*F.Pastore  (Royal Holloway Univ. of London)*
*francesca.pastore@cern.ch*

# FROM SMALL TO LARGE T/DAQ SYSTEMS

➡ **Examples of small experiments with their limits**

➡ **Overview of LHC experiments and their upgrade**

➡ **Future TDAQ systems (Dune/Proto-Dune)**

➡ **Examples of small experiments with their limits**

➡ Overview of LHC experiments and their upgrade

➡ Future TDAQ systems (Dune/Proto-Dune)

➡ **Data Size**

　➡ Summing up data from all Front-End channels

　　➡ 100 M channels of silicon detectors give few MB/event

　➡ depends on detector granularity (number of channels), on detector technology (single bit versus drift-time or TPC) and pile-up level

# SOME NOMENCLATURE

➡ **Data Size**

  ➡ Summing up data from all Front-End channels

    ➡ 100 M channels of silicon detectors give few MB/event

  ➡ depends on detector granularity (number of channels), on detector technology (single bit versus drift-time or TPC) and pile-up level

➡ **Data Rate**

  ➡ Front-End readout rate

  ➡ LHC clock gives about 40 M evt/s at 13 TeV

  ➡ Pierre Auger Observatory: about 1 evt/100years/km at EeV

# SOME NOMENCLATURE

➡ **Data Size**

   ➡ Summing up data from all Front-End channels

      ➡ 100 M channels of silicon detectors give few MB/event

   ➡ depends on detector granularity (number of channels), on detector technology (single bit versus drift-time or TPC) and pile-up level
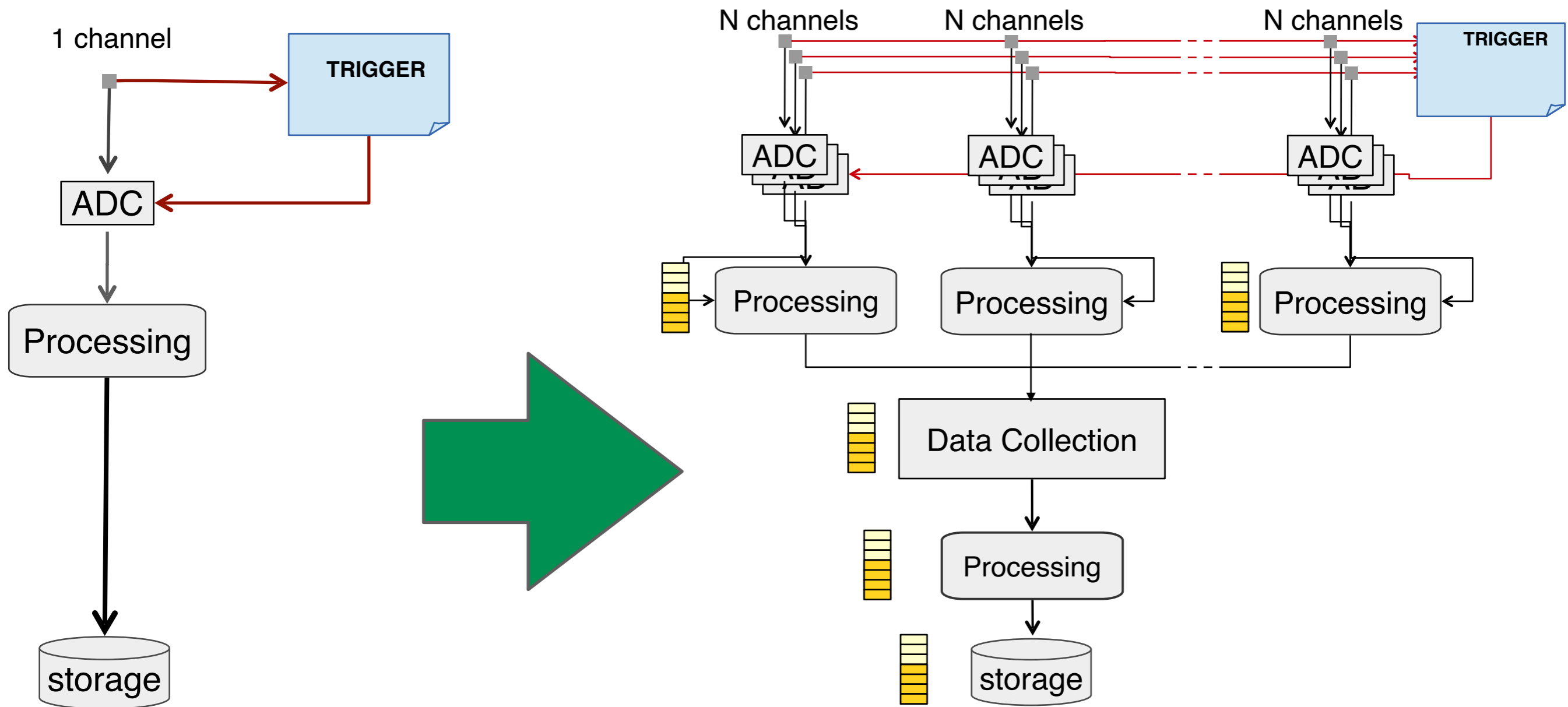
➡ **Data Rate**

   ➡ Front-End readout rate

   ➡ LHC clock gives about 40 M evt/s at 13 TeV

   ➡ Pierre Auger Observatory: about 1 evt/100years/km at EeV

➡ **DAQ bandwidth**

   ➡ 40MHz x 1 MB = 40TB/s

      ➡ too much data!

   ➡ select and record only the most important events

➡ **Two independent paths for trigger and DAQ**

➡ **Segmented Readout and trigger to allow parallel processing**

➡ **Included buffers at each stage to control dead-time**
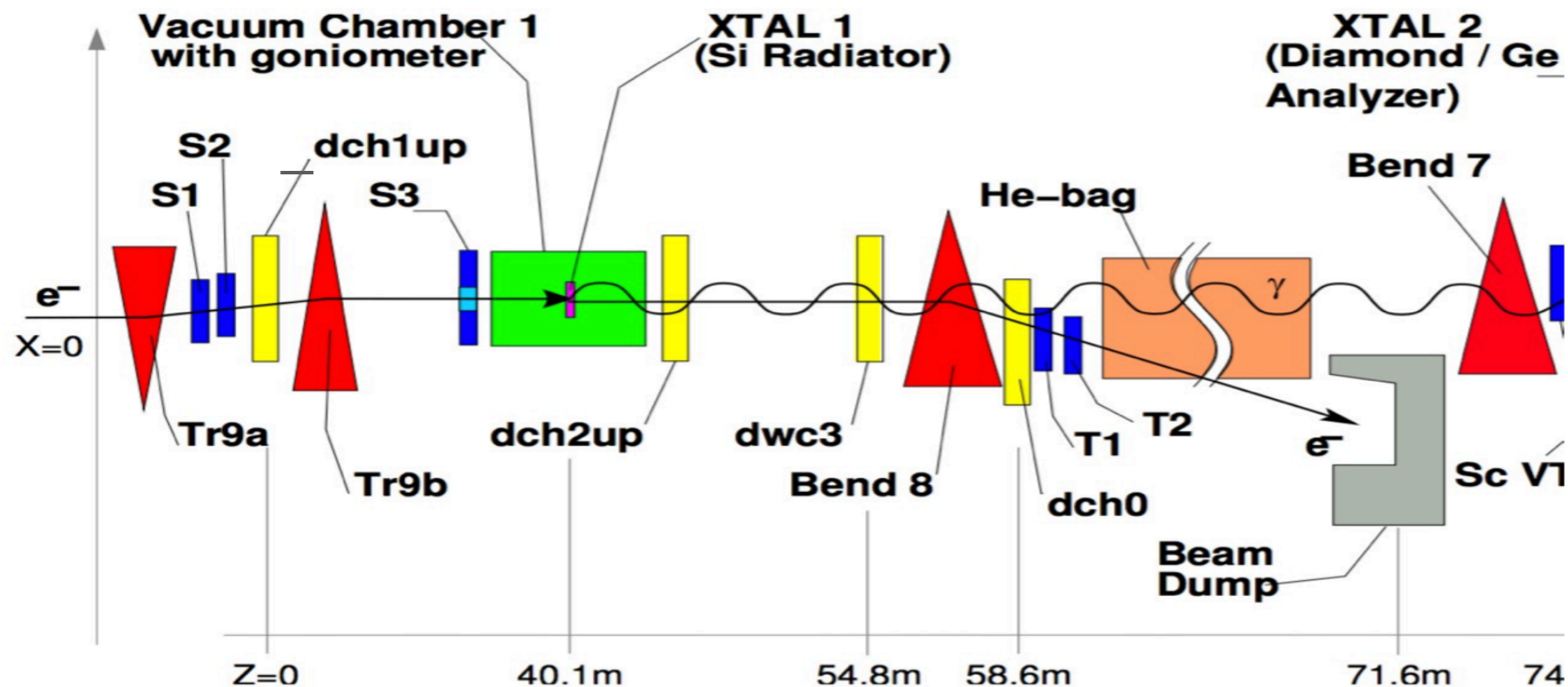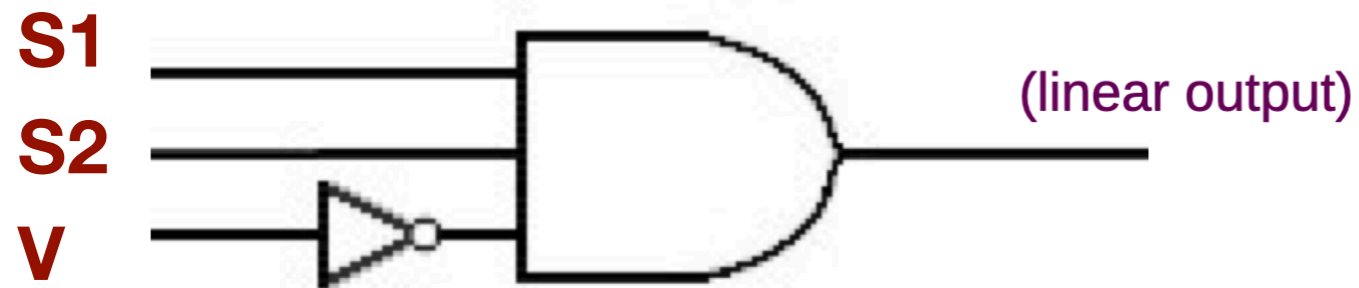
➡ **How to scale these systems?**
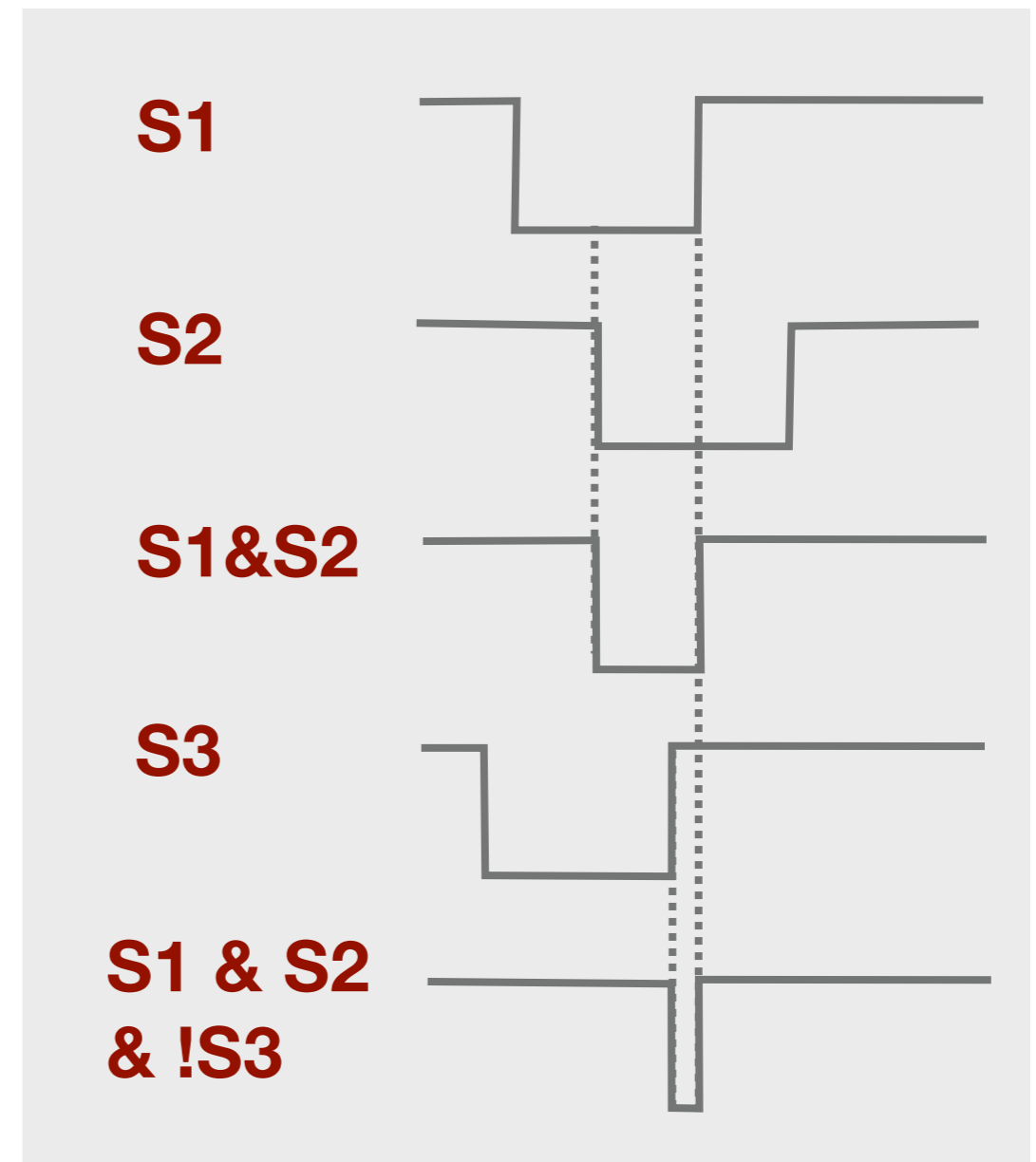
Fig. 1. Setup of the Na59 Experiment

➡ **Trigger event with an electron at the correct incident angle wrt crystal**

➡ **Three scintillators S1, S2 and S3 ensure the arrival of the beam within the acceptance of the crystal**
  - ➡ Input **N1 = S1 & S2 & !S3** ---> an electron is coming and it is not away from the central axis
  - ➡ use S3 as veto (anti-coincidence)

➡ **After the magnet, two scintillators to tag the electron out of the beam**
  - ➡ **N2 = N1 x (T1 || T2)** ---> the electron radiated a photon and was diverted by the magnet

➡ **Simple coincidence and veto logic can be broken if signals are not formed correctly**

S1
S2
V

(linear output)

➡ **Signals are random/independent**
➡ **Can fluctuate in duration and jitter**
  ➡ Need preliminary **timing alignment** between signals
    ➡ e.g adding delays to faster signals
  ➡ Need **forming output signals** with known width
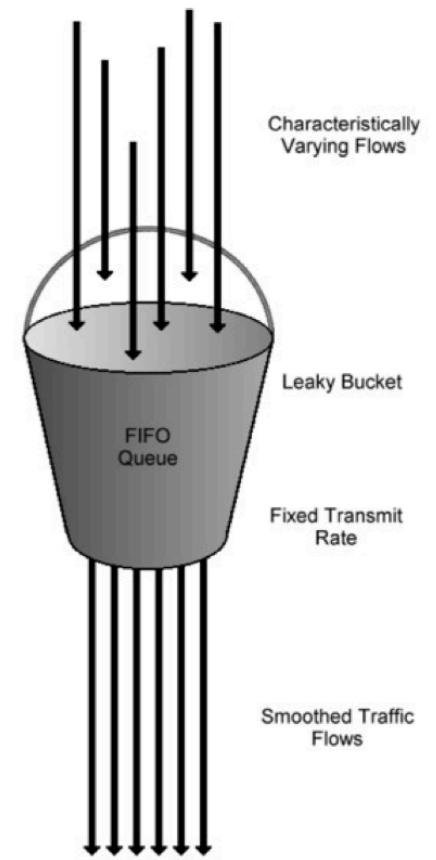    ➡ fix width of output signal at each step

S1

S2

S1&S2

S3

S1 & S2 & !S3

➡ **Step 1: Increasing rate**

➡ **Step 2: Increasing sensors**

➡ **Step 3: Multiple front-ends**

➡ **Step 4: Multi-level trigger**

➡ **Step 5: Data-flow control**

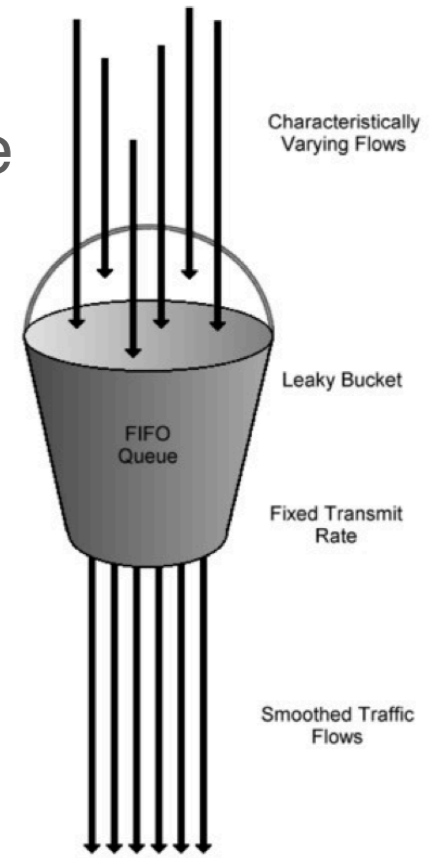Characteristically
Varying Flows

Leaky Bucket

FIFO
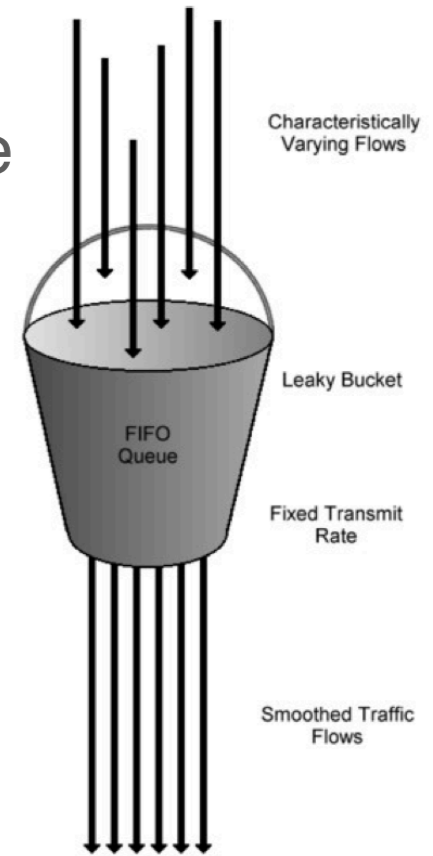Queue

Fixed Transmit
Rate

Smoothed Traffic
Flows

➡ **If two signals arrive very close in time**

  ➡ detector signals overlap (ask you detector expert, are you sure the detector is good at that rate? is your FE fast enough?)

  ➡ can have dead-time if not added any … FIFO!

Characteristically Varying Flows

Leaky Bucket

FIFO Queue

Fixed Transmit Rate

Smoothed Traffic Flows

➡ **If two signals arrive very close in time**

   ➡ detector signals overlap (ask you detector expert, are you sure the detector is good at that rate? is your FE fast enough?)

   ➡ can have dead-time if not added any … FIFO!

➡ **Is derandomization enough?**

   ➡ if FE readout windows overlap

      ➡ add artificial dead-time to protect the FrontEnd (**simple deadtime**)

   ➡ if FE buffers overflow in case of trigger bursts
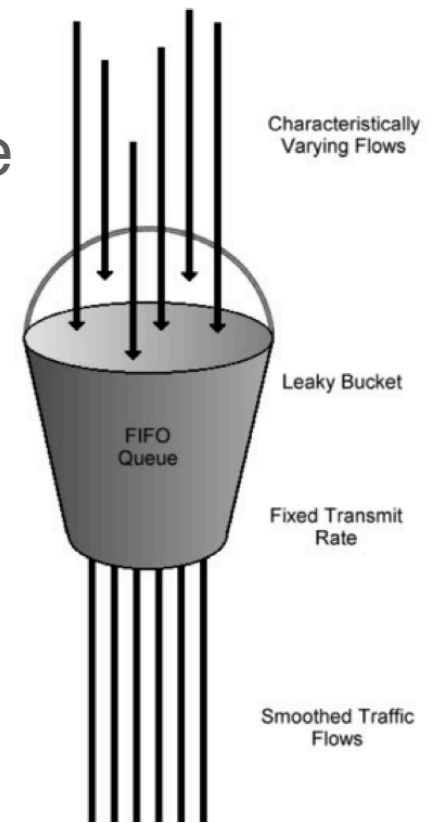
      ➡ add artificial dead-time (**complex deadtime**)

Characteristically
Varying Flows

Leaky Bucket

FIFO
Queue

Fixed Transmit
Rate

Smoothed Traffic
Flows

➡ **If two signals arrive very close in time**

  ➡ detector signals overlap (ask you detector expert, are you sure the detector is good at that rate? is your FE fast enough?)

  ➡ can have dead-time if not added any … FIFO!
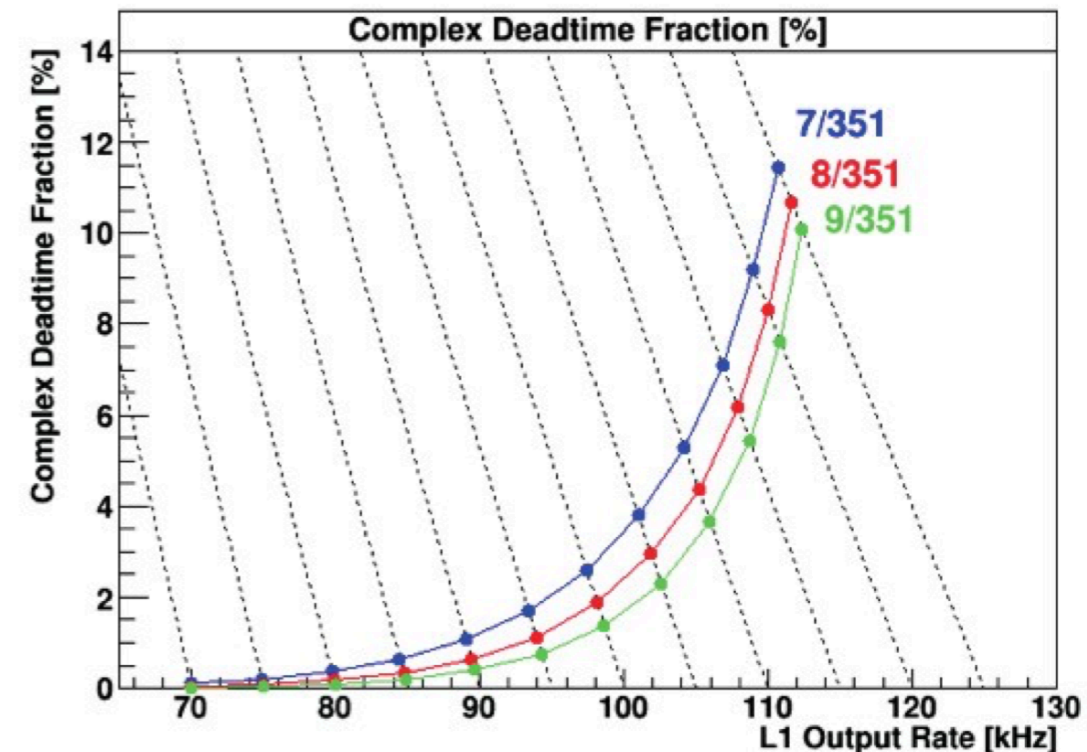
➡ **Is derandomization enough?**

  ➡ if FE readout windows overlap

    ➡ add artificial dead-time to protect the FrontEnd (**simple deadtime**)

  ➡ if FE buffers overflow in case of trigger bursts
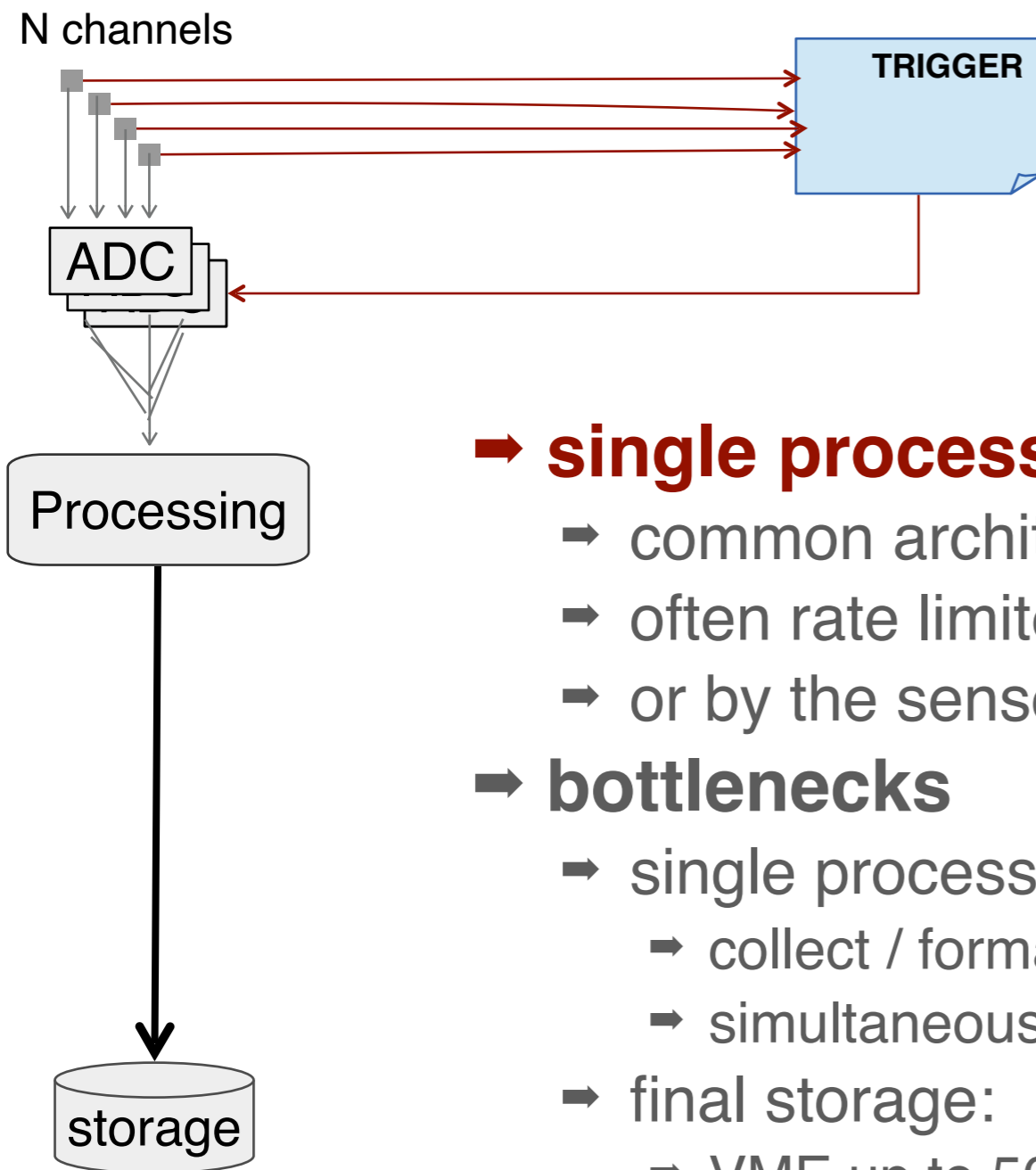
    ➡ add artificial dead-time (**complex deadtime**)



Characteristically Varying Flows

Leaky Bucket

FIFO Queue

Fixed Transmit Rate

Smoothed Traffic Flows

**Leaky bucket (LAr readout)**

➡ **Example in ATLAS @Run2: 90 kHz < 2%**

  ➡ Simple deadtime: 4 LHC BC [100 ns] after any L1 trigger

  ➡ Complex deadtime: leaky-bucket algorithms x4 detectors

    ➡ two params: bucket size (in number of events), /readout time (in BC units)

    ➡ i.e. 9 / 351 for LAr readout



Complex Deadtime Fraction [%]

7/351
8/351
9/351

Complex Deadtime Fraction [%] — L1 Output Rate [kHz]

N channels

TRIGGER

➡ **more sensors ==> more granularity**
➡ **multiple digitisers ==> parallelism**

ADC

Processing

storage

➡ **single processing system**
  ➡ common architecture in **test-beams and small experiments**
  ➡ often rate limited by (interesting) physics itself, not TDAQ
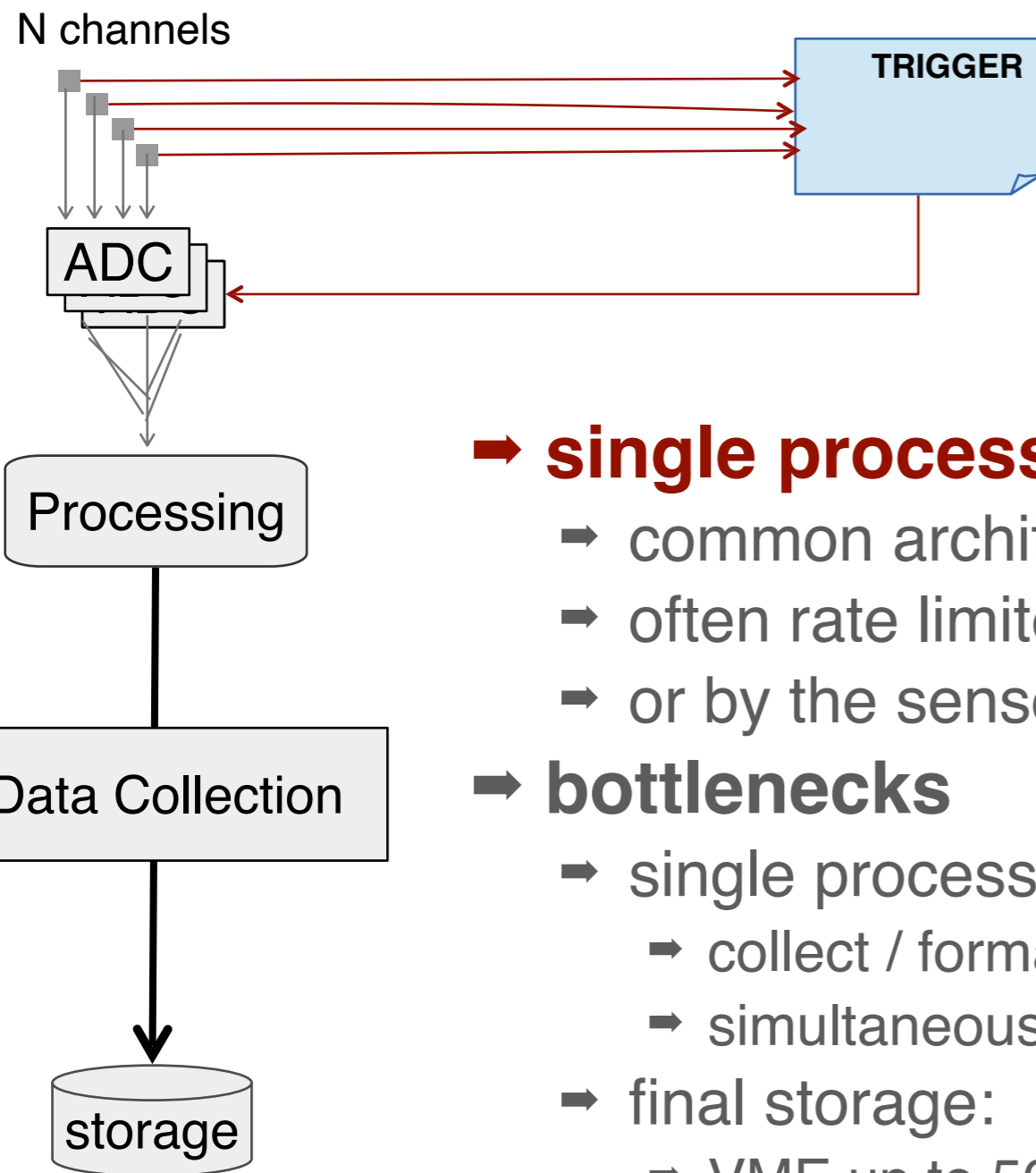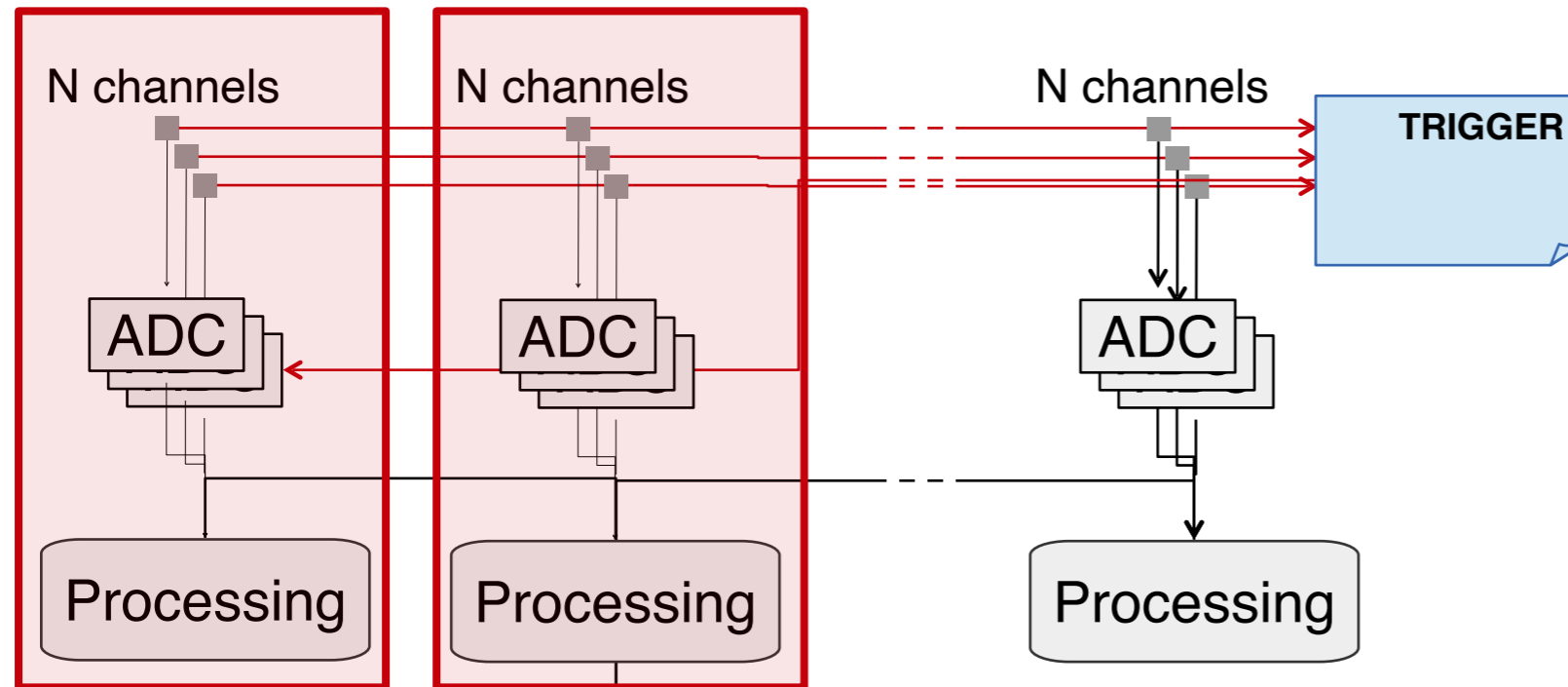  ➡ or by the sensors
➡ **bottlenecks**
  ➡ single processing unit
    ➡ collect / format / compress data can be heavy
    ➡ simultaneously writing storage
  ➡ final storage:
    ➡ VME up to 50MB/s → 1TB in 6h
    ➡ too many disks in one week!
➡ **decouple storage from processing unit (PU)**
  ➡ dedicated "Data Collection" unit to format, compress and store

