



# DAQ - from raw data to disk

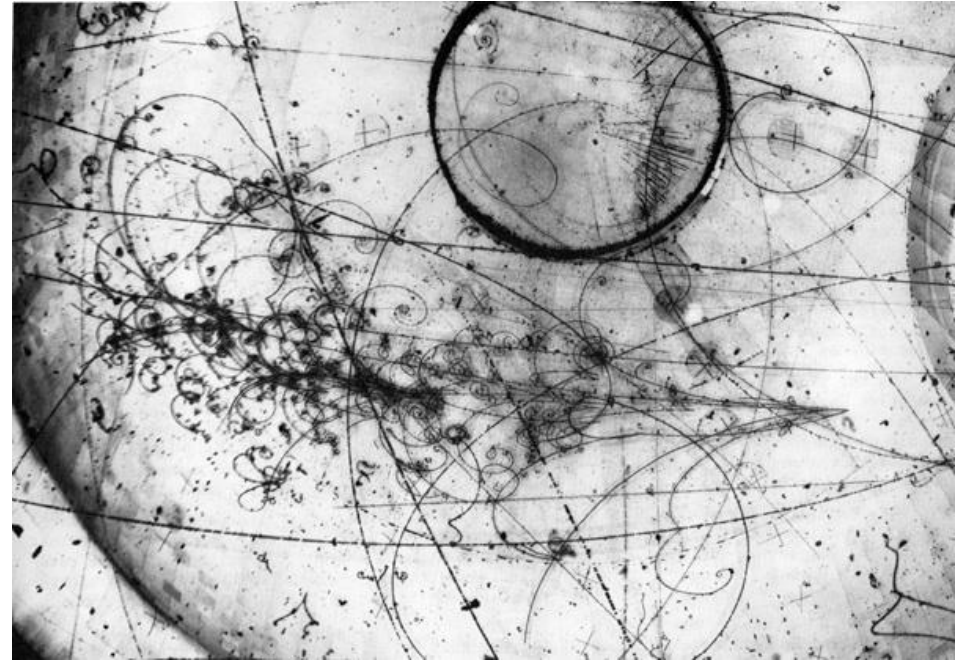
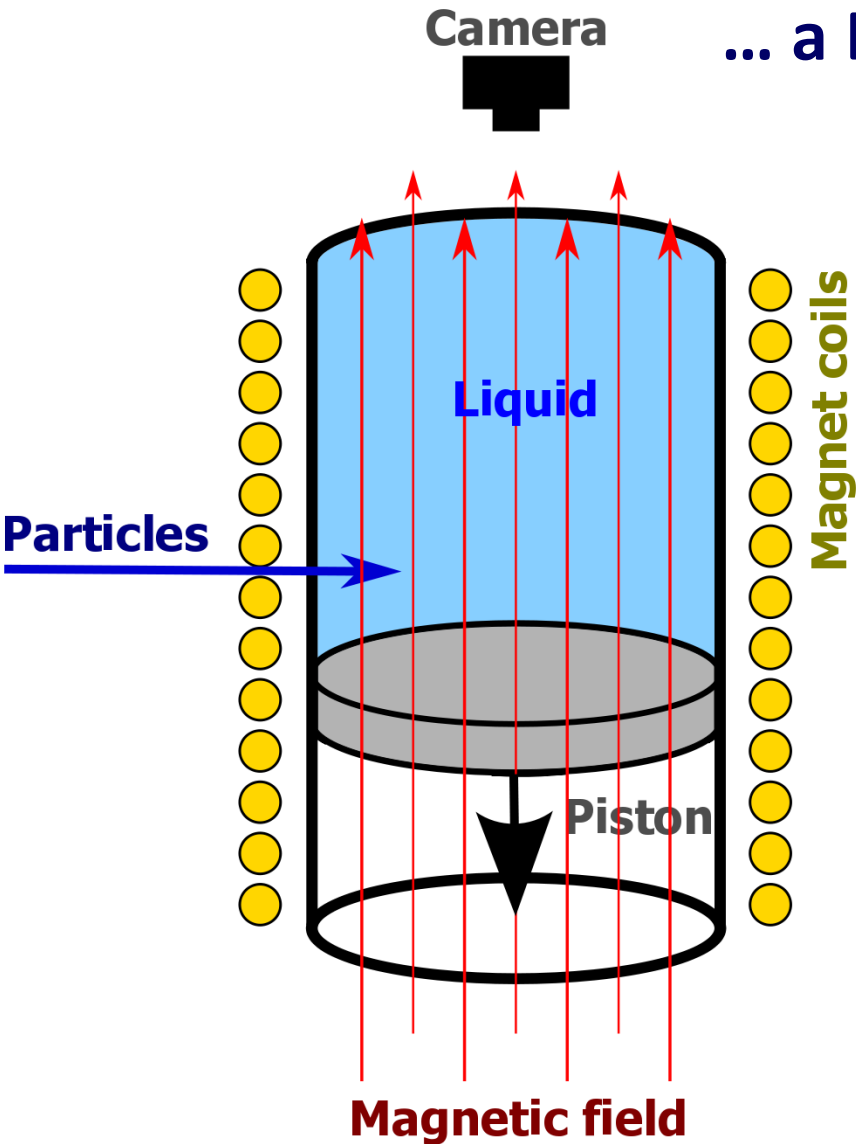
AIDA student tutorial, Vienna, 25<sup>th</sup> March, 2014

Hannes Sakulin, CERN/PH-CMD



# Once upon a time...

... a DAQ system looked like this

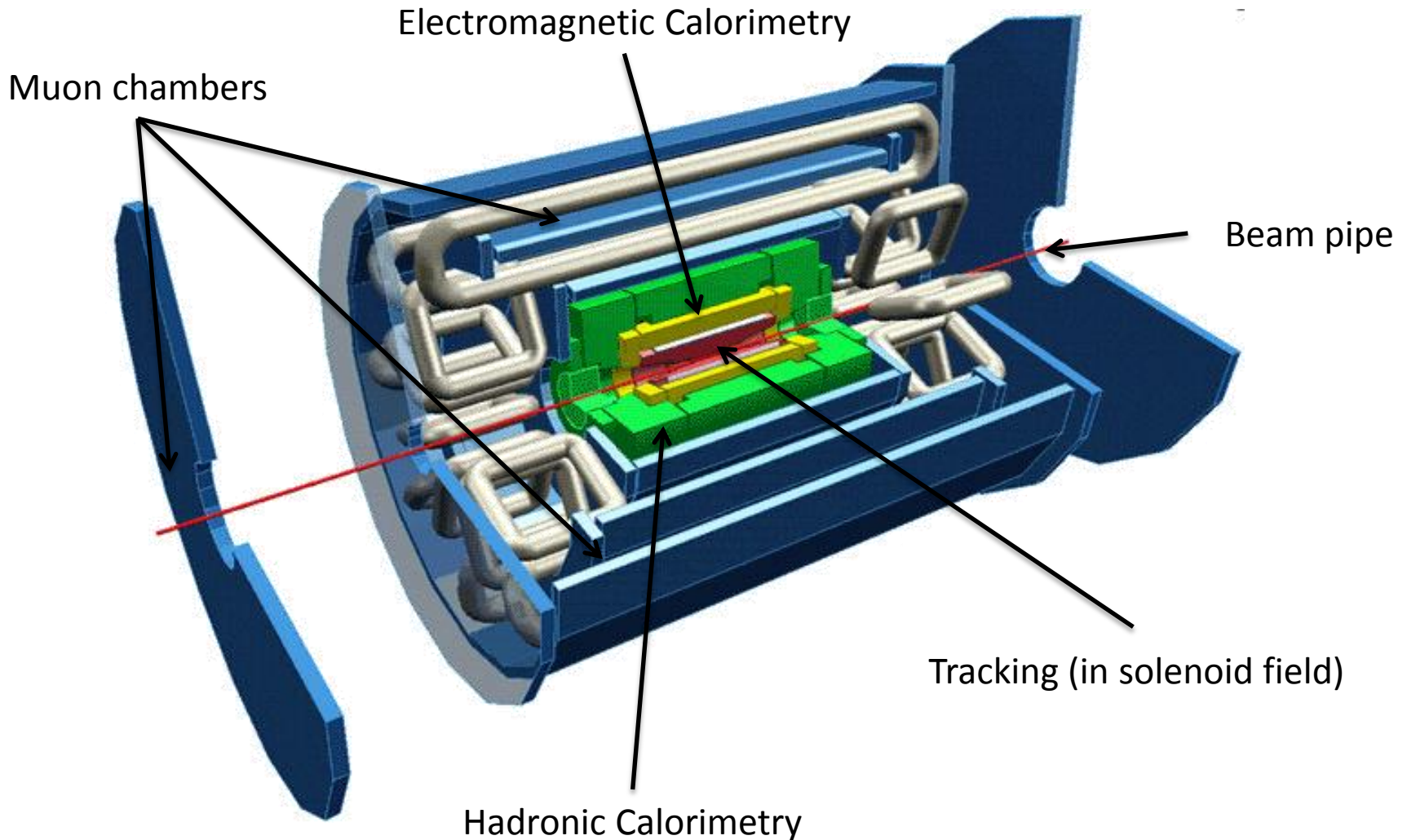


# And event selection was done like that ...





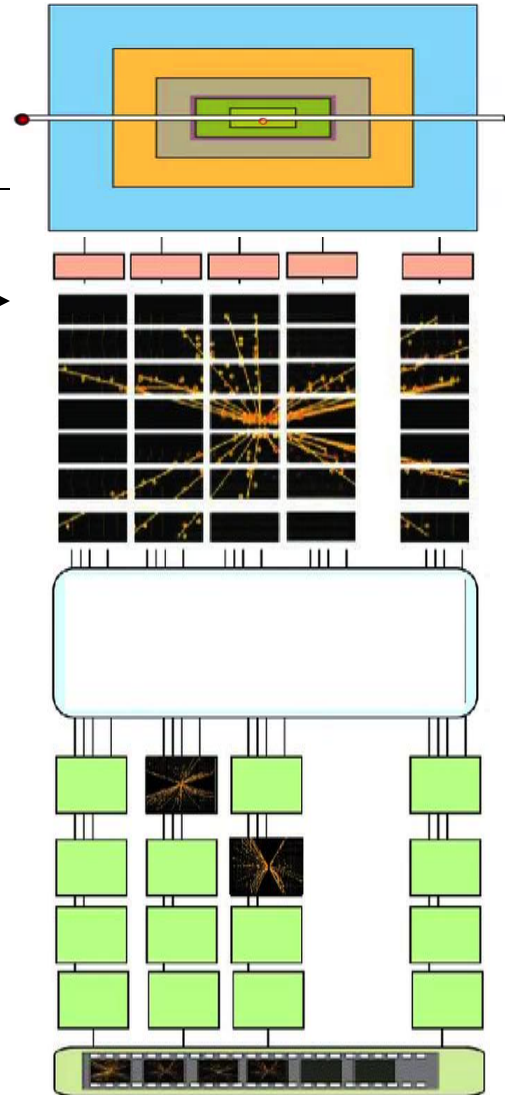
# How do we “read” a modern experiment ?



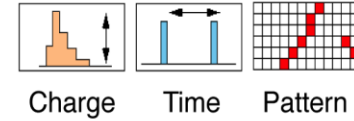
**$O(100 \times 10^6)$  detector channels**



# A typical DAQ system



## 0) Detectors



## 1) Front-end systems

Analog-digital conversion  
 Detector readout link  
 Feature extraction

## 2) DAQ readout Links

Interface from custom to commercial electronics

## 3) Event Building

Network

## 4) Event Filter (High Level Trigger)

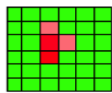
Thousands of CPU cores

## 5) Local mass storage

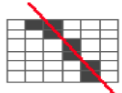


# Why do we need a trigger ?

**Beam crossing (bx) frequency given by the accelerator**  
**LHC: 40 MHz**



Energy



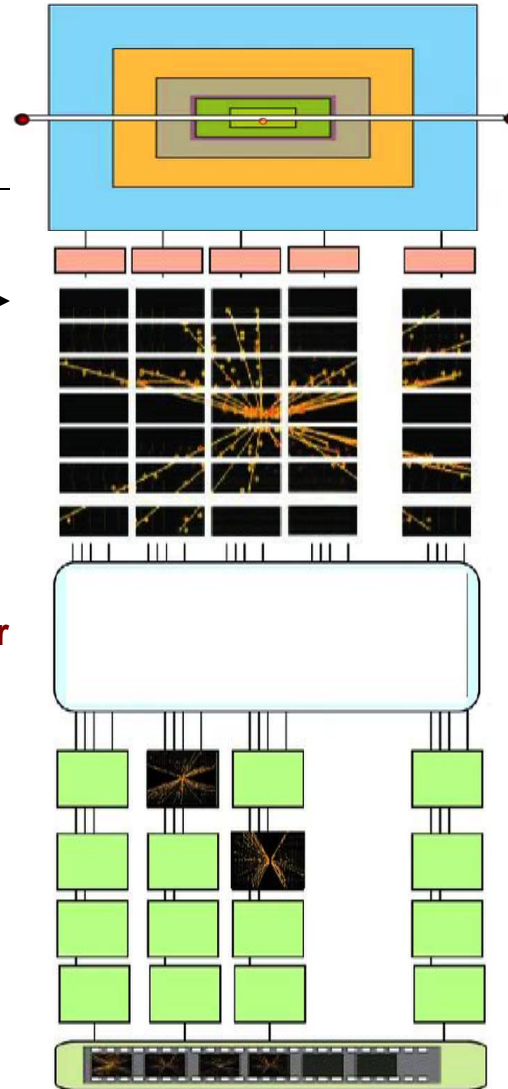
Tracks

**Level-1 Trigger**

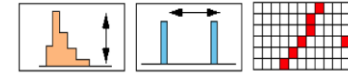
Level-1 accept

**Maximum readout rate limited by**

- B/W of readout links
- B/W of DAQ links
- Computing Infrastructure  
**Networks and Processing Power**



## 0) Detectors



Charge

Time

Pattern

## 1) Front-end systems

Analog-digital conversion  
 Feature extraction

## 2) DAQ readout Links

Interface from custom to commercial electronics

## 3) Event Building

Network

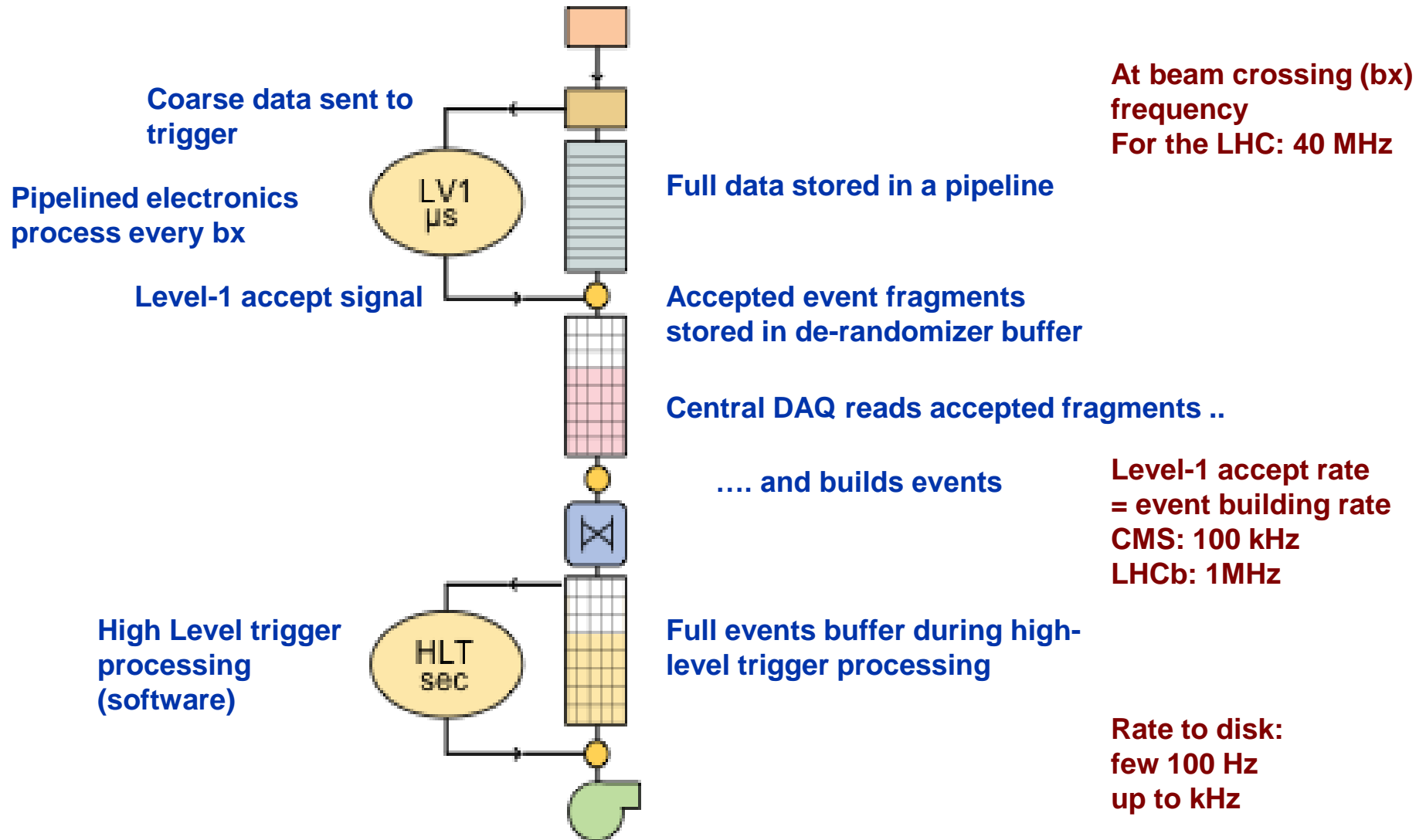
## 4) Event Filter (High Level Trigger)

Thousands of CPU cores

## 5) Local mass storage

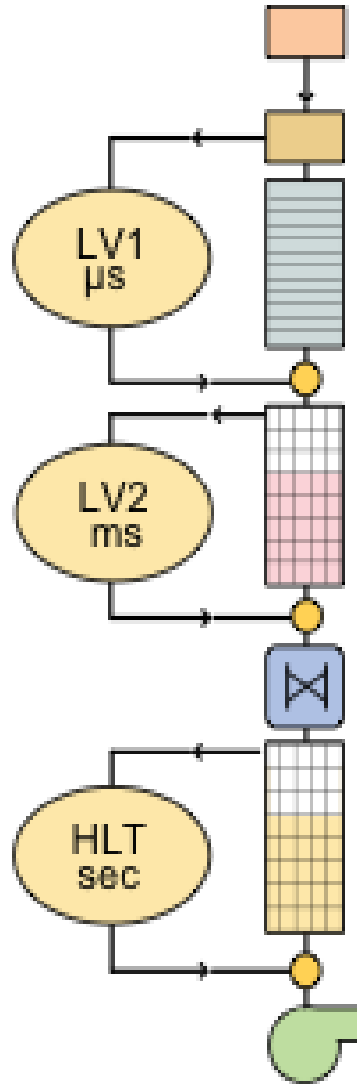


# Trigger and DAQ interaction





# Trigger and DAQ interaction



**Beam crossing (bx) frequency**  
**For the LHC: 40 MHz**

**Level-1 accept rate**  
**ATLAS: 65 KHz**

**Some experiments have multiple trigger levels**

**Level-2 accept rate = event building rate**  
**ATLAS: 6 KHz**

**Rate to disk:**  
**few 100 Hz**  
**up to kHz**



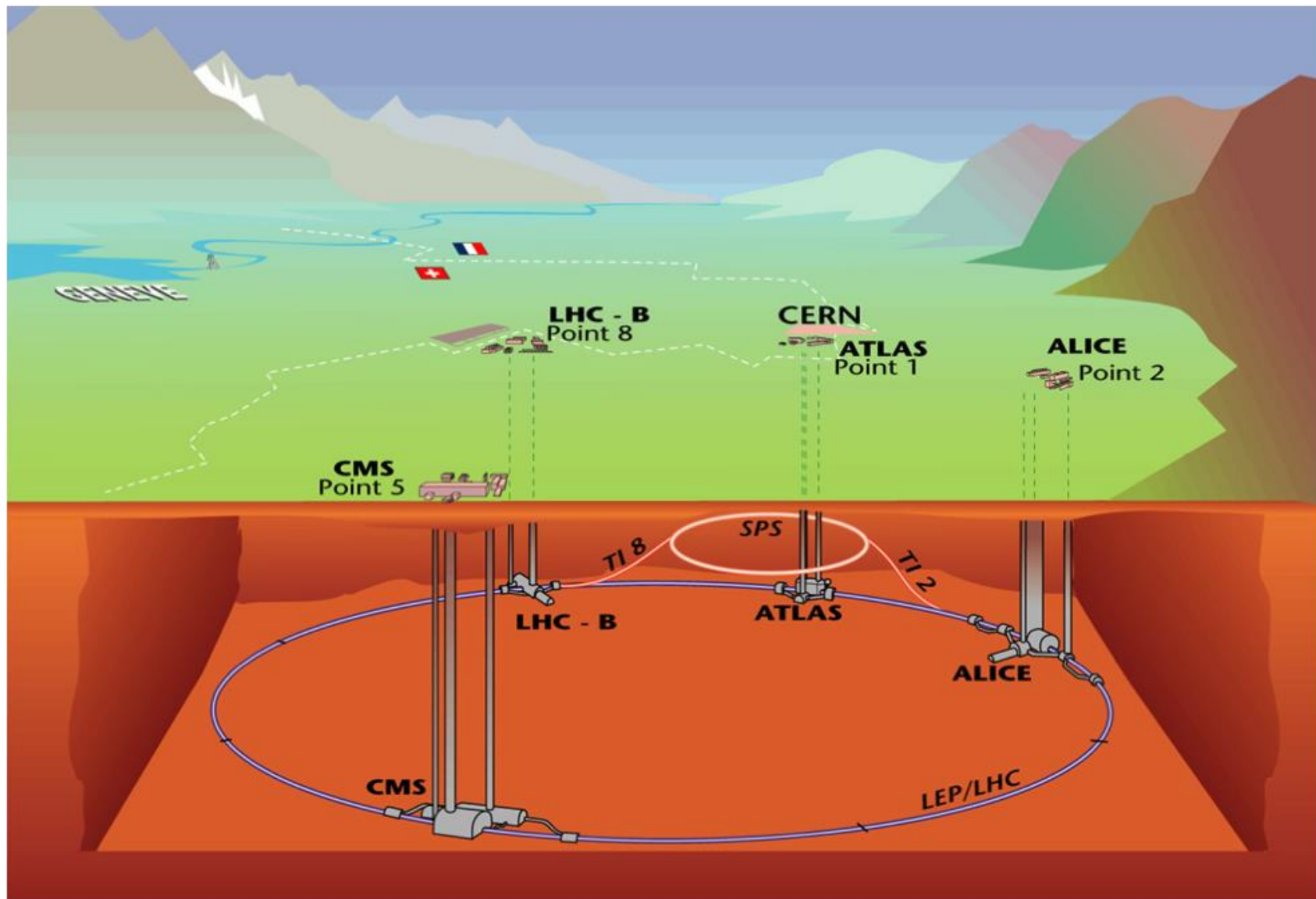


# Disclaimer

- This talk is based on the DAQ system of the LHC experiments. Many concepts are nevertheless generic.
- Thanks to all my colleagues who contributed material to this talk (actively and passively), especially to Christoph Schwick and Sergio Cittolin.

