April 15 – 19, 2024
Edinburgh, United Kingdom

**12th Beam Telescopes and Test Beams Workshop**

# INTRODUCTION TO TDAQ AND ITS SCALING PRINCIPLES

*F.Pastore (Royal Holloway Univ. of London)*

*francesca.pastore@cern.ch*

➡ **Aim of this lecture is to introduce the basic TDAQ concepts, avoiding as many technological details as possible**

➡ **Focus on High Energy Physics**

    ➡ But key concepts are common to other areas

*Credits to A.Negri and W.Vandelli whose material helped me preparing these slides*

# OUTLINE

➡ **Introduction**
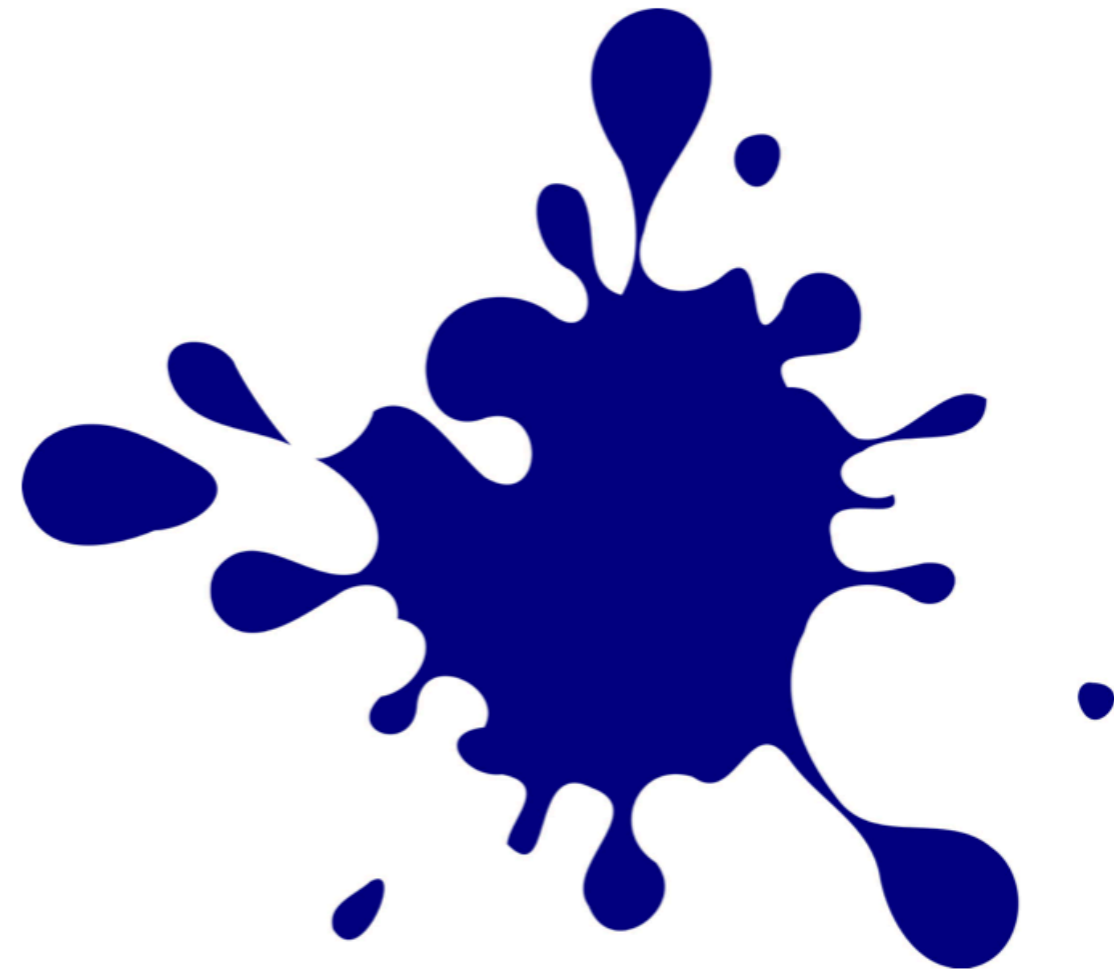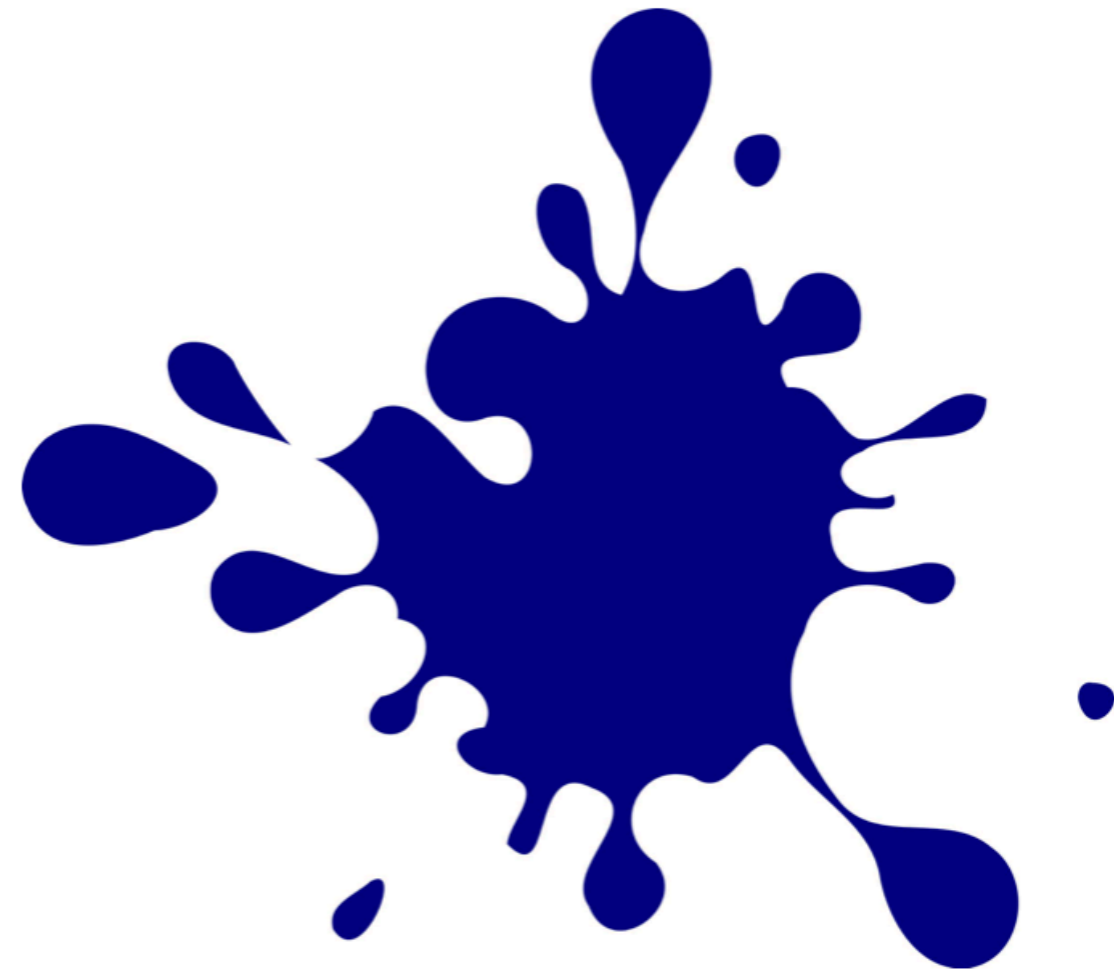  ➡ What is Trigger and DAQ?
  ➡ Overall TDAQ framework

➡ **Basic TDAQ concepts**
  ➡ Digitization, Latency
  ➡ Deadtime, Busy
  ➡ De-randomization

➡ **Scaling up**
  ➡ Readout and Event Building
  ➡ Buses vs Network

➡ **Fight bottlenecks**

➡ **Introduction**
  ➡ What is Trigger and DAQ?
  ➡ Overall TDAQ framework

➡ **Basic TDAQ concepts**
  ➡ Digitization, Latency
  ➡ Deadtime, Busy
  ➡ De-randomization

➡ **Scaling up**
  ➡ Readout and Event Building
  ➡ Buses vs Network

➡ **Fight bottlenecks**

# OUTLINE

➡ **Introduction**
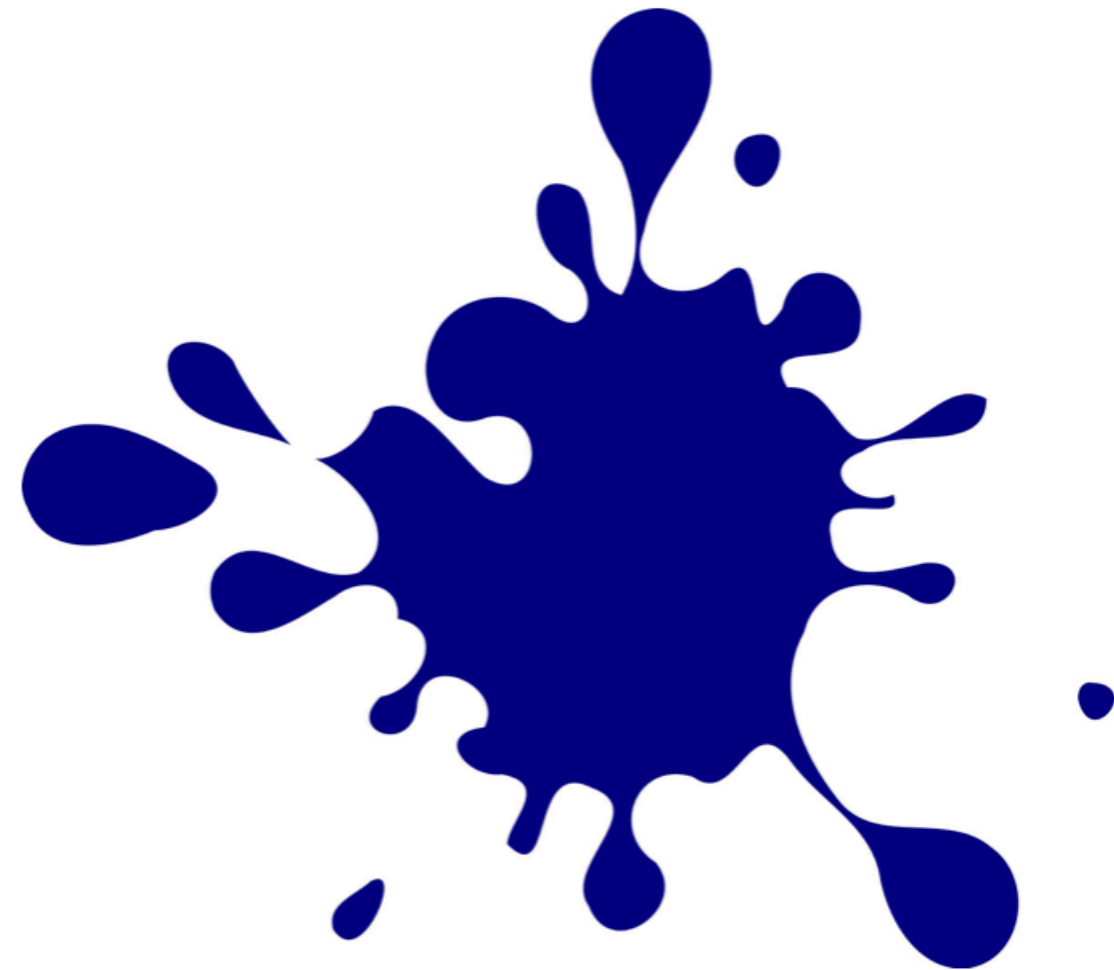  ➡ What is Trigger and DAQ?
  ➡ Overall TDAQ framework

➡ **Basic TDAQ concepts**
  ➡ Digitization, Latency
  ➡ Deadtime, Busy
  ➡ De-randomization
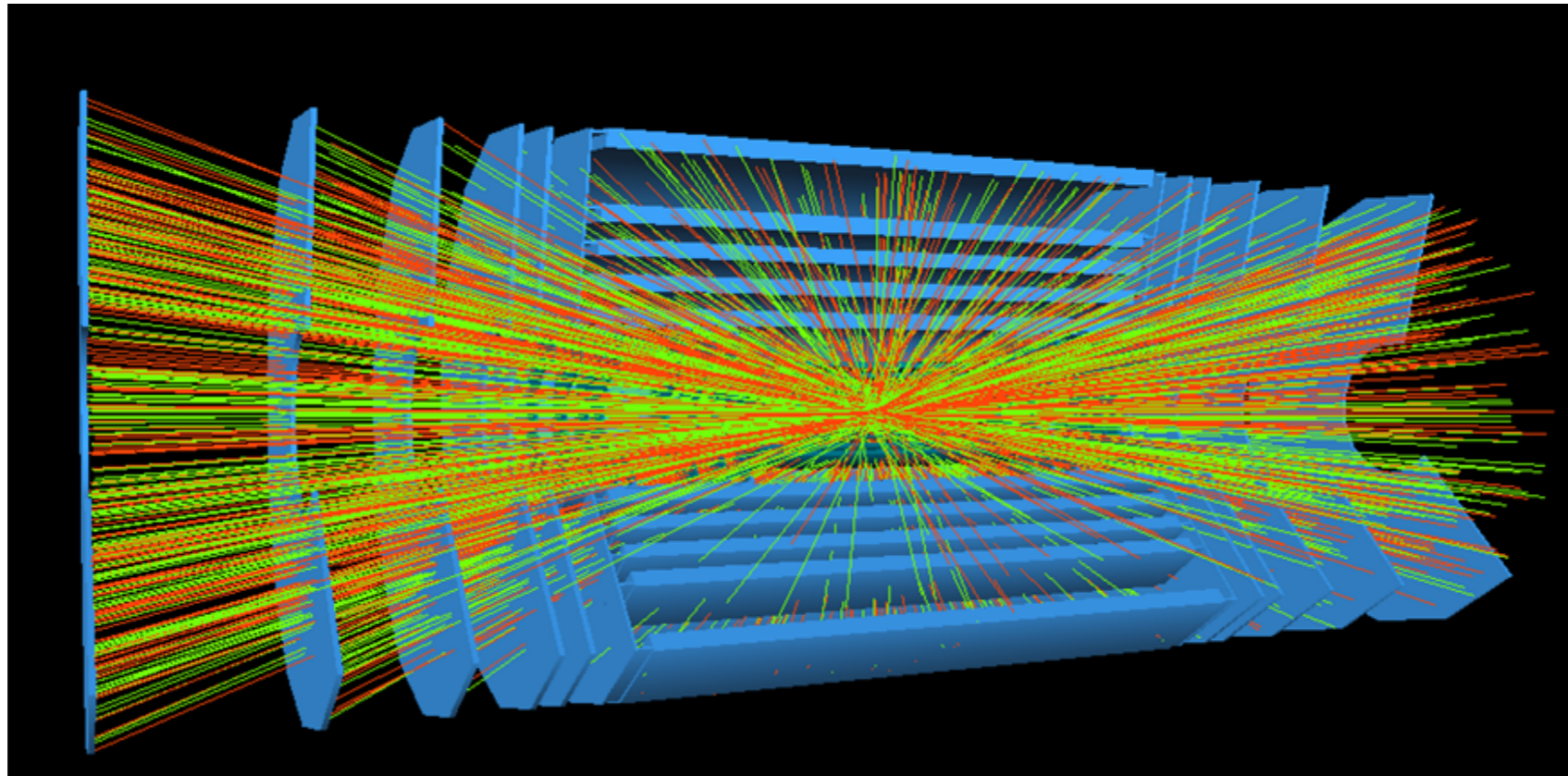
➡ **Scaling up**
  ➡ Readout and Event Building
  ➡ Buses vs Network
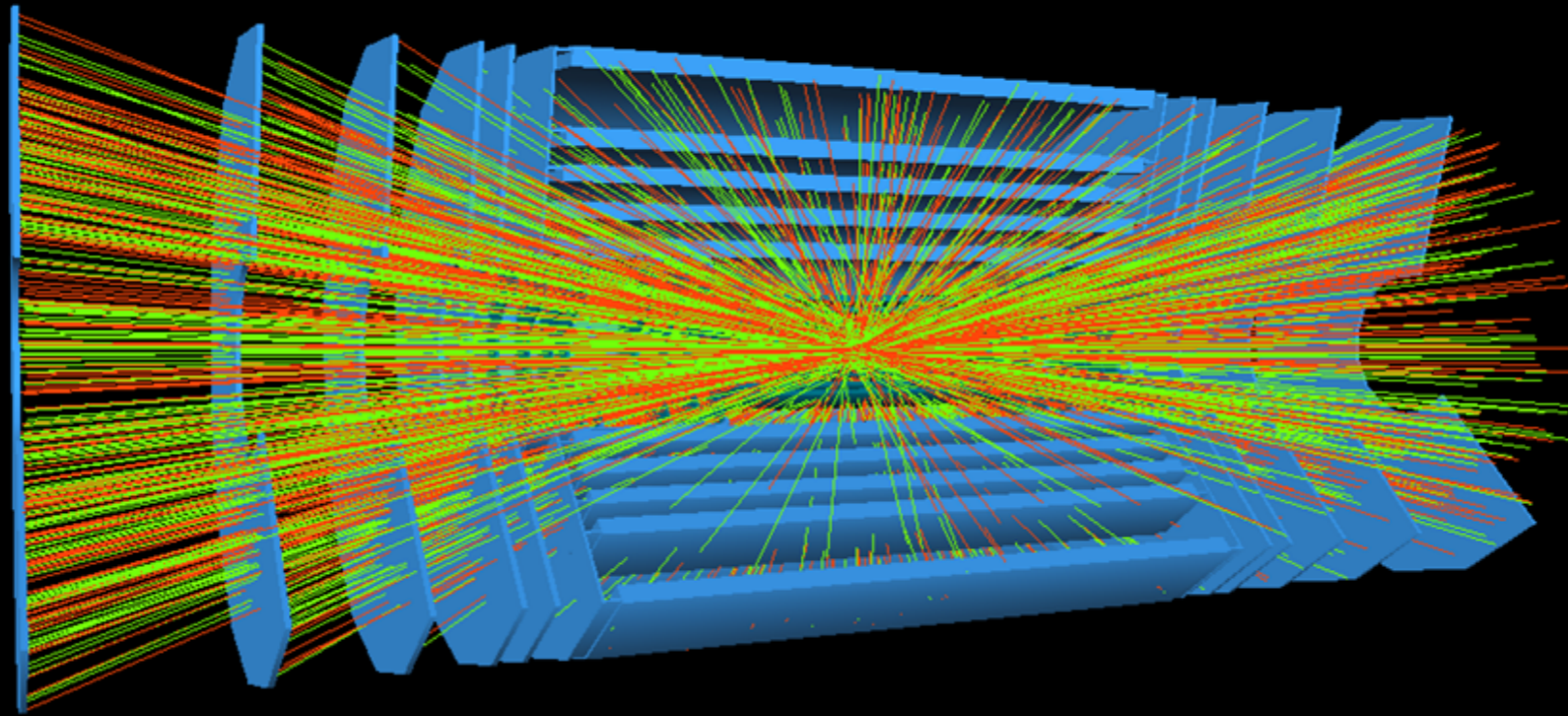
➡ **Fight bottlenecks**

# WHAT IS DAQ?



[Wikipedia]

➡ **Data AcQuisition (DAQ) is**

   ➡ the process of sampling signals that measure real world physical conditions

   ➡ and converting the resulting samples into digital numeric values that can be manipulated by a PC
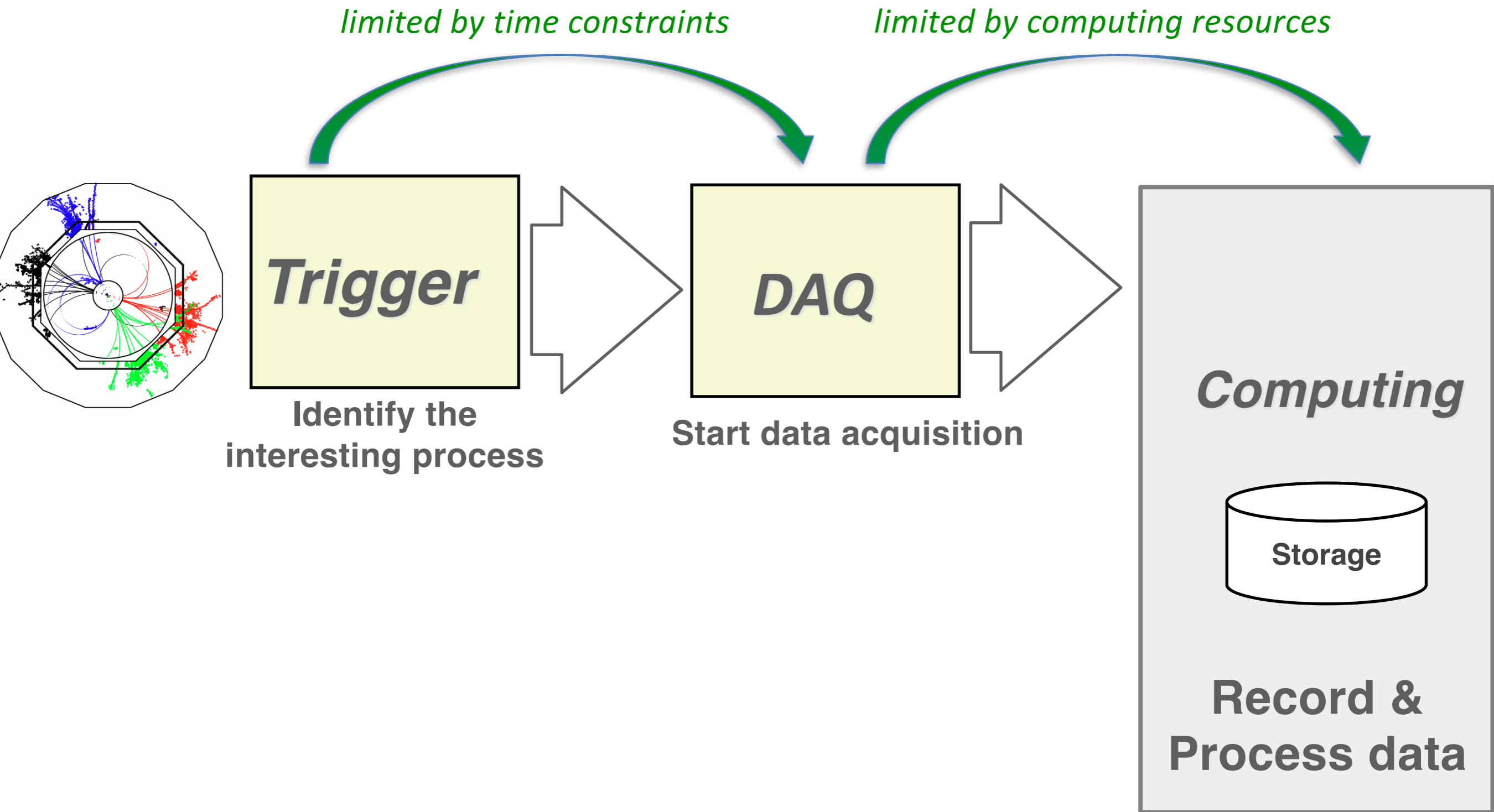
➡ **Main role of DAQ in HEP**

