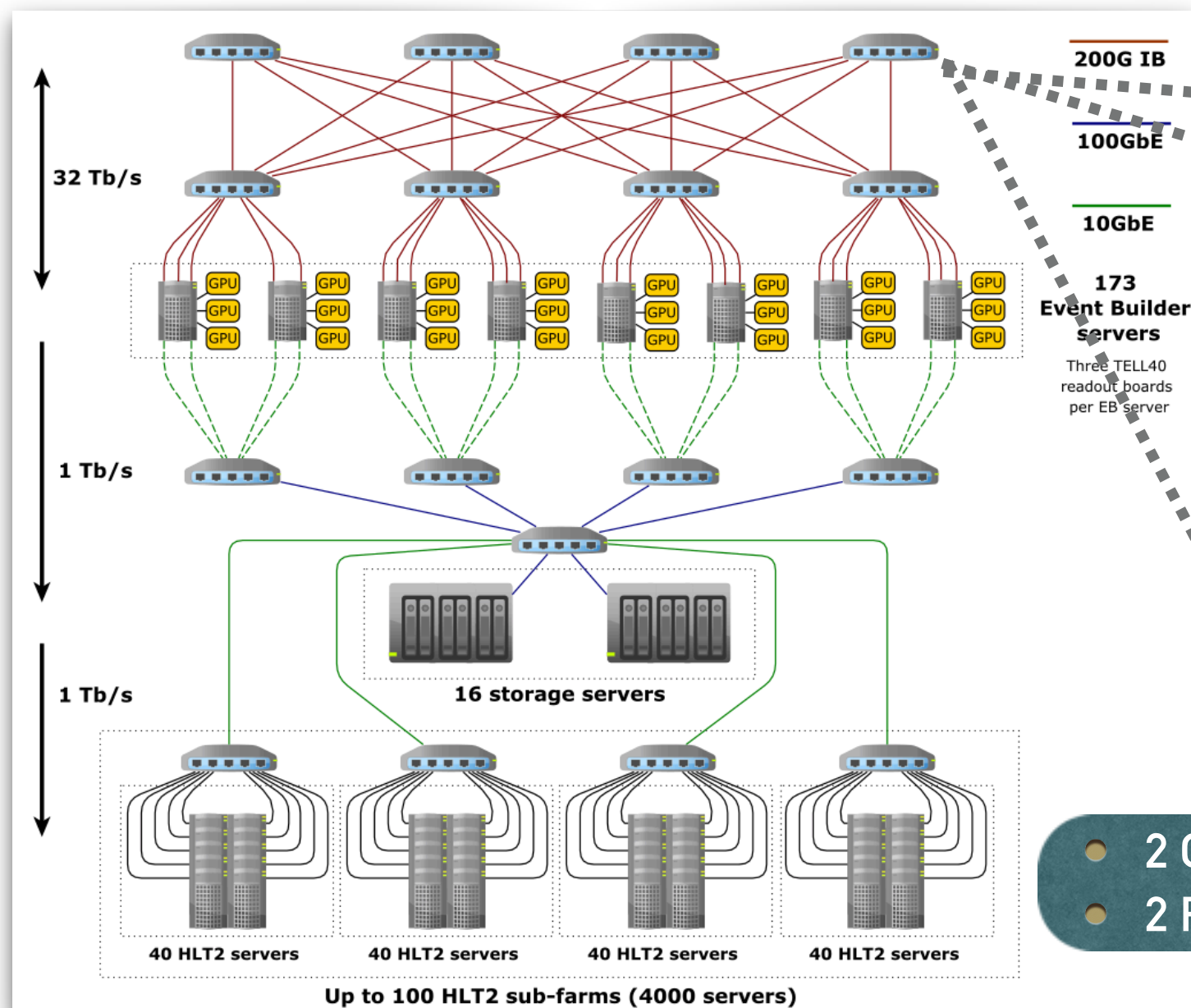➡ **Large rejection at HLT1: use O(200) GPU! throughput at ~100kHz**

➡ **Storage Buffer HLT1-HLT2 = 40 PB (3000 hard-disks) enough for days**

    ➡ SSD faster but have short lifetime wrt high read-write rate, so prefer hard-disks

## Large farm of equal nodes with 8 PCIe40 boards, specialised by firmware



*One node*

- 2 CPUs with large RAM (up to 512 GB!)
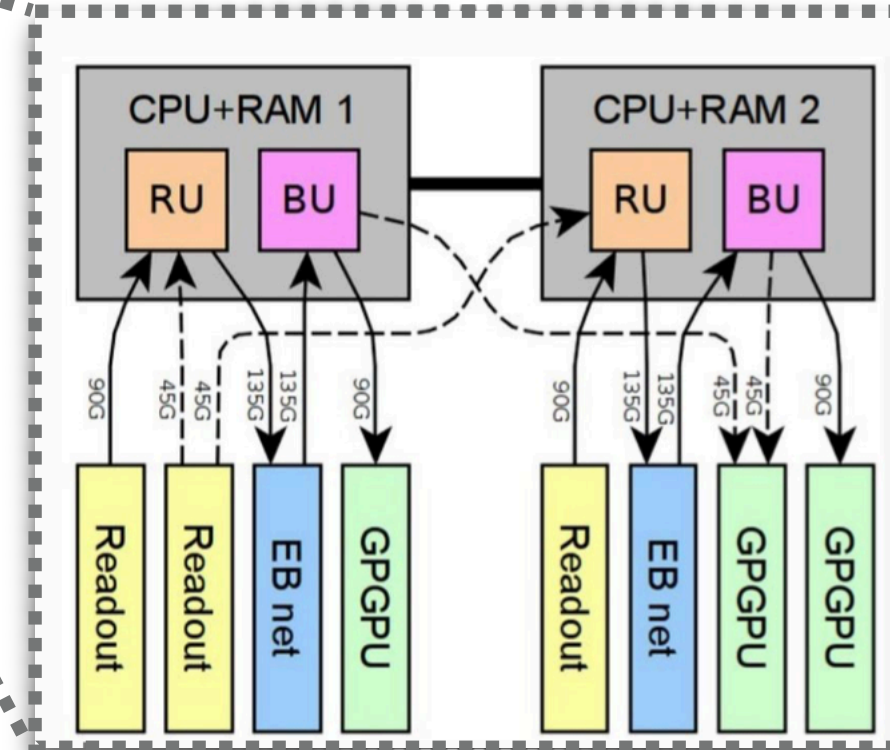- 2 RU, 2 BU, 2 infiniband NIC (200 Gb/s), 1-3 GPUs

➡ **EB network is oversized: able to manage 64Tb/s (320 network cards x 200Gb/s)**

➡ **Large rejection at HLT1: use O(200) GPU! throughput at ~100kHz**

➡ **Storage Buffer HLT1-HLT2 = 40 PB (3000 hard-disks) enough for days**

   ➡ SSD faster but have short lifetime wrt high read-write rate, so prefer hard-disks

Same data volume as ATLAS/CMS HL-LHC upgrades! But earlier and for less money

# ALICE: THE SMALL BIG-BANG
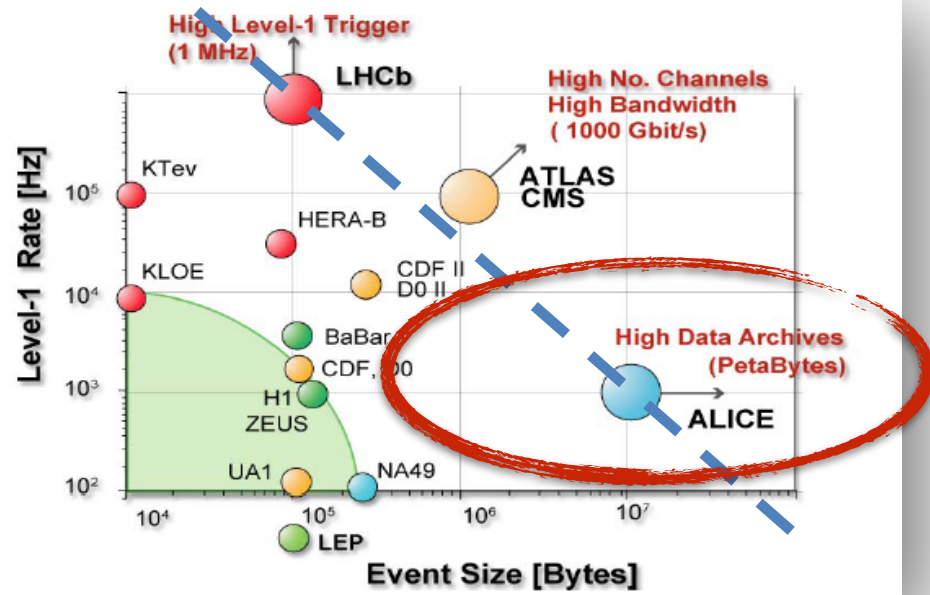
*Recording heavy ion collisions*

*http://alice-daq.web.cern.ch*

➡ **19 different detectors**

➡ With **high-granularity and timing information**

  ➡ in particular the Time Projection Chamber (**TPC**) has very high occupancy, and slow response

➡ **Large event size (> 40MB)**

  ➡ TPC producing 90% of data

➡ **Complex event topology**

  ➡ low trigger rate: max 3.5 kHz

cms = 5.5 TeV per nucleon pair
Pb–Pb collisions at L = $10^{27}$ cm$^{-2}$s$^{-1}$

➡**Challenges for TDAQ design:**

  ➡ detector readout: up to ~50 GB/s

  ➡ storage:  1.2 TB/s (Pb-Pb)

➡ **Dataflow with local (LDC) and global (GDC) data concentrators**

  ➡ Detector readout (~20 GB/s) with point-to-point optical links (DDL, max 6Gb/s)

  ➡ Rate to the LDCs can go above 13 GB/s

➡ **Transient Data Storage (TDS)**

  ➡ Before the Permanent Data Storage (PDS) and publish via the Grid

➡ **LHC heavy ion programme: <u>extend statistics by x100!</u>**
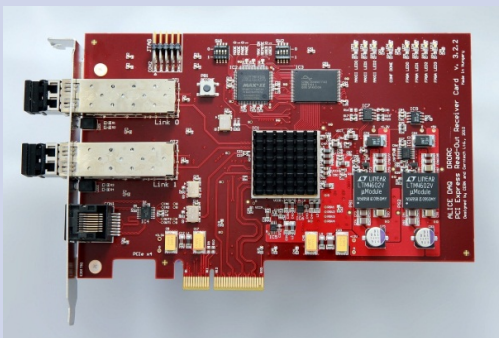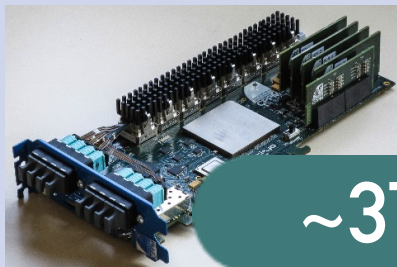
- ➡ Increase detector granularity (===> increase event size!)
- ➡ Increase storage bandwidth x O(100)
  - ➡ Offline reconstruction also challenging due to combinatorics
- ➡ Increase readout rates ~kHz ➡ 50 kHz (===> need new and faster electronics)
  - ➡ Rate very close to TPC readout !!

