



Physics at the high energy frontier - the Large Hadron Collider project

London 17 May 2011
Sergio Cittolin
ETHZ-UCSD-CERN



- **Introduction**

- Data handling requirements at LHC

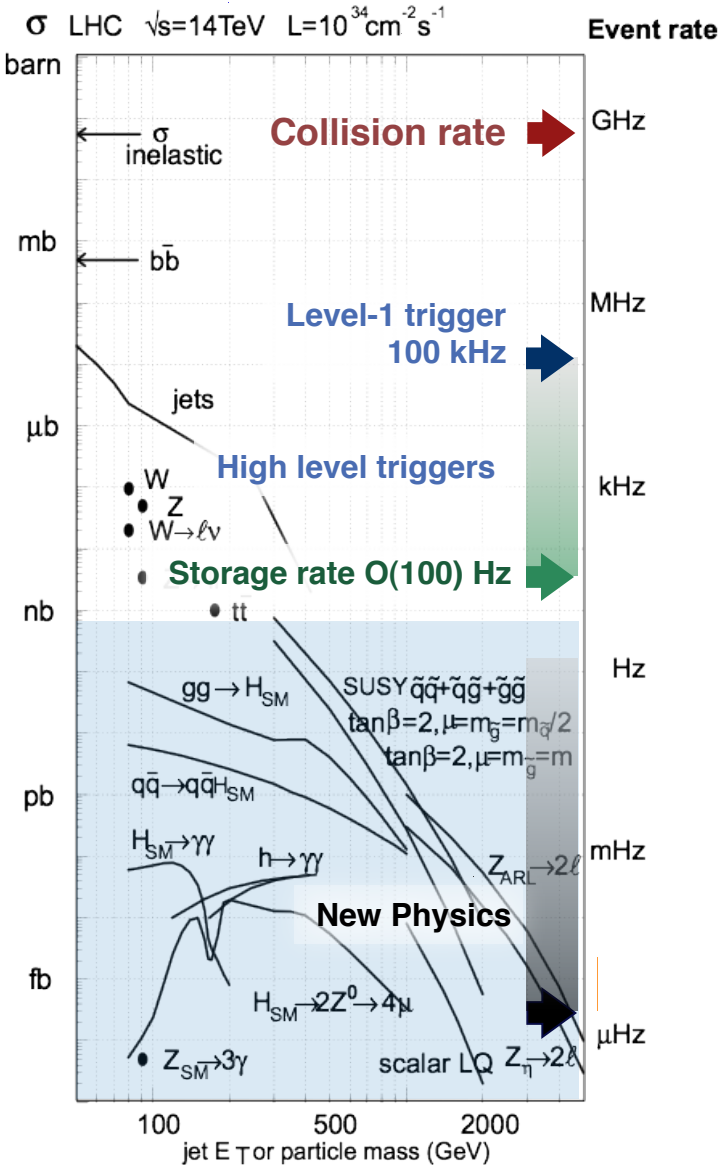
- **Design issues: Architectures**

- Front-end, event selection levels

- **Design issues: Technologies**

- Project history and technologies trends
- Predicted and unpredicted evolutions

- **Conclusion**



At the highest LHC beam intensities:

- **Observe every second** 40 million bunch crossings, each producing several (>20) p-p interactions resulting in events with 1000's particles every bunch crossing
- **Identify** and select single events out of 10 Trillion collisions
- **Locally digitize**, readout, transport and process 100's of TeraBits per second

- **Globally store**, retrieve and analyze efficiently tens of PetaBytes of data per year

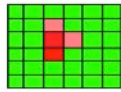
Collision rate	$\sim 10^9$ Hz
Detector granularity	$\sim 10^8$ sensors
Event size	~ 1 Mbyte
Event selection power	~ 1 in 10^{13}
Data readout bandwidth	\sim Terabit/s
Storage event rate	$\sim O(100\text{Hz})$
On/Off-line Processing power	\sim TeraFlops