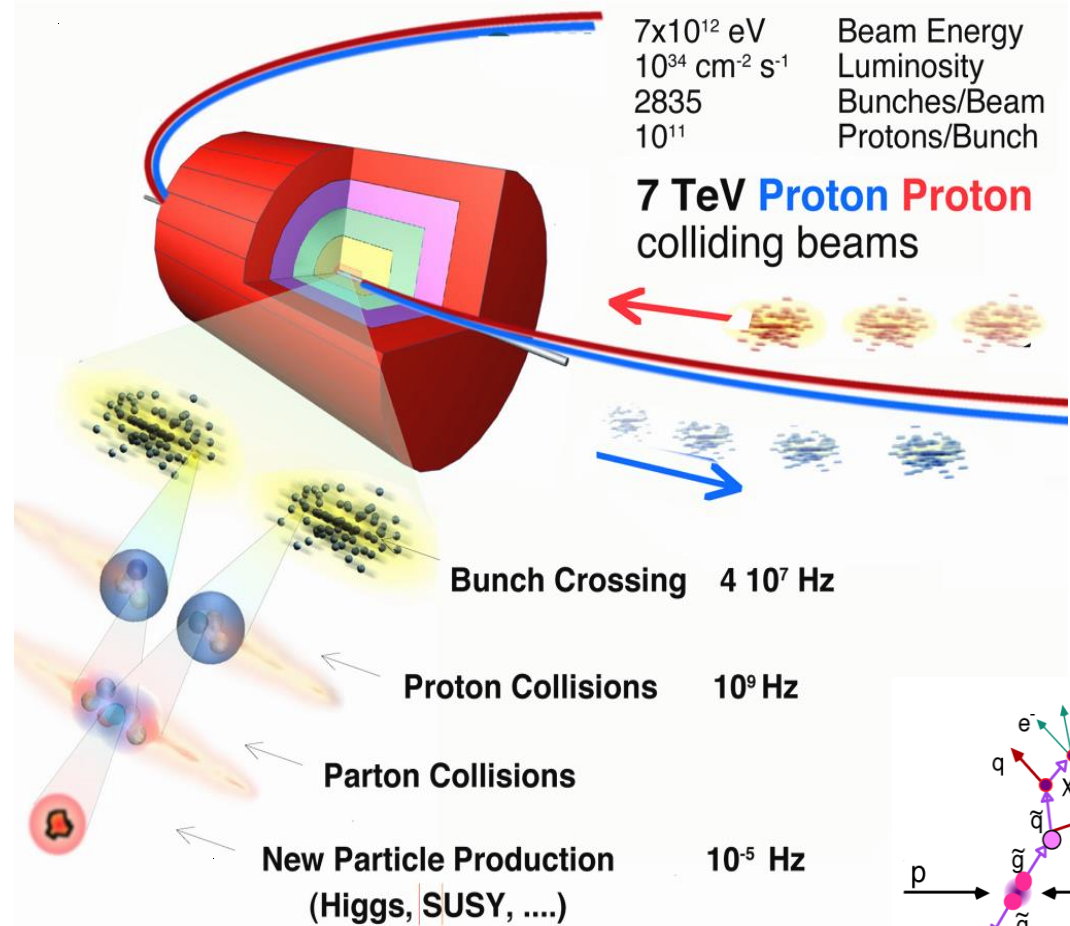
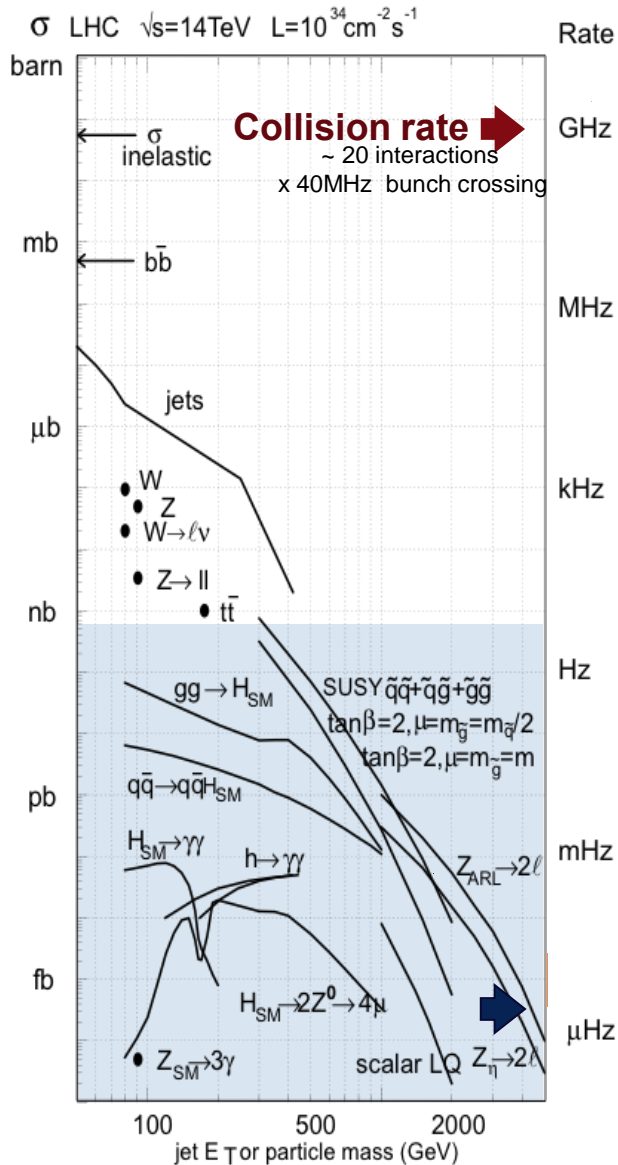


# DAQ at the LHC





# Proton-proton collisions at the LHC. Searching issue

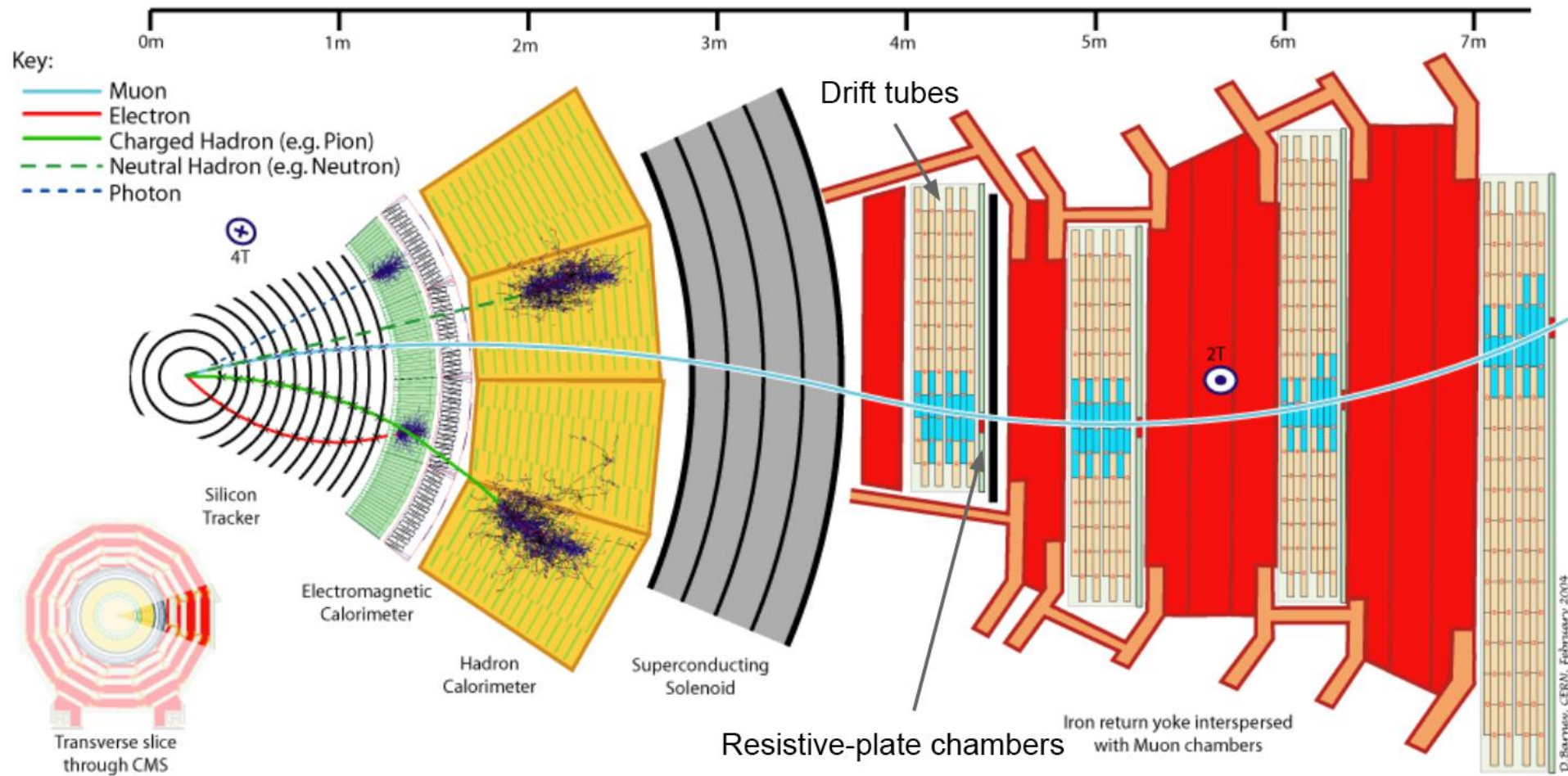


**Collision Rate:  $\sim 10^9$  Hz. Event Selection:  $\sim 1/10^{13}$**



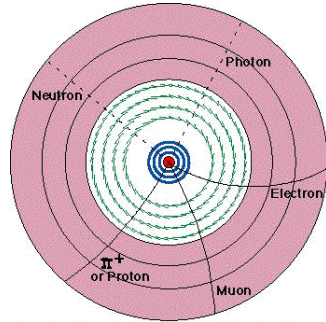
# 0) Detectors

# Typical detector layout



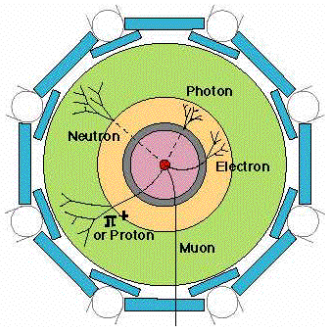
Combination of detectors allows particles to be identified

# Detectors are a science by themselves



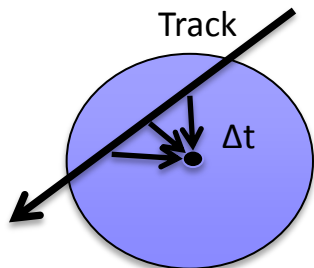
- Tracking detectors

- Charge collected in silicon semiconductor
- Precise position measurement – high number of channels



- Calorimeters

- Scintillation light collected by Photo-Detectors
- Use amplitude to measure deposited energy



- Muon detectors

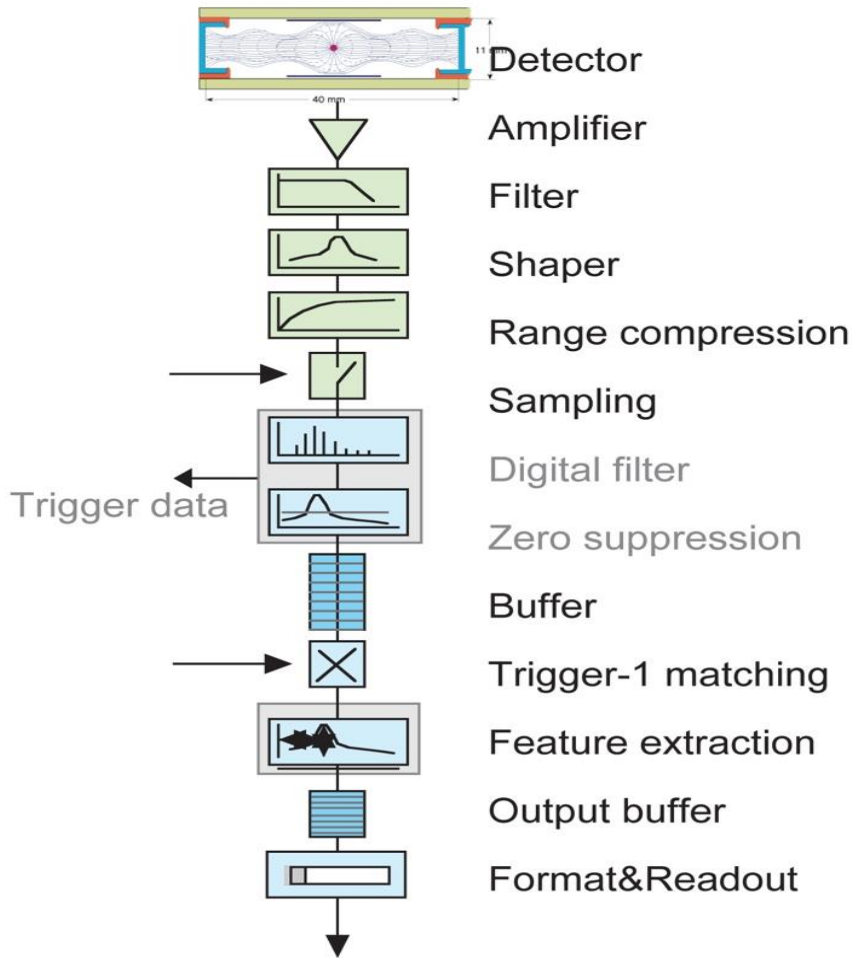
- Drift chamber. Electrons drift to central wire and cause avalanche.
- Timing determines position



# 1) Front-end systems



# Typical front-end electronics



## Signal

typically small current pulse

## Analog electronics

depending on detector type

## Analog/Digital conversion

## Digital Signal Processing (DSP)

### Front-end pipeline buffer

Store data while trigger is computing its decision

### De-randomizing buffer

Store data until read out by DAQ

## Readout-link to central DAQ

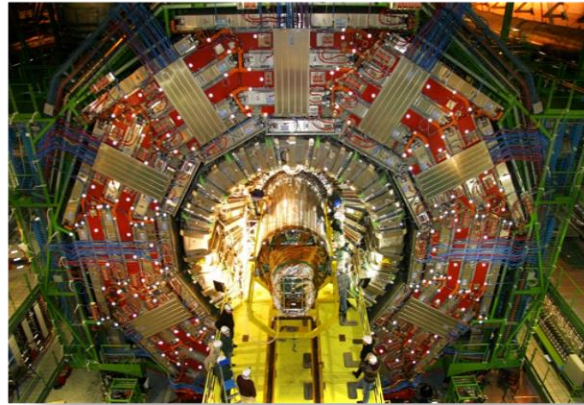




# Getting the data out of the detector

10K – 100M rad

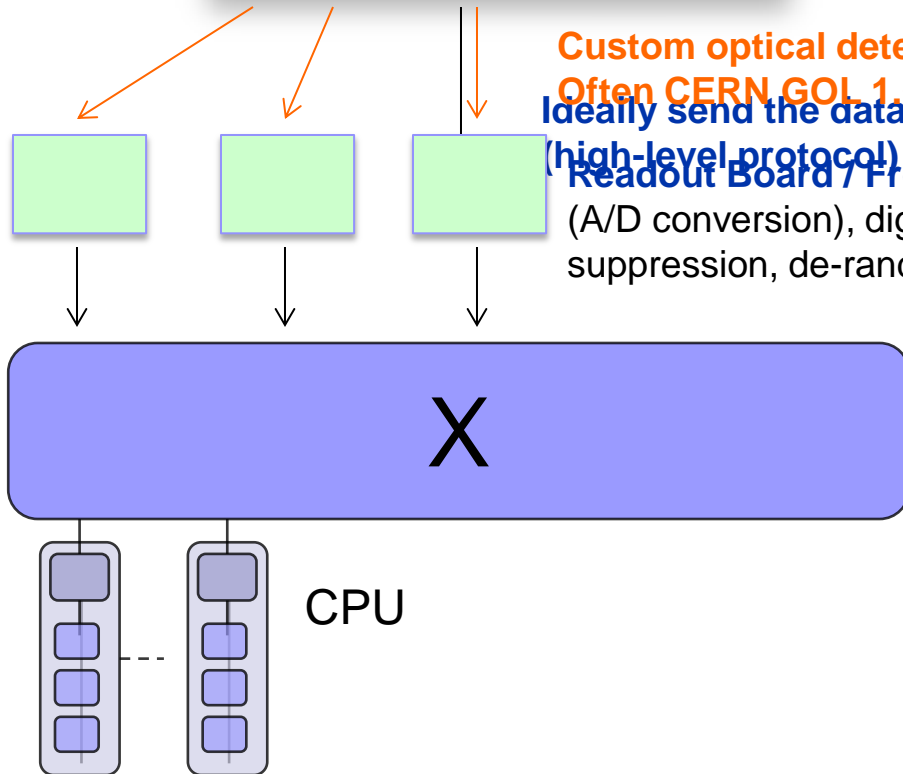
Max 4T



~zero mass



- Gigabit Optical Link (GOL)
- ASIC
- Radiation hard
- designed at CERN



Custom optical detector-readout link.

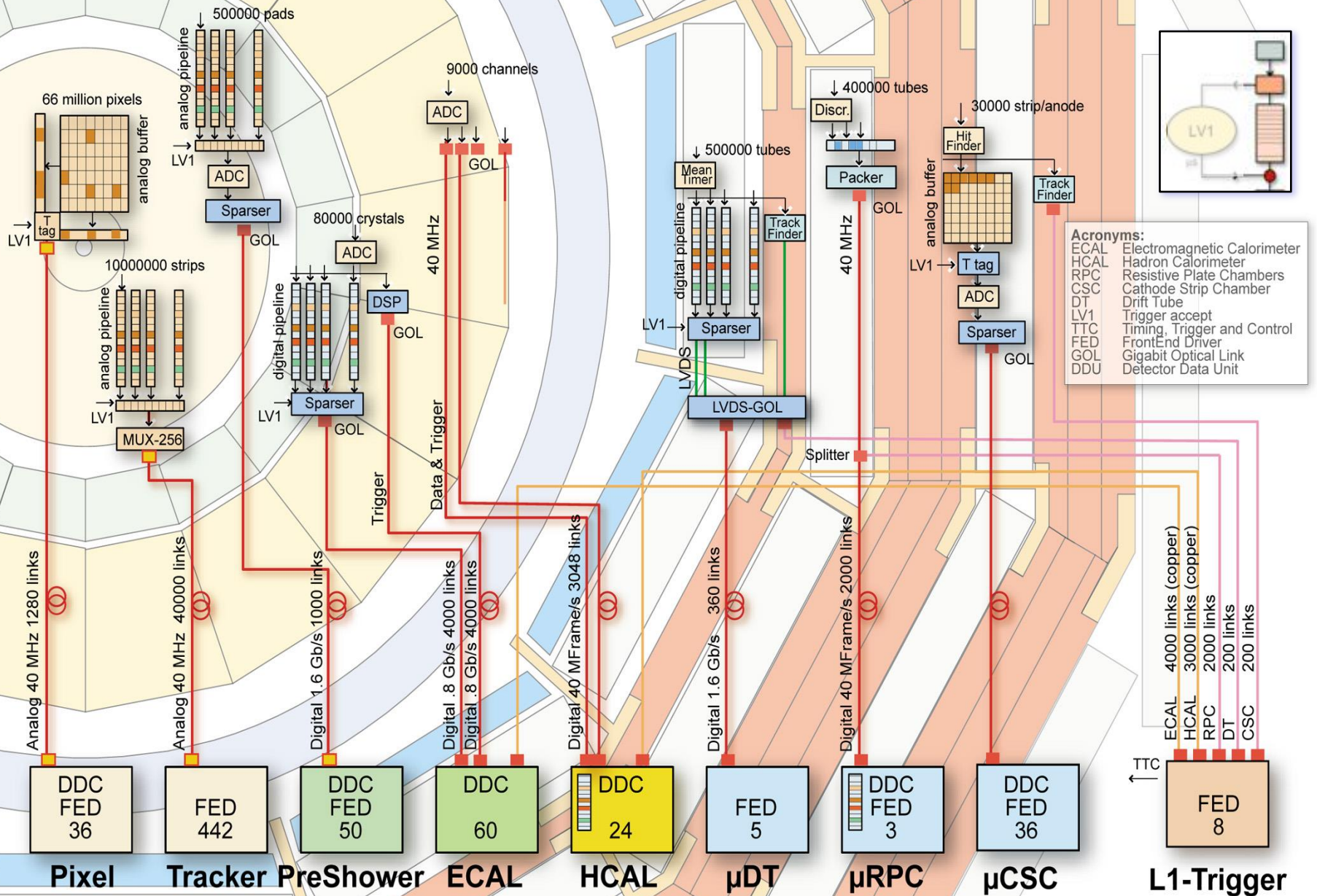
Often CERN GOL 1.6 Gb/s

Ideally send the data directly to the event builder (high-level protocol)

Readout Board / Frontend Driver / Readout Driver (A/D conversion), digital signal processing, zero suppression, de-randomizer buffer

CPU

# Example: front-end systems in CMS

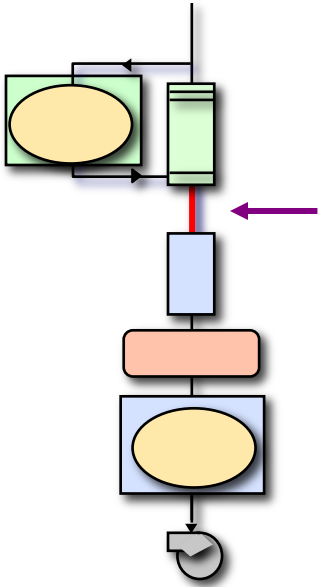




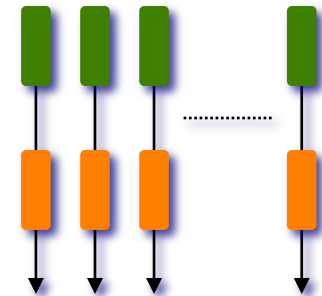
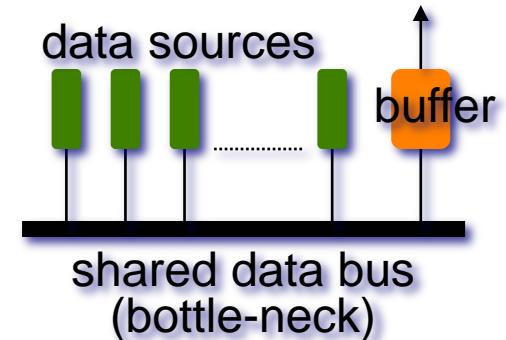
## 2) DAQ readout Links



# Data Flow: Data Readout



- Former times: Use of bus-systems
  - VME or Fastbus
  - Parallel data transfer (typical: 32 bit) on shared bus
  - One source at a time can use the bus
  
- Today: Point to point links
  - Optical or electrical
  - Data serialized
  - Custom or standard protocols
  - All sources can send data simultaneously



- Compare trends in industry market:
  - 198x: ISA, SCSI(1979),IDE, parallel port, VME(1982)
  - 199x: PCI( 1990, 66MHz 1995), USB(1996), FireWire(1995)
  - 200x: USB2, FireWire 800, PCIeexpress, Infiniband, GbE, 10GbE



# Example readout link : ATLAS



**High-Speed Optical Link for Atlas (HOLA) mezzanine board plugged onto Readout Driver electronics.**

**Detector electronics push 32-bit parallel data  
+ header + trailer**



**Optical fiber  
2 Gb/s**



**Readout-Buffer Input (ROBIN) PCI-x card  
plugged into a PC**

**Receiving up to 3 links**

**Data transfer to PC memory via DMA**

# Getting the data into the PC

## Problem:

Read data into PC with high bandwidth and low CPU load

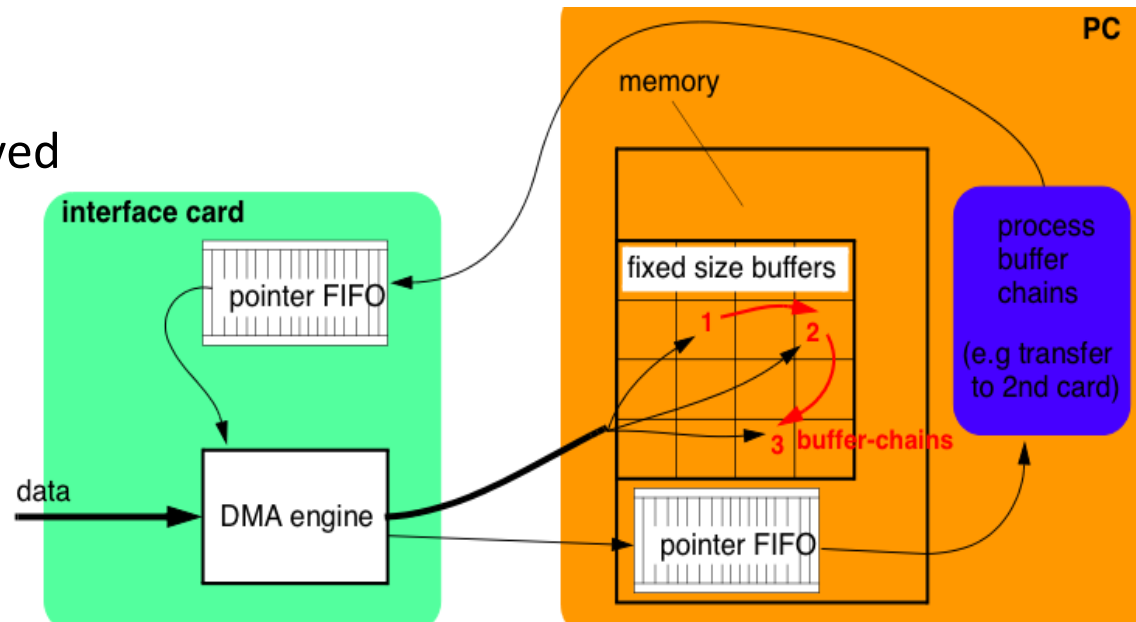
Note: copying data costs a lot of CPU time!