**New TDAQ challenges!**

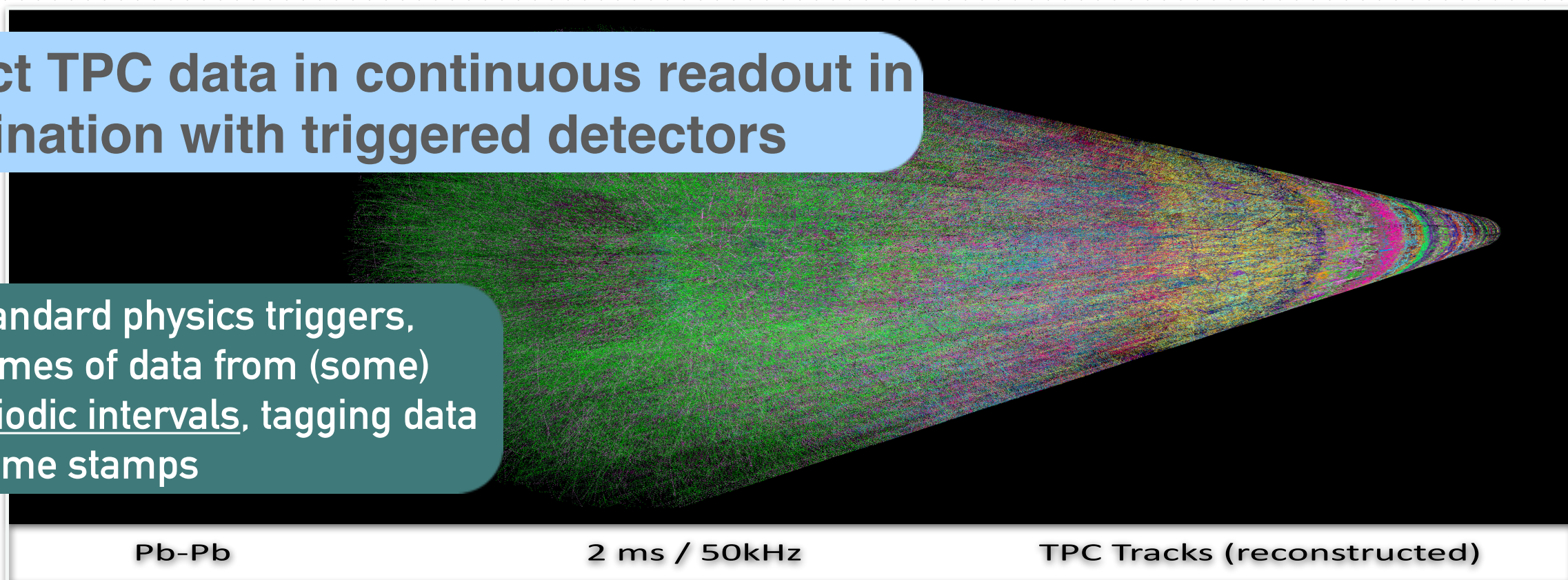| RORC 1 | C-RORC | CRU |
|---|---|---|
|  |  |  |
| 2 ch @ 2 Gb/s<br>PCIe gen.1 x4 (1 GB/s) | 12 ch @ up to 6 Gb/s<br>PCIe gen.2 x 8 (4 GB/s) | 24 ch @ 5 Gb/s<br>PCIe gen.3 X 16 (16 GB/s) |
| Custom DDL protocol | Custom DDL protocol<br>(same protocol but faster) | GBT |
| Protocol handling<br>TPC Cluster Finder | Protocol handling<br>TPC Cluster Finder | Protocol handling<br>TPC Cluster Finder<br>Common-Mode correction<br>Zero suppression |

**~3TB/s detector readout**

**New Common Readout Unit (CRU), based on PCIe40 card**

Run 1 | LS1 | Run 2 | LS 2 | Run 3

**Reconstruct TPC data in continuous readout in combination with triggered detectors**

In addition to standard physics triggers, DAQ collects frames of data from (some) detectors at <u>periodic intervals</u>, tagging data internally with time stamps



Pb-Pb      2 ms / 50kHz      TPC Tracks (reconstructed)

➡ **Heart Beat (HB) issued in continuous & triggered modes**

    ➡ subdivision of data into time intervals to allow synchronisation between different detectors

    ➡ 1 per LHC orbit, 89.4 $\mu$s: <u>~10 kHz</u>

➡ **Grouped in Time-Frames:**

    ➡ 1 every ~20 ms: <u>**~50 Hz**</u> (1 TF = ~256 HBF)

**CRU (& frontend)**

Time

**Heart Beat Frames (HBF):** data stream delimited by two HBs    Trigger data fragments

**FLP**

**Sub-Time Frame (STF) in FLP 0:** grouping of (~256) consecutive HBFs from one FLP    **FLP 1**    **FLP n**

**EPN**

**Time Frame (TF):** grouping of all STFs from all FLPs for the same time period from triggered or continuously read out detectors

➡ **Data compression in GPUs and FPGAs ==> x2 readout rate**

➡ **Network evolution: 2.5GB/s (2010) ⇒ 6GB/s (2015) ==> x2 DAQ throughput**



**Tracking processing based on GPUs since Run1!**

**Higher rates with smaller data?** → **Store reconstruction, discard raw data**

**Very heterogeneous system**

➡ **Synchronous, with continuous data**
  ➡ Data compression in FPGA/CPU
  ➡ 30s to analyse 20ms-time frame

➡ **Asynchronous, reconstruction in GPUs**
  ➡ 250 EPN servers with 8 GPU-cards
  ➡ Require large-memory GPUs!

## O² system

➡ **Common online/offline software**
  ➡ Same calibrations and resources

**Data reduction**
**Calibration 0**

**Data aggregation**
**Reconstruction**
**Calibration 1**

**More reconstruction**
**Calibration 2**

Detectors electronics

*3.4 TB/s  (over 8500 GBTs links)*

Base Line correction, zero suppr.
Readout
Data aggregation
Local data processing

CRU/FPGA
CPU
**FLP**

*500 GB/s*

Data aggregation
Synchronous global
data processing

CPU    GPU
**EPN**

*90 GB/s*

Data storage (60 PB)
1 year of compressed data
Write 170 GB/s, Read 270 GB/s

*20 GB/s*

Asynchronous (hours)
event reconstruction with
final calibration
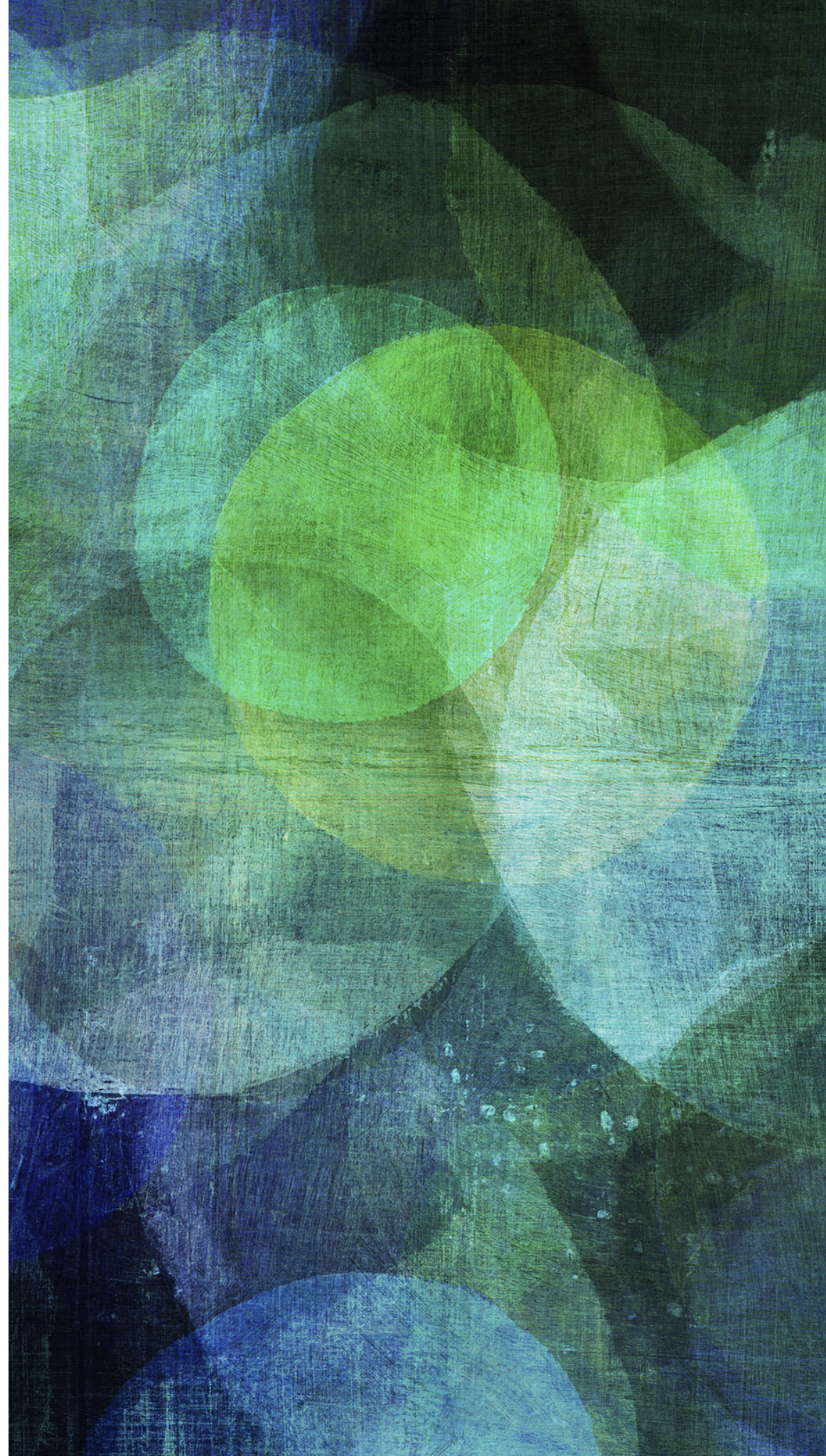
# SUMMARY OF THE SUMMARIES

➡ **LHC experiments are among the largest and most complex TDAQ systems in HEP, to cope with a very difficult environment (always top LHC Luminosity)**

➡ **Continuous upgrade following the LHC luminosity, with different approaches**

   ➡ **ATLAS/CMS** high-rate readout and Event Building, based on robust trigger selections

   ➡ **LHCb** pioneer online-offline merging with large data throughputs

   ➡ **ALICE** drives the GPU evolution and data compression

➡ **With a general trend, <u>towards higher bandwidths and comodity HW</u>**

   ➡ Scalability not obvious. Challenge remains for front-end and back-end technologies and efficient (cost, time, power) computing farms