➡ **more sensors ==> more granularity**
➡ **multiple digitisers ==> parallelism**

➡ **single processing system**
  ➡ common architecture in **test-beams and small experiments**
  ➡ often rate limited by (interesting) physics itself, not TDAQ
  ➡ or by the sensors
➡ **bottlenecks**
  ➡ single processing unit
    ➡ collect / format / compress data can be heavy
    ➡ simultaneously writing storage
  ➡ final storage:
    ➡ VME up to 50MB/s → 1TB in 6h
    ➡ too many disks in one week!
➡ **decouple storage from processing unit (PU)**
  ➡ dedicated "Data Collection" unit to format, compress and store

➡ **Multiple processing units**
  ➡ for data processing and storage

➡ **e.g.: CERN LEP experiments**
  ➡ complex detectors, moderate trigger rate, very little background
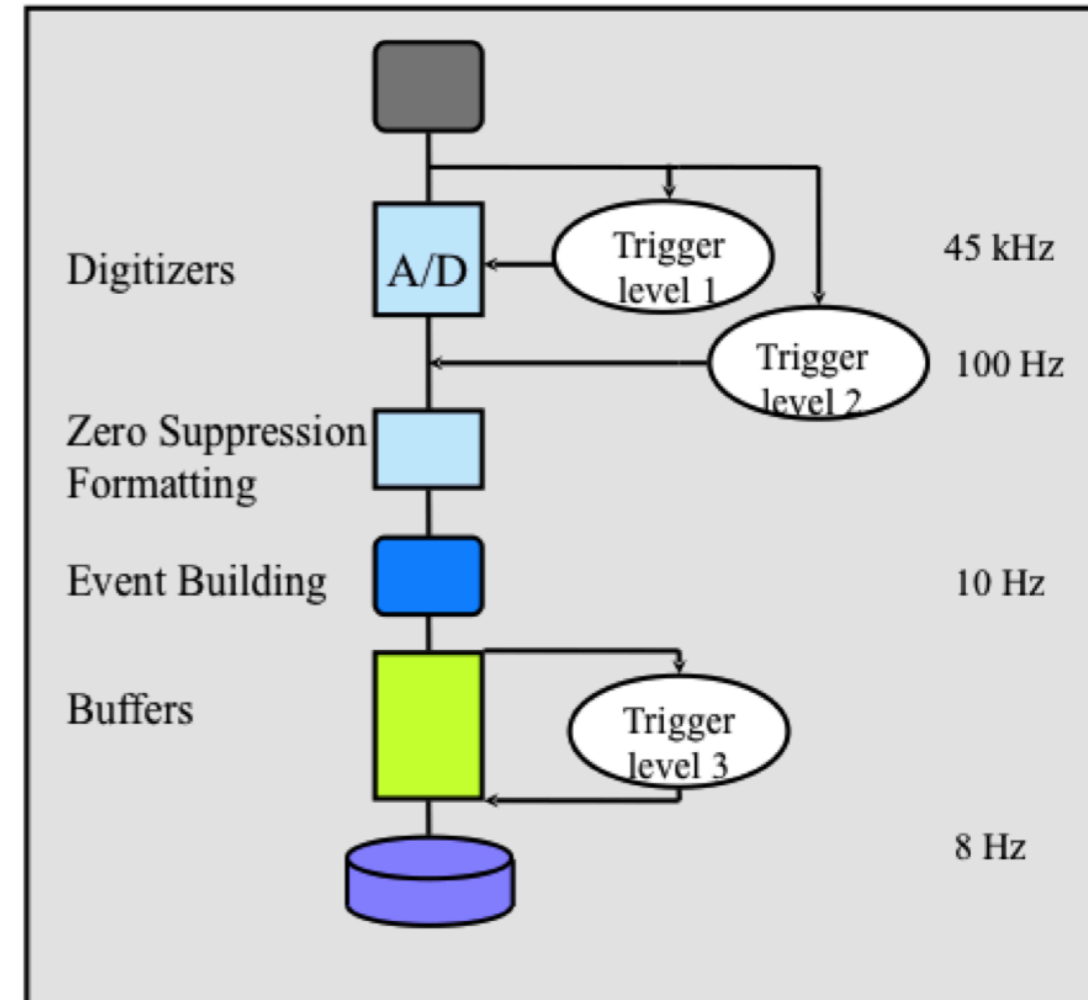  ➡ little pileup, limited channel occupancy

➡ **simpler, slow gas-based main trackers**

➡ **LEP**
➡ $10^5$ channels
➡ 22μs crossing rate – no event overlap
➡ single interaction

➡ **More channels + more rate + more data to process online ==> longer latency**

    ➡ single level trigger not enough

➡ **Add High level triggers with longer latency**

    ➡ more complex filters

    ➡ more data (for example silicon detectors)



➡ **Recall on trigger system architectures**

➡ **Real time system**

    ➡ must respond within some **fixed latency**

    ➡ → Latency = Max Latency

    ➡ → over fluctuations bad, will create deadtime

➡ **Non-real-time system**

    ➡ responds as soon as it's available

    ➡ → Latency = **Mean Latency**

    ➡ → over fluctuations fine, shouldn't create deadtime

➡ **LEP**
➡ $10^5$ channels
➡ 22µs crossing rate – no event overlap
➡ single interaction
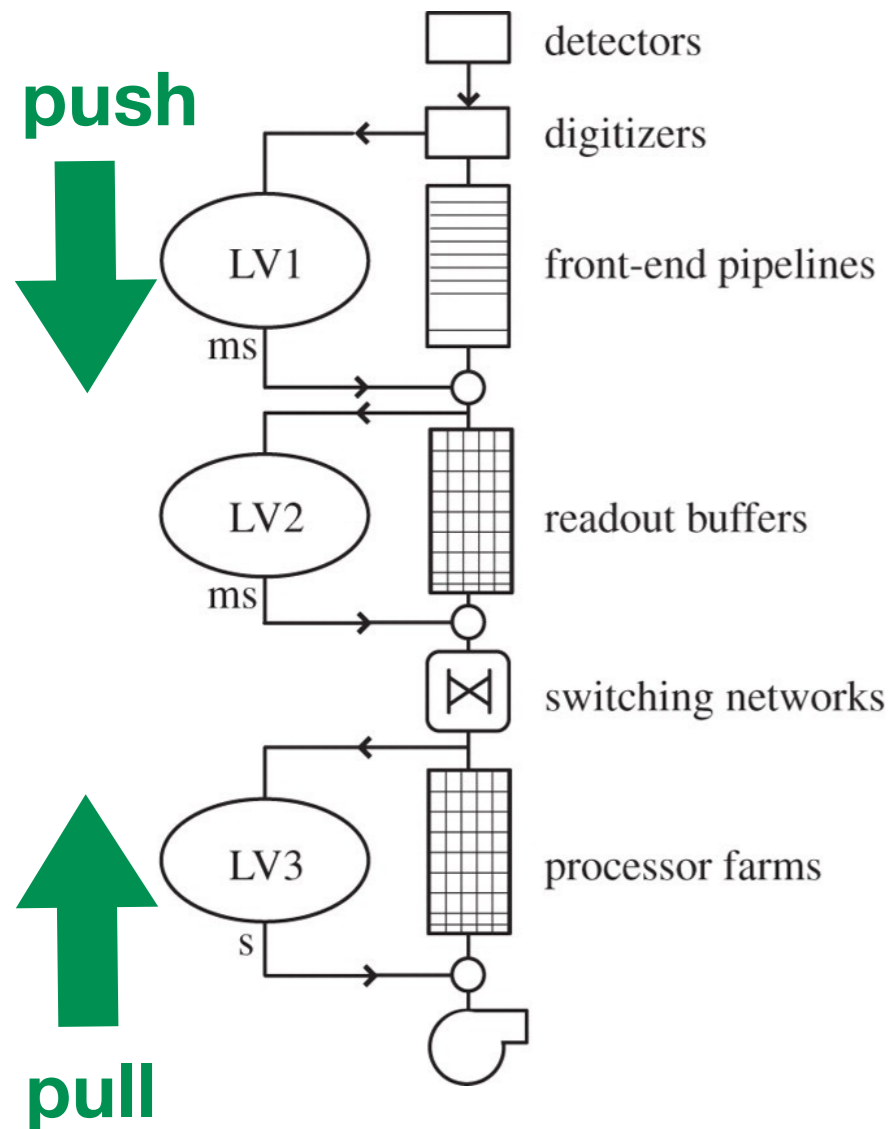➡ L1 ~$10^3$ Hz
➡ L2 ~$10^2$ Hz
➡ L3 ~10 Hz
➡ 100kB/ev → 1MB/s

**push**

**pull**

detectors

digitizers

LV1

ms

front-end pipelines

LV2

ms

readout buffers

switching networks

LV3

s

processor farms

**push**

**pull**
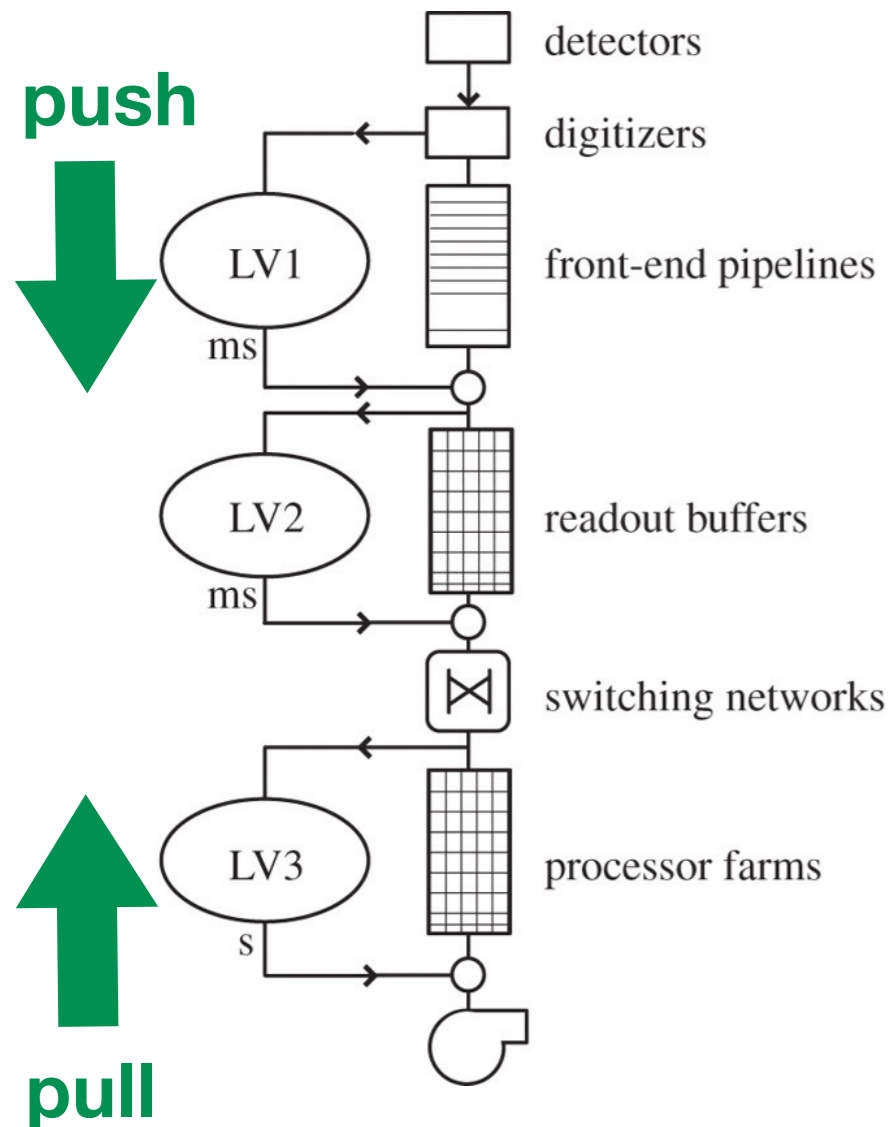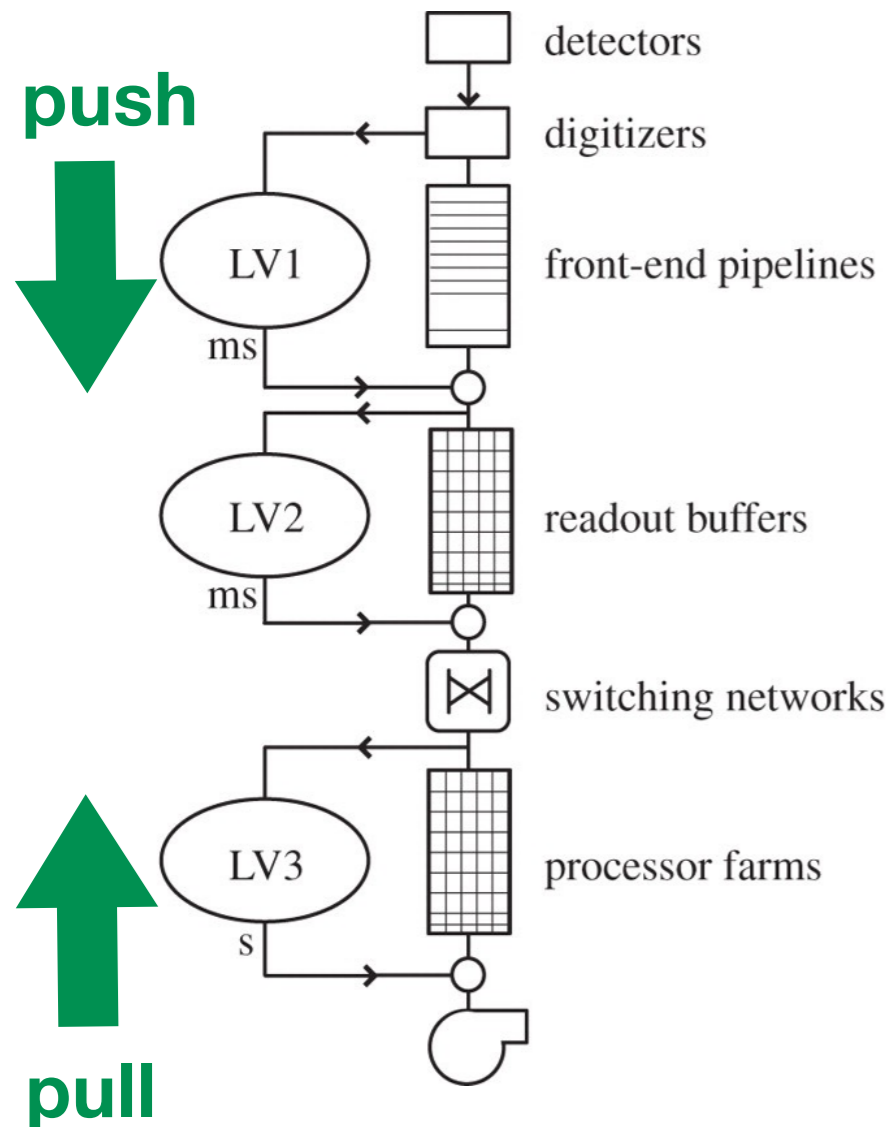

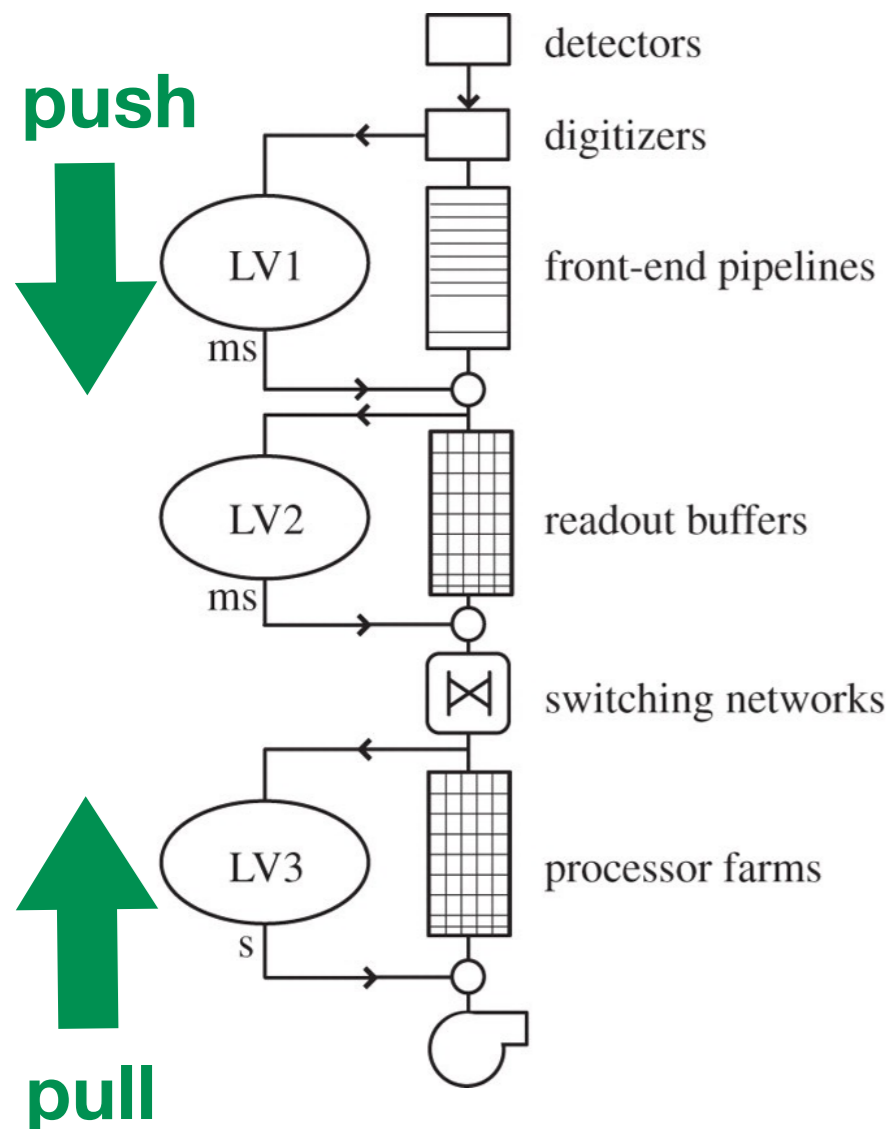
➡ **Buffers are not the "final solution"**

   ➡ Can overflow, with bursts and unusual event sizes

   ➡ In these cases

      ➡ discard data locally or

      ➡ exert "**back-pressure**", i. e. ask previous level(s) to block dataflow

**push**

**pull**

LV1
ms

LV2
ms

LV3
s

detectors

digitizers

front-end pipelines

readout buffers

switching networks

processor farms

➡ **Buffers are not the "final solution"**

   ➡ Can overflow, with bursts and unusual event sizes

   ➡ In these cases

      ➡ discard data locally or

      ➡ exert "**back-pressure**", i. e. ask previous level(s) to block dataflow

➡ **Throughput optimization means avoiding dead-time due to back-pressure**

   ➡ using knowledge of the input buffer state

**push**

**pull**

detectors

digitizers

front-end pipelines

readout buffers

switching networks

processor farms

LV1 ms

LV2 ms

LV3 s

➡ **Buffers are not the "final solution"**

➡ Can overflow, with bursts and unusual event sizes

➡ In these cases

➡ discard data locally or

➡ exert "**back-pressure**", i. e. ask previous level(s) to block dataflow

➡ **Throughput optimization means avoiding dead-time due to back-pressure**

➡ using knowledge of the input buffer state
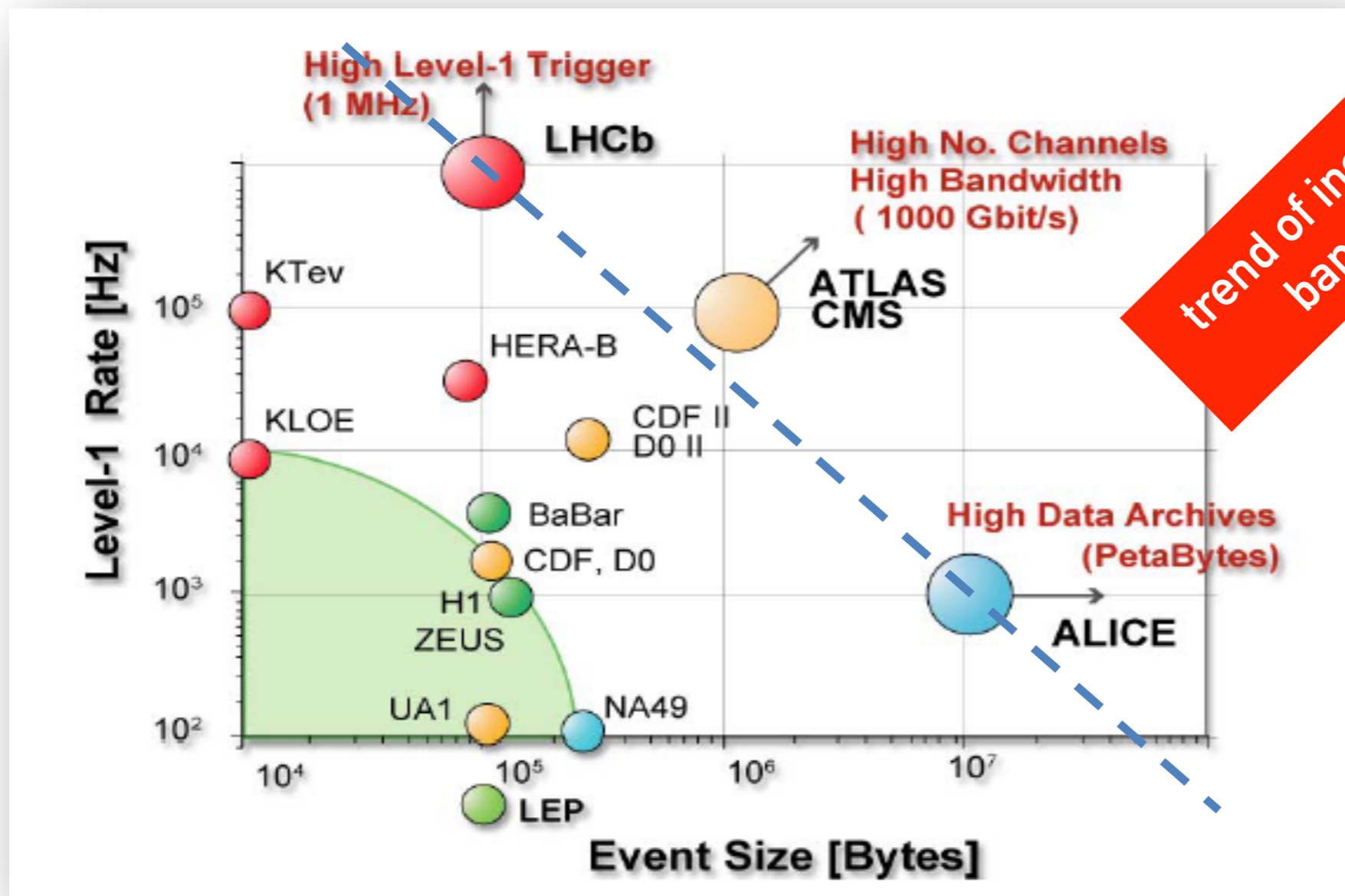
➡ **Who controls the flow?**

**push**

**pull**

(detectors, digitizers, front-end pipelines, readout buffers, switching networks, processor farms; LV1 ms, LV2 ms, LV3 s)

➡ **Buffers are not the "final solution"**
  ➡ Can overflow, with bursts and unusual event sizes
  ➡ In these cases
    ➡ discard data locally or
    ➡ exert "**back-pressure**", i. e. ask previous level(s) to block dataflow

➡ **Throughput optimization means avoiding dead-time due to back-pressure**
  ➡ using knowledge of the input buffer state

➡ **Who controls the flow?**

➡ **FE (push) or EB (pull)**
  ➡ **Push**: Events are sent as soon as data are available to the sender (for example round-robin algorithm) ==> Busy or Throttle
  ➡ **Pull** : events are required by a given destination processes (may need an event manager) ==> back-pressure
  ➡ **Push-Pull** ==> busy and back-pressure

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



*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

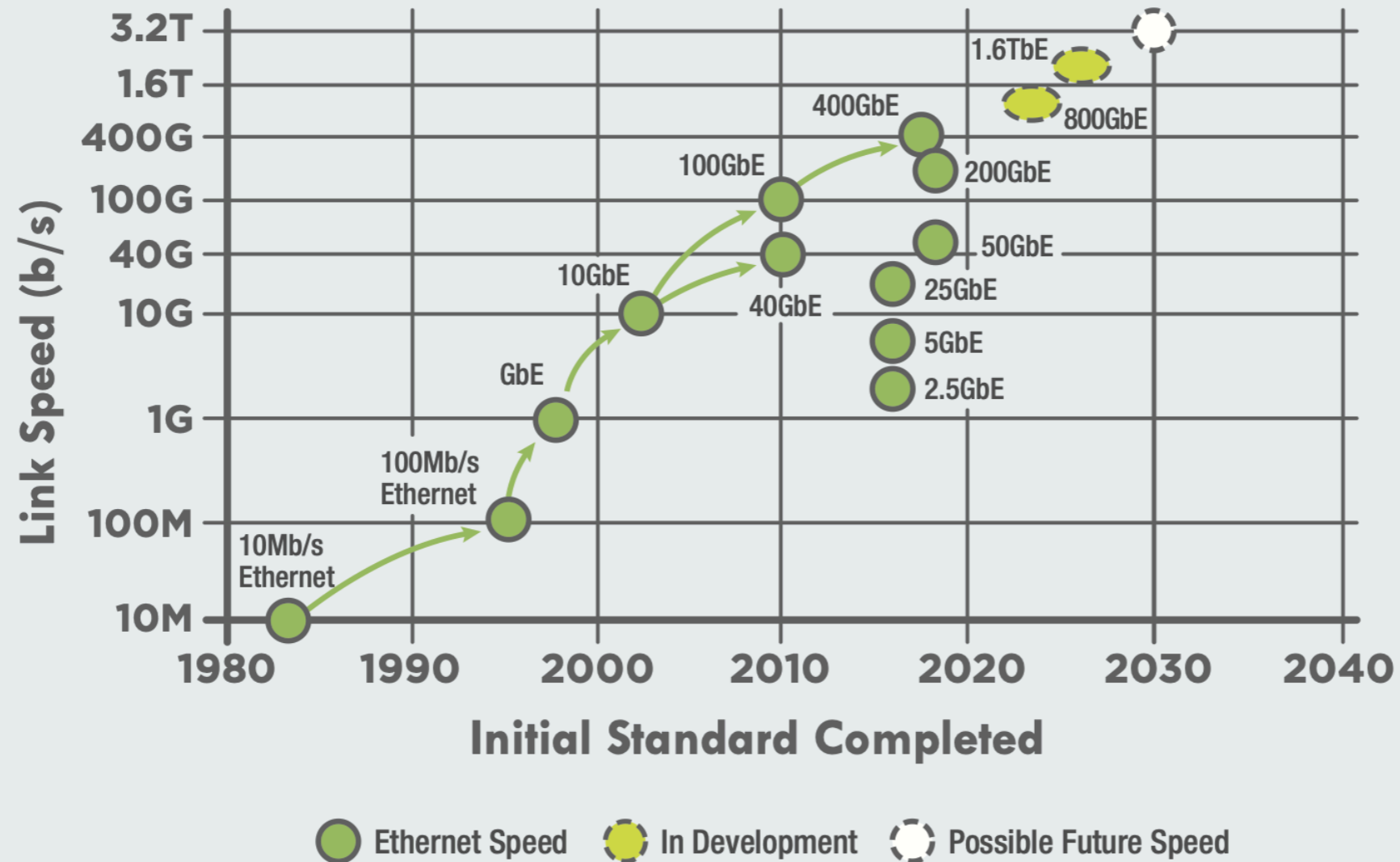**As the data volumes and rates increase, new architectures need to be developed**

*Courtesy of A.Cerri*

*Courtesy of A.Cerri*

# GENERAL T/DAQ TRENDS

➡ **Increasing readout channels, and front-end cards, distributed in multi-level three structure**

➡ **Integrate synchronous low-latency in Front-End**

  ➡ limitations do not disappear, but decouple (factorise)

➡ **Deal with dataflow instead of latency**

  ➡ **decouple** DAQ from High Level Triggers

  ➡ decouple dataflow from storage, with temporary buffers

  ➡ Use COTS network and processing

➡ **Use networks as soon as possible**

  ➡ toward commercial bidirectional point-to-multipoint architecture

➡ **Use "network" design already at small scale**

  ➡ easily get high performance with commercial components

➡ **Increase data aggregation at the Event Building**

  ➡ reducing request rates on DAQ software

  ➡ per time-frame, per orbit instead of per-event

42 Years of Microprocessor Trend Data

memory wall

Moore's law

power wall

multi-processing

Transistors (thousands)

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

Frequency (MHz)

Typical Power (Watts)

Number of Logical Cores

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!

