



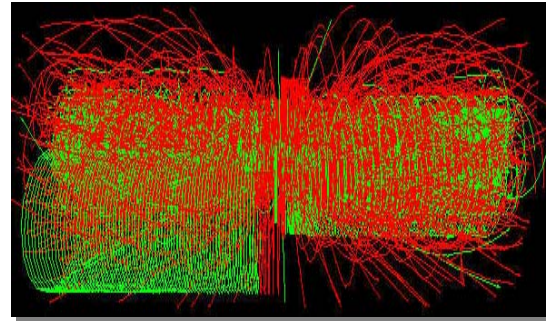
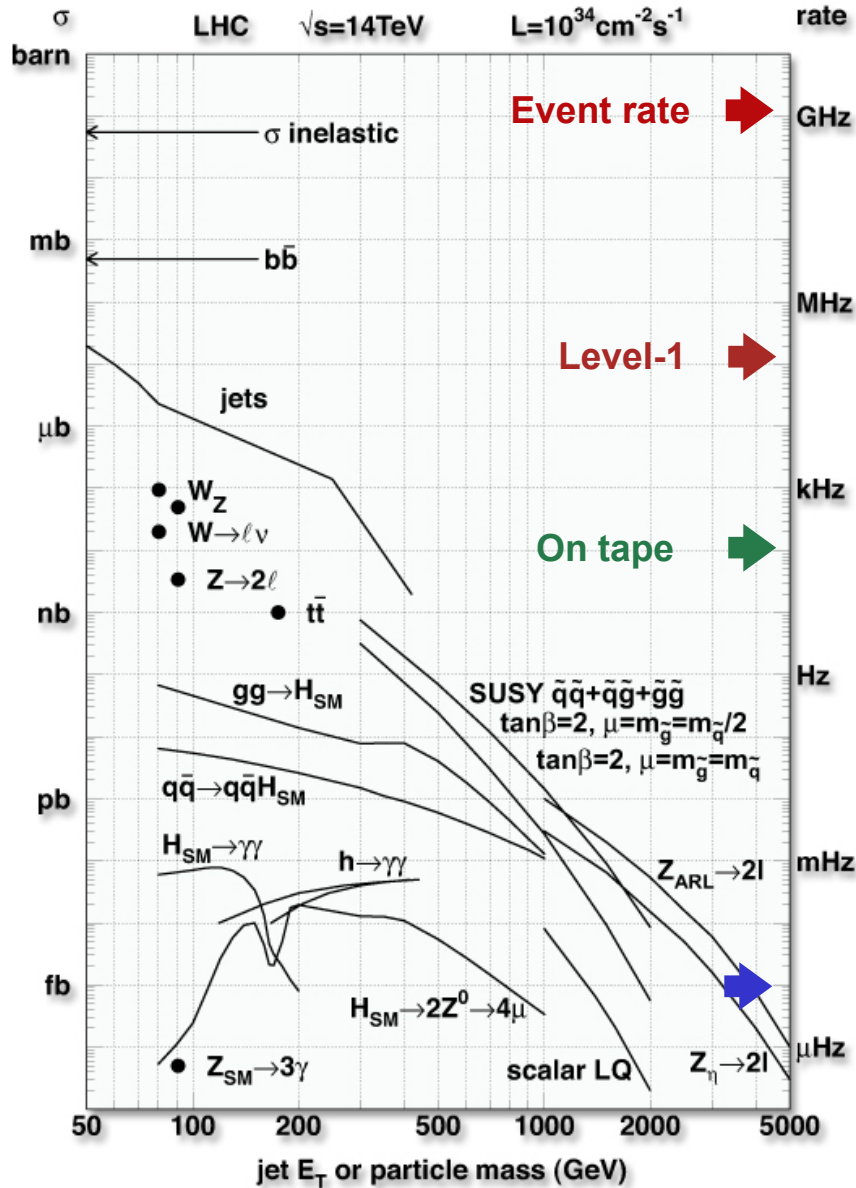
# **CMS Data Acquisition**

**S. Cittolin CERN/PH-CMD  
Oct. 7 2004 LHC DAYS IN SPLIT**

**Requirements and Architecture  
Trigger and Readout  
Event Builder and Technologies  
Plans and summary**



# p-p collisions at LHC



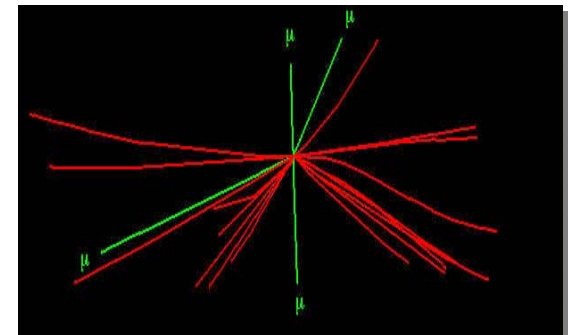
Operating conditions:  
Higgs in 4 muons  
+ ~20 minimum bias

All charged tracks  
with  $pt > 2\text{ GeV}$

**Event Rates:**  $\sim 10^9\text{ Hz}$   
**Event size:**  $\sim 1\text{ MByte}$

**Level-1 Output**  $100\text{ kHz}$   
**Mass storage**  $10^2\text{ Hz}$   
**Event Selection:**  $\sim 1/10^{13}$

Reconstructed tracks  
with  $pt > 25\text{ GeV}$





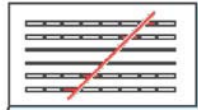
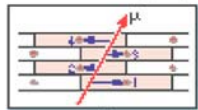
# Requirements and design parameters



## Detectors

### MUON BARREL

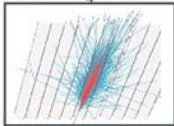
Drift Tube Chambers (DT)



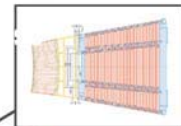
Resistive Plate Chambers (RPC)

### CALORIMETERS

ECAL Scintillating PbWO<sub>4</sub> Crystals



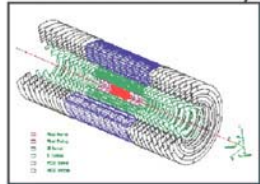
HCAL Scintillator brass sandwich



IRON YOKE

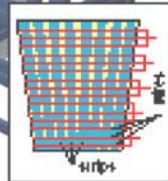
SUPERCONDUCTING COIL

### TRACKERS



Pixels  
Silicon Microstrips

MUON ENDCAPS



Cathode Strip Chambers (CSC)  
Resistive Plate Chambers (RPC)

Total weight : 12,500 t  
Overall diameter : 15 m

Overall length : 21.6 m  
Magnetic field : 4 Tesla

## Detector Channels Control Ev. Data

Detector	Channels	Control	Ev. Data
Pixel	60000000	1 GB	50 (kB)
Tracker	10000000	1 GB	650
Preshower	145000	10 MB	50
ECAL	85000	10 MB	100
HCAL	14000	100 kB	50
Muon DT	200000	10 MB	10
Muon RPC	200000	10 MB	5
Muon CSC	400000	10 MB	90
Trigger		1 GB	16

**Event size**

**1 Mbyte**

**Max LV1 Trigger**

**100 kHz**

**Online rejection**

**99.999%**

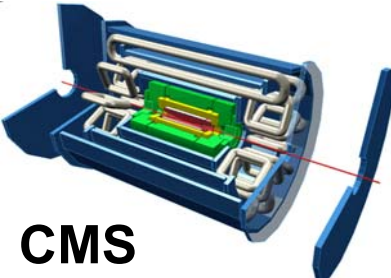
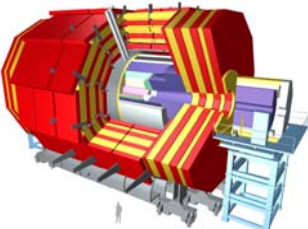
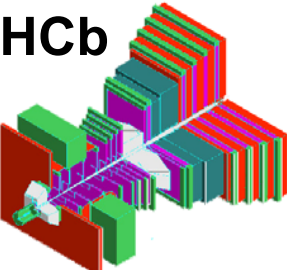

**System dead time**

**~ %**



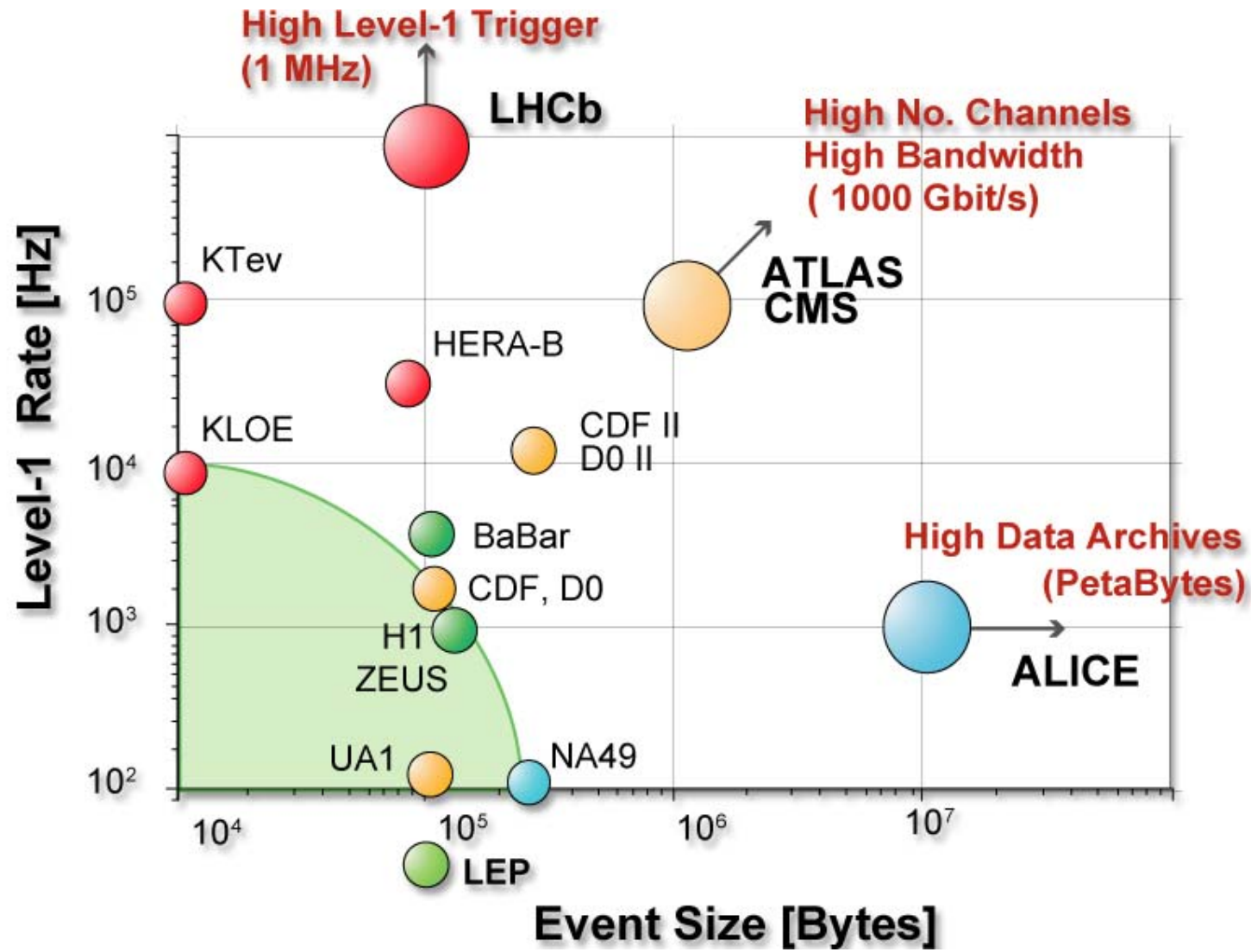
# LHC trigger and DAQ summary



	No. Levels Trigger	First Level Rate (Hz)	Event Size (Byte)	Readout Bandw. (GB/s)	Archive MB/s (Event/s)
<b>ATLAS</b> 	<b>3</b>	<b><math>10^5</math></b> LV-2 <b><math>10^3</math></b>	<b><math>10^6</math></b>	<b>10</b>	<b>100</b> ( $10^2$ )
<b>CMS</b> 	<b>2</b>	<b><math>10^5</math></b>	<b><math>10^6</math></b>	<b>100</b>	<b>100</b> ( $10^2$ )
<b>LHCb</b> 	<b>3</b>	LV-0 <b><math>10^6</math></b> LV-1 <b><math>4 \cdot 10^4</math></b>	<b><math>2 \times 10^5</math></b>	<b>4</b>	<b>40</b> ( $2 \times 10^2$ )
<b>ALICE</b> 	<b>4</b>	Pp-Pp <b>500</b> p-p <b><math>10^3</math></b>	<b><math>5 \times 10^7</math></b> <b><math>2 \times 10^6</math></b>	<b>5</b>	<b>1250</b> ( $10^2$ ) <b>200</b> ( $10^2$ )