   ➡ process the signals generated in a detector

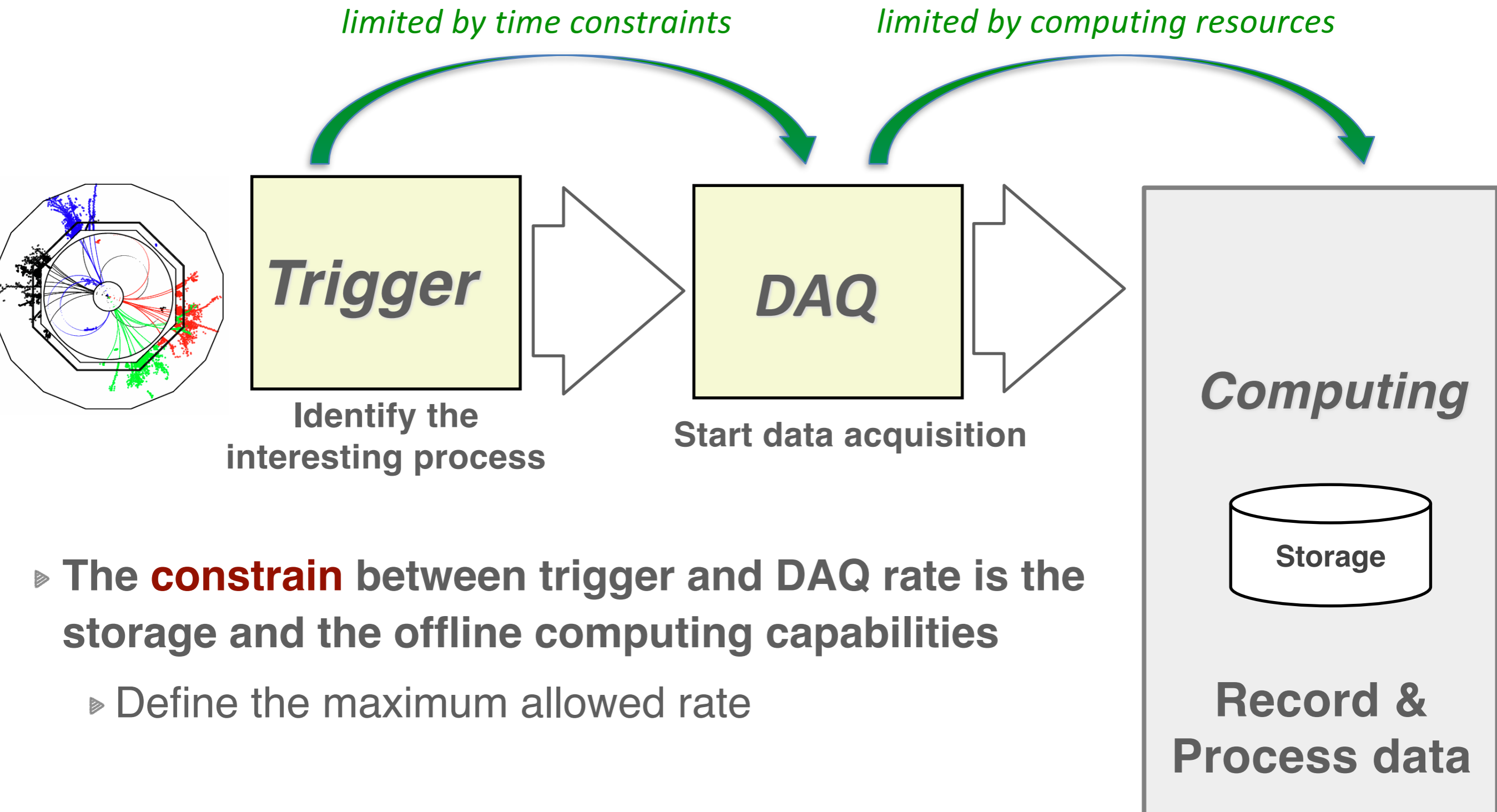   ➡ and saving the (interesting) information on a permanent storage

➡ **In many systems, like particle physics or astronomy experiments, to store all the possibly relevant data provided by the sensors is UNREALISTIC and often becomes also UNDESIRABLE**

➡ **Three approaches are possible:**
  ➡ Reduced amount of data (packing and/or filtering)   →*Trigger!*
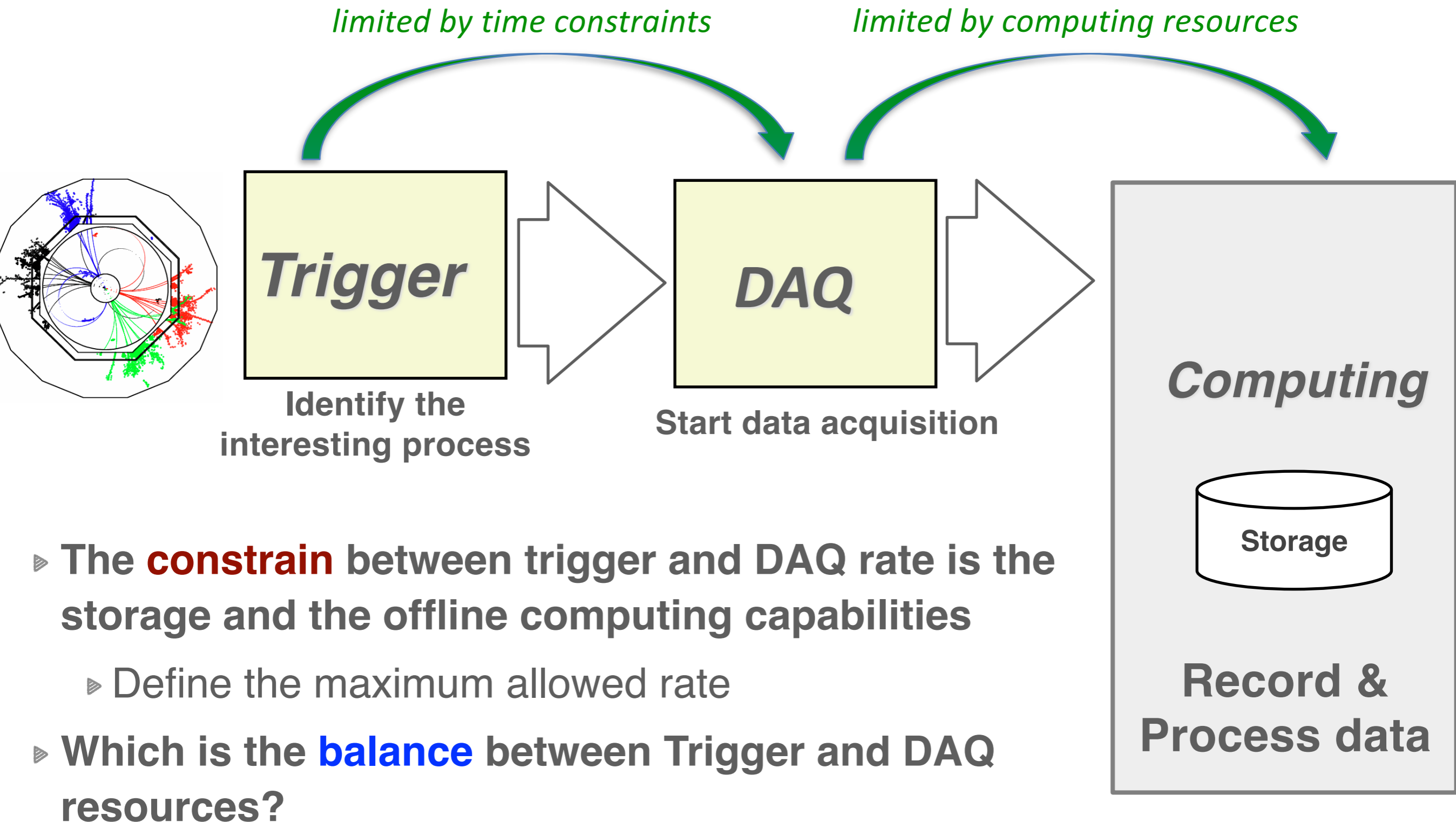  ➡ Faster data transmission and processing
  ➡ Both!

*limited by time constraints*

*limited by computing resources*

**Trigger**

**Identify the interesting process**

**DAQ**

**Start data acquisition**
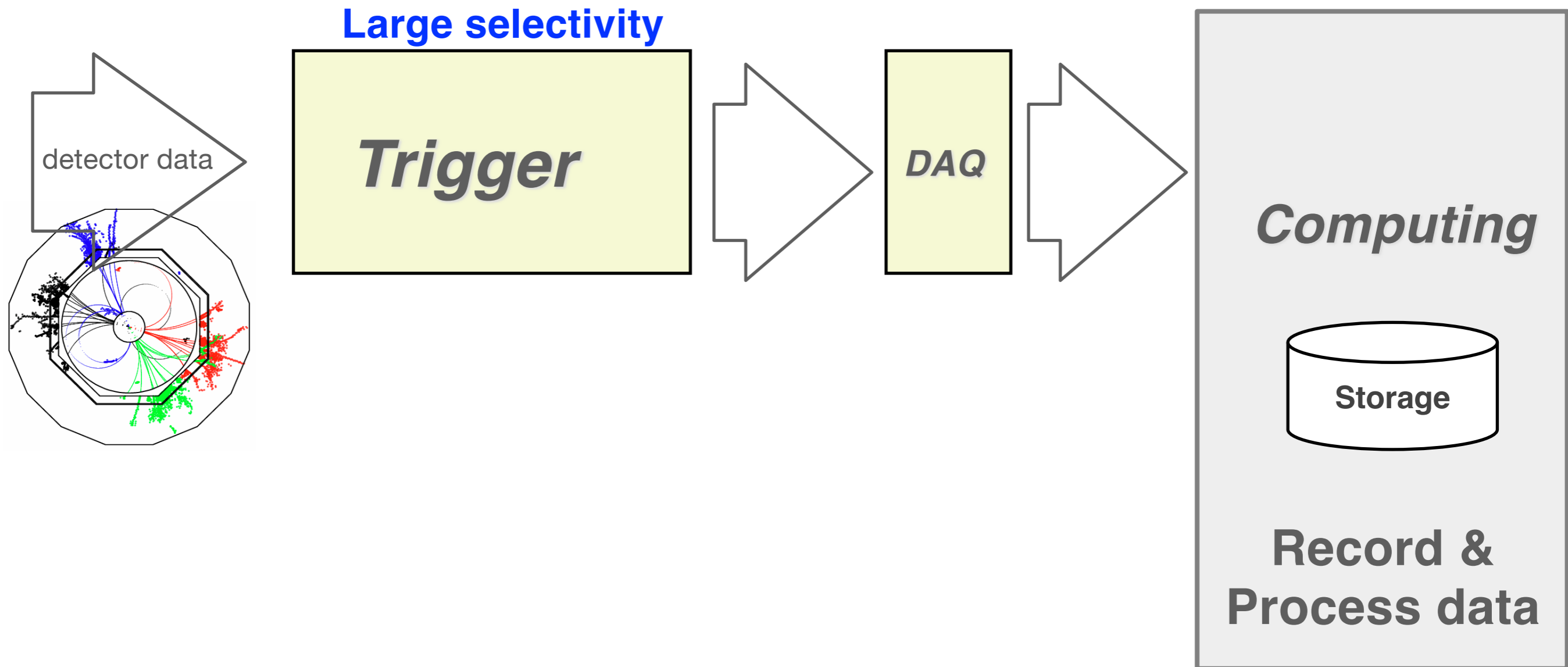
**Computing**

Storage

**Record & Process data**

▷ **The constrain between trigger and DAQ rate is the storage and the offline computing capabilities**

   ▷ Define the maximum allowed rate

*limited by time constraints*    *limited by computing resources*

**Trigger**

**Identify the interesting process**

**DAQ**

**Start data acquisition**

**Computing**

Storage

**Record & Process data**

▷ **The constrain between trigger and DAQ rate is the storage and the offline computing capabilities**
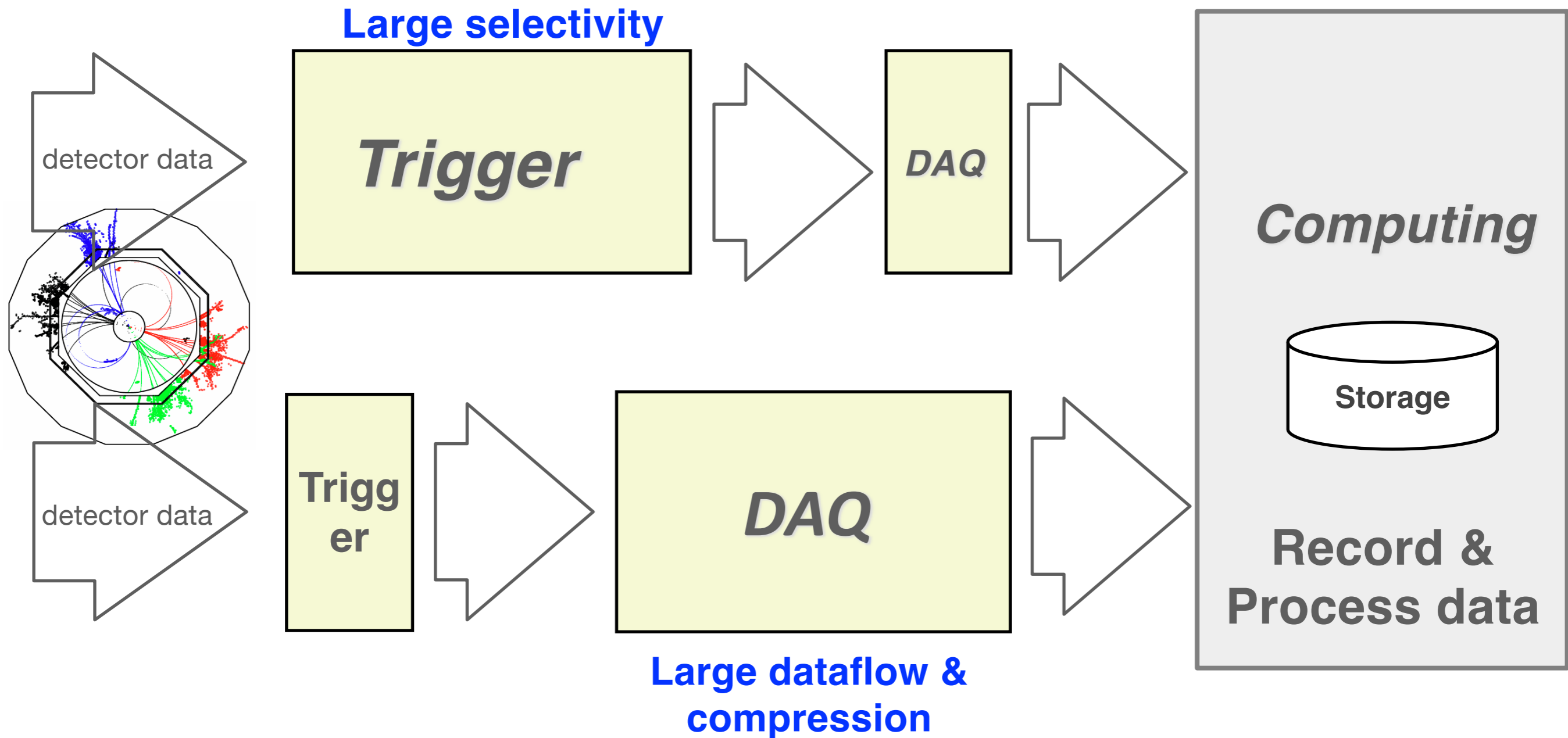
 ▷ Define the maximum allowed rate

▷ **Which is the balance between Trigger and DAQ resources?**

▷ **If the trigger is highly selective, one can reduce the size of the dataflow**

▷ **If the trigger is highly selective, one can reduce the size of the dataflow**

**Large selectivity**



▷ **If the selectivity of the trigger is not enough, due to large irreducible background, a large data flow (and data compression) is needed**

# TWO OPPOSITE EXAMPLES



▷ **LHC – ATLAS**

▷ **Project started in 1996**
▷ **Technology chosen in 2000**
▷ **Start data-taking 2008**

▷ **Full p-p collision rate: 40 MHz**
▷ **Average event size: 1.5 MB**
▷ **Full data rate:  ~60 PB/s**

▷ **Defined physics signal**
▷ **Selective trigger reduces 7 orders of magnitudes to ~kHz**

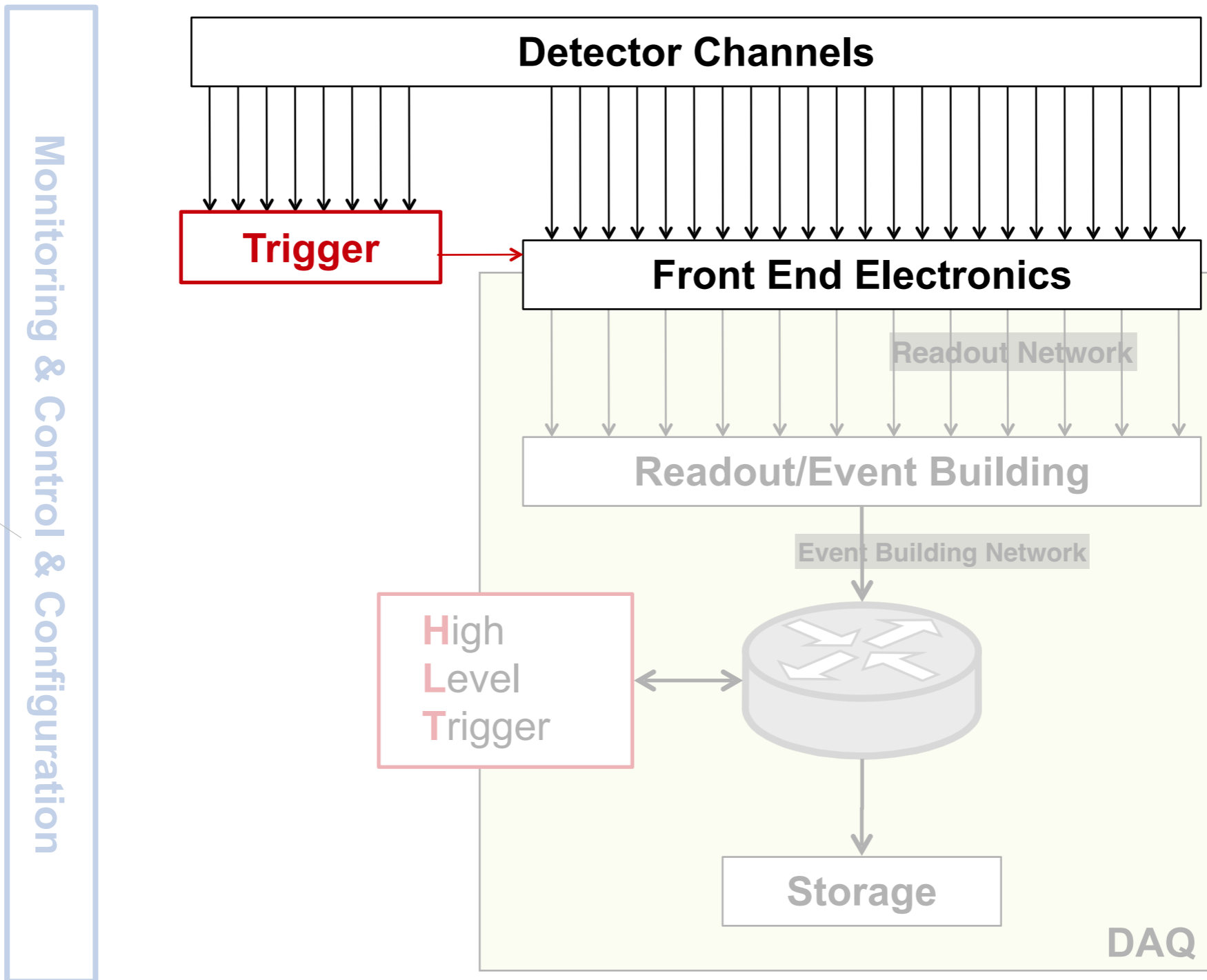▷ **Affordable DAQ rate: ~GB/s**
▷ **Data distribution (GRID)**



▷ **SKA (Square Km Array)**

▷ **Project started in 2011**
▷ **Technologies under evaluation**
▷ **Start operations in 2028**

▷ **Radio-photograph the sky continuously**
▷ **1.12 PB/s of photos collected**

▷ **EXASCALE system: $10^{18}$ operations for correlation and imaging**
▷ **Simple correlator : 10 TB/s**
▷ **Total Internet Traffic ≈ 8 TB/s in 2010**

▷ **Required large computing power**
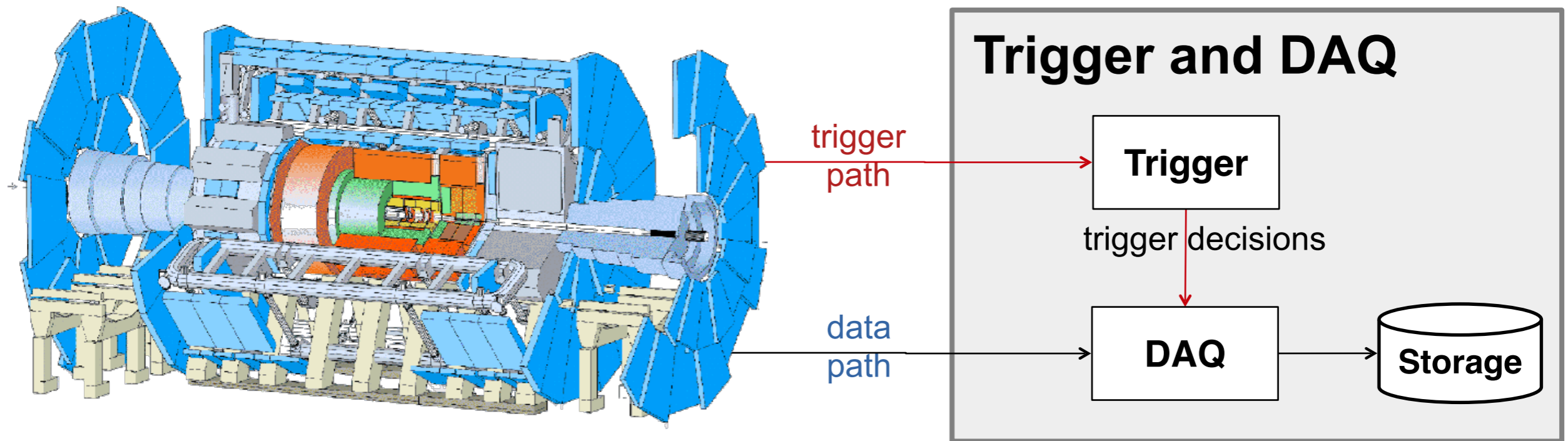▷ **Big-data and cloud-computing drive market**