**This requires fundamental re-write/ optimization of our software**
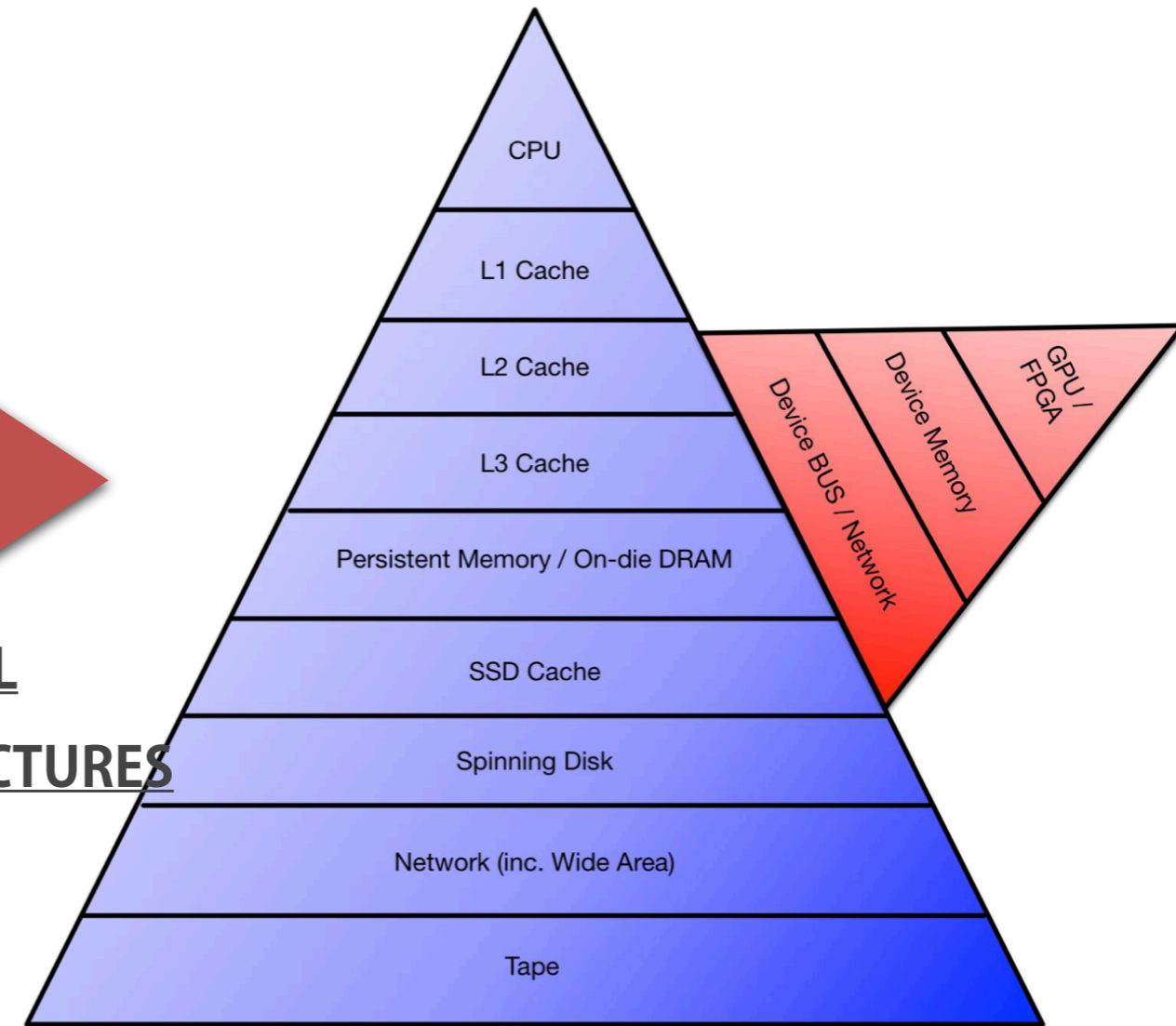
*Read: HPC computing*

"We're approaching the limits of computer power – we need new programmers now"
John Naughton, Guardian

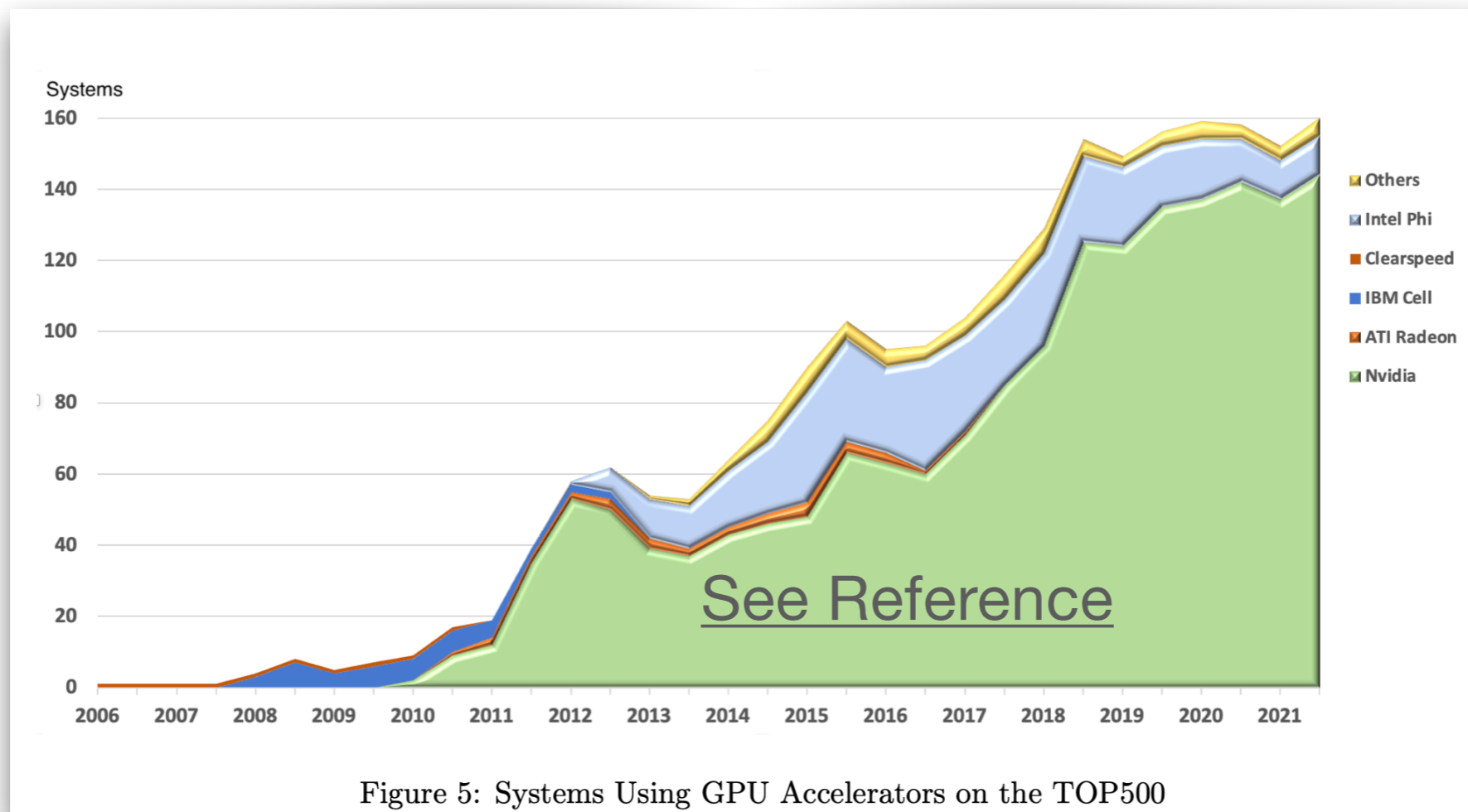*See LP-2022 slides from Graeme Stewart*



**EXPLOSION OF NOVEL**

**COMPUTER ARCHITECTURES**

➡ **Exploiting CPU hardware in new architectures**
  ➡ more complicated programming (vectorisation, memory sharing…)
➡ **Exploit more efficiently instruction level parallelism (ILP)**
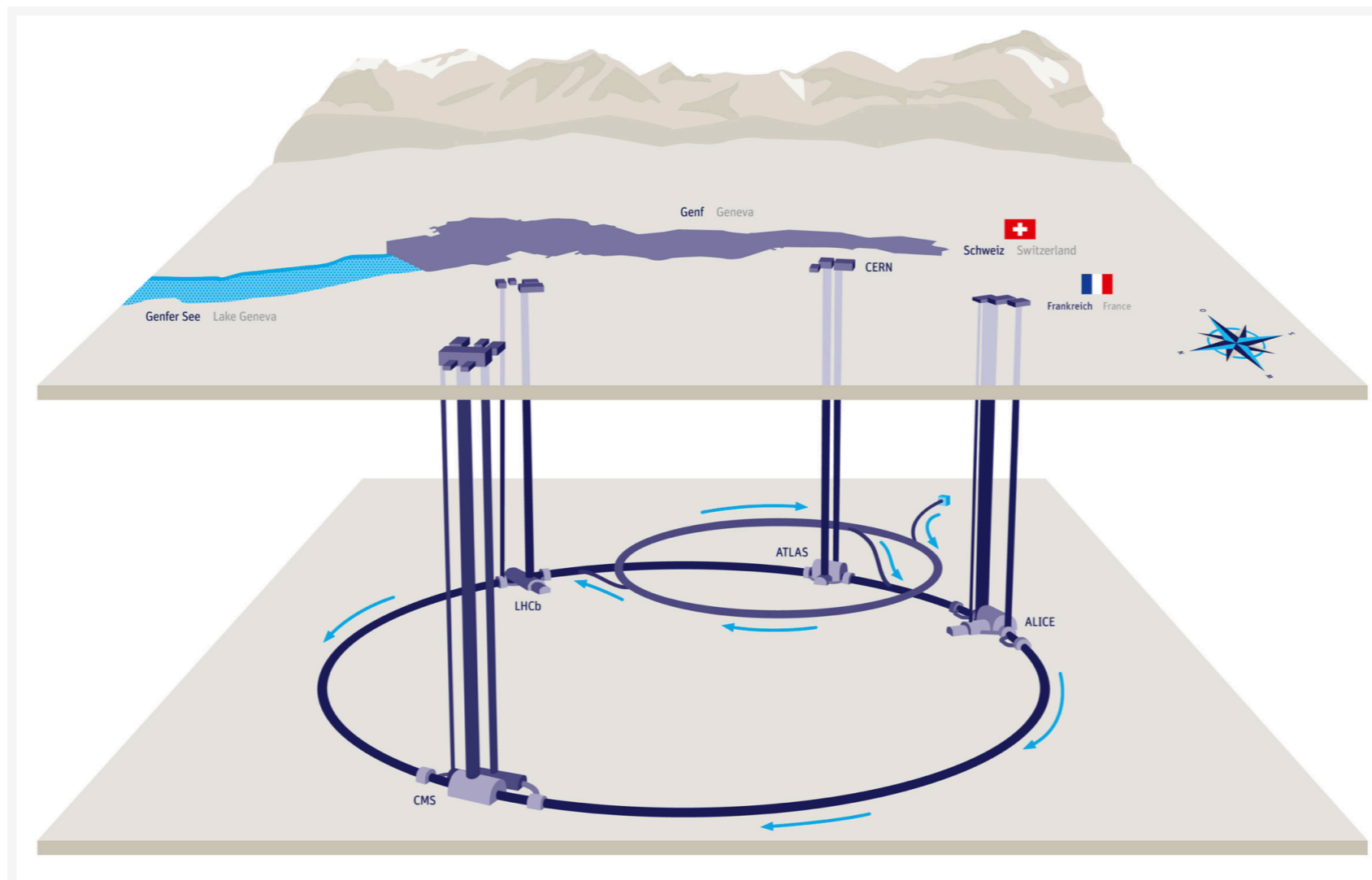
➡ **Scientific computing is the third paradigm, complementing theory and experiment**
- ➡ Global scientific facilities (e.g., LIGO, LHC, Vera Rubin Observatory, the Square Kilometer Array)

➡ **Future trends in HPC focusing on:**
- ➡ Rise of massive scale commercial clouds (Google Kubernetes, serverless computing,….)
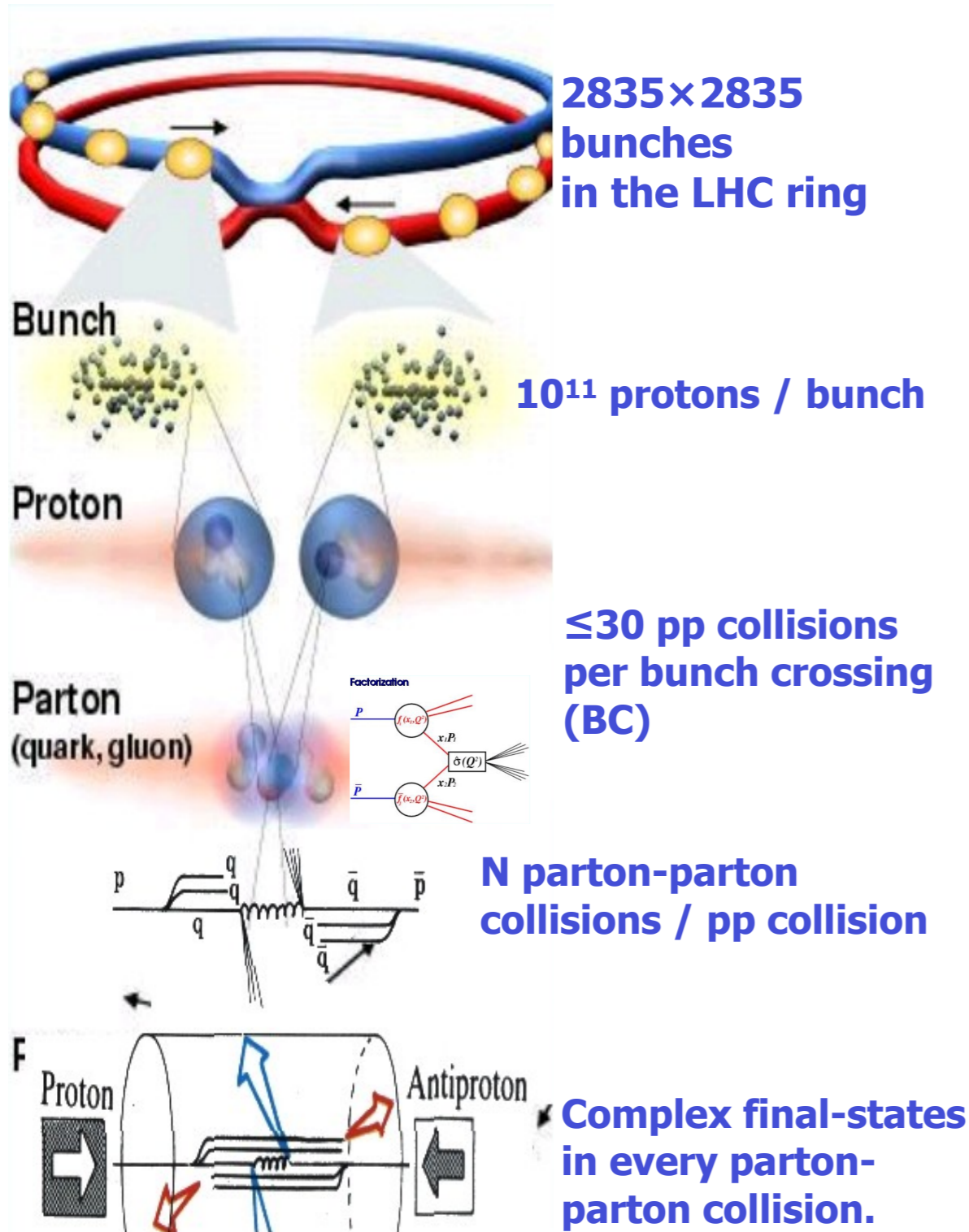- ➡ Evolution of semiconductor technology (chip size and packaging, see Amazon Graviton 3)



Figure 5: Systems Using GPU Accelerators on the TOP500

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

➡ **Examples of small experiments with their limits**

➡ **Overview of LHC experiments and their upgrade**

➡ **Future TDAQ systems (Dune/Proto-Dune)**

**2835×2835 bunches in the LHC ring**

**$10^{11}$ protons / bunch**

**≤30 pp collisions per bunch crossing (BC)**

**N parton-parton collisions / pp collision**

**Complex final-states in every parton-parton collision.**

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

**Search for rare events overwhelmed in abundant low-energy particles**

**Three major challenges for T/DAQ**

➡ **Face High Luminosity:**

➡ **fast electronics,** to resolve in time

➡ fine granularity detector, to resolve in space ➠ **high data volume**

➡ **Search for rare physics:**

➡ high rejection or large data collection

➡ **Be radiation resistant:**

➡ very costly for electronics ==> survive up to 100 Mrad= 1 MGy
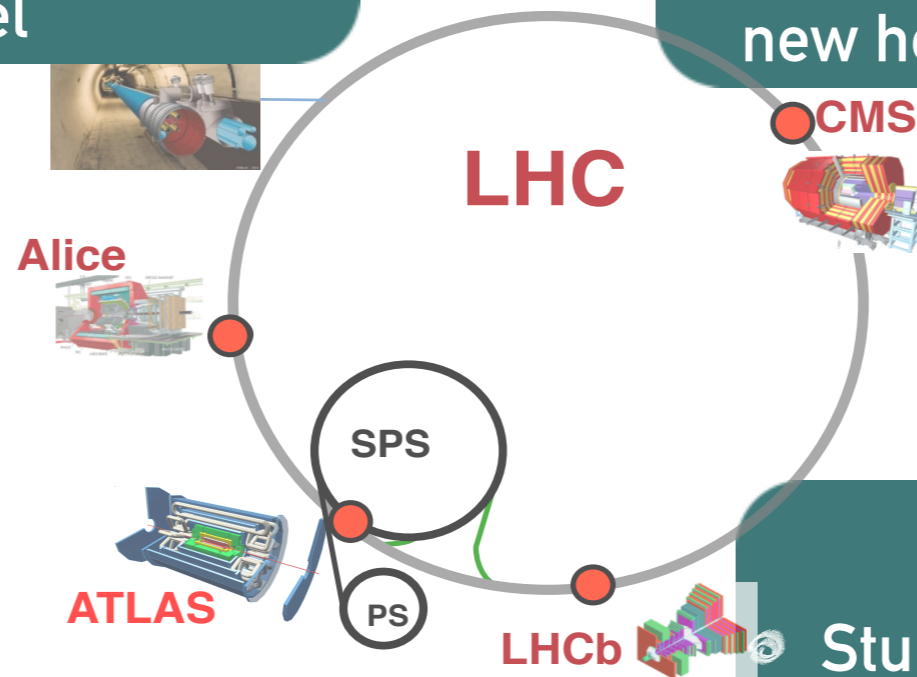
# LHC EXPERIMENTS FOR A DISCOVERY MACHINE

**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
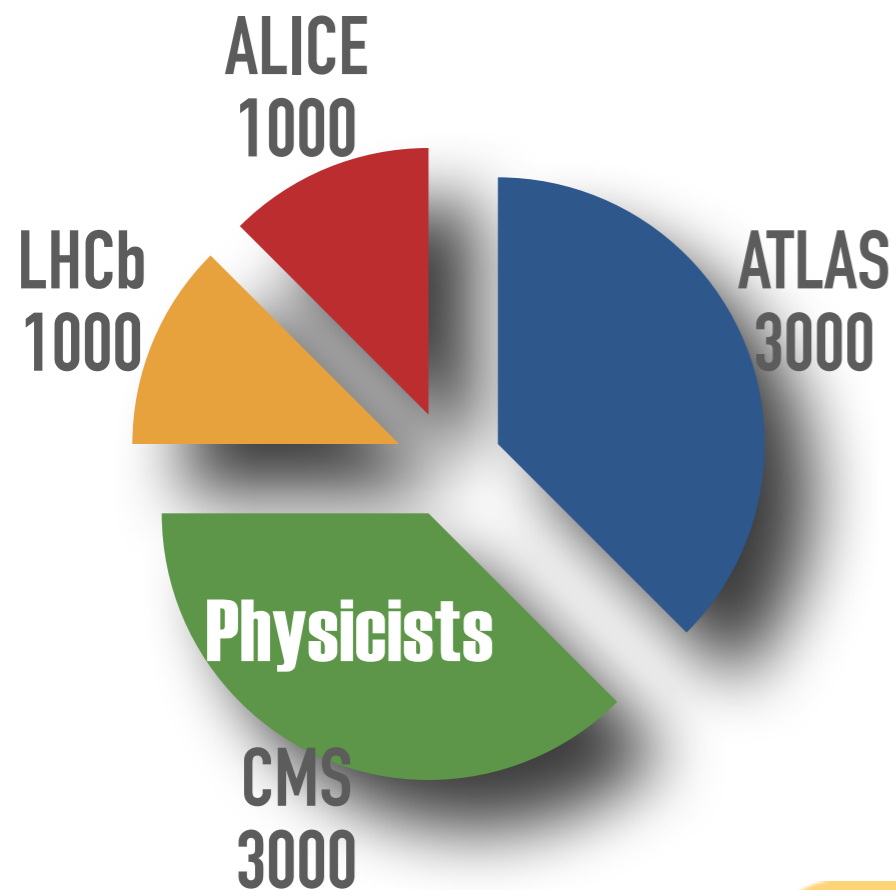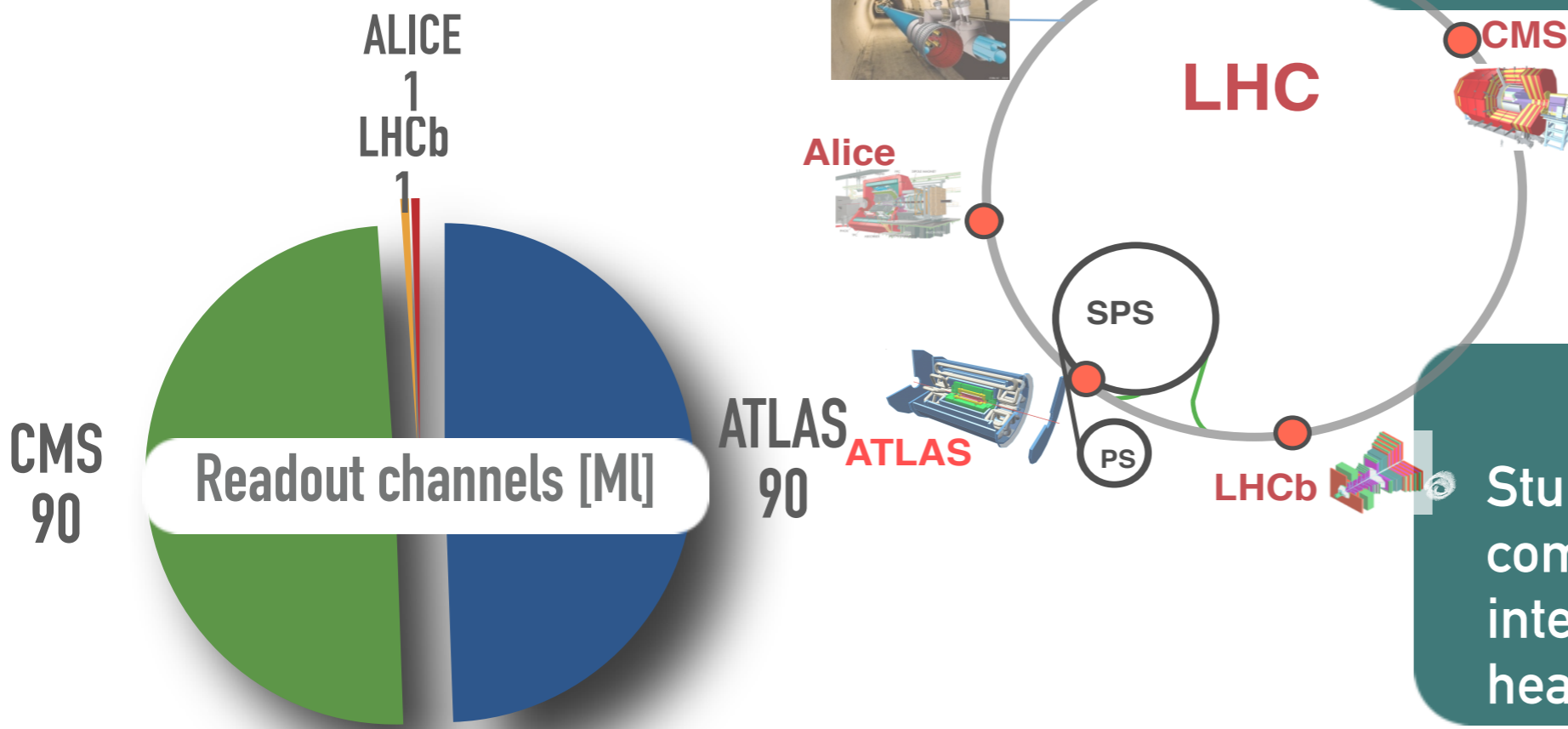
## LHCb

- 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

Alice
CMS
ATLAS
SPS
PS
LHCb

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

# 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**
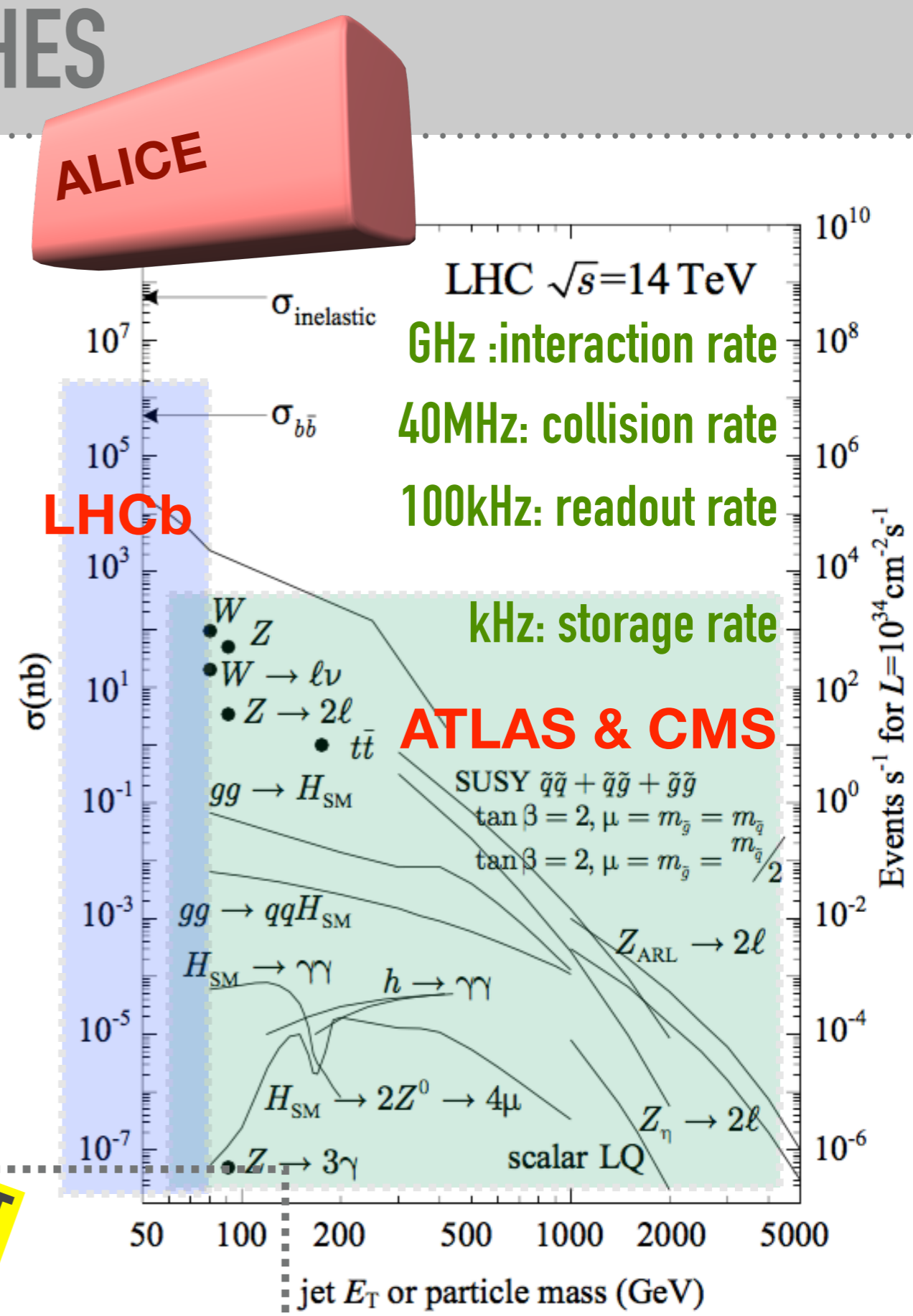  - search complex topologies of b-quark decays

- **ALICE: heavy-ion collisions ~2000 mb**
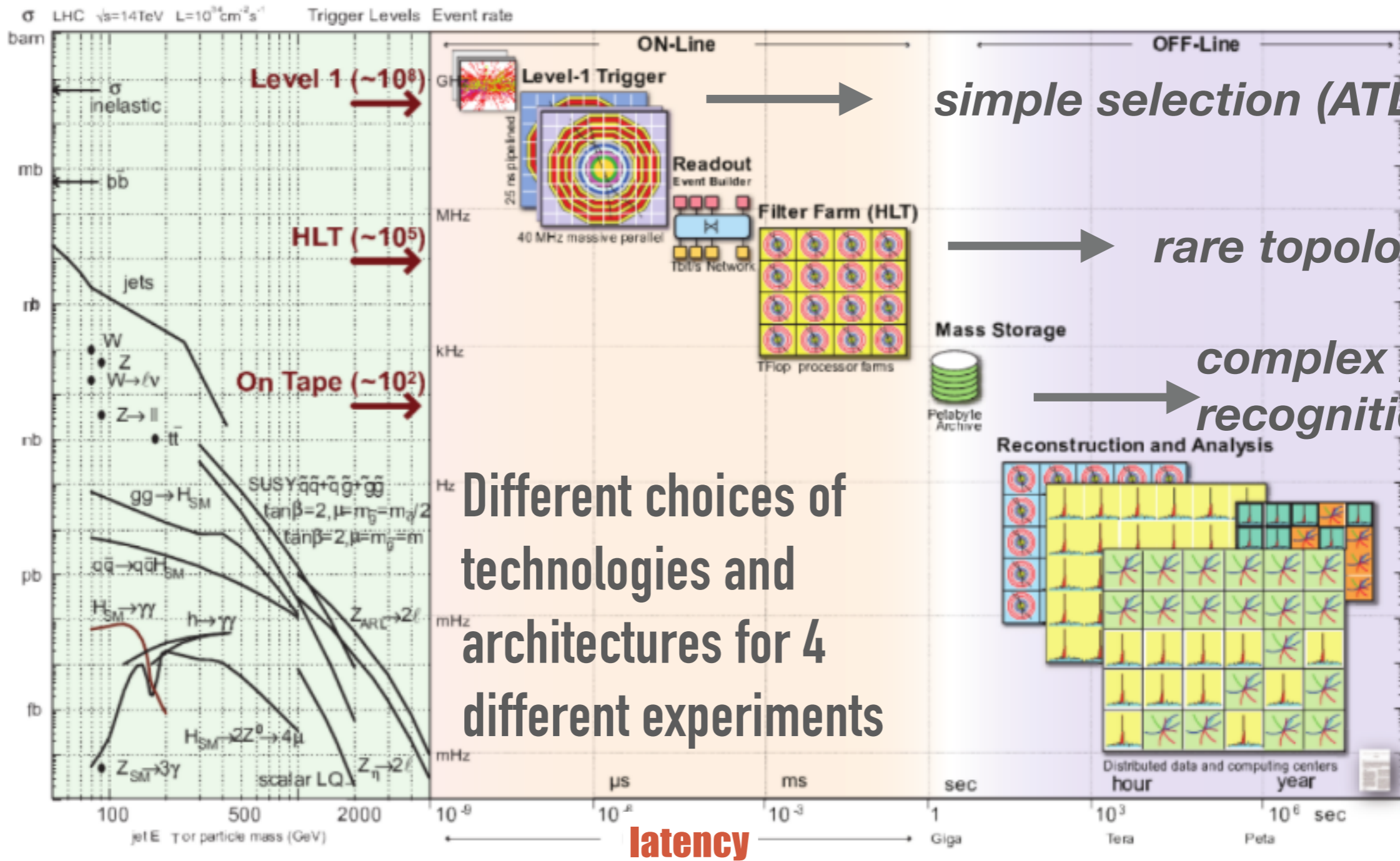  - search in high energy density

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

**DIFFERENT**



ALICE

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

**GHz :interaction rate**
**40MHz: collision rate**
**100kHz: readout rate**
**kHz: storage rate**

LHCb

ATLAS & CMS

$\sigma_{\text{inelastic}}$

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

$W$
$Z$
$W \to \ell\nu$
$Z \to 2\ell$
$t\bar{t}$
$gg \to H_{\text{SM}}$
$gg \to qqH_{\text{SM}}$
$H_{\text{SM}} \to \gamma\gamma$
$h \to \gamma\gamma$
$H_{\text{SM}} \to 2Z^0 \to 4\mu$
$Z \to 3\gamma$

SUSY $\tilde{q}\tilde{q} + \tilde{q}\tilde{g} + \tilde{g}\tilde{g}$
$\tan\beta = 2, \mu = m_{\tilde{g}} = m_{\tilde{q}}$
$\tan\beta = 2, \mu = m_{\tilde{g}} = m_{\tilde{q}}/2$
$Z_{\text{ARL}} \to 2\ell$
$Z_{\eta} \to 2\ell$
scalar LQ

$\sigma$(nb)

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

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

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

*simple selection (ATLAS, CMS)*

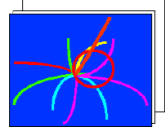*rare topology (LHCb)*

*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

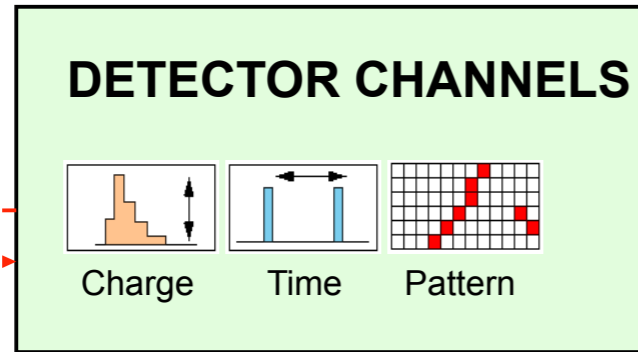# MANY PLAYERS, COMPLEX TDAQ ARCHITECTURES
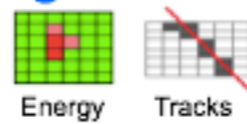
**Buffering and parallelism**
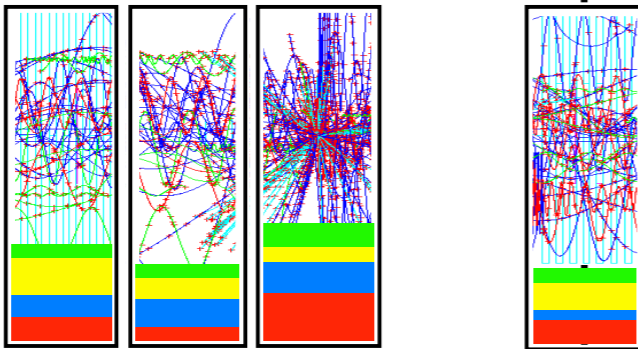
**Maximum 1-2% deadtime**