   ➡ Moore's law still valid for processors but needs more effort to be exploited

➡ **Each experiment trying to gain advantage from others' developments**

   ➡ joined efforts already started for hardware/software

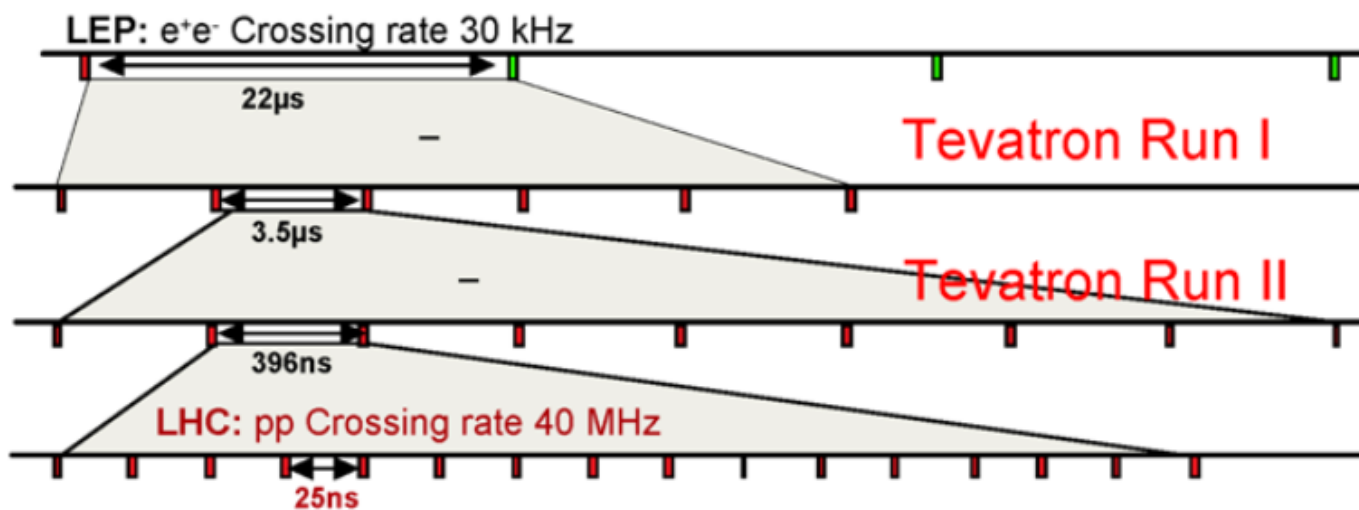   ➡ sometimes stealing ideas ("… but we can do better than that…")
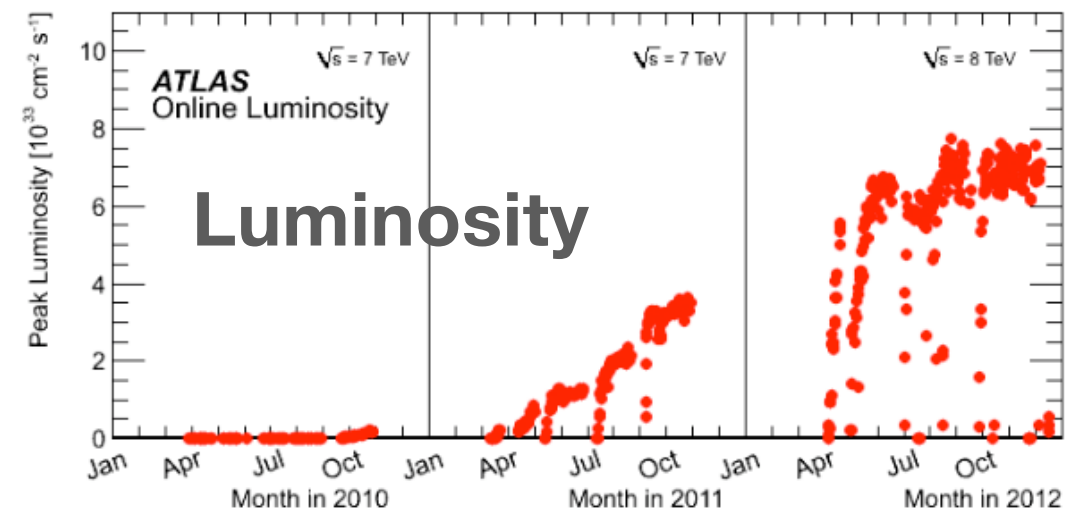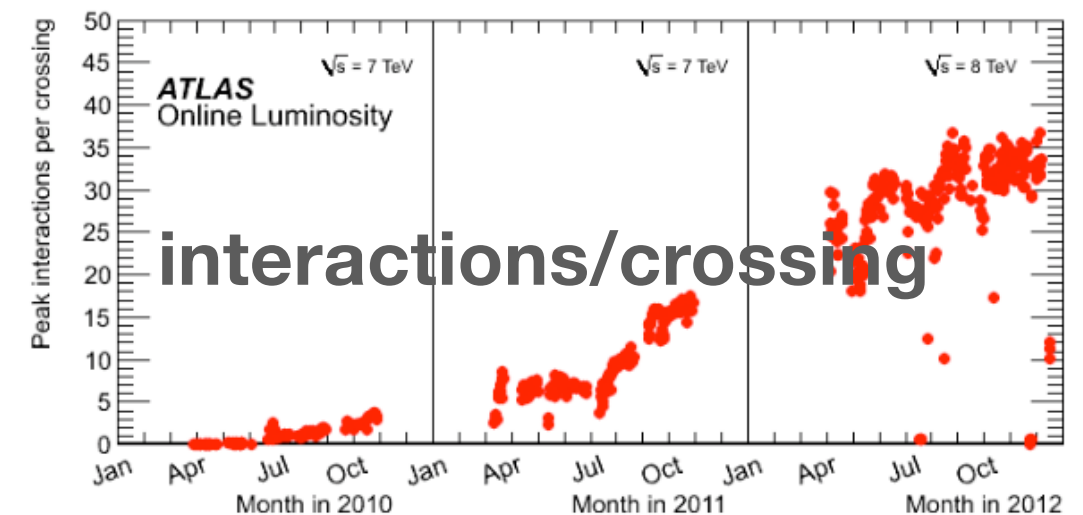
# BACK-UP SLIDES

## The clock source

- ➡ ~3600 bunches in 27km
- ➡ distance bw bunches: 27km/3600 = 7.5m
- ➡ distance bw bunches in time: 7.5m/c = 25ns



LEP: e⁺e⁻ Crossing rate 30 kHz
22μs
Tevatron Run I
3.5μs
Tevatron Run II
396ns
LHC: pp Crossing rate 40 MHz
25ns

**At full Luminosity, every 25ns, ~23 superimposed p-p interaction events**

## The pile-up source

- ➡ more collisions/bunch crossing: ~23 at design luminosity



**interactions/crossing**

**Luminosity**

➡ **Allow trigger decision longer than clock tick (and no deadtime)**

➡ Execute trigger selection in defined clocked steps (**fixed latency**)

➡ Intermediate storage in stacked buffer cells

➡ R/W pointers are moved by clock frequency

➡ **Tight design constraints for trigger/FE**

➡ **Analog/digital pipelines**

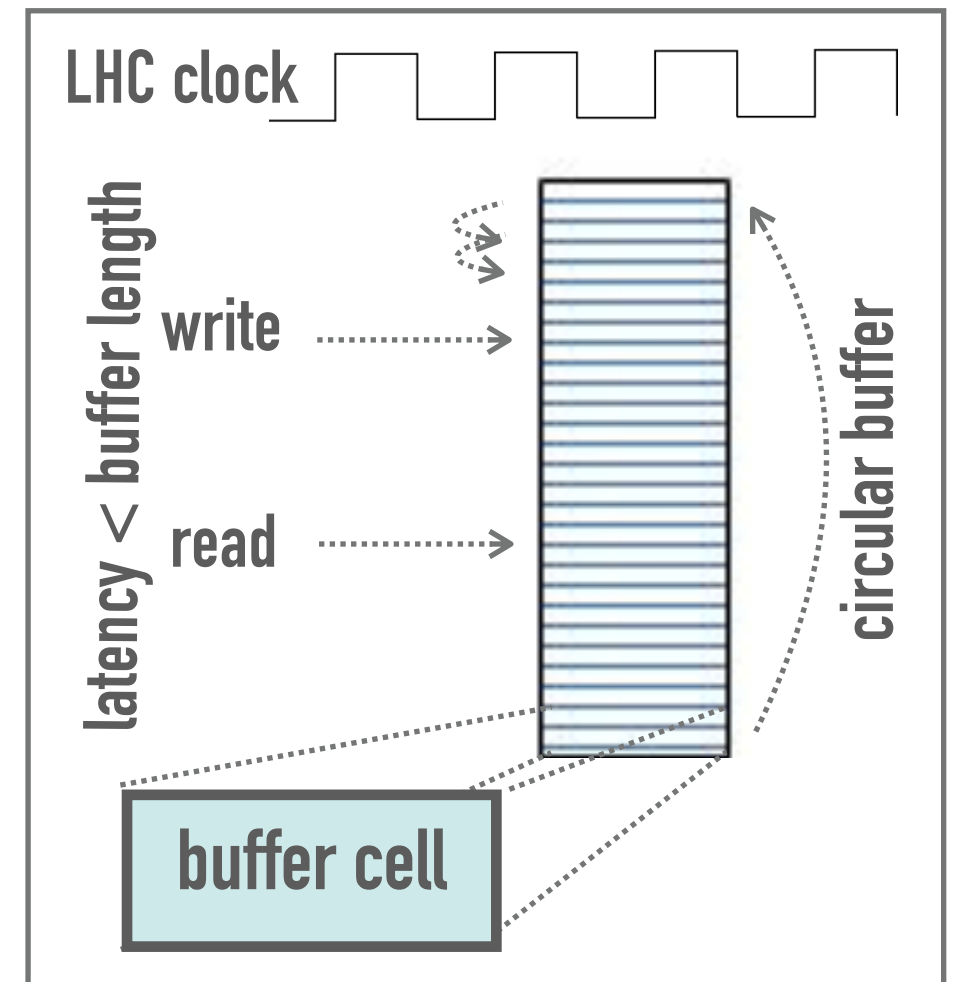➡ Analog: built from switching capacitors