Design issues: Architectures

- Data flow, Front-end, Event selection levels

Crossing rate: 40 MHz
Collision rate: 1 GHz

Raw data production
100's of Petabits/s



Energy



Tracks

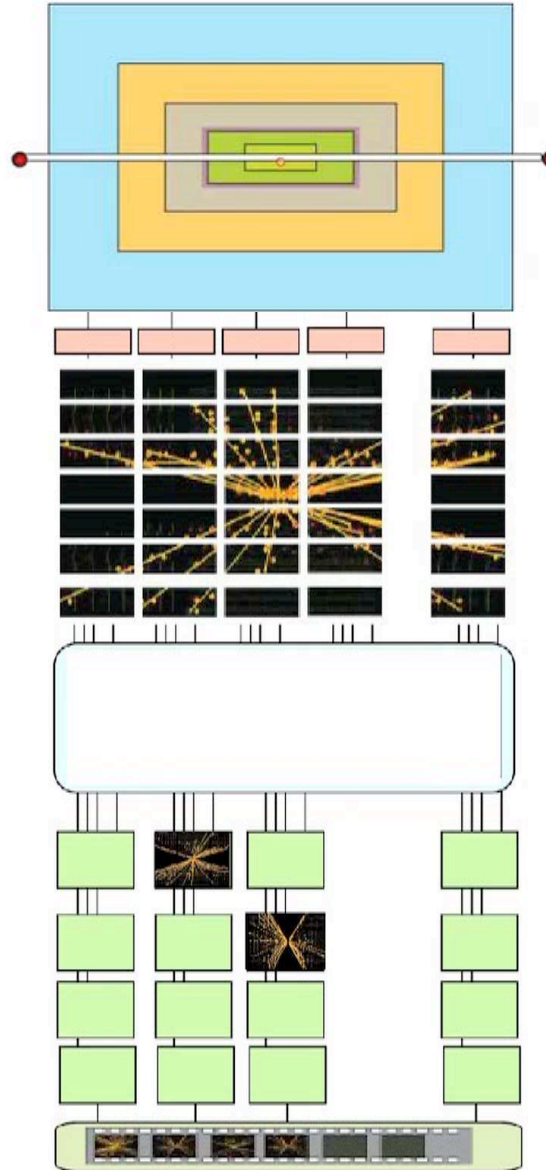
Level-1 Trigger

100 kHz output rate
50 Million fragments/s

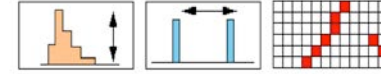
Readout network
2 Terabit/s

Build, process and select
10 TeraFlops
100000 event/s

Store O(100 Hz)
Tens of Gigabits/s



100 Millions instrumented sensors



Charge

Time

Pattern

Billions of (A/D) memory cells

Parallel readout

Hundreds of event fragment readers

Data to Surface

Thousands of Optical links

Event Builder

Switching system with thousands of ports
~1 Megabyte data per event

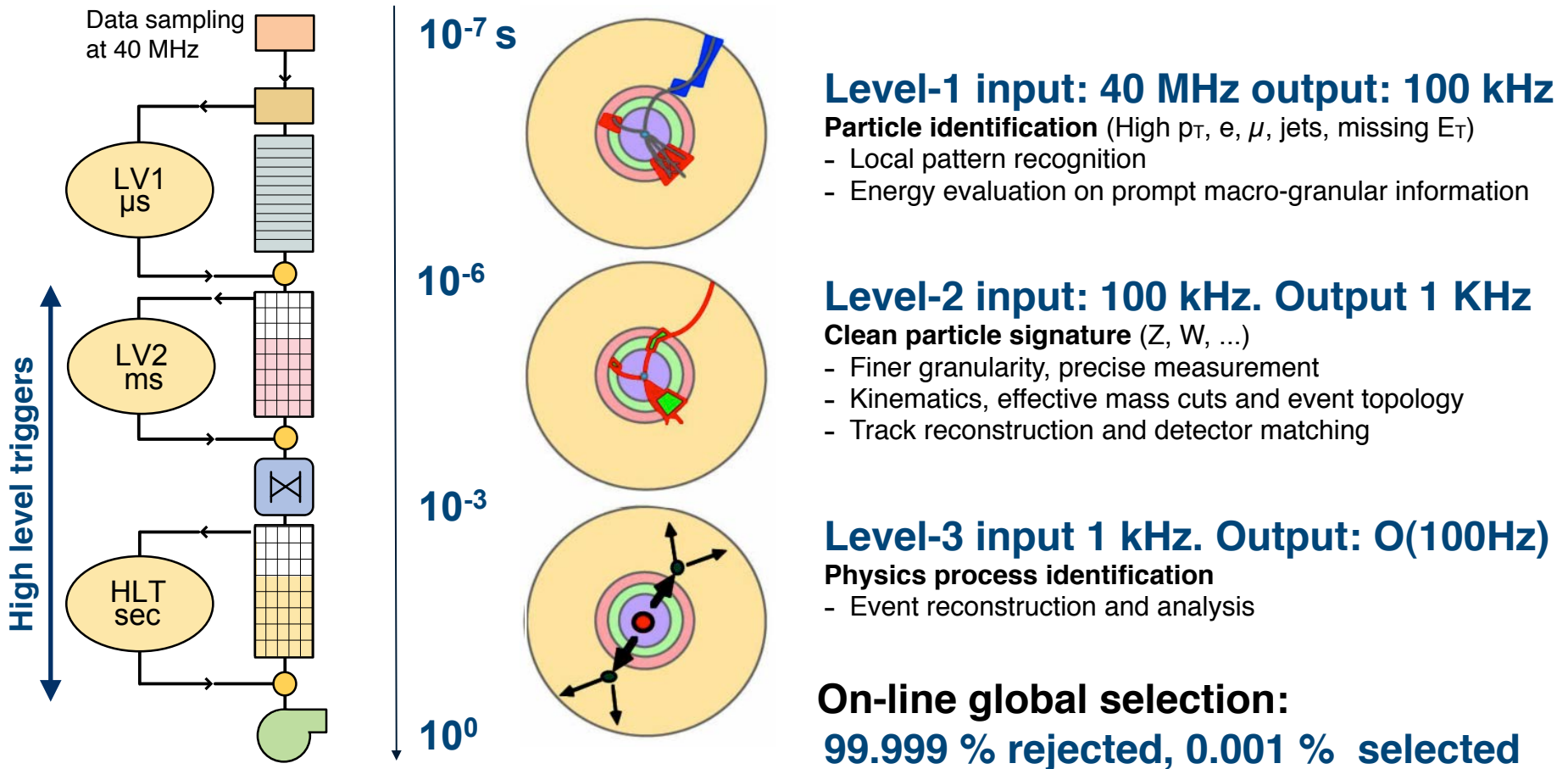
Event Filter

Thousands of CPU cores

Local mass storage

Hundreds of Terabits

Successively more complex decisions are made on successively lower data rates



Readout and trigger dead-time must be kept at minimum (typically of the order of few %)

The trigger system has to maximise the collection of data for physics process of interest at all levels, since **rejected events are lost for ever**