# Trigger and data acquisition trends

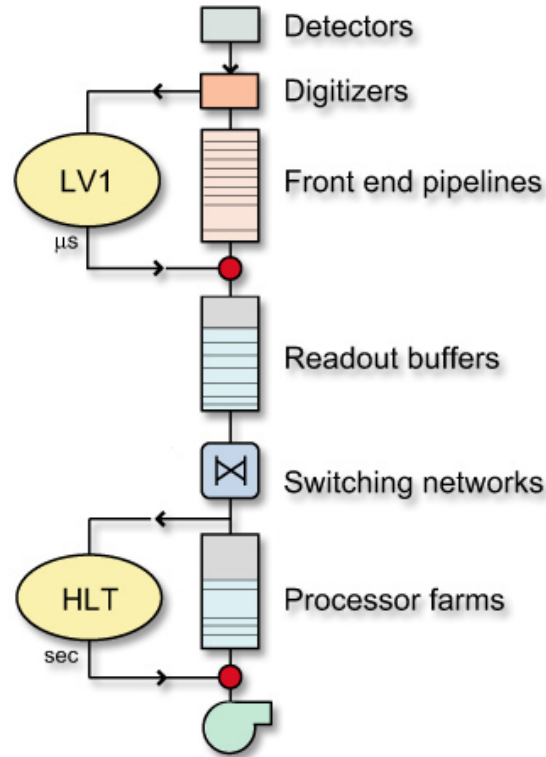




# CMS DAQ structure: 2 physical triggers



**40 MHz**  
Clock driven  
Custom processors



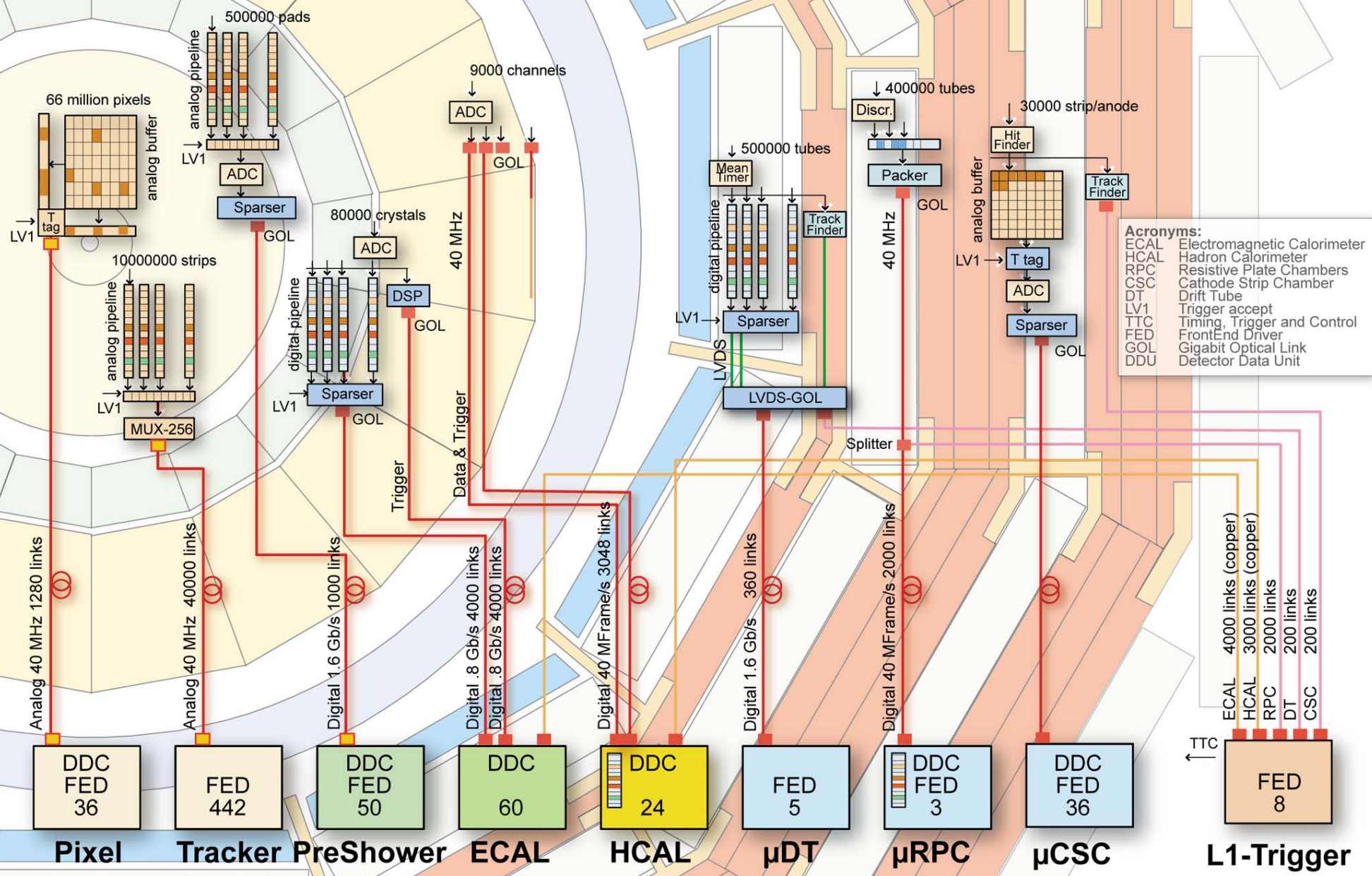
**Level-1 Trigger**  
**Custom design**

**100 kHz**  
Event driven  
PC network

**High-Level Trigger**  
**Industry products**

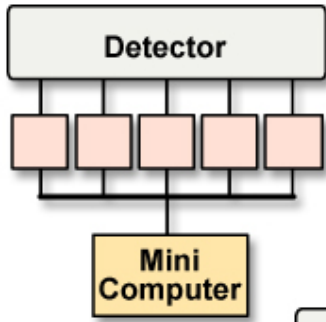
**Level-1 output / HLT input 100 kHz**  
**Network bandwidth 1 Terabit/s**  
**HLT output  $10^2$  Hz**  
**Invest in data transportation and CPU**

# CMS front-end readout systems



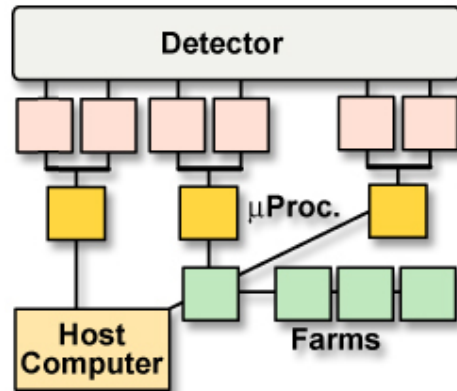


# Evolution of DAQ technologies and structures



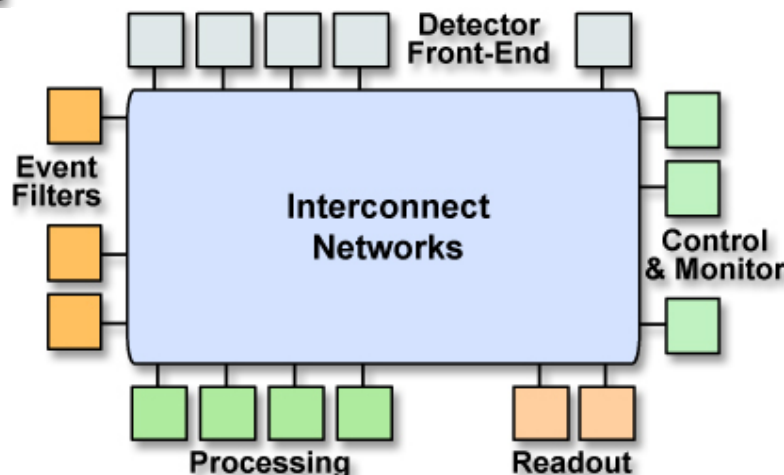
## 1970-80: Minicomputers

Readout custom design  
First standard: CAMAC  
• **kByte/s**



## 1980-90: Microprocessors

HEP standards (Fastbus)  
Embedded CPU, Industry standards (VME)  
• **MByte/s**



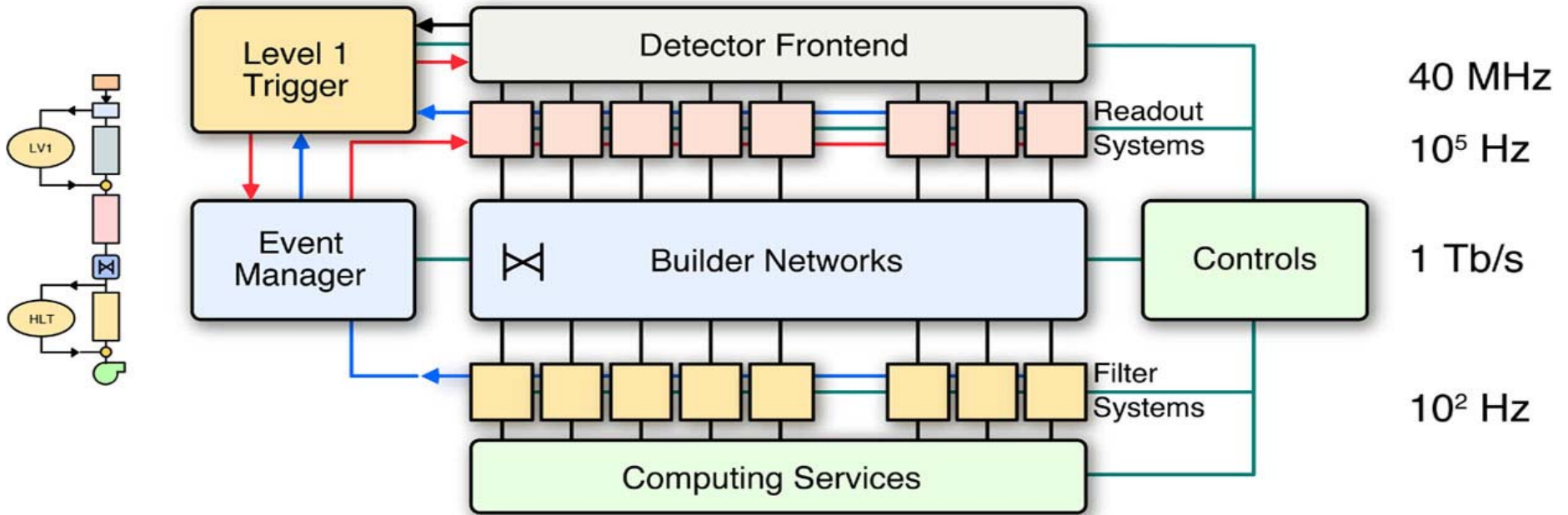
## 2000-xx: Networks

IT commodities, PC, Clusters  
Internet, Web, etc.  
• **GByte/s**





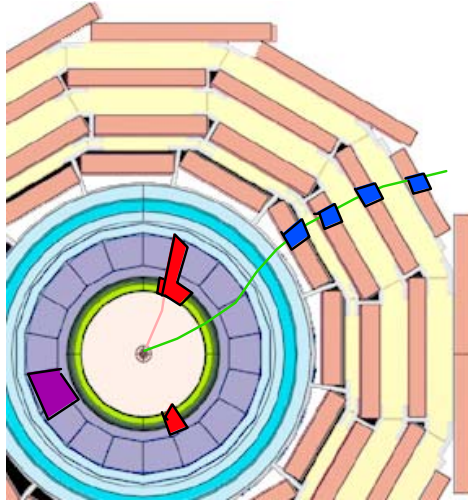
# DAQ baseline structure



Collision rate	40 MHz	No. of In-Out units	<b>512</b>
<b>Level-1 Maximum trigger rate</b>	<b>100 kHz</b>	<b>Readout network bandwidth</b>	<b>≈ 1 Terabit/s</b>
<b>Average event size</b>	<b>≈ 1 Mbyte</b>	<b>Event filter computing power</b>	<b>≈ 10<sup>6</sup> SI95</b>
Event Flow Control	≈ 10 <sup>6</sup> Mssg/s	Data production	≈ Tbyte/day
		No. of PC motherboards	≈ Thousands



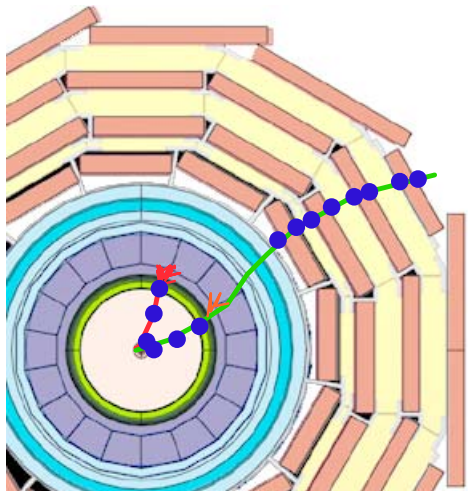
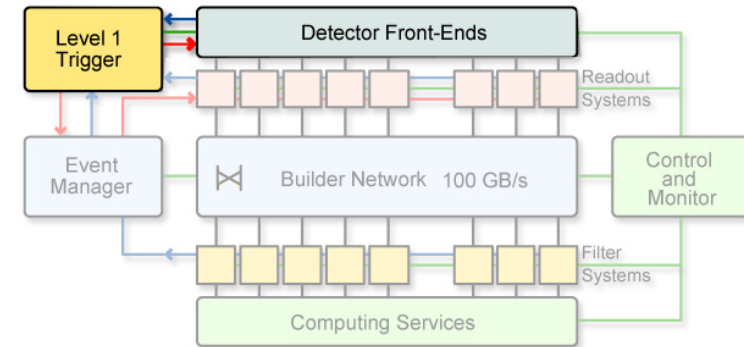
# Two trigger levels



## Level-1: Specialized processors 40 MHz synchronous

- Particle identification:
- high pT electron, muon, jets, missing ET
- Local pattern recognition and energy evaluation on prompt macro-granular information from calorimeter and muon detectors

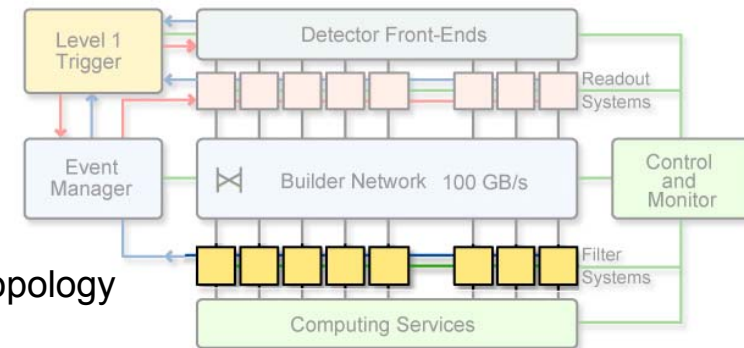
**99.99 % rejected 0.01 Accepted**



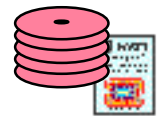
## High trigger levels: CPU farms 100 kHz asynchronous farms

- Clean particle signature
- Finer granularity precise measurement
- Kinematics. effective mass cuts and event topology
- Track reconstruction and detector matching
- Event reconstruction and analysis

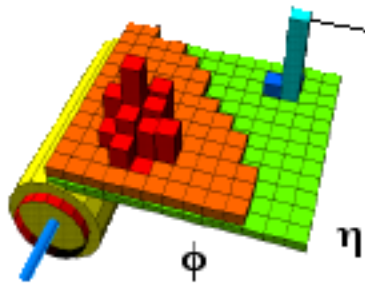
**99.9 % rejected 0.1 Accepted**



100-1000 Hz. Mass storage  
Reconstruction and analysis.