# DOUBLE PATHS

➡ **Trigger path**
- ➡ From dedicated detectors to trigger logic
- ➡ online selection

➡ **Data path**
- ➡ From all the detectors to storage
- ➡ On positive trigger decision

➡ **Use discriminating features within widely extended systems**

➡ **Reality is:**

 ➡ The trigger accepts events with features <u>similar</u> to the signal

 ➡ The final rate is often dominated by not interesting physics

# TRIGGER DUTIES

➡ **Either selects interesting events or rejects boring ones, in real time**

  ➡ Selective: **efficient** for "signal"
    and **resistant** to "background"

  ➡ Simple and robust

  ➡ Quick

➡ **Either selects interesting events or rejects boring ones, in real time**

  ➡ Selective: **efficient** for "signal"
     and **resistant** to "background"

  ➡ Simple and robust

  ➡ Quick

➡ **With minimal controlled latency**

  ➡ time it takes to form and distribute its decision

➡ **Either selects interesting events or rejects boring ones, in real time**

  ➡ Selective: **efficient** for "signal"
    and **resistant** to "background"

  ➡ Simple and robust

  ➡ Quick

➡ **With minimal controlled latency**

  ➡ time it takes to form and distribute its decision

➡ **Generates a prompt signal used to start the data-acquisition processes**

  ➡ To be distributed to front end electronics

- ➡ **Either selects interesting events or rejects boring ones, in real time**
  - ➡ Selective: **efficient** for "signal"
    and **resistant** to "background"
  - ➡ Simple and robust
  - ➡ Quick
- ➡ **With minimal controlled latency**
  - ➡ time it takes to form and distribute its decision
- ➡ **Generates a prompt signal used to start the data-acquisition processes**
  - ➡ To be distributed to front end electronics
- ➡ **Trigger and Front-End electronics have common design**
  - ➡ Data compression and formatting
  - ➡ Monitor and automatic fault detection

- *"Method of Registering Multiple Simultaneous Impulses of Several Geiger Counters"*
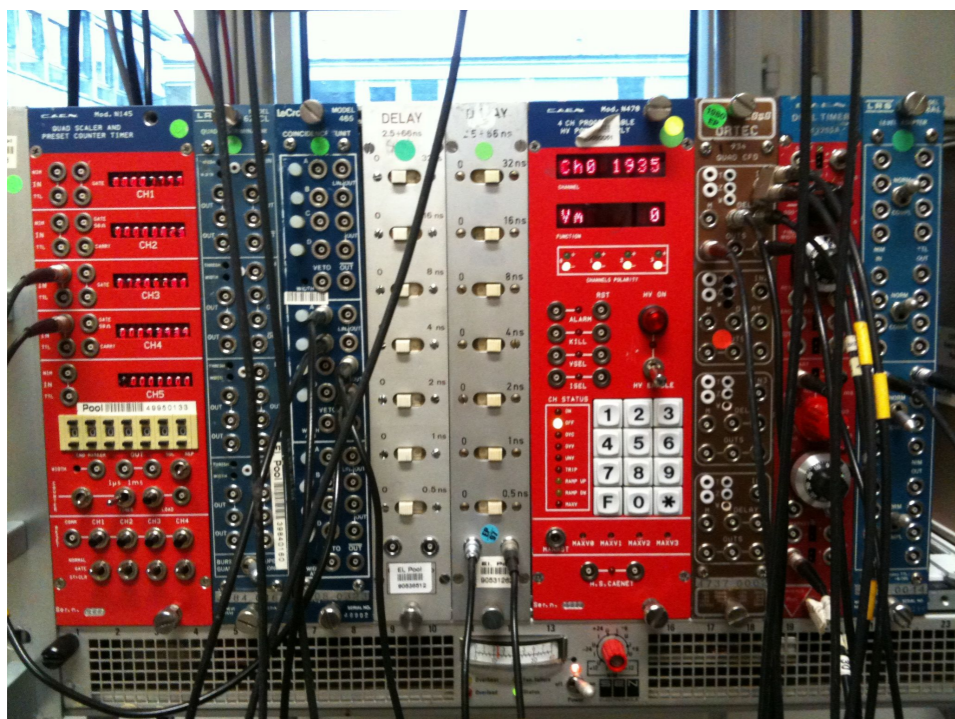
  **Bruno Rossi, Nature 1930**

  – Online coincidence of three signals



Geiger counters

Primary particle

Secondary particle

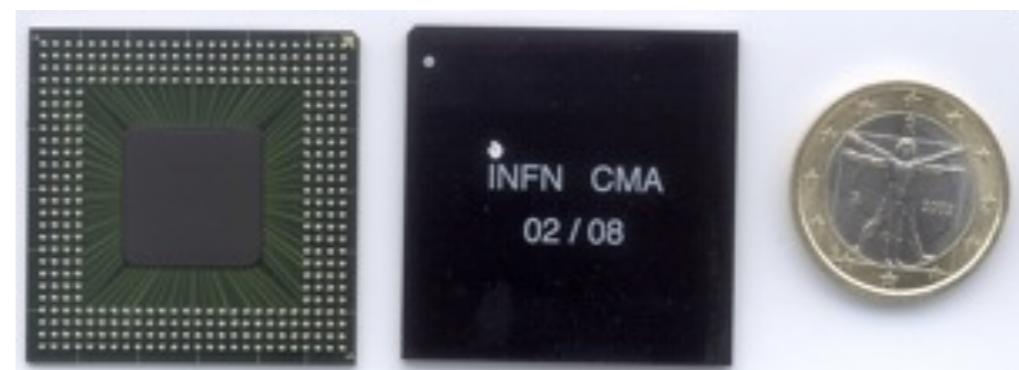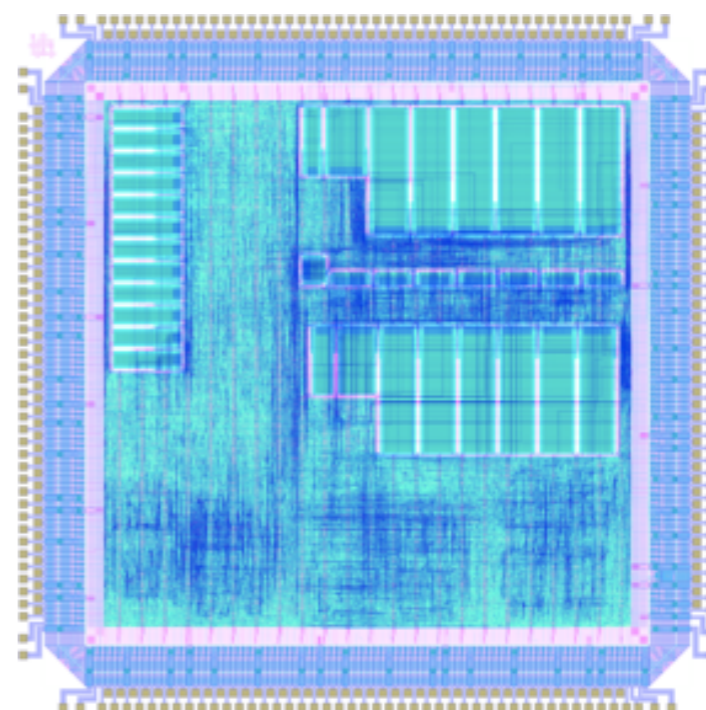Astronomia e Fisica a Firenze: dalla Specola ad Arcetri, Firenze Universiry Press, 2017

## Modular electronics

- Simple algorithms
- Low-cost
- Intuitive and fast use





## Digital integrated systems

- Highly complex algorithms
- Fast signals processing
- Knowledge of digital systems

# THREE TYPES OF TRIGGER

➡ **Global:**

    ➡ an external system identifies the "interesting" event, all the readout data is collected for that **event identifier**

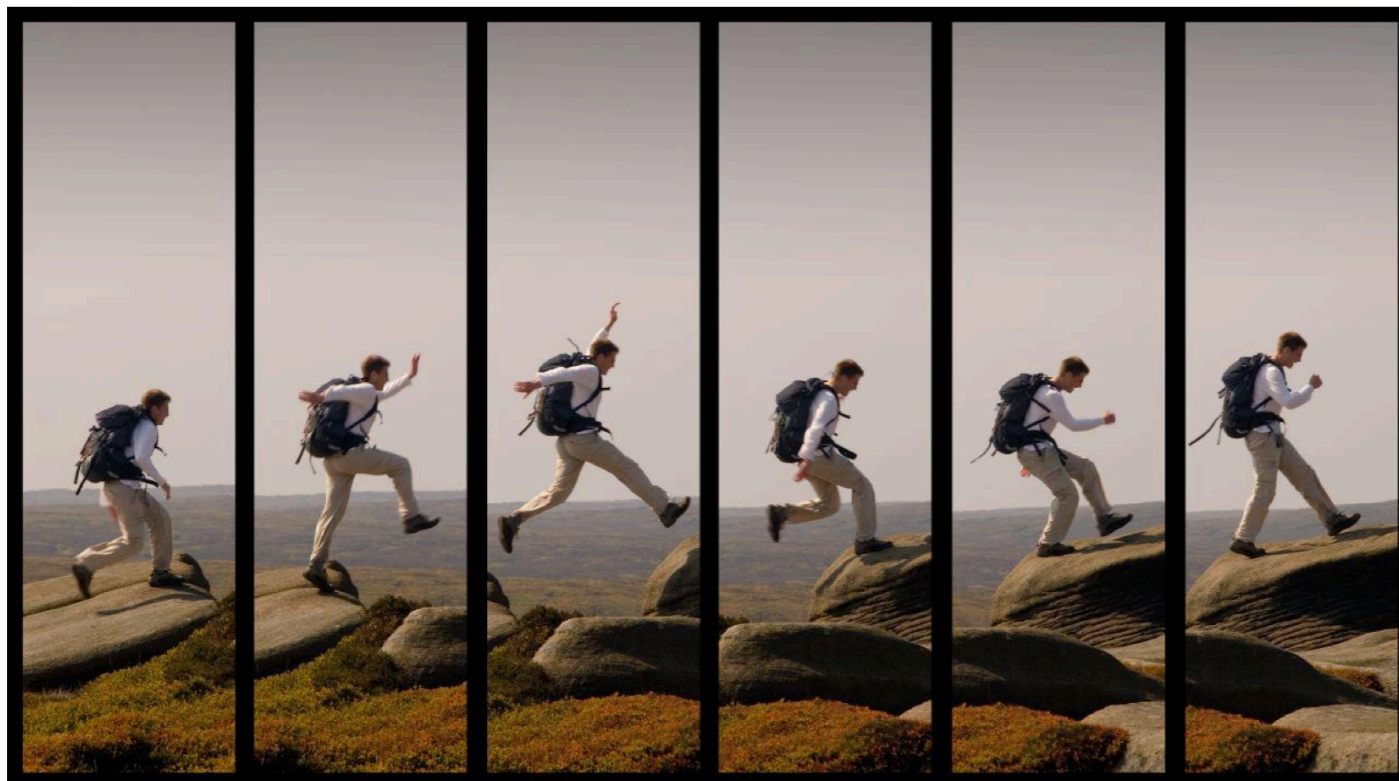# THREE TYPES OF TRIGGER

➡ **Global:**

   ➡ an external system identifies the "interesting" event, all the readout data is collected for that **event identifier**

➡ **Local:**

   ➡ local trigger decision to readout data on the **local front-end modules**, readout collects fragments corresponding to that trigger

# THREE TYPES OF TRIGGER

➡ **Global:**

  ➡ an external system identifies the "interesting" event, all the readout data is collected for that **event identifier**

➡ **Local:**

  ➡ local trigger decision to readout data on the **local front-end modules**, readout collects fragments corresponding to that trigger

➡ **Continuous readout:**

  ➡ front-end sends data continuously to the readout, at a fixed rate, regardless the data content. **Data size and rate are constant in size**. Readout cannot group fragments relative to an event



By Rick Harrison (license)

*not really a photo, almost a movie*

➡ **use cases:**

   ➡ **Colliders**: normally use global trigger: if something interesting has been seen somewhere, take all the data corresponding to that bunch crossing

   ➡ **Large distributed telescopes**: often use local trigger: readout data for the portions of the detector that have seen something

   ➡ **Very slow detectors**: sometimes use continuous readout: sample the analogue signals at a fixed rate and let the downstream DAQ decide whether there were any interesting signals



By Rick Harrison (license)

*not really a photo, almost a movie*

➡ **Gather data produced by detectors**

  ➡ Readout

➡ **Form complete events**

  ➡ Data Collection and Event Building

➡ **Possibly feed other trigger levels**

  ➡ High Level Trigger

➡ **Store event data**

  ➡ Data Logging

➡ **Manage the operations**

  ➡ Run Control, Configuration, Monitoring
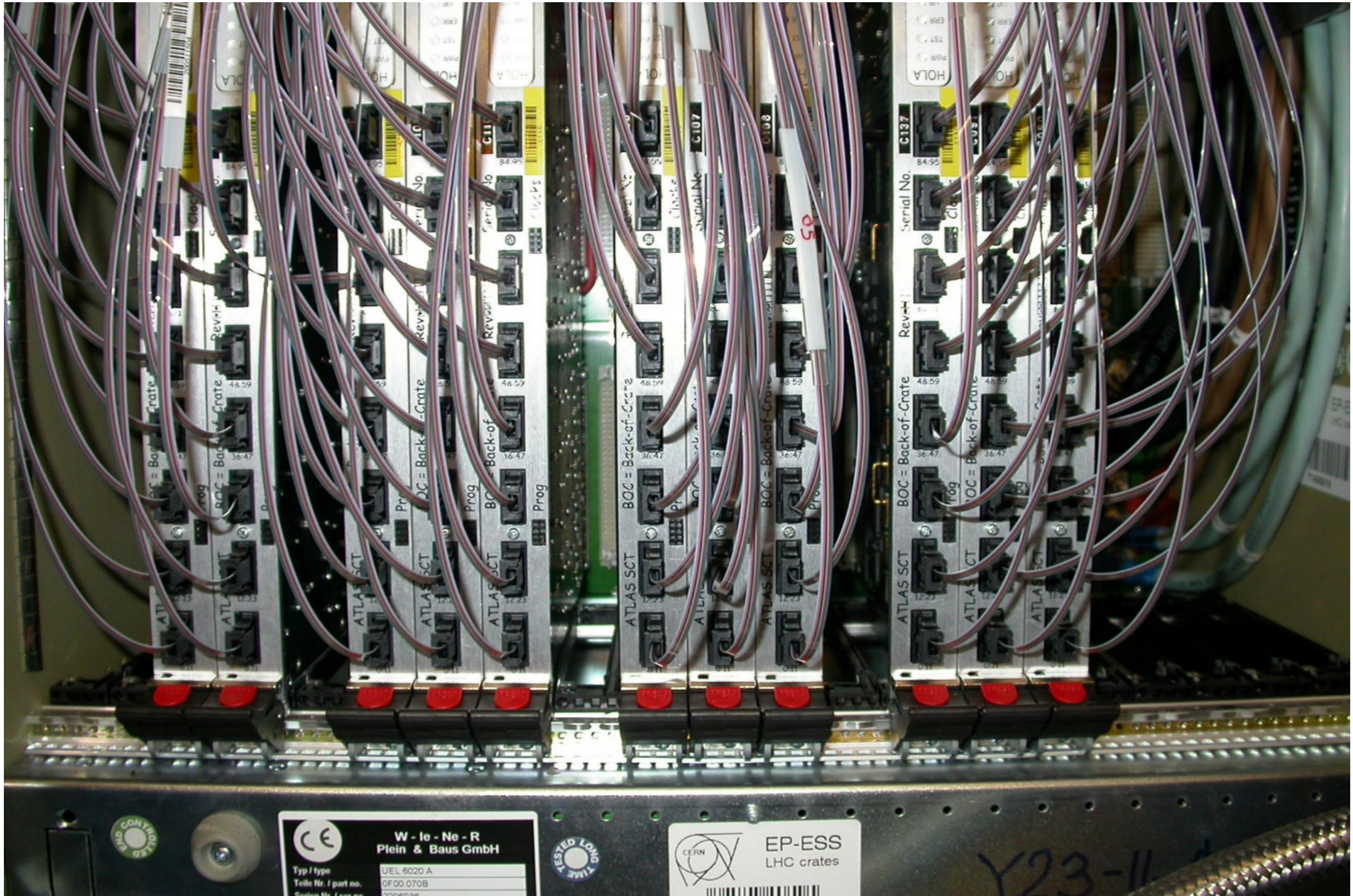
**Data Flow**

➡ **Intermediate crates off-detector to separate FE (long duration) and PCs**

➡ **Signal processing**, data formatting, parallelizable tasks (pattern recognition), machine learning, ...

   ➡ FPGAs are becoming the bread&butter of TDAQ

➡ **High-speed serial links**, electrical and optical, depending on distance

   ➡ Low-power LVDS, **400 Mbps**, < 10m

   ➡ Optical GHz-links for longer distances (up to 100 m)

➡ **High density backplanes** for data exchanges in crates

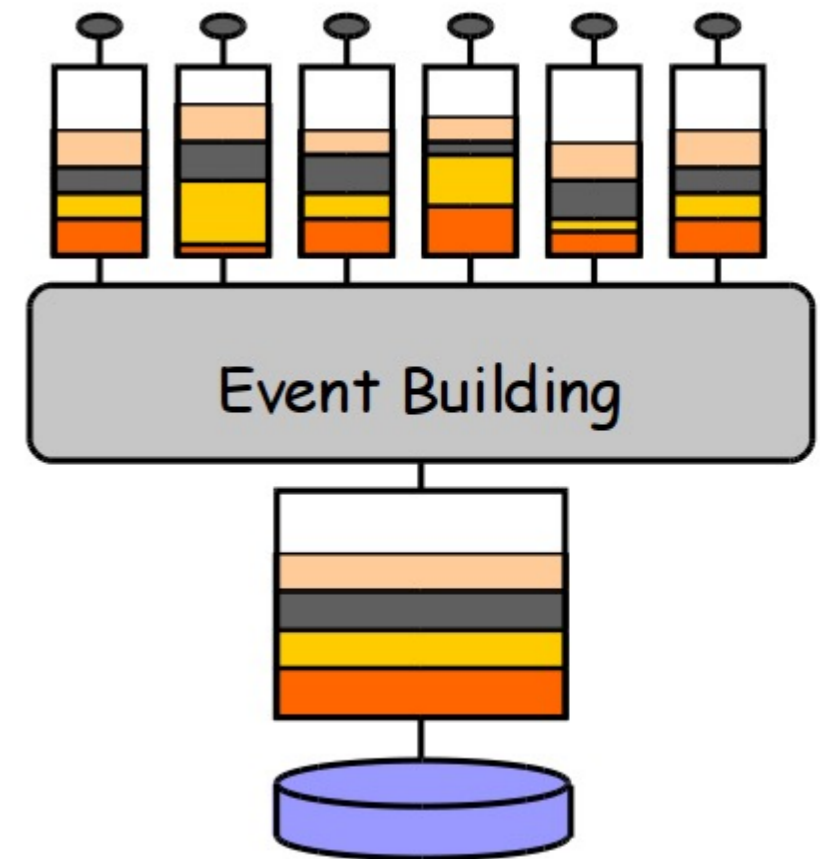   ➡ High pin count, with point-to-point connections up to **160 Mbit/s**

   ➡ Large boards preferred

➡ **Collects data from Front-End and associate fragments corresponding to**
  - ➡ the same event, e.g. same bunch crossing
  - ➡ the same accelerator orbit
  - ➡ the same time frame

➡ **Data must be marked with time-stamp**

➡ **Work done with: a distribution system (networks) and processing units (switch/PCs/custom board)**
  - ➡ Can be incremental on multiple networks
  - ➡ Or a separate network for data collection

➡ **Usually adopt the farm architecture**
  - ➡ assign one event per processor (node)
  - ➡ larger latency, but scalable

➡ **Network has intrinsic latency, so traffic control is critical**
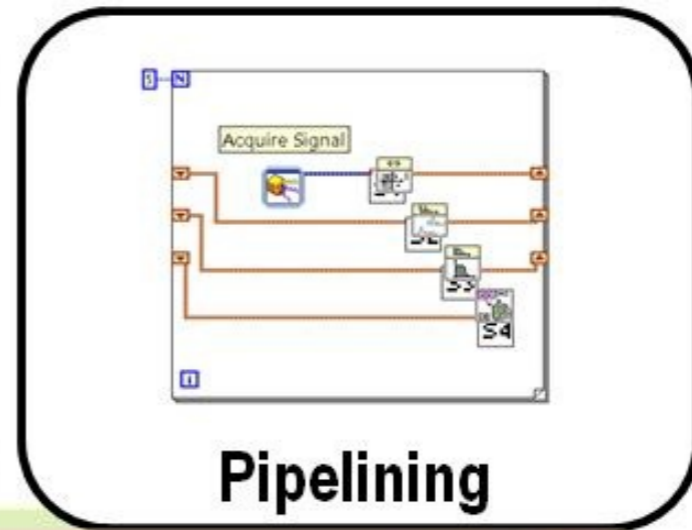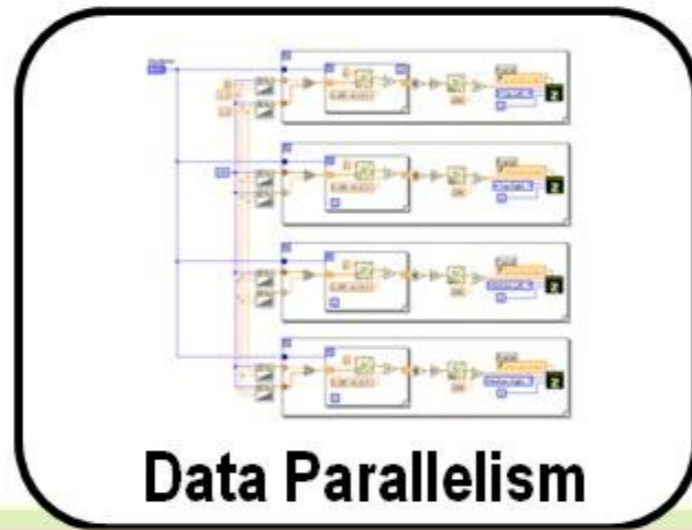  - ➡ can have one network only for traffic control