$Z \rightarrow \mu\mu + 2 \text{ jets}$

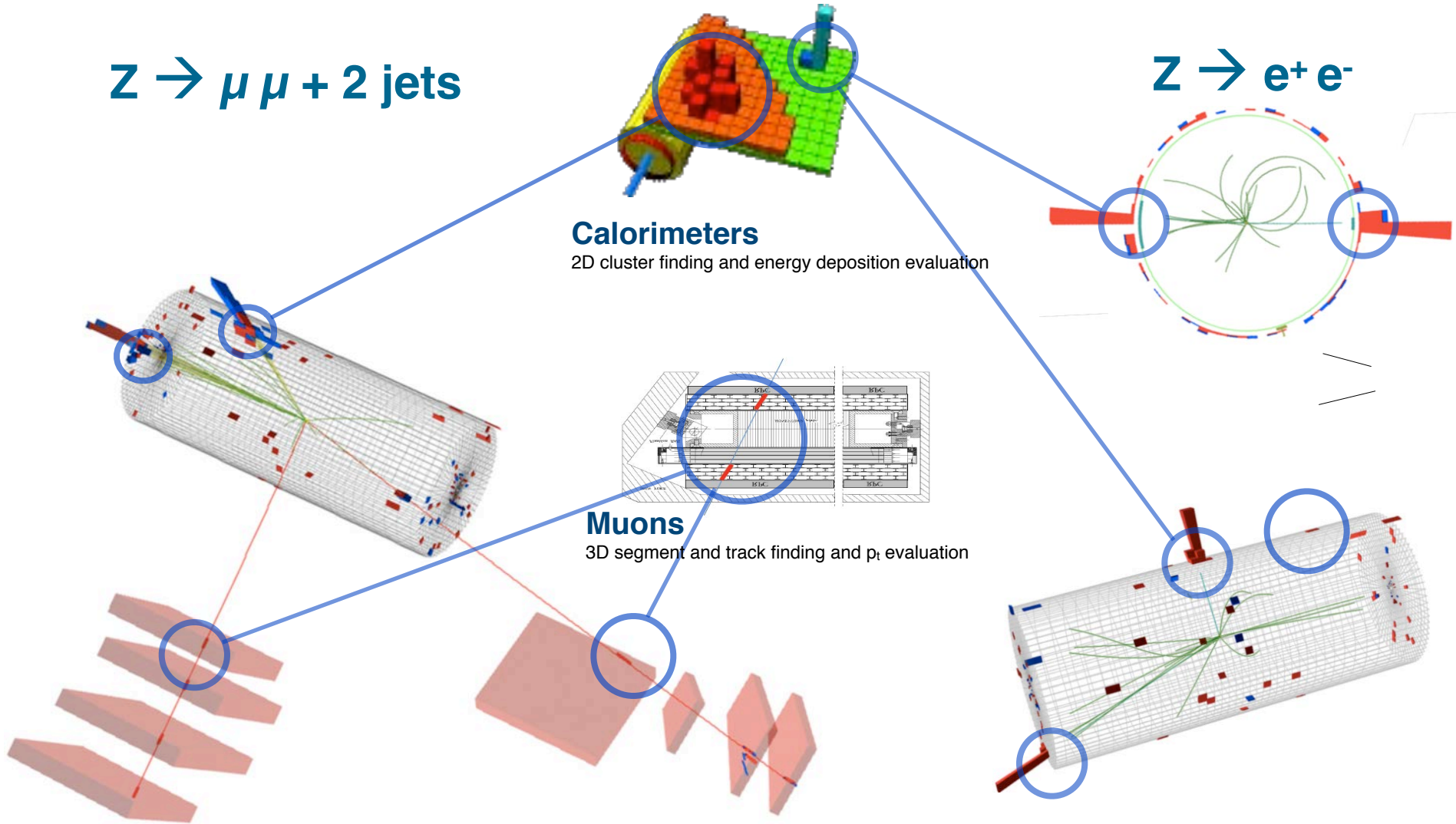
$Z \rightarrow e^+e^-$

Calorimeters

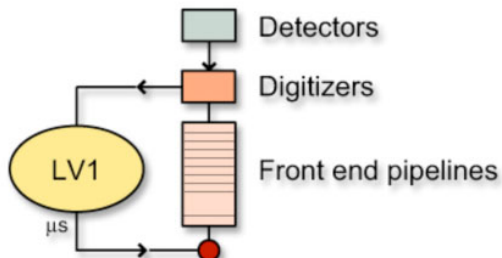
2D cluster finding and energy deposition evaluation

Muons

3D segment and track finding and p_T evaluation



Algorithms run on local **calorimeter and muon coarse data**. With **new data every 25 ns** and **decision latency $\sim \mu\text{s}$** . Special-purpose hardware reduces event rate (to be read out) **from 40 MHz to 100 kHz**.



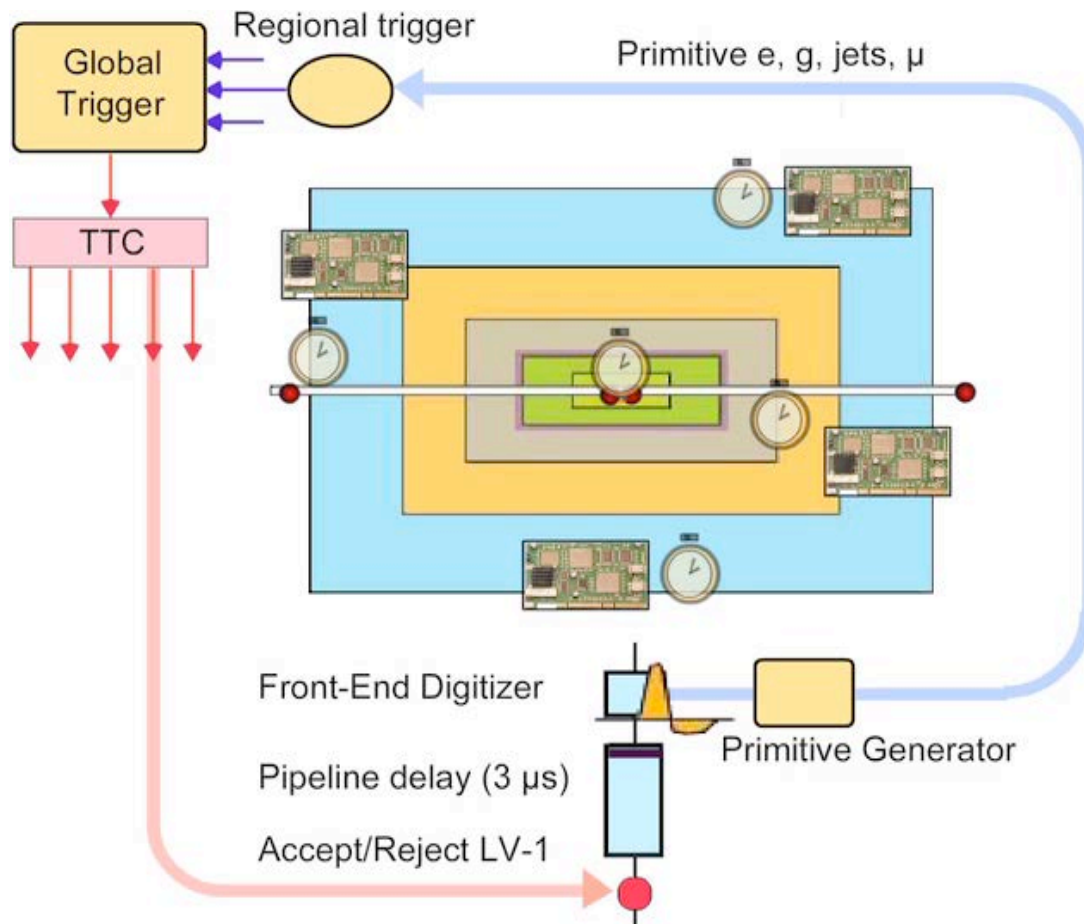
The **Front-end** system includes preamplifier, shaper, digitizer and the buffers to hold the data for the duration of the level-1 trigger decision and the signal propagation delay ($\sim 3 \mu\text{s}$)