➡ Digital: registers/FIFO/…

➡ **Full digitisation before/after L1A**

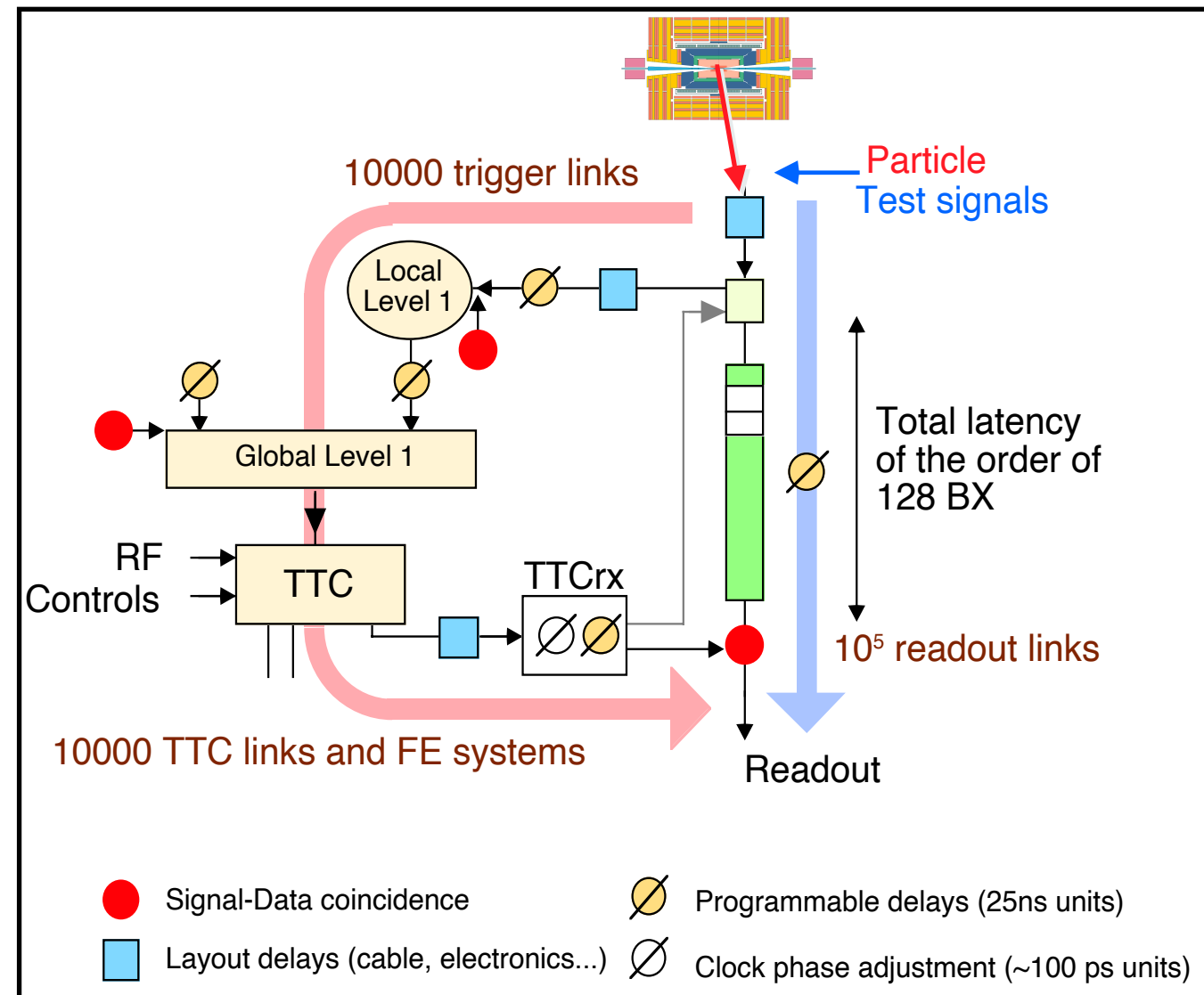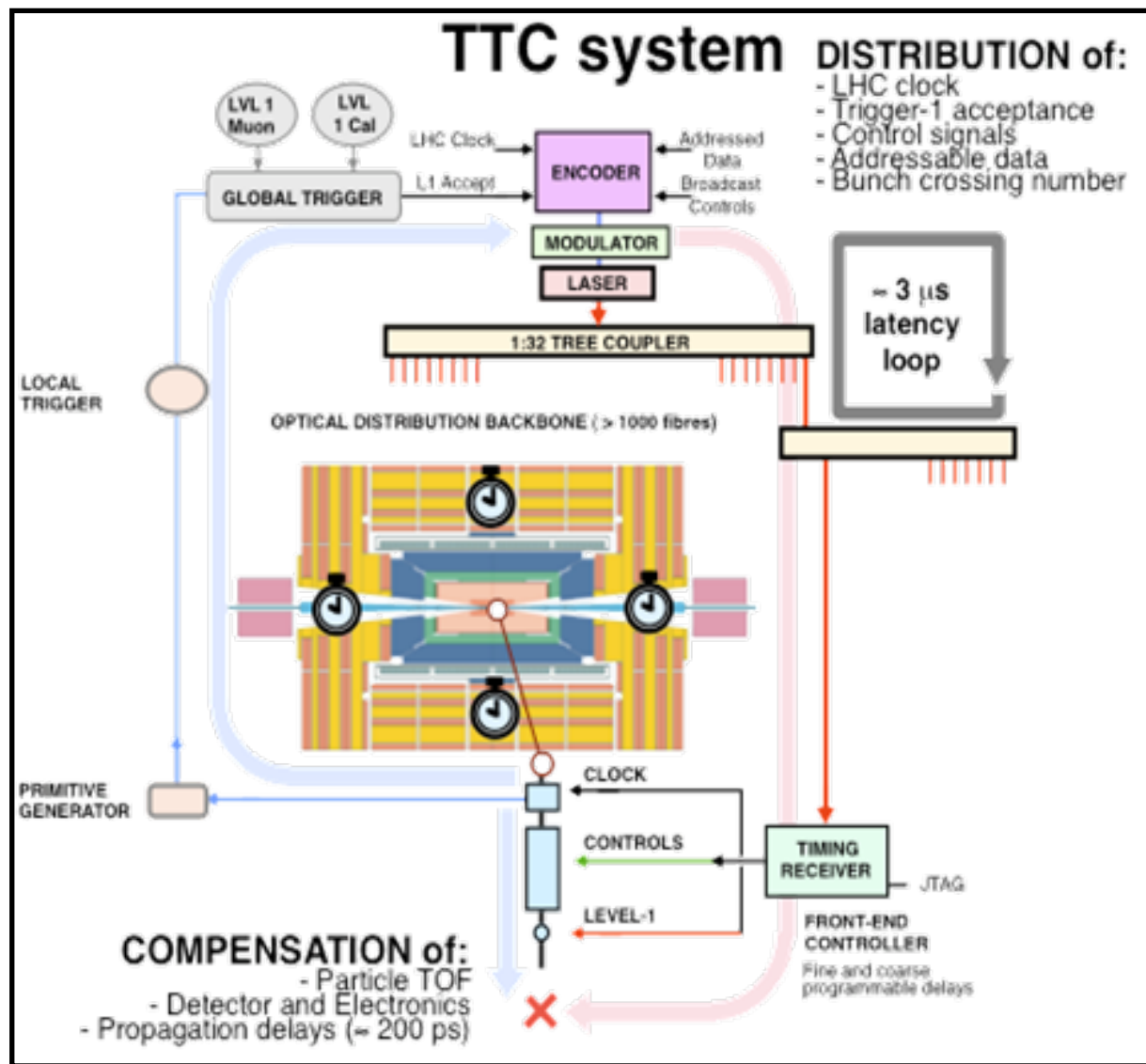➡ Fast DC converters (power consumption!)

➡ **Additional complication: synchronisation**

➡ BC counted and reset at each LHC turn

➡ large optical time distribution system



LHC clock

latency < buffer length

write

read

circular buffer

buffer cell

➡ **Common optical system: TTC**
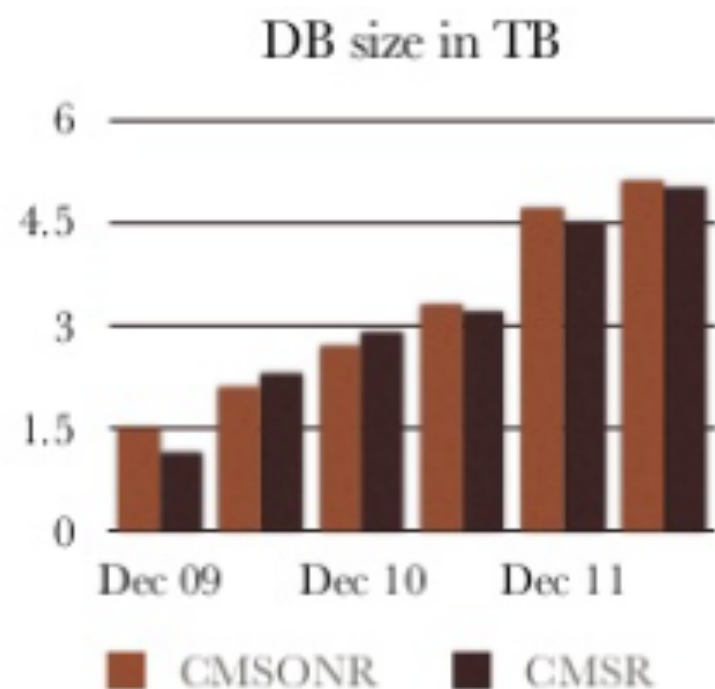  ➡ radiation resistance
  ➡ single high power laser
➡ **Large distribution**
  ➡ experiments with ~$10^7$ channels

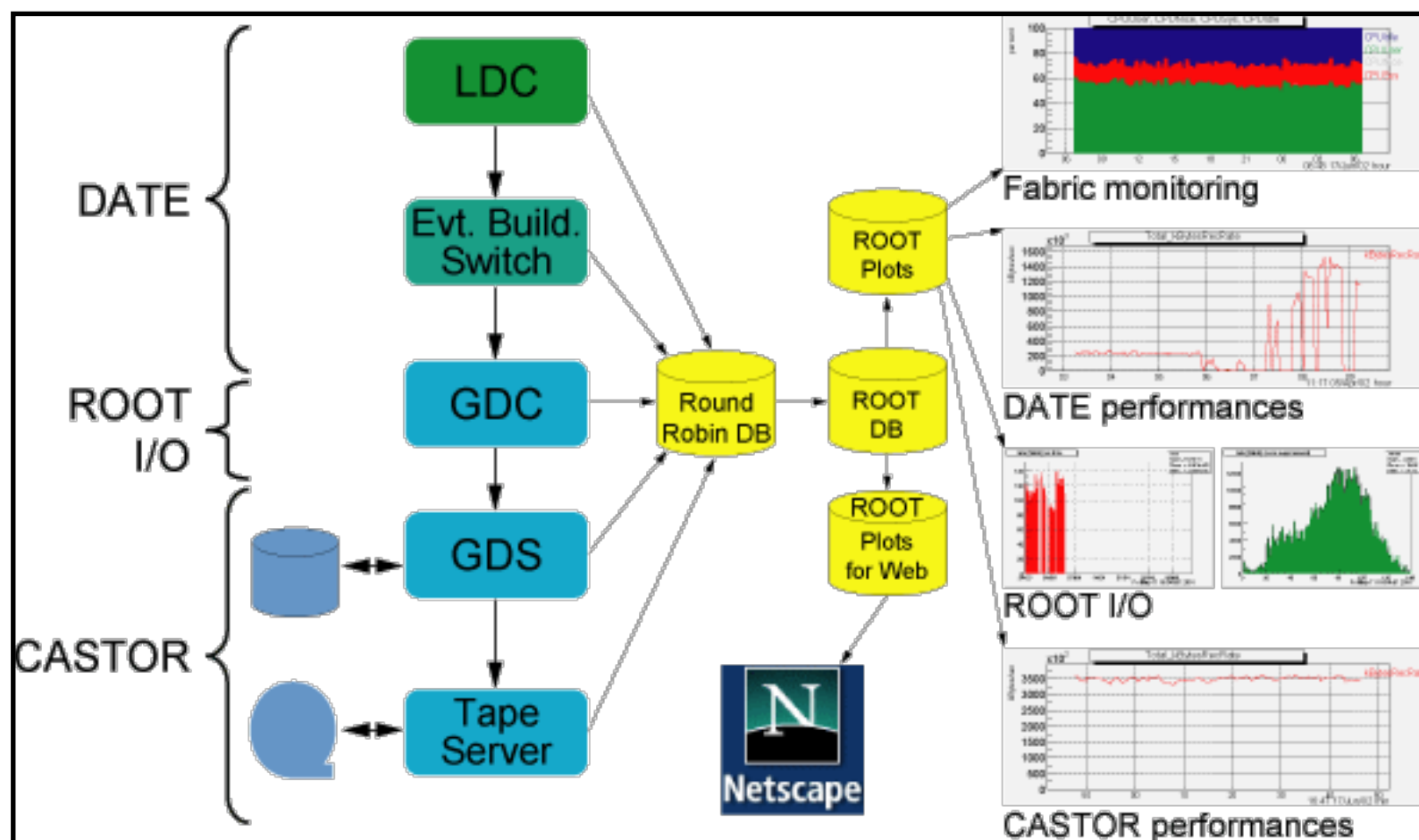➡ **Align readout & trigger at (better than) 25ns and correct for**
  ➡ time of flight (25 ns ≈ 7.5m)
  ➡ cable delays (10cm/ns)
  ➡ processing delays (~100 BCs)