**40 MHz COLLISION RATE**

**Level-1**

**Readout Buffers**

**Event building**

**Event filtering**

**Petabyte archive**

**DETECTOR CHANNELS**

Charge  Time  Pattern

Energy  Tracks

**SWITCH NETWORK**

**Computing Services**

**High speed electronics**

**Readout links and buffering**

**Large data network with dedicated technology**

**Dedicated PC farms**

**Readout**

**DAQ**

L1/Readout

HLT/DAQ

**Level-1 triggers**
➡ Set max Readout rate
➡ Hardware, synchronous
➡ Readout parallelism
➡ Latency ~ µsec/event

**Higher level triggers**
➡ Set max storage rate
➡ Software, asynchronous
➡ Event parallelism
➡ Latency < 1 sec/event

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

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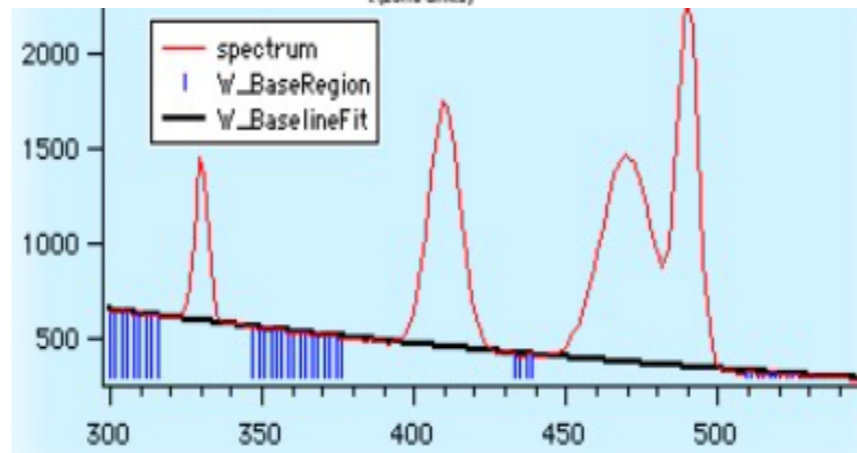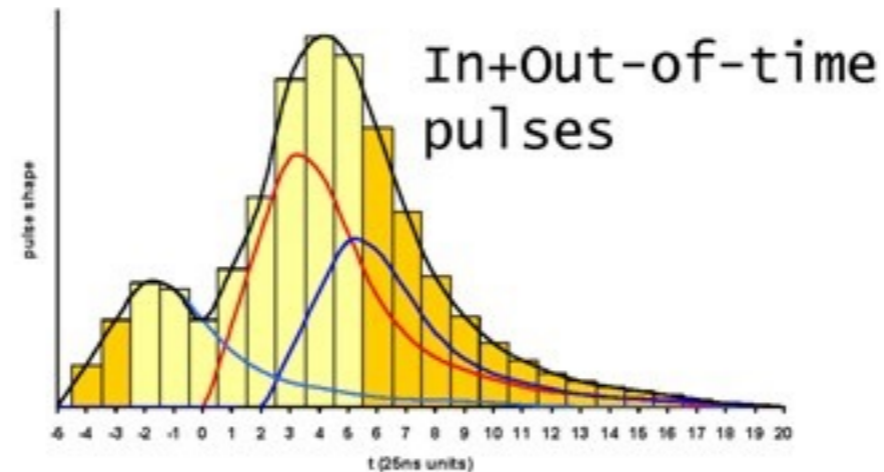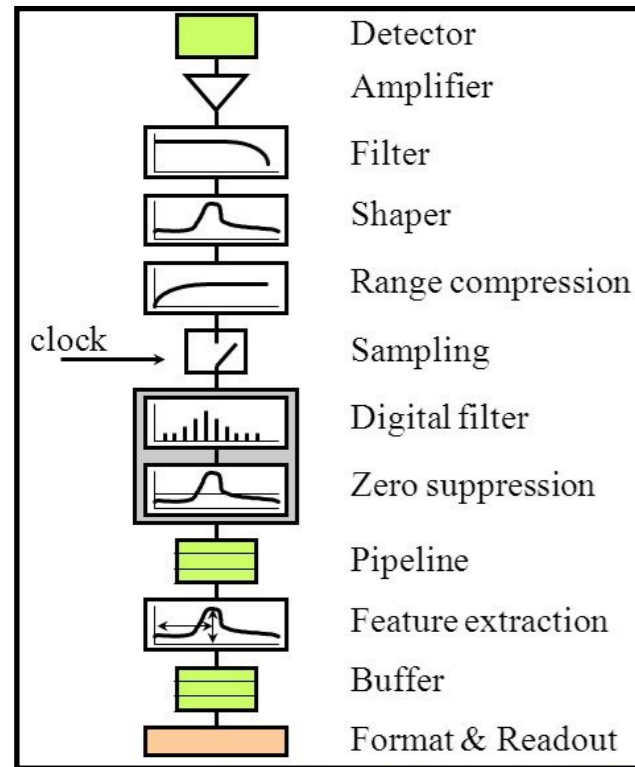➡ **Massively parallel**
➡ **Bunch Crossing identification capability**

# Fast, robust electronics

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

**Latency dominated by cable/transmission delay**

## Tight design constraints for trigger & FE
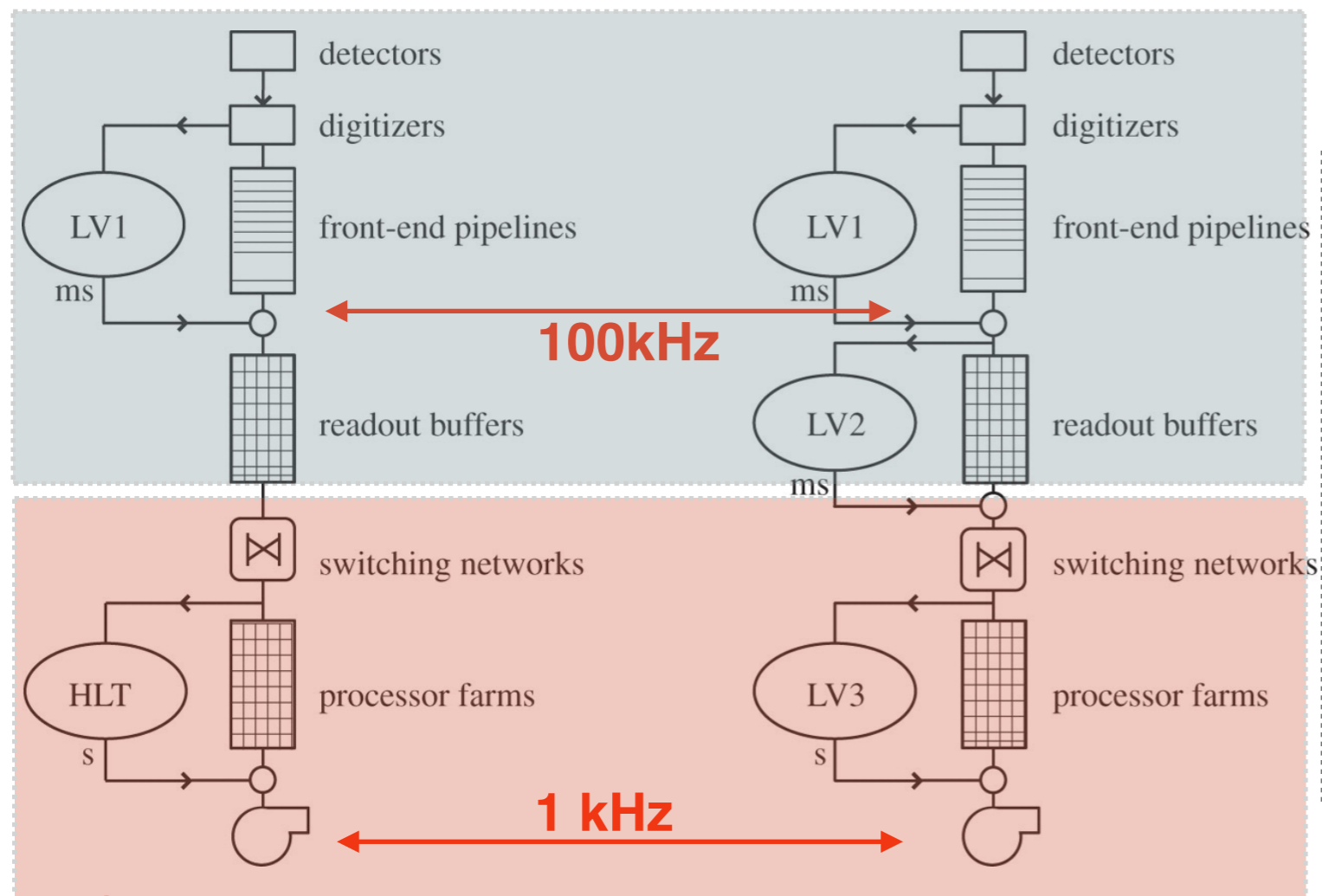


**ATLAS Liquid Argon calorimeter**



## Avoid

➡ **Electronic pile-up**
  ➡ source of dead-time
  ➡ distortion in pulse

➡ **In-time pile-up**
  ➡ more collisions/BC
  ➡ Baseline subtraction

➡ **Out-of-time pile-up**
  ➡ BC-identification capability
  ➡ peak finder algorithms

## Make it easier with fast, low occupancy and digital detectors

# HLT/DAQ REQUIREMENTS

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

**Prefer use of PCs (linux based), Ethernet protocols, standard LAN, configurable devices**



100kHz

1 kHz

**DAQ+HLT system**

### ATLAS/CMS Example

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

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

  ➡ HLT farm: 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)

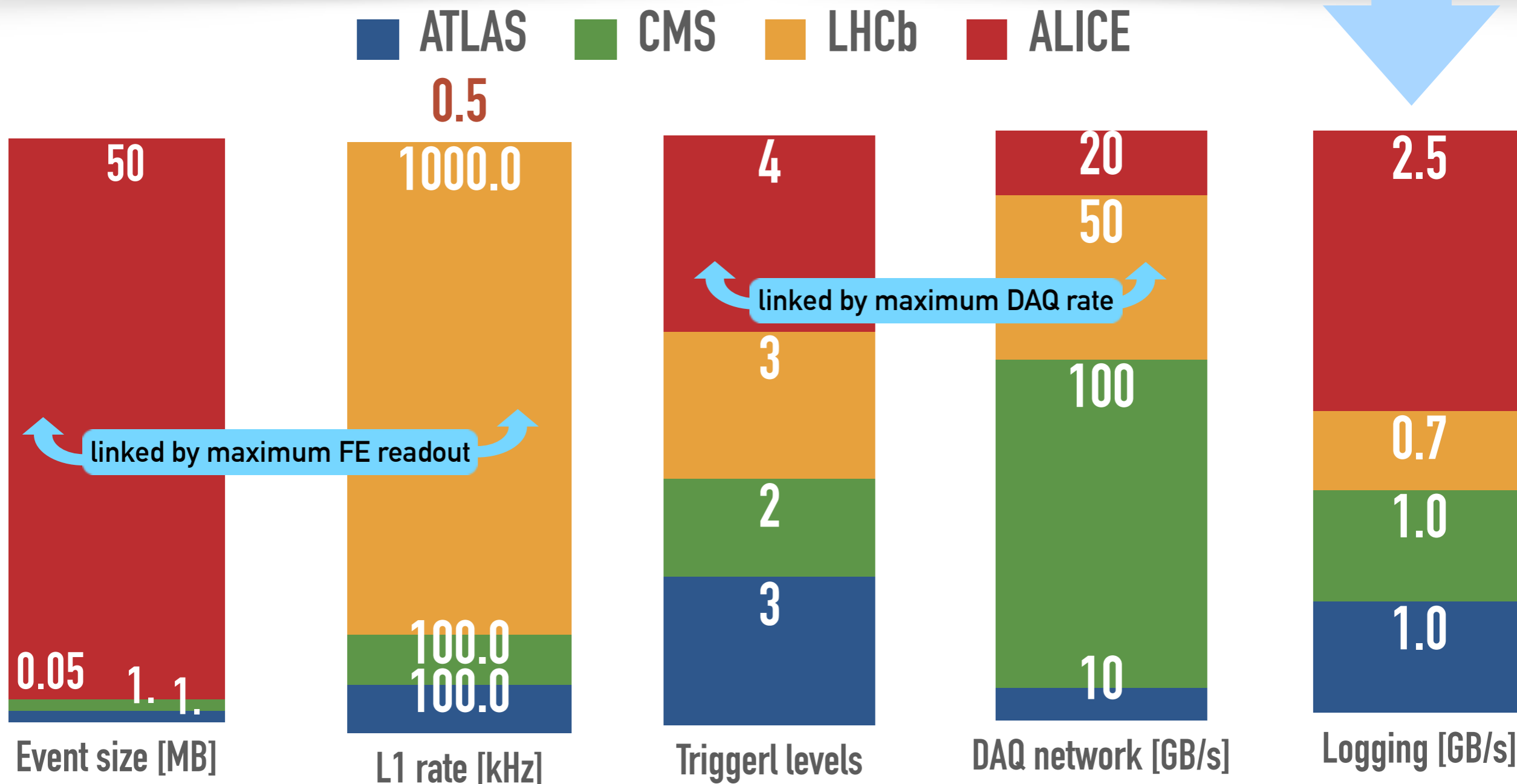**Very large uncertainties to take into account!**

# 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 $t\bar{t}$ event in ATLAS ITK
at $\langle\mu\rangle=200$

**Physics program for the future is towards more rare processes at the same energy scale**

## 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
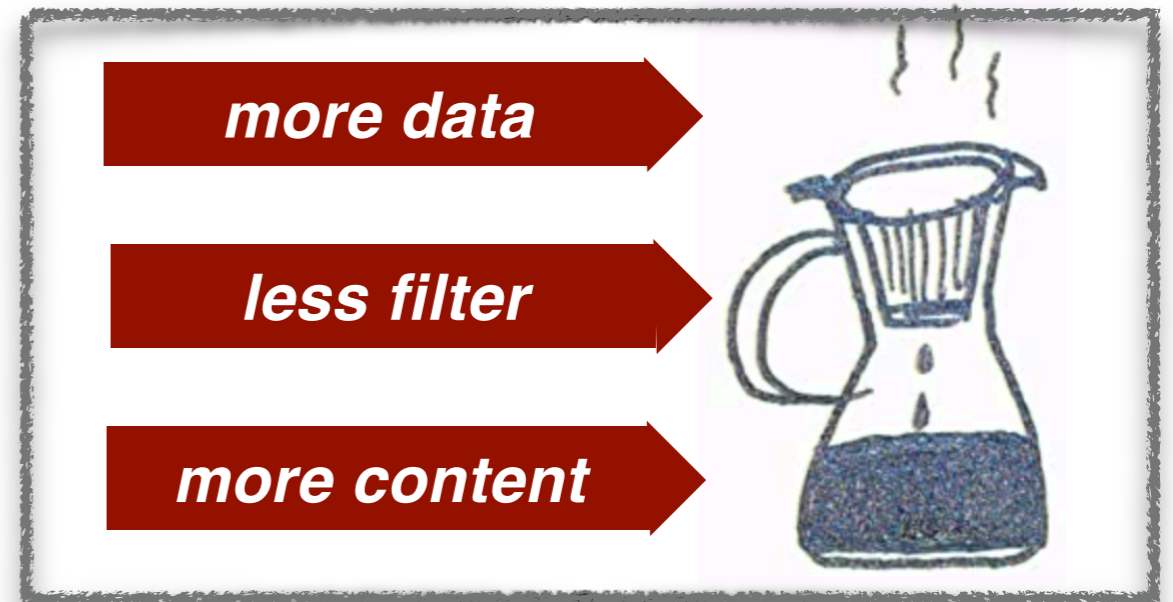
*ATLAS/CMS numbers*

*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] (y-axis: 0, 1000, 2000, 3000, 4000, 5000, 6000, 7000, 8000, 9000)

pileup interaction multiplicity (x-axis: 40, 60, 80, 100, 120, 140, 160)

# 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

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

Human Genome
8000 Eb/year

ATLAS/CMS
260 Eb/year

Global Internet
2800 Eb/year

LHCb
1000 Eb/year

2021    2023    2025    2027    2029    2031

*See Openlab workshop*

**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)

# COMPARE 4 EXPERIMENTS

*How to maximise physics acceptance*



spot the differences

➡ **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 leptons @L1**

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

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

## 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

## 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

Infiniband

Share (%)

2002    2014    2018

**Choose best prize/bitps!**

**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

➡ **Precision measurements and rare decays in the B system**
  ➡ Large production ($\sigma_{BB} \sim 500$ µb), but still $\sigma_{BB}/\sigma_{Tot} \sim 5\times10^{-3}$
  ➡ 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

# TRIGGER-LESS?



**Run1** / **Run3**

30 MHz inelastic event rate (full rate event building)

1 MHz / 30 MHz

*40Tbit/s*

Software High Level Trigger

Full event reconstruction, inclusive and exclusive kinematic/geometric selections — HLT-1

150 kHz / 1MHz

*1-2 Tbit/s*

Buffer events to disk, perform online detector calibration and alignment

Add offline precision particle identification and track quality information to selections — HLT-2

Output full event information for inclusive triggers, trigger candidates and related primary vertices for exclusive triggers

12.5 kHz / 50 kHz

2-5 GB/s to storage

*80 Gbit/s*

**From Run1 to Run3, TDAQ system evolved to handle more readout rate**

**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, δp/p ~ 15%
- Forward tracking — $p_T$ > 500 MeV, δp/p ~ 0.5%
- PV finding
- Rate reducing cuts — Output < 1 MHz
- Muon Identification
- Simplified Kalman fit
- Particle Identification

arXiv:2105.04031

**150kB x 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)

➡ **Data pushed to the Event Building with massive link usage:**
  - ➡ ~10,000 GBT (4.8 Gb/s, rad-hard)

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

**Inside Cavern**

**Surface data centre**



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

*Ref for PCIe40*

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

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

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)

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

➡ **19 different detectors**

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

  ➡ Time Projection Chamber (**TPC**): very high occupancy, and slow response

➡ **Large event size (> 40MB)**

  ➡ TPC producing 90% of data

➡ **Complex event topology**

  ➡ low trigger rate: ~ kHz



➡**Challenges for TDAQ and evolution:**

  ➡ detector readout: up to ~50 GB/s ==> x100 for Run3

  ➡ storage:  1.2 TB/s (Pb-Pb)  ==> x100 for Run3

**How can we increate the readout rate, when it's close to TPC readout?**

**Reconstruct TPC data in continuous readout**

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

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

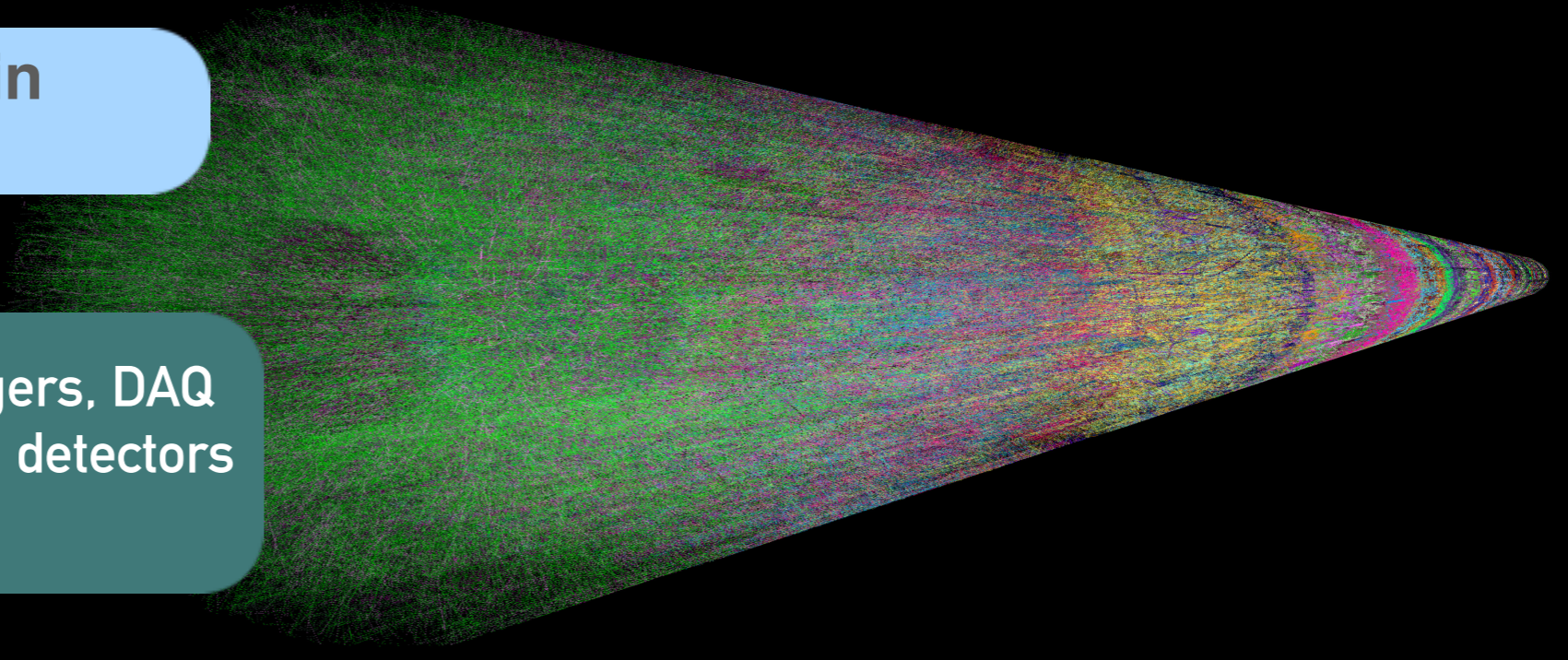  ➡ 1 per LHC orbit, 89.4 $\mu$s: <u>~10 kHz</u>

➡ **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!**

# OUTLINE

➡ Examples of small experiments with their limits

➡ Overview of LHC experiments and their upgrade

➡ **Future TDAQ systems (Dune/Proto-Dune)**

➡ **The next generation project for neutrino physics**

   ➡ the experiment does not exist (ready for 2030)

   ➡ the TDAQ of the experiment does not exist

➡ **Consider here design inputs:**

   ➡ have a broad understanding of what the experiment wants to achieve

   ➡ understand the detection principles and front-end electronics

   ➡ understand the constraints in which the TDAQ will live

➡ **http://dunescience.org**

➡ **DUNE Collaboration : 1317 members, 208 institutions, 33 Countries**

➡ **Strong International partnership to build a mega neutrino science project based in US**

➡ **see recent CERN colloquium**

A view of the ProtoDUNE cryostat at CERN (Image: CERN)

➡ **Two detectors on a muon-neutrino beam @Long-Baseline Neutrino Facility**
- ➡ One near the source of the beam, at Fermilab (**ND**), to characterise the beam & systematics
- ➡ One, much larger, 1300 km downstream, 1.48 km underground (**FD**)
  - ➡ Massive Liquid Argon Time Projection Chambers (70-kton, slow) + photon detectors (fast)
  - ➡ *the best particle imaging capability*

➡ **No quick access and no large host lab in the area !**

➡ **Prototypes at CERN Neutrino Platform (proto-DUNE)**
- ➡ 2 prototypes, 1/20th the size of planned DUNE
  - ➡ *the largest liquid-argon neutrino detector in the world!*
- ➡ Collected 4M events in 2018- 2020 from both cosmic rays and a beam

➡ **Extended physics cases:**

➡ Origin of matter: measure **neutrino oscillations** on large distances and unfold CPV from matter effects

➡ <u>trigger</u>: neutrino beam -> external trigger possible

➡ Unification of forces: search for **proton decay**

➡ <u>trigger</u>: very local, rare signature

➡ Black hole formation: observe neutrinos from **supernova collapse**

➡ Very distributed, rare signature

➡ **Extended physics cases:**

- ➡ Origin of matter: measure **neutrino oscillations** on large distances and unfold CPV from matter effects

  - ➡ <u>trigger</u>: neutrino beam -> external trigger possible

- ➡ Unification of forces: search for **proton decay**

  - ➡ <u>trigger</u>: very local, rare signature

- ➡ Black hole formation: observe neutrinos from **supernova collapse**

  - ➡ Very distributed, rare signature

➡ **TDAQ active at "all" times, mixing readout strategies**

- ➡ local readout for photon detectors, sampling @ 150 MHz

- ➡ continuous readout for TPC, sampling @ 2 MHz

- ➡ post-readout system combines data fragments into time windows of interesting detector regions

  - ➡ data reordering appears to be the biggest CPU consumer

➡ **Extended physics cases:**

- ➡ Origin of matter: measure **neutrino oscillations** on large distances and unfold CPV from matter effects
  - ➡ <u>trigger</u>: neutrino beam -> external trigger possible
- ➡ Unification of forces: search for **proton decay**
  - ➡ <u>trigger</u>: very local, rare signature
- ➡ Black hole formation: observe neutrinos from **supernova collapse**
  - ➡ Very distributed, rare signature

➡ **TDAQ active at "all" times, mixing readout strategies**

- ➡ local readout for photon detectors, sampling @ 150 MHz
- ➡ continuous readout for TPC, sampling @ 2 MHz
- ➡ post-readout system combines data fragments into time windows of interesting detector regions
  - ➡ data reordering appears to be the biggest CPU consumer

➡ **Adding all up, TDAQ has to sustain readout of ~5 TB/s**

- ➡ TPC: 384 k channels (12 bit ADC) @ 2 MHz = 9.2 Tb/s (dominates)

➡ **Extended physics cases:**

- ➡ Origin of matter: measure **neutrino oscillations** on large distances and unfold CPV from matter effects
  - ➡ <u>trigger</u>: neutrino beam -> external trigger possible
- ➡ Unification of forces: search for **proton decay**
  - ➡ <u>trigger</u>: very local, rare signature
- ➡ Black hole formation: observe neutrinos from **supernova collapse**
  - ➡ Very distributed, rare signature

➡ **TDAQ active at "all" times, mixing readout strategies**

- ➡ local readout for photon detectors, sampling @ 150 MHz
- ➡ continuous readout for TPC, sampling @ 2 MHz
- ➡ post-readout system combines data fragments into time windows of interesting detector regions
  - ➡ data reordering appears to be the biggest CPU consumer

➡ **Adding all up, TDAQ has to sustain readout of ~5 TB/s**

- ➡ TPC: 384 k channels (12 bit ADC) @ 2 MHz = 9.2 Tb/s (dominates)

➡ **Sounds very much like HL-LHC...**

➡ **Differently from LHC, time frames varies a lot**

  ➡ from few ms to ~100s for the supernova core collapse

  ➡ Data corresponding to a trigger can have size ranging << 1 GB to ~100 TB!

# DIFFERENCE WITH COLLIDER EXPERIMENTS

➡ **Differently from LHC, time frames varies a lot**

    ➡ from few ms to ~100s for the supernova core collapse

    ➡ Data corresponding to a trigger can have size ranging << 1 GB to ~100 TB!

➡ **The rate of events varies widely from few Hz to <<1/month**

➡ **Differently from LHC, time frames varies a lot**

- ➡ from few ms to ~100s for the supernova core collapse

- ➡ Data corresponding to a trigger can have size ranging << 1 GB to ~100 TB!

➡ **The rate of events varies widely from few Hz to <<1/month**

➡ **The trigger selection need to accumulate data from detectors over several seconds**

- ➡ readout needs very large buffers to accommodate the long decision latency

- ➡ fast storage of 3.5 TB with a sequential write performance @ 25 GBps

➡ **Differently from LHC, time frames varies a lot**

  ➡ from few ms to ~100s for the supernova core collapse

  ➡ Data corresponding to a trigger can have size ranging << 1 GB to ~100 TB!

➡ **The rate of events varies widely from few Hz to <<1/month**

➡ **The trigger selection need to accumulate data from detectors over several seconds**

  ➡ readout needs very large buffers to accommodate the long decision latency

  ➡ fast storage of 3.5 TB with a sequential write performance @ 25 GBps

➡ **Complexity and size are similar, but uptime is much larger (100% instead of ~30%)**

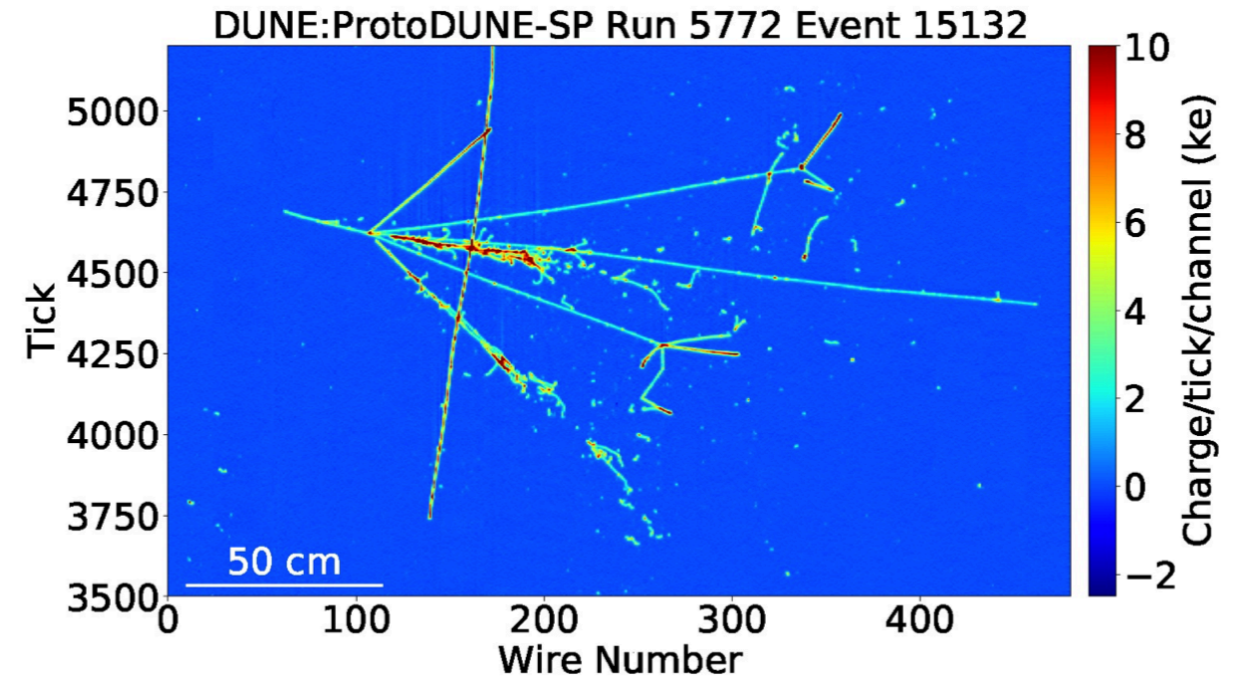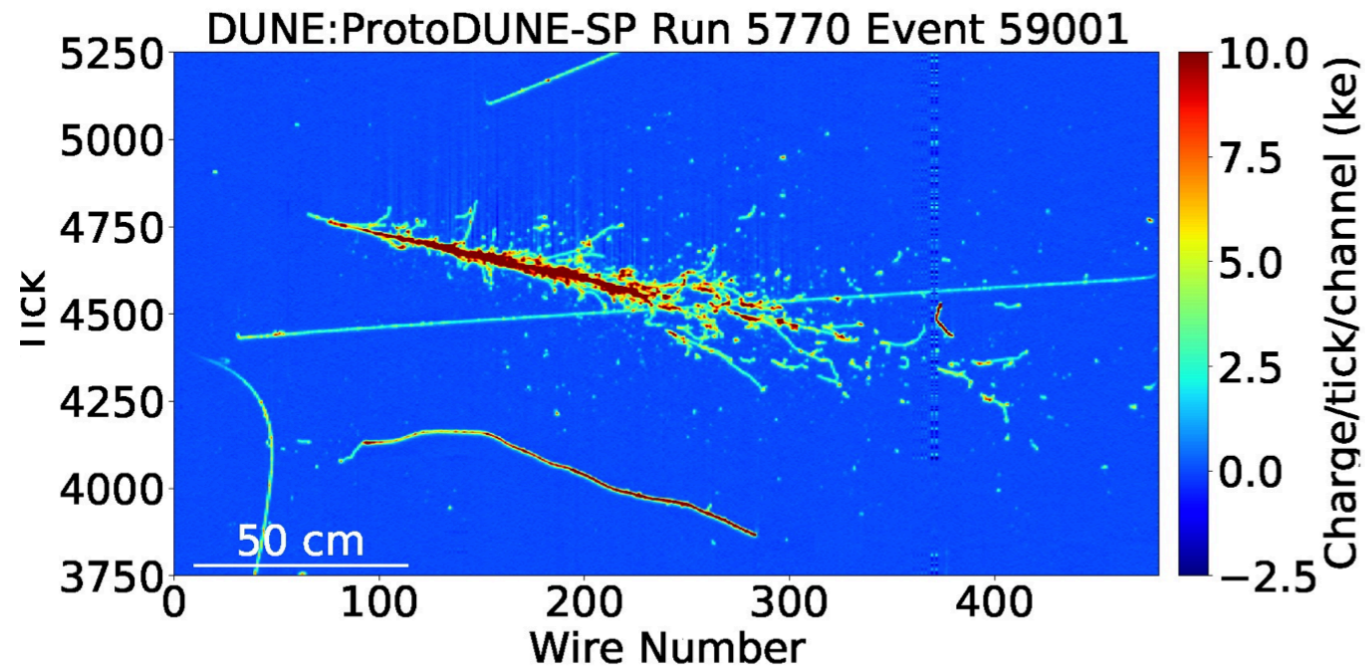➡ **Differently from LHC, time frames varies a lot**

- ➡ from few ms to ~100s for the supernova core collapse
- ➡ Data corresponding to a trigger can have size ranging << 1 GB to ~100 TB!