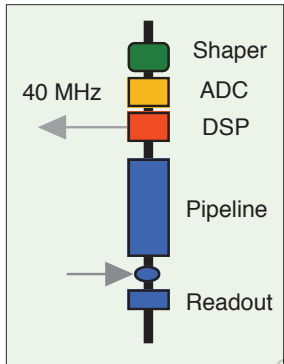
The **Level-1** trigger is a centralized multi-step logic system collecting, buffering and processing sub-set of detector data every 25 ns.

TTC. A multichannel **optical distribution system broadcasts the LHC 40 MHz clock** and the Global Trigger signals to several thousand destinations

Front-end time budget (128 Bunch crossings)

Signals to central logic	18bx	450 ns
Muon/Calorimeter logic	90bx	2250 ns
Back to detectors	18bx	450 ns
total latency	$\leq 128\text{bx}$	$3.2 \mu\text{s}$

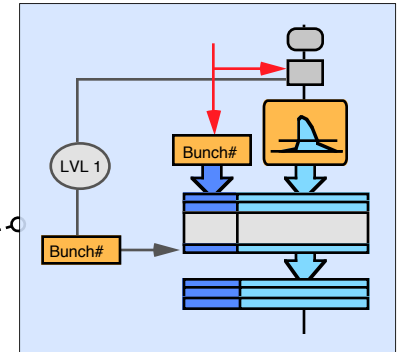
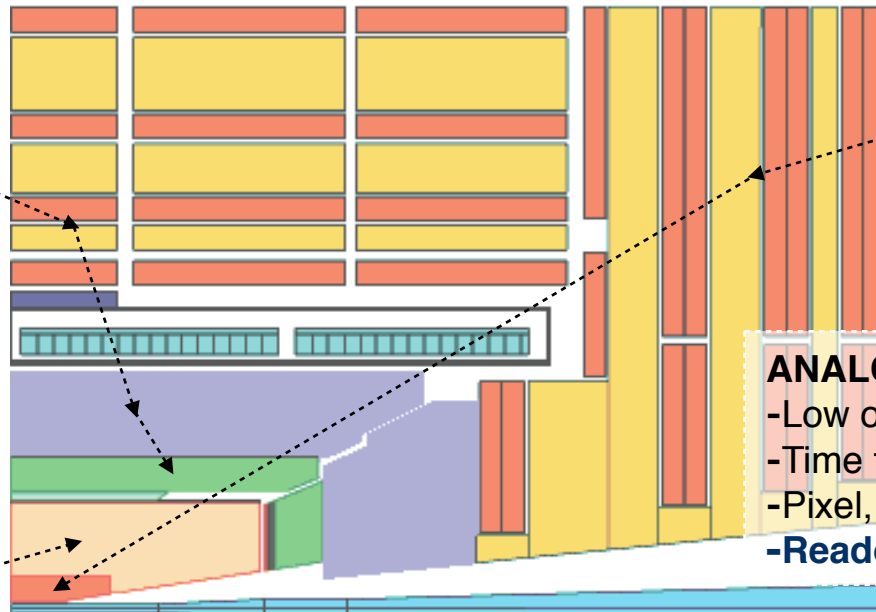




500.000 ch.

DIGITAL (25 ns) synchronous pipeline

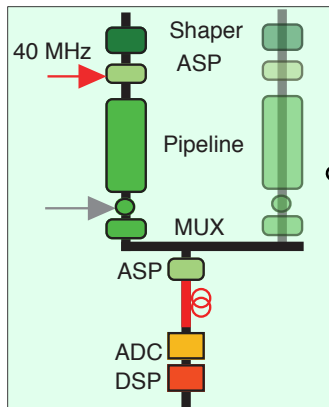
- Large dynamic range (15bits)
- Digital filtering (high power consumption)
- ECAL, HCAL Calorimeters, DT, RPC Barrel Muons
- Readout: OPTICAL digital links**



60.000.000 ch.

ANALOG/DIGITAL asynchronous buffer

- Low occupancy. High No.channels
- Time tag identification system
- Pixel, CSC forward muons
- Readout: OPTICAL&COPPER**

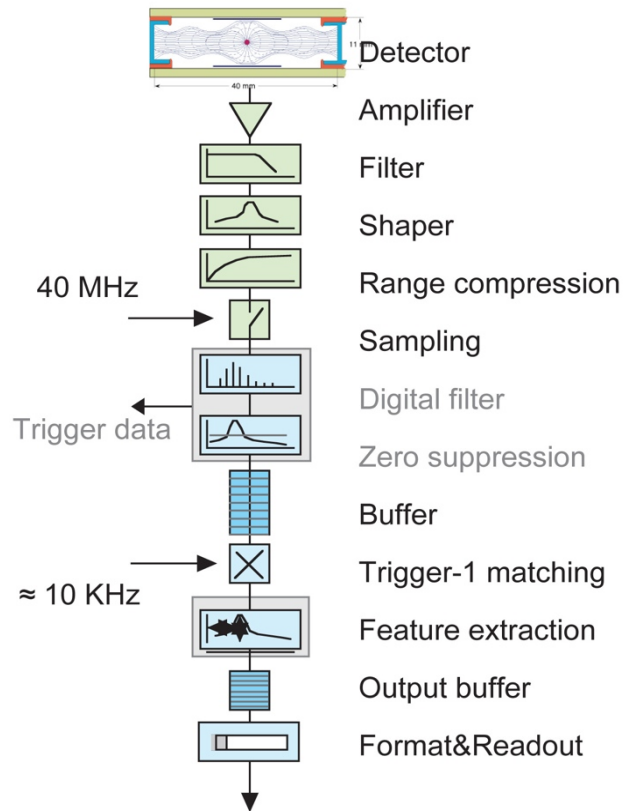


10.000.000 ch.