➡ **Multiple Databases: configuration, condition, both online and offline**

　➡ Use (<u>Frontier</u>) caches to minimise access to Oracle servers

➡ **Monitoring and system administration**

　➡ thousands of nodes and network connections

　➡ advanced tools of monitoring and management

　➡ support software updates and rolling replacement of hardware



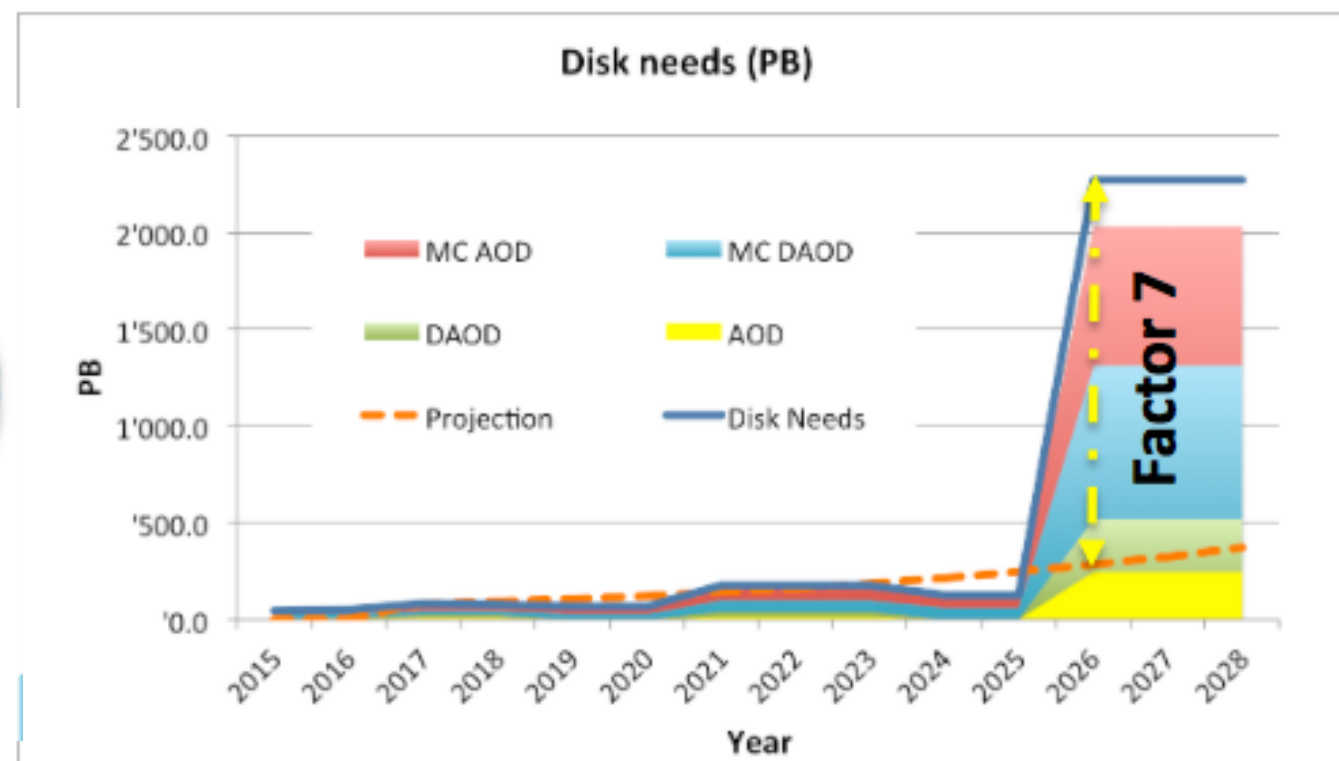CMS DB grows about 1.5TB/year, condition data only a small fraction

➡ **Re-thinking of distributed data management, distributed storage and data access.**

➡ **A network driven data model allows to reduce the amount of storage, particularly for disk**
  - ➡ Tape today costs 4 times less than disk

➡ **Computing infrastructure in HL-LHC**
  - ➡ Network-centric infrastructure
  - ➡ Storage and computing loosely coupled
  - ➡ Storage on fewer data centers in WLCG
  - ➡ Heterogeneous computing facilities (Grid/Cloud/HPC/ ...) everywhere

Projection of available resources in HL–LHC: 20% more CPU/year, 15% more storage/year

**electrons,
photons, taus,
jets,
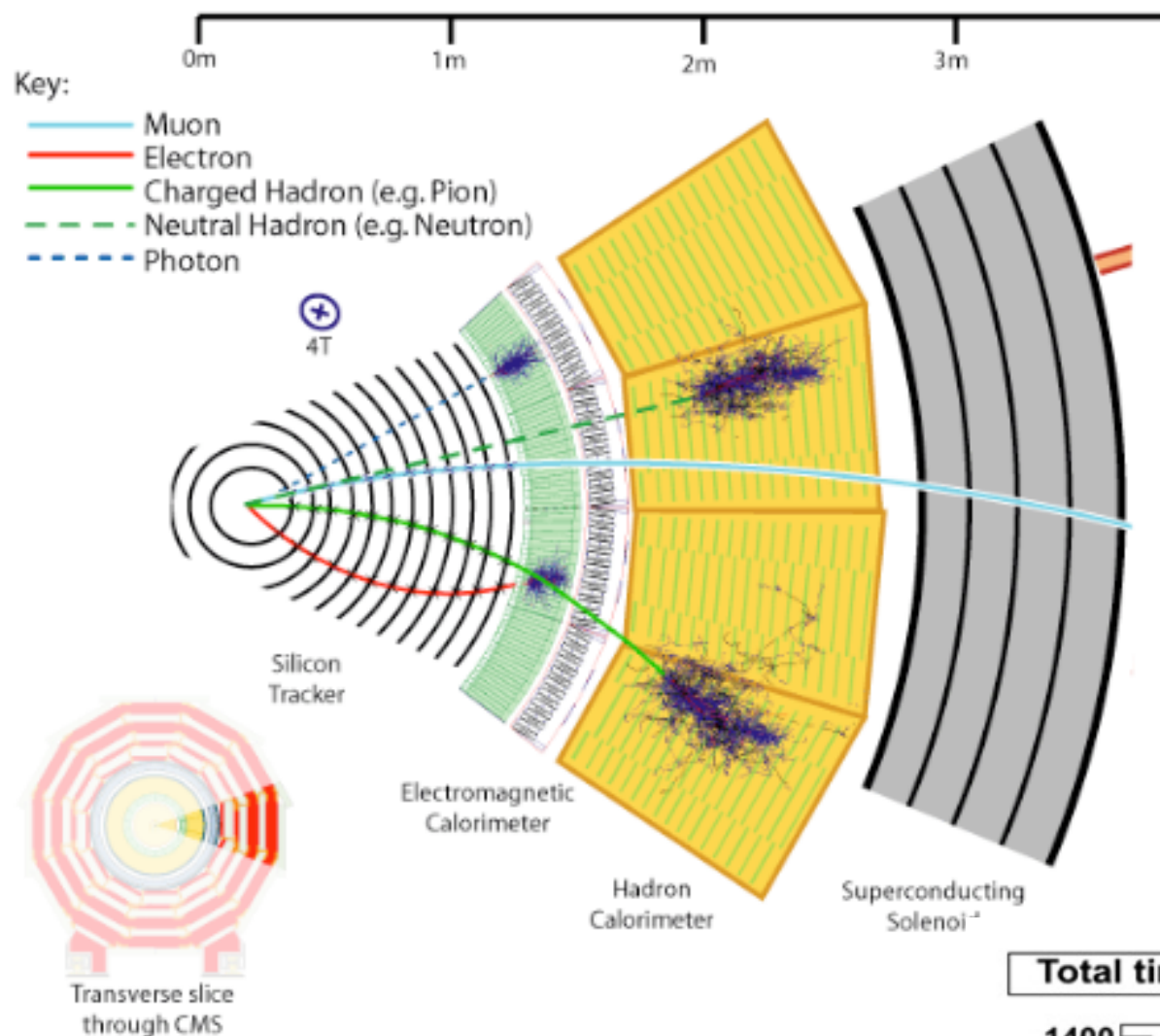total energy,
missing energy
Isolation**



Key:
— Muon
— Electron
— Charged Hadron (e.g. Pion)
--- Neutral Hadron (e.g. Neutron)
····· Photon