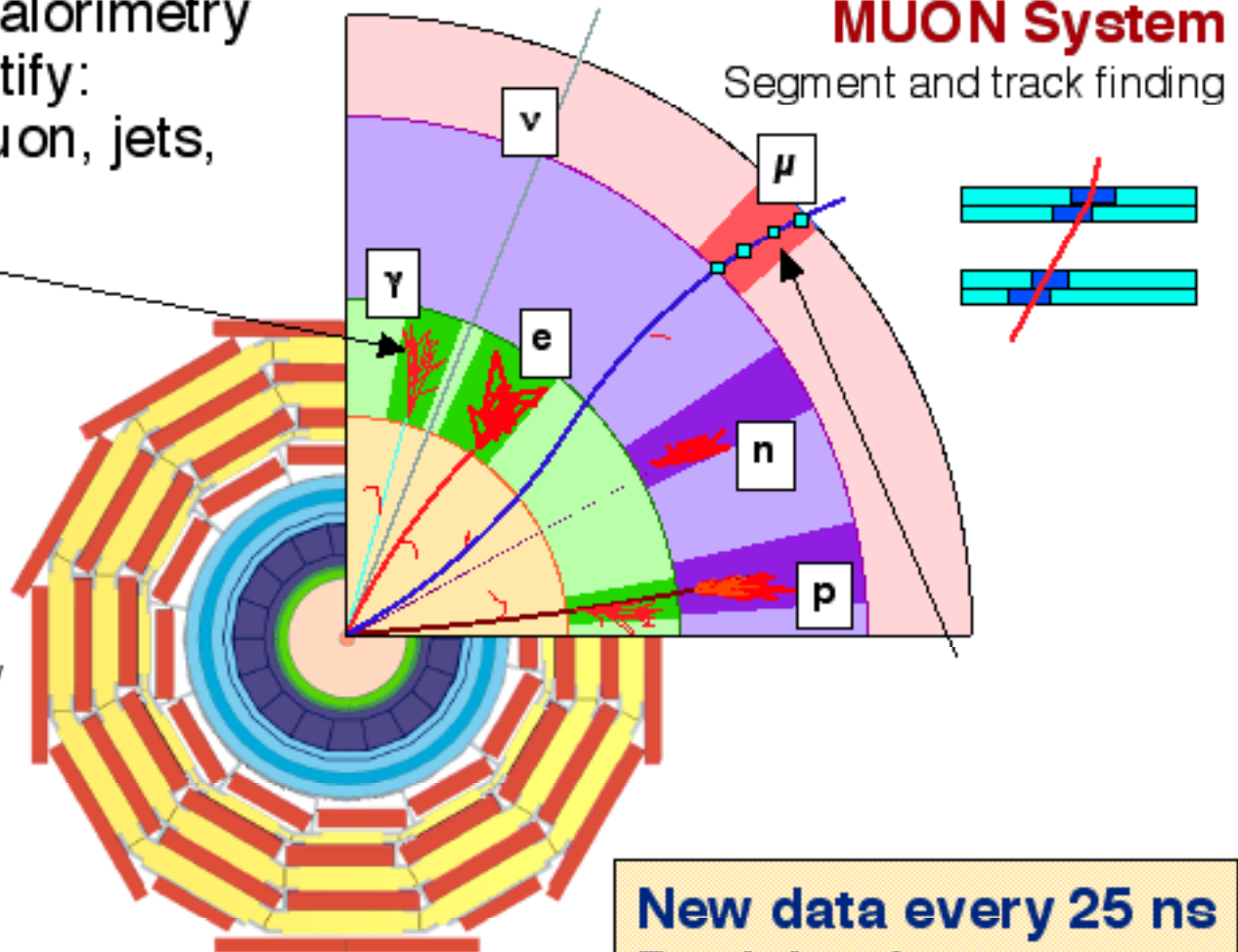


Use prompt data (calorimetry and muons) to identify:  
 High  $p_t$  electron, muon, jets,  
 missing  $E_T$

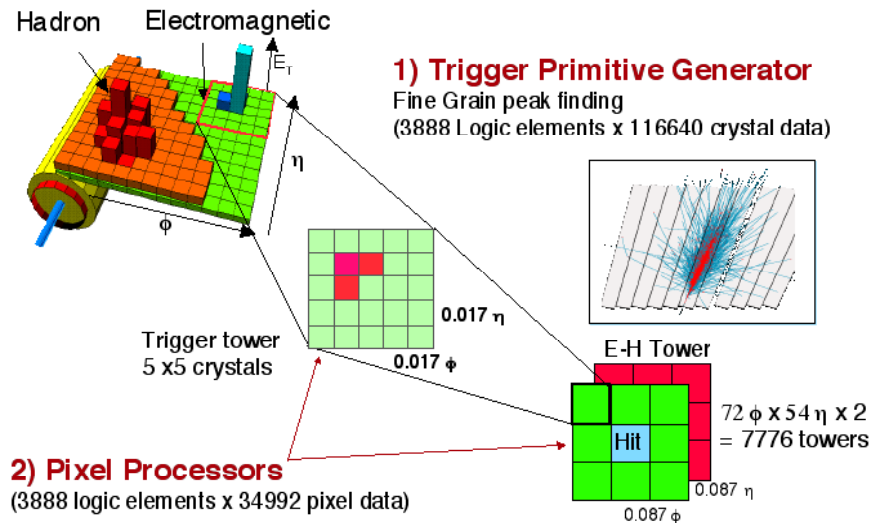


## CALORIMETERS

Cluster finding and energy deposition evaluation

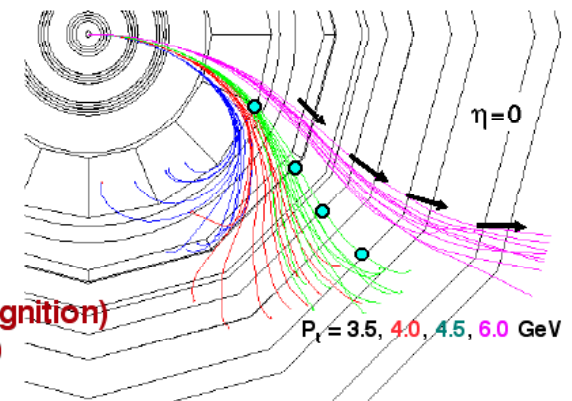


**New data every 25 ns**  
**Decision latency  $\sim \mu\text{s}$**



**Trigger based on tracks in external muon detectors that point to interaction region**

- Low- $p_t$  muon tracks don't point to vertex
- Multiple scattering
- Magnetic deflection
- Two detector layers
- Coincidence in "road"



**Detectors:**

- RPC (pattern recognition)**
- DT(track segment)**

## Trigger Primitive Generator

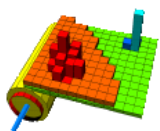
Fine grain Flag Max of ( )

## Pixel Processor

$E_t$  cut + Max ( ) > Threshold

Longitudinal cut (H/E) / < 0.05

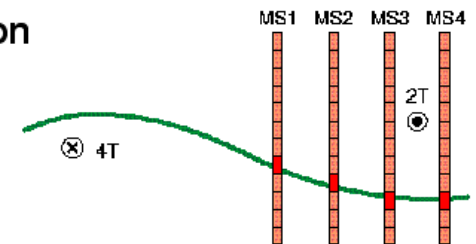
Neighbors longitudinal cut / < 2 GeV



One of ( )  
↓  
**ISOLATED ELECTRON**

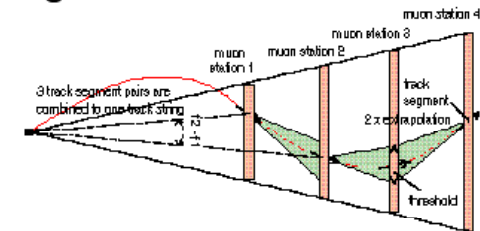
## RPC pattern recognition

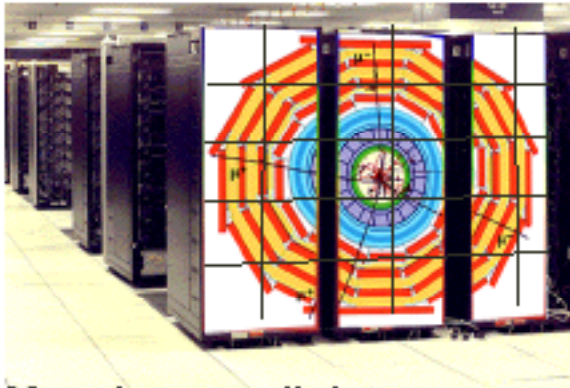
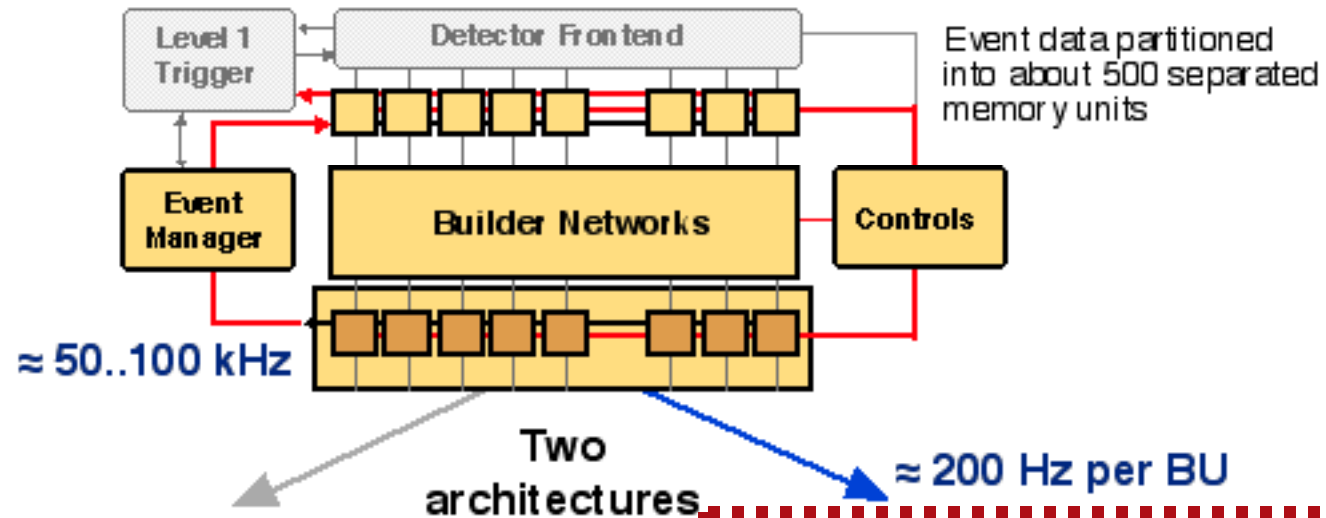
- Pattern catalog
- Fast logic



## DT and CSC track finding:

- Finds hit/segments
- Combines vectors
- Formats a track
- Assigns  $p_t$  value