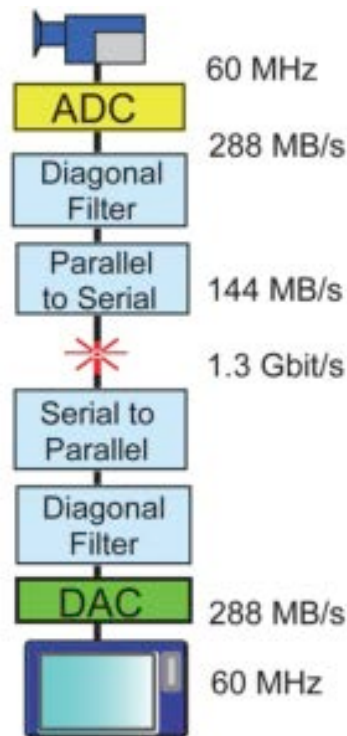
ANALOG (25 ns) synchronous pipeline

- Low dynamic range and resolution (≤ 10 -bits)
- Low power consumption, radiation tolerant,
- High number of channels. Tracker, Pre-shower
- Readout: OPTICAL analog links**

1990. LHC detector channel



1990. HDTV chain



One HDTV = One LHC channel

Analog bandwidth ~ 100 MHz
 Digital resolution 12_14 bits
 Digital bandwidth ~ 1 Gb/s

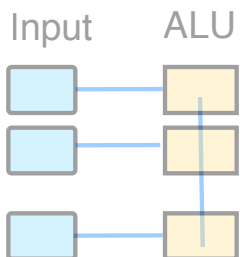
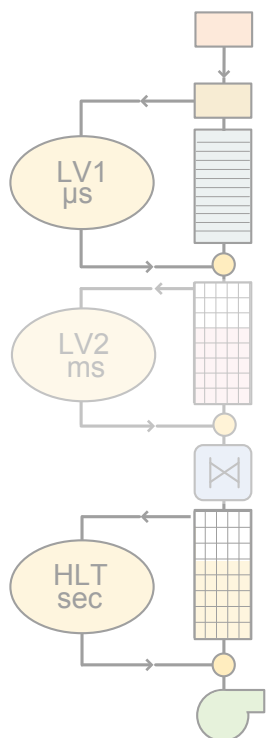
Since early 80's:

- Digital Signal Processing (**DSP**) has become pervasive at all levels in our society.
- **DSP** as a technology has become the primary growth driver for the entire semiconductor market.
- The telecommunication industry has been one of the major customers for the development of this technology.
- Analog to digital converters (**ADC**)
- Multiply accumulator (**MAC**)
- GHz **optical links** and Laser LED
- Finite Impulse Response (**FIR**) digital filters and vector processing are today the **building blocks of any LHC detector readout chain** as well.

Level-1 trigger architecture. Massively parallel: One event -> Multi-processors

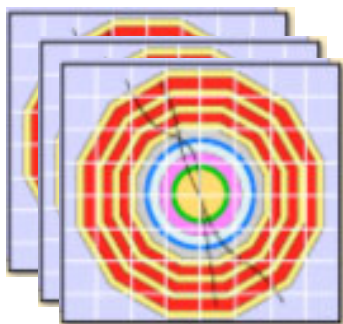
- High (fixed interconnections), Short (fixed) latency. Pipelined simple ALUs, **Clock driven**

MEMORY (I/O) access



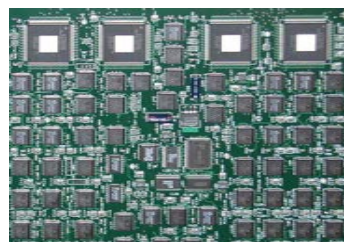
DATA PROCESSING

Synchronous pipelined



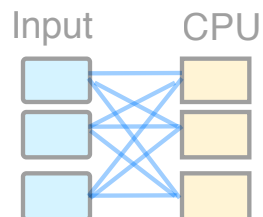
IMPLEMENTATION

Custom design
(ASIC, FPGA, LVDS..)