## Solution: **Buffer-Loaning**

- Hardware shuffles data via DMA (Direct Memory Access) engines
- Software maintains tables of buffer-chains

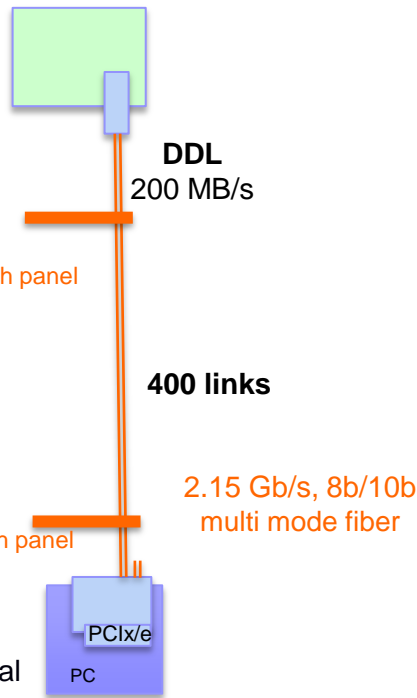
## Advantage:

- No CPU copy involved





# ALICE

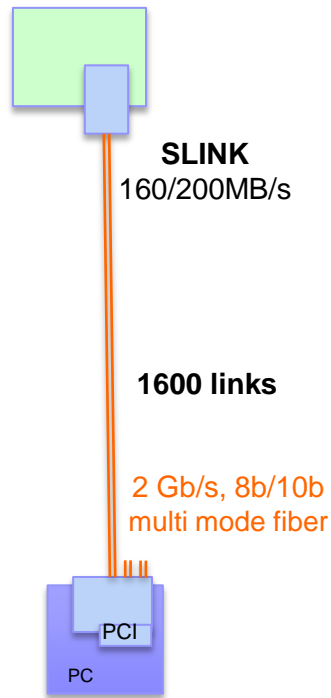


### SIU



### DIU

# ATLAS

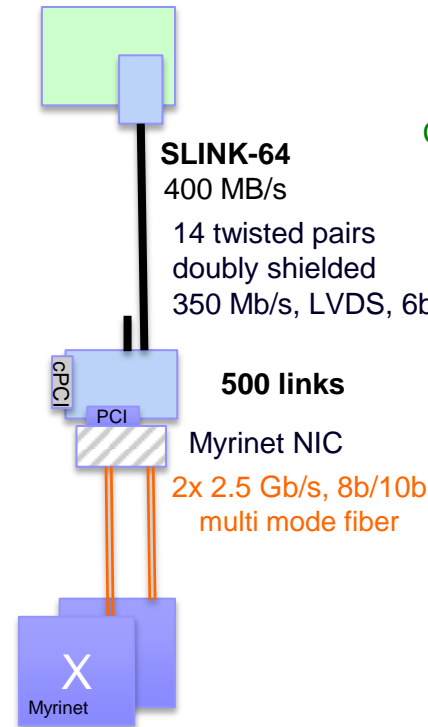


### HOLA



### ROBIN

# CMS

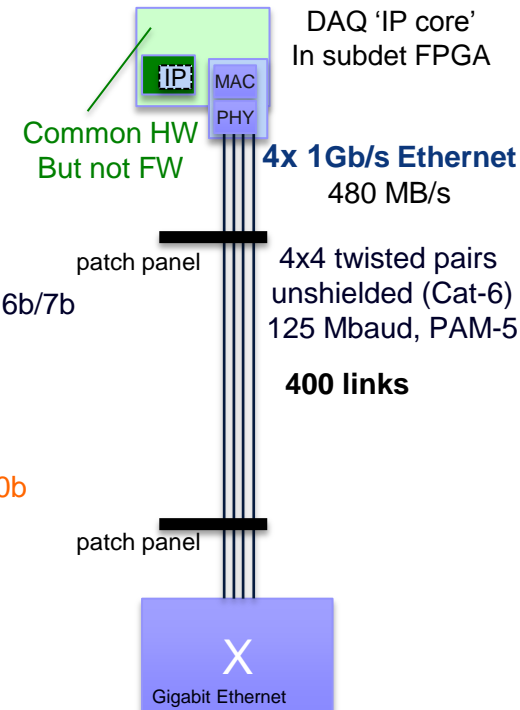


### SLINK CMC



### FRL

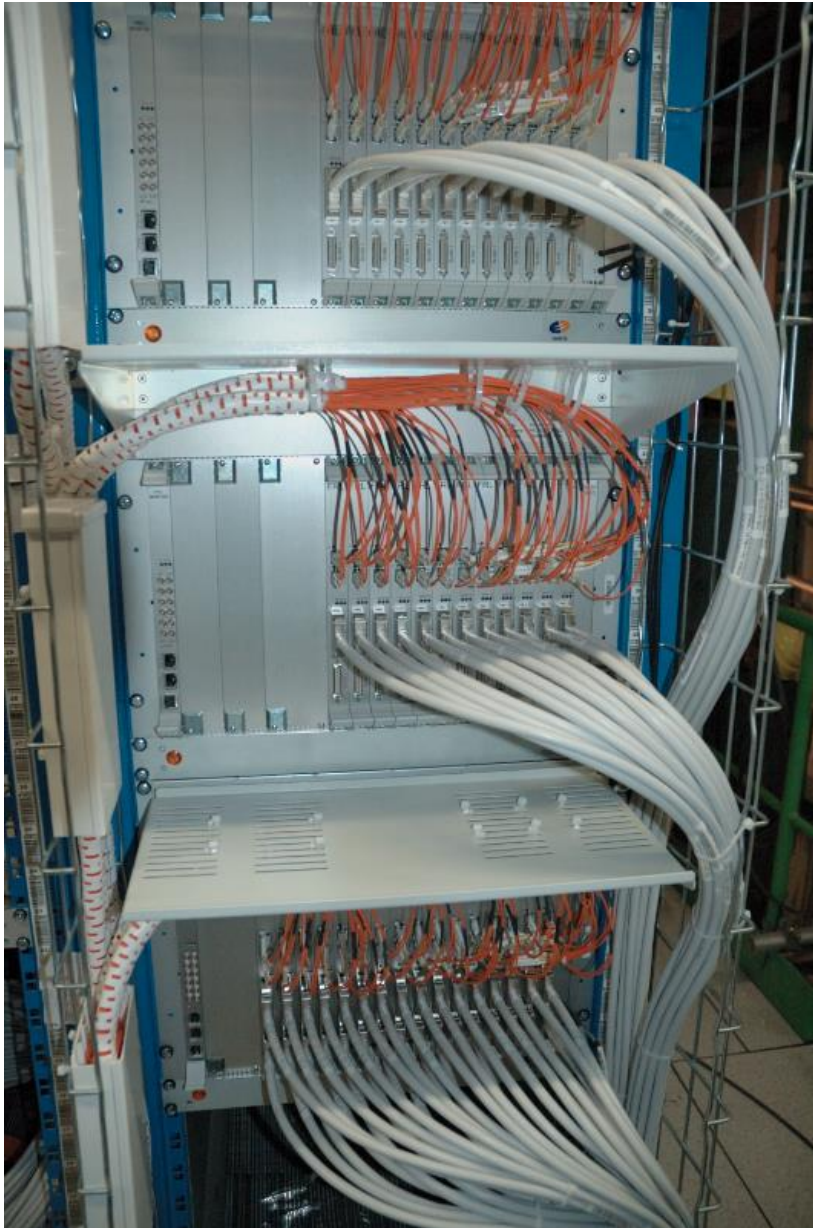
# LHCb



### TELL-1

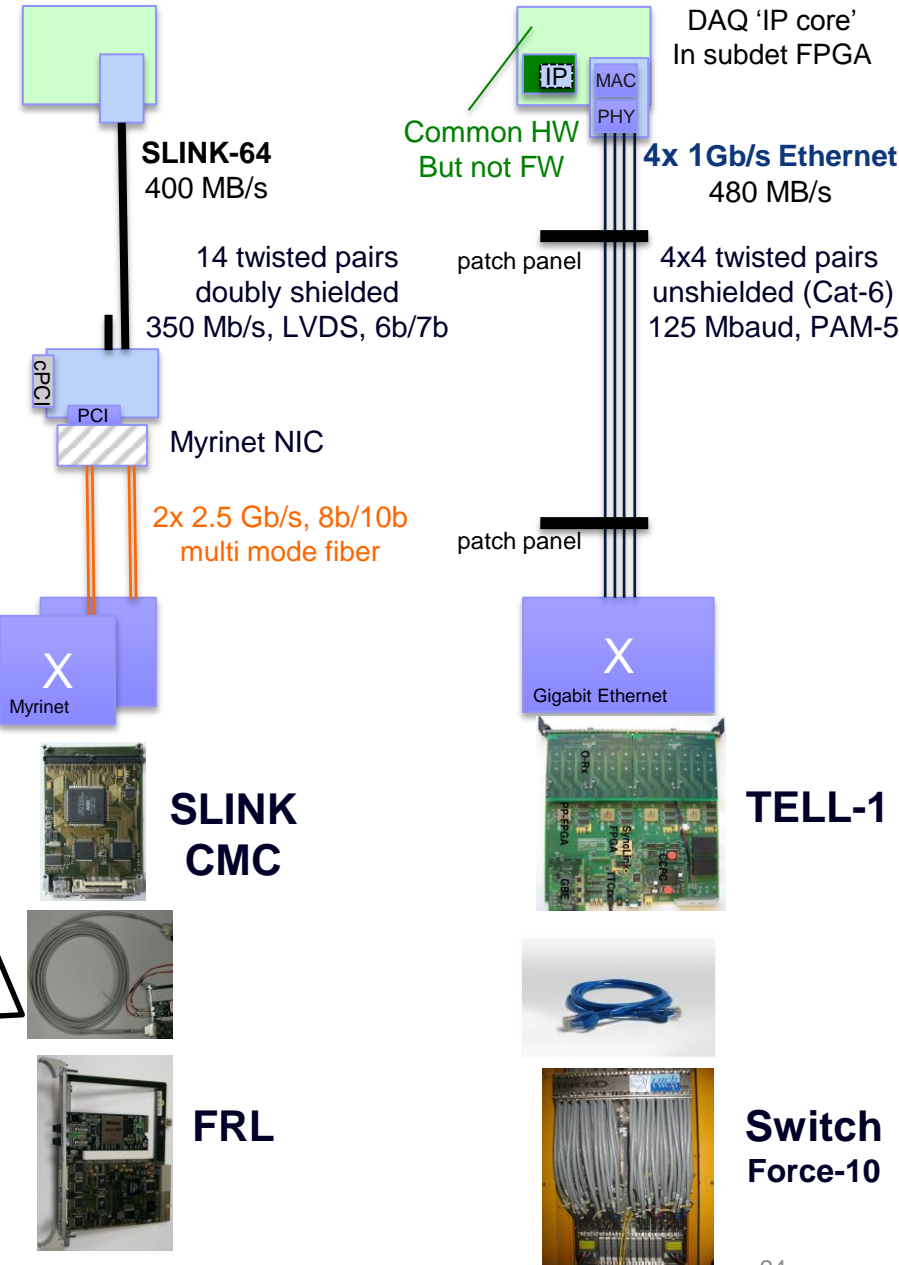


### Switch Force-10



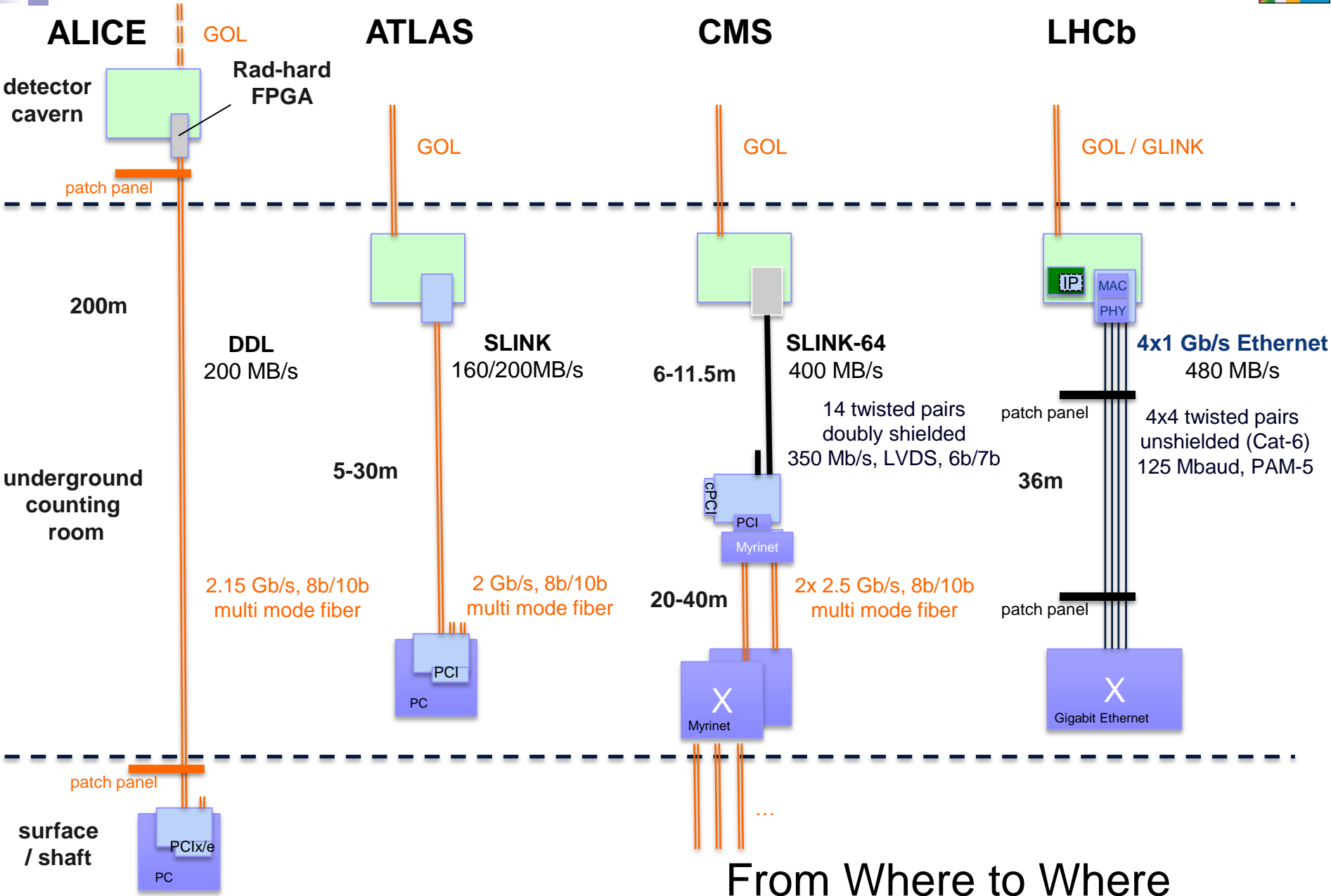
**CMS**

**LHCb**



IN





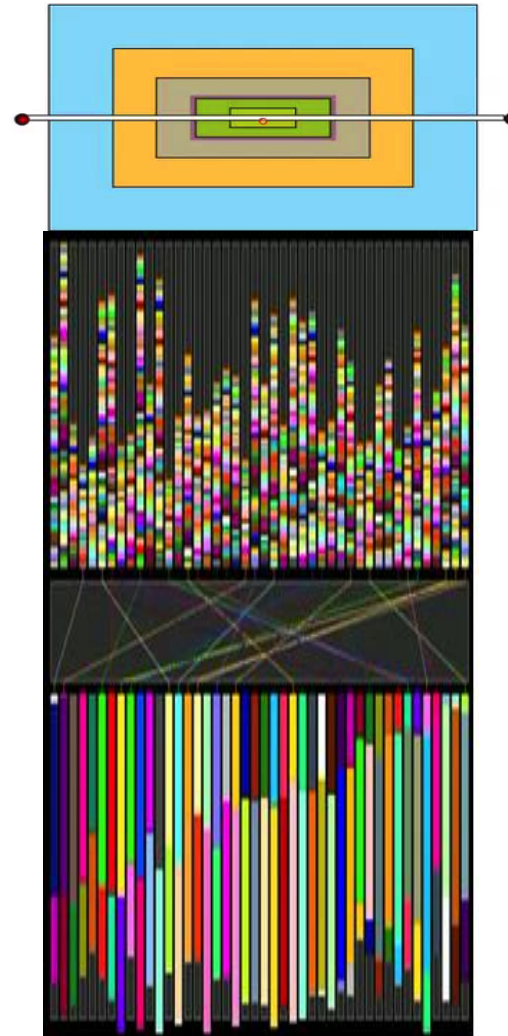
From Where to Where



# 3) Event-Building



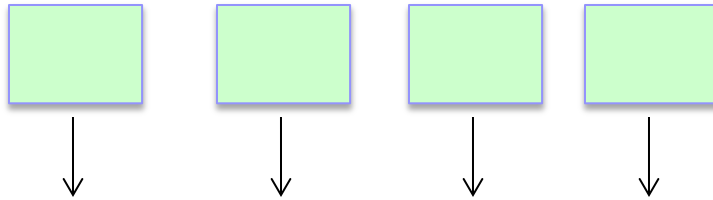
# Event Building



**3) Event Building**  
Network

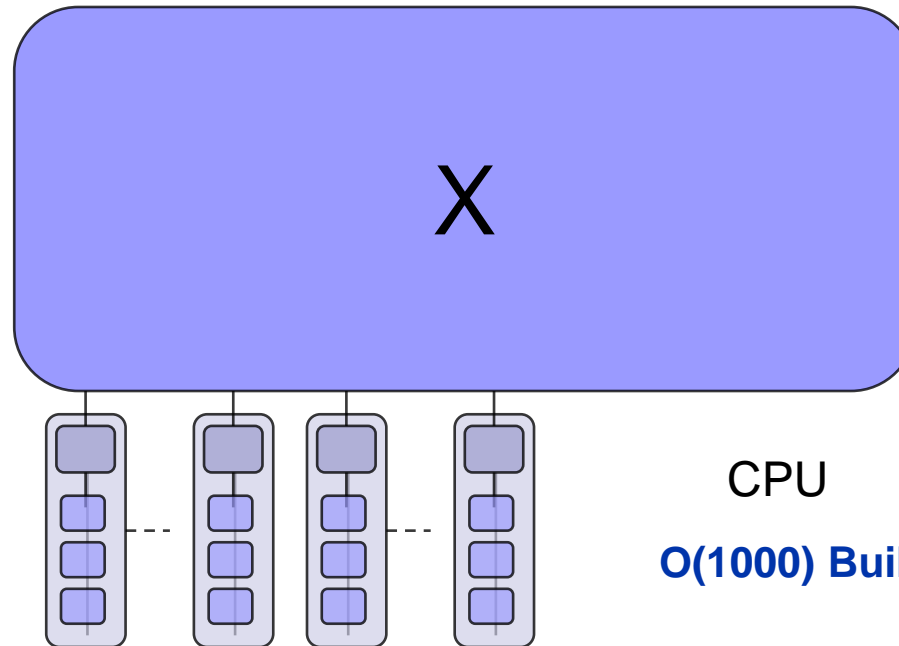


# Principle of event building



Readout Board / Frontend Driver / Readout Driver

$O(1000)$  data sources



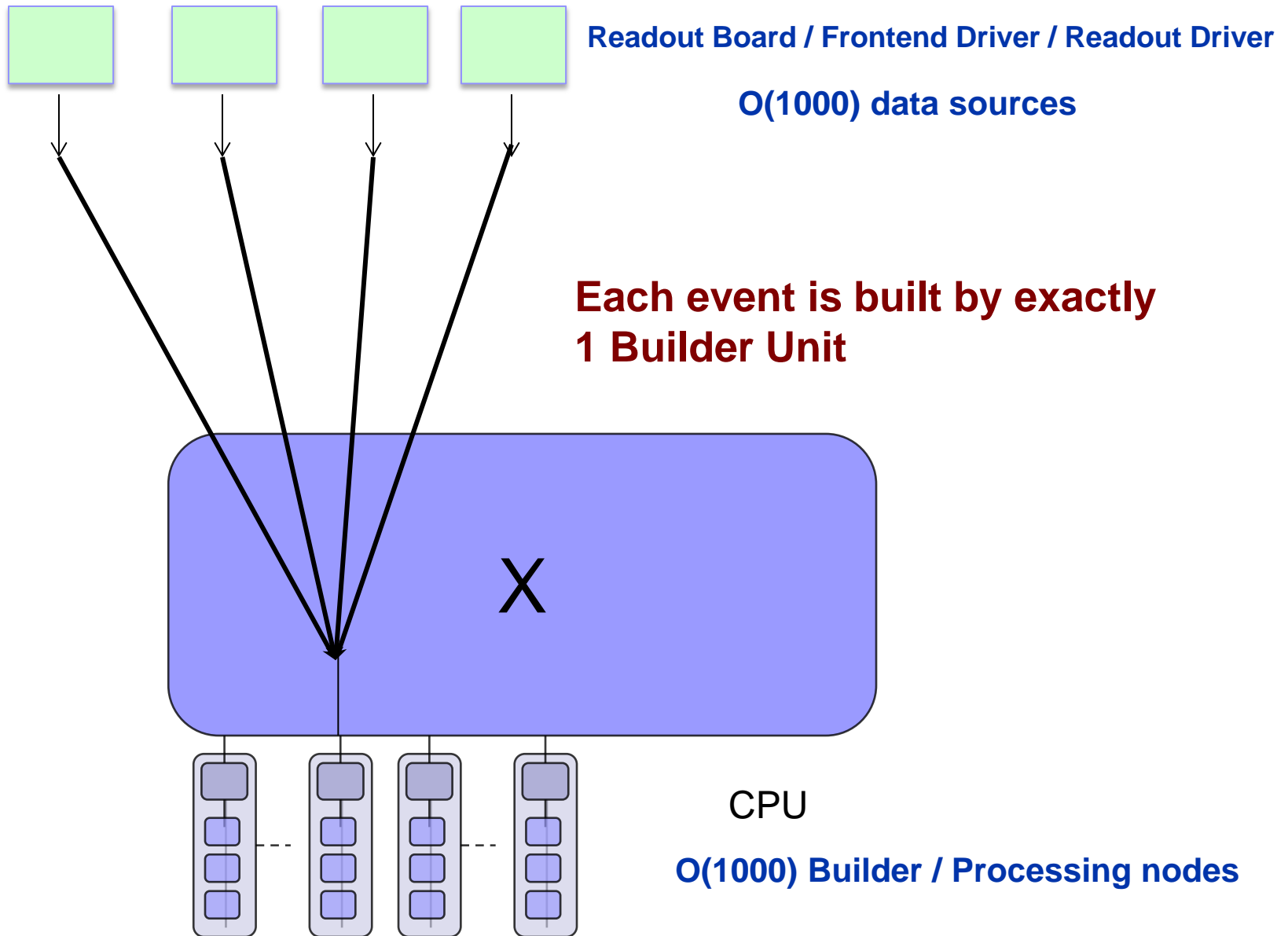
Network switch

CPU

$O(1000)$  Builder / Processing nodes



# Principle of event building

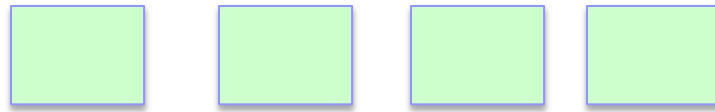


# Congestion



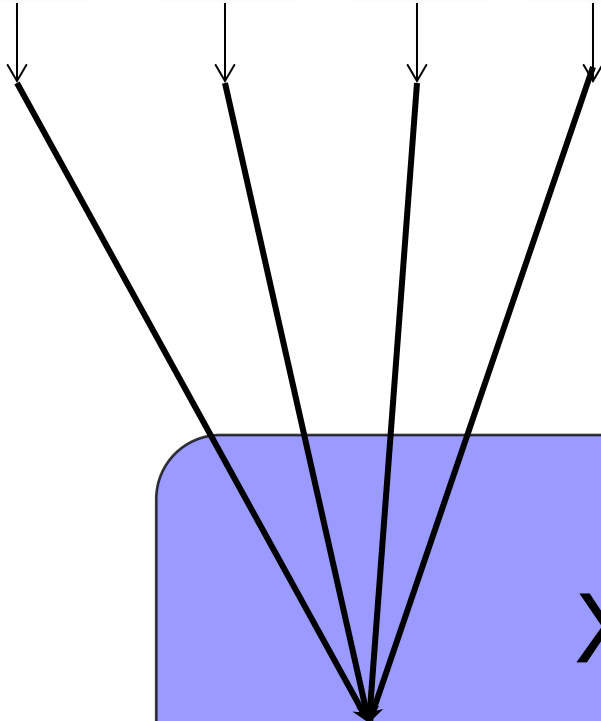


# Building the events



Readout Board / Frontend Driver / Readout Driver

$O(1000)$  data sources

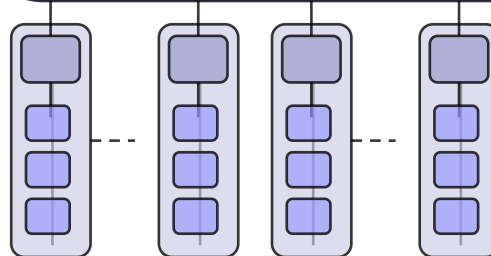
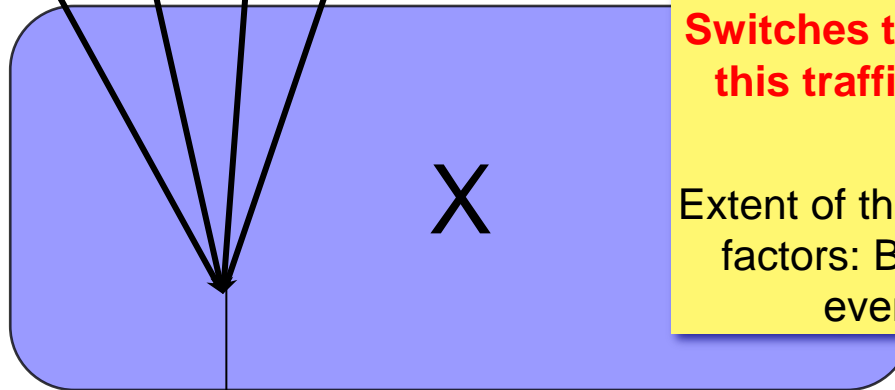


## EVB traffic

all sources send to the same destination (almost) concurrently.