⊕
4T

Silicon
Tracker

Electromagnetic
Calorimeter

Hadron
Calorimeter

Superconducting
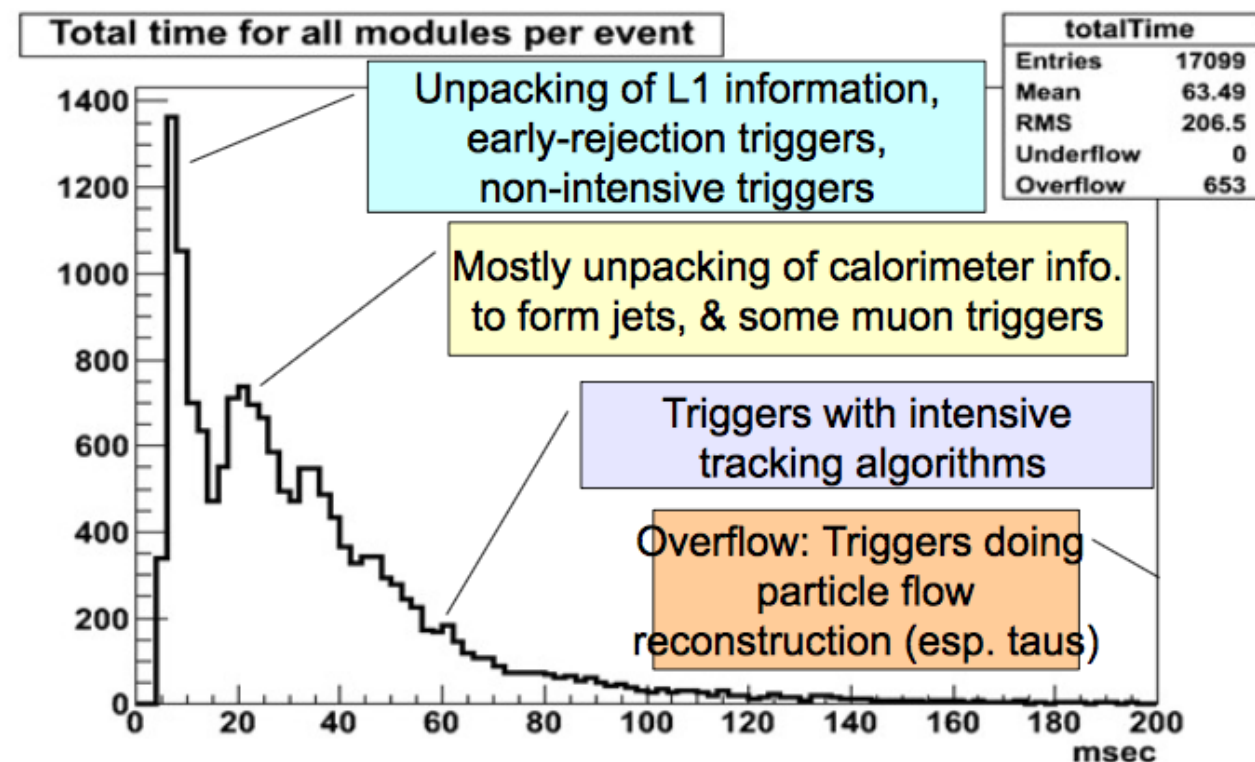Solenoi ⁻¹

Transverse slice
through CMS

➡ **Fast and good resolution (LArg, PbW$_4$ for e-m)**

➡ **First-level processing (40MHz)**
  ➡ "trigger towers" to reduce data (10-bit range)
  ➡ sliding-window technique for local maxima
  ➡ parallel algorithms for cluster shape and energy distribution
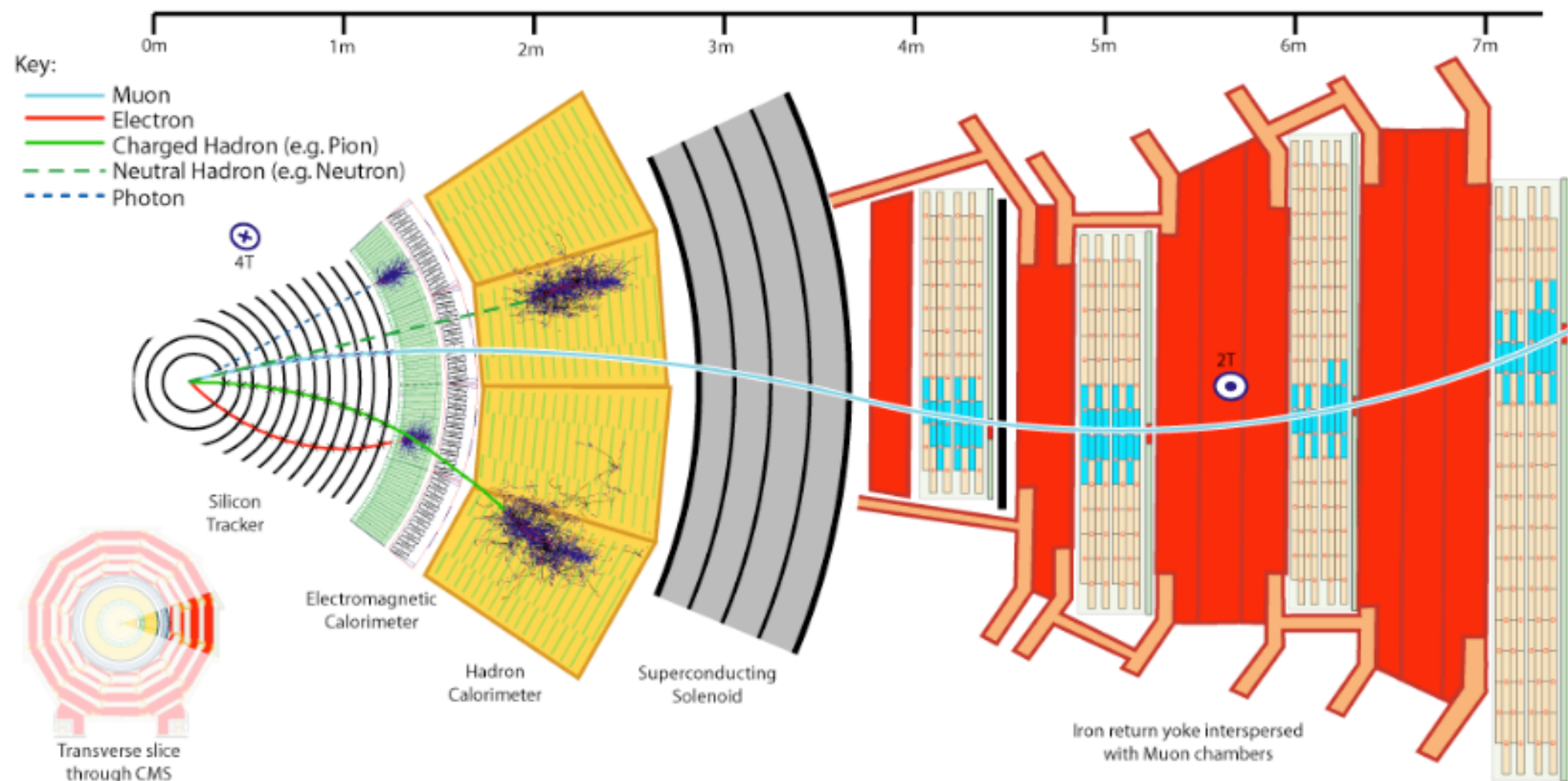
➡ **High-level processing (100 kHz)**
  ➡ regional tracking in the inner detectors
  ➡ bremsstrahlung recovery
  ➡ measure activity in cones (with tracks/ clusters) to isolate e/jets
  ➡ jet algorithms



**Total time for all modules per event**

Unpacking of L1 information,
early-rejection triggers,
non-intensive triggers

Mostly unpacking of calorimeter info.
to form jets, & some muon triggers

Triggers with intensive
tracking algorithms

Overflow: Triggers doing
particle flow
reconstruction (esp. taus)

| totalTime | |
|---|---|
| Entries | 17099 |
| Mean | 63.49 |
| RMS | 206.5 |
| Underflow | 0 |
| Overflow | 653 |

msec
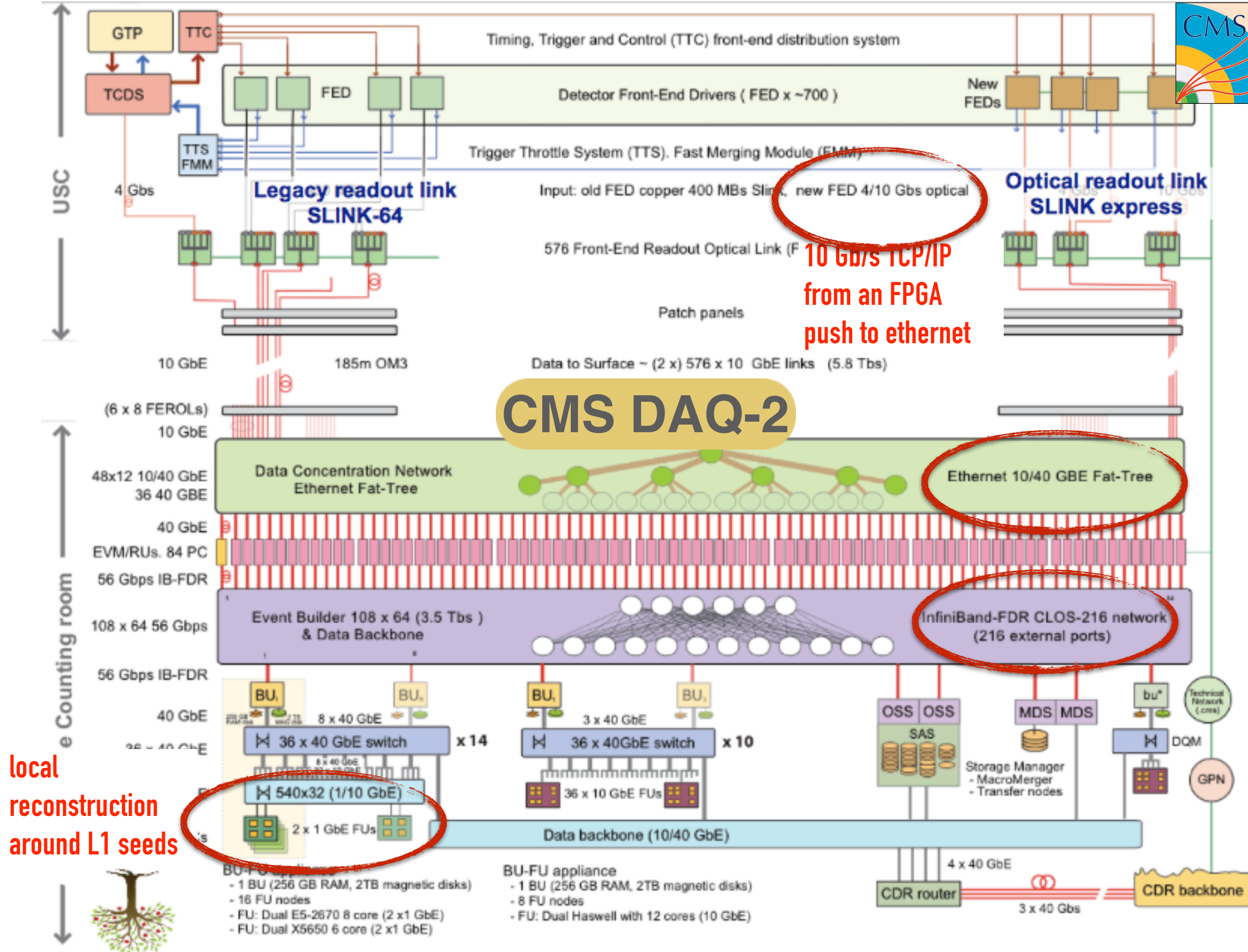
Transverse slice through CMS

➡ **Dedicated detectors:**

➡ low occupancy for fast pattern recognition

➡ optimal time-resolution for BC-identification

➡ **L1 processing (40 MHz)**
   ➡ pattern matching with patterns stored in buffers
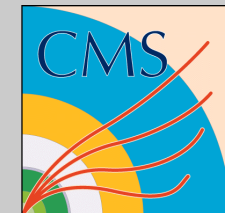   ➡ simplified fit of track segments

➡ **High level processing (100 kHz)**
   ➡ full detector resolutions
   ➡ match segments with tracks in the ID
   ➡ isolation
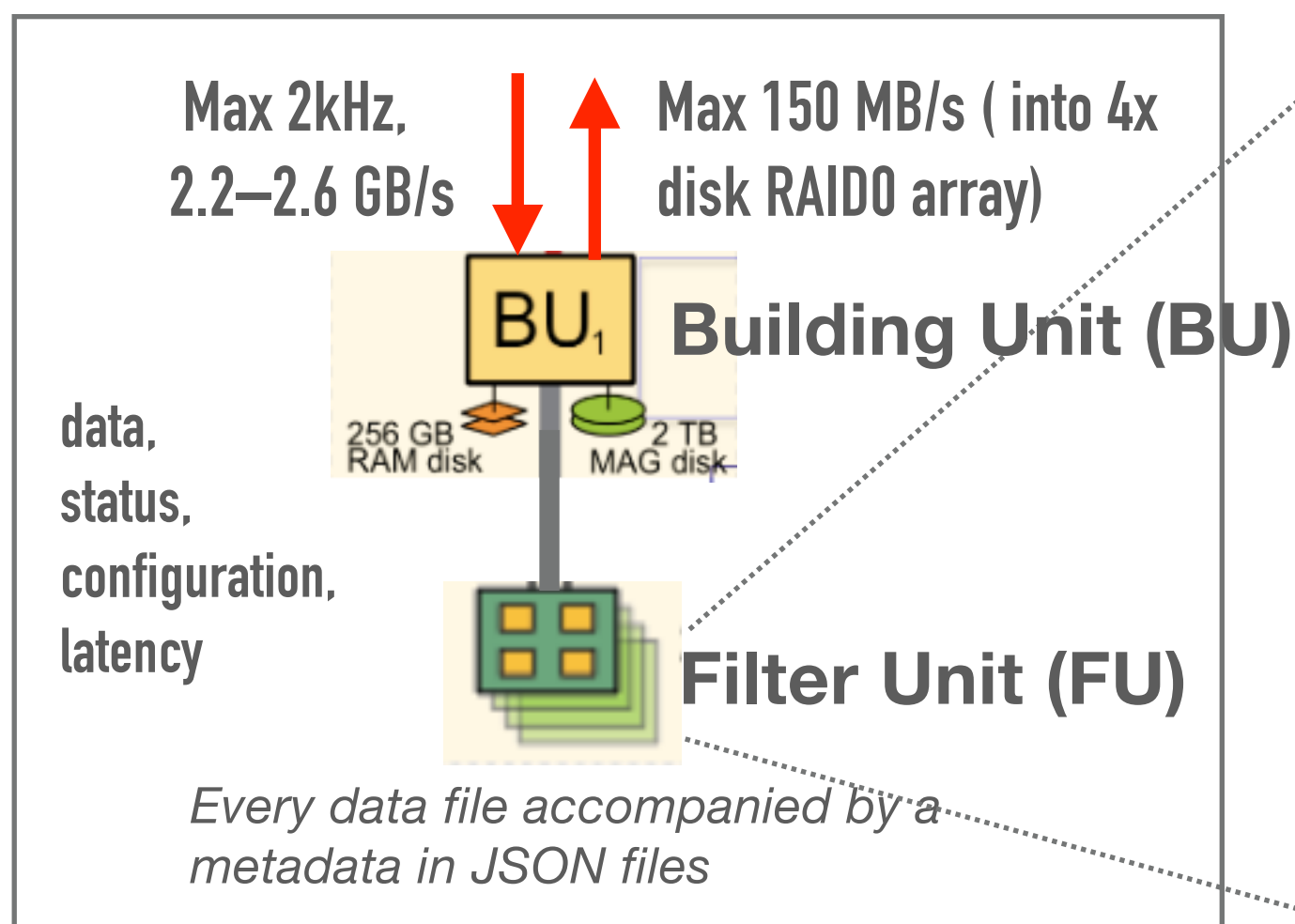
## Full readout, but <u>regional</u> <u>reconstruction</u> in HLT seeded by L1 trigger objects