**Massive parallel system**  
**ONE event, ALL processors**

- Low latency
- Complex I/O
- Parallel programming

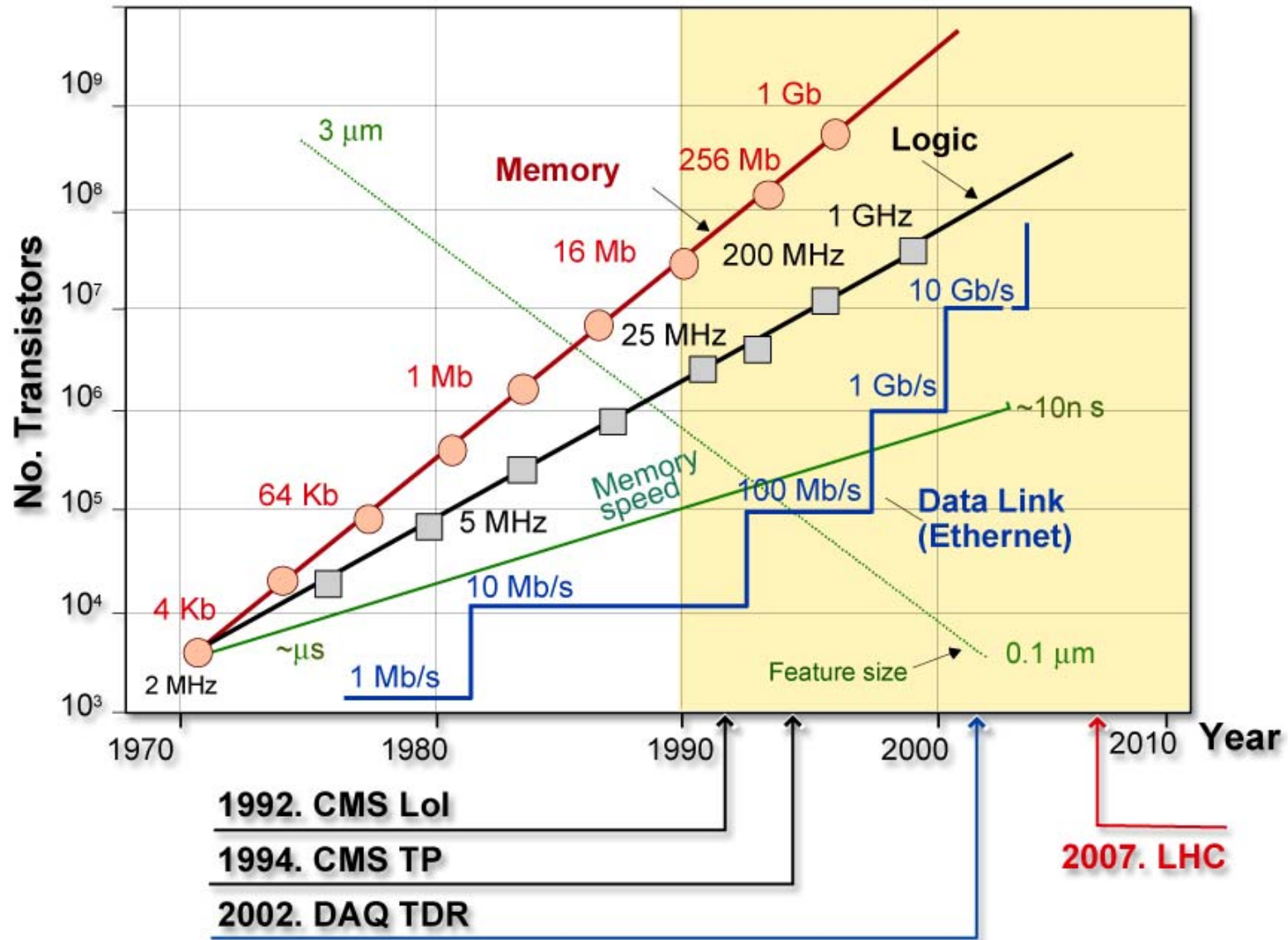


**Farm of processors**  
**ONE event, ONE processor**

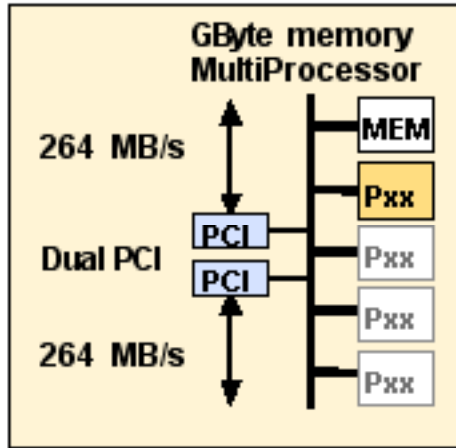
- High latency (larger buffers)
- Simpler I/O
- Sequential programming



# Technology trends



## 1990' PCI



IO and Processing systems : Commercial PCs  
 Operating systems : Unix(Linux),  
 Interfaces standards : PC IO systems (e.g. PCI)

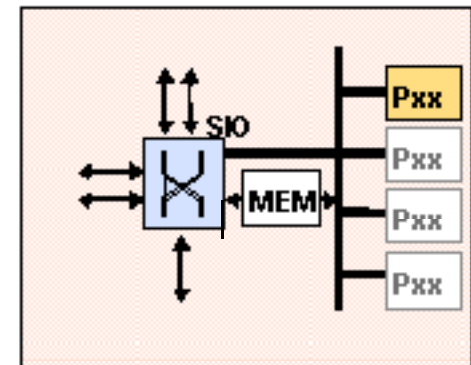
Desktop/Server  
 current architecture  
 Peripheral IO bus PCI:  
 33/66 MHz x 32/64 bit  
**100/200/400 MB/s**



## 200X: PCI-X ...

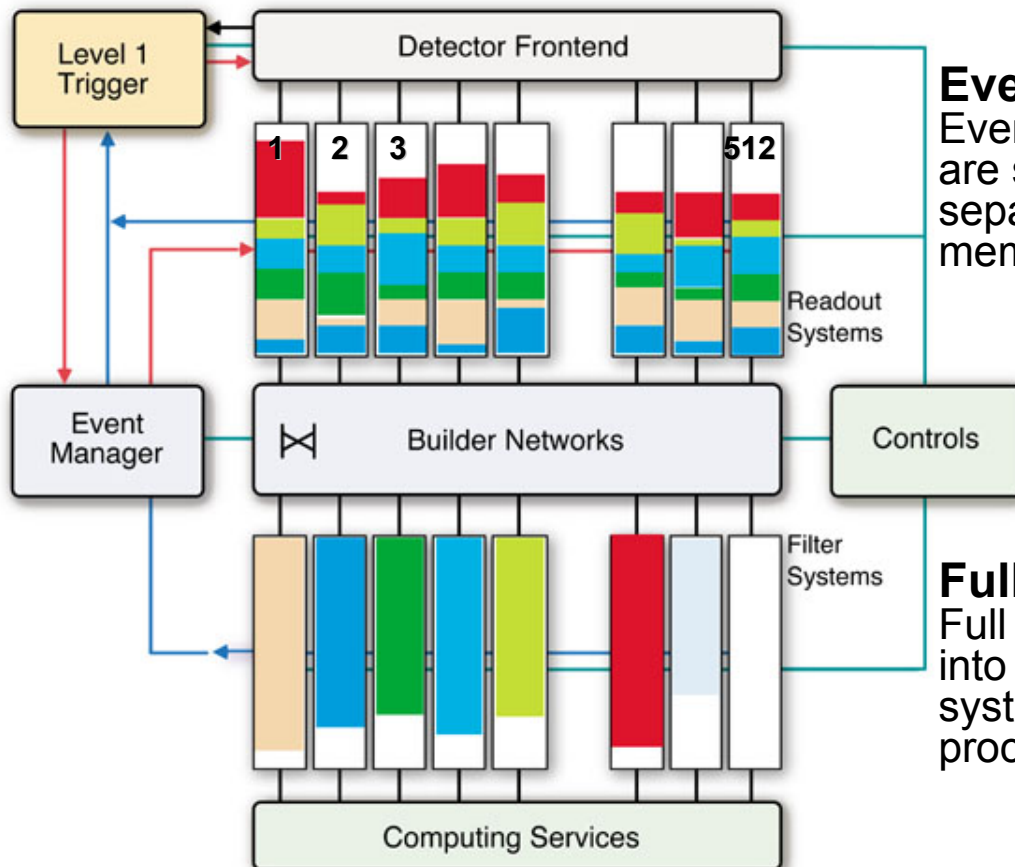
### 2002 PC mother boards:

- 2 GHz dual processors
- 4 PCI-X ports at 1GB/s
- 3 GB/s memory bandwidth
- Suitable for all DAQ readout applications

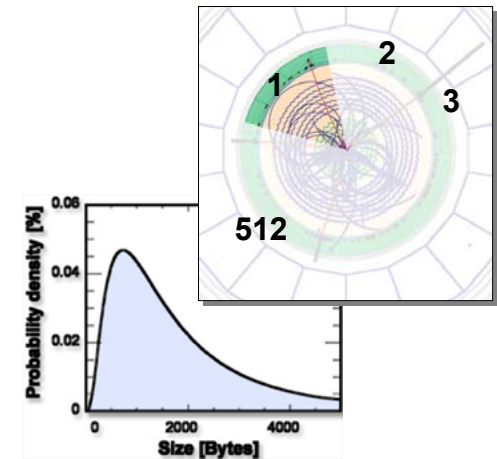


## Event builder :

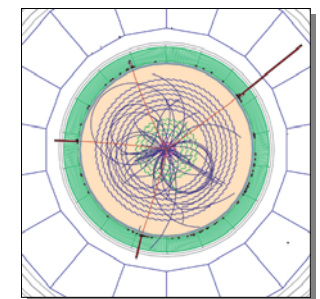
Physical system interconnecting data sources with data destinations. It has to move each event data fragments into a same destination



**Event fragments :**  
Event data fragments are stored in separated physical memory systems



**Full events :**  
Full event data are stored into one physical memory system associated to a processing unit



**512 Data sources for 1 Mbyte events**  
**~1000s HTL processing nodes**



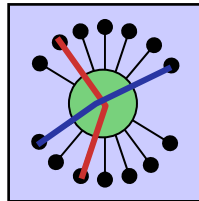


# EVB and switch technologies

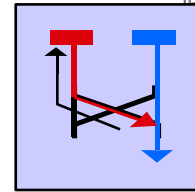


## Myricom Myrinet 2000

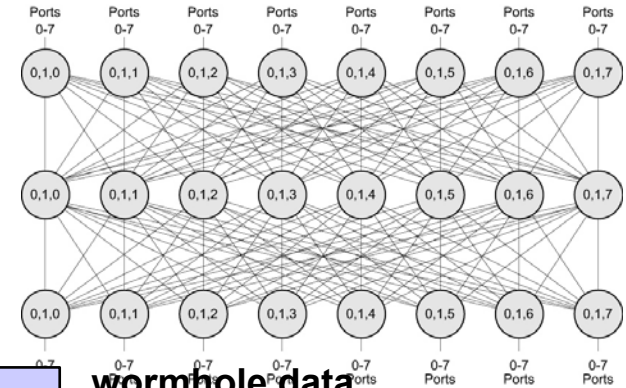
- Switch: **Clos-128 x 2.5 Gb/s port**
- NIC: M3S-PCI64B-2 (**LANai9**)
- **Custom Firmware**



**Implementation :**  
16x16 port X-bar capable of channeling data between any two ports.

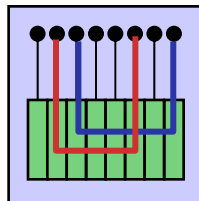


**wormhole data**  
transport with flow control at all stages

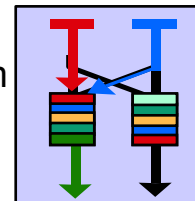


## Gigabit Ethernet

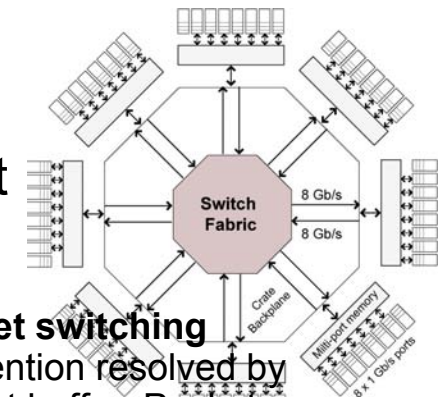
- Switch: Foundry **FastIron 64 x 1.2 Gb/s port**
- NIC: **Alteon** (running standard firmware)



**Implementation:**  
Multi-port memory system of R/W access bandwidth greater than the sum of all port speeds



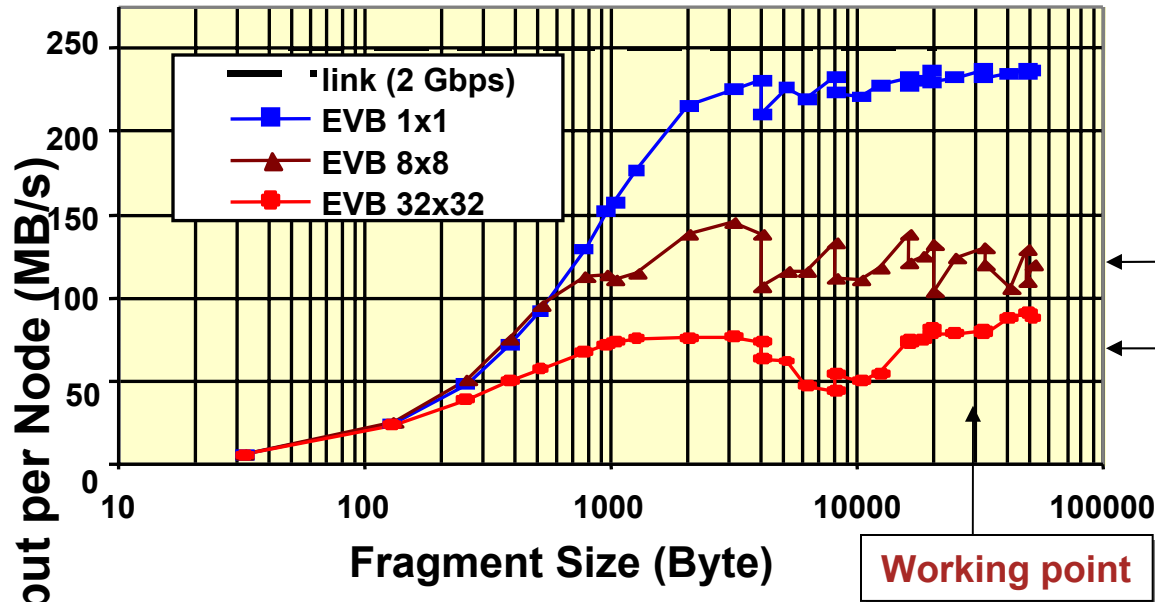
**Packet switching**  
Contention resolved by Output buffer. Packets can be lost.