Event Building

**Task Parallelism**

**Data Parallelism**

**Pipelining**

**NVIDIA GPUS:**
**3.5 B TRANSISTORS**

**VIRTEX-7 FPGA:**
**6.8 B TRANSISTORS**

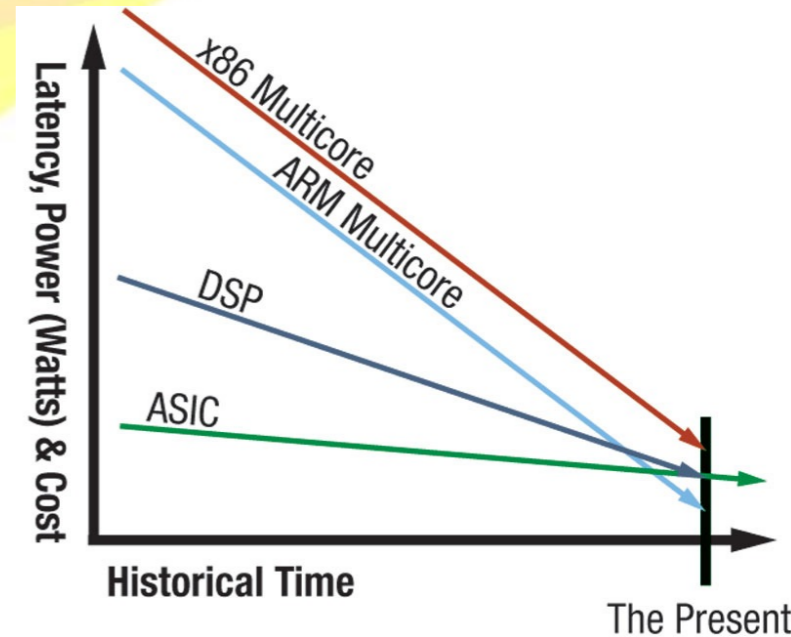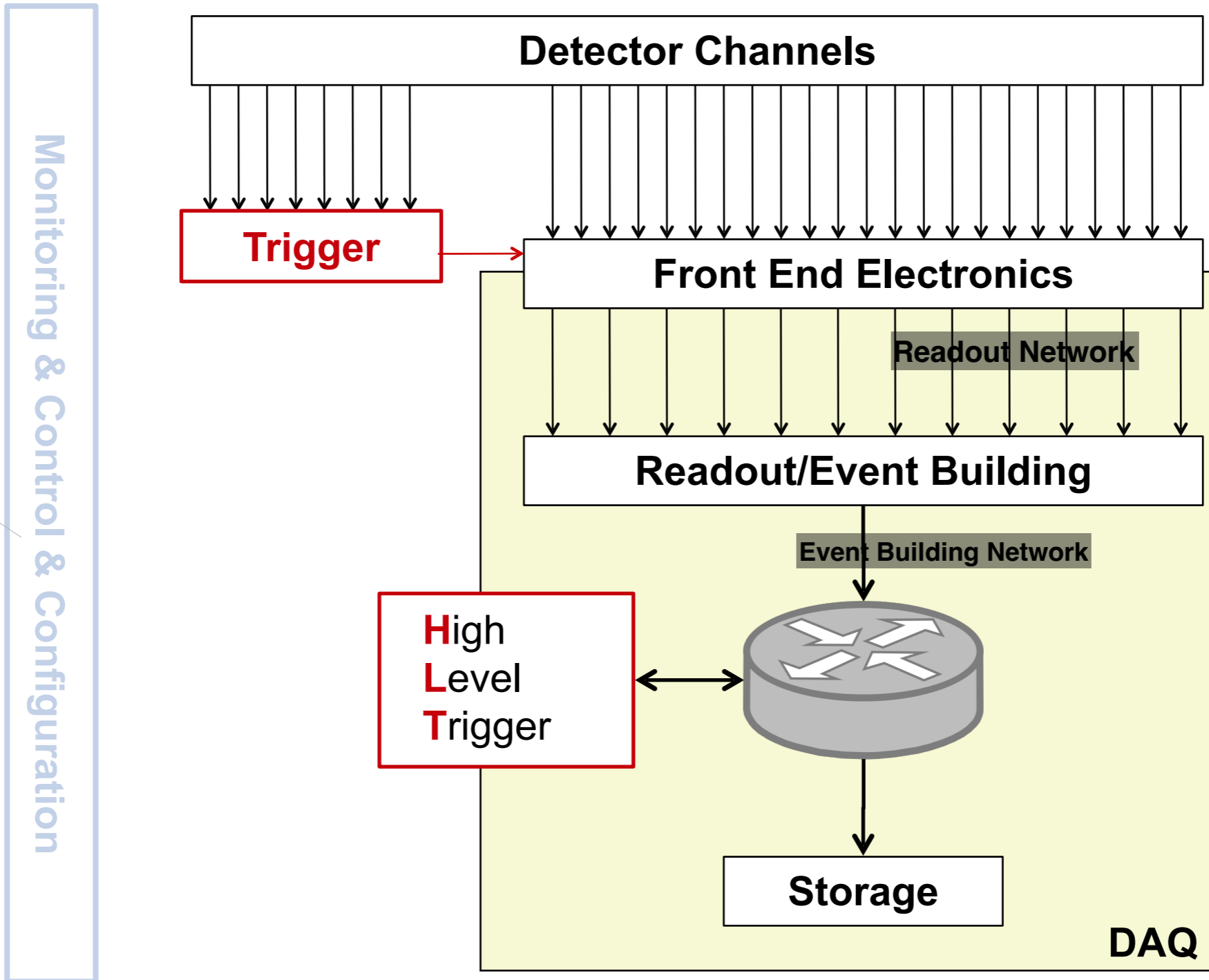**Multicore Processors**

**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.**
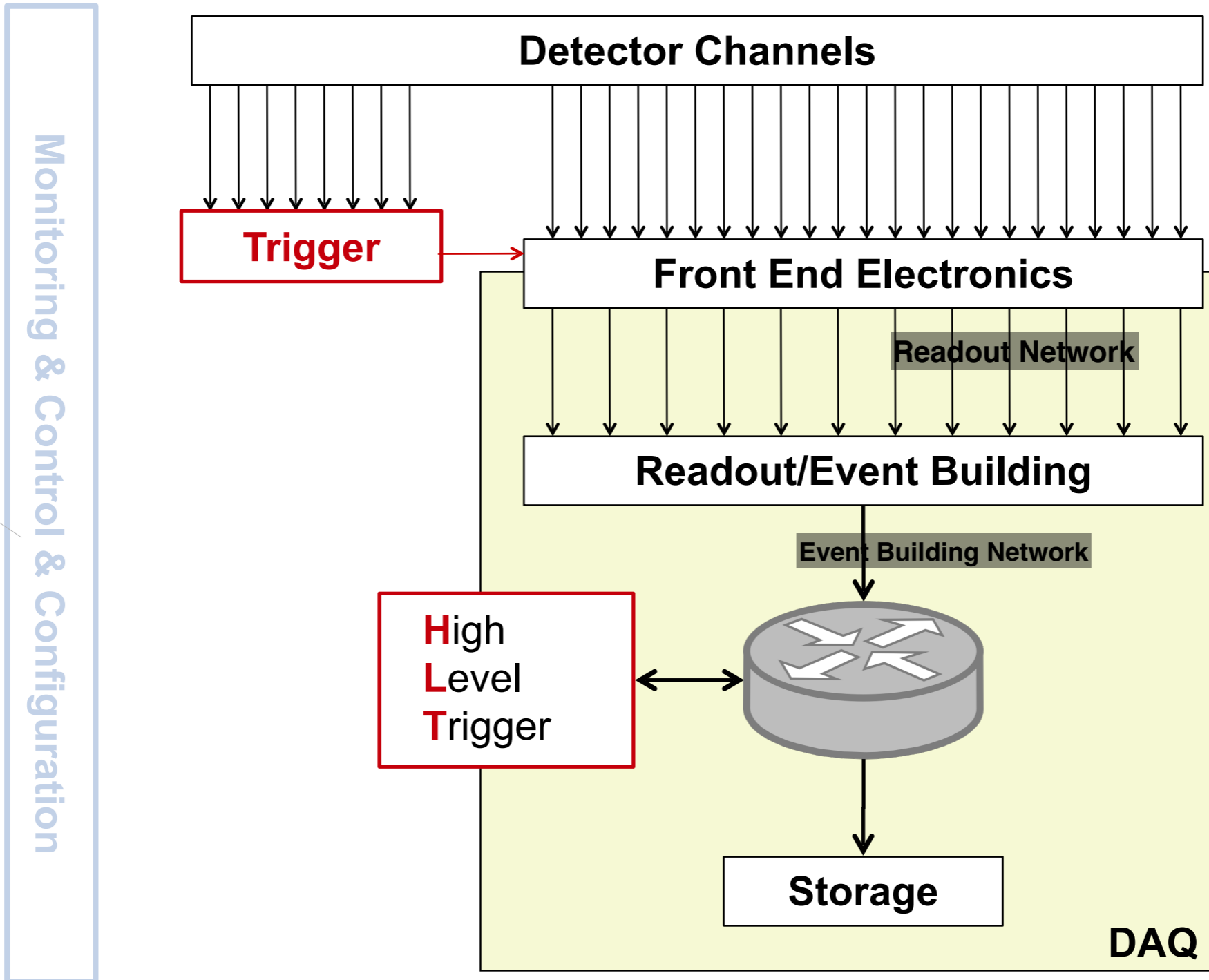
➡ **Storage device technologies gaining importance in HEP**

    ➡ Storage data rate increasing with luminosity

    ➡ Distributed file systems being used as data-flow frameworks

        ➡ CMS, ATLAS run 4 (?), …

    ➡ Also use large temporary buffers with high rate access

        ➡ LHCb: 40 PB (3000 hard-disks) enough for days

        ➡ SSD faster but have short lifetime wrt
high read-write rate, so prefer hard-disks
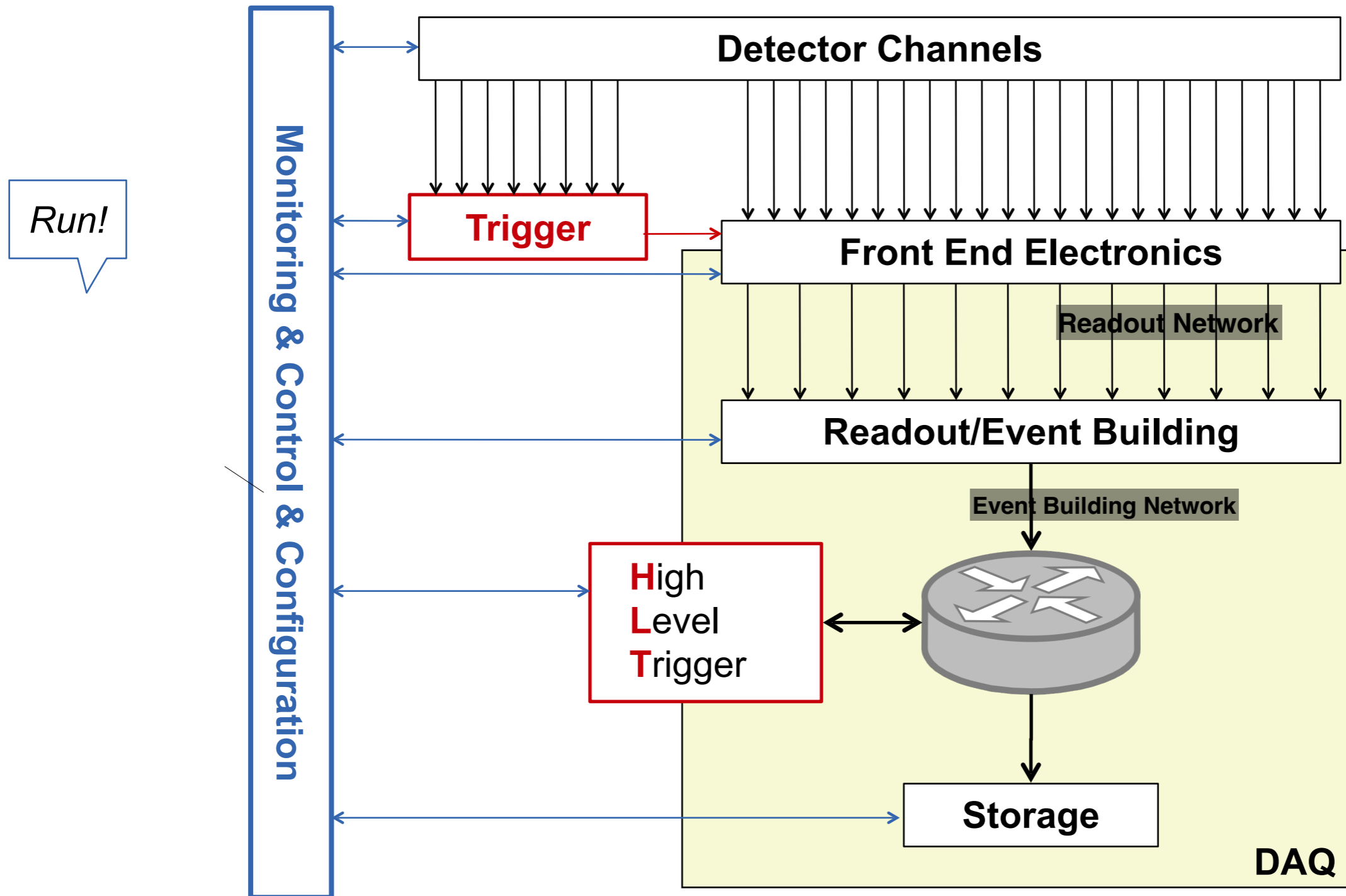
# THE GLUE OF YOUR EXPERIMENT

➡ **Configuration**

    ➡ data taking or test?

➡ **Control**

    ➡ Orchestrate applications participating to data taking

    ➡ Via distributed **Finite State Machine**

➡ **Monitoring**

    ➡ What is going on?

    ➡ What happened?

    ➡ When? Where?

➡ **Introduction**

   ➡ What is Trigger and DAQ?

   ➡ Overall TDAQ framework

➡ **Basic TDAQ concepts**

   ➡ Digitization, Latency

   ➡ Deadtime, Busy

   ➡ De-randomization

➡ **Scaling up**

   ➡ Readout and Event Building

   ➡ Buses vs Network

**Via a toy model**

➡ **Eg: measure temperature at a fixed frequency**

  ➡ Clock trigger

➡ **ADC performs analog to digital conversion, digitization (on front-end electronics)**

  ➡ Encoding analog value into binary representation

➡ **CPU does**

  ➡ Readout, Processing, Storage



**Physical View**

T sensor — ADC Card — CPU ← disk

➡ **System clearly limited by the time $\tau$ to process an "event"**

   ➡ ADC conversion + CPU processing + Storage

➡ **The DAQ maximum sustainable rate is simply the inverse of $\tau$, e.g.:**

   ➡ E.g.: $\tau$ = 1 ms ® R = $1/\tau$ = 1 kHz

**Physical View**

T sensor

ADC Card

CPU

disk

TRIGGER

ADC

Processing

disk

$\tau$ = 1 ms

➡ **Events are asynchronous and unpredictable**

   ➡ E.g.: beta decay studies

➡ **A physics trigger is needed**

   ➡ **Discriminator**: generates an output digital signal if amplitude of the input pulse is greater than a given threshold

➡ NB: delay introduced to compensate for the trigger latency

   ➡ Signal split in trigger and data paths

➡ **Events are asynchronous and unpredictable**

➡ E.g.: beta decay studies

➡ **A physics trigger is needed**

➡ **Discriminator**: generates an output digital signal if amplitude of the input pulse is greater than a given threshold

➡ **NB: delay introduced to compensate for the trigger latency**

➡ Signal split in trigger and data paths

➡ **Stochastic process**

   ➡ Fluctuations in time between events

➡ **Let's assume for example**

   ➡ physics rate f = 1 kHz, i.e. $\lambda$ = 1 ms

   ➡ and, as before, $\tau$ = 1 ms

➡ **Stochastic process**

 ➡ Fluctuations in time between events

➡ **Let's assume for example**

 ➡ physics rate f = 1 kHz, i.e. $\lambda$ = 1 ms

 ➡ and, as before, $\tau$= 1 ms



Probability of time (in ms) between events for average decay rate of f=1kHz → $\lambda$=1ms

Trigger path

Data path

TRIGGER

discriminator

delay

start

interrupt

ADC

Processing

disk

$\tau$ = 1 ms

➡ **Stochastic process**

➡ Fluctuations in time between events

➡ **Let's assume for example**

➡ physics rate f = 1 kHz, i.e. $\lambda = 1$ ms

➡ and, as before, $\tau = 1$ ms



Probability of time (in ms) between events for average decay rate of f=1kHz → $\lambda$=1ms

**Trigger path**

**Data path**

**TRIGGER**

**discriminator**

delay

start

ADC

interrupt

Processing
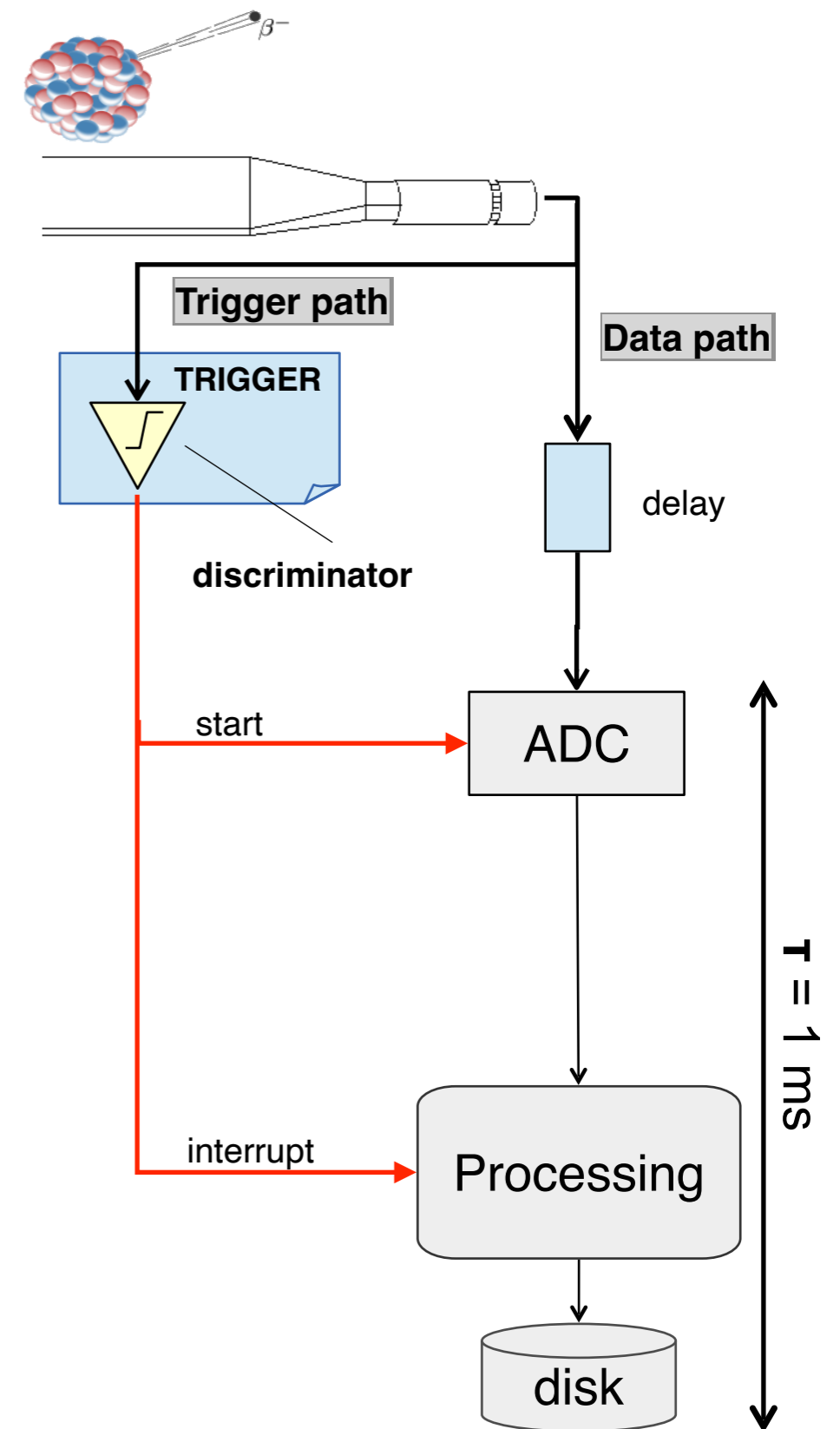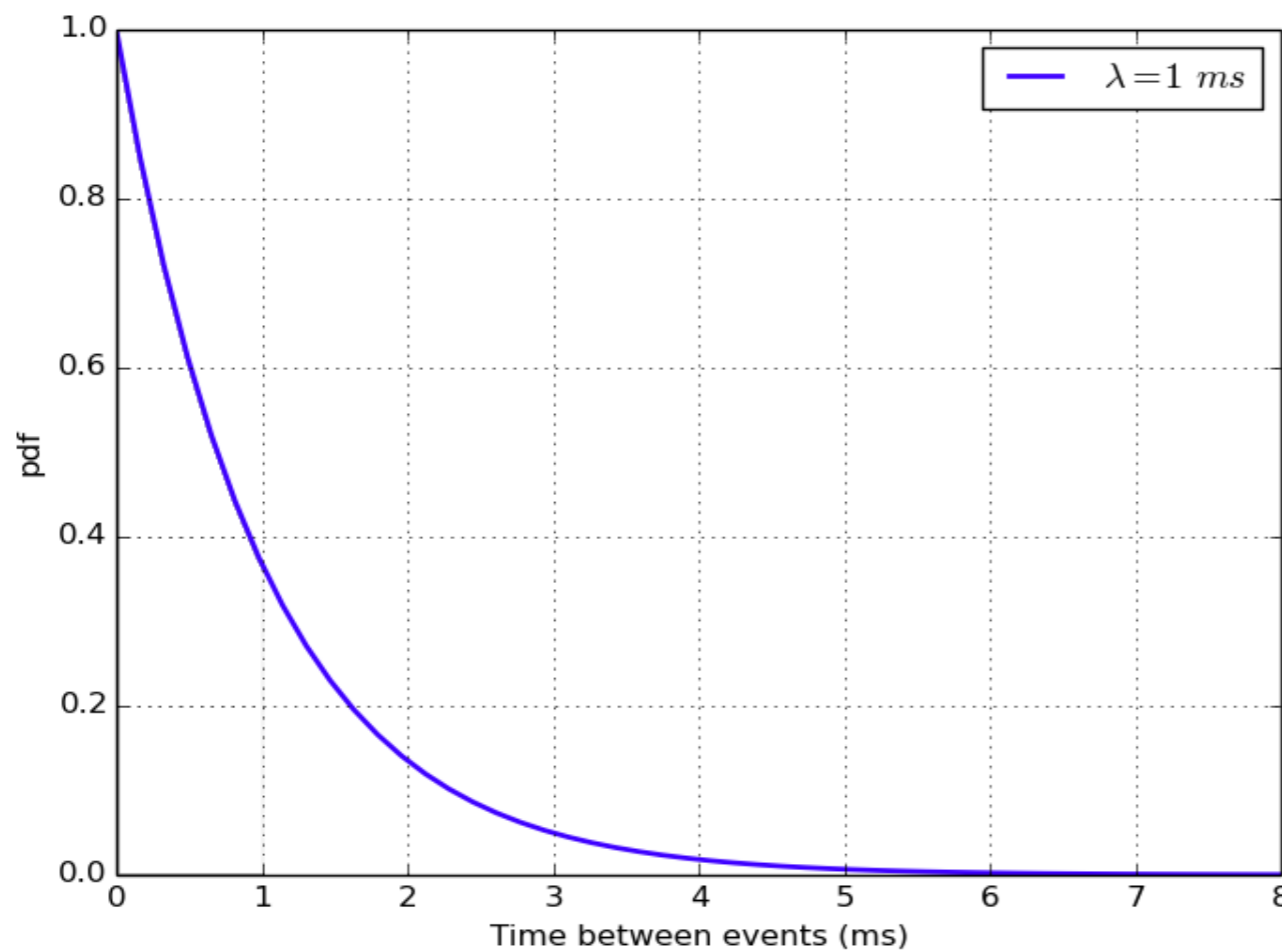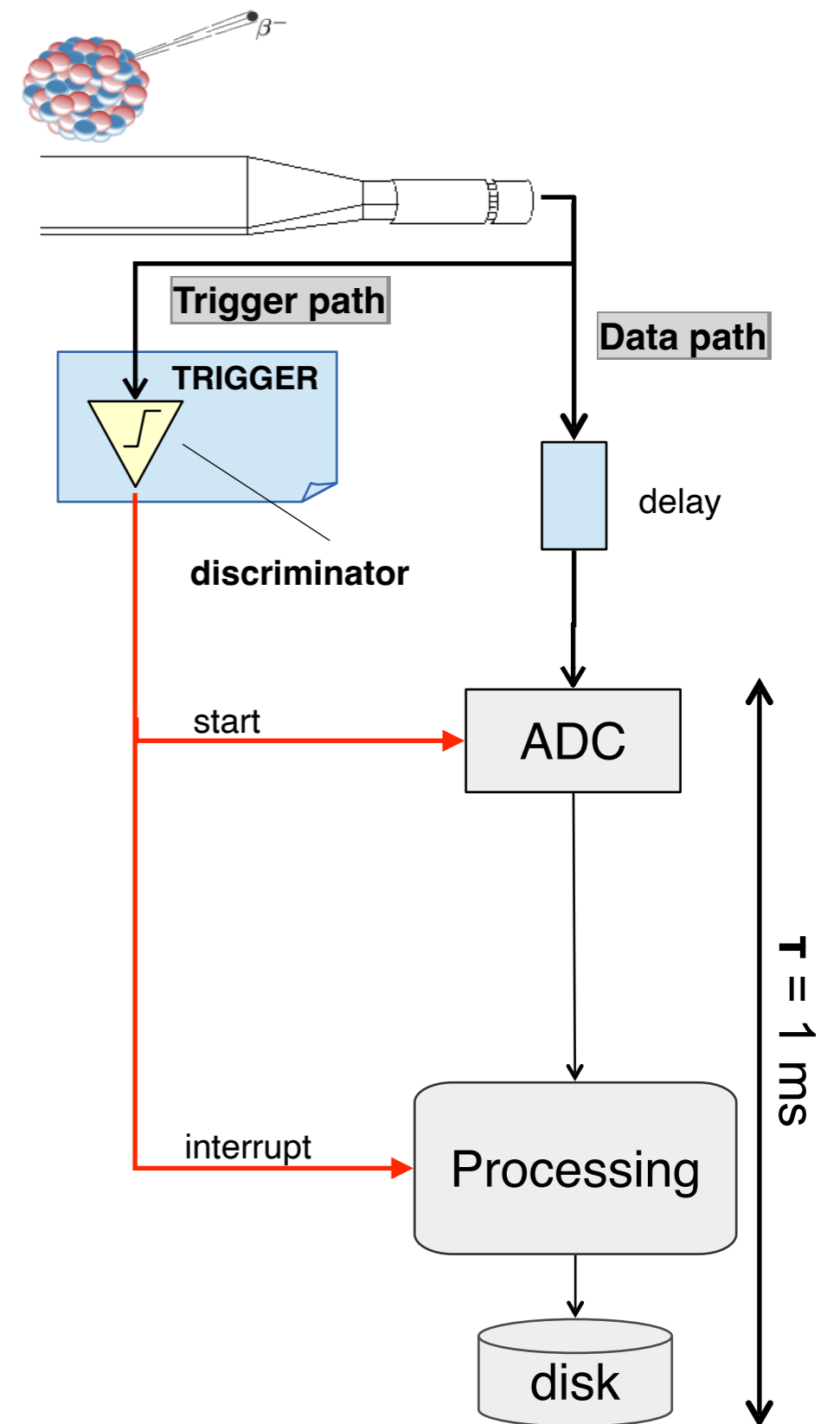
$\tau$ = 1 ms