➡ **The rate of events varies widely from few Hz to <<1/month**

➡ **The trigger selection need to accumulate data from detectors over several seconds**

- ➡ readout needs very large buffers to accommodate the long decision latency
- ➡ fast storage of 3.5 TB with a sequential write performance @ 25 GBps

➡ **Complexity and size are similar, but uptime is much larger (100% instead of ~30%)**

➡ **Limited accessibility makes things more complex**

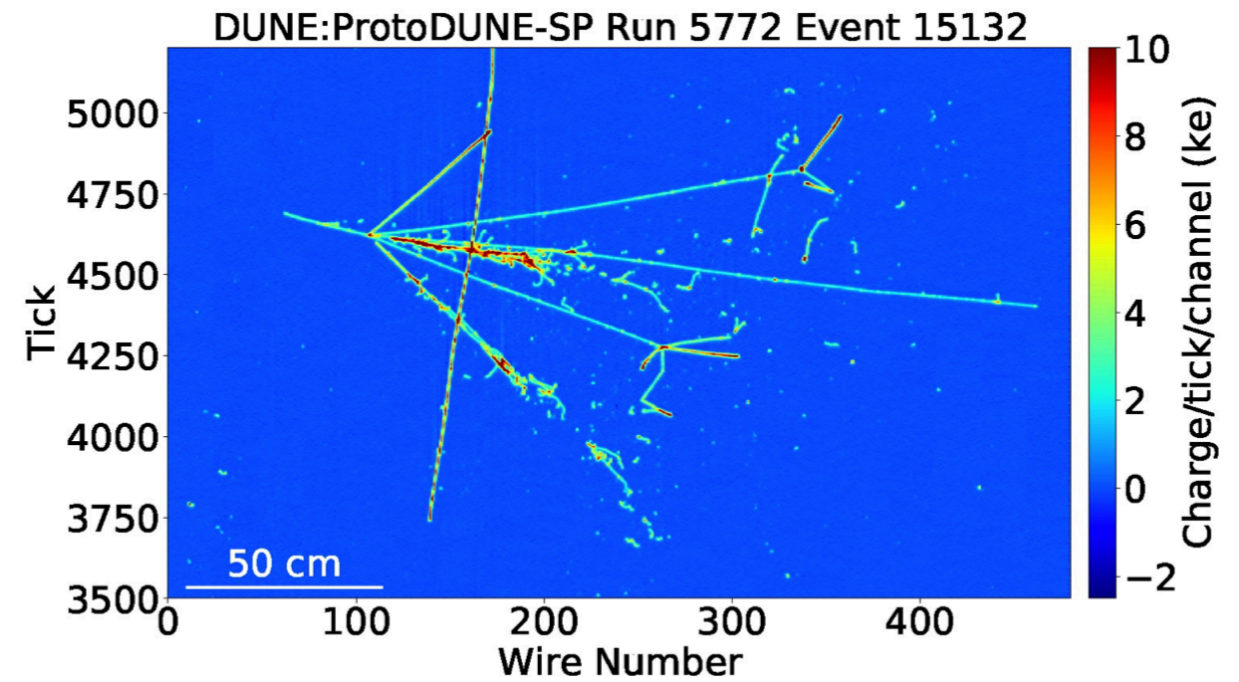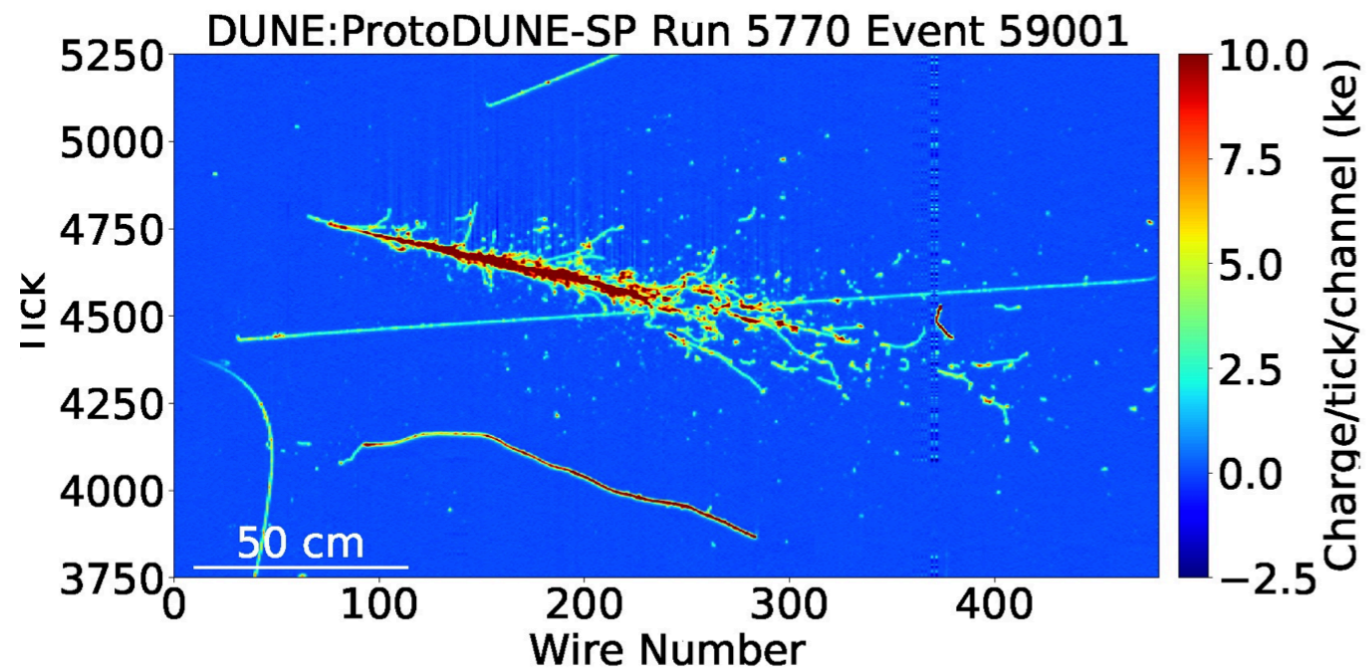➡ **Differently from LHC, time frames varies a lot**
  - ➡ from few ms to ~100s for the supernova core collapse
  - ➡ Data corresponding to a trigger can have size ranging << 1 GB to ~100 TB!

➡ **The rate of events varies widely from few Hz to <<1/month**

➡ **The trigger selection need to accumulate data from detectors over several seconds**
  - ➡ readout needs very large buffers to accommodate the long decision latency
  - ➡ fast storage of 3.5 TB with a sequential write performance @ 25 GBps

➡ **Complexity and size are similar, but uptime is much larger (100% instead of ~30%)**

➡ **Limited accessibility makes things more complex**

➡ **The control and monitoring system will have a predominant role for the success of the DUNE TDAQ**
  - ➡ Automated anomaly detection and recovery
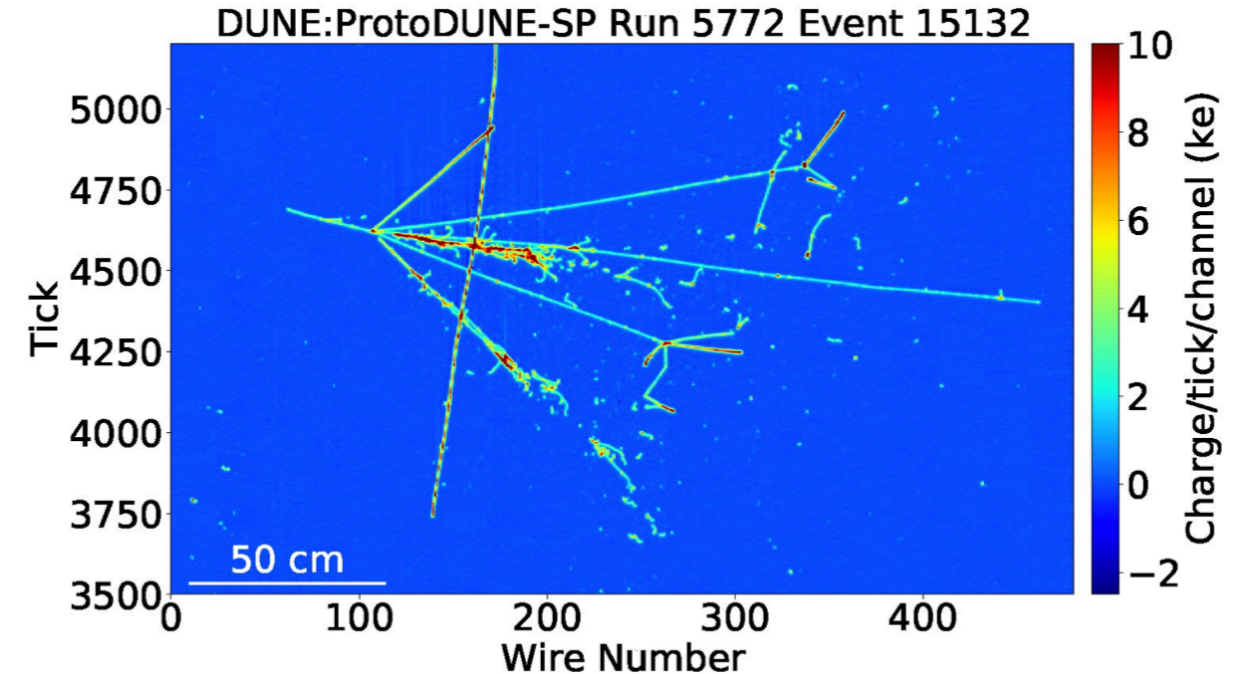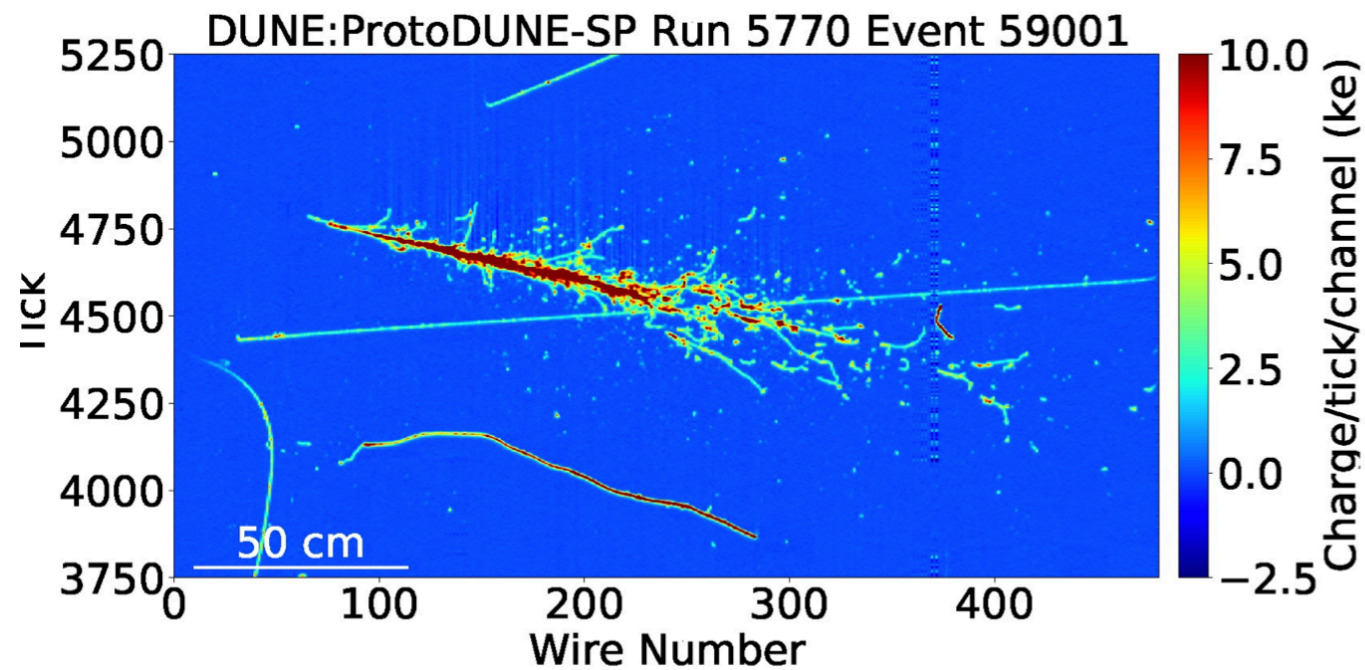  - ➡ Remote monitoring and control

DUNE:ProtoDUNE-SP Run 5770 Event 59001

DUNE:ProtoDUNE-SP Run 5772 Event 15132

➡ **Why?**

DUNE:ProtoDUNE-SP Run 5770 Event 59001

DUNE:ProtoDUNE-SP Run 5772 Event 15132

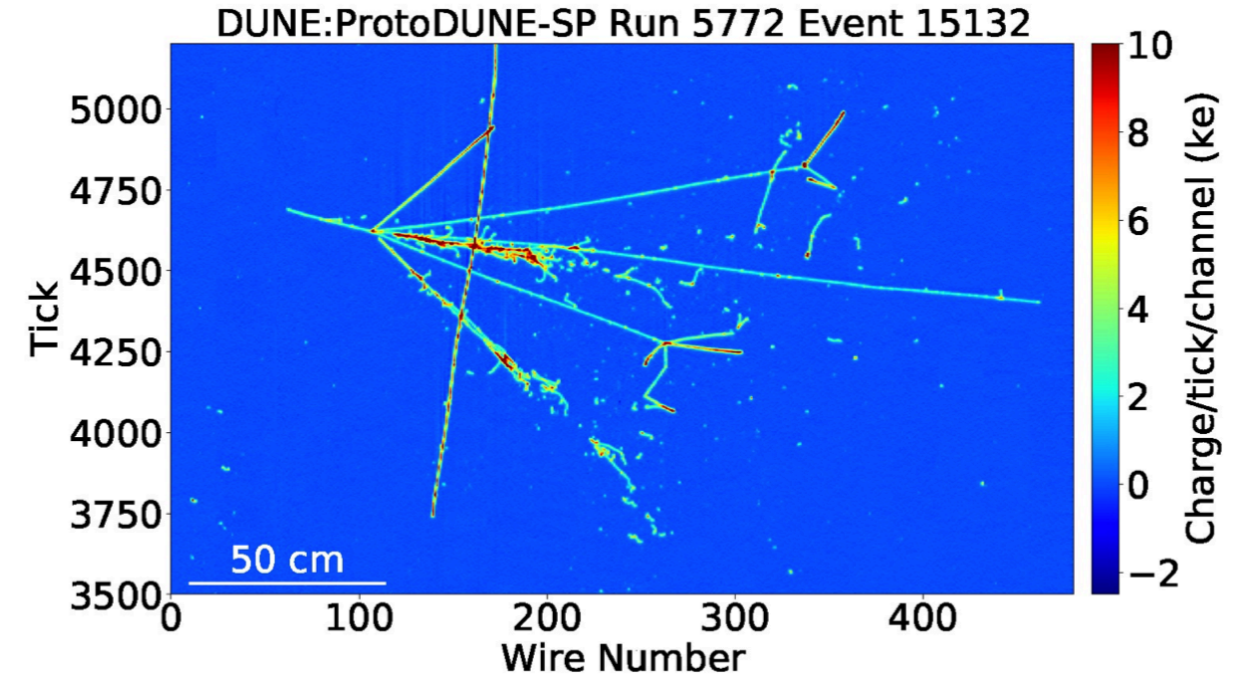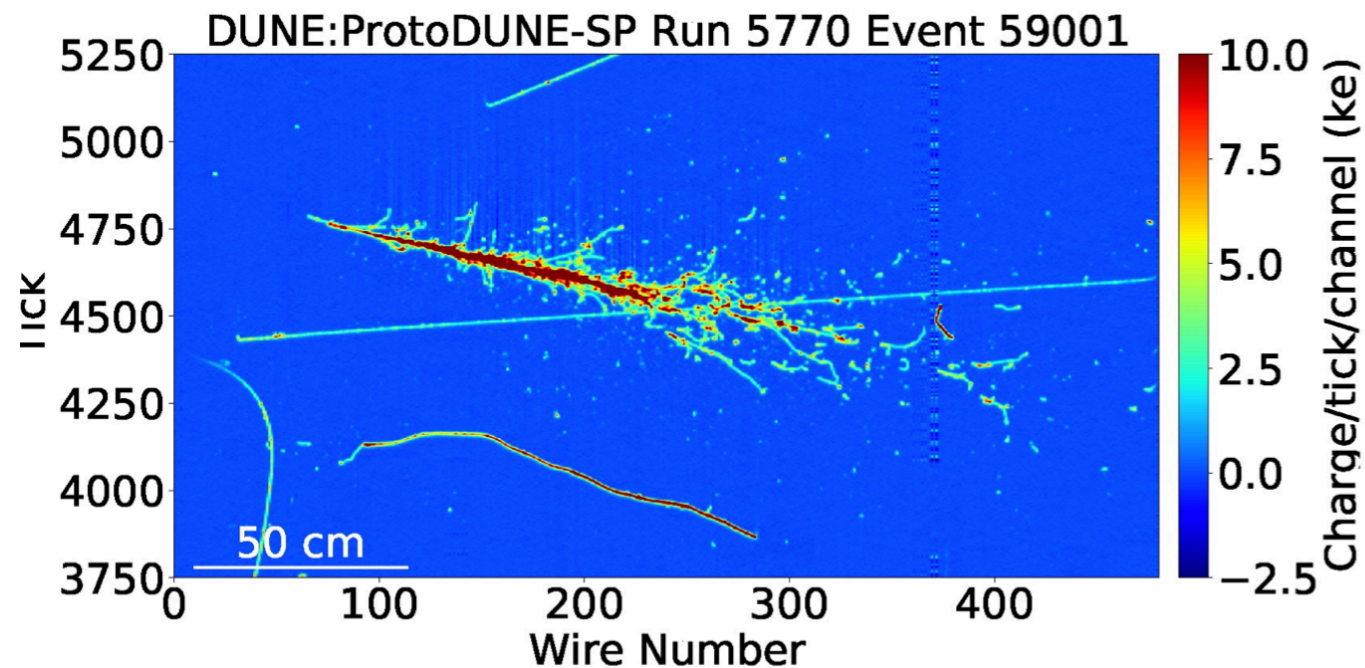➡ **Why?**

➡ **TPC Information is very rich**

    ➡ triggering algorithms are more sophisticated than what a hardware trigger could do

DUNE:ProtoDUNE-SP Run 5770 Event 59001

DUNE:ProtoDUNE-SP Run 5772 Event 15132
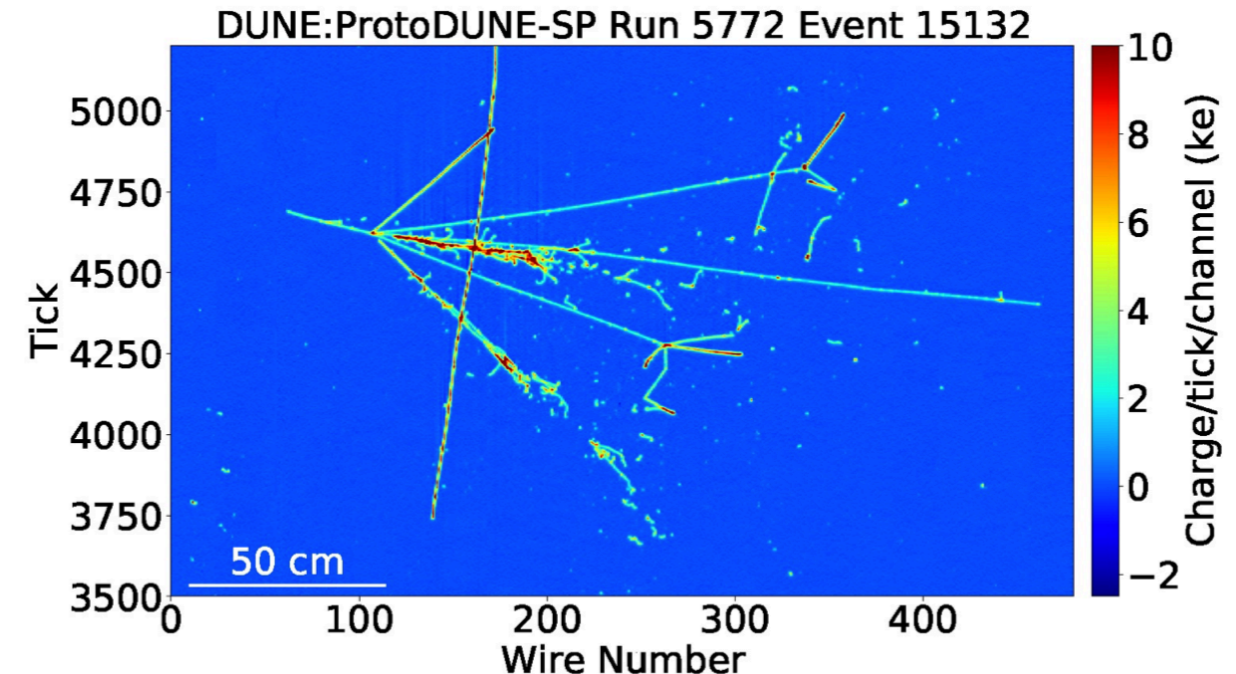
➡ **Why?**

➡ **TPC Information is very rich**

  ➡ triggering algorithms are more sophisticated than what a hardware trigger could do

➡ **TPC is also very slow and u/g rates are very low...**

  ➡ Plenty of time to make decisions, large buffers add more time

  ➡ Not naturally "friendly" to a hardware approach

DUNE:ProtoDUNE-SP Run 5770 Event 59001

DUNE:ProtoDUNE-SP Run 5772 Event 15132
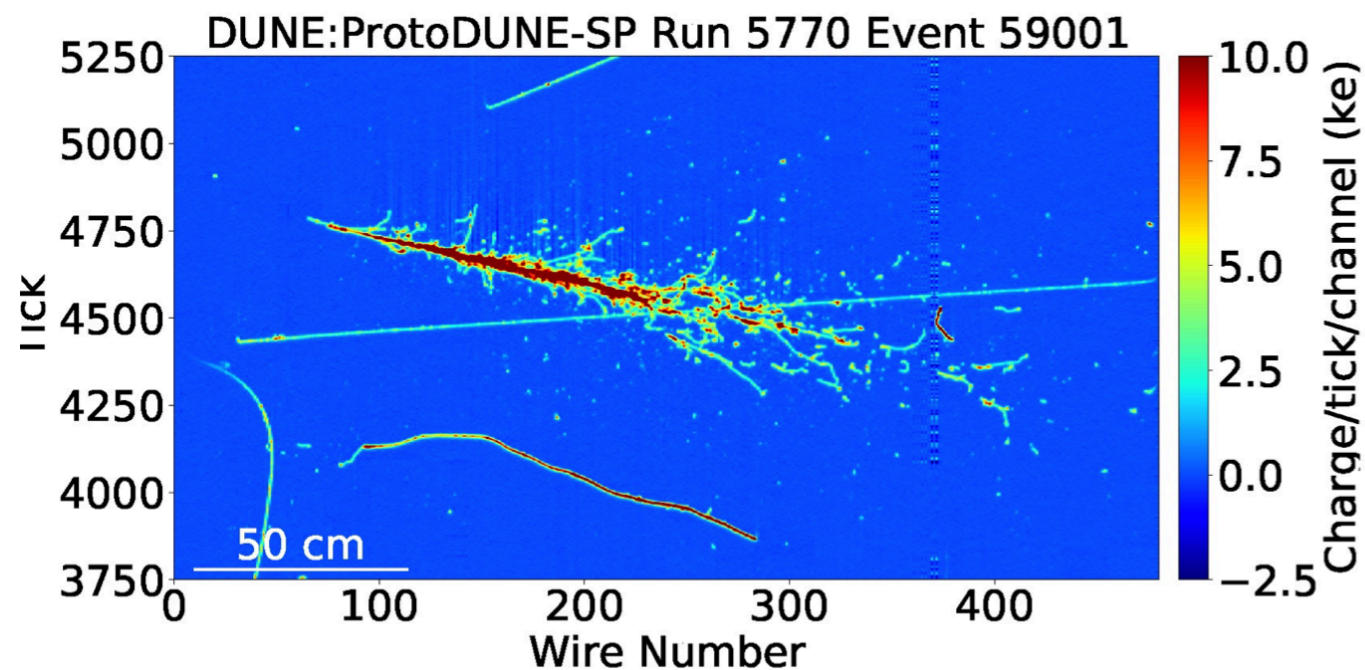
➡ **Why?**

➡ **TPC Information is very rich**

   ➡ triggering algorithms are more sophisticated than what a hardware trigger could do

➡ **TPC is also very slow and u/g rates are very low...**

   ➡ Plenty of time to make decisions, large buffers add more time

   ➡ Not naturally "friendly" to a hardware approach

➡ **Want out-of-beam triggering for broad program**

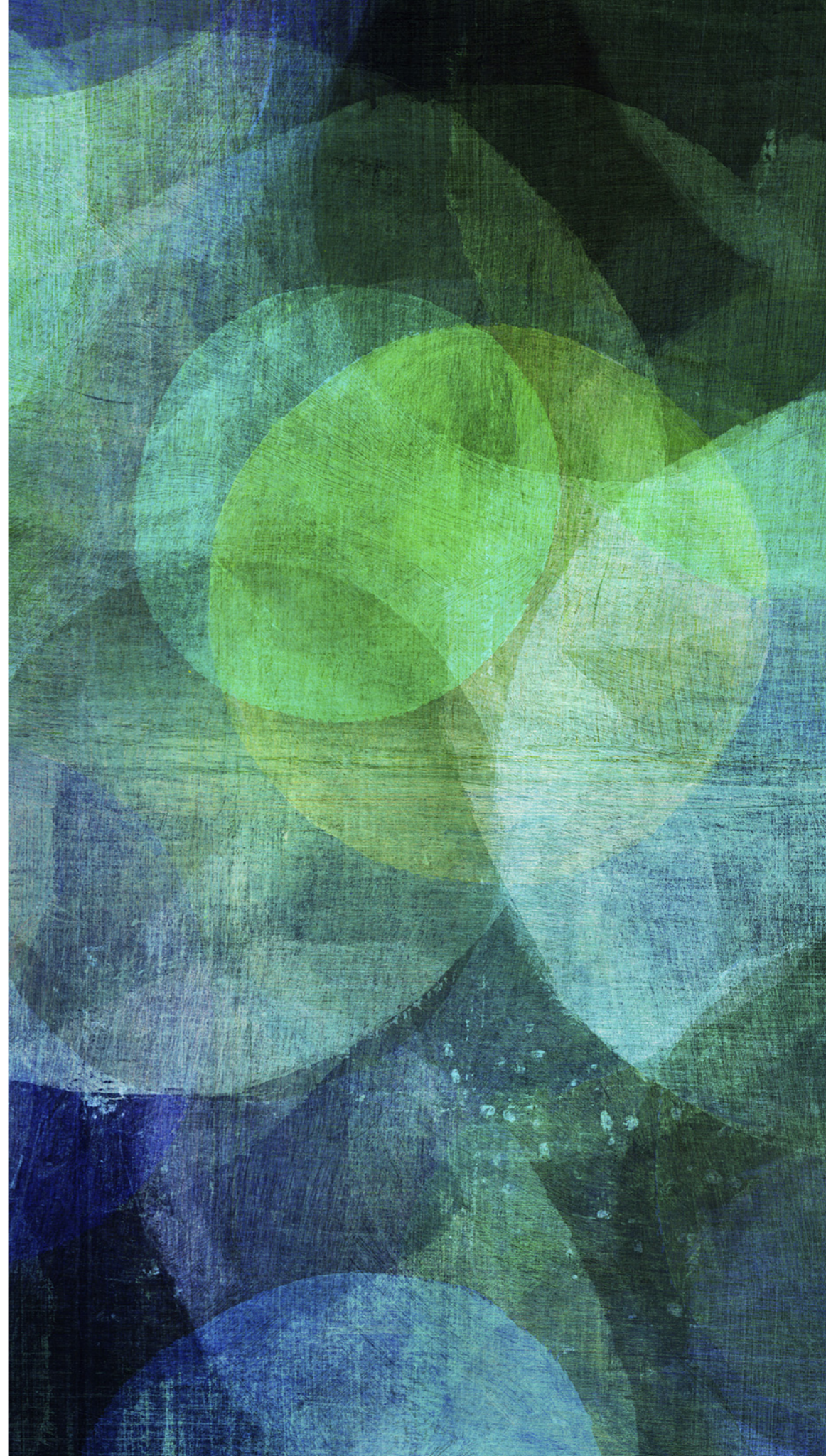   ➡ And beam information may be slow to arrive anyway

# IT'S ALL ABOUT PHYSICS

➡ **The knowledge of hardware and software technologies is becoming critical in our community**

    ➡ thanks to this school we try to keep a high level

➡ **The physics goals depends on technology and innovation**

    ➡ Particle physicists must monitor technological trends and make innovation  (especially true in TDAQ field)

➡ **Not always easy to make extrapolations for the future**

➡ **[Snowmass 2022 report]**

    ➡ *"Modern computing architectures and emerging technologies are changing the way we do particle physics"*

    ➡ *"Machine learning was essentially not a part of the 2013 Snowmass report"*

➡ **[ATLAS TDR, 2003]**

    ➡ *"Thanks to the Moore law, in 2007 our event selection farm will be based on 8 GHz CPUs"*

➡ **[Ken Olsen, Founder of DEC, 1977]**

    ➡ *"There is no reason anyone would want a computer at home."*
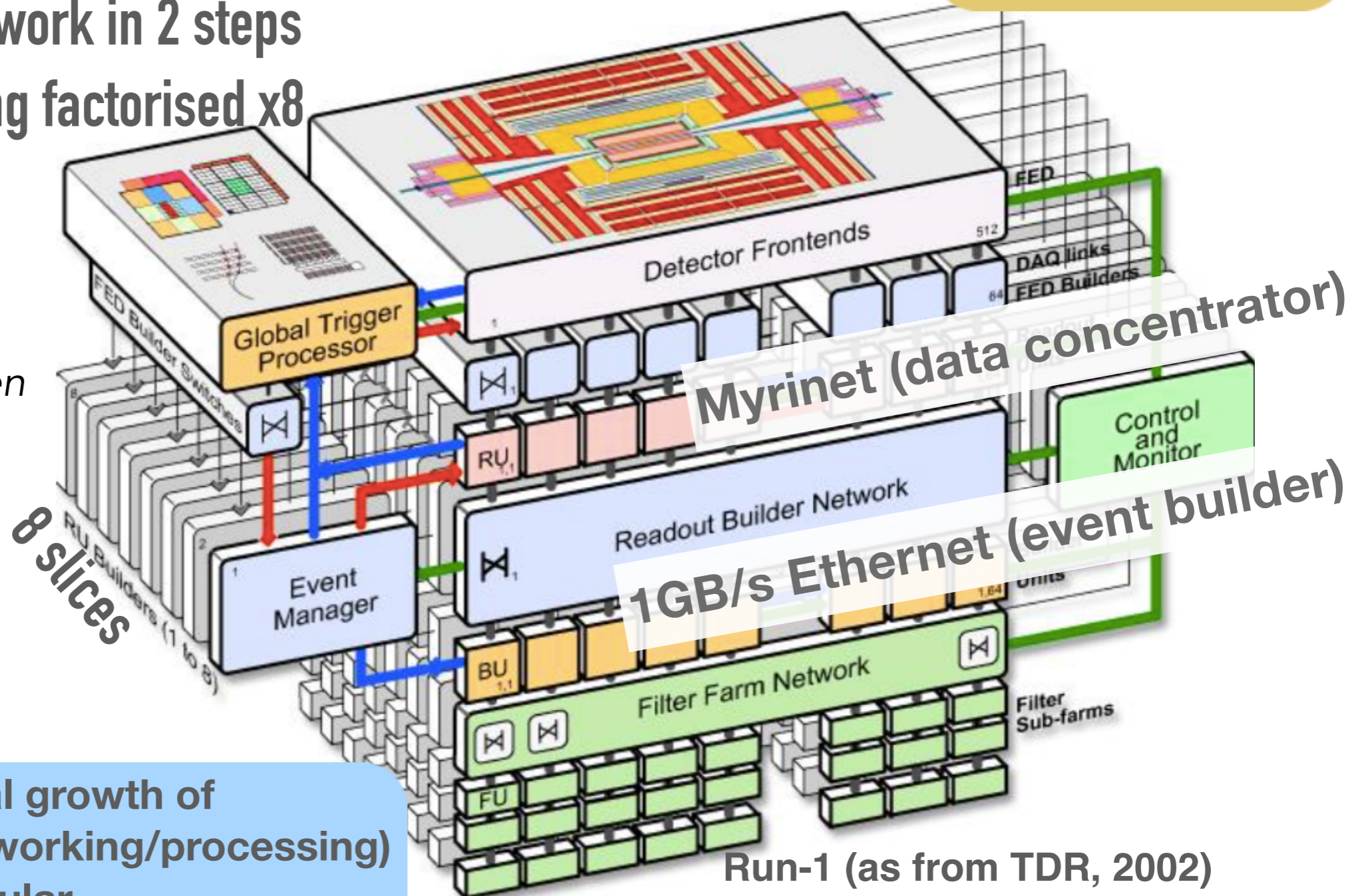
# BACK-UP SLIDES

**Cannot do Event Building at 100 kHz**

**CMS DAQ-1**

100 GB/s readout network in 2 steps

100 kHz Event Building factorised x8

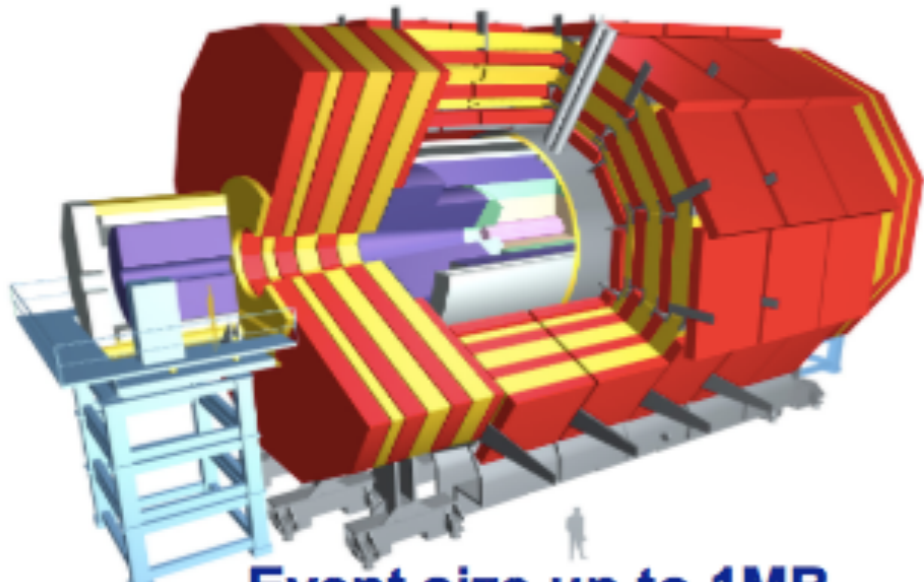*2 EB networks in blu*

*Filter network in green*



Myrinet (data concentrator)

1GB/s Ethernet (event builder)