**Infiniband** • 2.5 Gb/s demo product. Tests ongoing with a small 2x2 setup



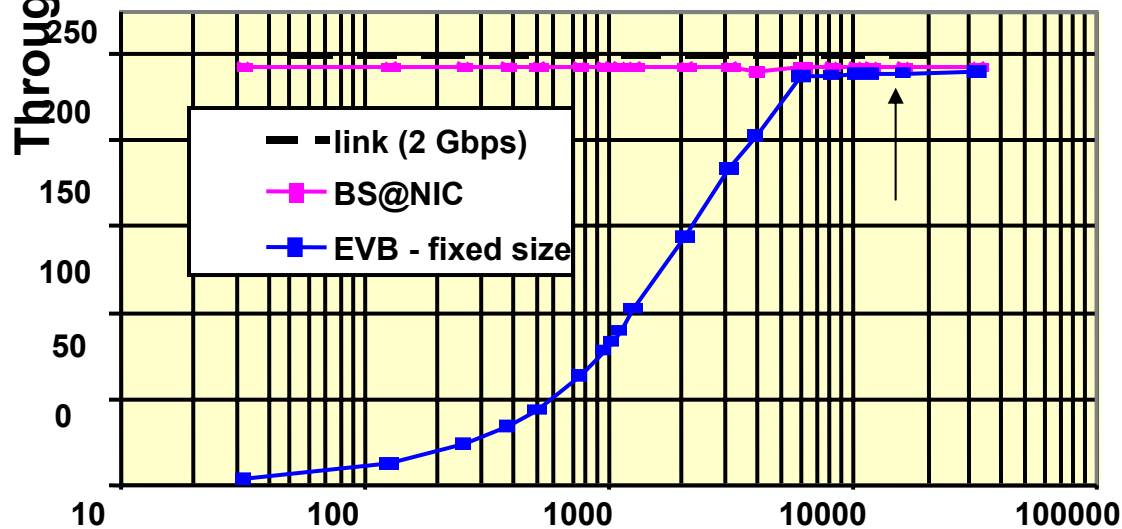
# 32x32 Myrinet EVB protocols results



## Random traffic

**8x8**: single stage:  
max. **utilization**:  $\approx 50\%$

**32x32**: two stage network  
max. **utilization**  $\approx 30\%$

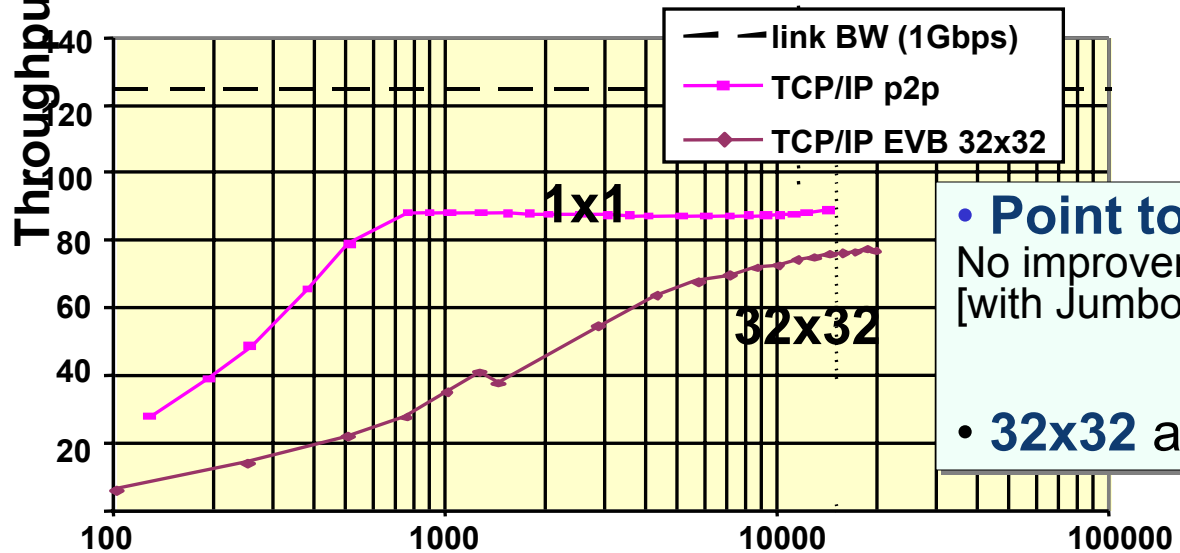
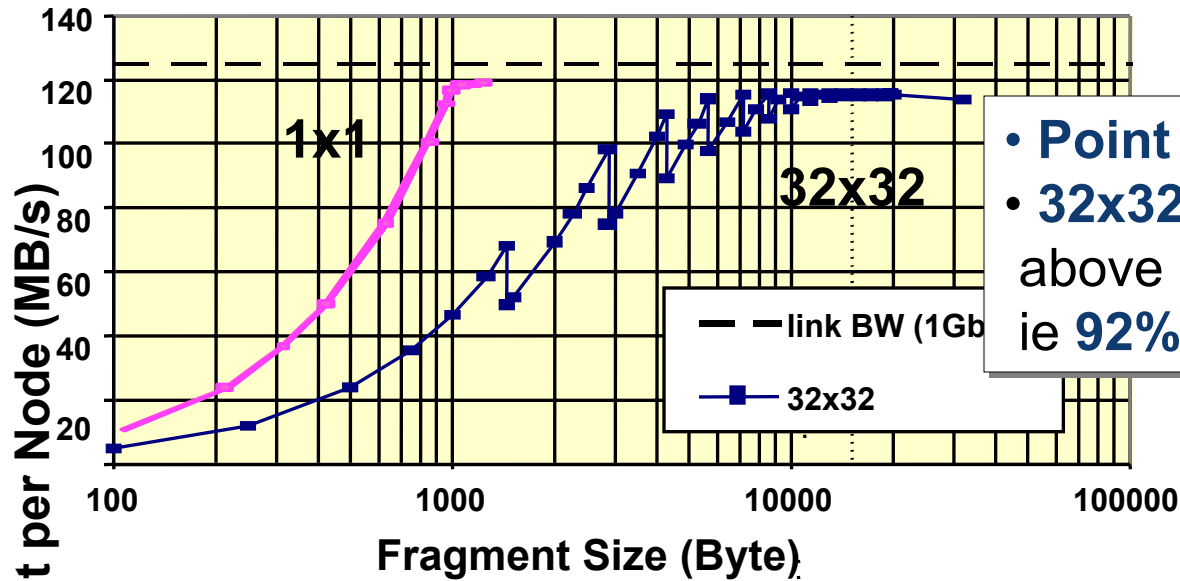


## Barrel shifter

- Fixed size event fragments  
below 4k: Fragment < BS carrier  
above 4k: Fragment > BS carrier
- Throughput at **234 MB/s**  
= **94% of link Bandwidth**



# 32x32 GbE EVB protocols results





# EVB demonstrators summary



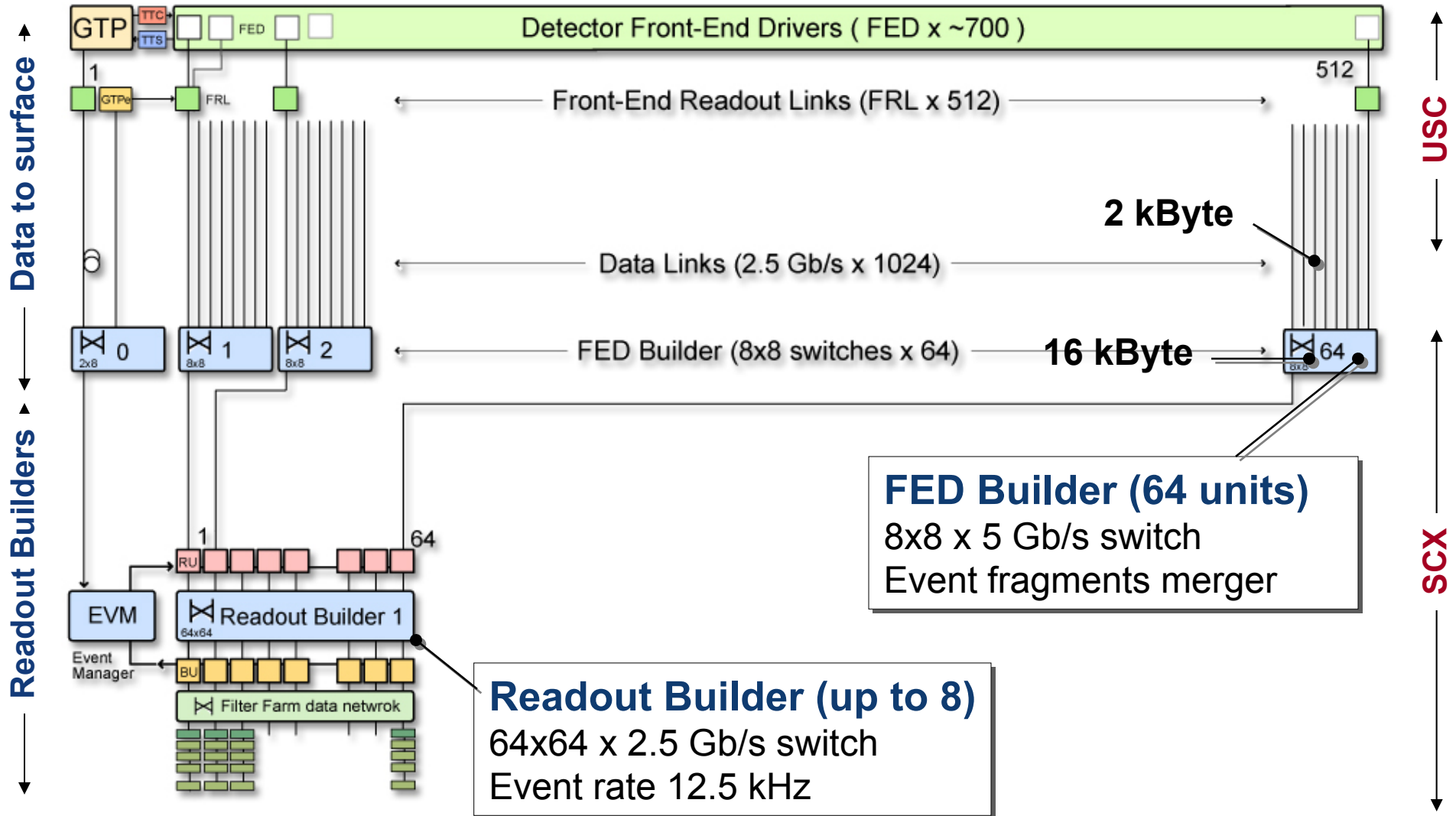
	<b>Myrinet 2000</b>	<b>GbE raw packet</b>	<b>GbE TCP/IP</b>
<b>Test bench</b>	<b>32x32</b>	<b>32x32</b>	<b>32x32</b>
<b>Port speed</b>	<b>2.5 Gbit/s</b>	<b>1.2 Gbit/s</b>	<b>1.2 Gbit/s</b>
<b>Random traffic</b>	<b>30-50%</b>	<b>50%, 92% (*)</b>	<b>30%, 60% (*)</b>
<b>Barrel switch</b>	<b>94%</b>	-	-
<b>CPU load</b>	<b>Low</b>	<b>High</b>	<b>High</b>
<b>1 Tbit/s EVB</b>	<b>512x512</b>	<b>1024x1024</b>	<b>2048x2048</b>
<b>No. switches</b>	<b>8 128-Clos</b>	<b>16 256-port</b>	<b>32 256-port</b>