disk

# BASIC DAQ: "REAL" TRIGGER

➡ **Stochastic process**
  ➡ Fluctuations in time between events
➡ **Let's assume for example**
  ➡ physics rate f = 1 kHz, i.e. $\lambda$ = 1 ms
  ➡ and, as before, $\tau$ = 1 ms



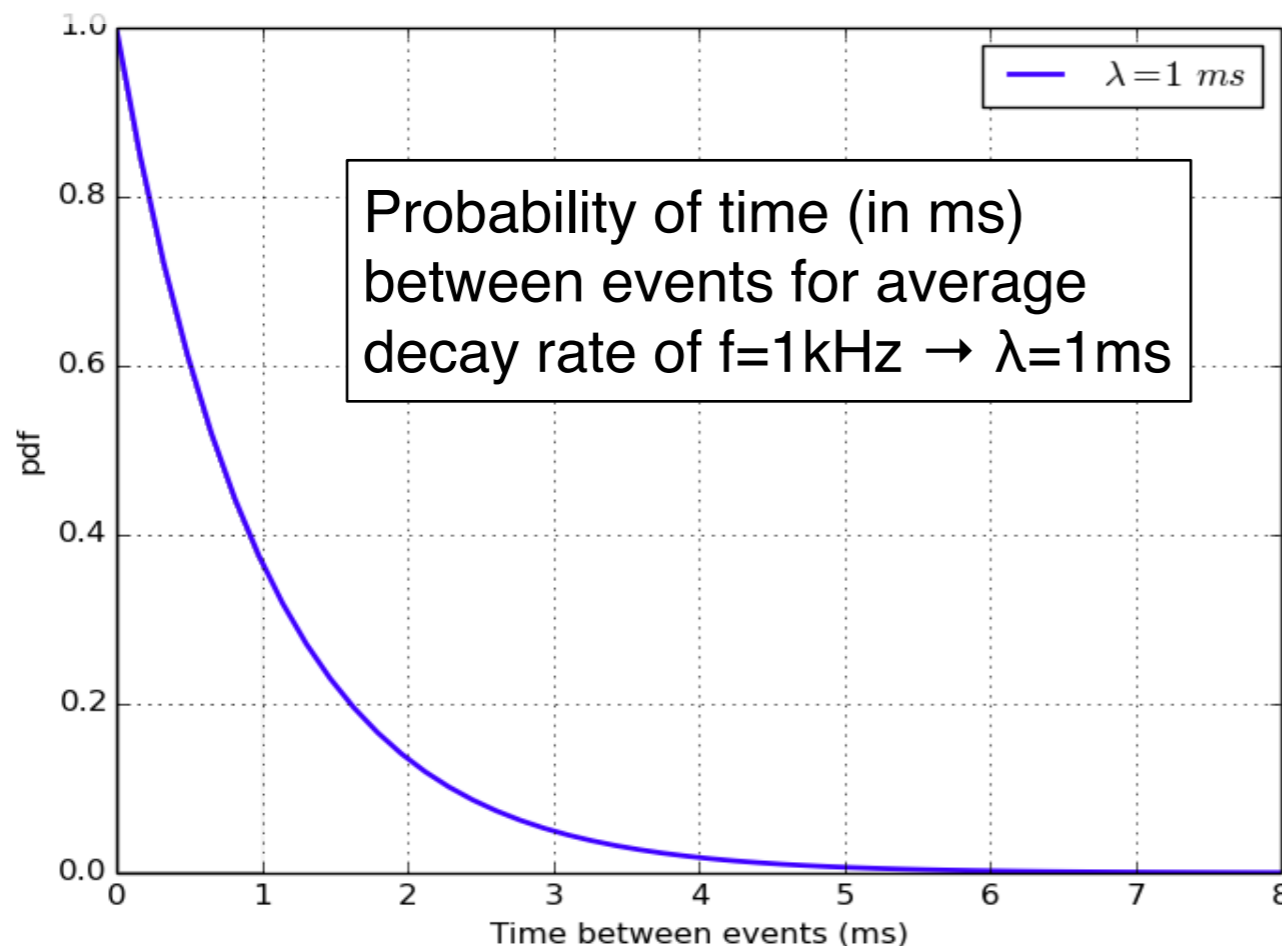What if a trigger is created when the system is busy?

➡ **If a new trigger arrives when the system is still processing the previous event**

➡ The processing of the previous event could be screwed up

What if a trigger is created when the system is busy?

➡ **Need a feedback mechanism, to know if the data processing pipeline is free to process a new event: the** **busy logic**



f = 1 kHz
λ= 1 ms

Trigger path

Data path

TRIGGER

discriminator

delay

start

ADC

interrupt

Processing

disk

τ = 1 ms

➡ **Need a feedback mechanism, to know if the data processing pipeline is free to process a new event: the busy logic**

➡ **A minimal busy logic can be implemented with**

  ➡ an **AND** gate

  ➡ a **NOT** gate

  ➡ a **flip-flop**



f = 1 kHz
λ = 1 ms

Trigger path

Data path

TRIGGER

discriminator

delay

start

ADC

interrupt

Processing

disk

τ = 1 ms

➡ **Need a feedback mechanism, to know if the data processing pipeline is free to process a new event: the busy logic**

➡ **A minimal busy logic can be implemented with**

  ➡ an **AND** gate

  ➡ a **NOT** gate

  ➡ a **flip-flop**

**Any new trigger is inhibited by the AND gate (busy)**

f = 1 kHz
λ = 1 ms

Trigger path

Data path

TRIGGER

delay

NOT

AND

start

ADC

BUSY LOGIC

flip-flop

Q   SET
    CLEAR

ready

Processing

T = 1 ms

disk

# DEADTIME AND EFFICIENCY

➡ **The busy mechanism protects our electronics from unwanted triggers**

  ➡ During the busy time, no signals are accepted, cause of **inefficiency**

  ➡ this is a source of **dead-time**

➡ **Due to stochastic fluctuations**

  ➡ DAQ rate always < physics rate

  ➡ Efficiency always < 100%

➡ **To cope with the input signal fluctuations, we have to over-design our DAQ system**

  ➡ can we mitigate this effect?

*f* average rate of physics (input)

$\nu$ average rate of DAQ (output)

**τ: deadtime**

$$\nu = \frac{f}{1 + f\tau}$$

➡ **The busy mechanism protects our electronics from unwanted triggers**

   ➡ During the busy time, no signals are accepted, cause of **inefficiency**

   ➡ this is a source of **dead-time**

➡ **Due to stochastic fluctuations**

   ➡ DAQ rate always < physics rate

   ➡ Efficiency always < 100%

➡ **To cope with the input signal fluctuations, we have to over-design our DAQ system**

   ➡ can we mitigate this effect?



*f* average rate of physics (input)

$\nu$ average rate of DAQ (output)

**τ: deadtime**

$$\nu = \frac{f}{1 + f\tau}$$

➡ **The busy mechanism protects our electronics from unwanted triggers**

   ➡ During the busy time, no signals are accepted, cause of **inefficiency**
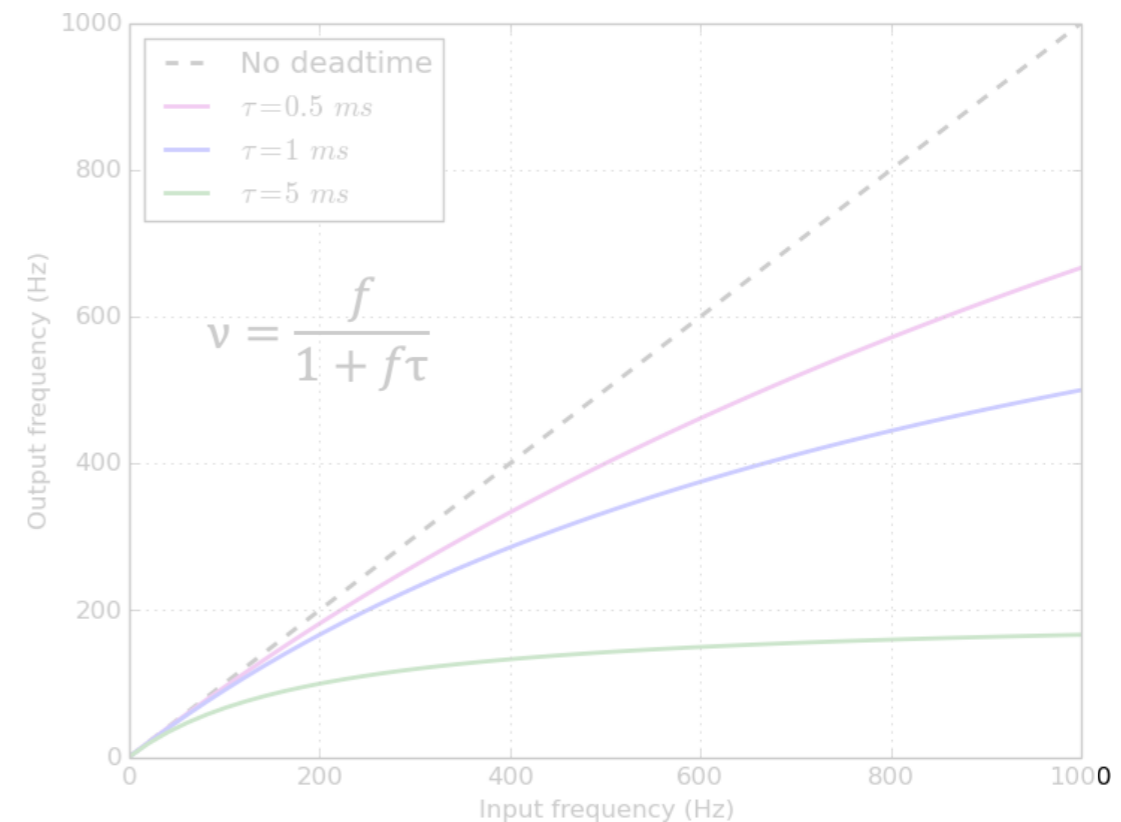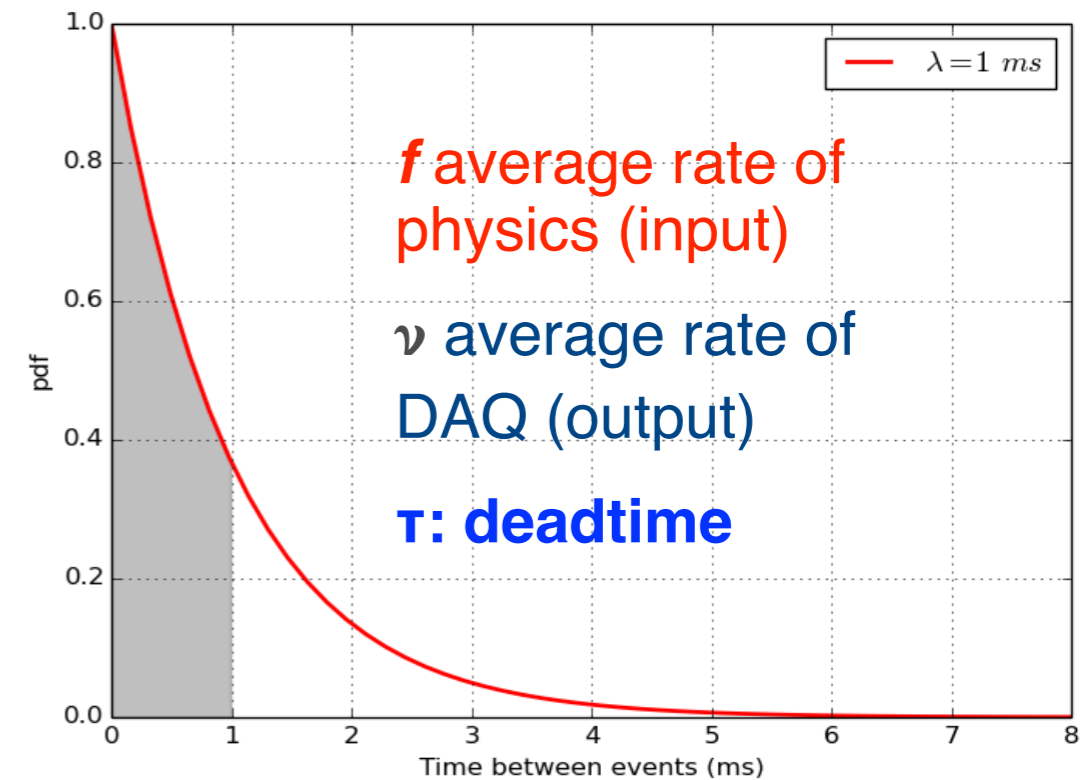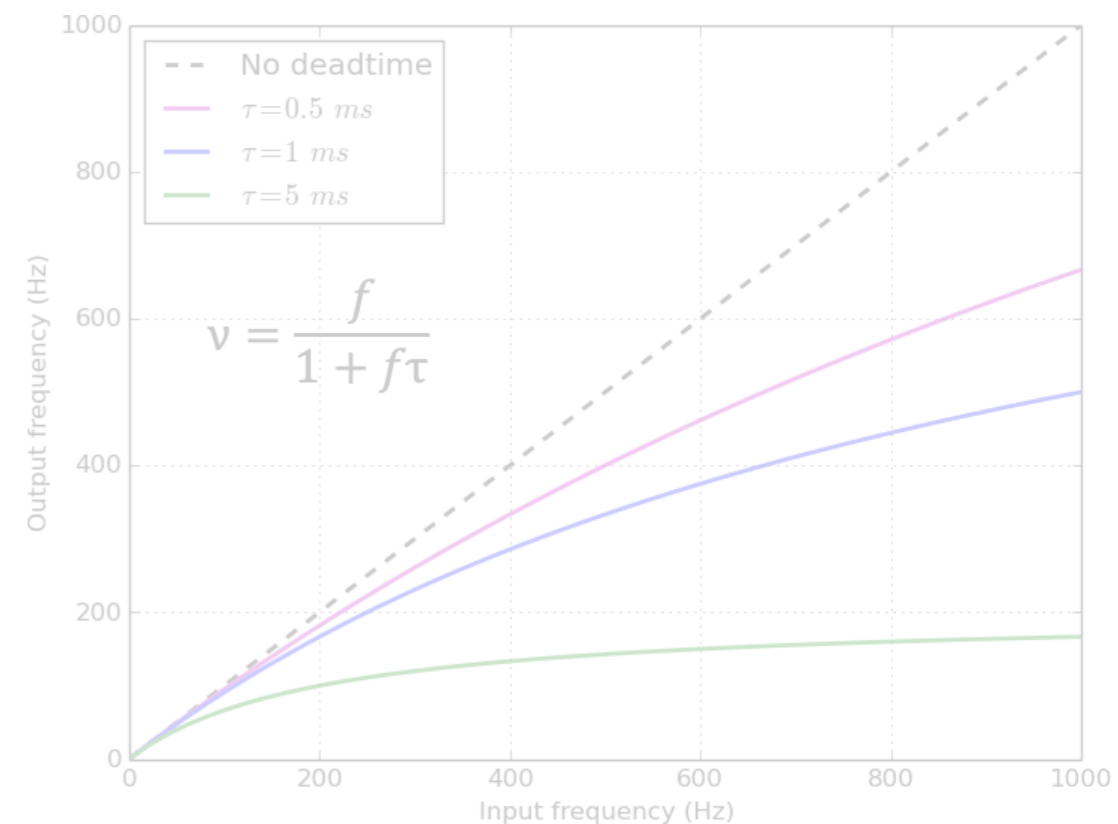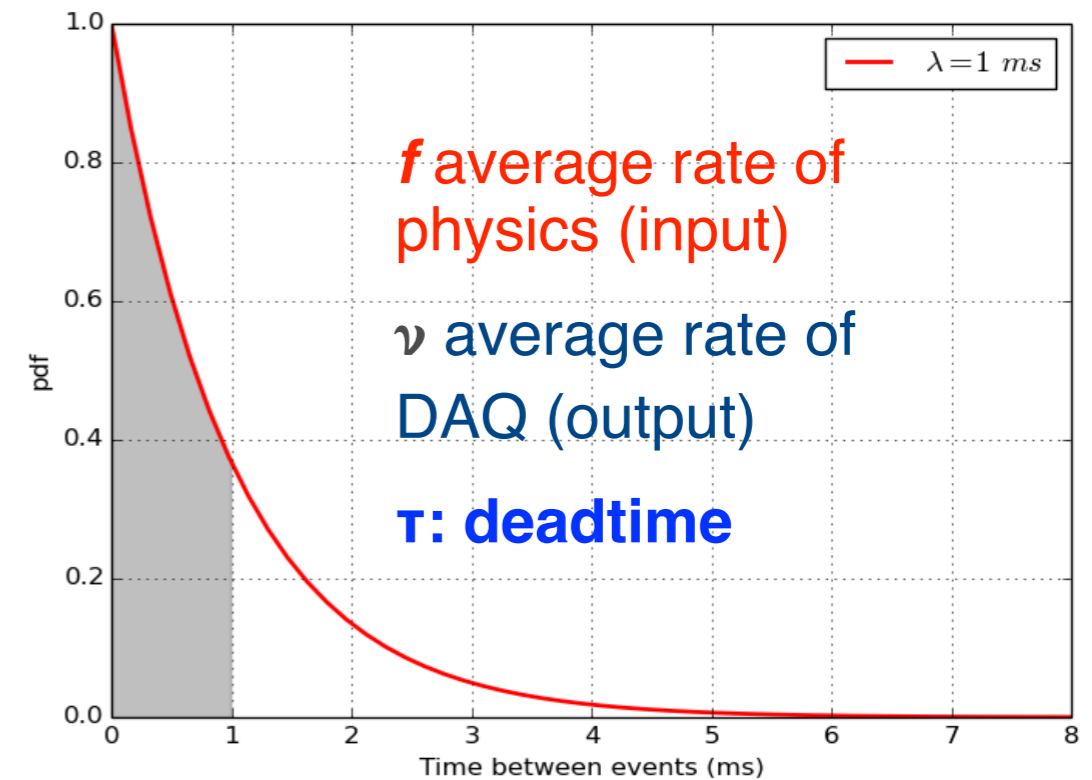
   ➡ this is a source of **dead-time**

➡ **Due to stochastic fluctuations**

   ➡ DAQ rate always < physics rate

   ➡ Efficiency always < 100%

➡ **To cope with the input signal fluctuations, we have to over-design our DAQ system**

   ➡ **can we mitigate this effect?**

$f$ average rate of physics (input)

$\nu$ average rate of DAQ (output)

**τ: deadtime**

$$\nu = \frac{f}{1 + f\tau}$$

➡ **What if we were able to make the system more deterministic and less dependent on the arrival time of our signals?**

  ➡ Then we could ensure that events don't arrive when the system is busy

  ➡ This is called **de-randomization**

➡ How can be achieved?