➡ **Bet on exponential growth of technologies (networking/processing)**
➡ **Scalable and modular**
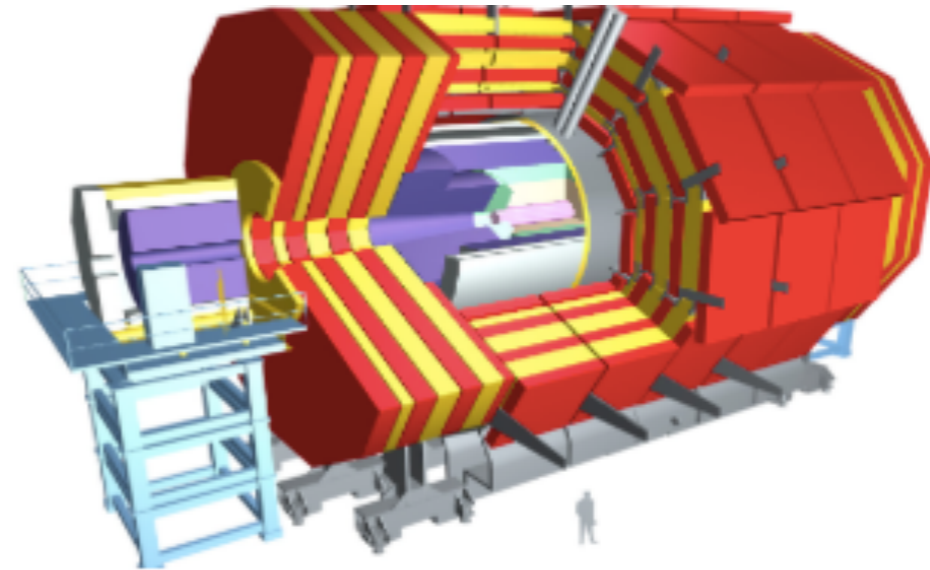  ➡ Independent development of two network technologies

**Run-1 (as from TDR, 2002)**
➡ Myrinet + 1GBEthernet
➡ 1-stage building: 1200 cores (2C)
➡ HLT: ~13,000 cores
➡ 18 TB memory @100kHz: ~90ms/event

**Event size up to 1MB**

**Event size up to 2MB**

**100 kHz L1 rate**

**100 kHz L1 rate**

**Myrinet**

**1 Gb/s Ethernet**

**10/40 Gb/s Ethernet**

**56 Gb/s Infiniband**

**100 GB/s 8 slices**

**~200 GB/s**

**1 slice**

**CMS DAQ 1**

**CMS DAQ 2**

13000 core, 1260 host filter farm

16000+ core, 900 host filter farm

**max. 1.2 GB/s to storage**

**~ 3-6 GB/s to storage**

**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**

# 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

**LHCb 2015 Trigger Diagram**

**40 MHz bunch crossing rate**

L0 Hardware Trigger
readout, high Et

**450 kHz**
**h±**

**μ/μμ**

**150 kHz**
**e/γ**

**NO L0 trigger**

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**

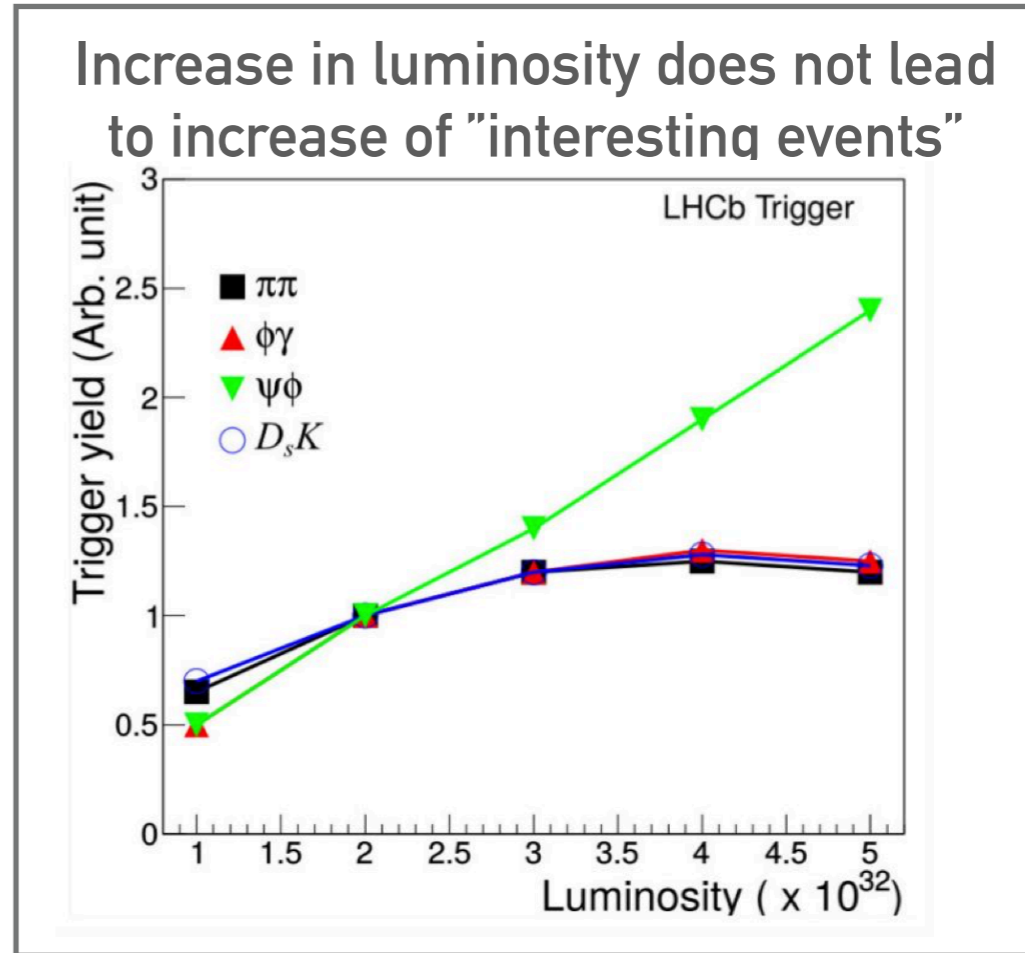**NO offline analysis**

**12.5 kHz Rate**

*See Phase-I upgrade TDR*

**Can increase luminosity x10 ?**
**Can increase b-hadron efficiency x2?**

**YES, remove limit from L0 –1MHz readout!**

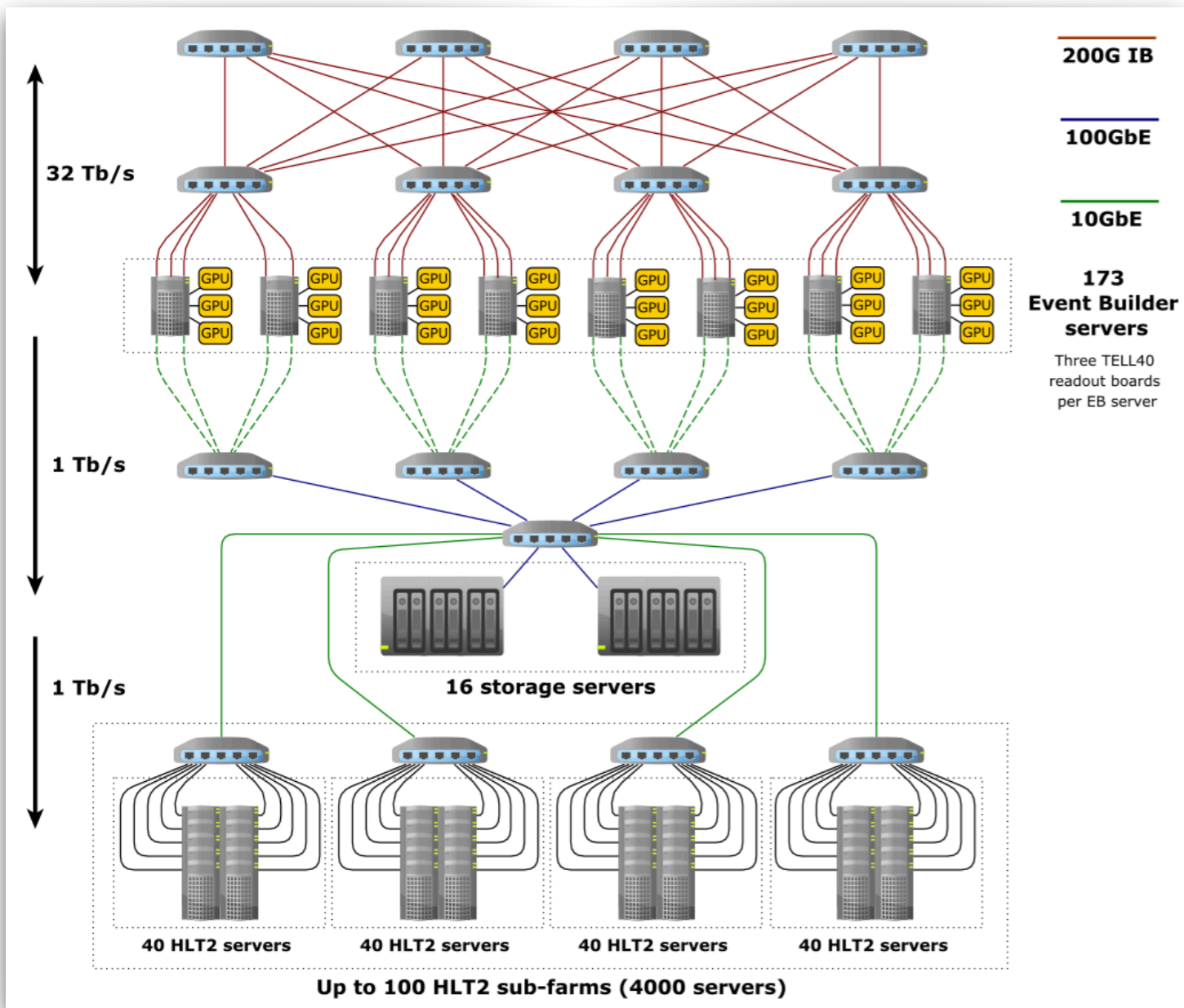Increase in luminosity does not lead to increase of "interesting events"



LHCb Trigger

- ■ $\pi\pi$
- ▲ $\phi\gamma$
- ▼ $\psi\phi$
- ○ $D_sK$

Trigger yield (Arb. unit)

Luminosity ( $\times 10^{32}$ )

**Allow detector readout and reconstruction at unprecedented rate: 30MHz !!**
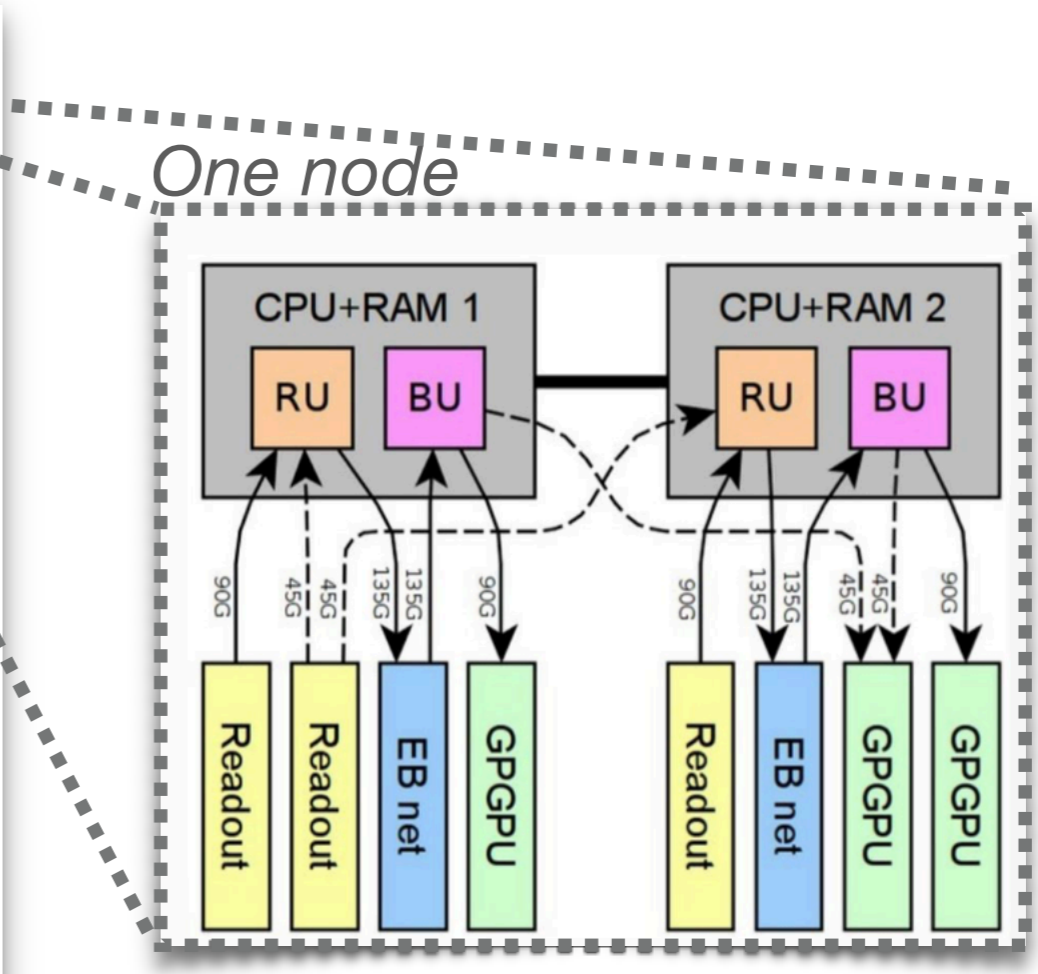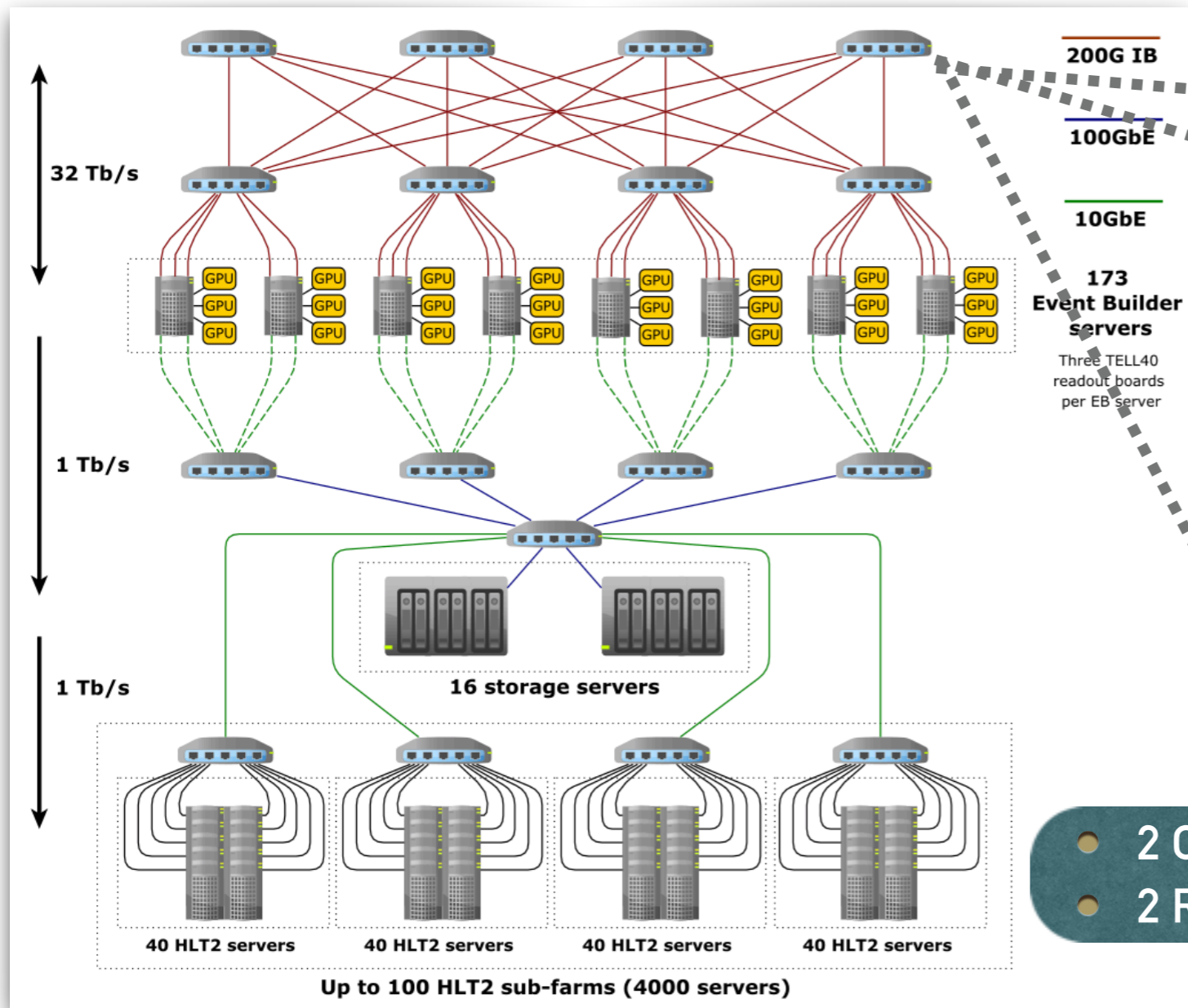
**Large farm of equal nodes with 8 PCIe40 boards, specialised by firmware**



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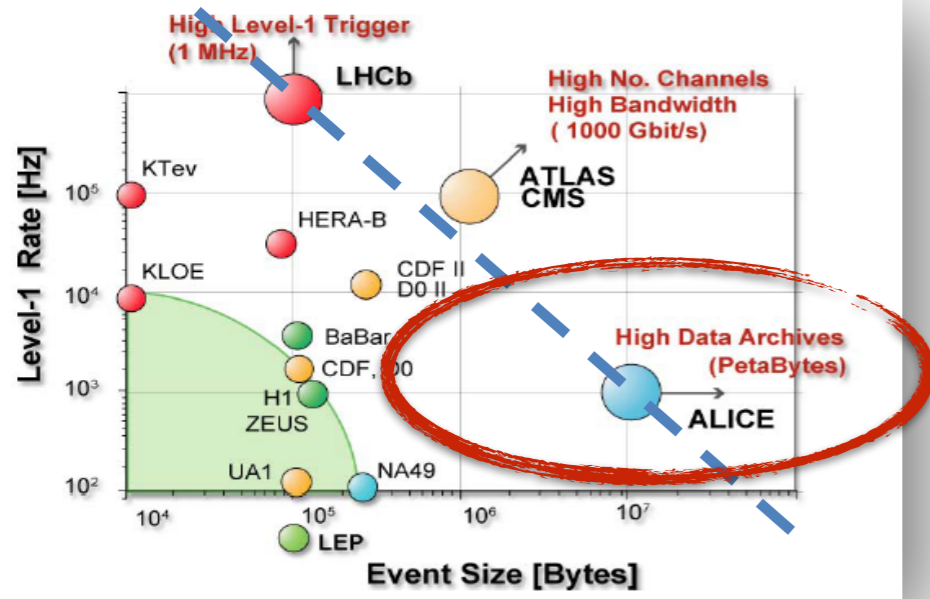
**Large farm of equal nodes with 8 PCIe40 boards, specialised by firmware**



*One node*

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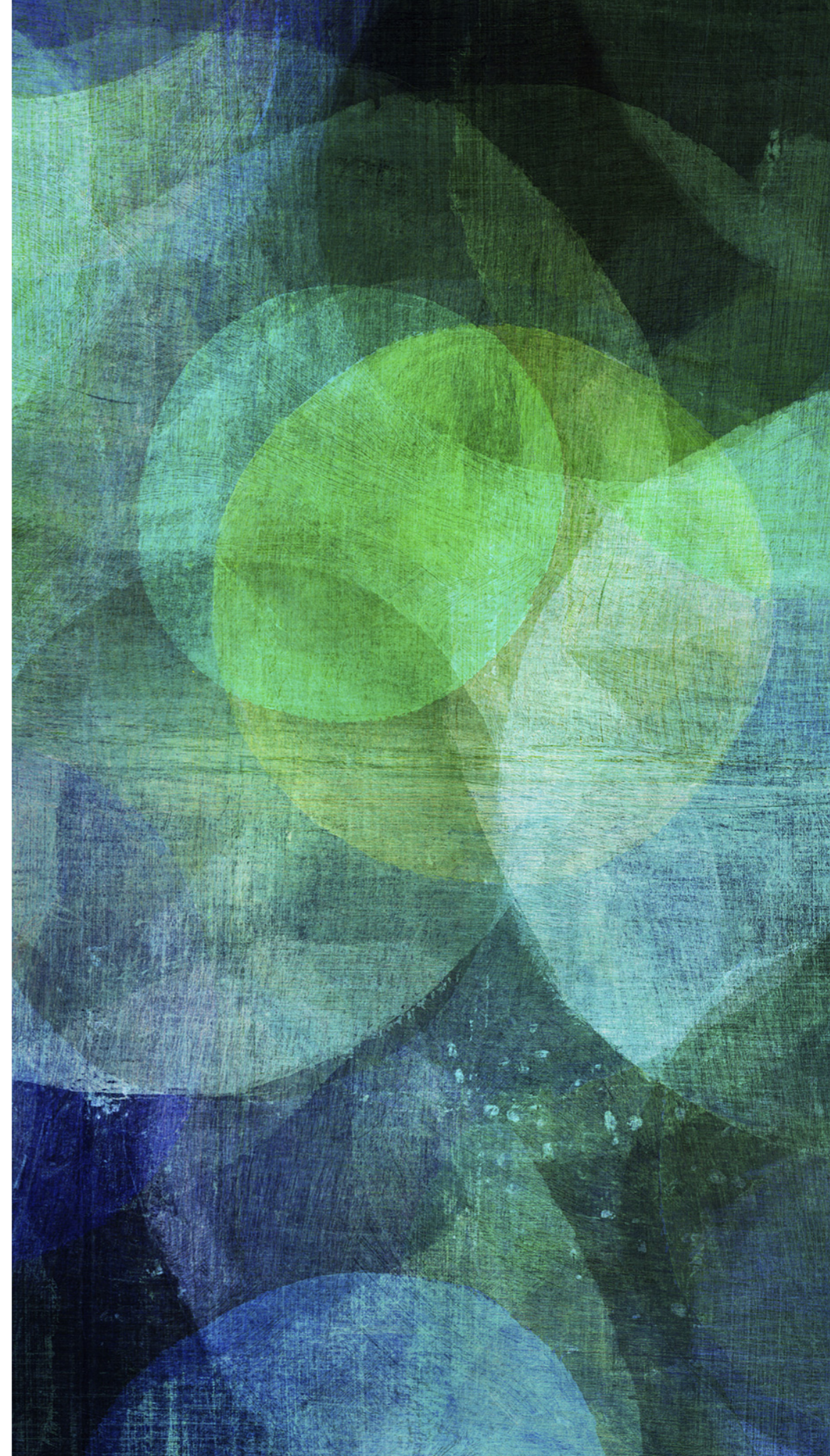- 2 RU, 2 BU, 2 infiniband NIC (200 Gb/s), 1-3 GPUs

➡ **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**

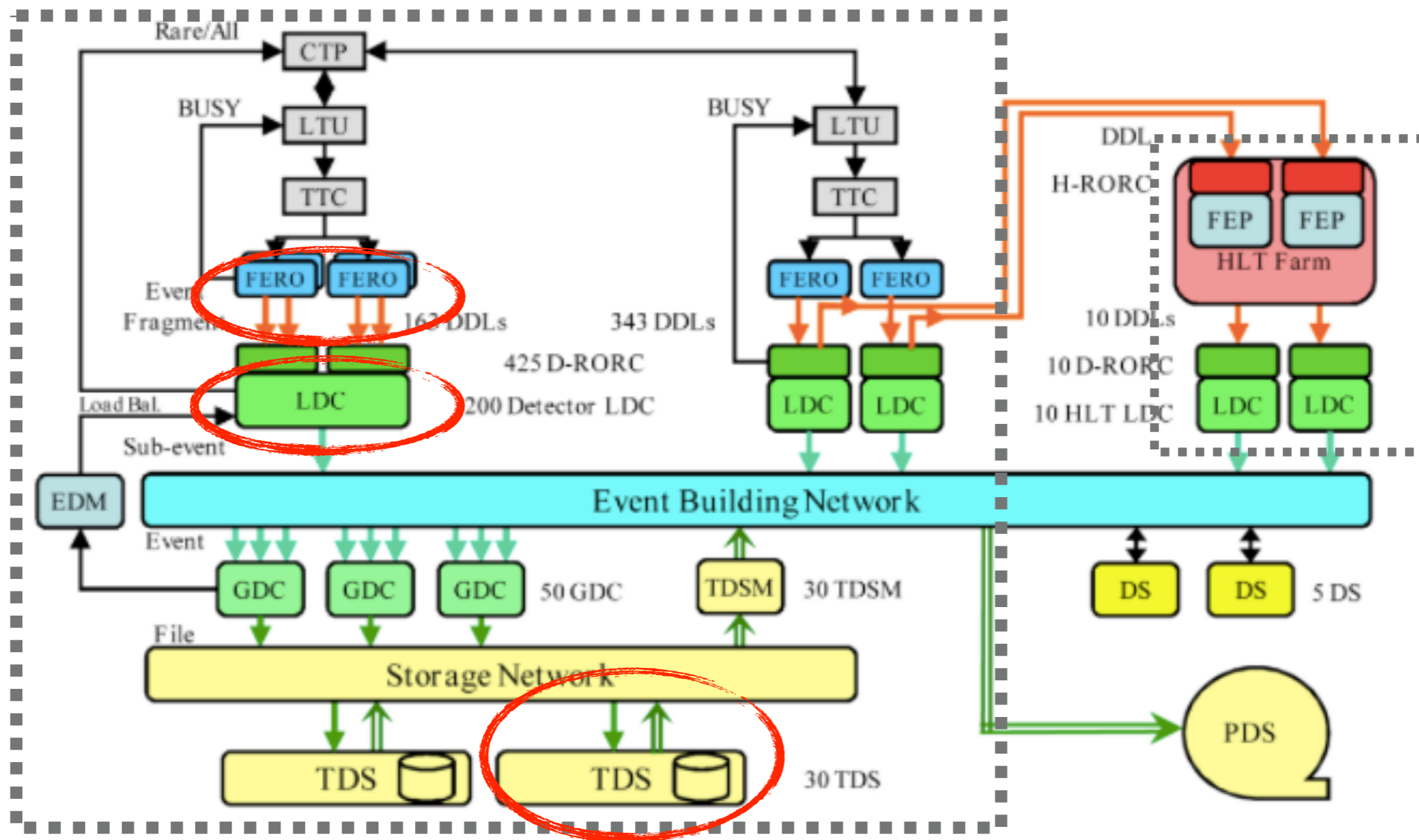  ➡ SSD faster but have short lifetime wrt high read-write rate, so prefer hard-disks

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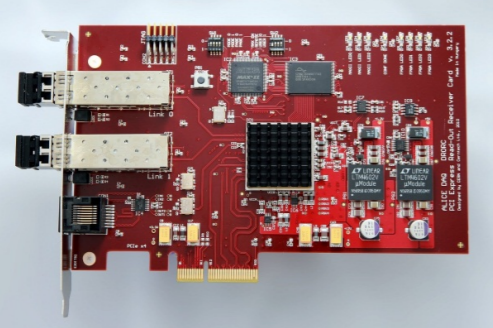
*Recording heavy ion collisions*

*http://alice-daq.web.cern.ch*
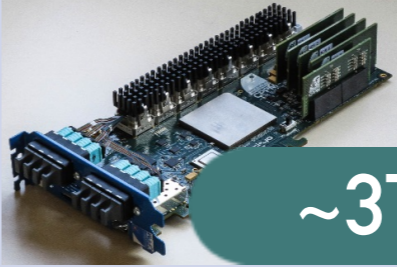
➡ **Dataflow with local (LDC) and global (GDC) data concentrators**

   ➡ 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!**

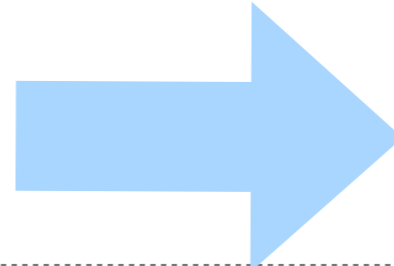| RORC 1 | C-RORC | CRU |
|---|---|---|
| 2 ch @ 2 Gb/s<br>PCIe gen.1 x4 (1 GB/s) | 12 ch @ up to 6 Gb/s<br>PCIe gen.2 x 8 (4 GB/s) | 24 ch @ 5 Gb/s<br>PCIe gen.3 X 16 (16 GB/s) |
| Custom DDL protocol | Custom DDL protocol<br>(same protocol but faster) | GBT |
| Protocol handling<br>TPC Cluster Finder | 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
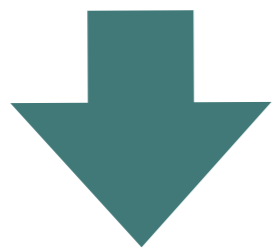
**Higher rates with smaller data?**

**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**
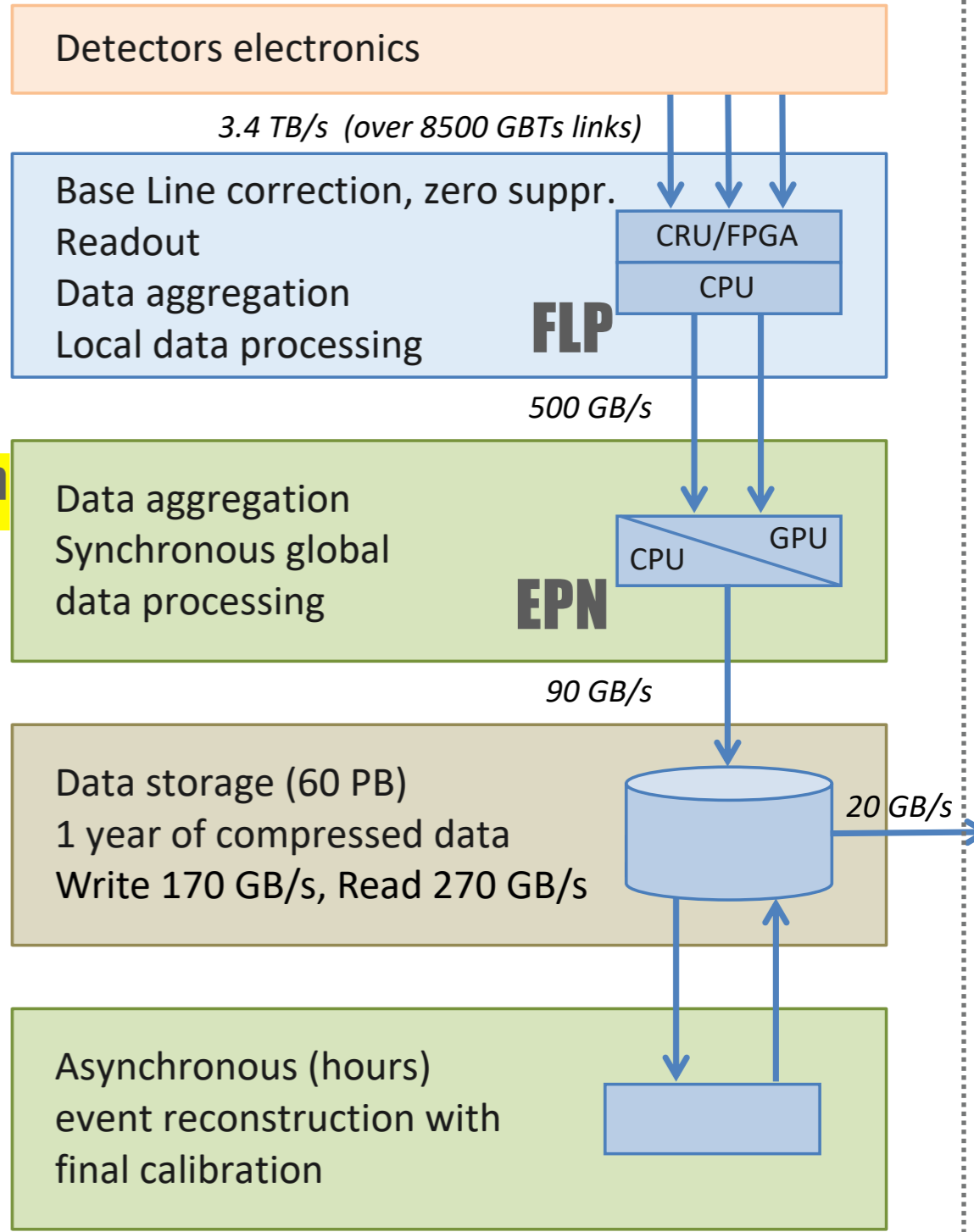  - ➡ 250 EPN servers with 8 GPU-cards
  - ➡ Require large-memory GPUs!