(\*) with fragment sizes **larger than 16kB**

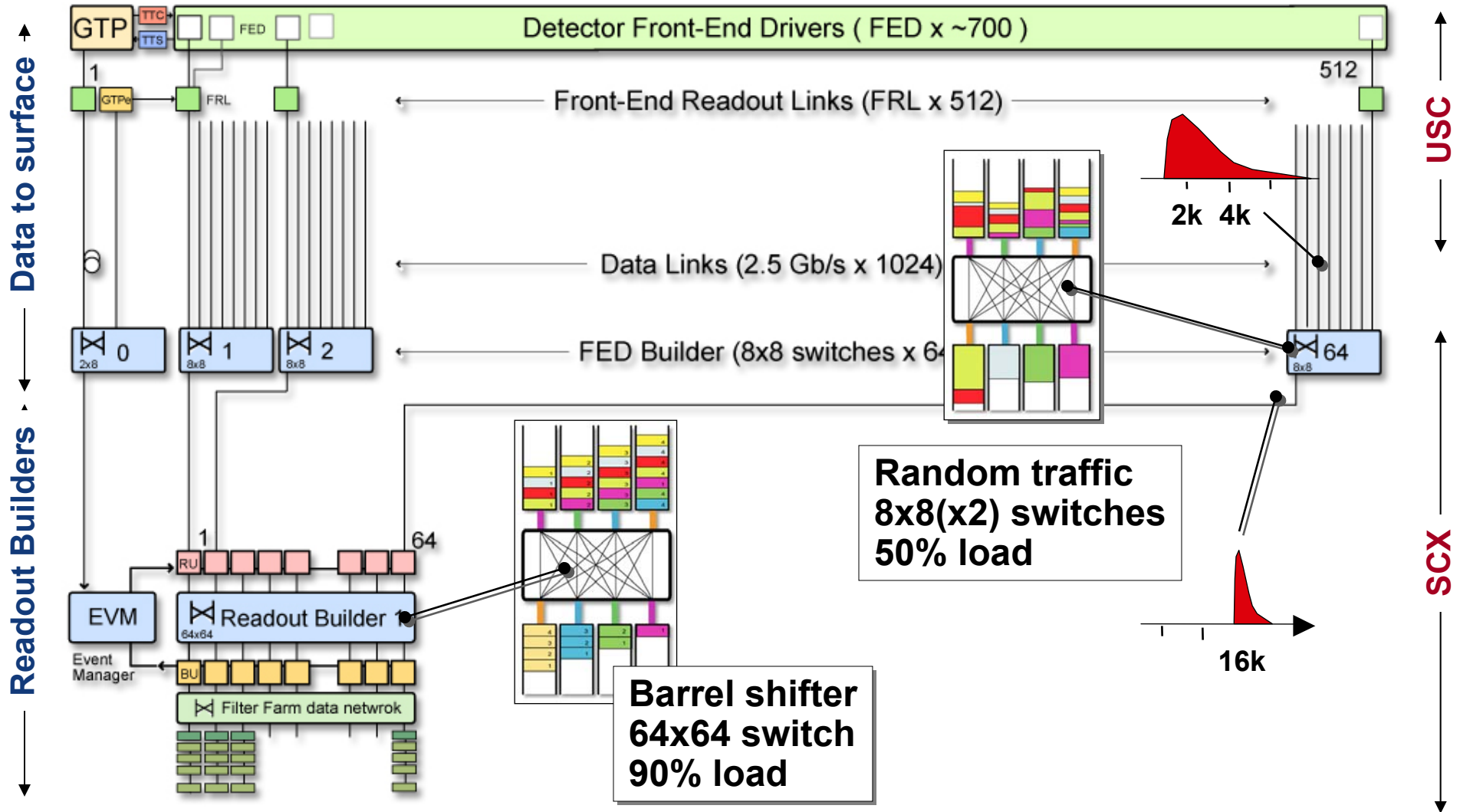


# 2 stages: Data to surface & Readout Builder



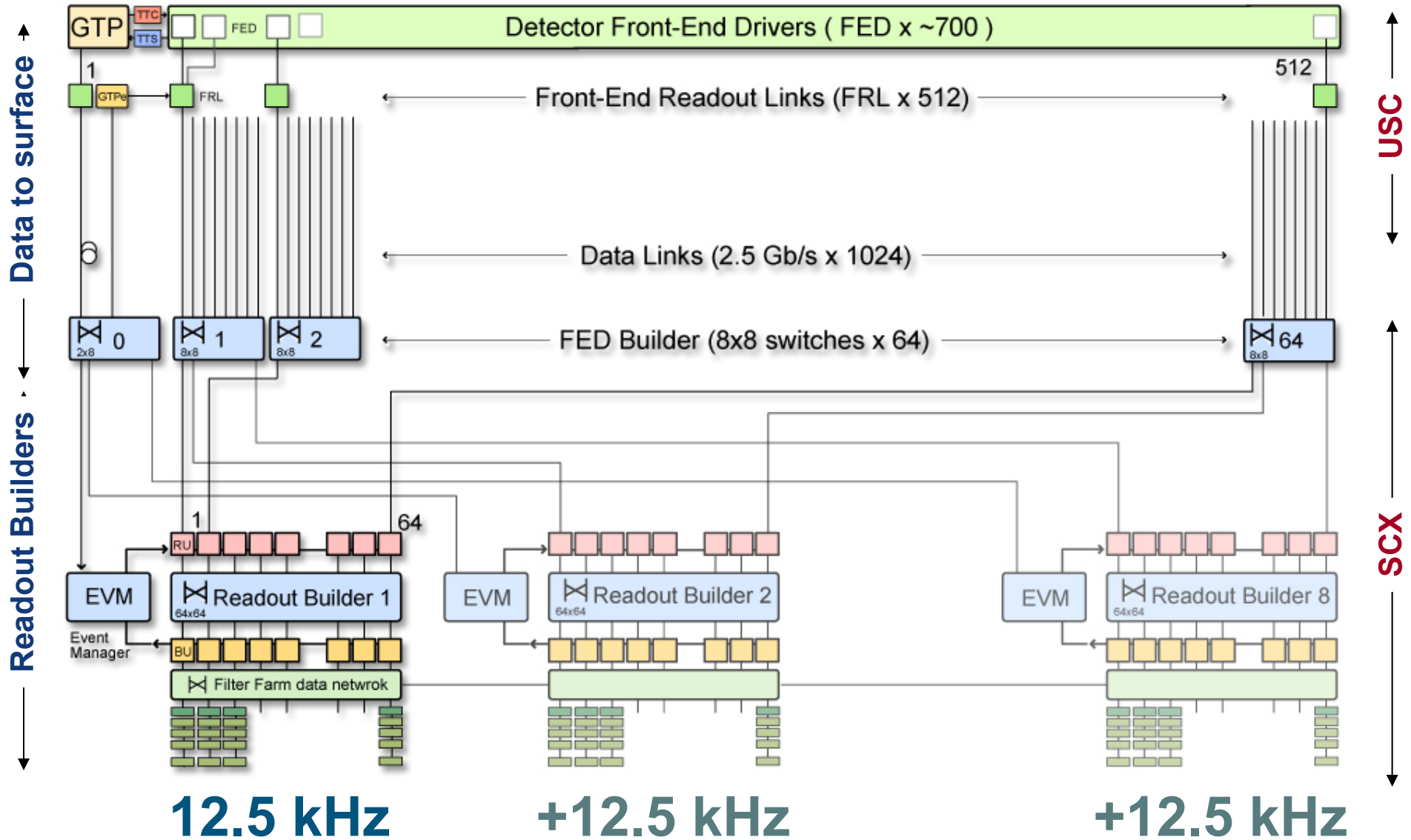


# Builders protocols (e.g. Myrinet)



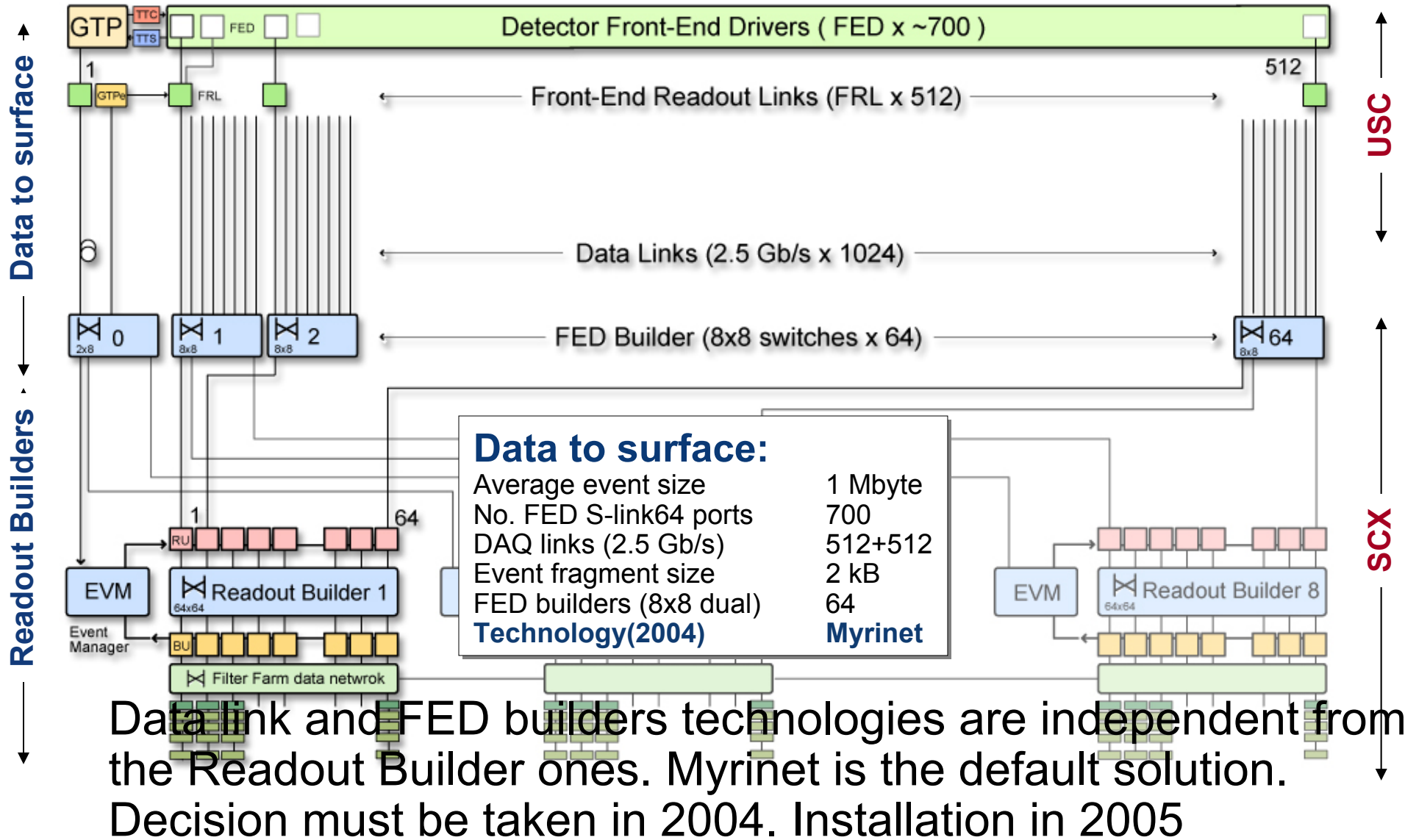


# DAQ staging : 1 to 8 RBs = 100 kHz





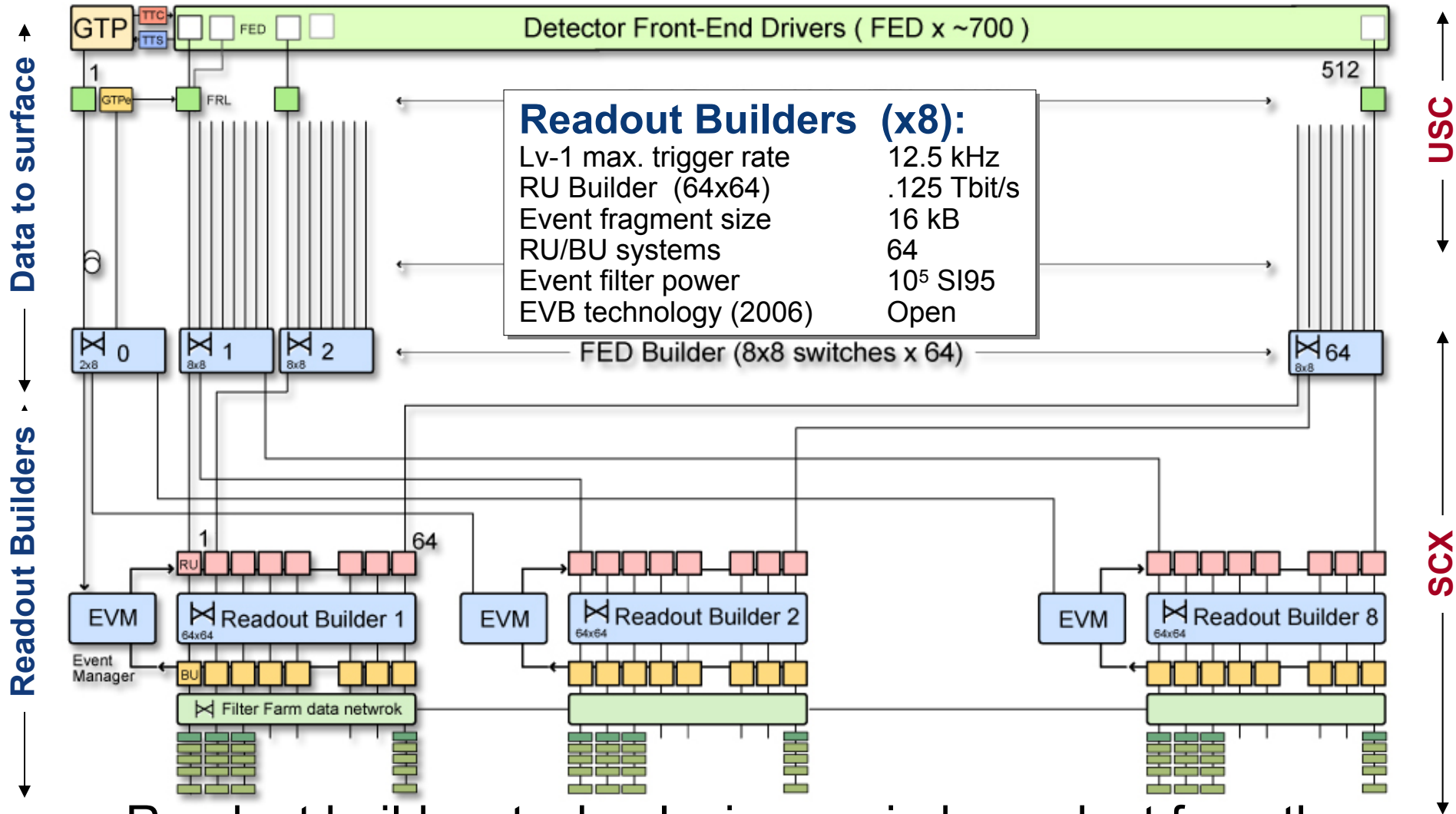
# I) Data to surface (D2S)