  ➡ by **buffering** the data (having a holding queue where we can slot it up to be processed)

  ➡ Maintaining $\tau \sim \lambda$ (traffic intensity), high efficiency can be obtained even with moderate depth of FIFOs

**Inter-arrival time distribution**

**ms**

**Data access time distribution**

**ms**

# DE-RANDOMIZATION

➡ **What if we were able to make the system more deterministic and less dependent on the arrival time of our signals?**

➡ Then we could ensure that events don't arrive when the system is busy

➡ This is called **de-randomization**

➡ **How can be achieved?**

➡ by **buffering** the data (having a holding queue where we can slot it up to be processed)

➡ Maintaining $\tau \sim \lambda$ (traffic intensity), high efficiency can be obtained even with moderate depth of FIFOs

**Inter-arrival time distribution**

ms

**Data access time distribution**
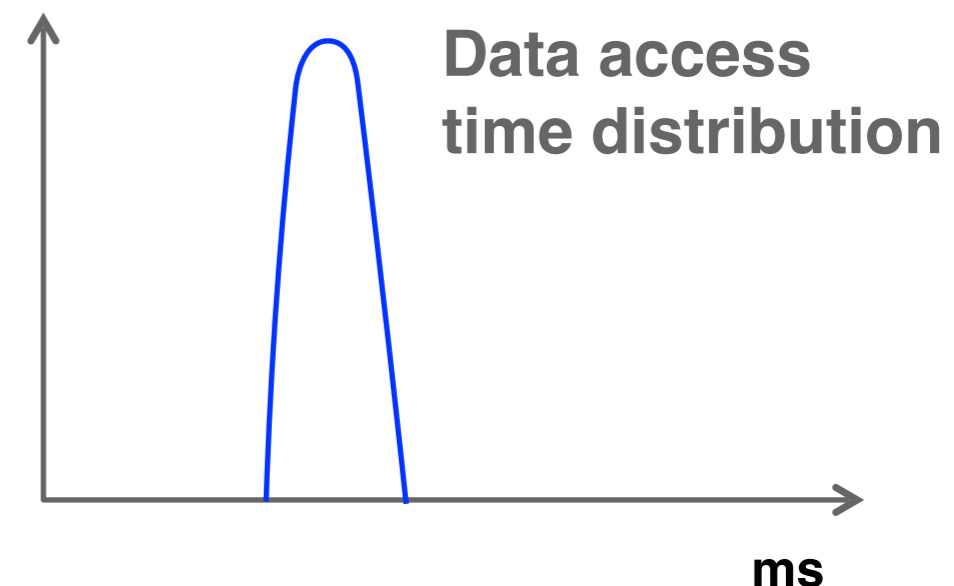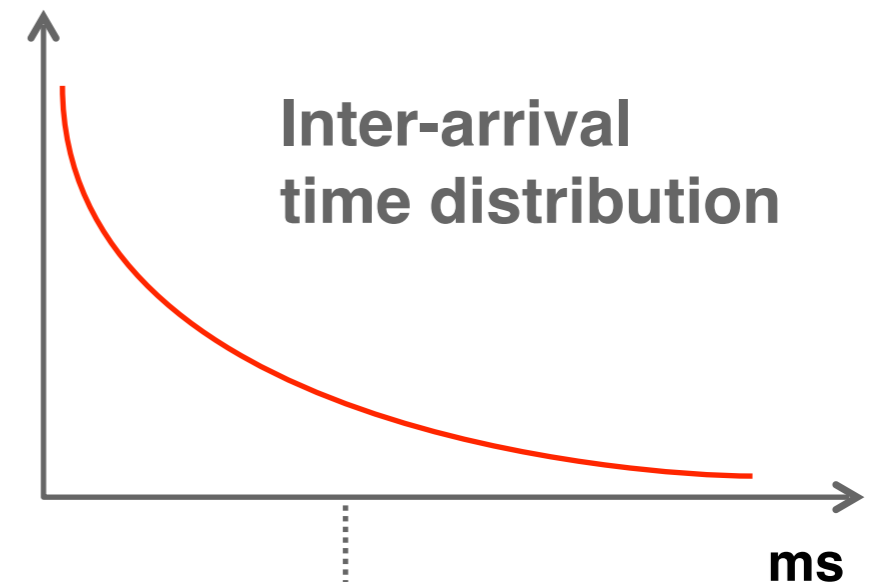
ms

➡ **What if we were able to make the system more deterministic and less dependent on the arrival time of our signals?**

    ➡ Then we could ensure that events don't arrive when the system is busy

    ➡ This is called **de-randomization**

➡ **How can be achieved?**

    ➡ by **buffering** the data

    ➡ having a holding **queue** where we can slot it up to be processed

    ➡ Maintaining $\tau \sim \lambda$ **(traffic intensity)**, high efficiency can be obtained even with moderate depth of FIFOs

Inter-arrival time distribution

ms

λ (ms)  f (Hz)

FIFO

τ (ms)  ν (Hz)

Data access time distribution

ms

➡ **What if we were able to make the system more deterministic and less dependent on the arrival time of our signals?**

    ➡ Then we could ensure that events don't arrive when the system is busy

    ➡ This is called **de-randomization**

➡ **How can be achieved?**

    ➡ by **buffering** the data

    ➡ having a holding **queue** where we can slot it up to be processed

    ➡ Maintaining $\tau \sim \lambda$ **(traffic intensity~1)**, high efficiency can be obtained even with moderate depth of FIFOs
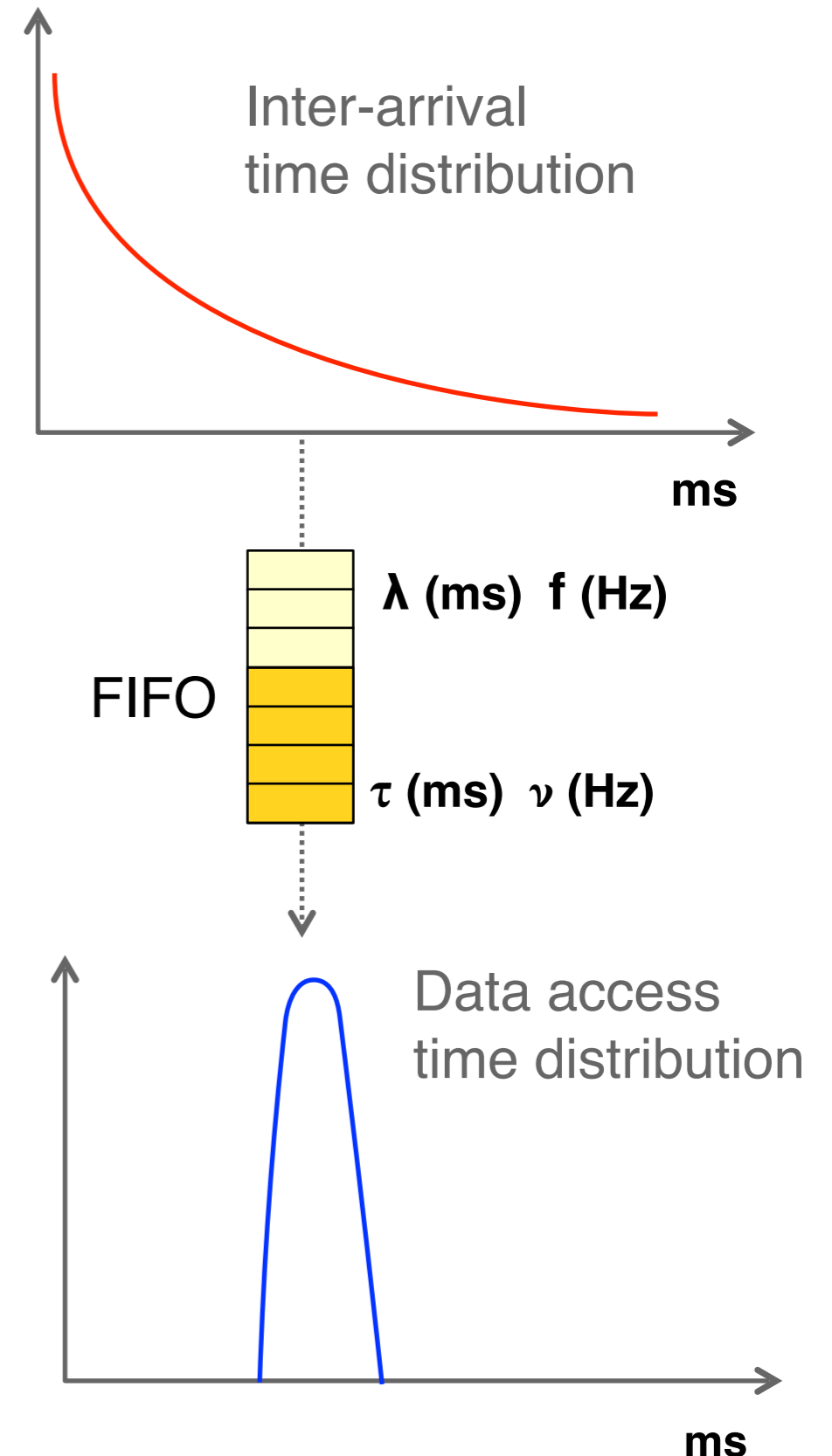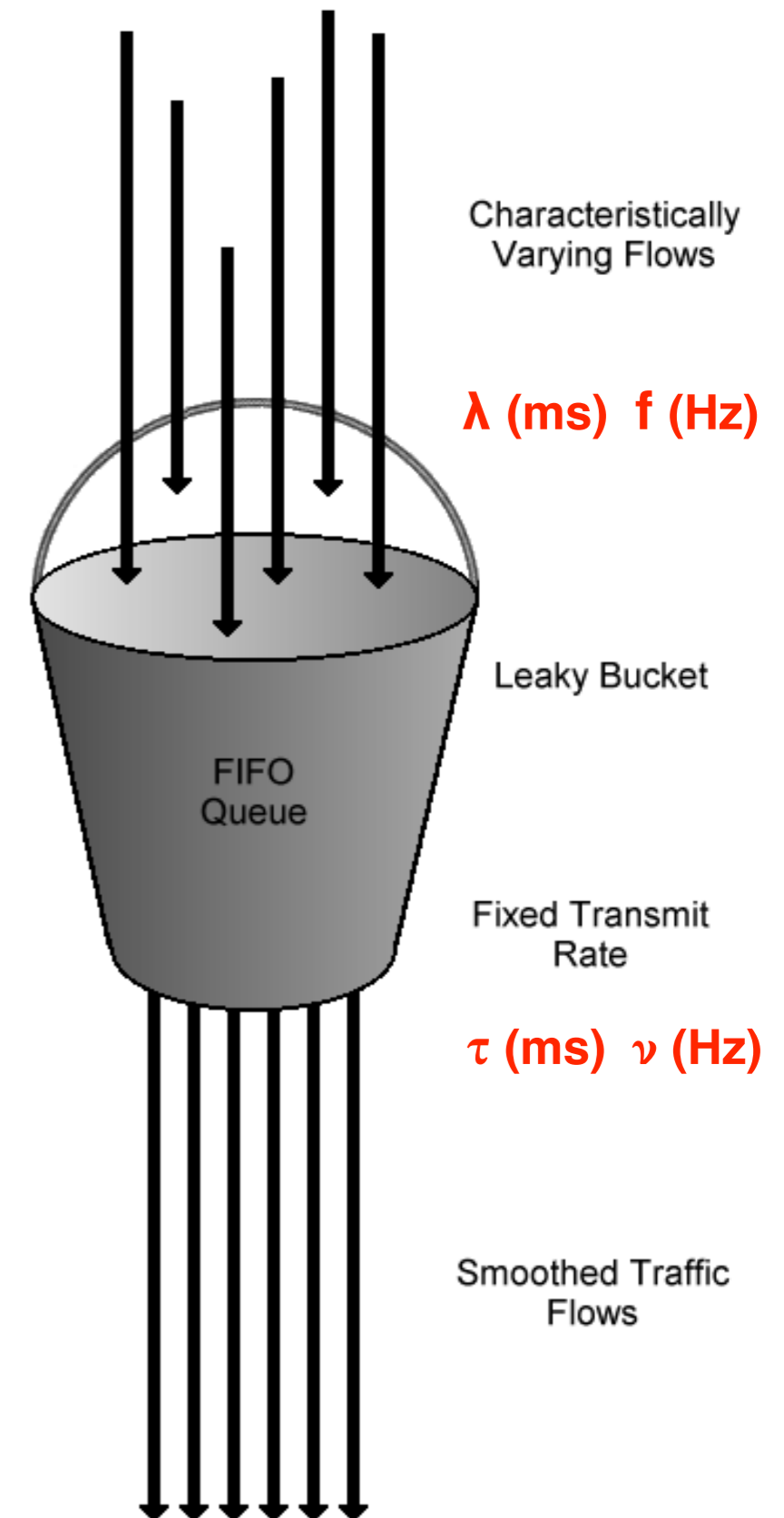
*—> search for "Queuing theory"*



Characteristically Varying Flows

$\lambda$ **(ms)  f (Hz)**

Leaky Bucket

FIFO Queue

Fixed Transmit Rate

$\tau$ **(ms)  $\nu$ (Hz)**

Smoothed Traffic Flows

➡ **Input fluctuations can be absorbed and smoothed by a queue**

    ➡ A FIFO can provide a ~steady and **de-randomized** output rate

➡ **Busy is now defined by the buffer occupancy**

    ➡ Processor pulls data from the buffer at fixed rate, separating the event receiving and data processing steps



$f = 1\ kHz$
$\lambda = 1\ ms$

Trigger path

Data path

TRIGGER

delay

BUSY LOGIC

AND

start

ADC

FIFO

busy (full)

data ready

Processing

disk

$T = 1\ ms$

➡ **The FIFO decouples the low latency front-end from the data processing**

- ➡ Minimize the amount of "unnecessary" fast components

➡ **~100% efficiency with minimal deadtime achievable if**

- ➡ ADC can operate at rate >> f
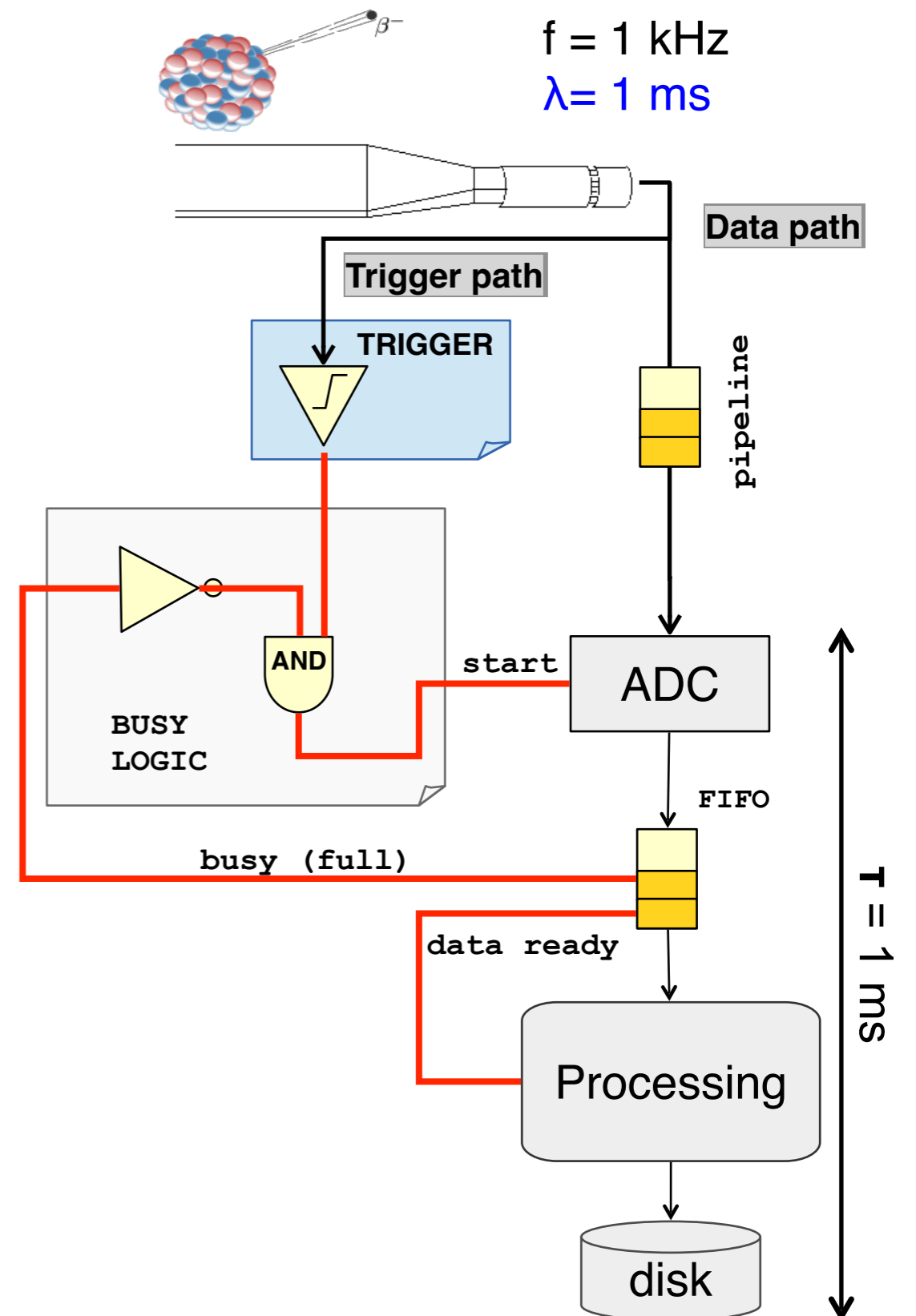- ➡ Data processing and storing operate at a rate ~ f

➡ **Could the delay be replaced with a "FIFO"?**

- ➡ Analog pipelines, heavily used in LHC DAQs



$f = 1$ kHz
$\lambda = 1$ ms

Trigger path

Data path

TRIGGER

delay

AND

BUSY LOGIC

start

ADC

FIFO

busy (full)

data ready

Processing

$\tau = 1$ ms

disk

➡ **The FIFO decouples the low latency front-end from the data processing**

➡ Minimize the amount of "unnecessary" fast components

➡ **~100% efficiency with minimal deadtime achievable if**

➡ ADC can operate at rate >> f

➡ Data processing and storing operate at a rate ~ f

➡ **Could the delay be replaced with a "FIFO"?**
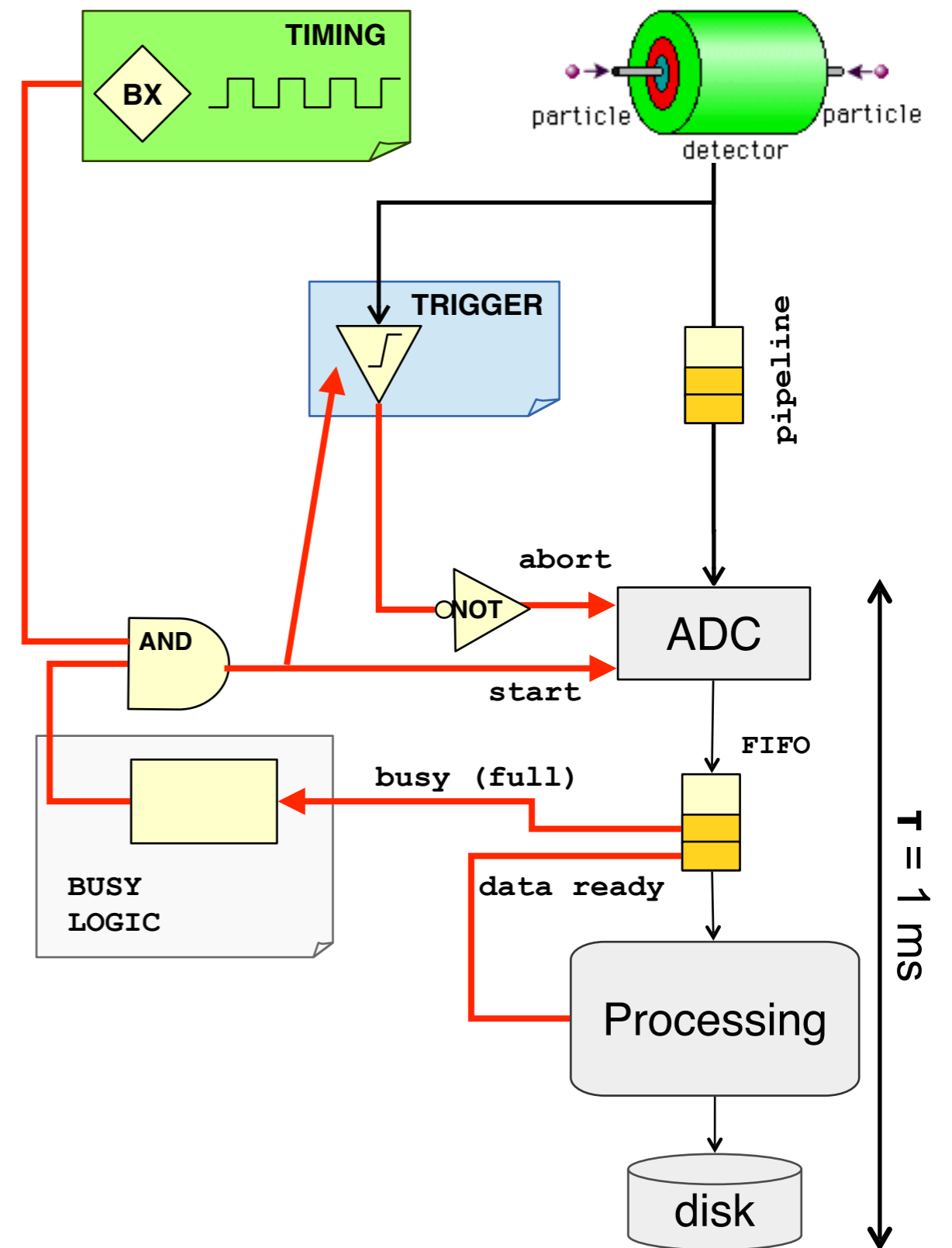
➡ Analog pipelines, heavily used in LHC DAQs

➡ **Do we need de-randomization buffers also in collider setups?**

- ➡ Particle collisions are synchronous
- ➡ But the time distribution of triggers is random: good events are unpredictable



**TIMING**

BX

**TRIGGER**

pipeline

abort

NOT

ADC

AND

start

BUSY LOGIC

busy (full)