**Switches typically are not built for this traffic pattern. Congestion.**

Extent of the effect depends on many factors: Buffer space, event size, event building rate, ...



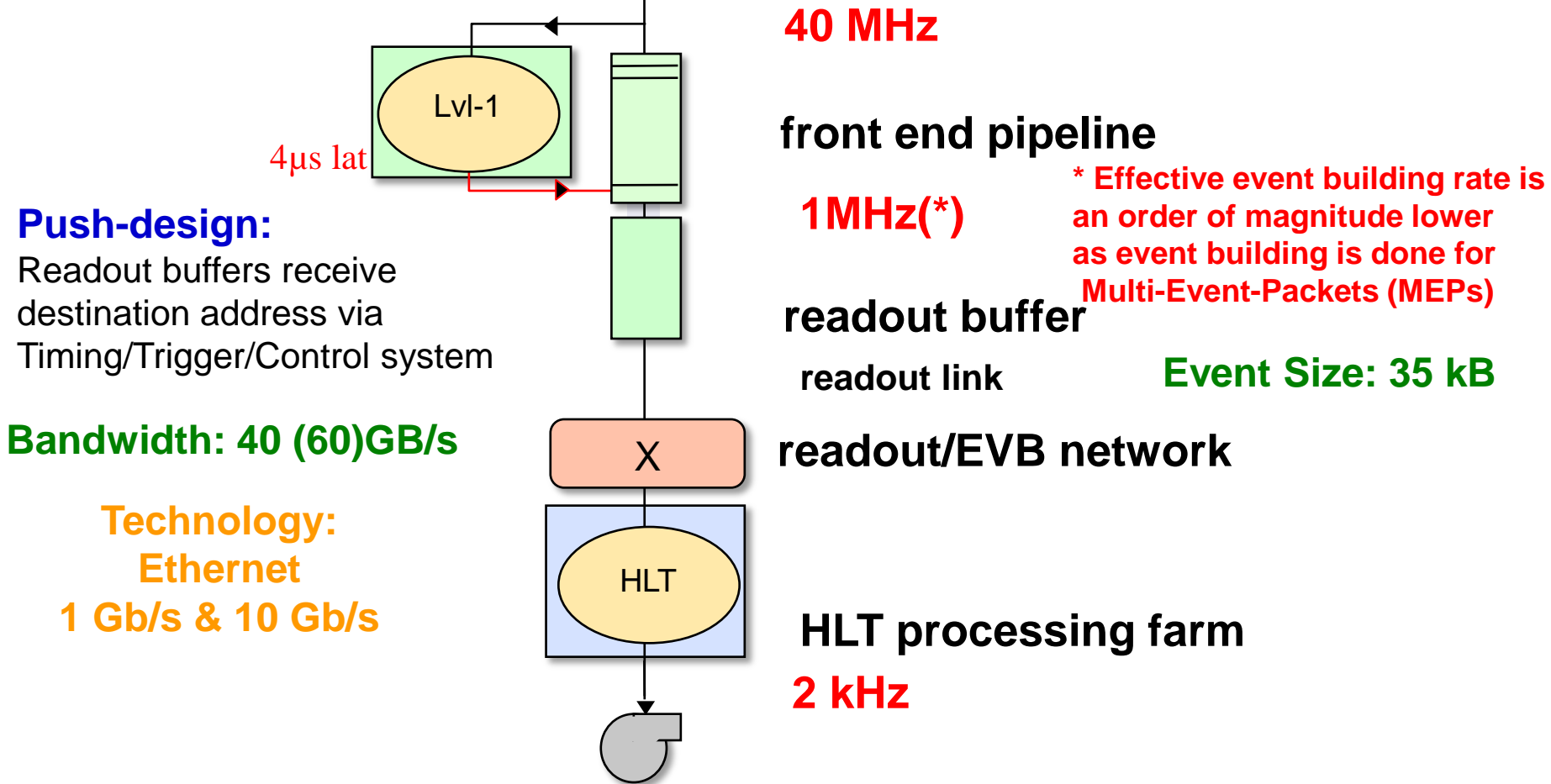
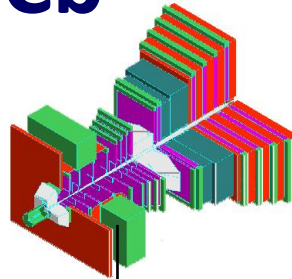
CPU

$O(1000)$  Builder / Processing nodes



# Data Flow: LHCb

- custom hardware
- PC
- network switch



**Push-design:**

Readout buffers receive destination address via Timing/Trigger/Control system

**Bandwidth: 40 (60)GB/s**

**Technology:**

**Ethernet**

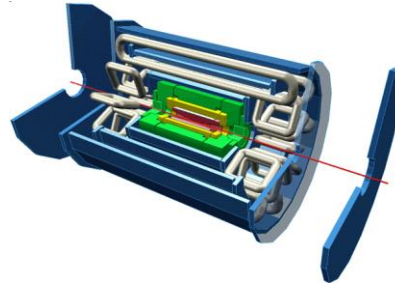
**1 Gb/s & 10 Gb/s**



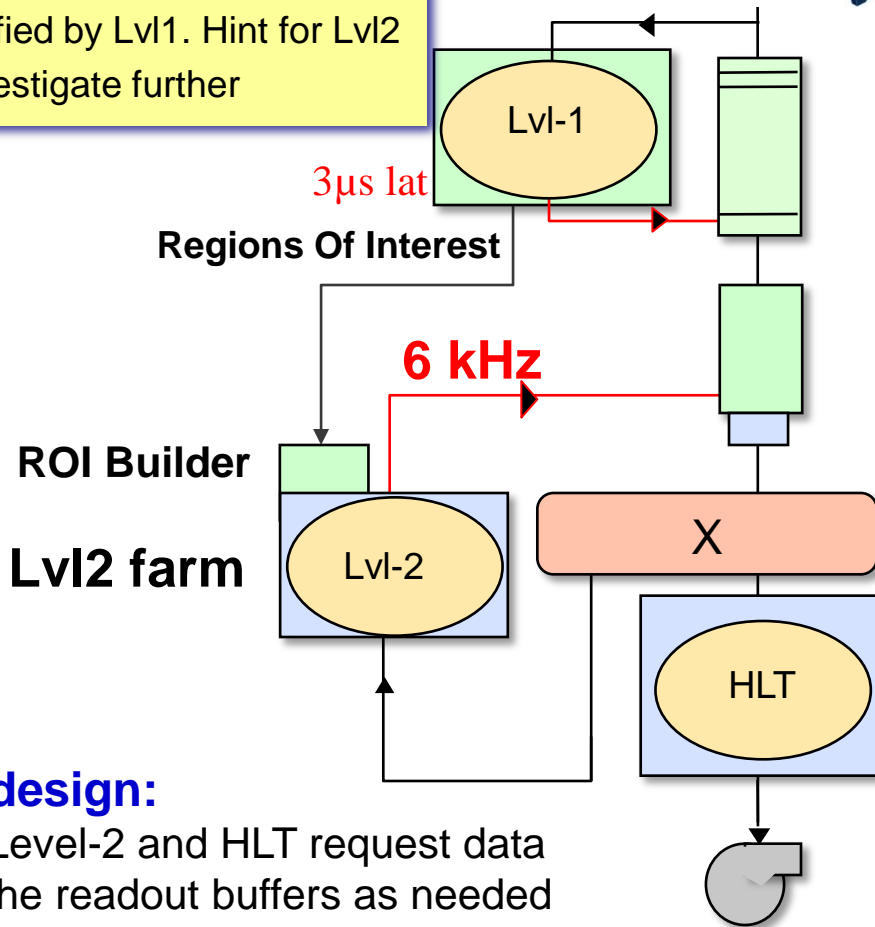


# Data Flow: ATLAS

- custom hardware
- PC
- network switch



**Region Of Interest (ROI):**  
Identified by Lvl1. Hint for Lvl2 to investigate further



**40 MHz**

**front end pipeline**

**65 kHz**

**readout link**

**readout buffer**

**Event Size: 1.5 MB**

**event builder**

**Bandwidth: 5.25 (10) GB/s**

**HLT farm**

**Technology:  
Ethernet**

**200 Hz**

**1 Gb/s & 10 Gb/s**

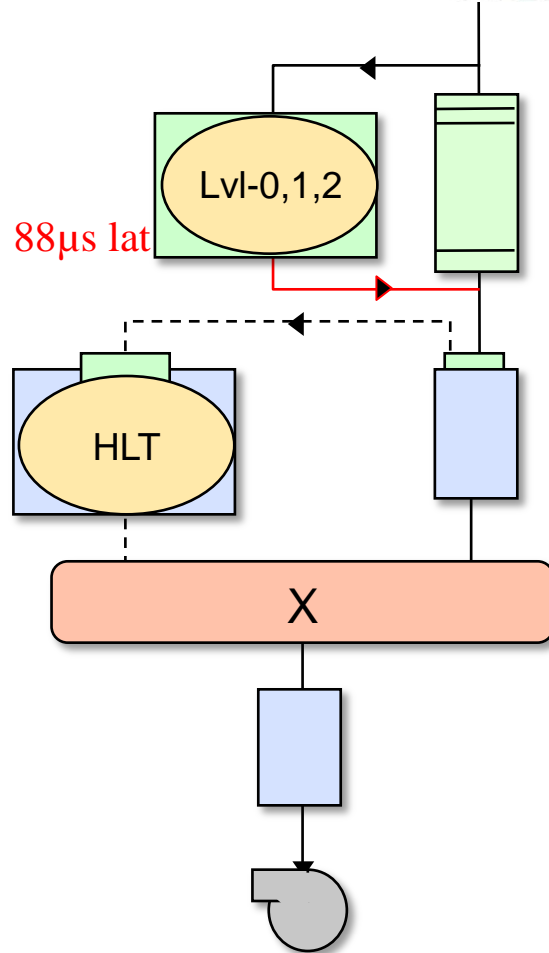
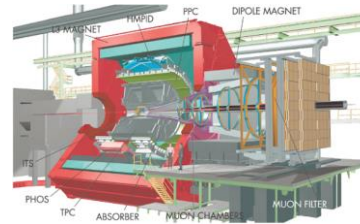
**Pull design:**

Both Level-2 and HLT request data from the readout buffers as needed



# Data Flow: ALICE

- custom hardware
- PC
- network switch



**front end pipeline**

**500 Hz**  
readout link

**Event Size: 70 MB**

**readout buffer**

**100 Hz**  
**event builder**

**Bandwidth: 2GB/s**

**Technology:**  
**Ethernet**

**1 Gb/s & 10 Gb/s**

**event buffer**

**100 Hz**

**HLT performs data compression by factor 4-5**

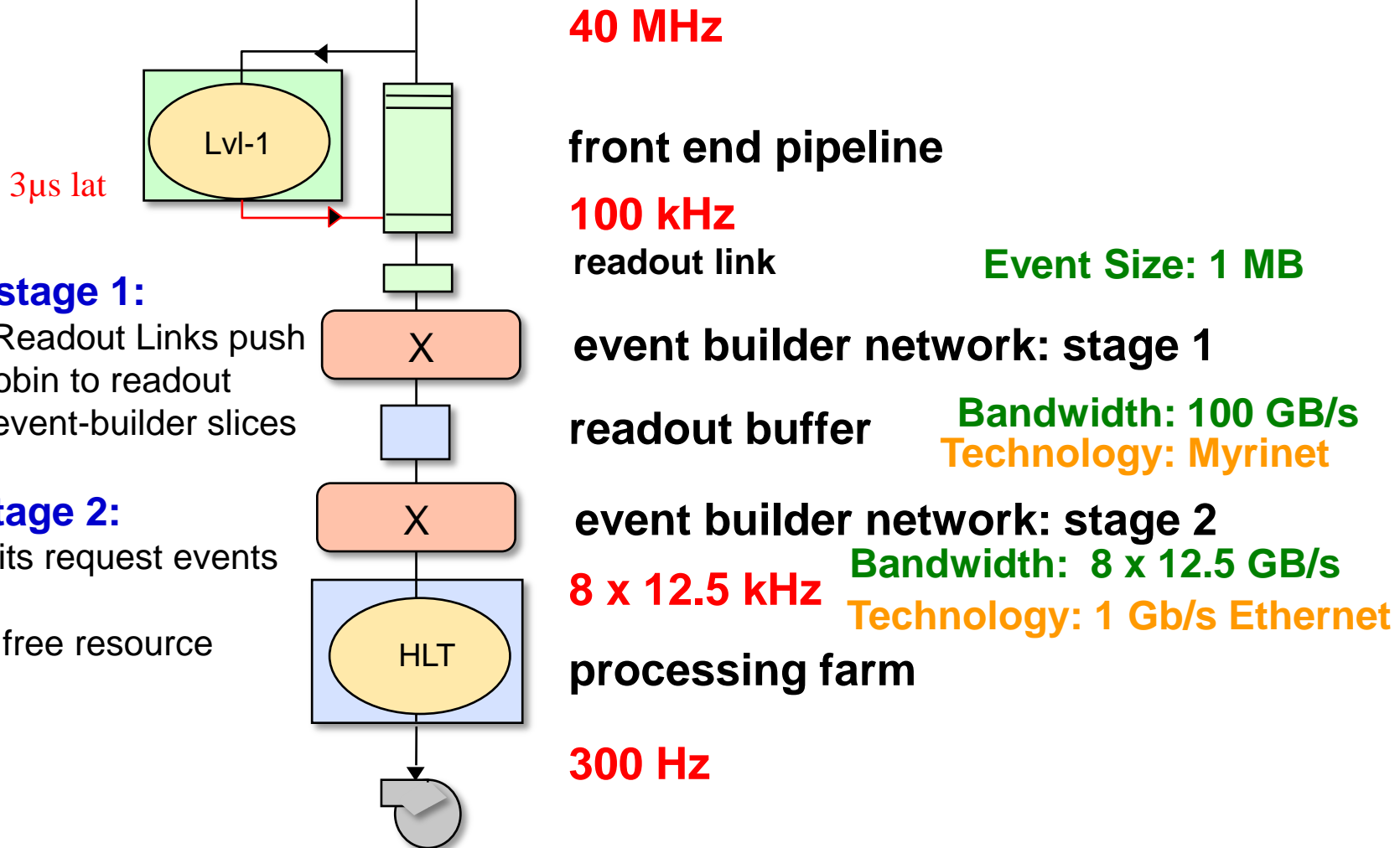
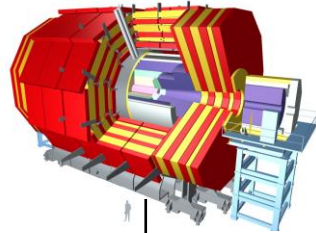
**HLT farm**

**Push design:**  
Push over TCP/IP



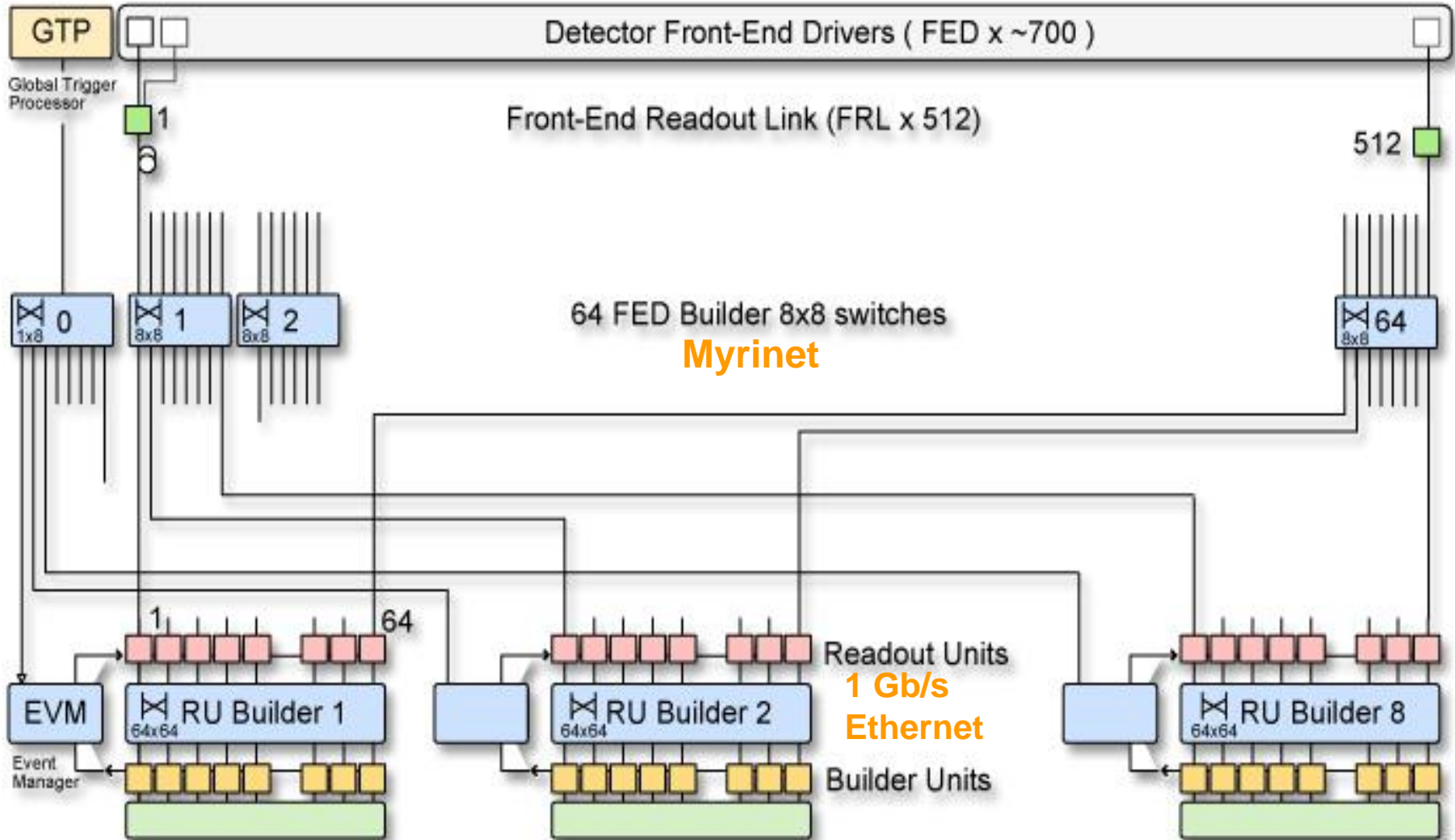
# Data Flow: CMS

- custom hardware
- PC
- network switch



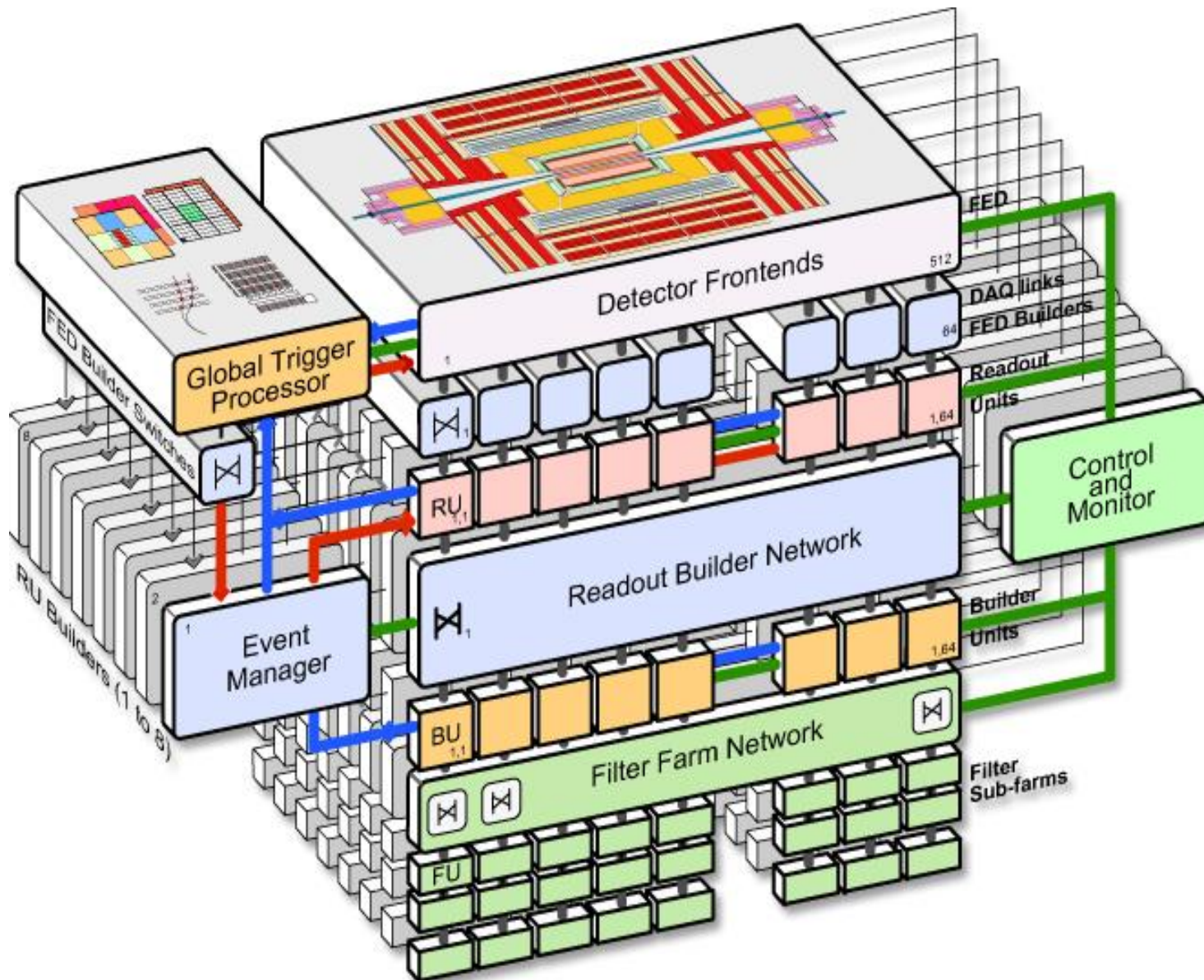


# EVB CMS: 2 stages



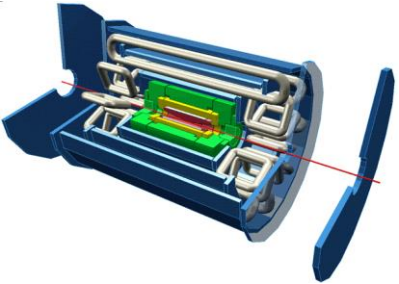
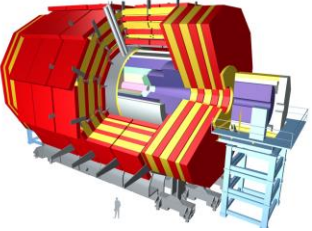
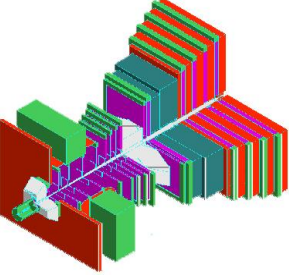
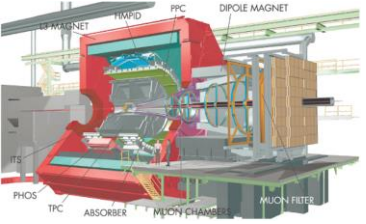