## 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
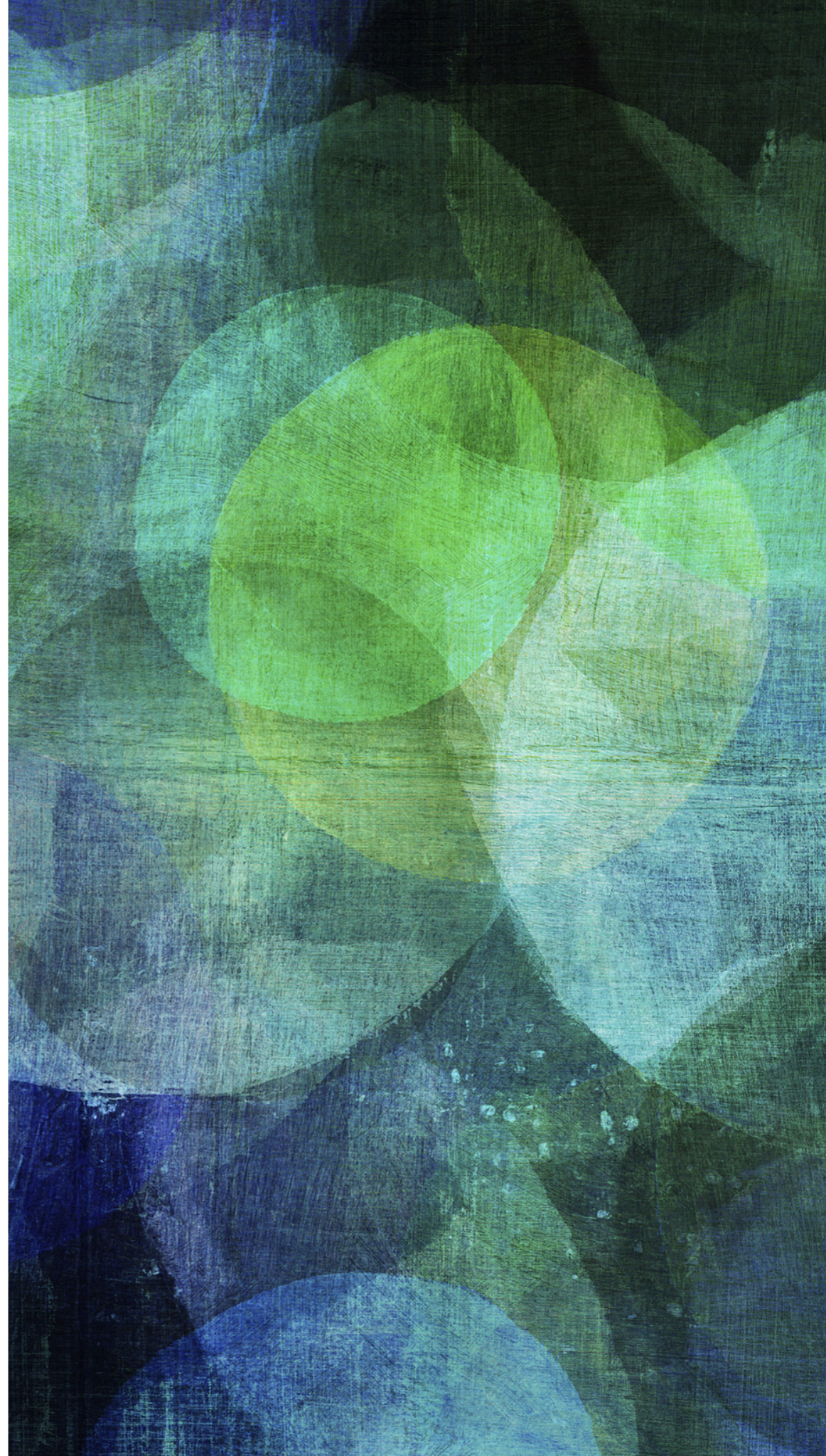
*20 GB/s*

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

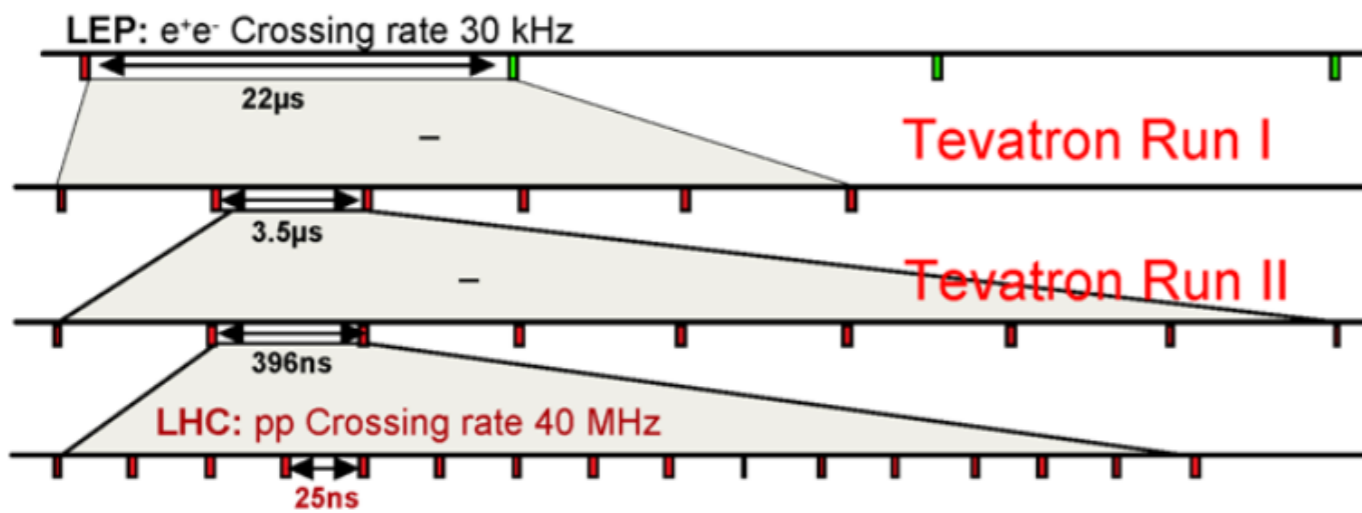  ➡ sometimes stealing ideas ("… but we can do better than that…")

# BACK-UP SLIDES

## The clock source

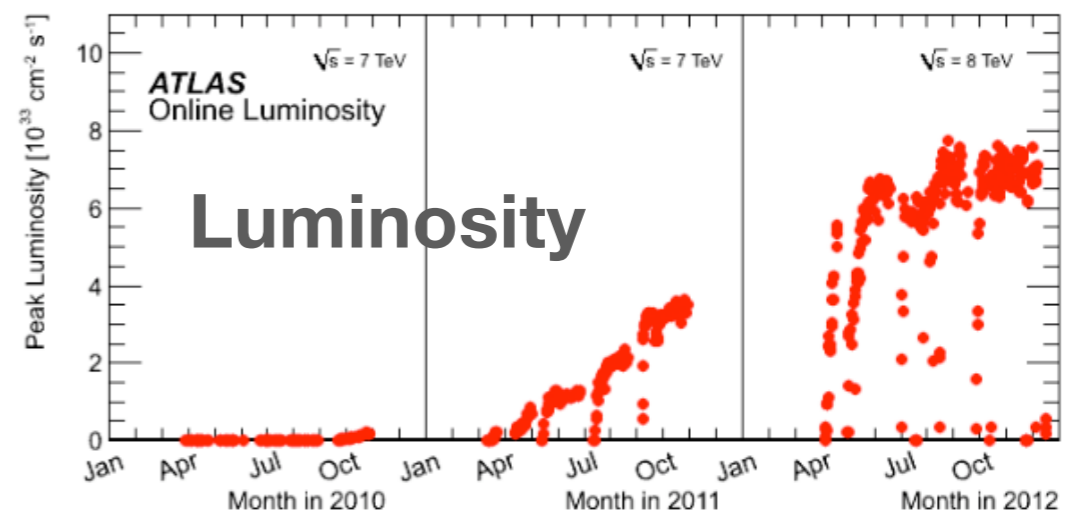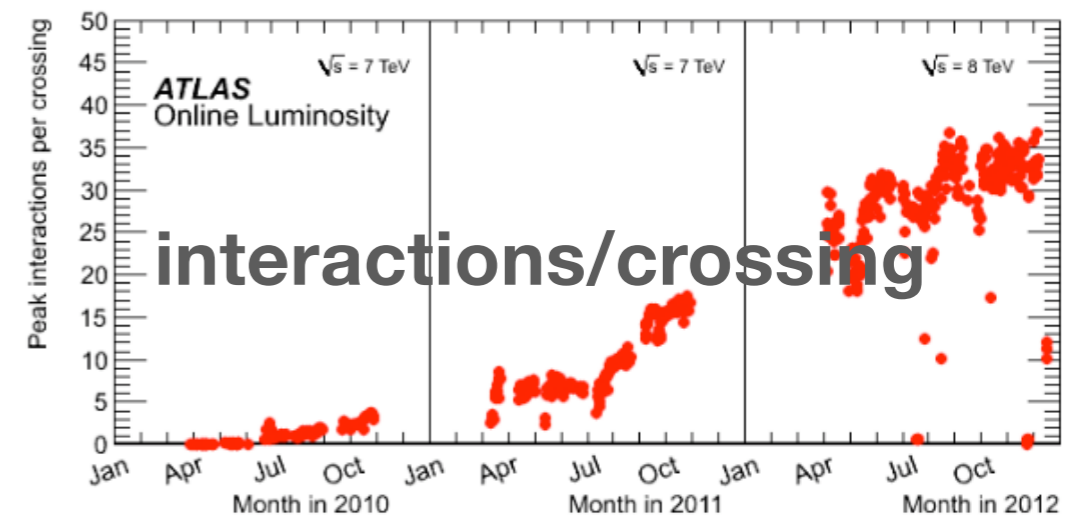- ➡ ~3600 bunches in 27km
- ➡ distance bw bunches: 27km/3600 = 7.5m
- ➡ distance bw bunches in time: 7.5m/c = 25ns



## The pile-up source

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



**interactions/crossing**

**Luminosity**

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

➡ **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

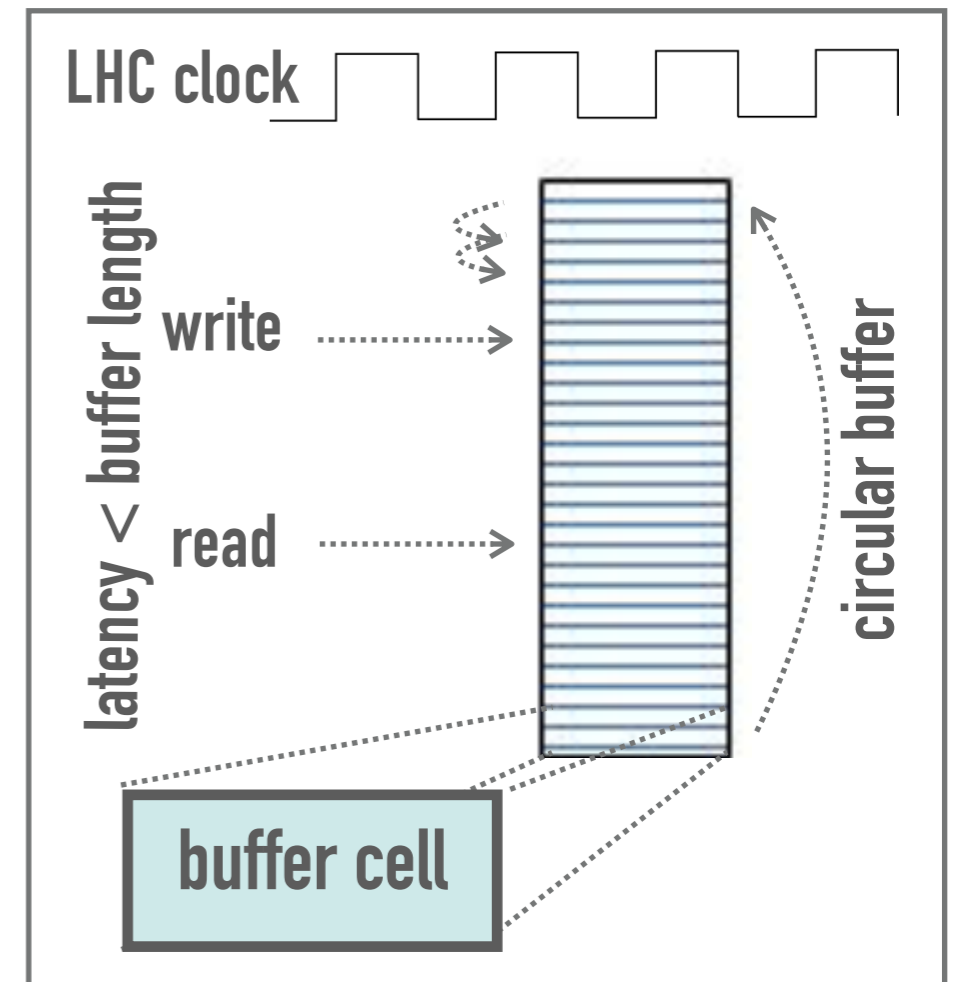➡ **Tight design constraints for trigger/FE**

➡ **Analog/digital pipelines**

- ➡ Analog: built from switching capacitors
- ➡ Digital: registers/FIFO/…

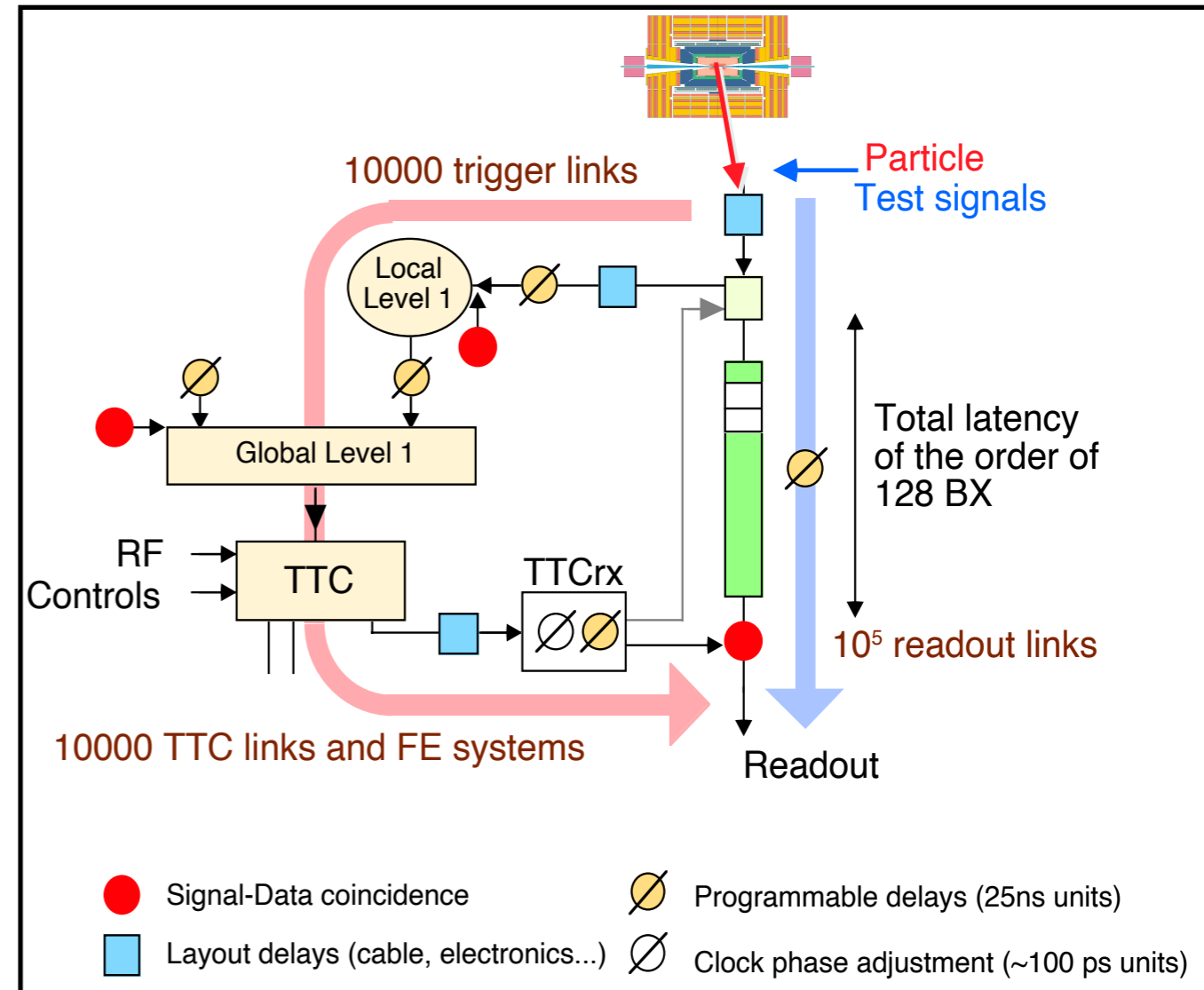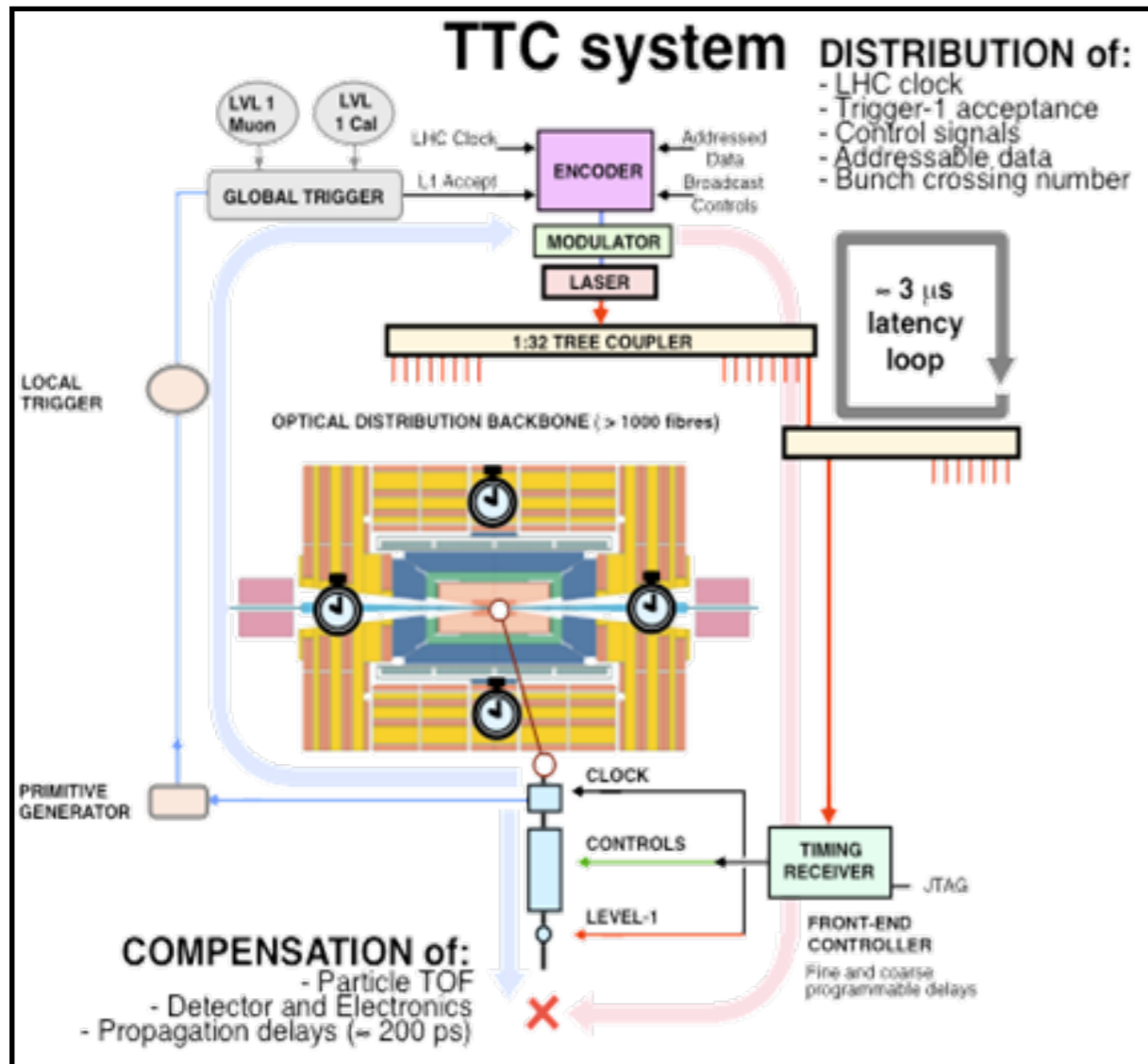➡ **Full digitisation before/after L1A**

- ➡ Fast DC converters (power consumption!)

➡ **Additional complication: synchronisation**

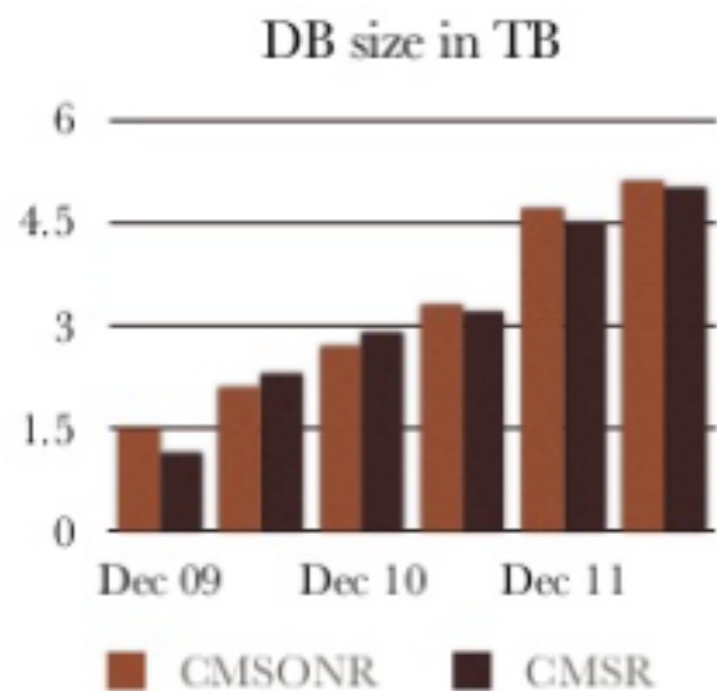- ➡ BC counted and reset at each LHC turn
- ➡ large optical time distribution system

➡ **Common optical system: TTC**
  ➡ radiation resistance
  ➡ single high power laser
➡ **Large distribution**
  ➡ experiments with ~$10^7$ channels
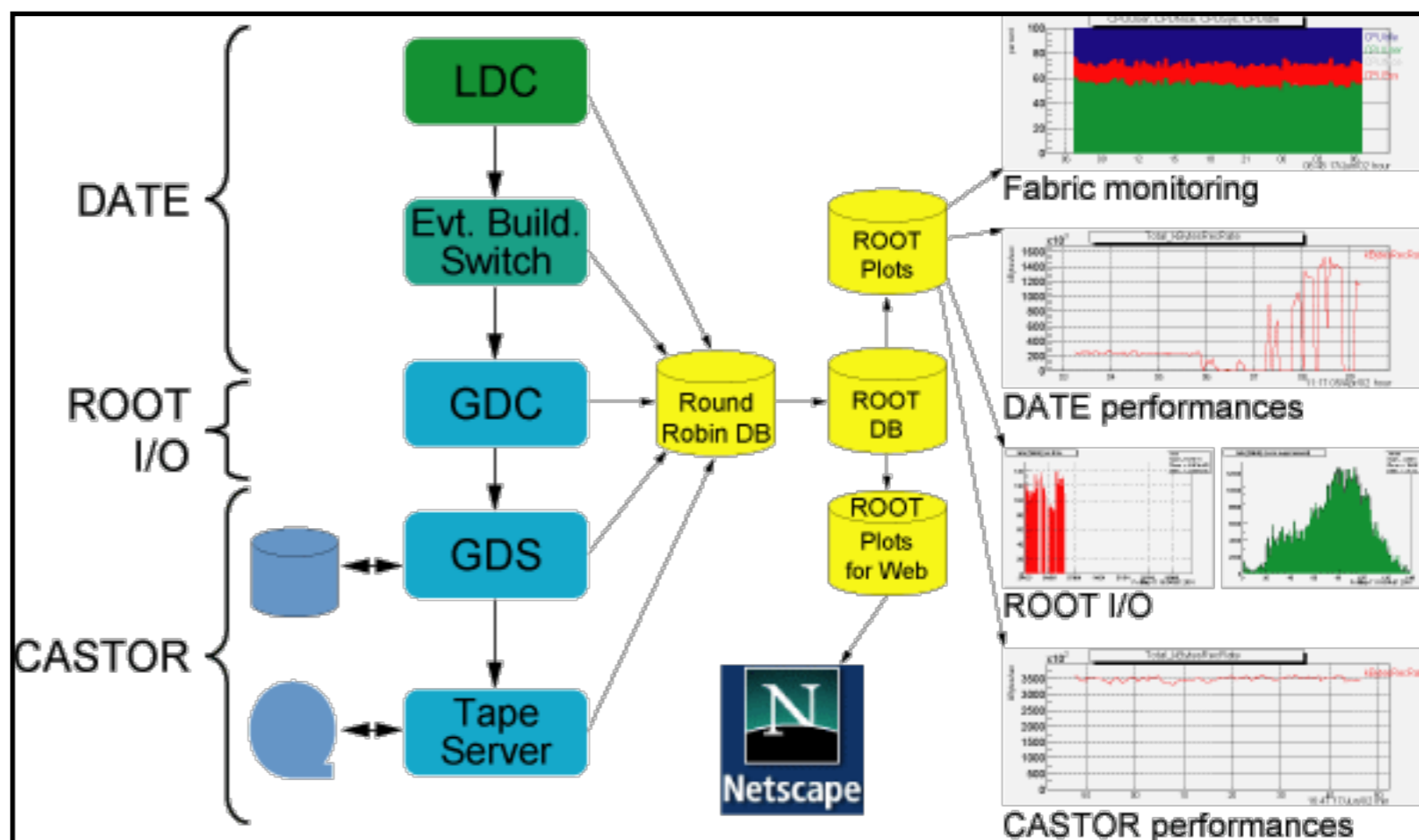
➡ **Align readout & trigger at (better than) 25ns and correct for**
  ➡ time of flight (25 ns ≈ 7.5m)
  ➡ cable delays (10cm/ns)
  ➡ processing delays (~100 BCs)

➡ **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

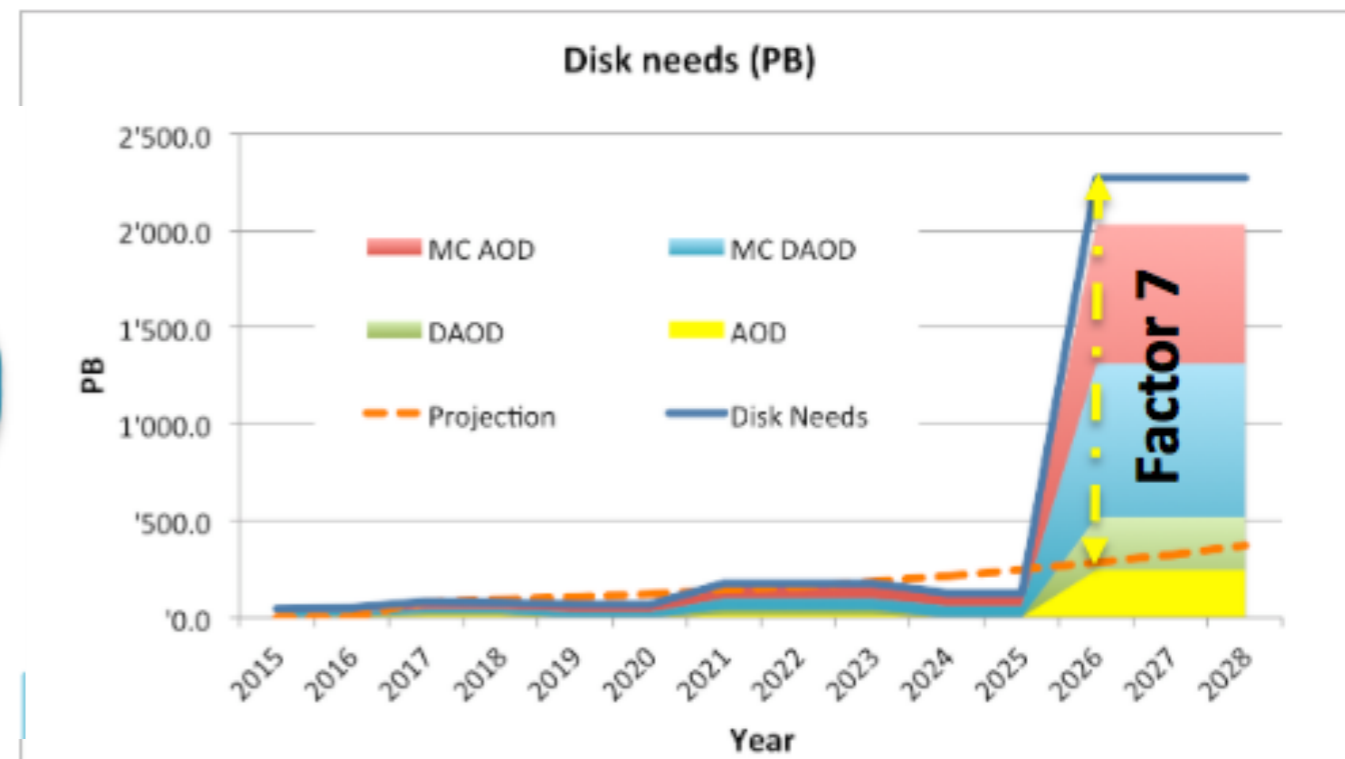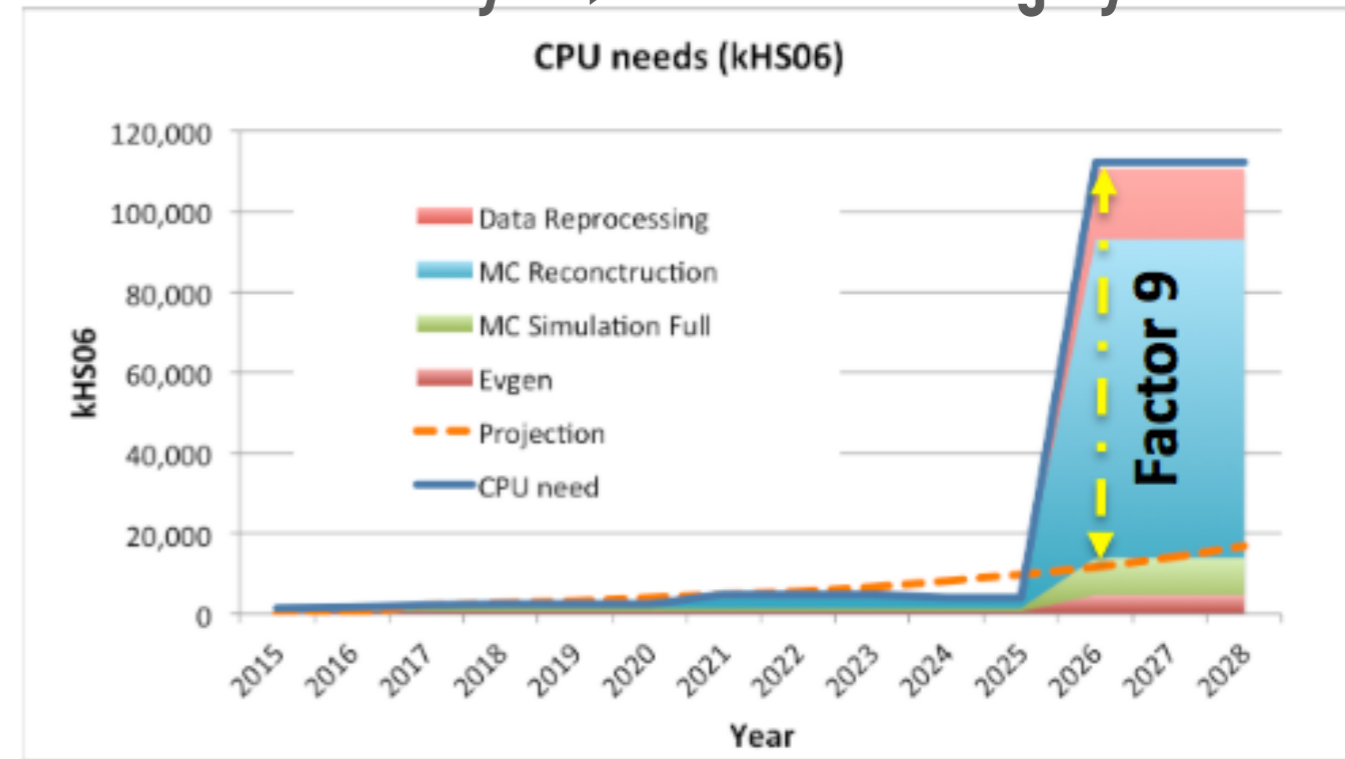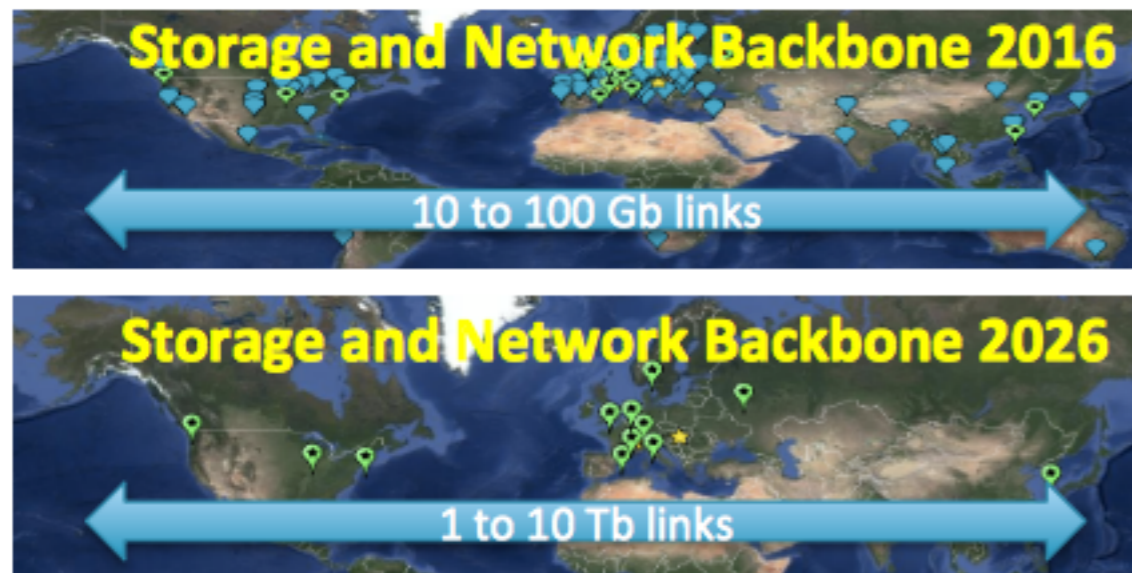  ➡ support software updates and rolling replacement of hardware