**Max 2kHz, 2.2–2.6 GB/s**

**Max 150 MB/s ( into 4x disk RAID0 array)**

BU₁  **Building Unit (BU)**

data, status, configuration, latency

256 GB RAM disk   2 TB MAG disk

**Filter Unit (FU)**

*Every data file accompanied by a metadata in JSON files*

**Integrated Cloud capability (New!)**
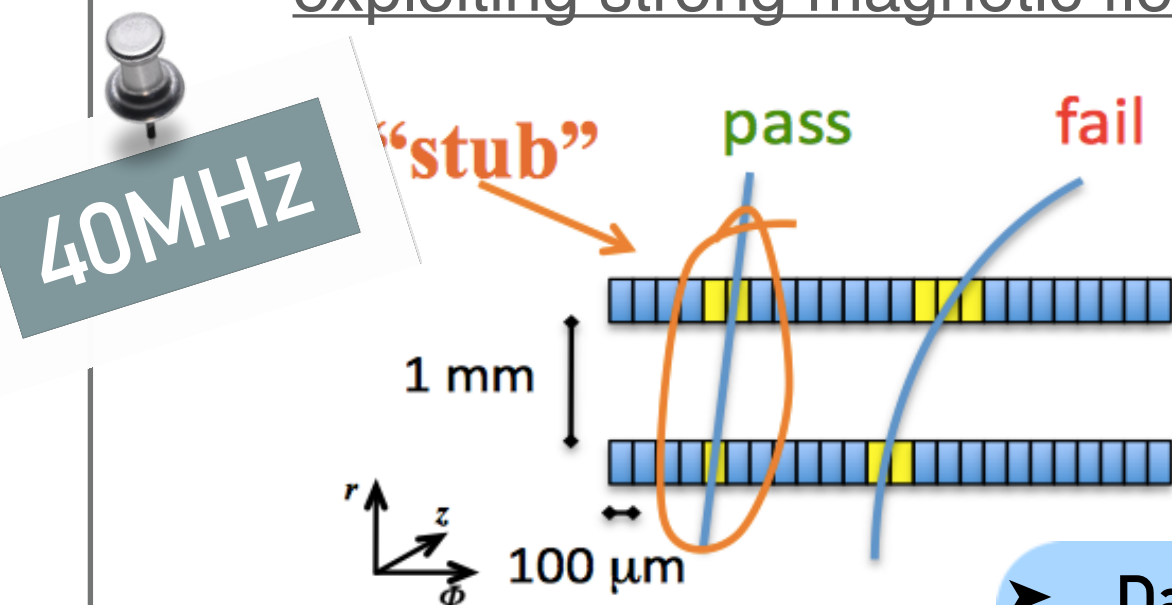➡ Added ability to run WLCG grid jobs in FUs during stops/interfill



HLT contribution

## File-based communication
➡ HLT and DAQ completely decoupled
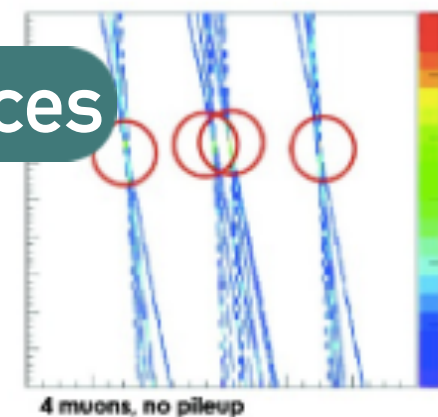➡ Network filesystem used as transport (and resource arbitration) protocol (LUSTRE FS)

## Track filtering (low p$_T$)

## Track finding options

**Reduce readout 40 ⇒ 1MHz by detector coincidences**
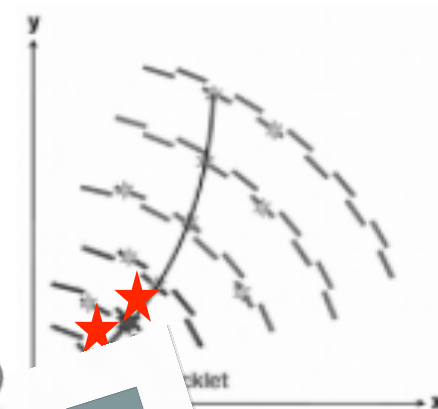
➡ **Special outer tracker modules**

  ➡ two layers of silicon at few mm

  ➡ using cluster width and stacked trackers

➡ **Design tracker to have coherent p$_T$ threshold in the full volume**
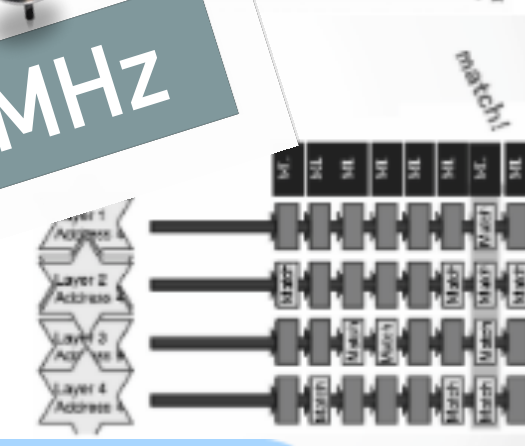
  ➡ exploiting strong magnetic field of CMS



"stub"   pass   fail

40MHz

1 mm

100 µm

Hough Transform

4 muons, no pileup

Tracklets

1MHz

Associative Memories

➤ Data rates > 50–100 Tbps
➤ Latency: 4+1 µs
➤ Three R&D efforts: FPGA/ASIC

Readout: 40 MHz
Event size: 100kB
DAQ: 40 Tbit/s
Record: 100 kHz

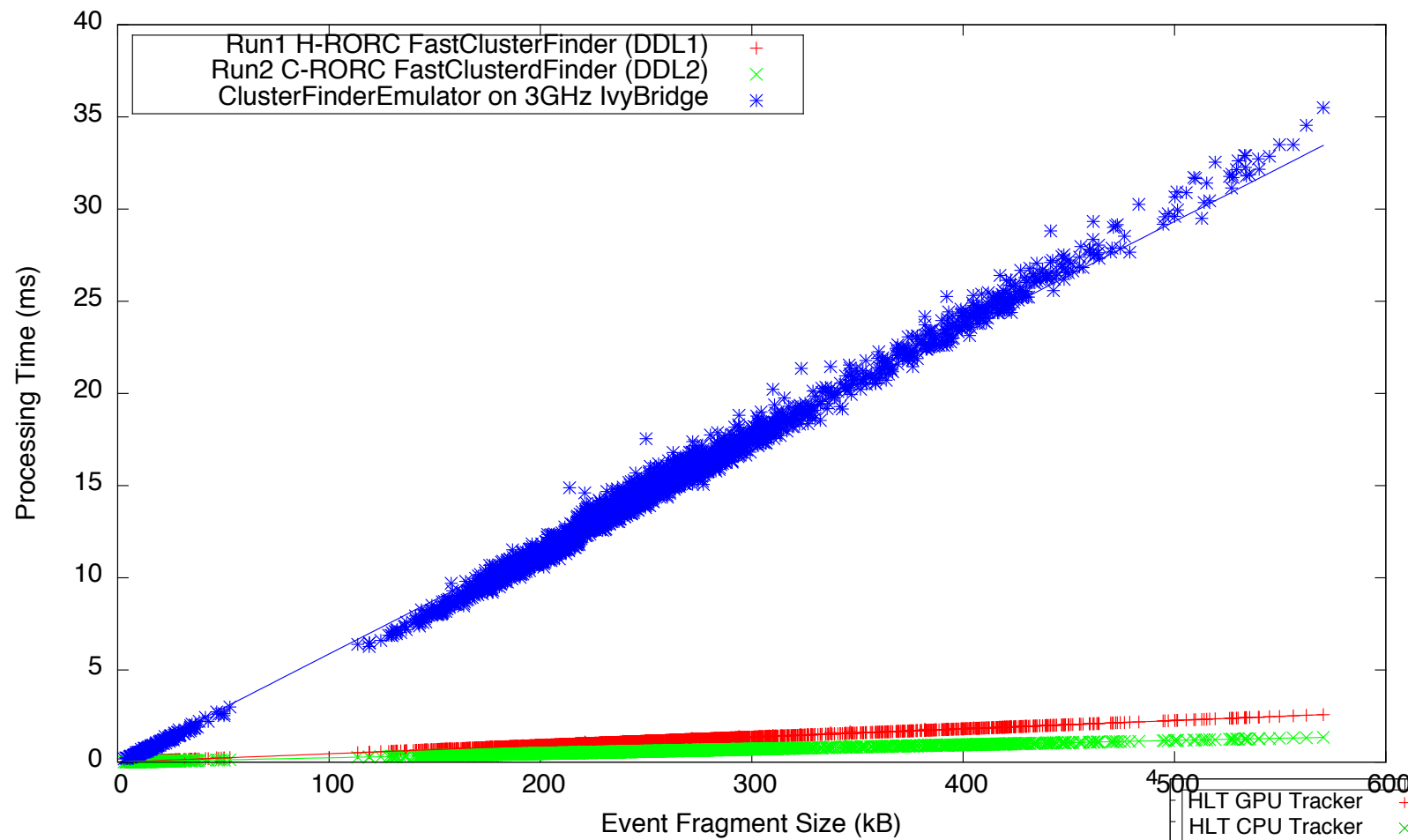➡ **Need zero-suppressing on front-end electronics**
➡ **A single, high performance, custom FPGA-card (PCIe40)**
  ➡ 8800 (# VL) * 4.48 Gbit/s (wide mode) => 40 Tbps
➡ **Single board up to 100 Gbits/s (to match DAQ links in 2018)**
➡ **Event-builder with 100 Gbit/s technology and data centre-switches**

**Detector**

VELO | ST | OT | RICH | ECal | HCal | Muon

FE Electronics

Readout Board

**Front - End**

**READOUT NETWORK**

~60 GB/s

**Event Building**

SWITCH | SWITCH | SWITCH | SWITCH | SWITCH | SWITCH

CPU CPU CPU CPU

**HLT farm**

**Experiment Control System**

L0 Trigger

L0 trigger
LHC clock

**TFC System**

MEP Request

~700 MB/s

STORAGE

SWITCH

CPU CPU CPU CPU

**MON farm**

**10GB Ethernet**

Deep buffering in the readout network (overloaded x300 at L0A)

**PUSH**

**PUSH**

62 sub-farms, total 1780 nodes, with edge-routers (12 Gbps)

—— Event data
- - - Timing and Fast Control Signals
—— Control and Monitoring data