# CMS: 3D - EVB





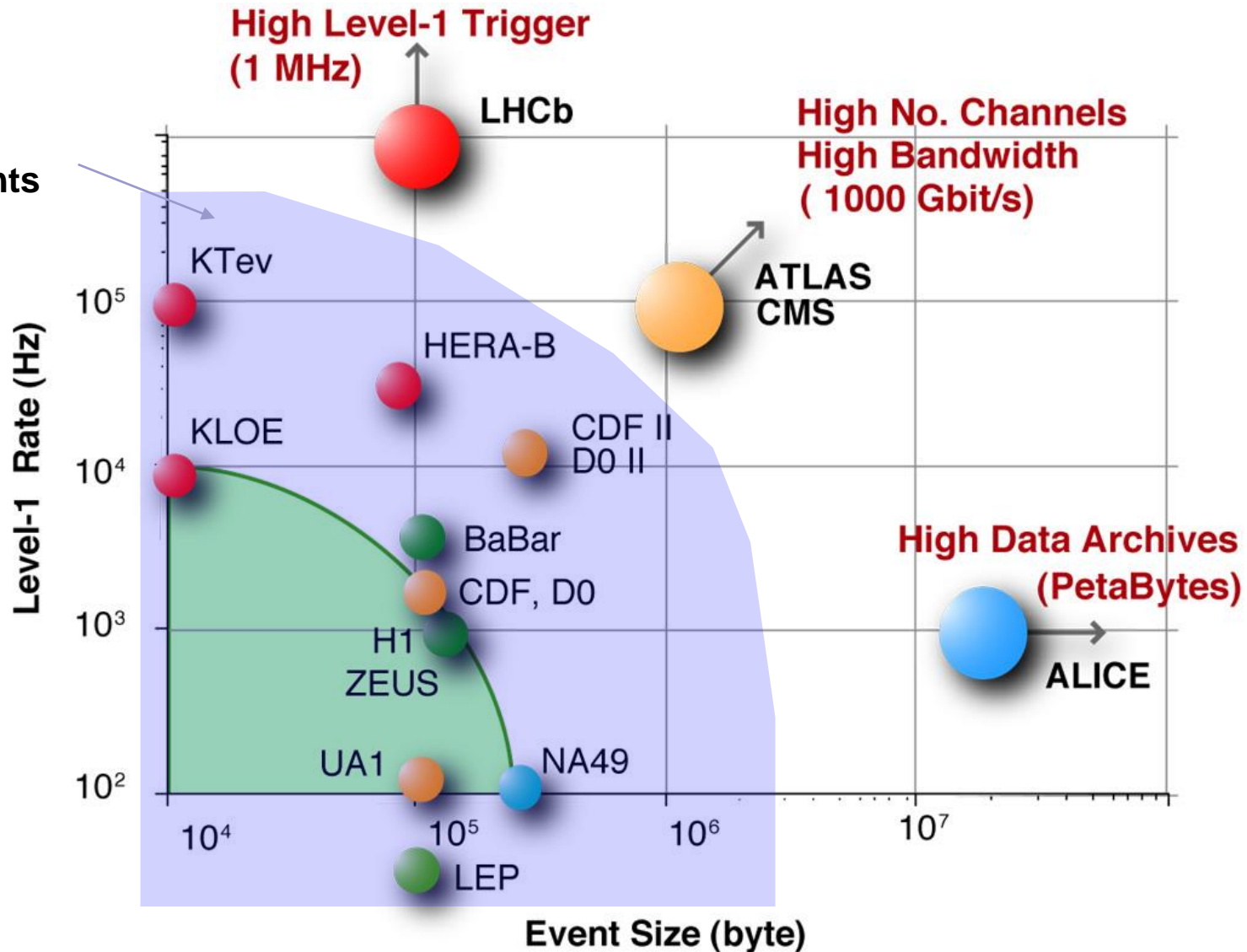
# Trigger/DAQ parameters (Run 1)

	No.Levels	Lvl 0,1,2	Event	Evt Build.	High Level Trigger HLT Out
	Trigger	Rate (Hz)	Size (Byte)	Bandw.(GB/s)	(design) MB/s (Event/s)
	<b>3</b>	LV-1 <b><math>10^5</math></b> LV-2 <b><math>6 \times 10^3</math></b>	<b>1.5 MB</b>	<b>5.25</b>	<b>300</b> (200)
	<b>2</b>	LV-1 <b><math>10^5</math></b>	<b>1.0 MB</b>	<b>100</b>	<b>300</b> (200)
			<b>Pb-Pb 1500MB/s</b>		
	<b>2</b>	LV-0 <b><math>10^6</math></b>	<b>30 kB</b>	<b>40</b>	<b>60</b> (2 kHz)
	<b>4</b>	Pb-Pb <b>500</b>	<b>70 MB</b>	<b>2</b>	<b>1250</b> (100)



# LHC experiments: Lvl 1 rate vs size

pre-LHC experiments





# 4) High-Level Trigger





# High-level trigger

- Reconstruct event data – often with full off-line software (minus calibration)
- To save time: multiple levels in software
  - E.g. First look only at muon and calorimeter data
  - Only if this looks interesting, spend time on reconstructing tracks in the trackers
- Issue of coupling of Online Software and Offline Software
- Multi-core machines
  - Run 1 or 2 copies of filter process per core
  - Processes forked from a mother-process to benefit from Copy-on-Write



# High-level trigger farm sizes

Number of..	Boxes	CPU cores	Filter procs	Logical Grouping
ALICE	~ 200	~ 5000 <sup>(1)</sup>	~ 3000	
ATLAS	~ 1600	~ 17000	1 per core <sup>(2)</sup>	49 Racks
CMS	~ 1600	~ 16000	~ 35000 <sup>(2)</sup>	O(20)BUs in 8 slices
LHCb	~ 1600	~ 16000	~ 30000 <sup>(2)</sup>	57 Racks

(1) 2300 CPU cores + 54 FPGA + 64 GPU cards (estimated to 100-200% of the CPU)

(2) Overcommitment if hyper-threading is supported by worker node



# 5) Storage to disk

# Local mass-storage

- Task
  - store a few days worth of data on-site.
  - handle peak recording rates (higher than network bandwidth to outside)
  - keep data safe
- Implementation
  - Disk arrays (Storage Area Network)
  - Mostly fiber channel (ATLAS: direct)

	ATLAS	CMS	ALICE	LHCb
GB/s max	1.6	0.8 pp 2 HI	2.2	0.25
space	160 TB	225 TB	610 TB	120 TB
raid	RAID-5	RAID-6	RAID-6	RAID-6
# writers	9	16	80+	6

## ALICE Storage System





# Run Control and Monitoring



# Run control tasks

- Start, configure and control  $O(10000)$  processes on farms of  $O(1000)$  machines
- Parallelize or execute in the right order according to dependencies
- Configure and monitor frontend electronics (millions of channels)
- Database access
- Graphical user interface to be used by non-expert shifters (integrated guidance)
- Automation

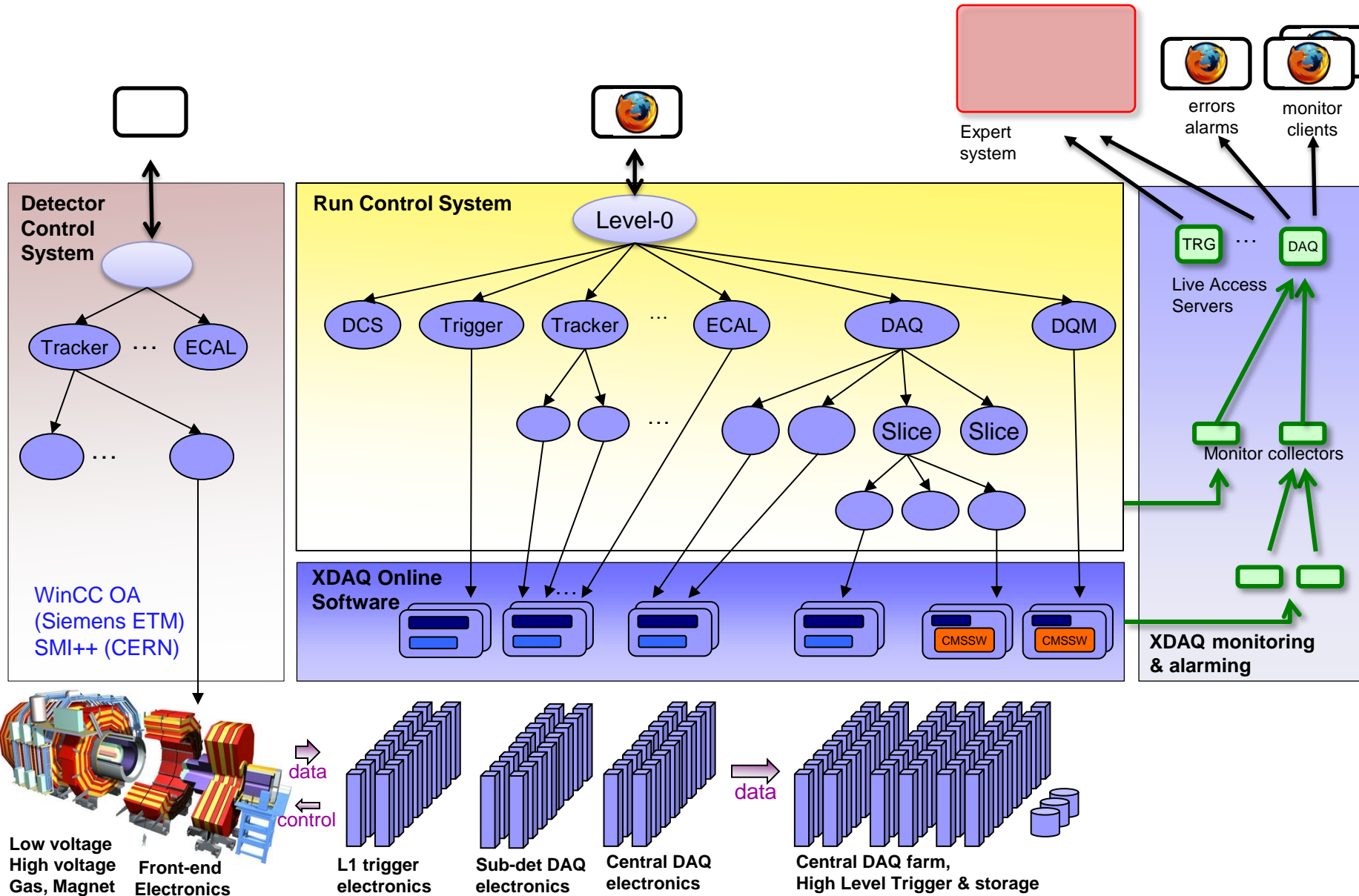


# Technologies used

- Communication
  - CORBA (ATLAS)
  - HTTP/SOAP (CMS)
  - DIM (LHCb, ALICE)
- Behavior and Automation
  - SMI++ (Alice, LHCb)
  - CLIPS (ATLAS)
  - RCMS Web Applications / Java (CMS)
- Process control
  - Based on XDAQ online software, CORBA
  - FMC/PVSS



# Run Control / Monitoring / Online software (CMS)







# Run Control GUI (LHCb)

Vision\_1: fwFSM\FSMOperat...
LHCb: TOP

**System** LHCb

**State** RUNNING

**Auto Pilot** ON

Thu 07-Feb-2013 09:52:49

root

Sub-System	State
DCS	READY
DAI	READY
DAQ	RUNNING
RunInfo	RUNNING
TFC	RUNNING
HLT	RUNNING
Storage	RUNNING
Monitoring	RUNNING
Reconstruction	RUNNING
Calibration	RUNNING
HV	READY

**Run Info**

Run Number: 136863

Run Start Time: 07-Feb-2013 09:37:45

Run Duration: 000:15:03

Nr. Events: 15502054

Step Nr: 0 To Go: 0

L0 Rate: 17192.04 Hz

HLT Rate: 5584.61 Hz

Dead Time: 0.78 %

Overflow: 0.00 %

Data Destination: Offline Data Type: IONPROTON13 Automatic

File: /daqarea/lhcb/data/2013/RAW/FULL/LHCb/IONPROTON13/136863 Run DB

**Deferred HLT Info**

LHCb\_Deferred NOT\_ALLOCATED

Runs/Files: 0 / 0

Processing: 0%

Disk Usage: 29%

Efficiency Trigger Rates TFC Control TELL1s LHCb Elog

**Sub-Detectors:**

TDET	VELOA	VELOC	TT	IT	OTA	OTC	RICH1	RICH2	PRS
HOT_READY	RUNNING	RUNNING	RUNNING	RUNNING	RUNNING	RUNNING	RUNNING	RUNNING	RUNNING

ECAL	HCAL	MUONA	MUONC
RUNNING	RUNNING	RUNNING	RUNNING

**Trigger Components:**

L0DU	TCALO	TMUA	TMUC	TPU
RUNNING	RUNNING	RUNNING	RUNNING	RUNNING

**Messages:**

```

07-Feb-2013 09:03:41 - LHCb executing action START_RUN
07-Feb-2013 05:03:42 - LHCb in state ACTIVE
07-Feb-2013 05:03:42 - LHCb in state RUNNING
07-Feb-2013 05:28:21 - *** INFO - VELO Closed, Changing RUN...
07-Feb-2013 05:28:21 - LHCb executing action CHANGE_RUN
                    
```

Close

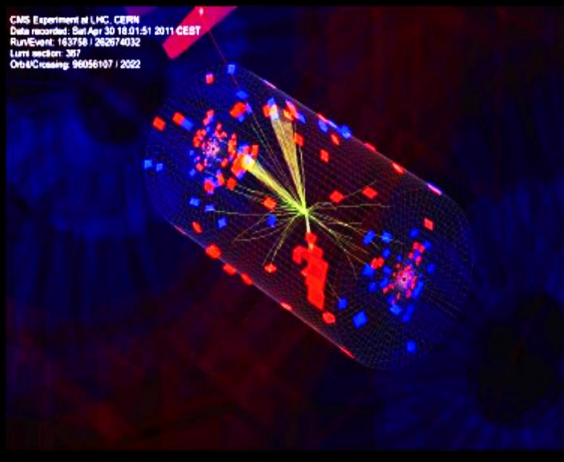
View All Owners



# Monitoring (CMS)

CMS 30/04/11 PROTON PHYSICS DAQ state Running Run Number 163758 **Lv1 rate 30.929 kHz** Ev. <Size> kB 401.0 [224.8] DeadTime(AB) 0.673 % Acc. Hz (%) 30991.1 (100.0%) **HLT <CPU> 21.96 %**

CMS Experiment at LHC, CERN  
 Date recorded: Sat Apr 30 18:04:31 2011 CEST  
 Run/Ev#s: 163758 / 262074032  
 Luma section: 357  
 Data Crossing: 90656107 / 2022



### Data to Surface

Sub-System	State	FRL	FEB	IN
TRG	Running	3	3	3
CSC	Running	9	9	9
DAQ	Running	0	0	0
DQM	Running	0	0	0
DT	Running	6	6	6
ECAL	Running	54	54	54
ES	Running	39	39	39
HCAL	Running	26	26	26
HFLUMI	Running	6	6	6
PIXEL	Running	40	40	40
RPC	Running	3	3	3
SCAL	Running	1	1	1
TRACKER	Running	250	438	438
CASTOR	Running	3	3	3

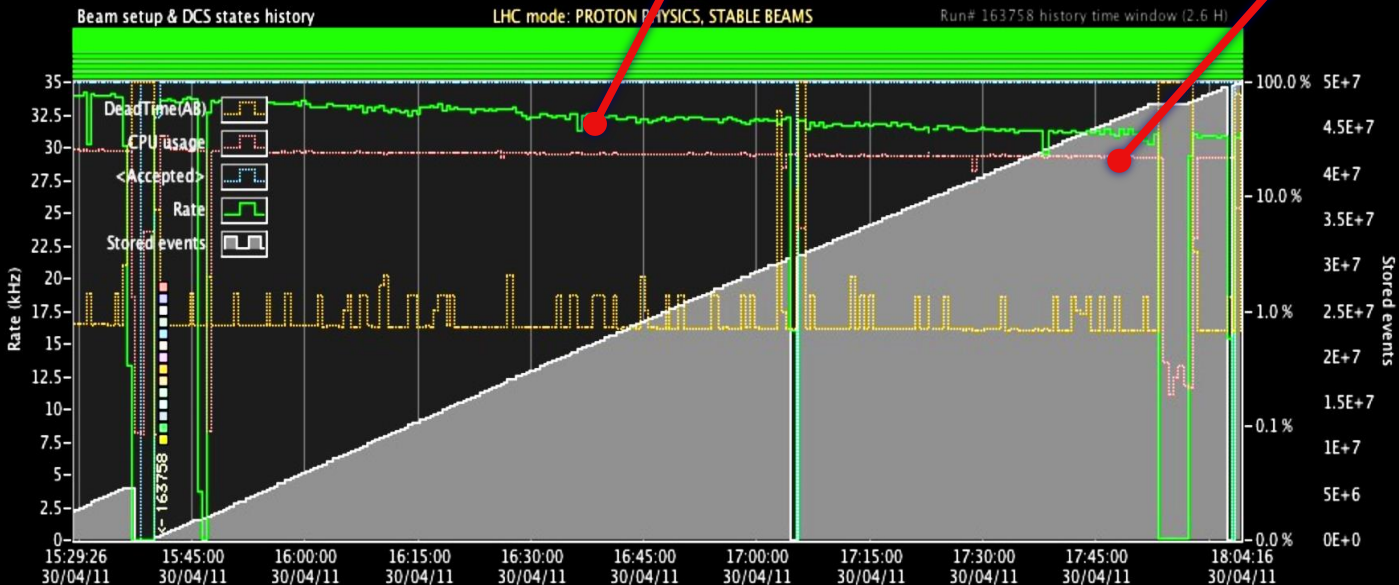
### SM streams

Stream	No.Events	Rate (Hz)	BnW (MB/s)
NanoDST	25.985E+6	3003.71	5.79
ALCAPO	9.281E+6	1076.24	8.95
RPCMON	7.323E+6	847.27	11.55
A	2.720E+6	311.82	68.48
ALCAPHISYM	2.592E+6	313.51	1.36
Calibration	827.060E+3	100.40	9.14
EcalCalibrati	827.060E+3	100.40	2.91
Express	177.089E+3	19.51	3.86
HLTMON	24.976E+3	3.00	0.74
OnlineErrors	1.762E+3	0.27	0.08
FaultyEvents	0.000E+0	0.00	0.00
Error	0.000E+0	0.00	0.00

### Data Flow

/cdaq/physics/Run2011/Se32/v8.3/HLT/V4  
 #LS 374  
 PreScaleIndex 1  
 #L1(GT) 266915432  
**Lv1 Rate 30.929 kHz**  
 Pending Lv1 111887  
 #Frag. in RU 584  
 Min 382  
 BnW (MB/s) 1.2E+4  
 Events in BU 6  
 <Ev.> 0  
 Pending Req. 21701  
 <#P> 22  
 #Running FUs 7996  
 100.00%  
 A 318.8 Hz  
 BnW MB/s 115  
 EventRate Hz 5811.1  
 Stored 50002951  
 Time to fill disk 3 of srv-c2c06-14 > week  
 TIER0 TRANSFER ON

Random ON  
 Physics ON  
 CalibCyc ON  
 55 FEDCRC  
 FBI occ. %  
 Max 2  
 Min 0  
 FBO occ. %  
 Max 13  
 Min 0  
 EvSize (kB) 401  
 Rcv.-Disc. 2160  
 53179 P.M-m  
 53183 A.M-m  
 <FU-CPU> 22 %  
 100  
 0  
 <SM-CPU> 13.2 %  
 Free space TB 227.7



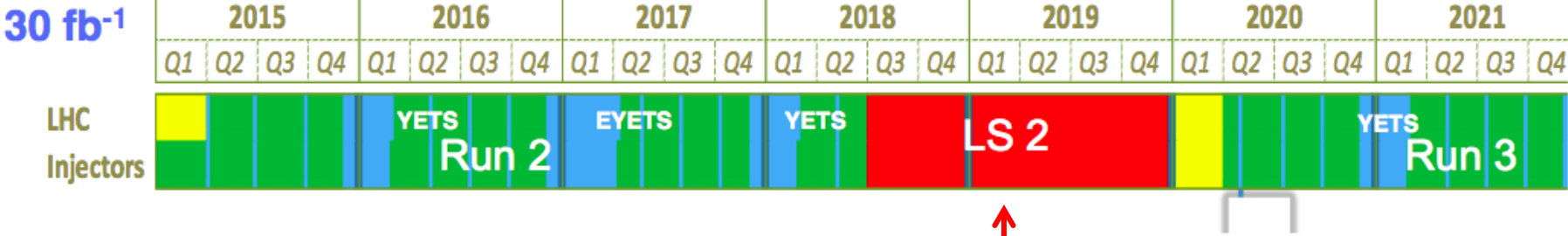
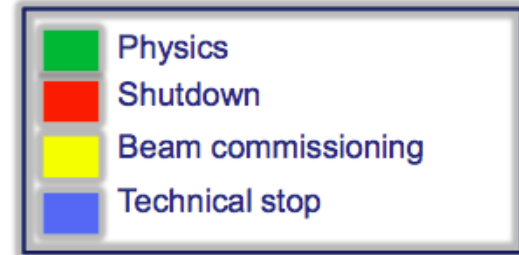


# **Upgrades. around the corner and further away**



# LHC plans

**Run-2** : 13 TeV center-of-mass energy  
 Targeting 40 fb<sup>-1</sup> / year  
 higher pile-up



↑  
**Switch over to Linac-4**



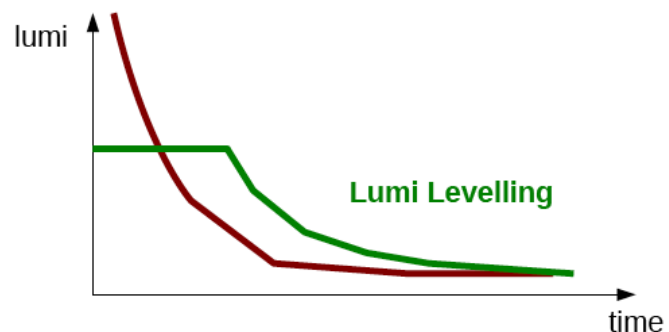
(Extended) Year End Technical Stop: (E)YETS

**LS3 : HL-LHC installation**

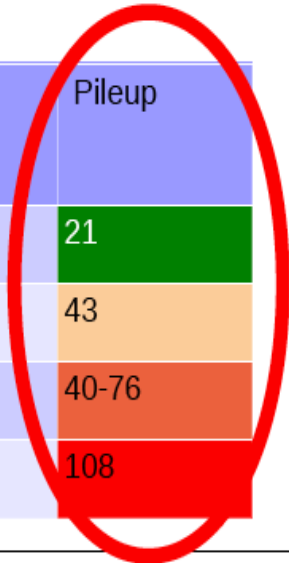


# Upgrade of the LHC

- Event sizes increase with the number of underlying events
  - LHC experiments were designed for a pile-up of 24 (with some margin)
  - In 2012 we had a pile-up of > 30
  - After LS1 we must be ready for a pile-up of 50
    - LHC may do “lumi-levelling” to keep pile-up below 50
    - Beams separated at the start of the fill

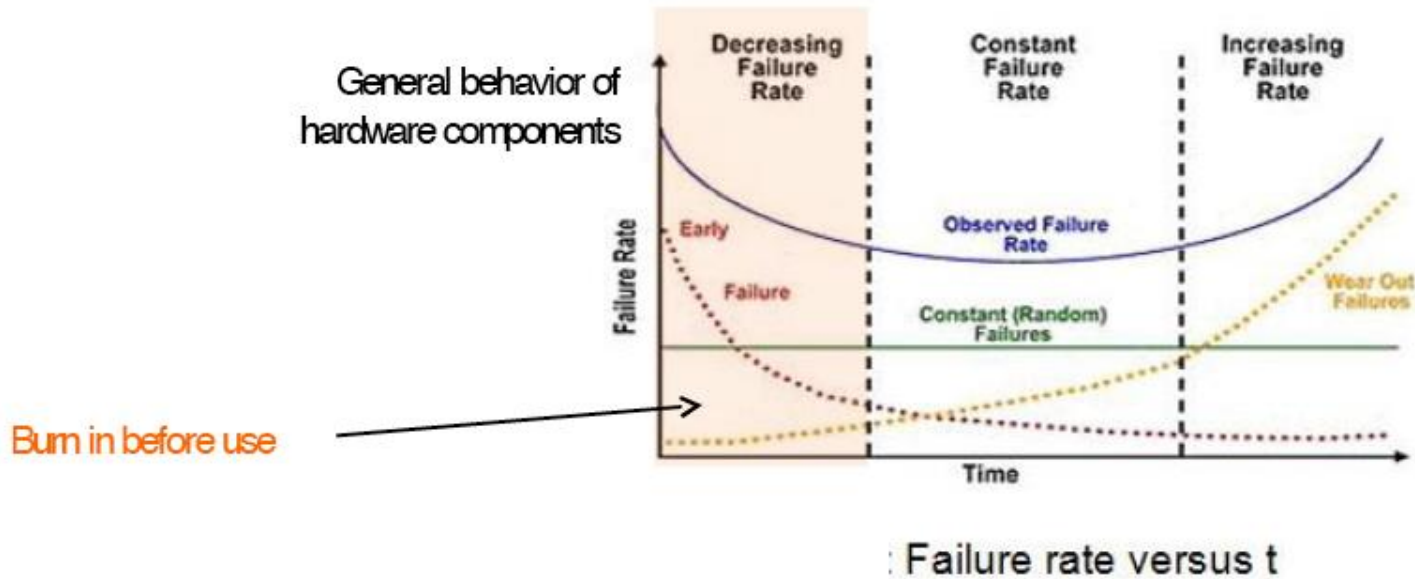


BX spacing [ns]	Beam current [ $\times 10^{11}$ e]	Emittance [ $\mu\text{m}$ ]	Peak Lumi [ $\times 10^{34} \text{cm}^{-2} \text{s}^{-1}$ ]	Pileup
25	1.15	3.5	0.92	21
25	1.15	1.9	1.6	43
50	1.6	2.3	0.9-1.7	40-76
50	1.6	1.6	2.2	108



Plan for startup after LS1 (2015) →

# Hardware life-time ... another reason to upgrade



- System reliability decreases after a few years
  - It makes sense to replace PCs and networking equipment every ~5 years
  - Can then also benefit from advances in technology  
(Get about a factor 10 higher performance according to Moore's Law)
- Custom hardware is usually kept much longer
  - At some point it also starts becoming unreliable ...