FIFO

data ready

Processing

disk

T = 1 ms

particle — detector — particle

➡ **Do we need de-randomization buffers also in collider setups?**

   ➡ Particle collisions are synchronous

   ➡ But the time distribution of triggers is random: good events are unpredictable
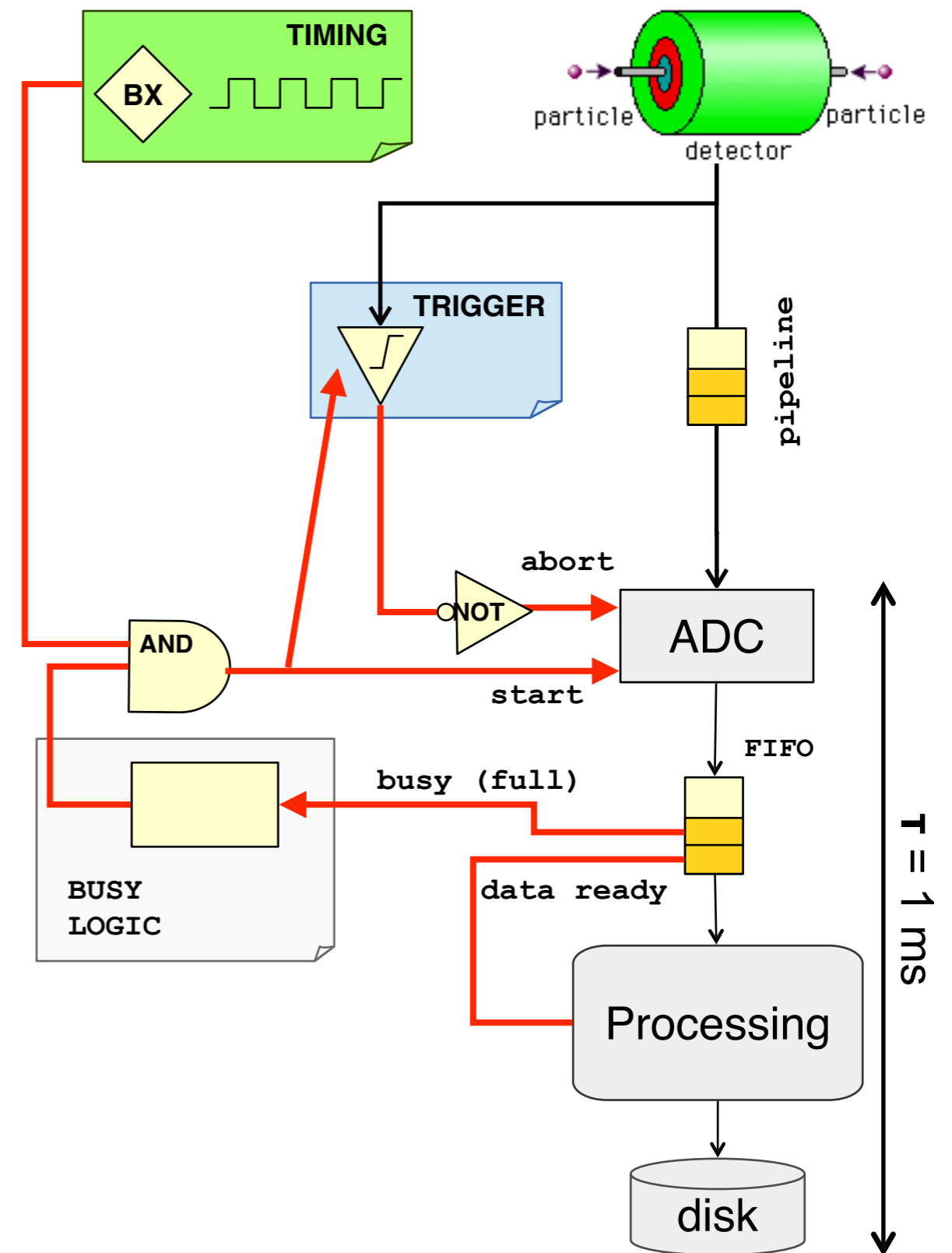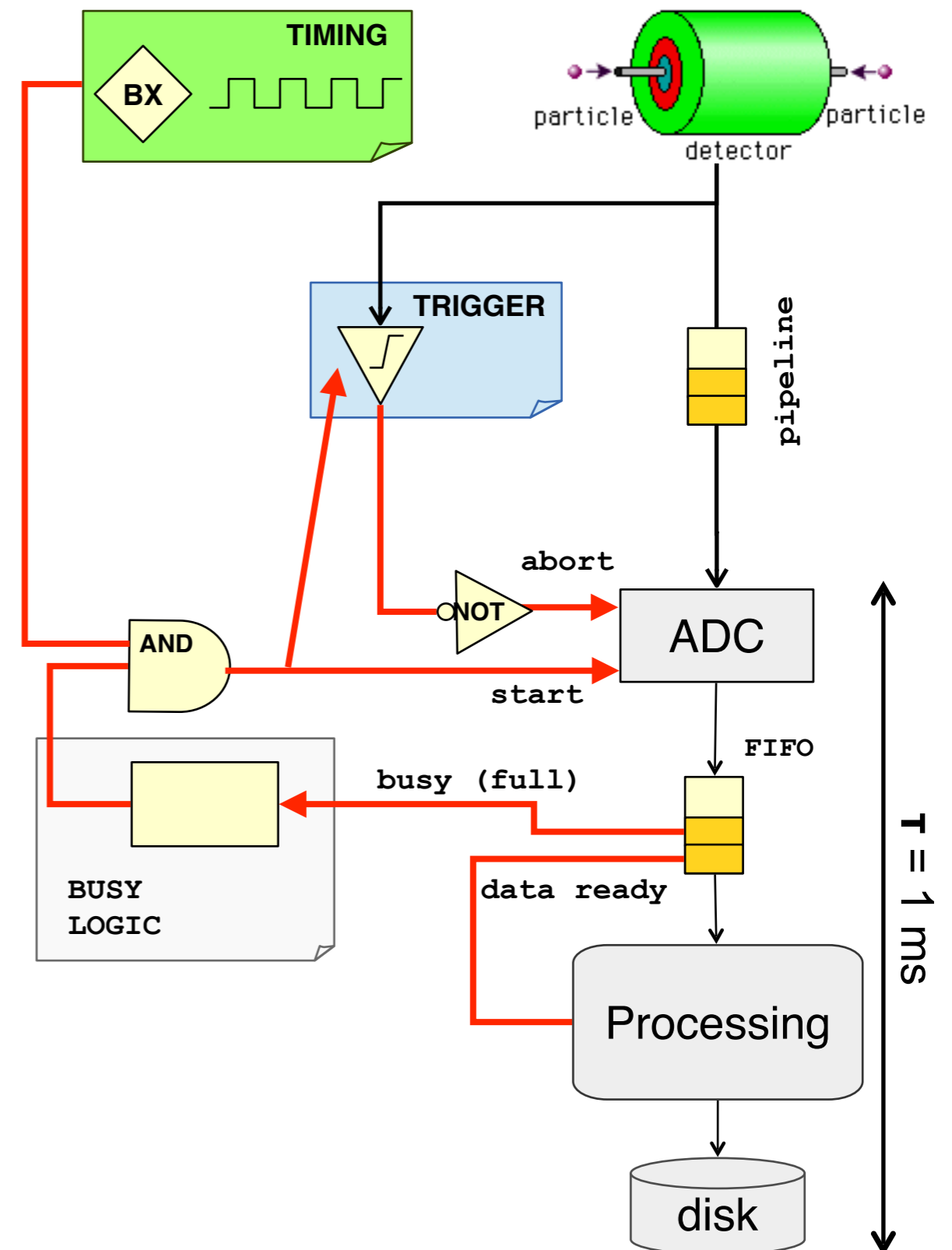
➡ **De-randomization is still needed**

➡ **Do we need de-randomization buffers also in collider setups?**

➡ Particle collisions are synchronous

➡ But the time distribution of triggers is random: good events are unpredictable

➡ **De-randomization is still needed**

➡ **More complex busy logic to protect buffers and detectors**

➡ Eg: accept n events every m bunch crossings

➡ Eg: prevent some dangerous trigger patterns

➡ **Introduction**

  ➡ What is Trigger and DAQ?

  ➡ Overall TDAQ framework

➡ **Basic TDAQ concepts**

  ➡ Digitization, Latency

  ➡ Deadtime, Busy

  ➡ De-randomization

➡ **Scaling up**

  ➡ Readout and Event Building
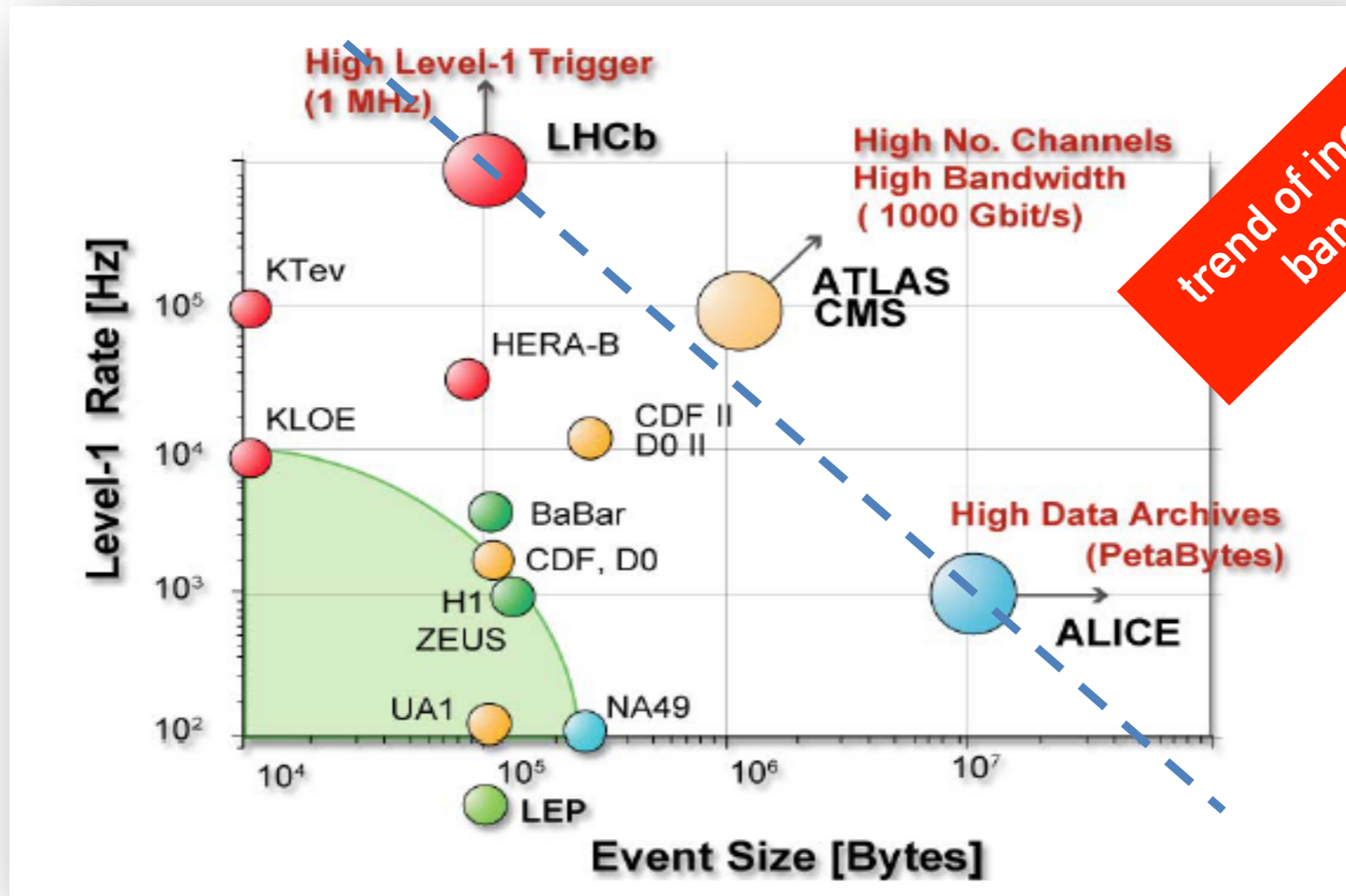
  ➡ Buses vs Network
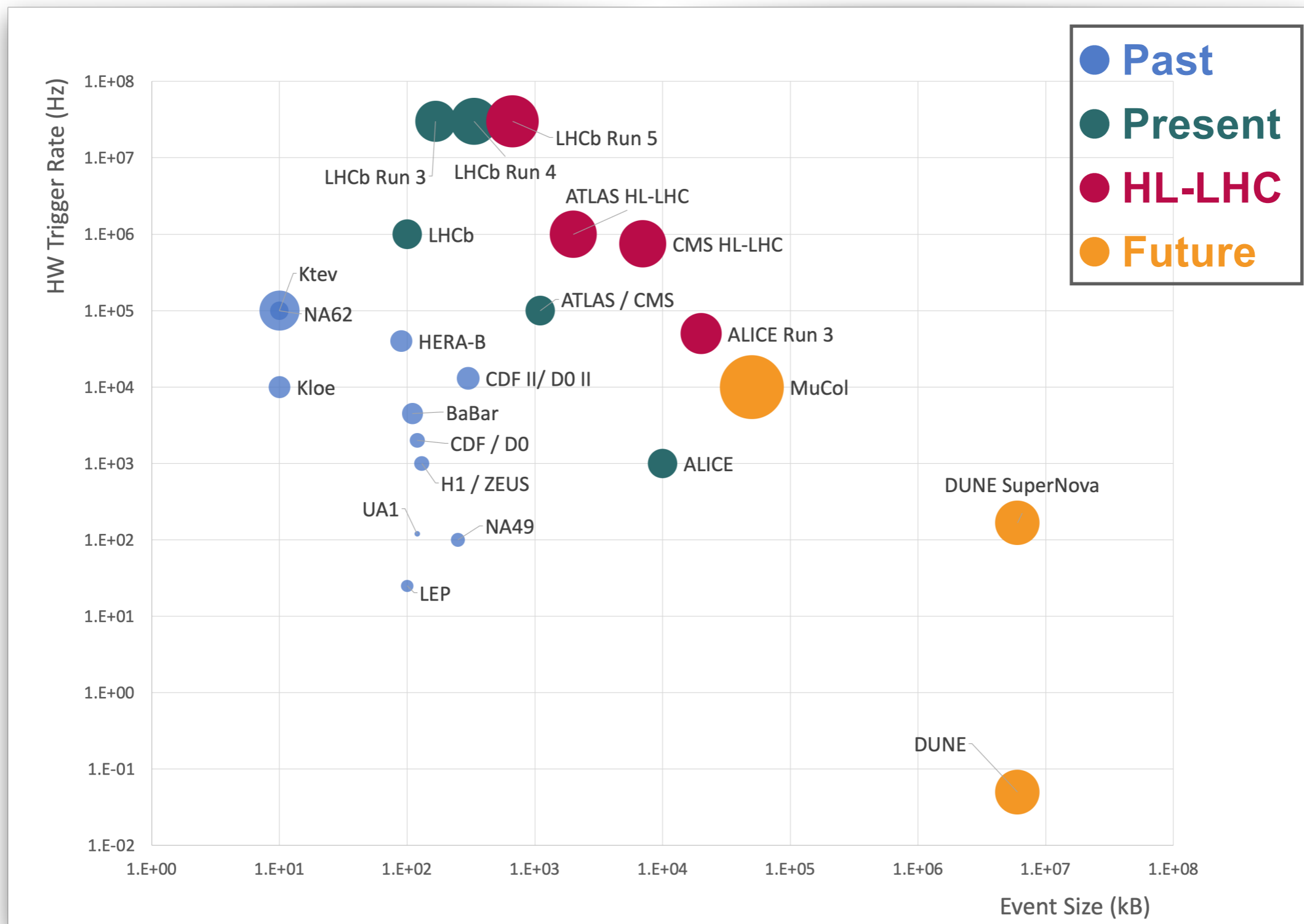
➡ **Fight bottlenecks**

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



*faster L1 electronics*

*more channels, more complex events*

➡ **As the data volumes and rates increase, new architectures are needed**

*Courtesy of A.Cerri*

- **Adding more channels requires a hierarchical structure committed to the data handling and conveyance**

- **Adding more channels requires a hierarchical structure committed to the data handling and conveyance**

- **Adding more channels requires a hierarchical structure committed to the data handling and conveyance**

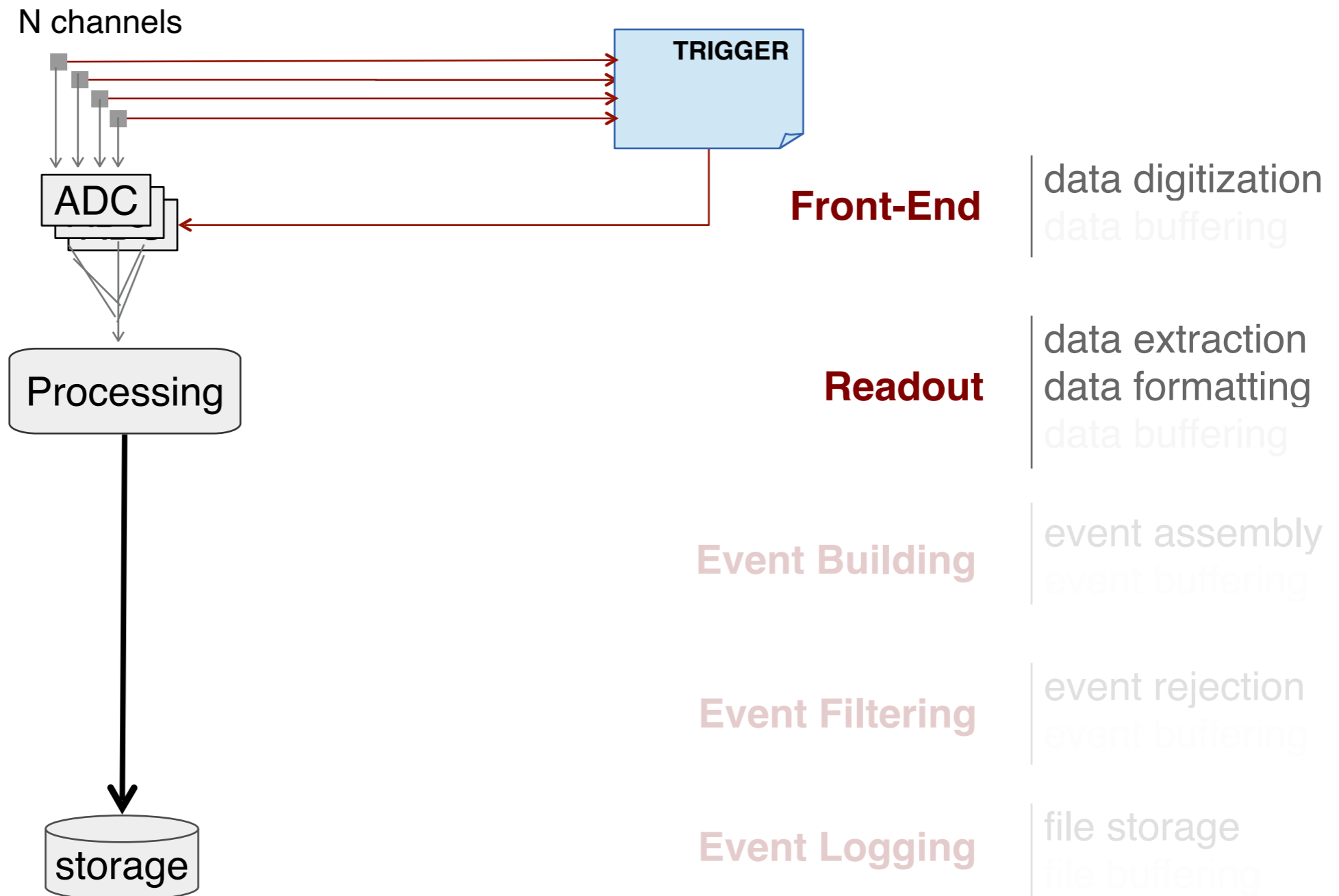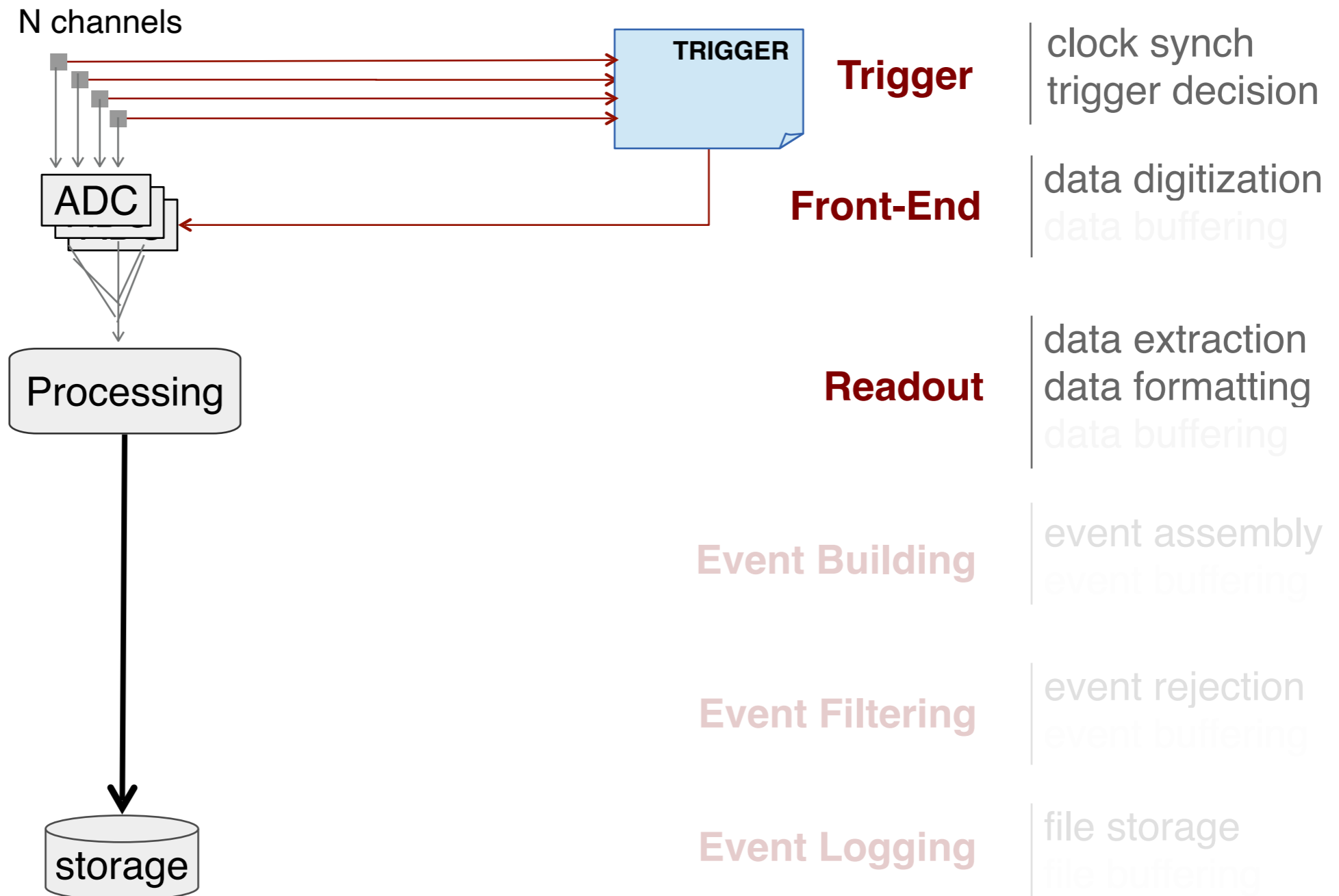- **Adding more channels requires a hierarchical structure committed to the data handling and conveyance**