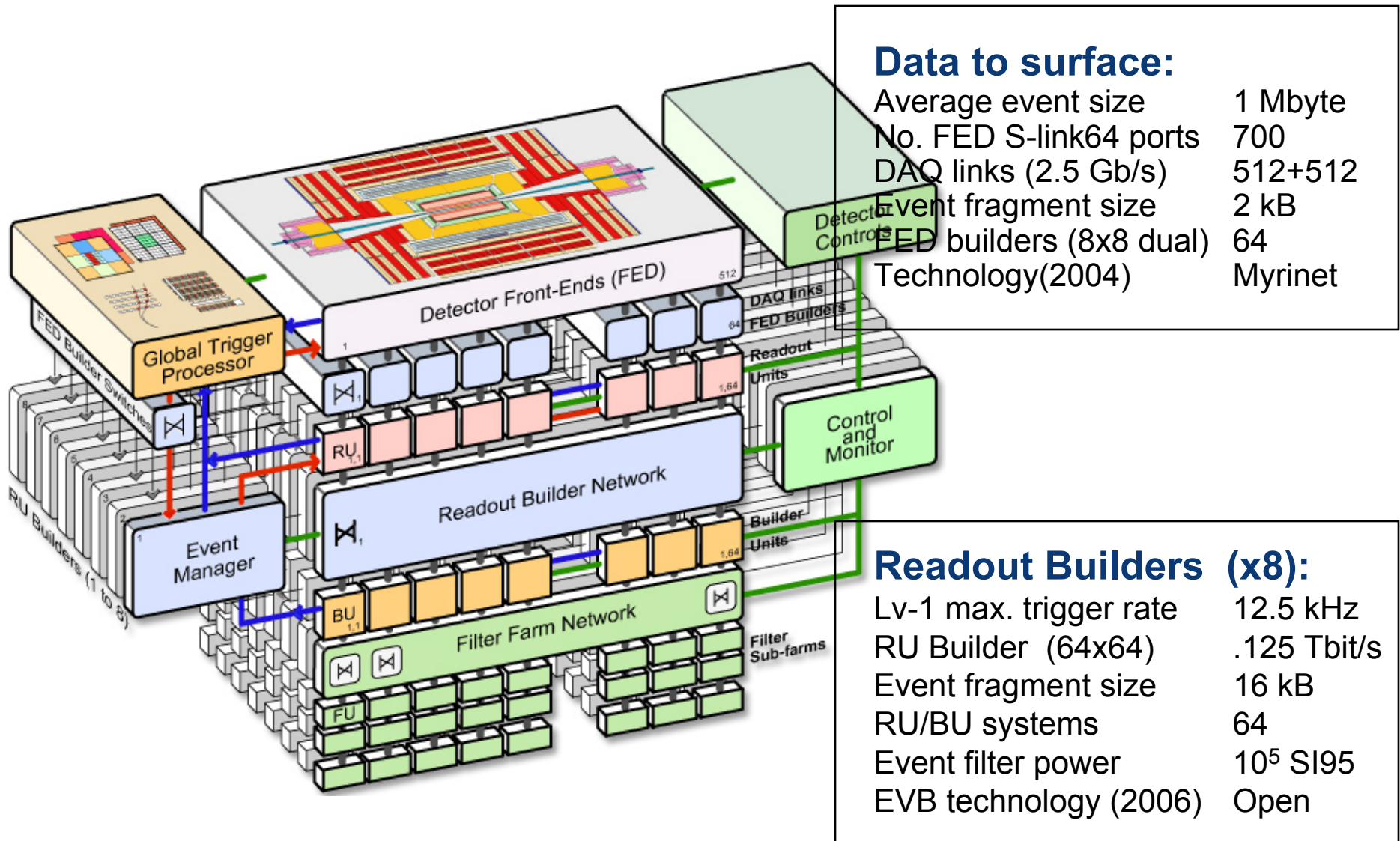


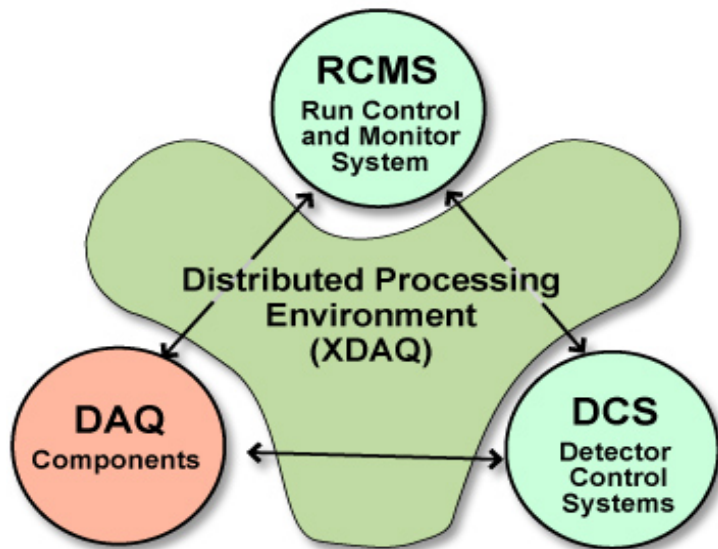


# II) Readout Builders (RB)



Readout builders technologies are independent from the D2S one. Decision will be taken later in 2006





## Online software architecture :

- Cross-platform DAQ framework: XDAQ
- Data acquisition components
- Run Control and Monitor System (RCMS)
- Detector Control System (DCS)

The RCMS, DCS and data acquisition components interoperate through a distributed processing environment called XDAQ (cross-platform DAQ framework)

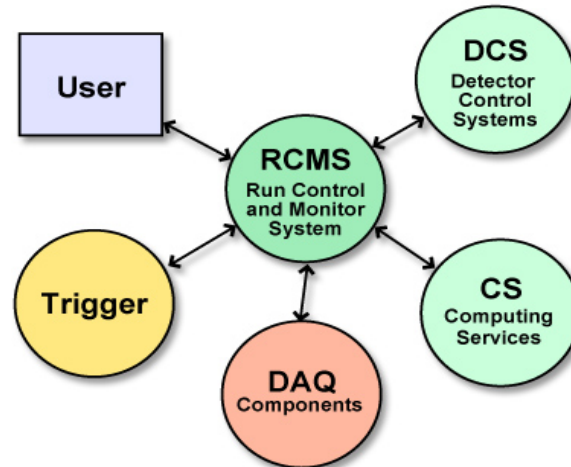
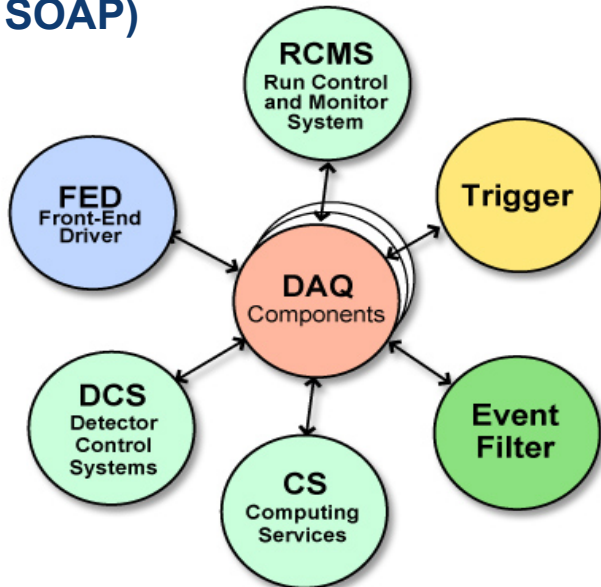


# Configuration, operation and monitoring



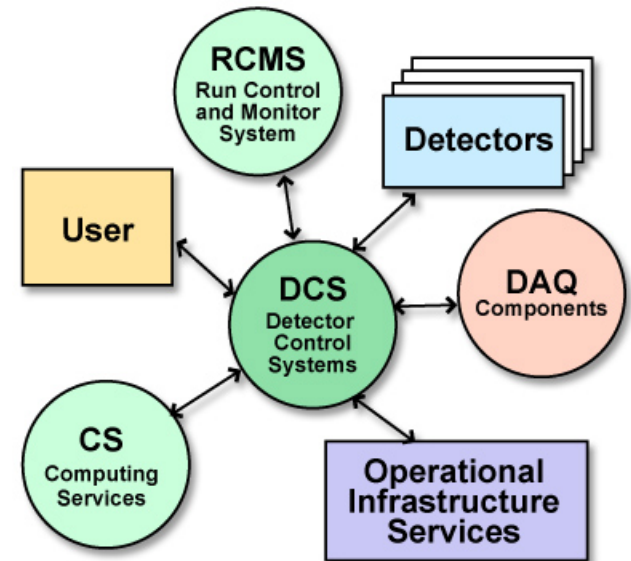
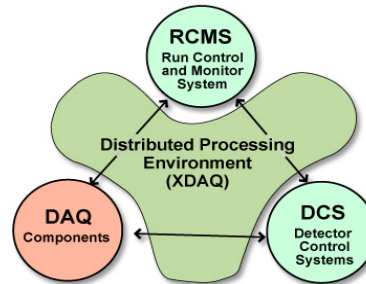
## XDAQ: on-line framework and DAQ components:

- Services and tools for local and remote inter-process communication, configuration and control and data storage
- Components to build data acquisition systems (RU, BU, EVM,..)
- (C++, JAVA, I2O, http, XML, SOAP)



## RCMS: Run Control and Monitoring System

Based on open protocols, web services and emerging e-tools tools (JAVA, http, XML, MySQL, ....)

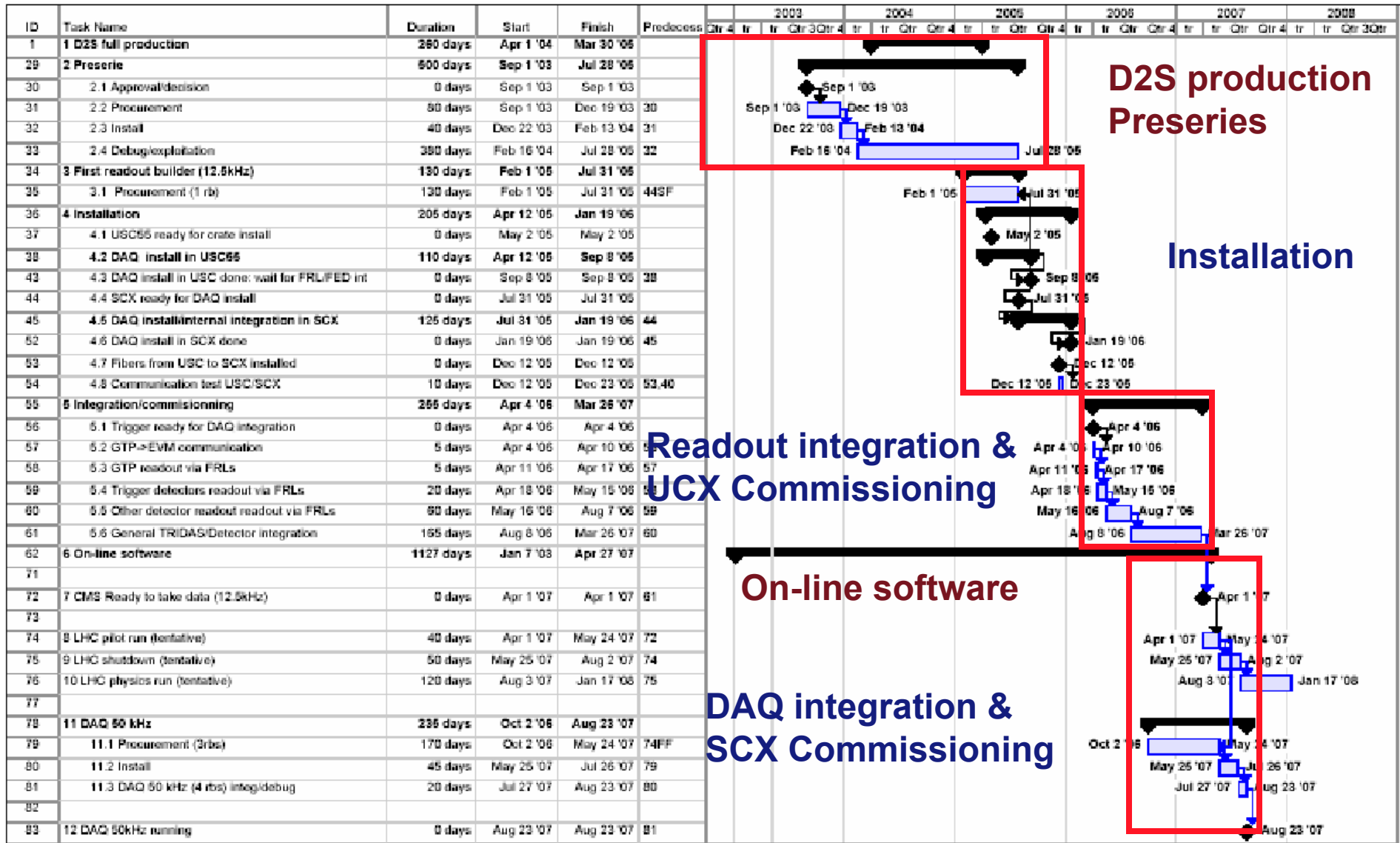


## DCS: Detector Control Systems

Based on industry supported hardware and software (PLC, field buses, PVSS and JCOP tools)



# DAQ raw schedule

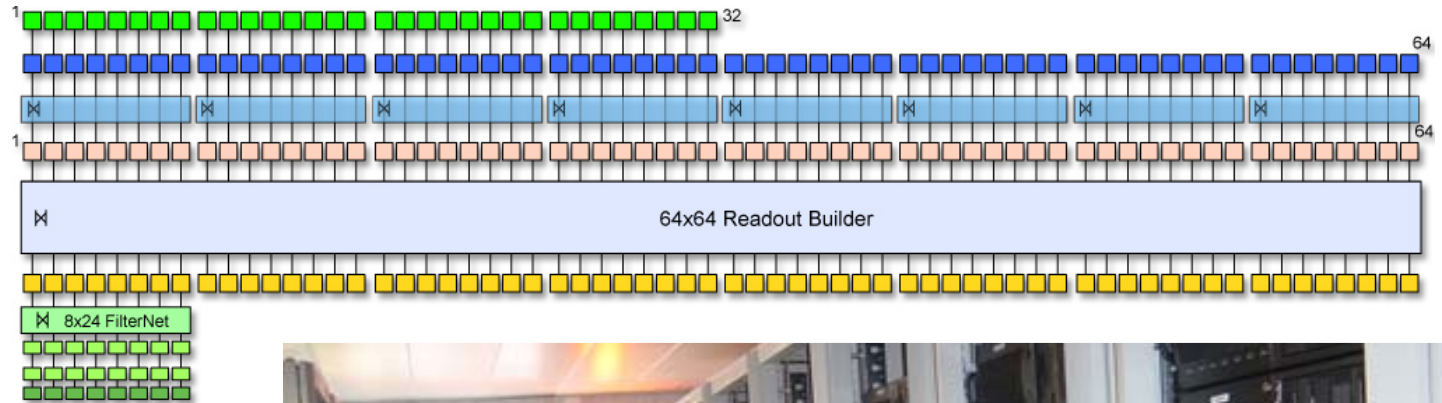




# 2004 P5 green barrack: Pre-series RB



- 32 GII-FED emulators
- 64 FRLs
- 13 Water cooled racks
- 93 PC dual-CPU
- D2S Myrinet equipment
- Readout Builder Myrinet
- 16 PC Filter Farm