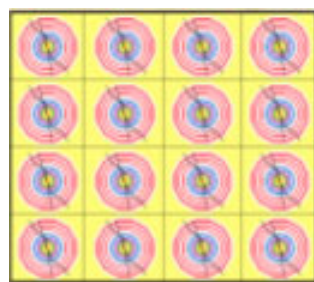


INTERCONNECTION

Distributed



Asynchronous clusters



Commodities
PCs and networks

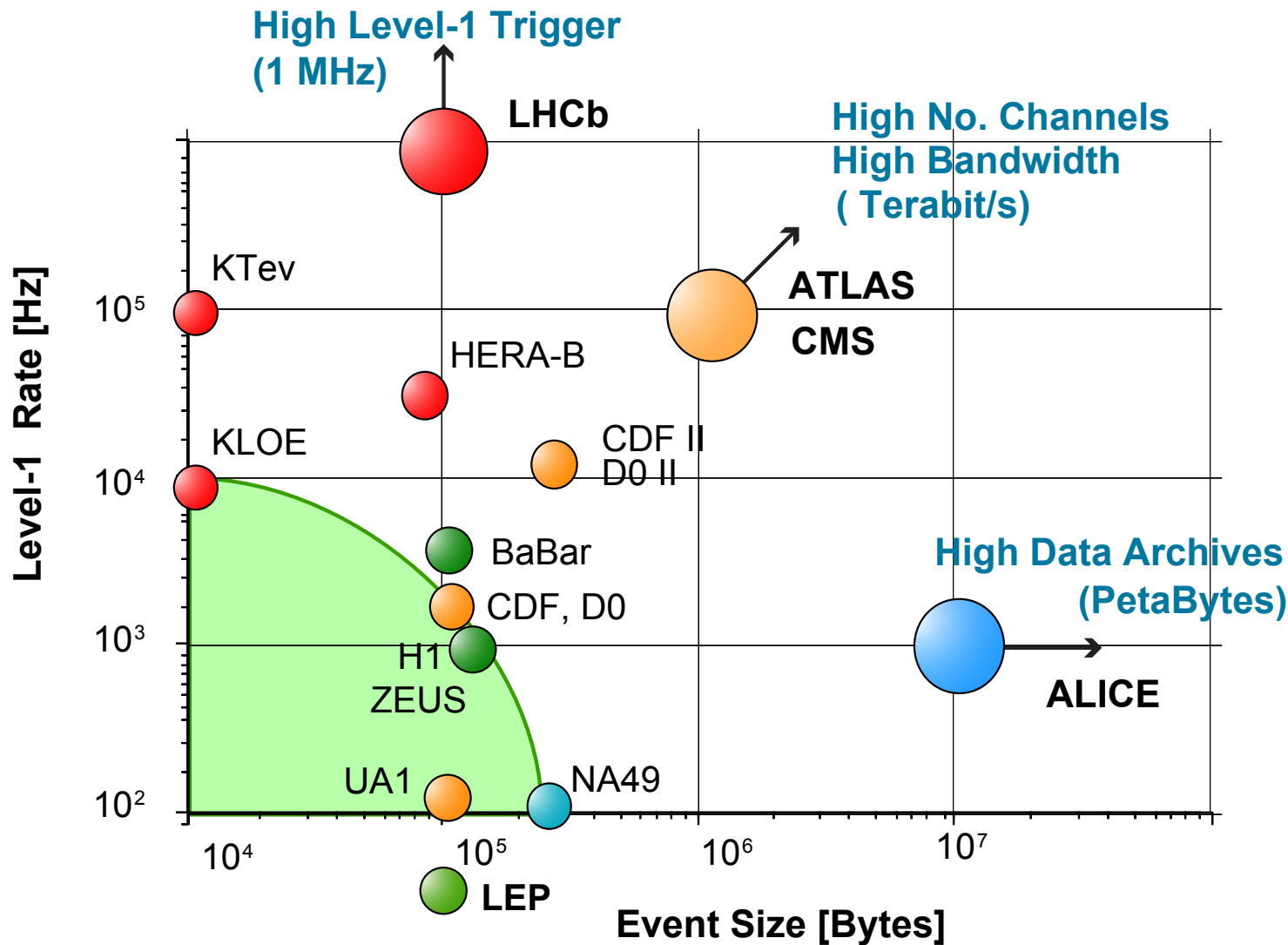


Decentralised



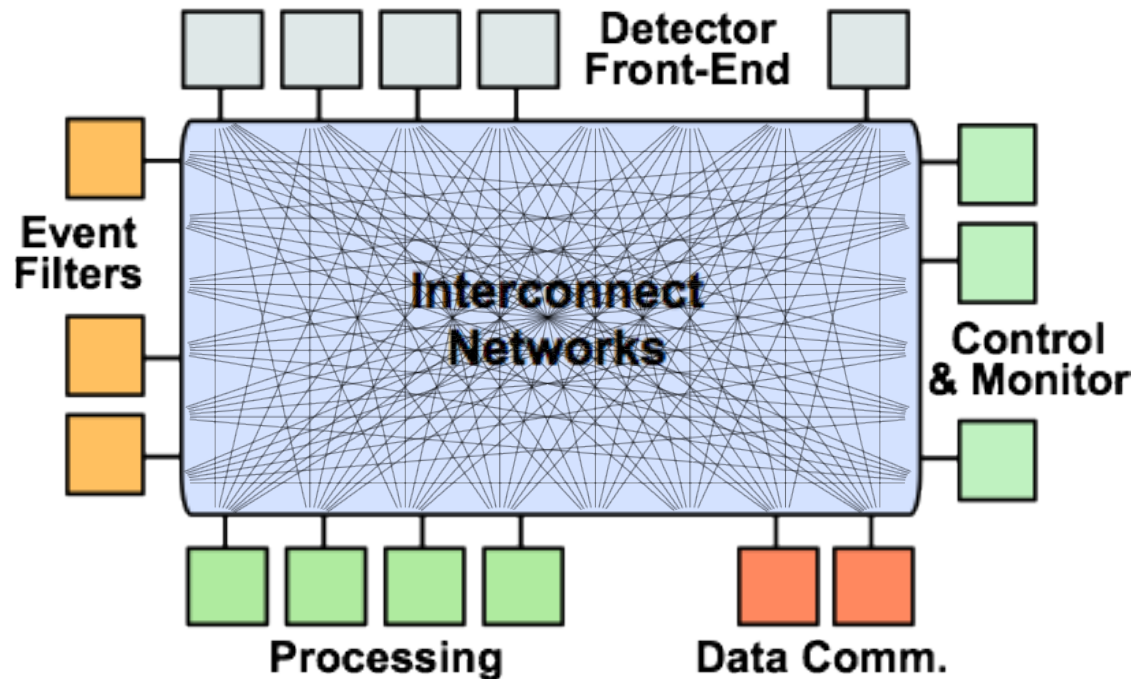
High Level Triggers architecture. Cluster structure: One event -> One processor

- Loose coupling, large latency. Node high power, **Event driven**

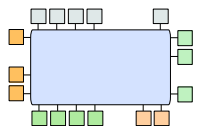


2000's On&Off-line processing and communication model

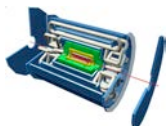
Consists of buffer memories, processors, communication links, data-flow supervisors, storage and data analysis units. Conceptually, the On/Off-line systems can be seen as a global **network interconnecting all** the data-flow, control and processing units