- **Adding more channels requires a hierarchical structure committed to the data handling and conveyance**

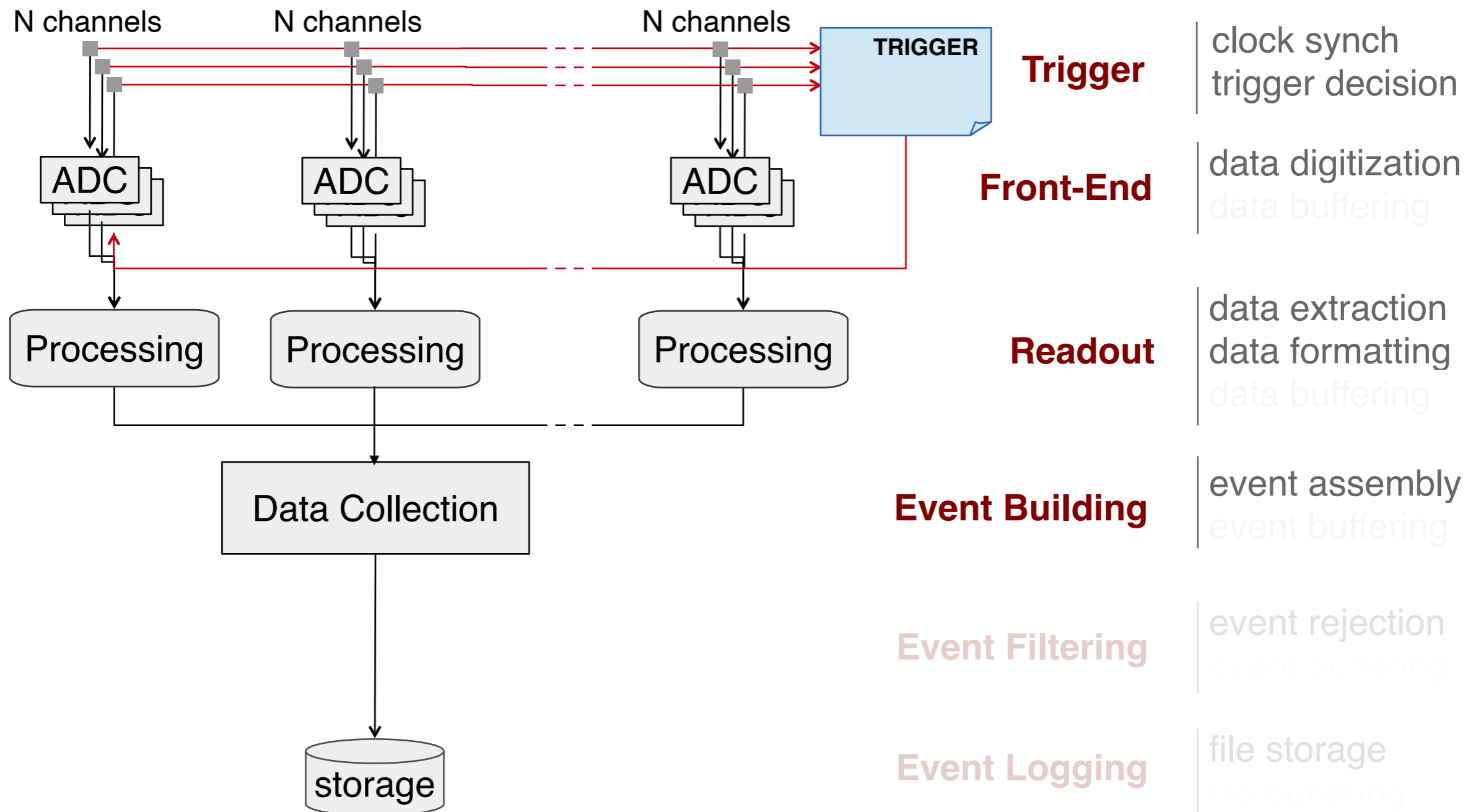- **Adding more channels requires a hierarchical structure committed to the data handling and conveyance**

- **Buffering** usually needed at every level
- DAQ can be seen as a multi level buffering system

➡ **Reading out data or building events out of many channels requires many components**



In the design of our hierarchical data-collection system, we have better define "**building blocks**"

➤ Readout crates

➤ HLT racks

➤ Event Building groups

➤ DAQ slices

➜ **How to organize the interconnections inside the building blocks and between building blocks?**

   ➜ How to connect data sources and data destinations?

   ➜ Two main classes: **bus** or **network**



data sources

bus

bus

bus

network

data processors

➡ **How to organize the interconnections inside the building blocks and between building blocks?**
  ➡ How to connect data sources and data destinations?
  ➡ Two main classes: **bus** or **network**

# BUSES

➡ **Devices connected via a shared bus**

    ➡ Bus → group of electrical lines

➡ **Sharing implies arbitration**

    ➡ Devices can be **master** or **slave**

    ➡ Devices can be addresses (uniquely identified) on the bus

➡ **E.g.: SCSI, Parallel ATA, VME, PCI …**

    ➡ local, external, crate, long distance, ...

➡ **Simple** :-)

   ➡ Fixed number of lines (bus-width)

   ➡ Devices have to follow well defined interfaces

      ➡ Mechanical, electrical, communication, ...

➡ **Scalability issues** :-(

   ➡ Bus bandwidth is shared among all the devices

   ➡ Maximum bus width is limited

   ➡ Maximum number of devices depends on bus length

   ➡ Maximum bus frequency is inversely proportional to the bus length

   ➡ **On the long term, other "effects" might limit the scalability of your system**

On the long term, other "effects" might limit the scalability of your system

➡ **VME Modular electronics**

➡ **VME bus programming**

➡ **µATCA**

➡ **PCI express**

# NETWORK

➡ **All devices are equal**

  ➡ They communicate directly with each other via messages

  ➡ No arbitration, simultaneous communications

➡ **Eg: Telephone, Ethernet, Infiniband, …**

➡ **In switched networks, <span style="color:blue">switches</span> move messages between sources and destinations**

   ➡ Find the right path

➡ **How <span style="color:darkred">congestions</span> (two messages with the same destination at the same time) are handled?**

   ➡ The key is .... buffering

➡ **In switched networks, switches move messages between sources and destinations**

➡ Find the right path

➡ **How congestions (two messages with the same destination at the same time) are handled?**

➡ The key is .... **buffering**

➡ **Networks scale well (and allow redundancy)**
   ➡ They are the backbones of the LHC DAQ systems

➡ **Very Front End**

➡ does the analog part

➡ ADC, low-level calibration, zero suppression, lossless compression

➡ low-power, rad-tolerant

➡ **Quasi Front End**

➡ medium scale aggregation, local reconstruction, "lossy" compression, transition to standard protocol on optical links

➡ **Commodity Of The Shelf**

➡ COTS switched networks

➡ COTS servers, with co-processors (GPU, FPGA)



*E. Meschi, Summer Student Lectures 2022*

➡ **Artificial deadtime**

➡ **Data collection**

➡ **Multi-level trigger**

➡ **Data-flow control**

➡ **Data reconstruction**

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 parameters: bucket size (in number of events) / readout time (in BC units)

      ➡ i.e. 9 / 351 for LAr readout

**Complex Deadtime Fraction [%]**

Complex Deadtime Fraction [%]

7/351
8/351
9/351

L1 Output Rate [kHz]

**more sensors ⇒ more granularity**

**multiple digitisers ⇒ more 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 to storage
  ➡ final storage:
    ➡ VME up to 50MB/s → 1TB in 6h
    ➡ too many disks in one week!
➡ **Data Collection unit decouples storage from processing**
  ➡ dedicated to format, compress and store

N channels

TRIGGER

**more sensors ⇒ more granularity**

**multiple digitisers ⇒ more parallelism**

ADC

Processing

Data Collection

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 to storage
- ➡ final storage:
  - ➡ VME up to 50MB/s → 1TB in 6h
  - ➡ too many disks in one week!

➡ **Data Collection unit decouples storage from processing**
- ➡ dedicated to format, compress and store

➡ **Reduce the rate at each stage, with limited buffer size and no deadtime**
  - ➡ $\tau \sim \lambda$ (traffic intensity~1)

➡ **High level triggers with longer latency**
  - ➡ more complex filters
  - ➡ more data (for example silicon detectors)



Diagram labels: detectors, digitizers, LV1 ms, front-end pipelines, LV2 ms, readout buffers, switching networks, LV3 s, processor farms

---

**Recall on trigger architectures**

➡ **Real time system**
  - ➡ must respond within some **fixed latency**
  - ➡ Latency = max Latency
  - ➡ over fluctuations is bad, will create deadtime

➡ **Non-real-time system**
  - ➡ responds as soon as it's available
  - ➡ Latency = **Mean Latency**
  - ➡ over fluctuations is fine, shouldn't create deadtime

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

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

   ➡ Can overflow, with bursts and unusual event sizes

   ➡ In these cases, can

      ➡ discard data locally or

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

detectors

digitizers

LV1

ms

front-end pipelines

LV2

ms

readout buffers

switching networks

LV3

s

processor farms

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

  ➡ Can overflow, with bursts and unusual event sizes

  ➡ In these cases, can

    ➡ discard data locally or

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

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

  ➡ using knowledge of the input buffer state

detectors

digitizers

LV1
ms

front-end pipelines

LV2
ms

readout buffers

switching networks

LV3
s

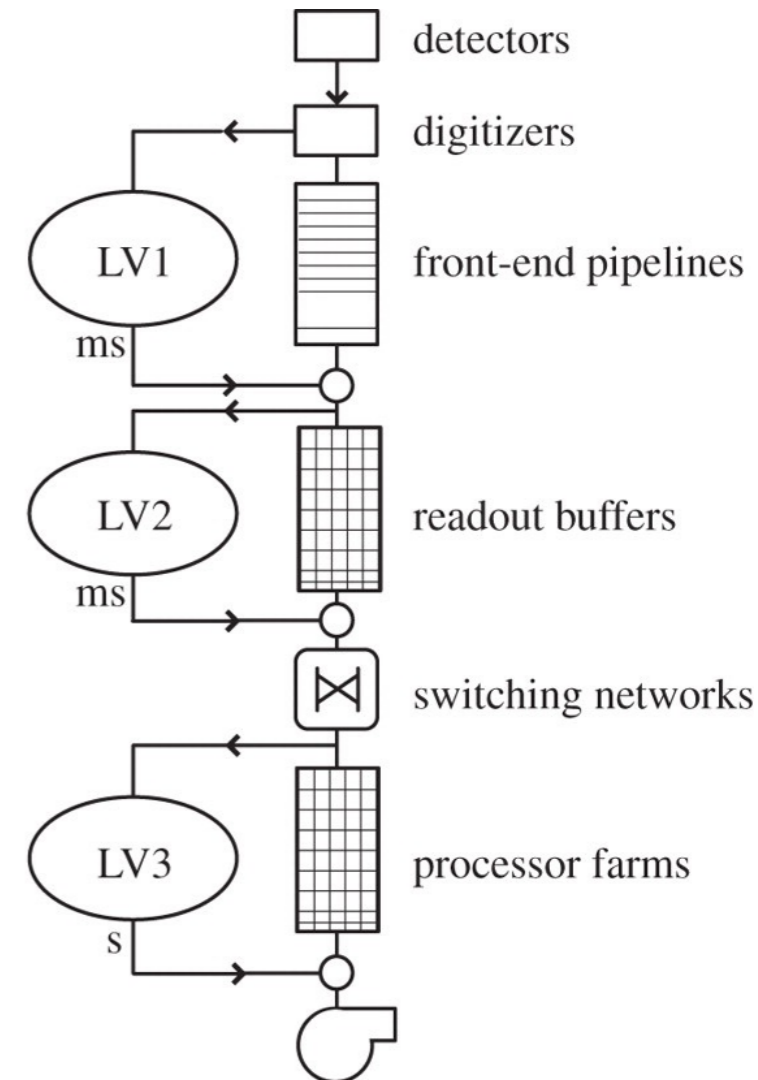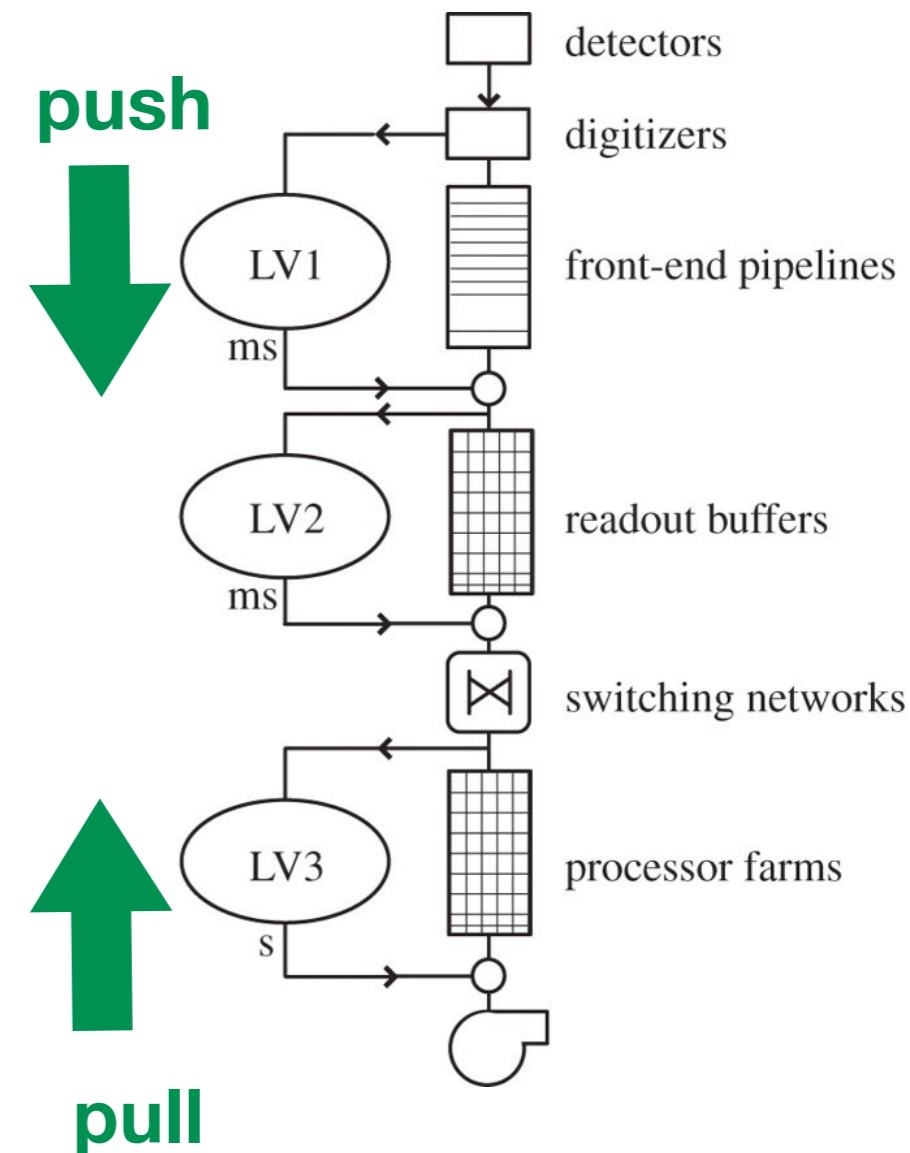processor farms

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

   ➡ Can overflow, with bursts and unusual event sizes

   ➡ In these cases, can

      ➡ discard data locally or

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

   ➡ **FE Push**: Events are sent as soon as data are available to the sender (e.g. round-robin algorithm) $\implies$ Busy or Throttle (block trigger)

   ➡ **EB Pull** : events are required by a given destination processes (may need an event manager) $\implies$ back-pressure (block dataflow)

   ➡ **Push-Pull** $\implies$ busy and back-pressure

**push**

**pull**

detectors

digitizers

LV1

ms

front-end pipelines

LV2

ms

readout buffers

switching networks

LV3

s

processor farms

➡ **Can play with data size and delayed reconstruction to overcome limitations**

➡ **Trigger Level Analysis / data scouting: data compressed via full event reconstruction, avoid to save raw detector data**

  ➡ if the <u>bandwidth</u> to write to the permanent storage is limited

➡ **Data parking: data saved on temporary storage and reconstructed when resources are available (during fills,….)**

  ➡ if the <u>resources</u> to promptly reconstruct the data in the computing center are limited



CMS paper, 2024

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

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

➡ **Deal with dataflow instead of latency**

    ➡ **decouple** DAQ from High Level Triggers

    ➡ decouple dataflow from storage, with temporary buffers

    ➡ Use COTS network and processing

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

➡ **Deal with dataflow instead of latency**

　➡ **decouple** DAQ from High Level Triggers

　➡ decouple dataflow from storage, with temporary buffers

　➡ Use COTS network and processing

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

　➡ reducing request rates on DAQ software

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

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

➡ **Deal with dataflow instead of latency**

  ➡ **decouple** DAQ from High Level Triggers

  ➡ decouple dataflow from storage, with temporary buffers

  ➡ Use COTS network and processing

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

  ➡ reducing request rates on DAQ software

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

➡ **Use networks as soon as possible**

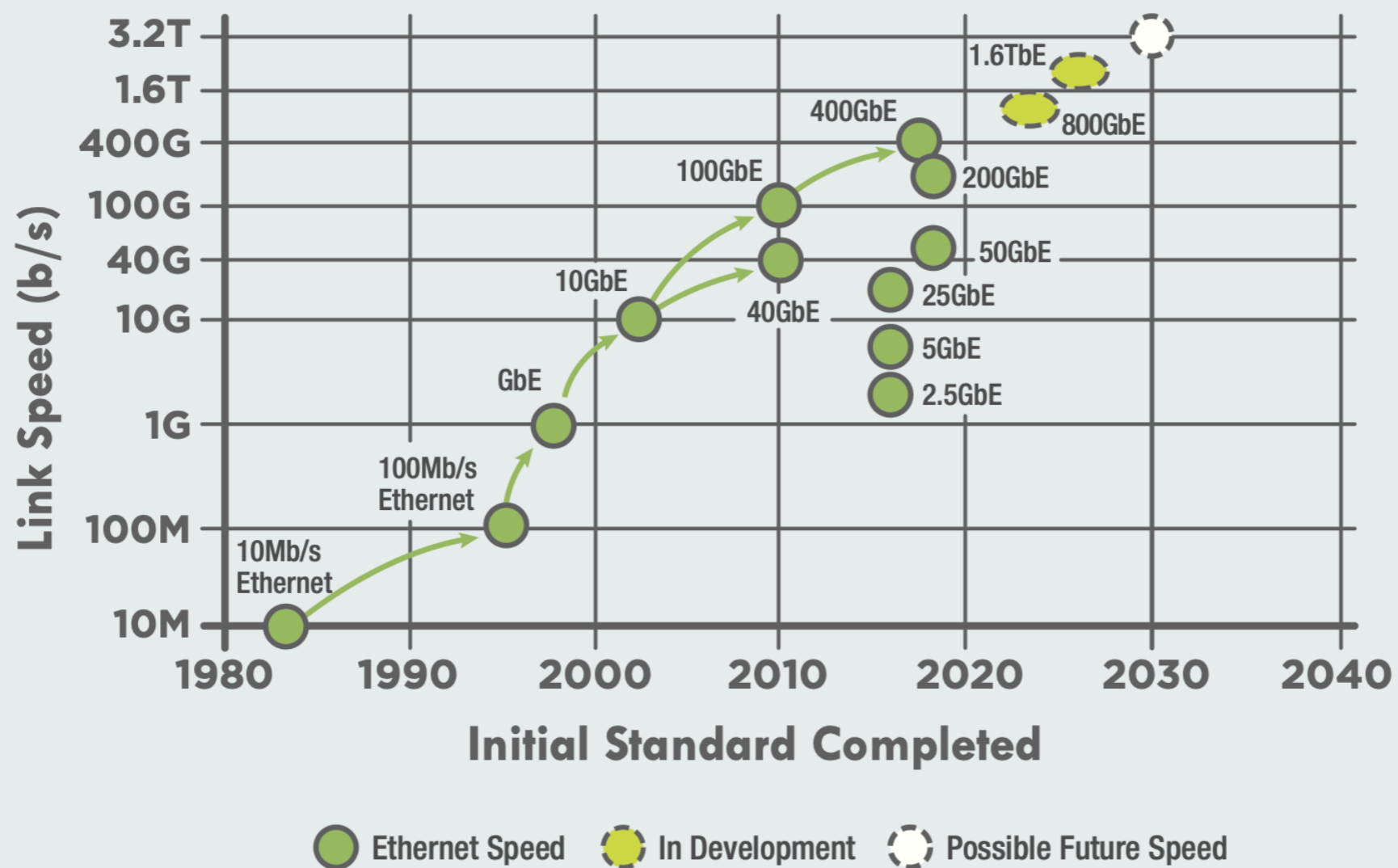  ➡ toward commercial bidirectional point-to-multipoint architecture

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

➡ **Deal with dataflow instead of latency**

- ➡ **decouple** DAQ from High Level Triggers

- ➡ decouple dataflow from storage, with temporary buffers

- ➡ Use COTS network and processing

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

- ➡ reducing request rates on DAQ software

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

➡ **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

ETHERNET SPEEDS

# ISOTDAQ 2024

## 14th International School of Trigger and Data Acquisition

### Jun. 19th – Jun. 28th 2024
### University of Science and Technology of China, Hefei, China

**Target Audience:** Physicists, engineers and computer scientists with an interest in trigger and data acquisition systems

**Places are limited:** Acceptance is by a selection committee

**Registration Deadline:** 31 January 2024

**Email:**
isotdaq.schools@cern.ch

**Website:**
https://indico.cern.ch/event/1337180

### ☑Trigger

- Modular Electronics
- Front-End Electronics
- Associative memories

### ☑DAQ

- ADC, TDC, Detector Readout
- Micro Controllers
- VMEBus, xCTA, PCI, PCIe

### ☑Application Examples

- General Concepts for TDAQ
- Insight on LHC TDAQ
- Non-LHC Systems