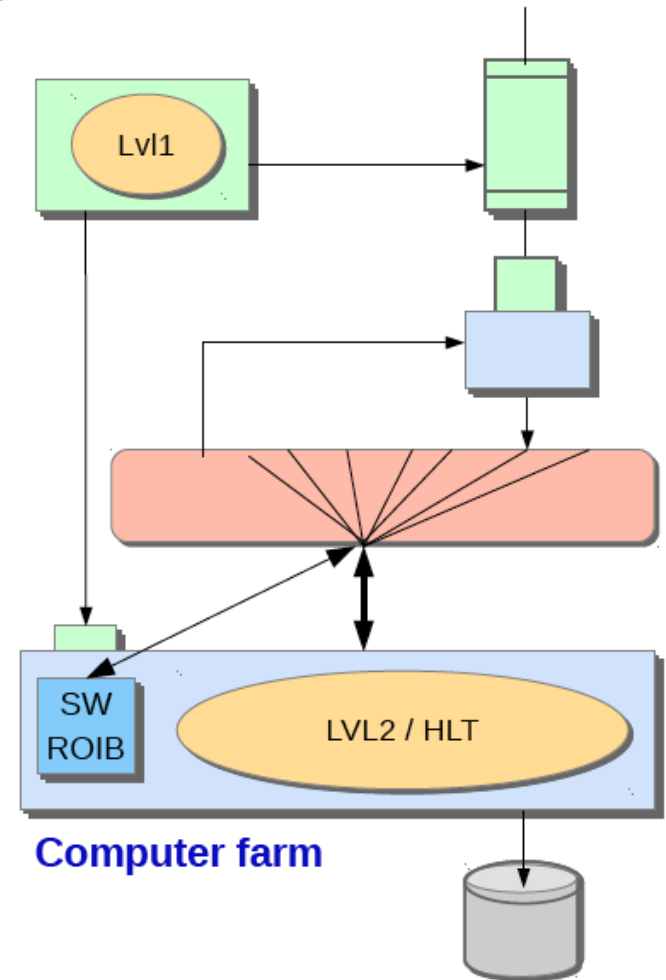
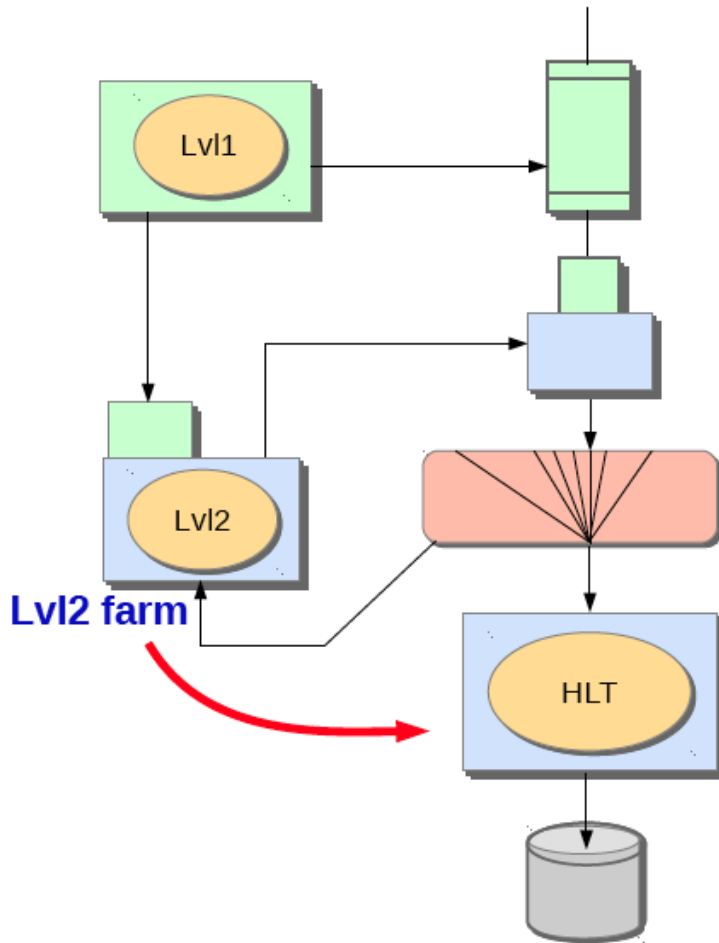


# Upgrade highlights of the LHC experiments



# ATLAS DAQ upgrade (for Run2)

NOW | FUTURE



One integrated farm for Level-2 and High-Level Trigger.  
 Region-of-Interest principle is kept, but ROIs are built in software.

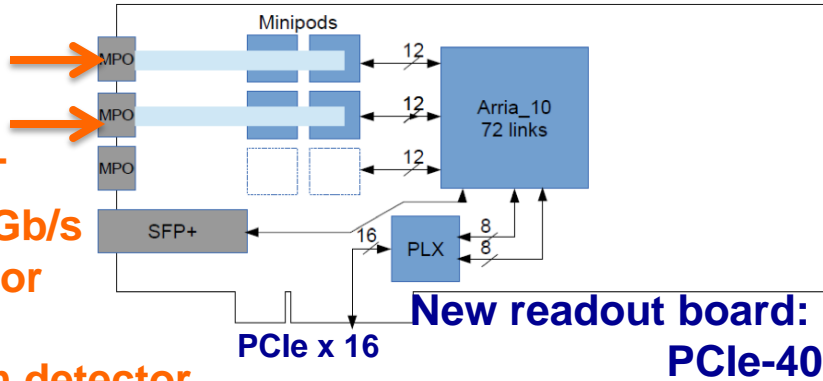




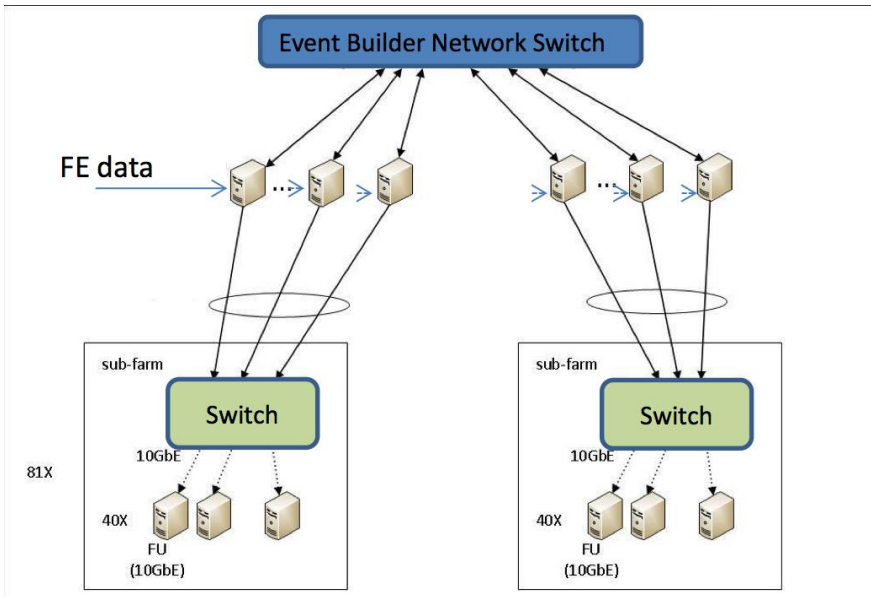
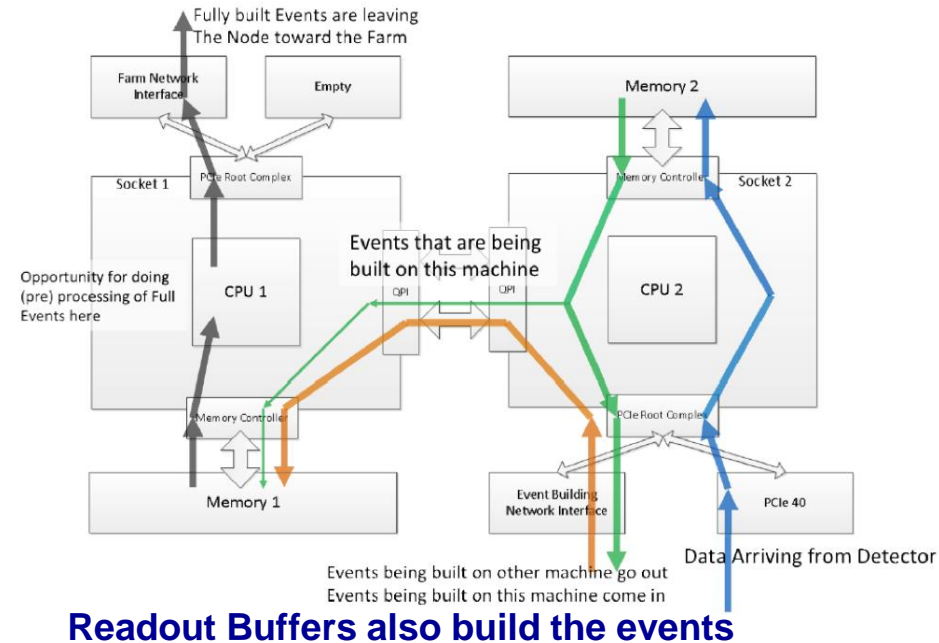


# LHCb to be read out at up to 40 MHz (for Run 3)?

**GBT**  
**4.8 Gb/s**  
**x24 or**  
**X36**  
**from detector**



## Investigating DAQ in a container



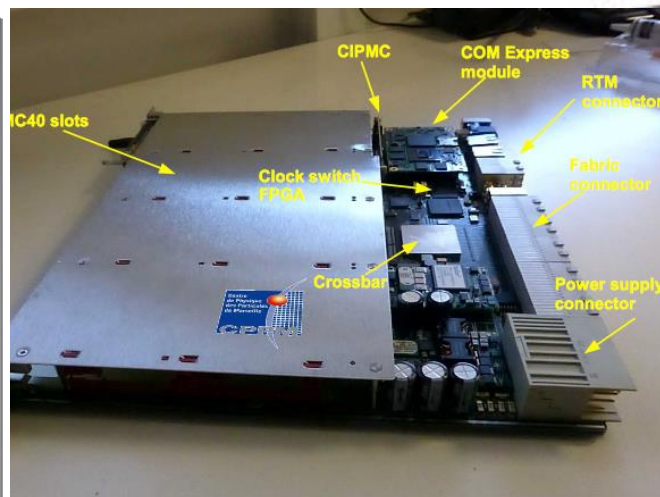
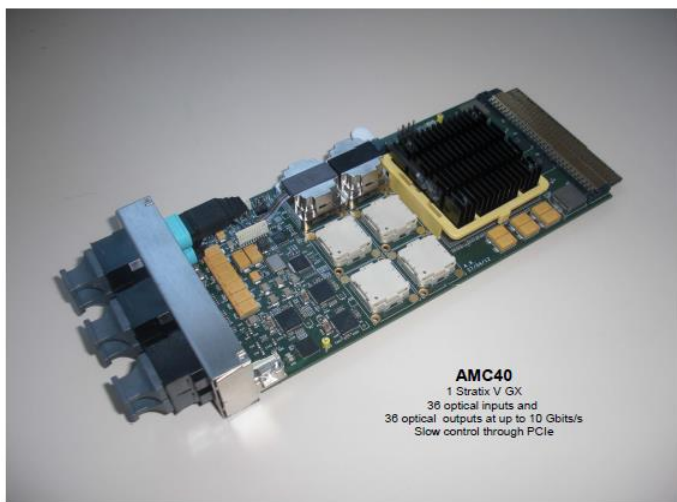
**Readout Buffers connected to event-builder and Filter-farm networks**



# ALICE (for Run3)

- Read-out all Pb-Pb interactions at 50 kHz (= interaction rate). 1 TB/s event building.
- Continuous read-out of Time-Projection-Chamber
- No on-line filtering. Just on-line data reconstruction ( clustering )
- Common cluster for Online + Offline - estimated 250000 cores

## Investigating a Common Readout Unit based on Advanced TCA: (also considering PCI-e)



**Tell-40 AMC card**  
(originally developed by LHCb)  
**in: GBT x 24**  
**Out: DDL3xn (commercial 10 Gbs)**

**ATCA carrier**  
carrying 4 Tell-40

**10 ATCA shelves would be needed**



# Summary

- DAQ is an interesting part of a HEP experiment
  - No DAQ system alike
  
- Exciting times ahead for the LHC DAQ systems
  - Upgrade projects in all experiments
  
- If you want to learn more ...
  - ISOTDAQ school
    - Feb 2013, Budapest  
<http://indico.cern.ch/event/274473/other-view?view=standard>
    - Feb 2014, Rio (?)
  - DAQ@LHC workshop (March 2013)
    - <http://indico.cern.ch/event/217480/other-view?view=standard>



# Backup



# Example of network performance

## CMS Stage1 FED-Builder

### ■ FED Builder functionality

- Receives event fragments from approx. 8 Readout Links (FRLs).
- FRL fragments are merged into “super-fragments” at the destination (Readout Unit).

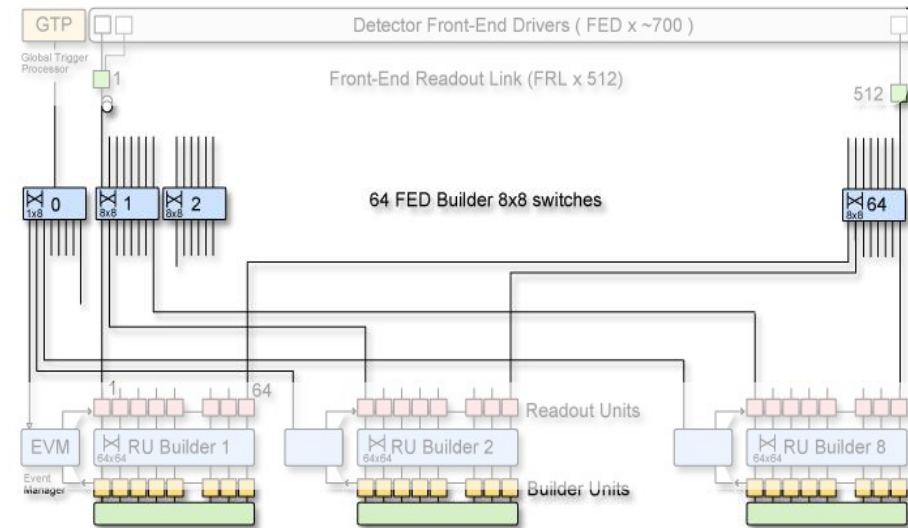
### ■ FED Builder implementation

#### □ Requirements:

- Sustained throughput of 200MB/s for every data source (500 in total).
- Input interfaces to FPGA (in FRL) -> protocol must be simple.

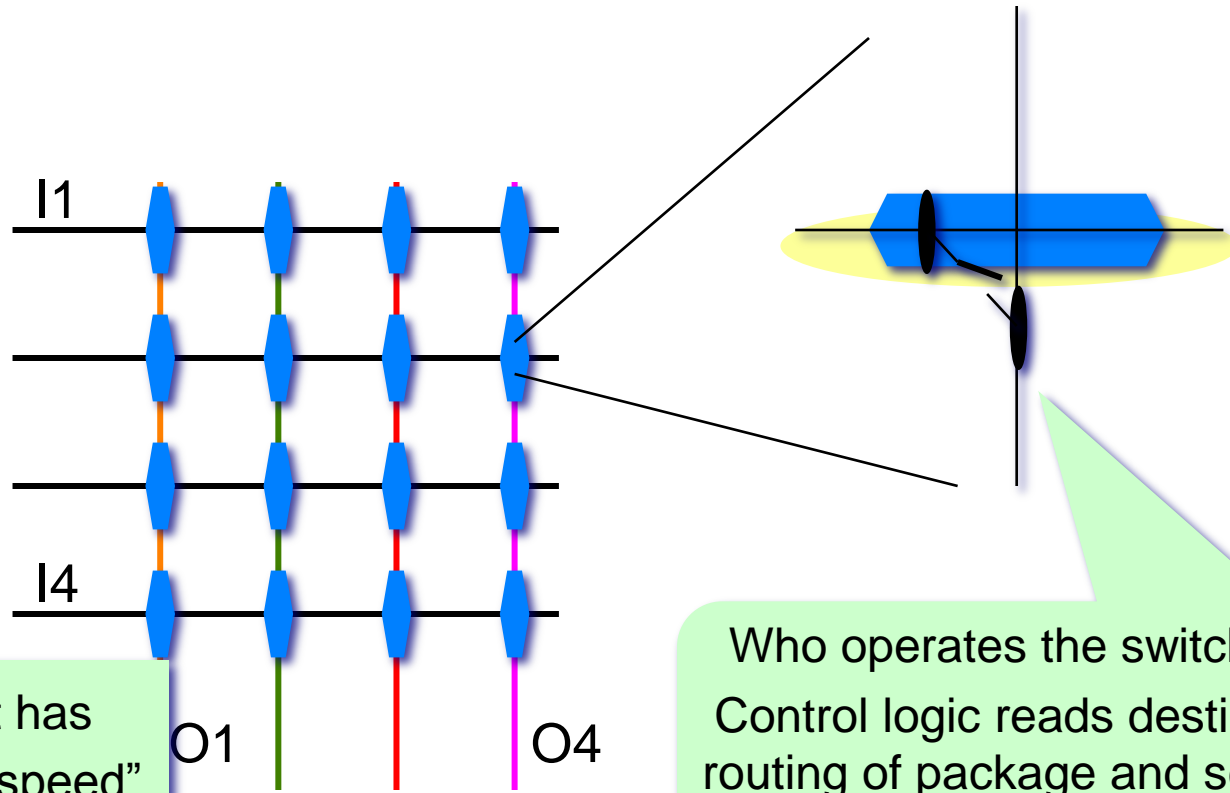
#### □ Chosen network technology: Myrinet

- NICs (Network Interface Cards) with 2x2.5 Gb/s optical links ( $\approx 2 \times 250$  MB/s)
- Full duplex with flow control (no packet loss).
- NIC cards contain RISC processor. Development system available.  
Can be easily interfaced to FPGAs (custom electronics: receiving part of readout links)
- Switches based on cross bars (predictable, understandable behavior).
- Low cost!





# Switch implementation: crossbar

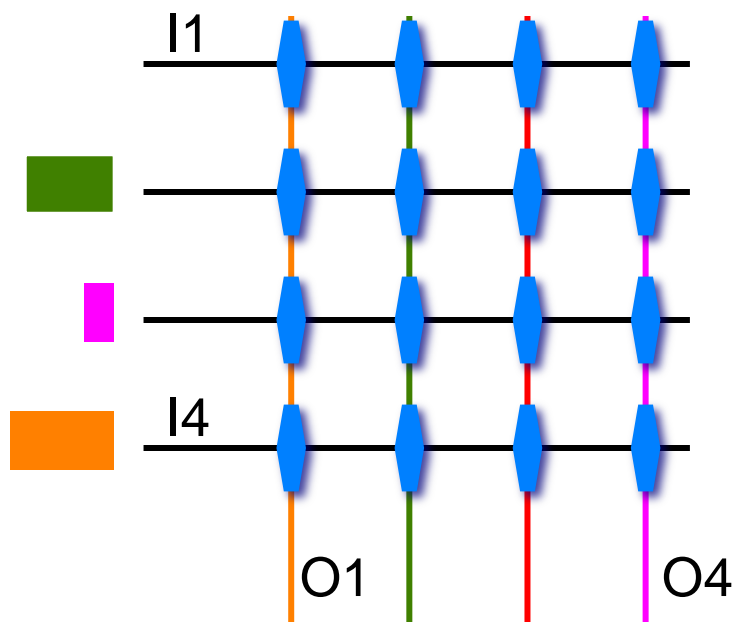


Every input / output has  
 A given max. “wire-speed”  
 (e.g. 2Gbits/sec).  
 Internal connections are  
 much faster!

Who operates the switches ?  
 Control logic reads destination  
 routing of package and sets the  
 switches appropriately.



# Switch implementation: crossbar



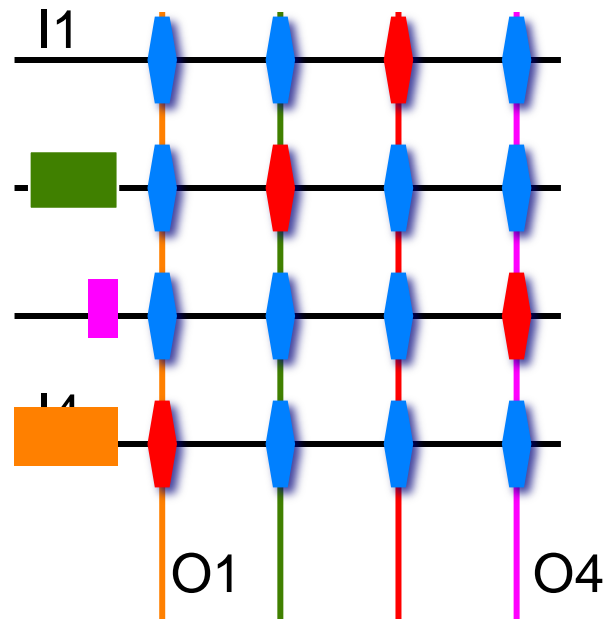
Paradise scenario:

All inputs want to send data to different destinations





# Switch implementation: crossbar

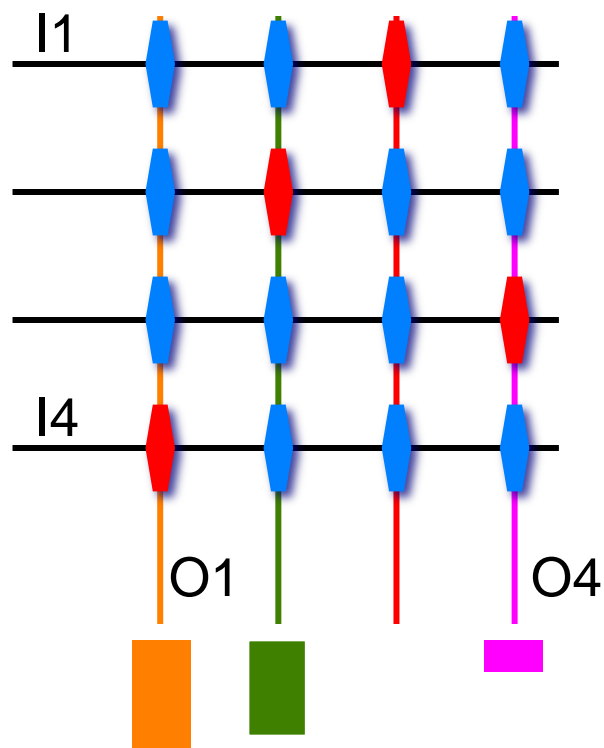


Paradise scenario:

No congestion, since every data package finds a free path through the switch.