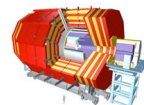
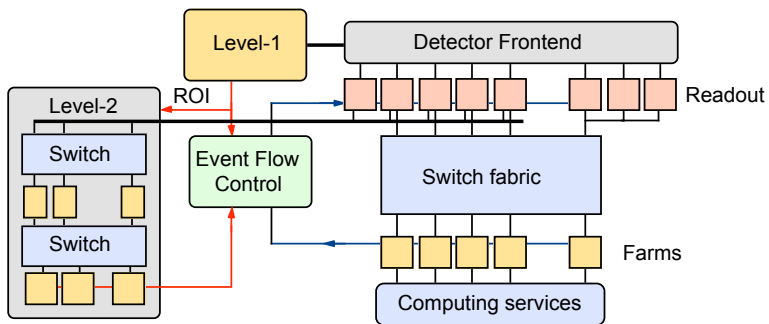
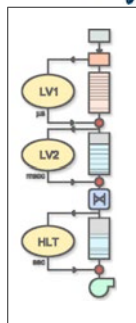
At the time of the finalization of the system design (2002-03), **a single network technology could not satisfy at once all the LHC requirements. The LHC DAQ designs had to adopt multiple specialized networks instead.**



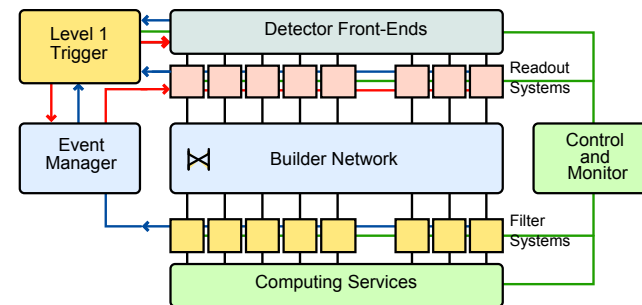
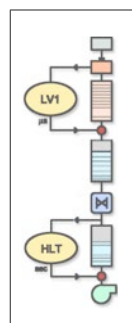
Each LHC experiment developed its own scheme to cut the rate, to process events online and/or optimize the throughput. In a sense, the systems designed and built are “approximations” of the basic architecture/conceptual design



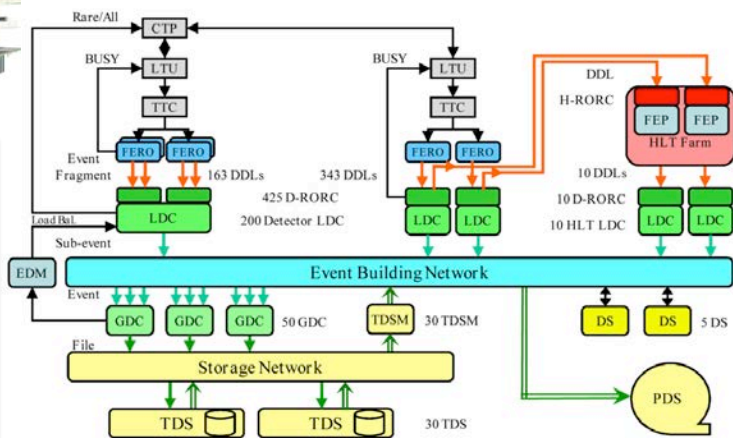
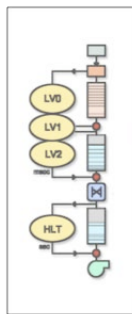
ATLAS



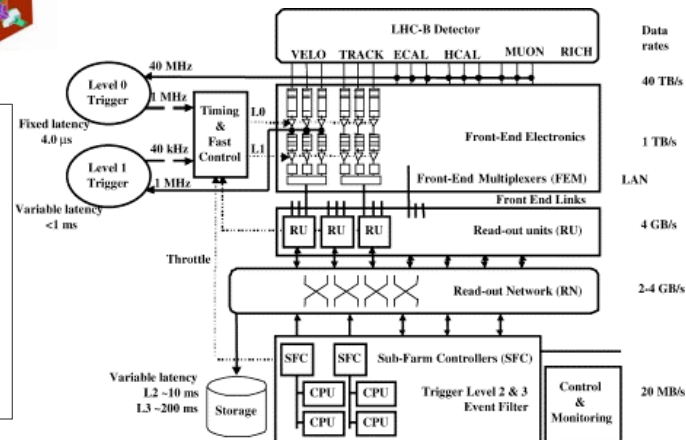
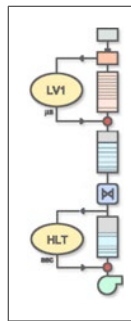
CMS