**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
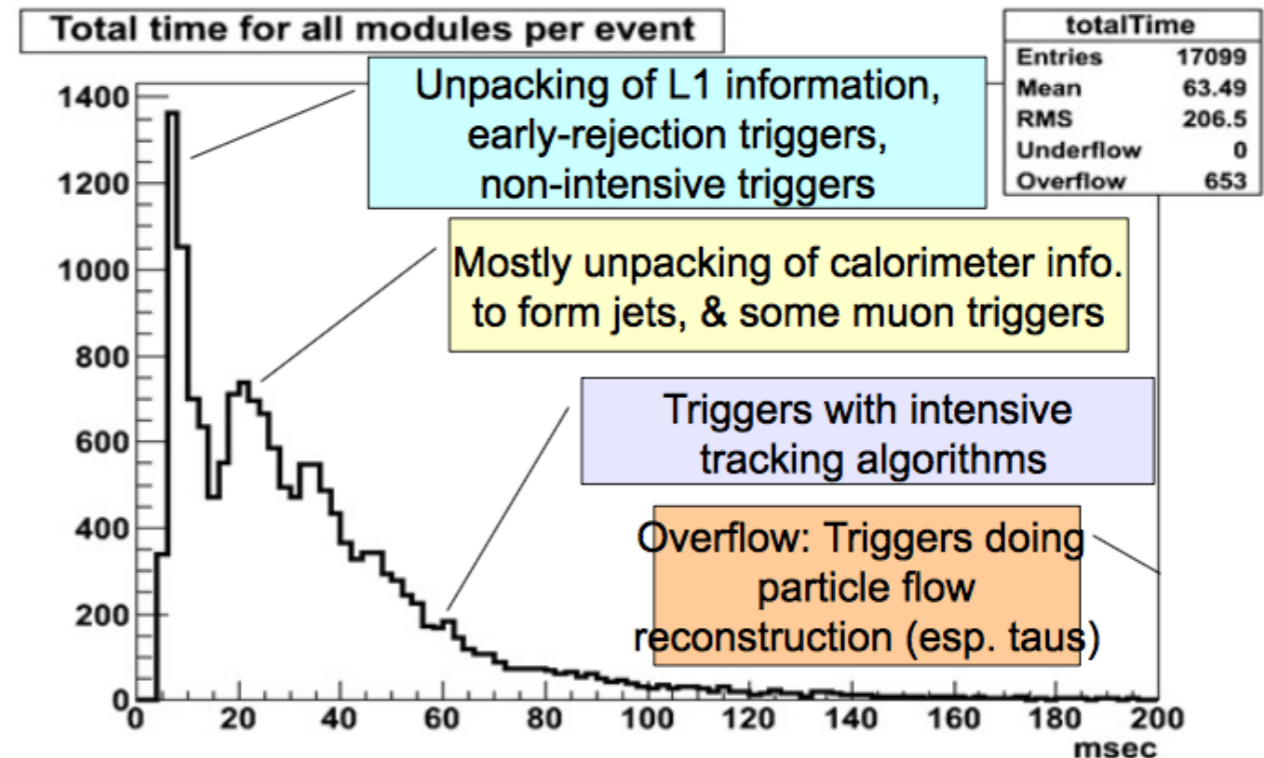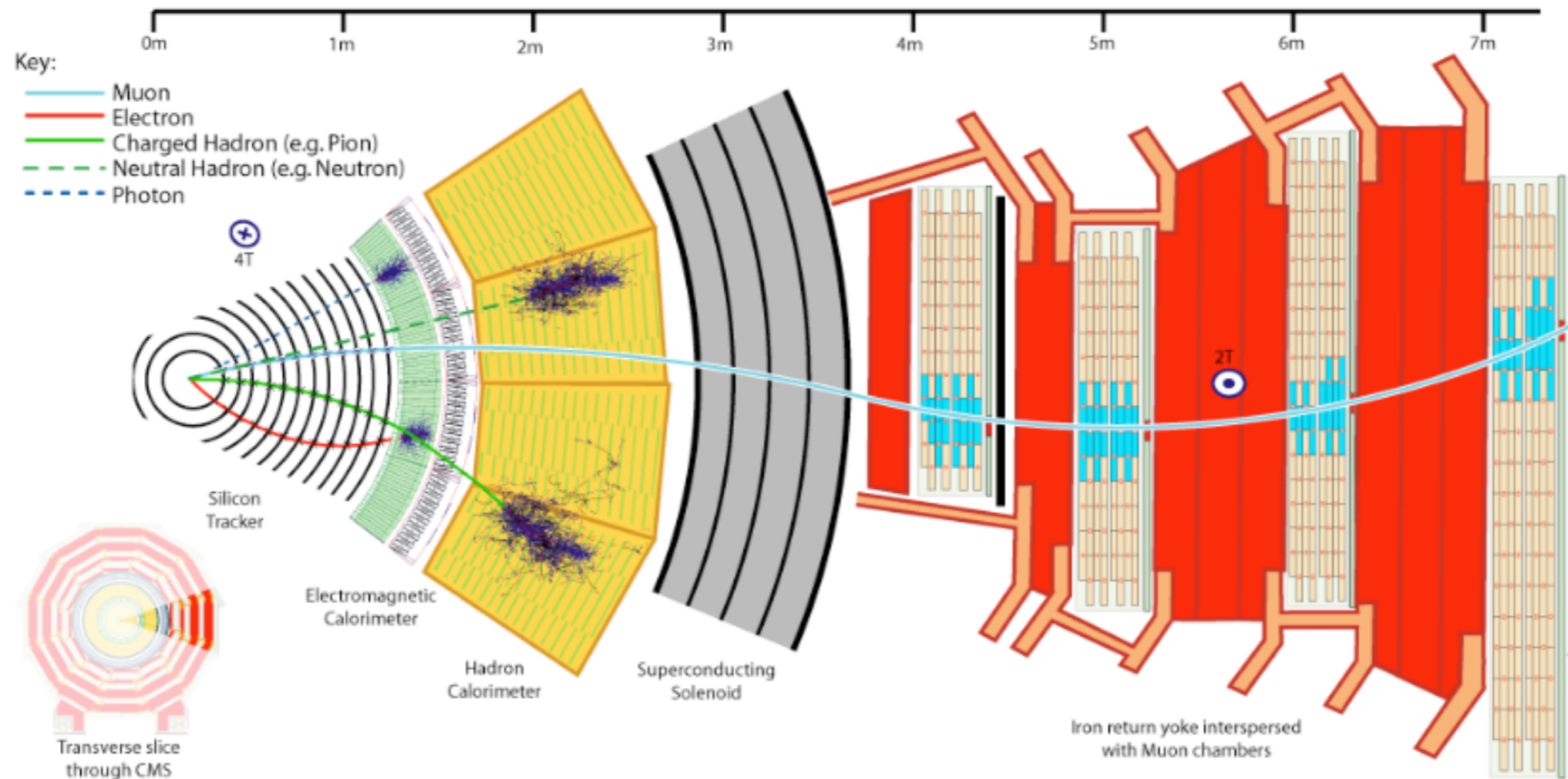— Electron
— Charged Hadron (e.g. Pion)
--- Neutral Hadron (e.g. Neutron)
····· Photon

4T

Silicon Tracker

Electromagnetic Calorimeter

Hadron Calorimeter

Superconducting Solenoi⁻

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 |

Transverse slice through CMS

Silicon Tracker

Electromagnetic Calorimeter

Hadron Calorimeter

Superconducting Solenoid

Iron return yoke interspersed with Muon chambers

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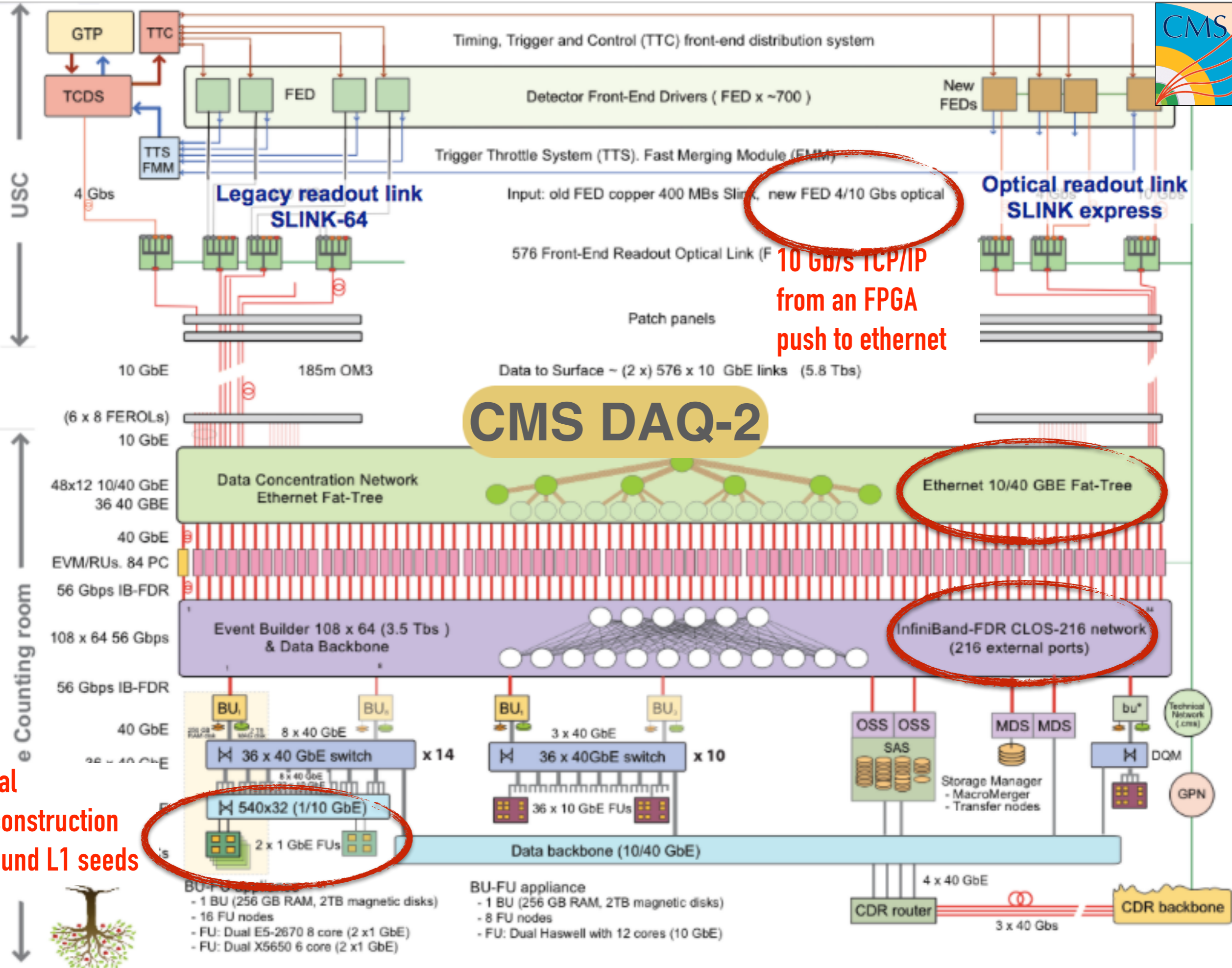- Neutral Hadron (e.g. Neutron)
- Photon

➡ **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

CMS DAQ-2

## 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)**

256 GB RAM disk    2 TB MAG disk

data, status, configuration, latency

**Filter Unit (FU)**

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

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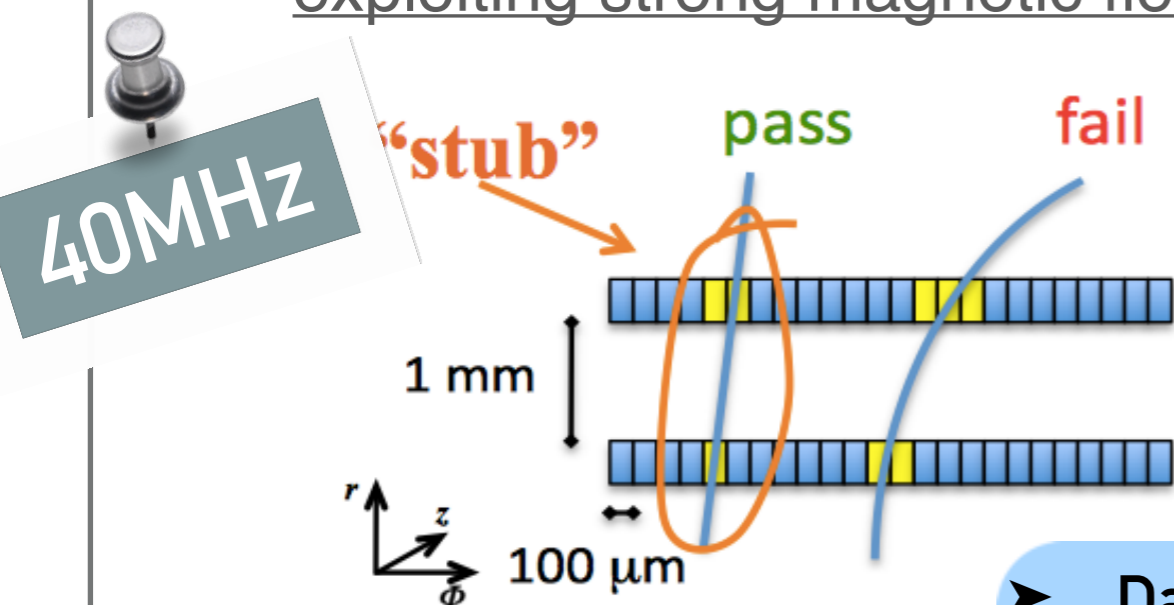➡ Added ability to run WLCG grid jobs in FUs during stops/interfill



HLT contribution

## File-based communication
➡ HLT and DAQ completely decoupled
➡ Network filesystem used as transport (and resource arbitration) protocol (LUSTRE FS)
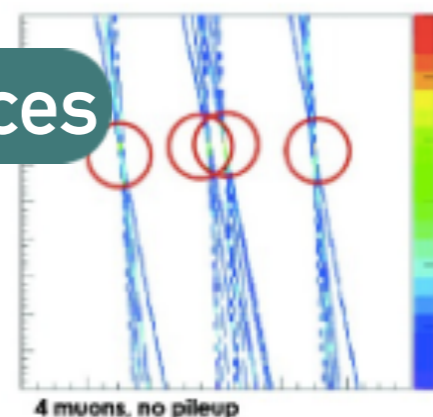
## Track filtering (low p$_T$)

## 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**
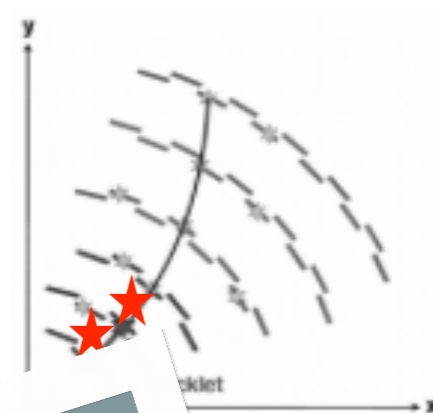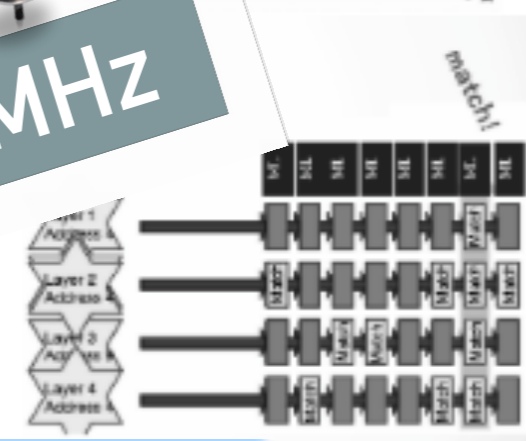
➡exploiting strong magnetic field of CMS



"stub"   pass   fail

1 mm

100 μm

40MHz

1MHz

**Hough Transform**

4 muons, no pileup

**Tracklets**

**Associative Memories**

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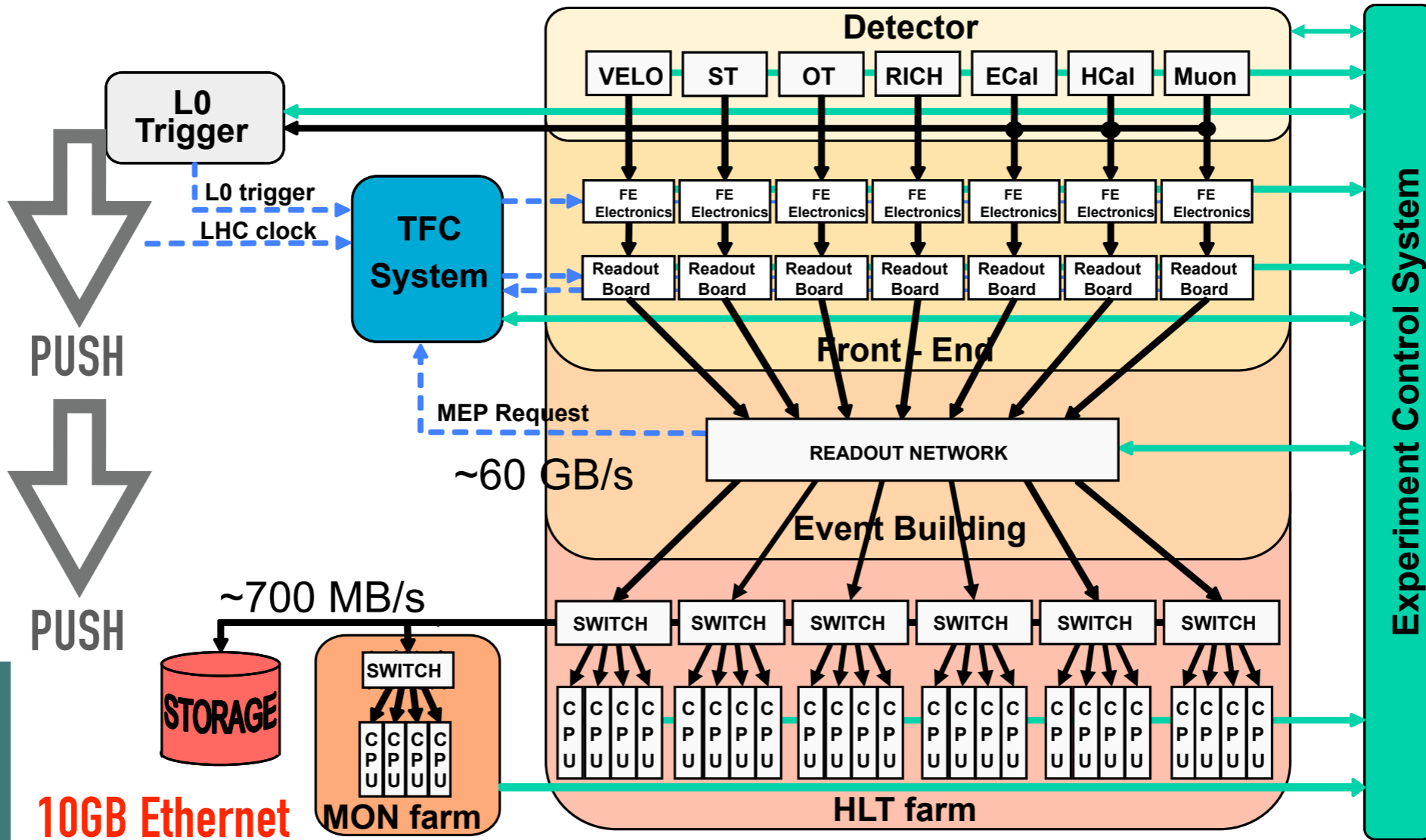➤ Three R&D efforts: FPGA/ASIC

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**

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

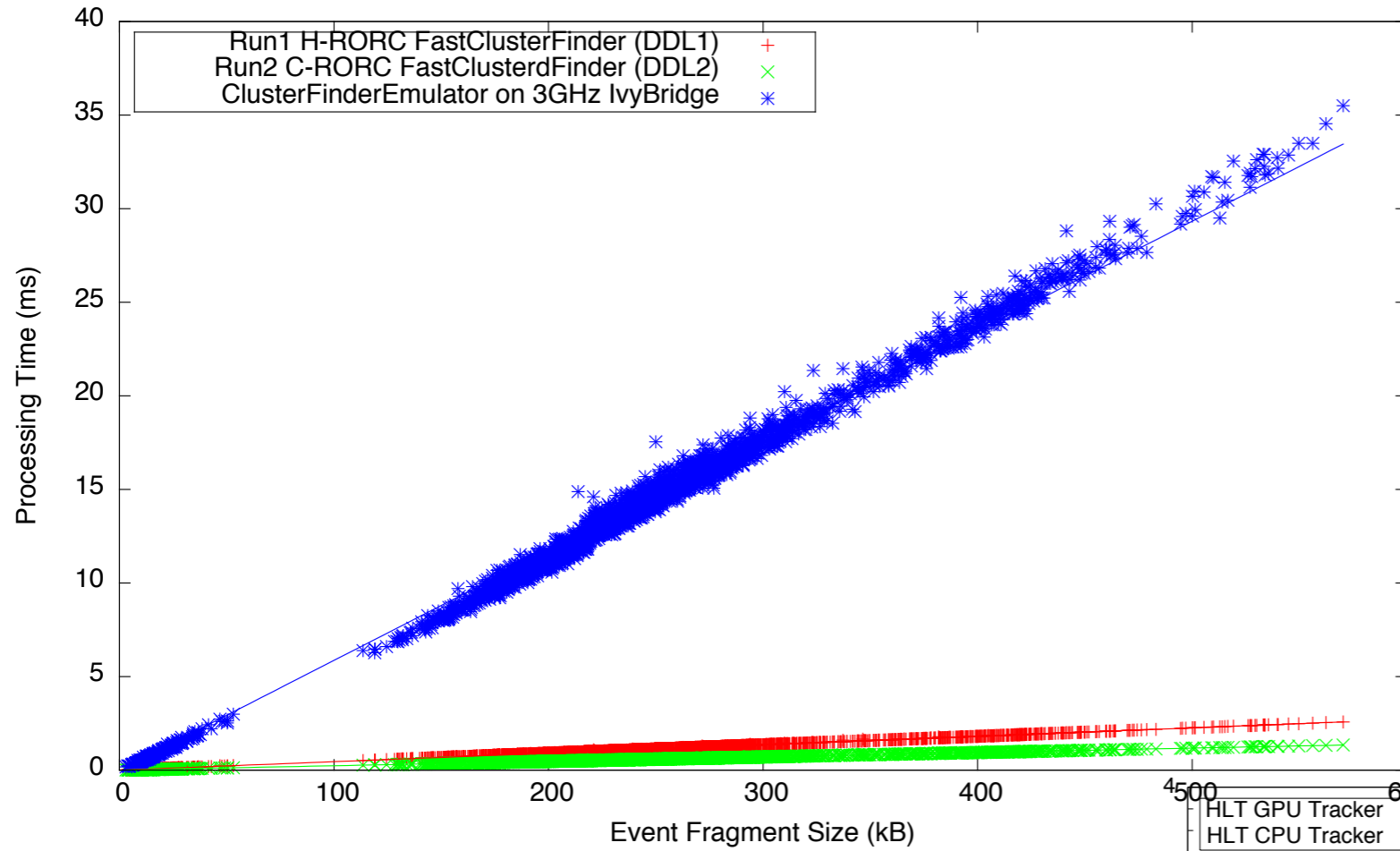**62 sub-farms, total 1780 nodes, with edge-routers (12 Gbps)**

PUSH

PUSH

~700 MB/s

**10GB Ethernet**

Detector

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

**L0 Trigger**

L0 trigger
LHC clock

**TFC System**

FE Electronics

Readout Board

MEP Request

**Front - End**

READOUT NETWORK

~60 GB/s

**Event Building**

SWITCH

STORAGE

SWITCH

MON farm

CPU

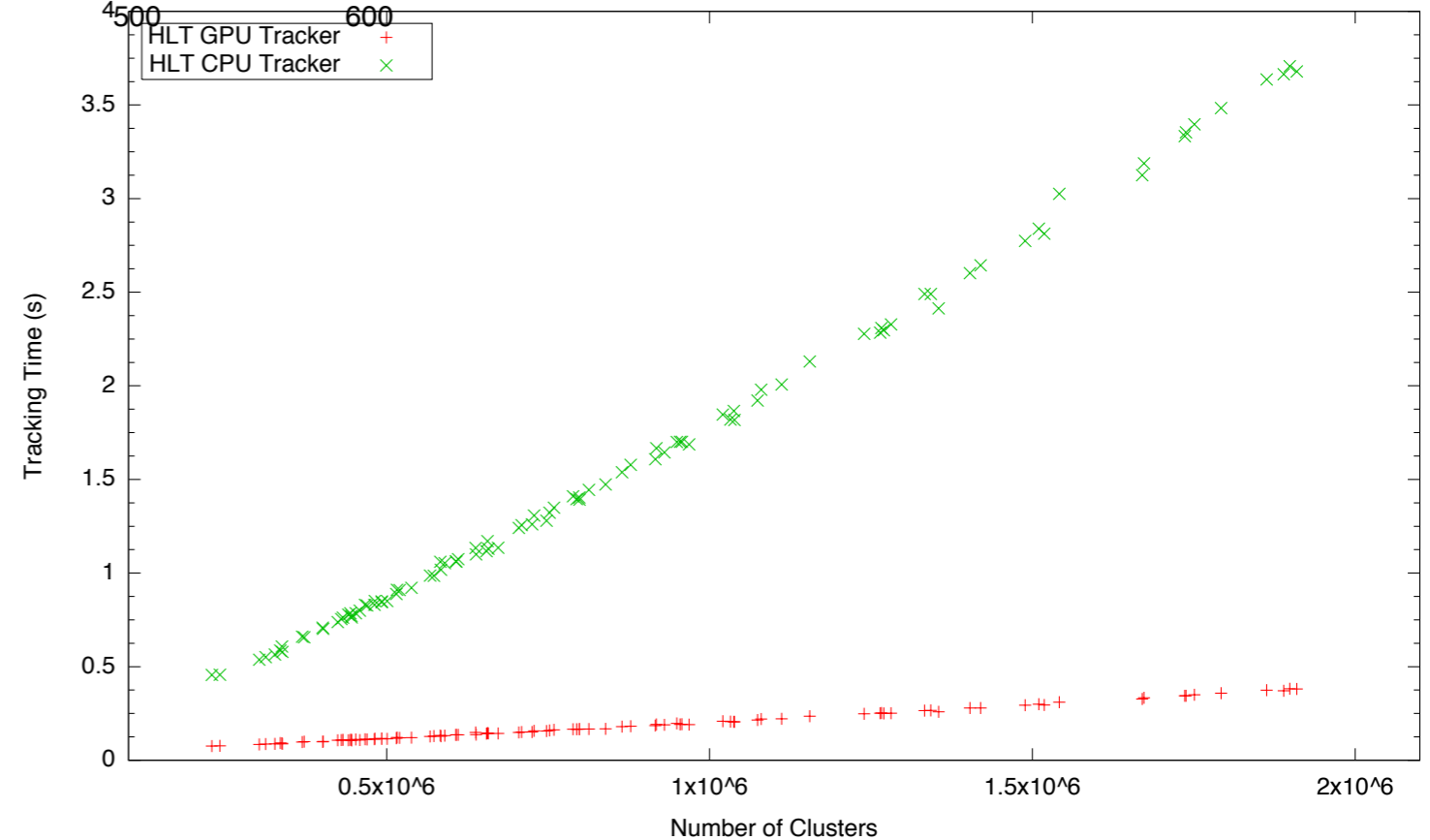**HLT farm**

**Experiment Control System**

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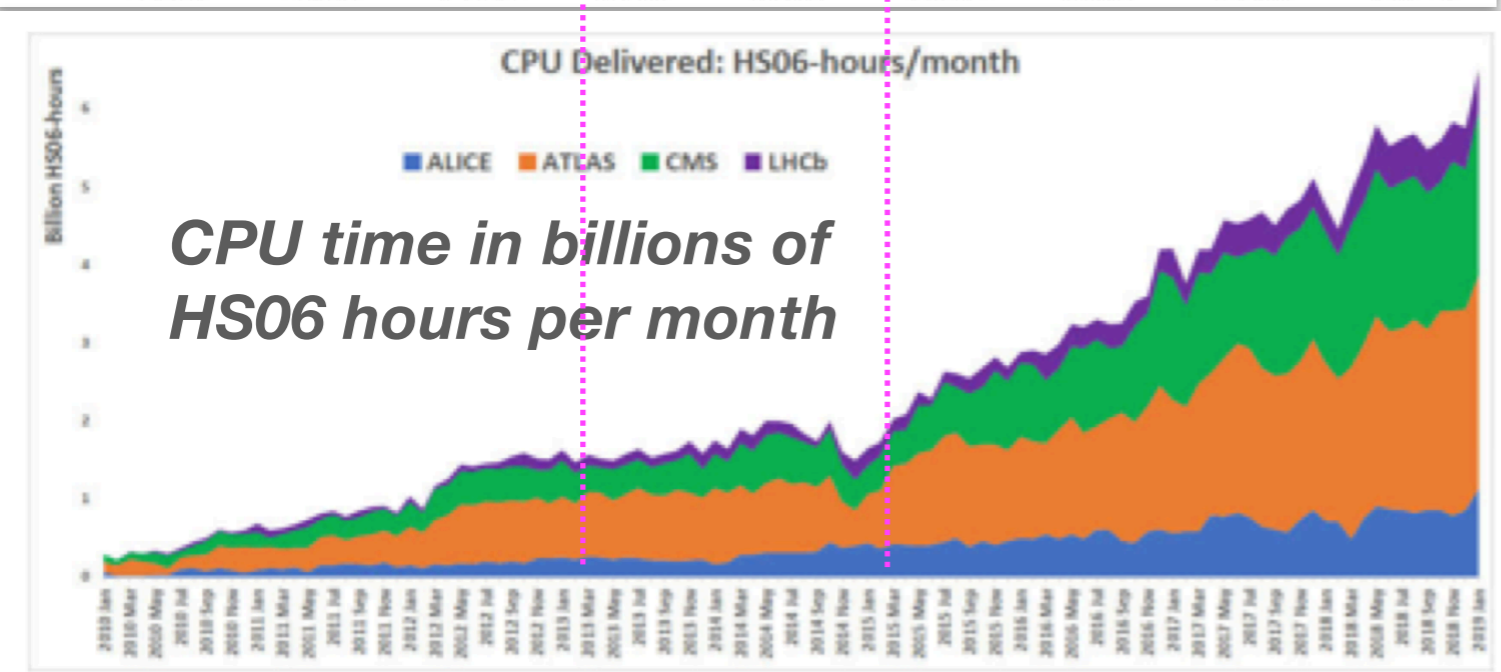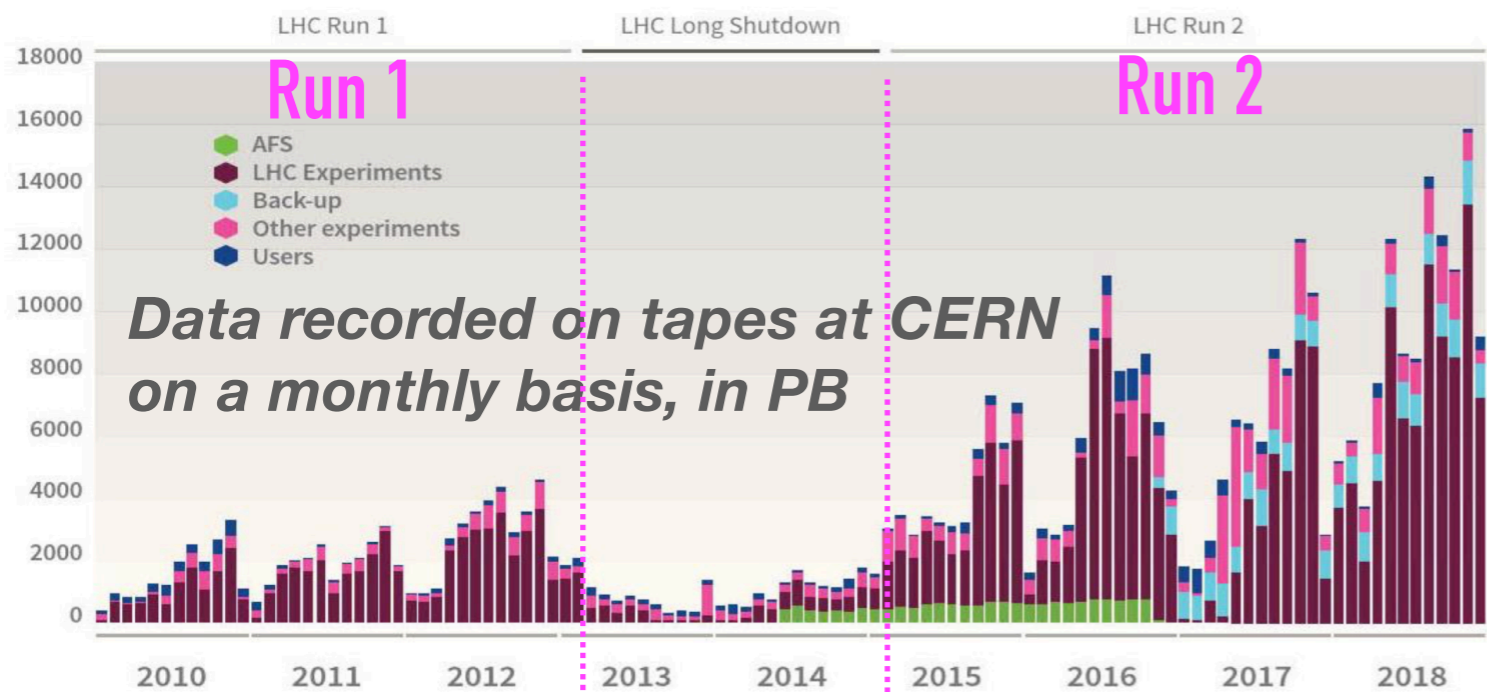  ➡ Ethernet UDP, with 10-15 events packed ⇒ ~ **80 kHz**

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*

*see [Ref]*

## 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