**Average event size 60 kB**
**Average rate into farm 1 MHz**
**Average rate to tape ~12 kHz**
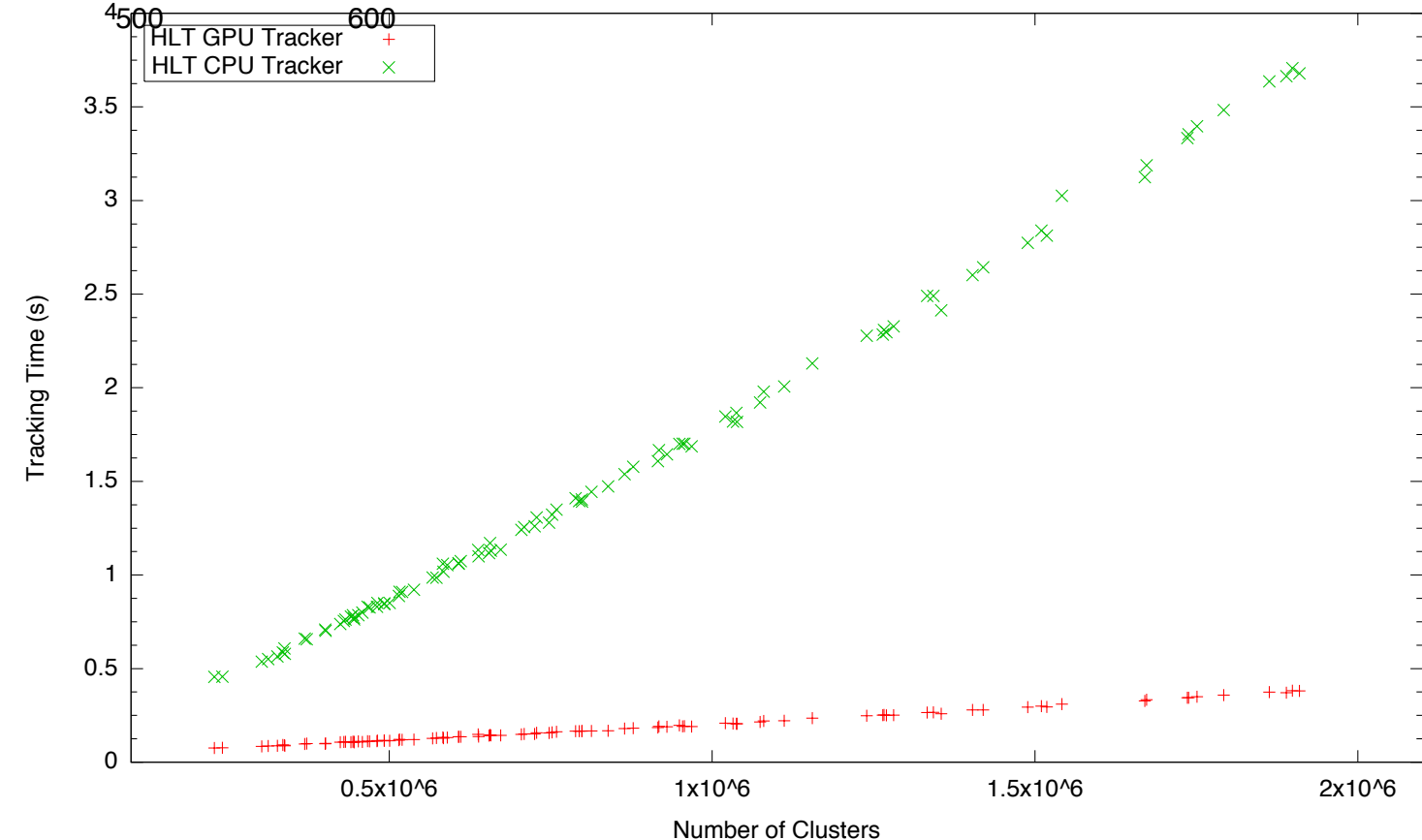
➡ **Small event, at high rate: ask for optimized transmission**

  ➡ TTC system is used to assign IP addresses to RO boards

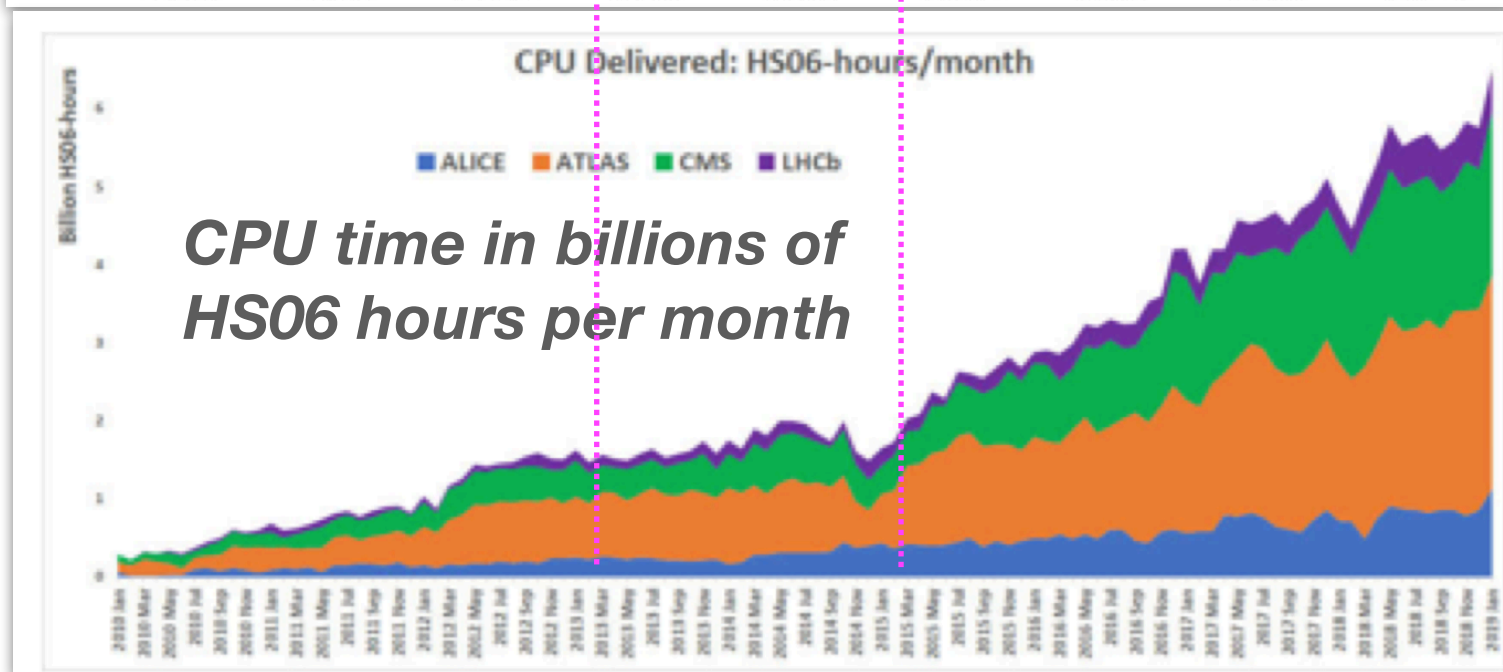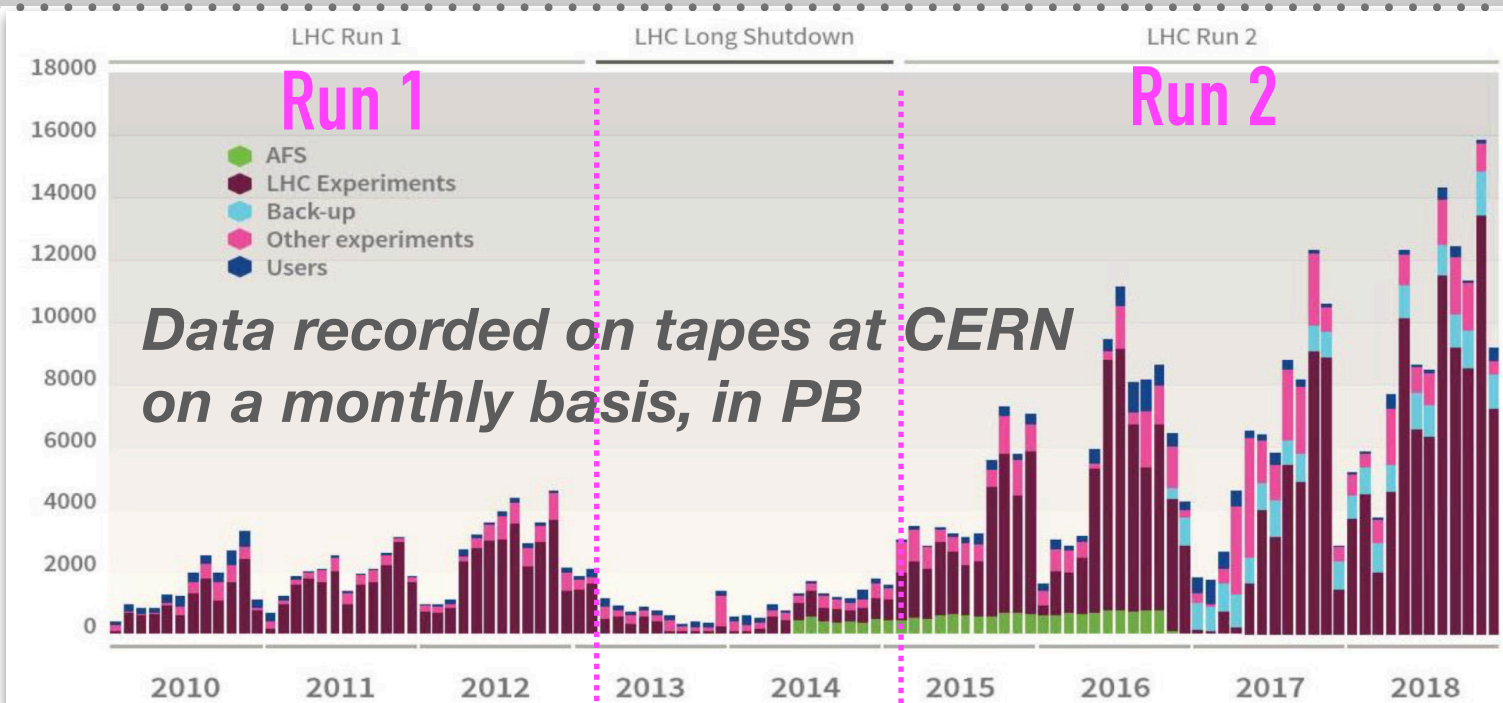  ➡ Ethernet UDP, with 10-15 events packed ⇒ ~ **80 kHz**

Tracking time of HLT TPC Cellular Automata tracker on Nehalem CPU (6Cores) and NVIDIA Fermi GPU.

Performance of the FPGA-based FastClusterFinder algorithm for DDL1 (Run1) and DDL2 (Run2) compared to the software implementation on a recent server PC.

# LHC COMPUTING TOWARDS NEW PARADIGMS



Run 1

Run 2

*Data recorded on tapes at CERN on a monthly basis, in PB*

*CPU time in billions of HS06 hours per month*
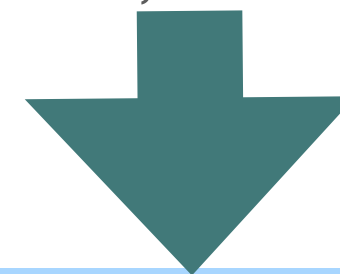
## Run1 + Run2

➡ **Data storage**
- ➡ 339 PB on tapes, 173 PB on disks

➡ **Global CPU time delivered by Worldwide LHC Computing Grid (WLCG)**
- ➡ about 900,000 cores

## Run 3

➡ **Evolution of current technologies and current (flat) funding is ok**

## Run 4

➡ **Linear increase of digitisation time**
➡ **Factorial increase of reconstruction time**
➡ **Larger events, lots of more memory**

➡**Need factor 2-3 more storage and computing resources for HL-LHC**
- ➡ new developments and R&D projects for data management and processing, SW multithreading, new computing models and data compression