# Switch implementation: crossbar



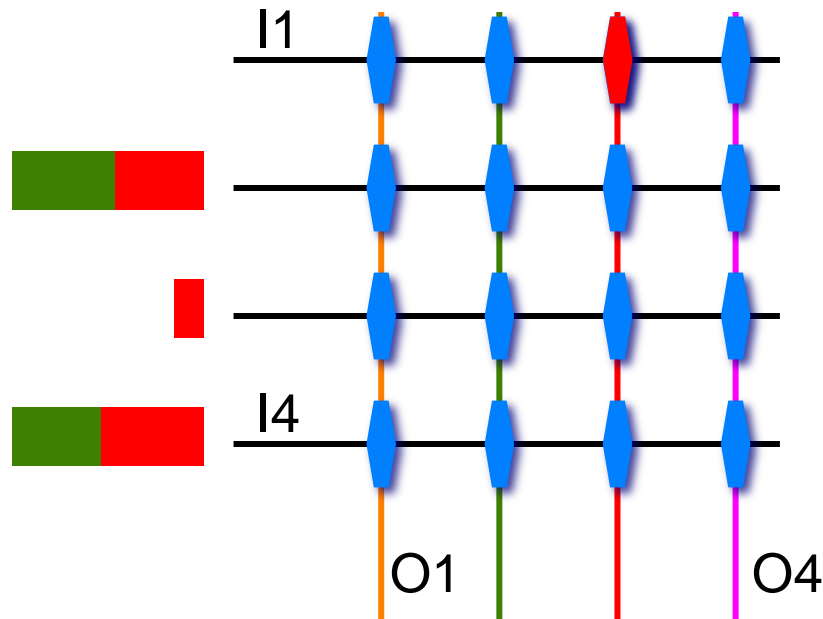
Paradise scenario:

Data traffic performs with  
“wire speed” of switch



# Switch implementation

Crossbar switch: **Congestion in EVB traffic**

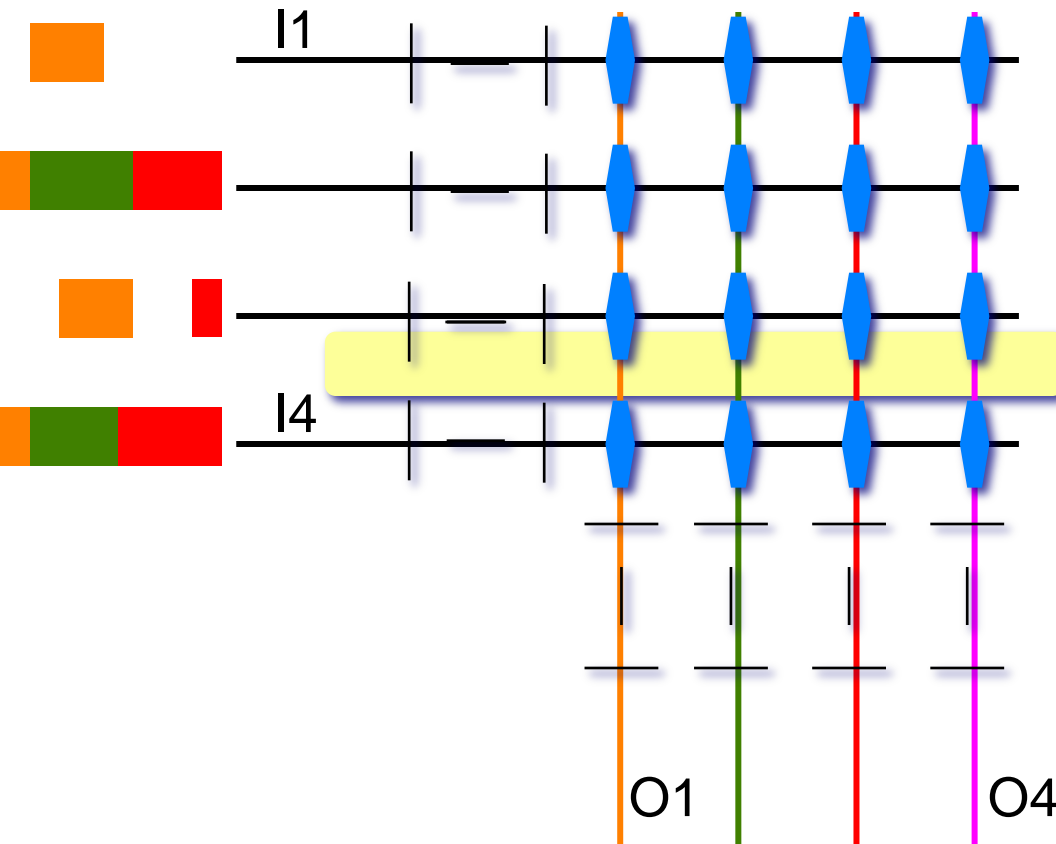


Only one packet at a time can be routed to the destination.  
"Head of line" blocking



# Switch implementation

Crossbar switch: **Improvement : additional input FIFOs**

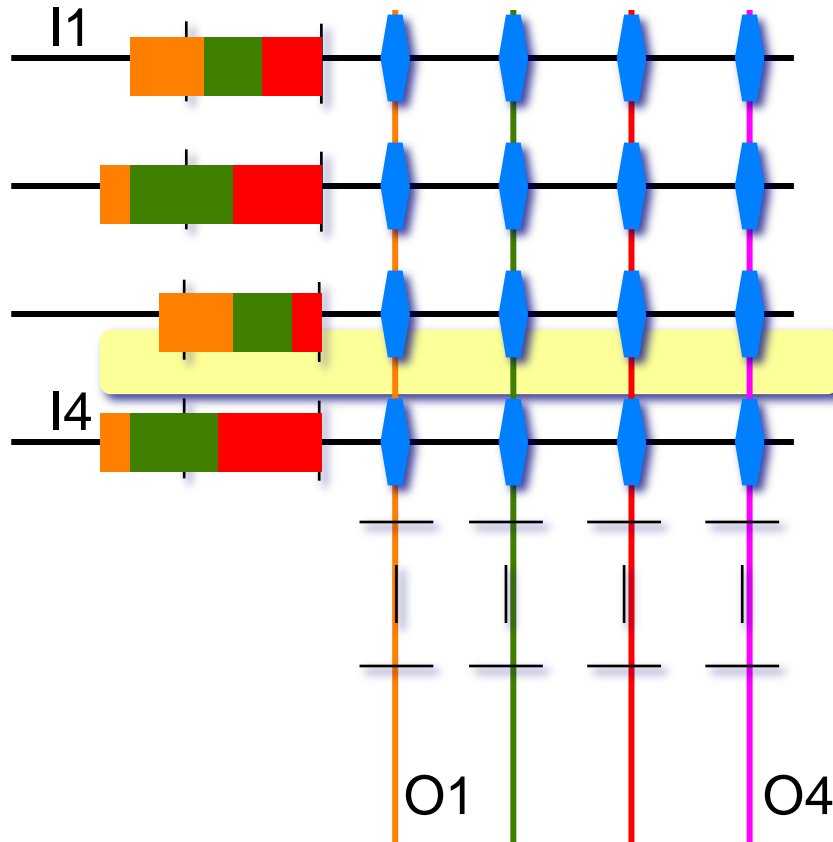


Fifos can “absorb” congestion ...until they are full.



# Switch implementation

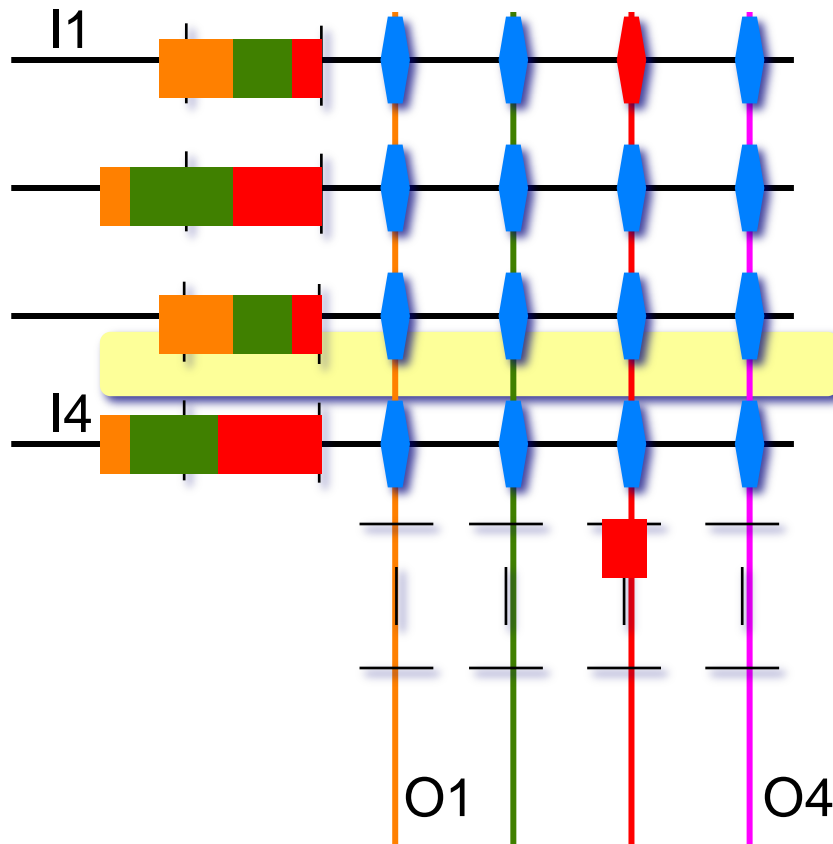
Crossbar switch: **Improvement : additional input FIFOs**





# Switch implementation

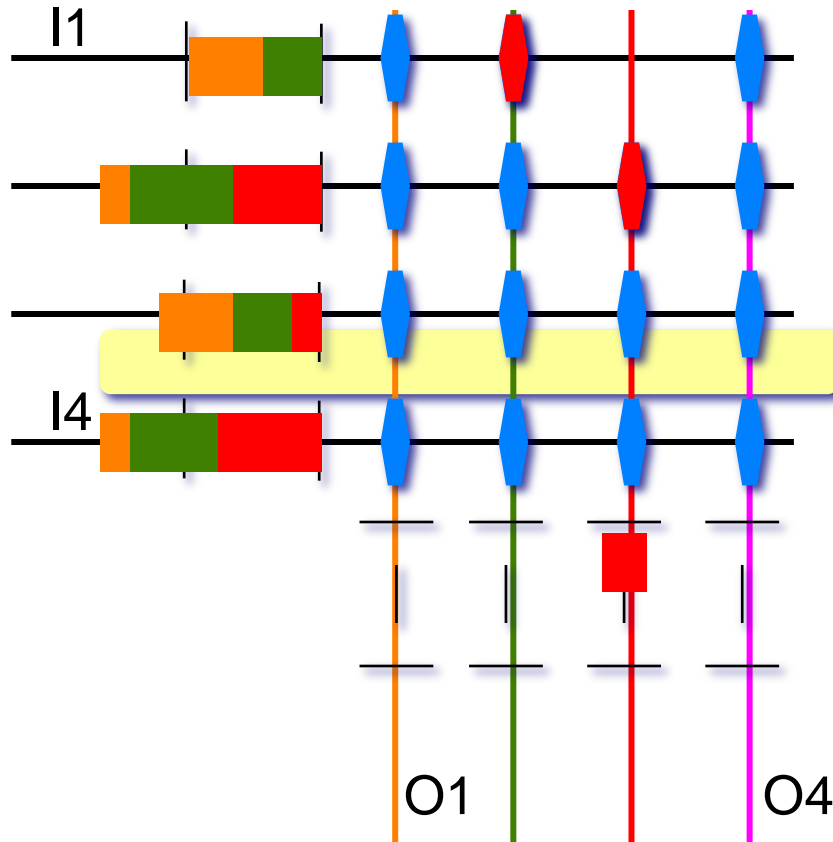
Crossbar switch: **Improvement : additional input FIFOs**





# Switch implementation

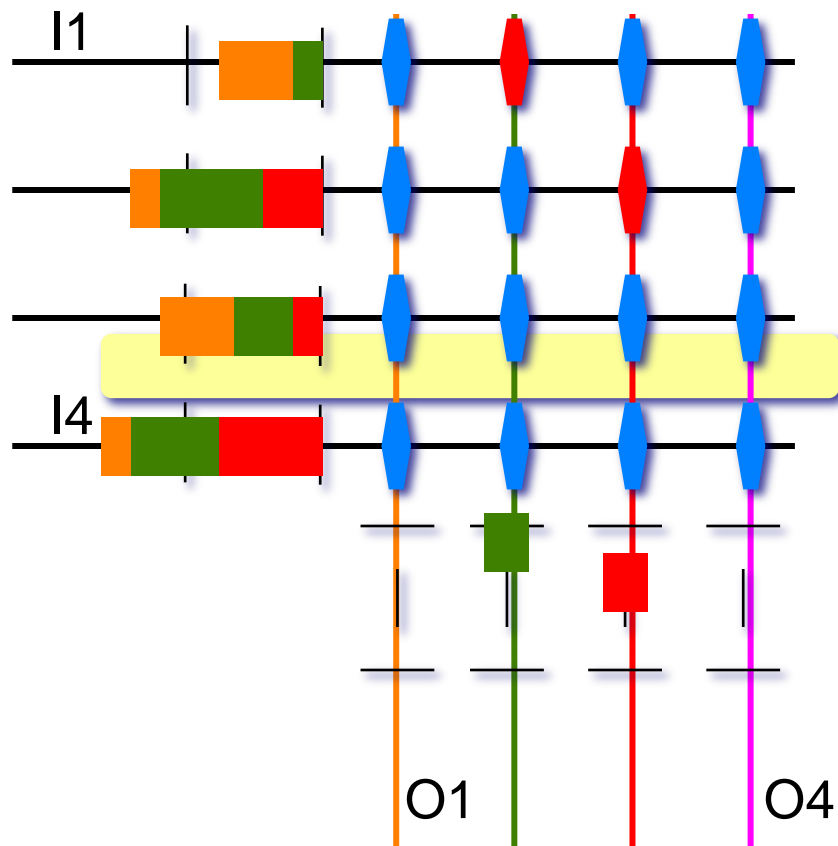
Crossbar switch: **Improvement : additional input FIFOs**





# Switch implementation

Crossbar switch: **Improvement : additional input FIFOs**

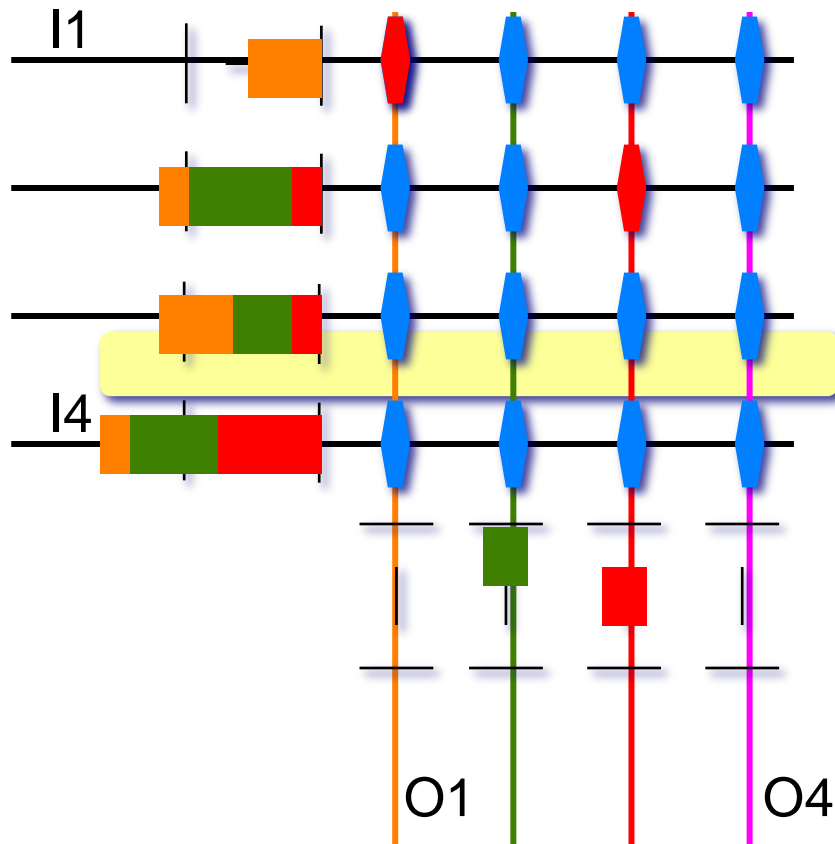






# Switch implementation

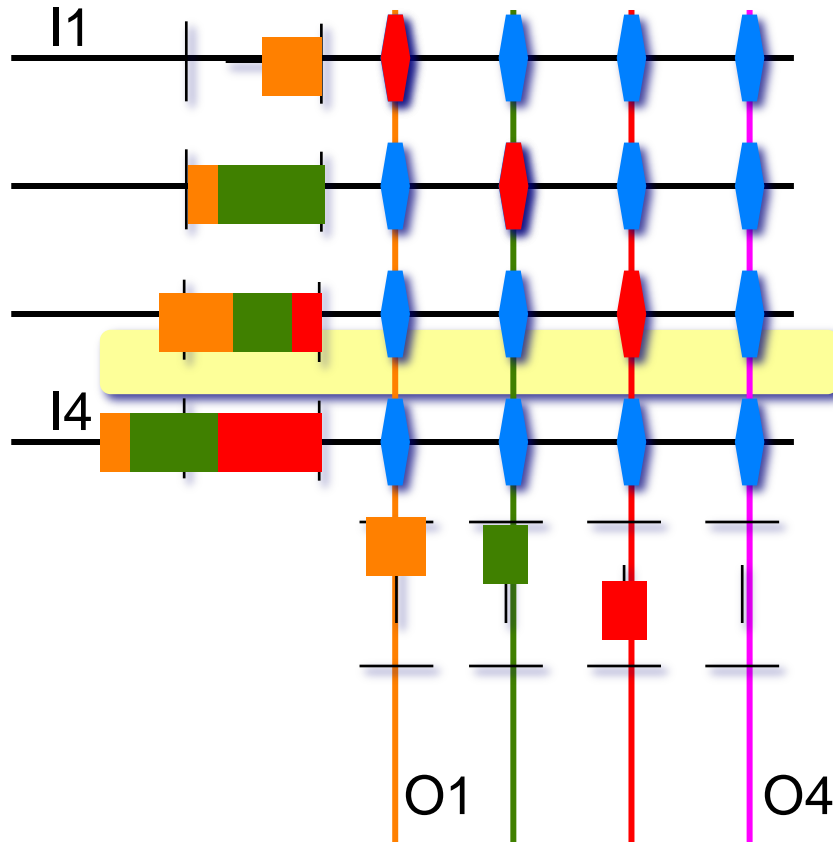
Crossbar switch: **Improvement : additional input FIFOs**





# Switch implementation

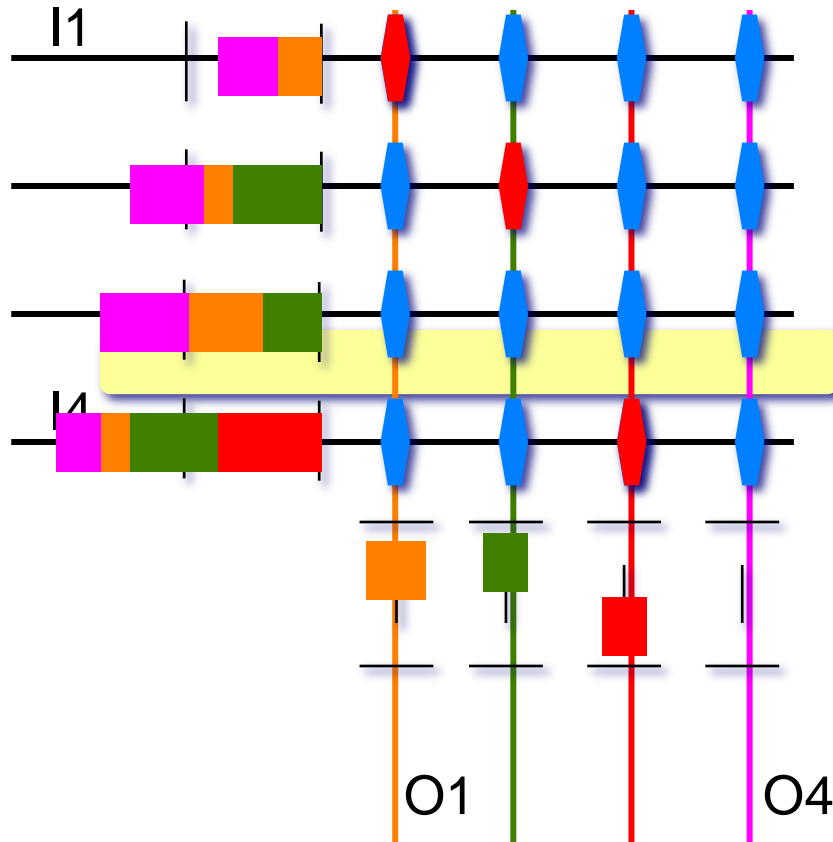
Crossbar switch: **Improvement : additional input FIFOs**





# Switch implementation

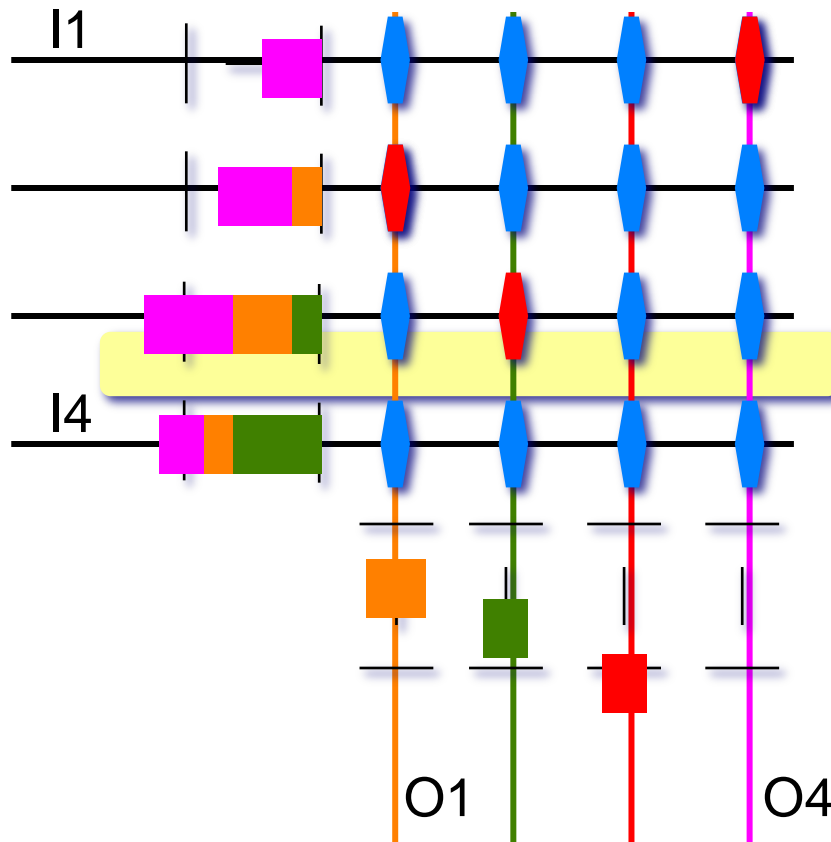
Crossbar switch: **Improvement : additional input FIFOs**





# Switch implementation

Crossbar switch: **Improvement : additional input FIFOs**

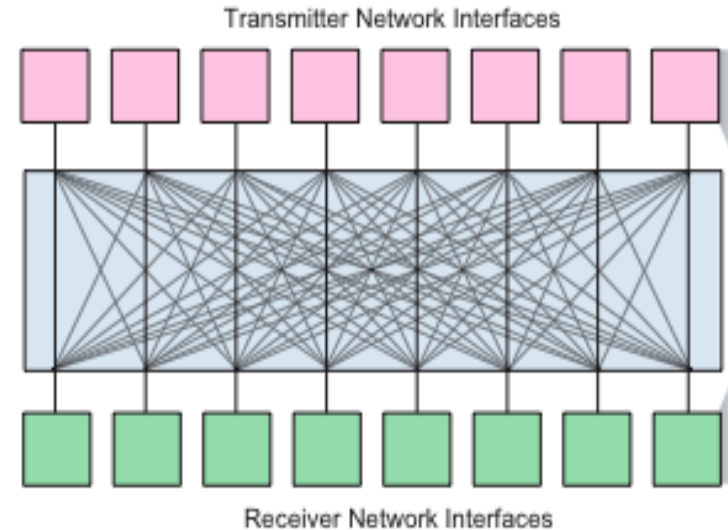


**Still problematic:**  
 Input Fifios can absorb data fluctuations until they are full.  
 How good it works depends on:  
 Fifos capacity  
 event size distribution  
 Internal speed of the switch  
 EVB traffic: blocking problem  
 remains to some extend



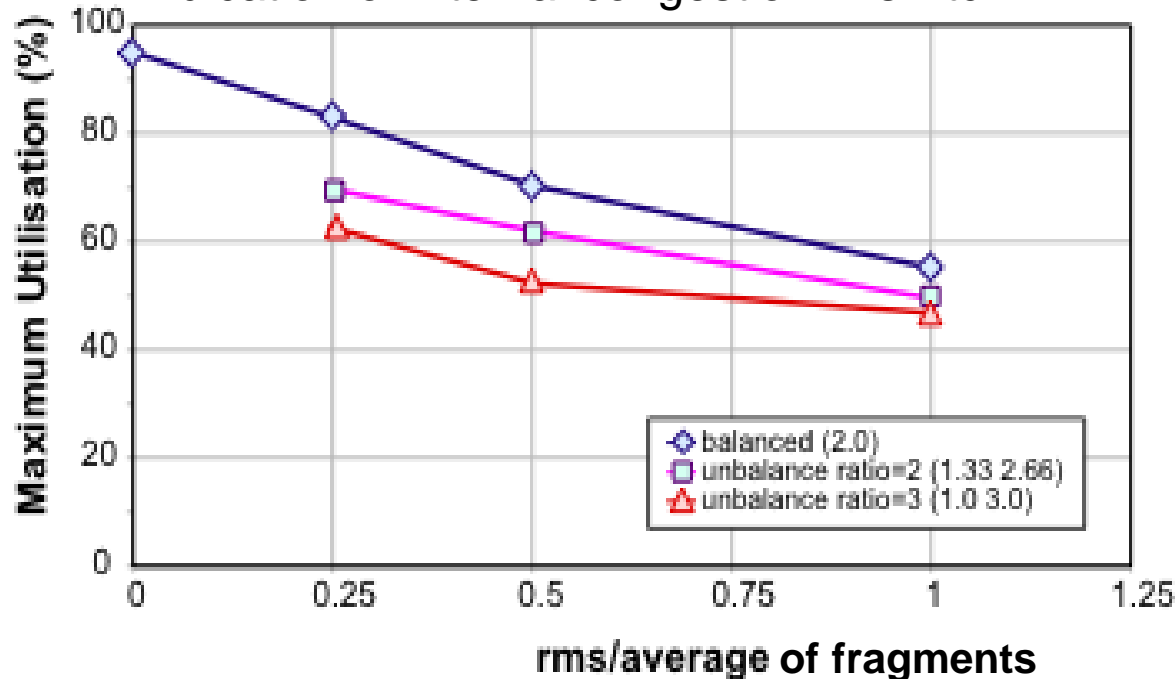
# Performance of “1 rail” FEDBuilder

Measurement configuration:  
8 sources to 8 destinations



% of wire-speed

Indication of internal congestion in switch:



Measured switch utilization:

**Blue:** all inputs 2 kB avg

**Magenta:** 4 x 1.33 kB

4 x 2.66 kB

**Red:** 4 x 1 kB

4 x 3 kB

**≈ 50 %**



# Conclusion: EVB traffic and switches

- EVB network traffic is particularly hard for switches
  - The traffic pattern is such that it leads to congestion in the switch.
  - The switch either “blocks” (= packets at input have to “wait”) or **throws away** data packets (Ethernet switches)

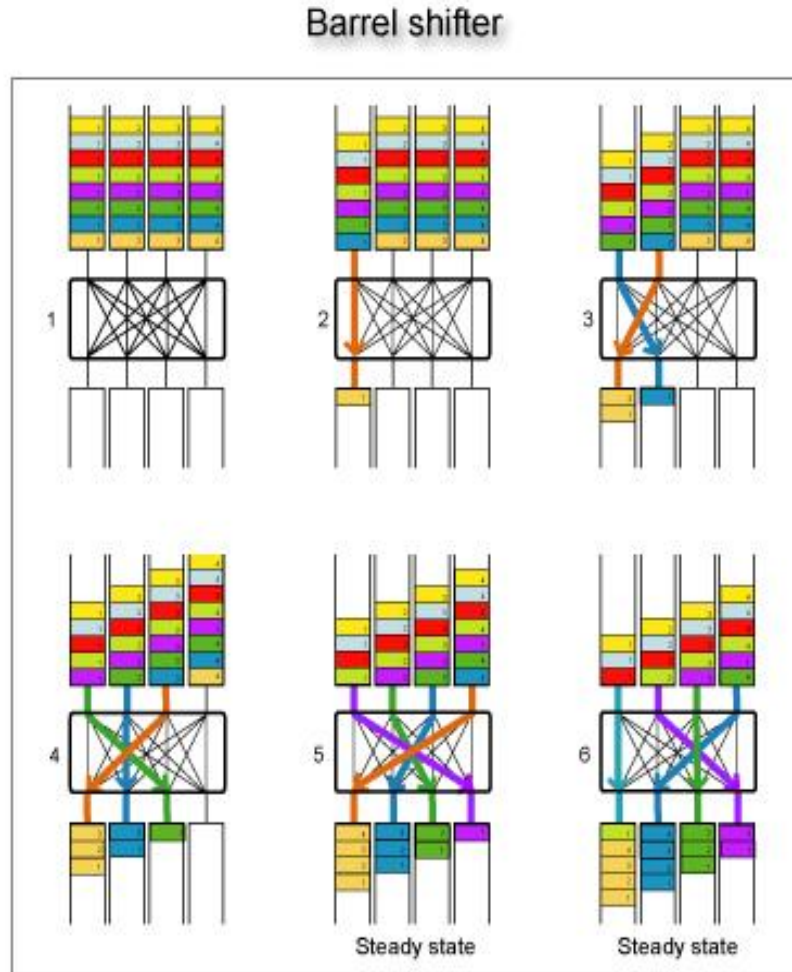
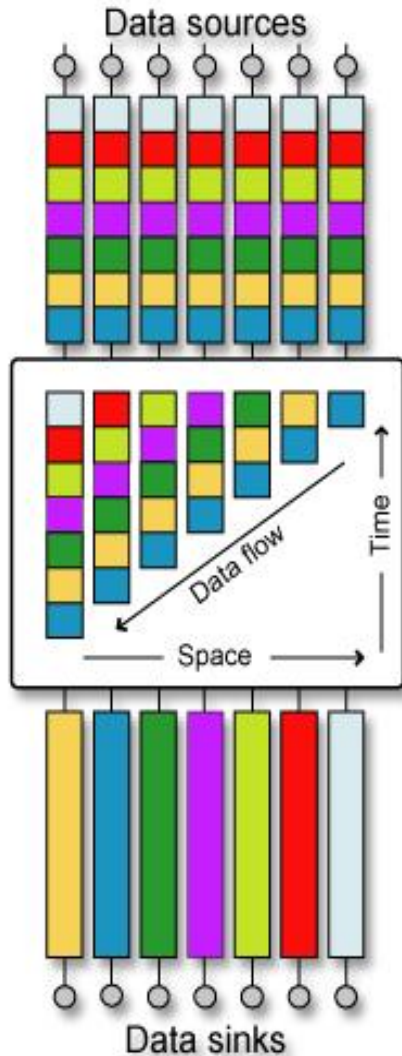
How to deal with this ???

→ 2 possible solutions



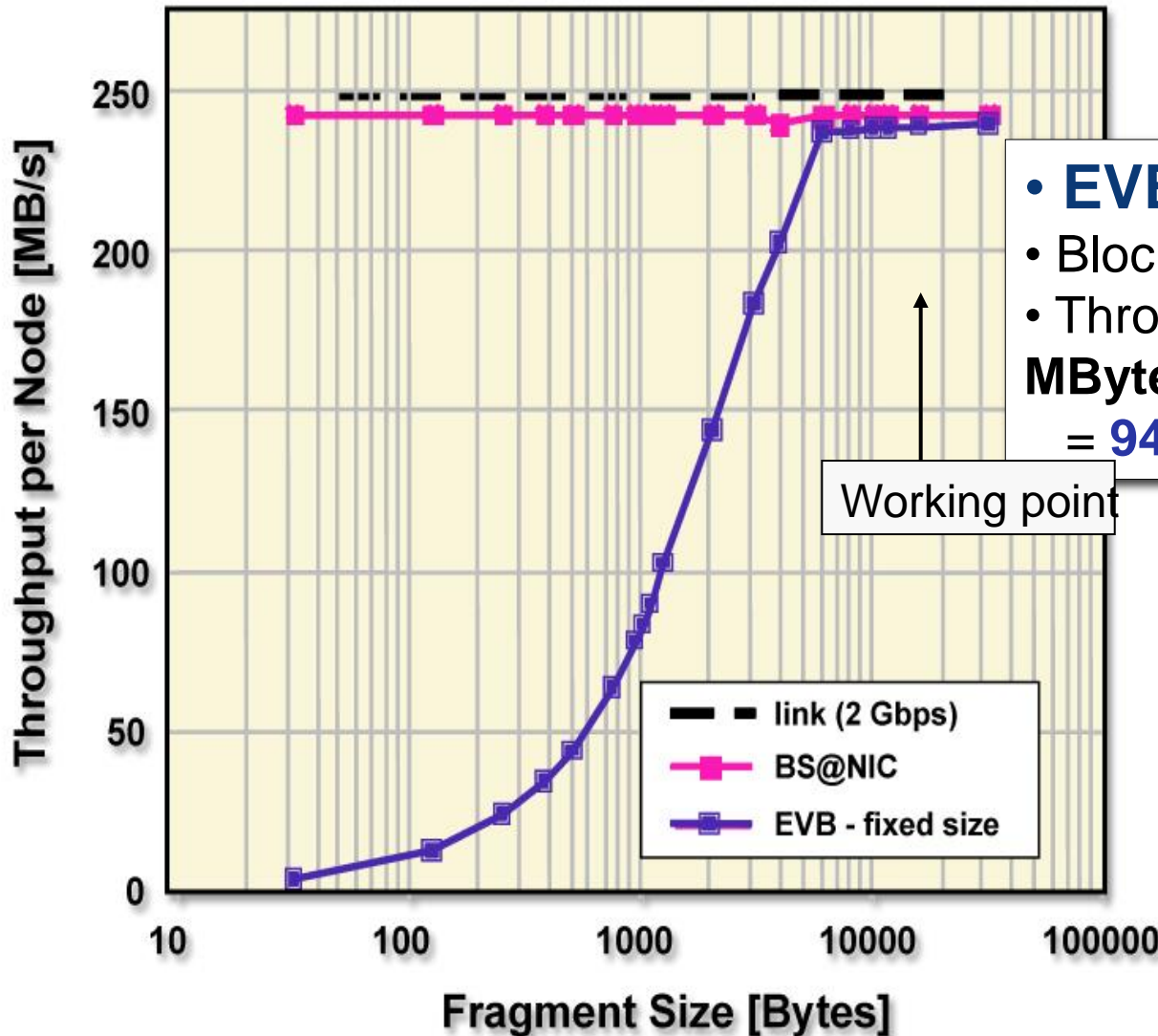
# 1<sup>st</sup> : the clever solution: traffic shaping

## Example: Barrel Shifter





# Barrel Shifter: Measured Performance



- **EVB - Demo 32x32**
- Blocksize 4kB
- Throughput at **234 MByte/s**  
= **94% of link Bandwidth**

Working point

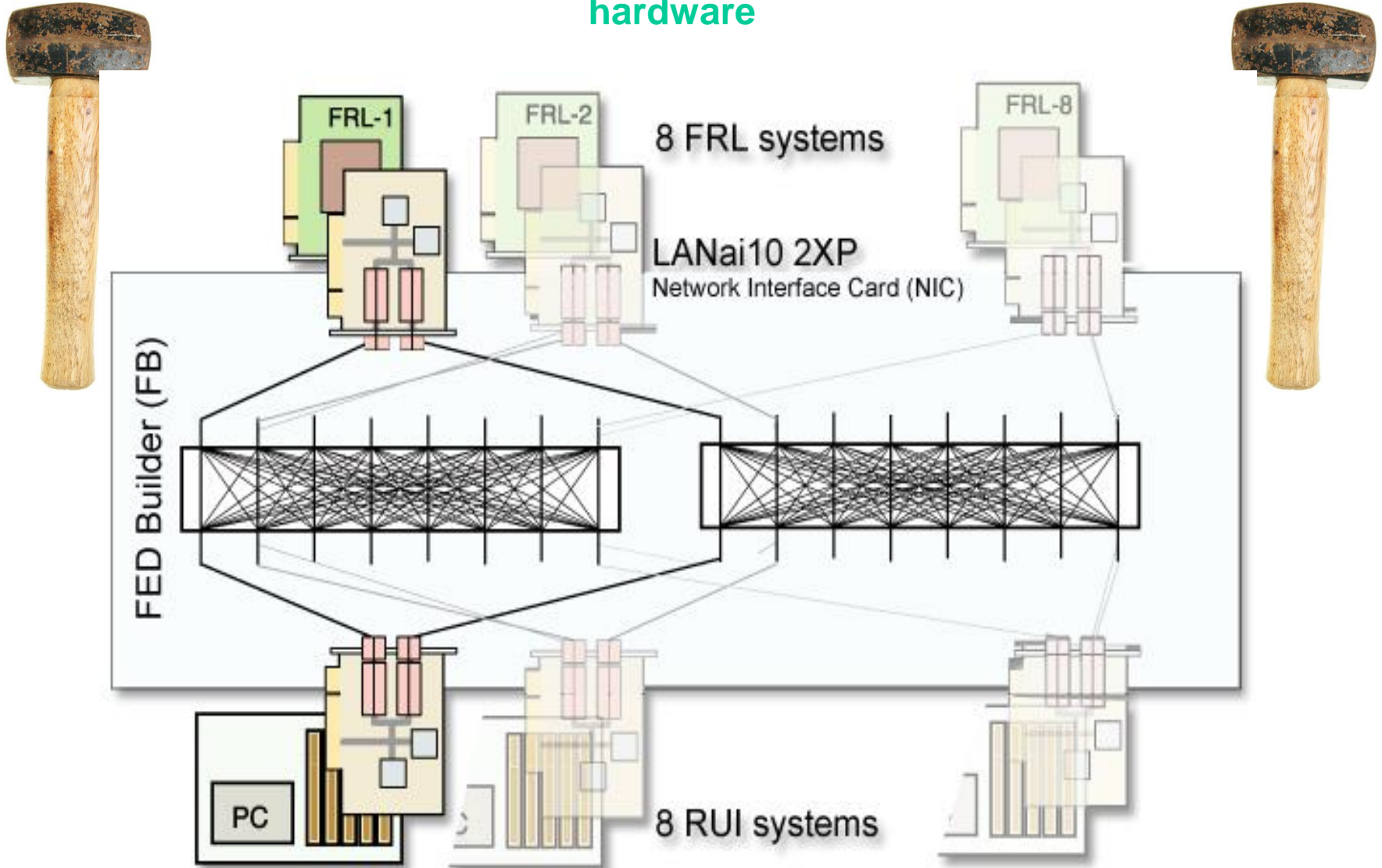
Measurement 2003  
(still valid)





# 2<sup>nd</sup> Solution: “take the hammer”

Over-dimension the system: buy twice as much hardware





# What did CMS do???

☺ Of course: we took the hammer ☺

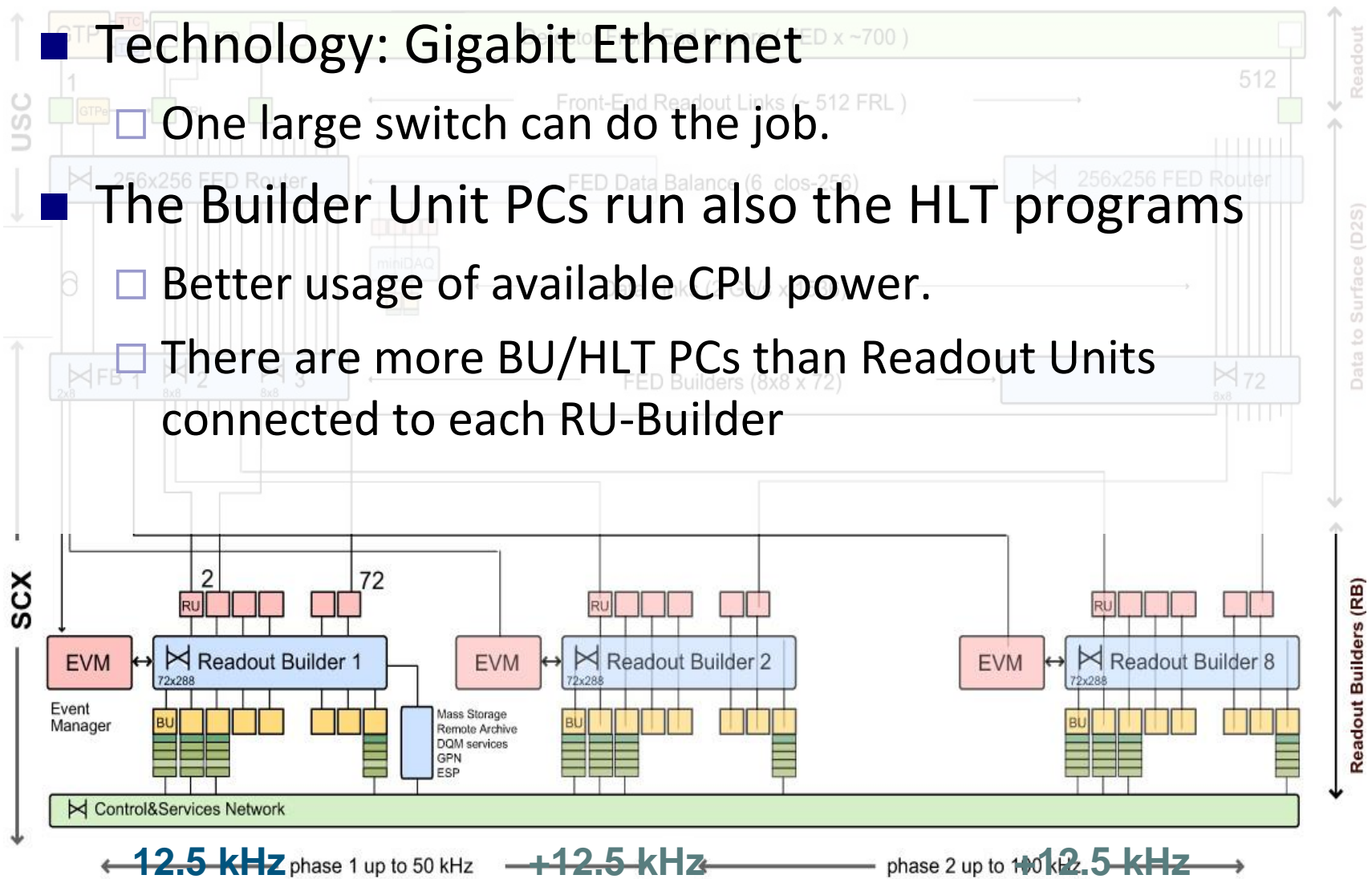
## ■ Advantages:

- Much less development work
- No dependence on internal working of the switch
- Much less maintenance work
  
- Most important: redundancy**
  - If one rail fails: continue to run with one leg ( → less performance but still taking data !!!!)



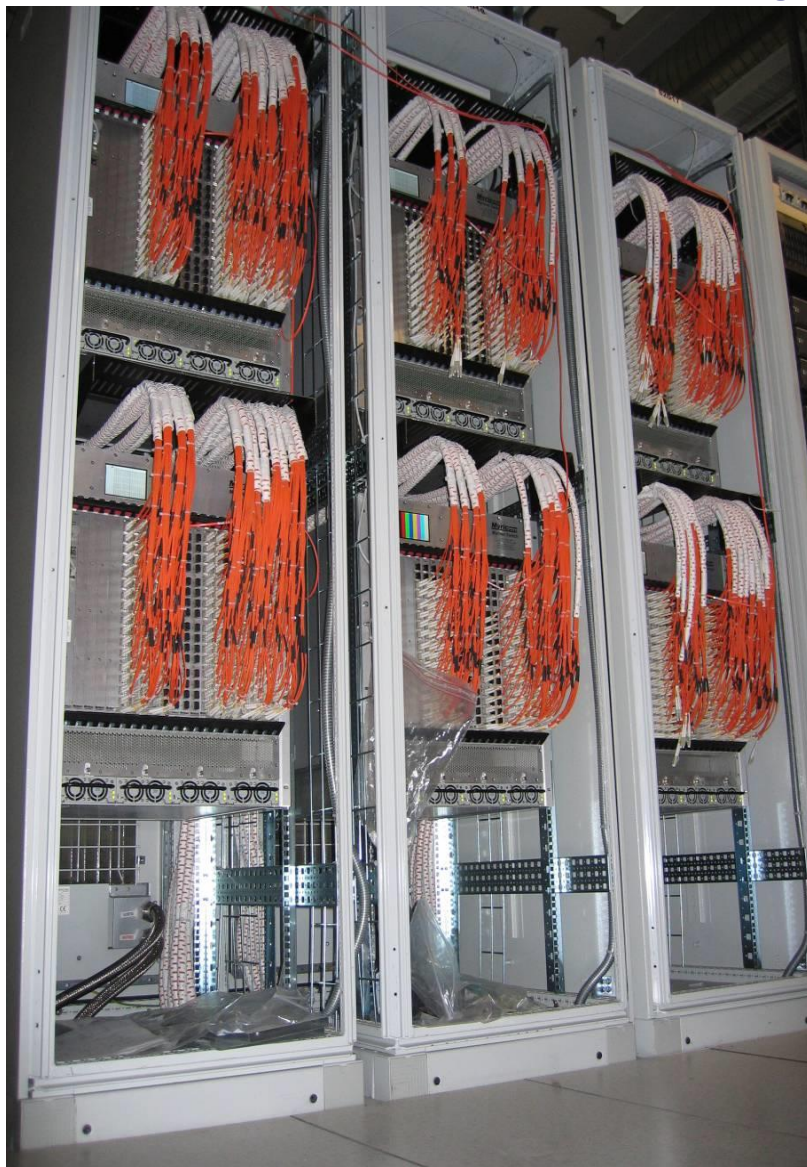
# 2<sup>nd</sup> stage Event Builder: “bread and butter”

- Technology: Gigabit Ethernet
  - One large switch can do the job.
- The Builder Unit PCs run also the HLT programs
  - Better usage of available CPU power.
  - There are more BU/HLT PCs than Readout Units connected to each RU-Builder





# Event Builder Components



Half of the CMS FED Builder

One half of the FEDBuilder is installed close to the experiment in the underground.

The other half is on the surface close to the RU-Builder and the Filter Farm implementing the HLT.

The FEDBuilder is used to transport the data to the surface.

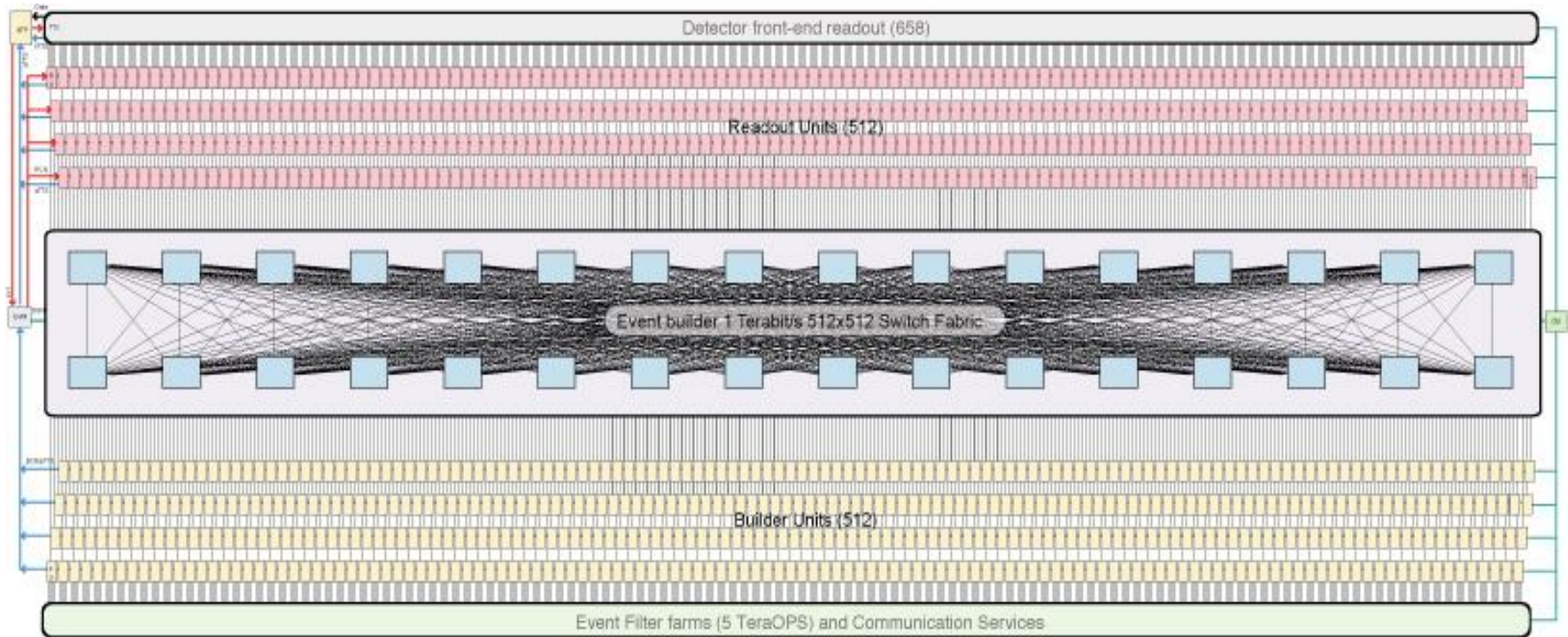
# Event Builder Components

The RU-Builder Switch  
(for 2 slices)





# EVB example: CMS



<b>Level-1 maximum trigger rate</b>	<b>100 kHz</b>	No. Readout systems	≈ 512
<b>Average event size</b>	<b>1 Mbyte</b>	No. Filter Subfarms	≈ 512 x n
<b>Builder network</b>	<b>1 Terabit/s</b>	No. (C&D) network ports	≈ 10000
<b>Event filter computing power</b>	<b>5 10<sup>6</sup> MIPS</b>	No. programmable units	≈ 10000
<b>Event flow control</b>	≈ 10 <sup>6</sup> Mssg/s	System dead time	≈ %

## Achronyms

- TPO Trigger Primitive Generator
- RTT Regional Trigger Processor
- L1T Level-1 Trigger Processor
- GTP Global Trigger Processor
- TTC Timing, Trigger and Control
- STTS Synchronous Trigger Transfer System
- STTS asynchronouse Trigger Transfer System
- FCS FrontEnd System
- FCD FrontEnd Controller
- FCC FrontEnd Controller
- DSF Data to Station
- FRL FrontEnd Readout Link
- FRJ FrontEnd Readout Link
- BU Builder Unit
- FS Filter Subfarm
- EVM Event Manager
- FMM FrontEnd Manager
- EM Event Manager
- EVB Event Builder
- FCN FrontEnd Control Network
- SCN Subferm Control Network
- CSN Computing Service Network
- CCN Central Control Network
- CSN C&D Service Network
- CCSN Central Control System
- FCS Farm Control System