# Preseries integration programme



32 GIII FED emulators

Up to 32 detector FED

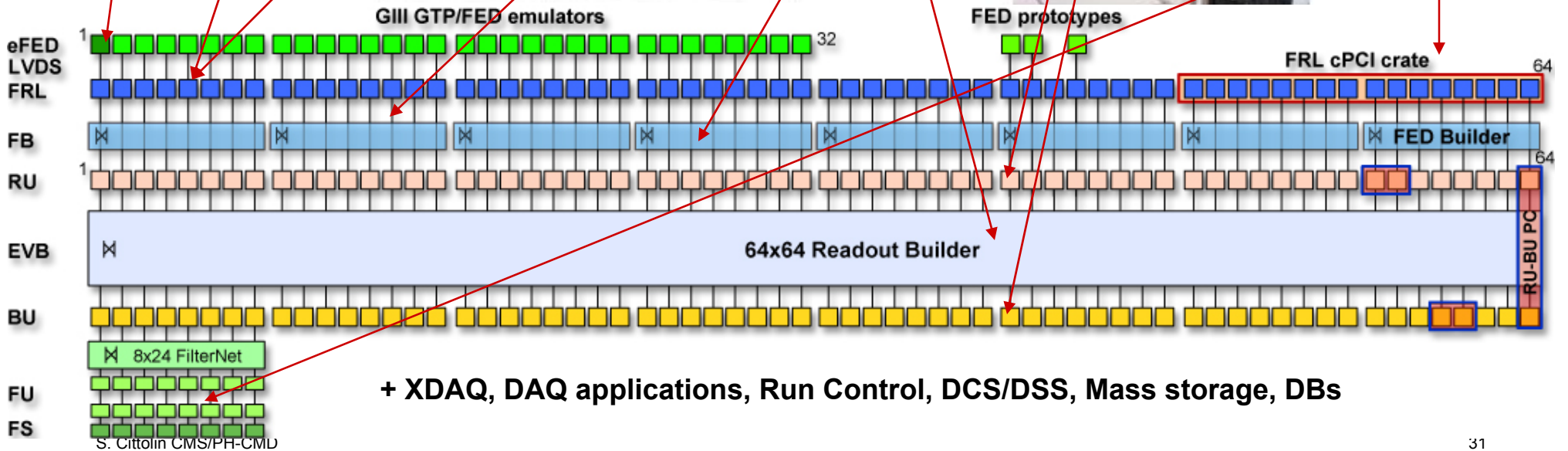
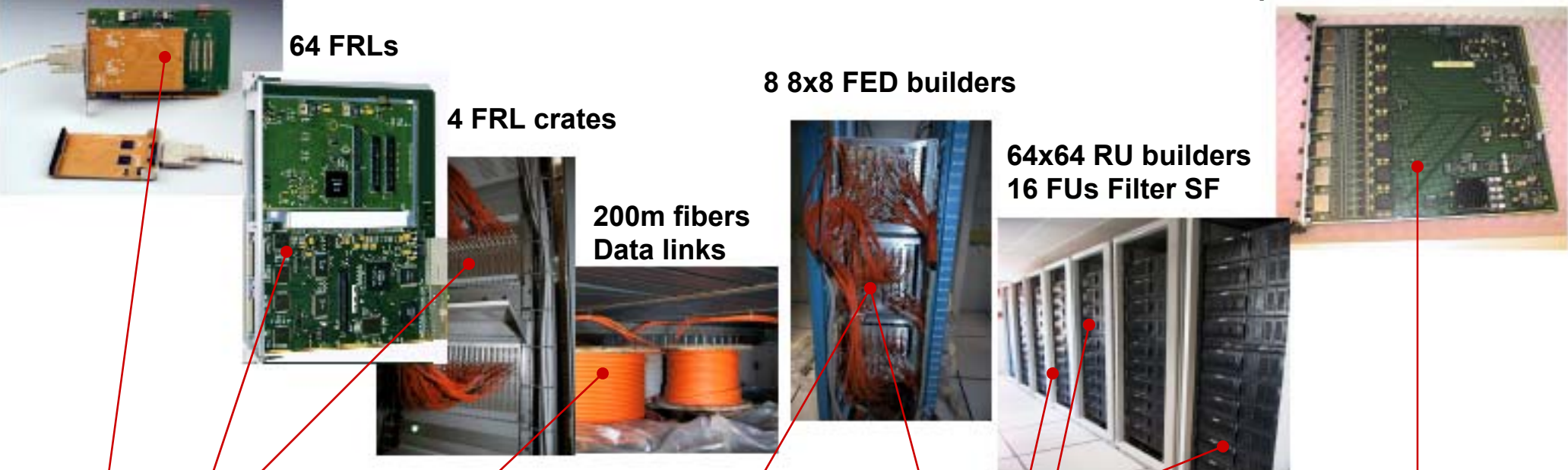
64 FRLs

4 FRL crates

8 8x8 FED builders

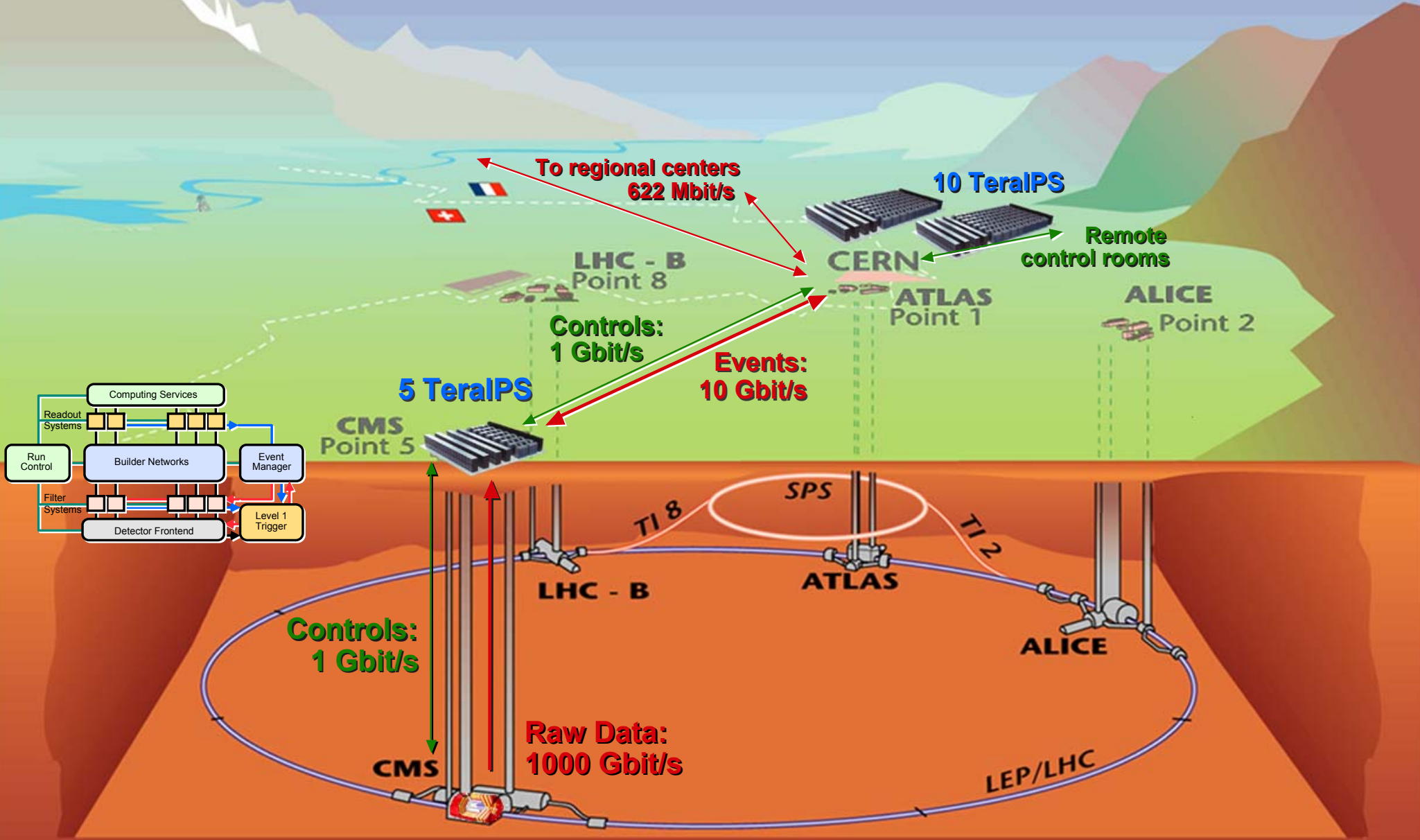
64x64 RU builders  
16 FUs Filter SF

200m fibers  
Data links



+ XDAQ, DAQ applications, Run Control, DCS/DSS, Mass storage, DBs

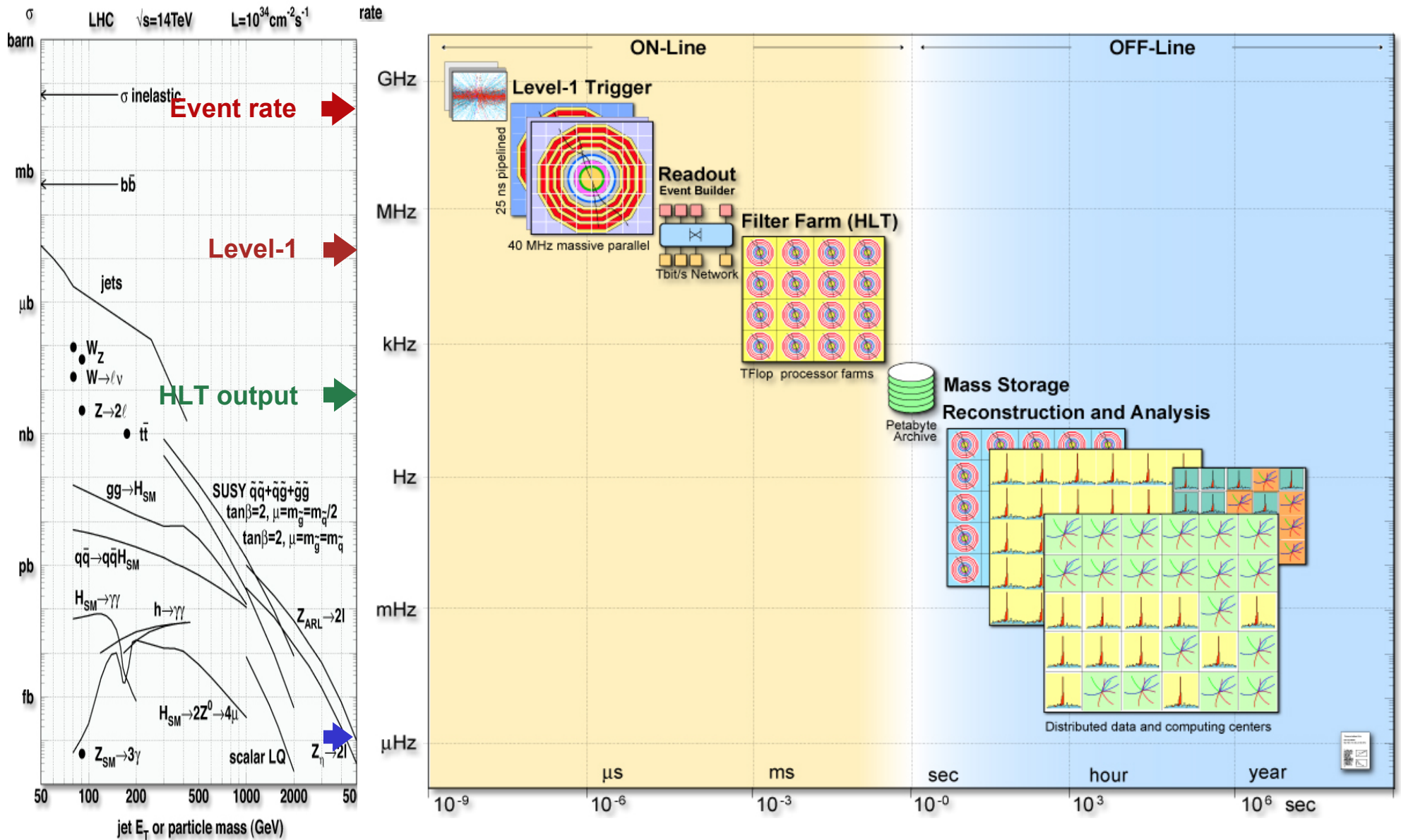
# CMS data flow and on(off) line computing







# DAQ data flow and computing model





# Summary design principles



- **Invest in the advance of communication and processing technologies**
  - Computing (**100 kHz Readout, HLT by PC farms**)
  - Communication (**Terabit/s networks, GB/s memories**)
- **Maximally scaling architecture**
  - Exploit technology **evolution**
  - Cost optimization via **staged installation**
  - Modular system** (simpler controls, error handling, smaller basic units)
- **Rely on hardware and software industry standards**
  - Custom/**standards** (PCI, Ethernet, C++, JAVA, http, XML,..)



# Conclusions



The CMS design fulfils the major requirements:

- ✓ **100 KHz level-1 readout**
  - ✓ **Event builder:**
    - a scalable structure that can go up to 1 Terabit/s**
  - ✓ **High-Level Trigger by fully programmable processors**
- 
- This **design** should be considered **complete**, but not final.
    - e.g. switches procured in 2008-09 can be different from those of the startup system
  - It is a **system** that is **expected to change** with time, accelerator and experiment conditions. And it has been designed to do so
  - It is conceived to provide the **maximum possible flexibility** to execute a physics selection on-line