Alice



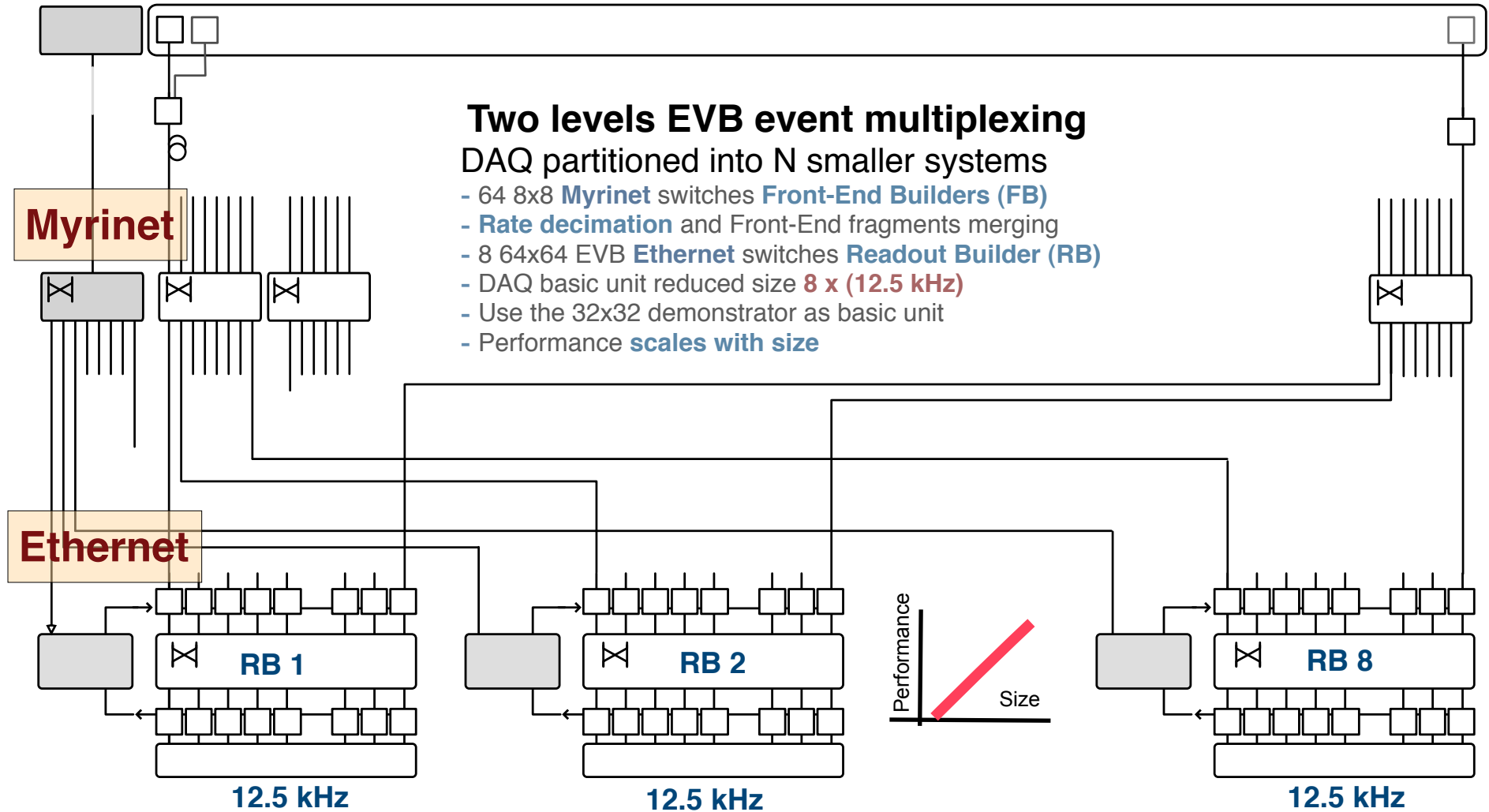
LHCb



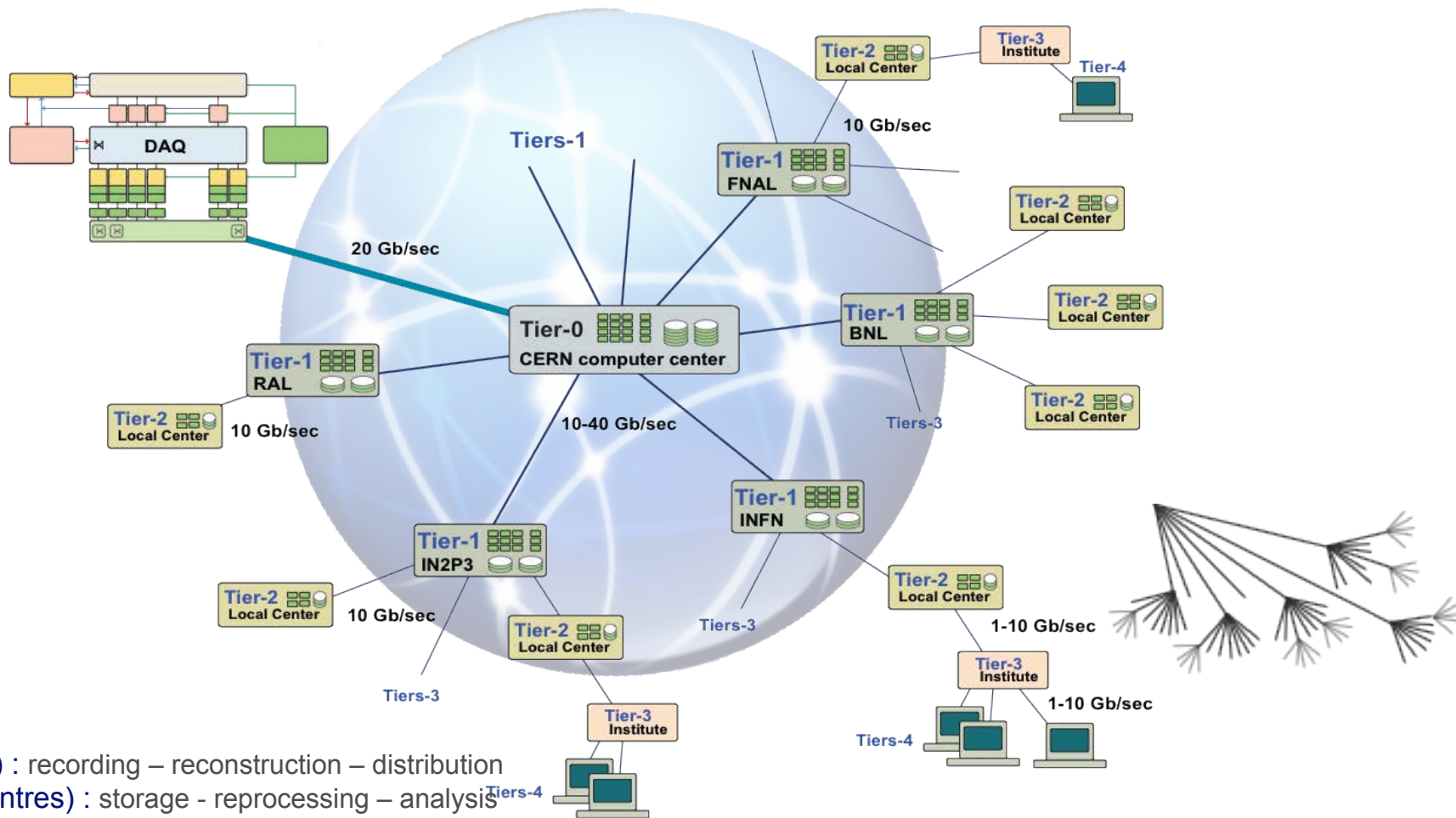
Two levels EVB event multiplexing

DAQ partitioned into N smaller systems

- 64 8x8 **Myrinet** switches **Front-End Builders (FB)**
- **Rate decimation** and Front-End fragments merging
- 8 64x64 EVB **Ethernet** switches **Readout Builder (RB)**
- DAQ basic unit reduced size **8 x (12.5 kHz)**
- Use the 32x32 demonstrator as basic unit
- Performance **scales with size**



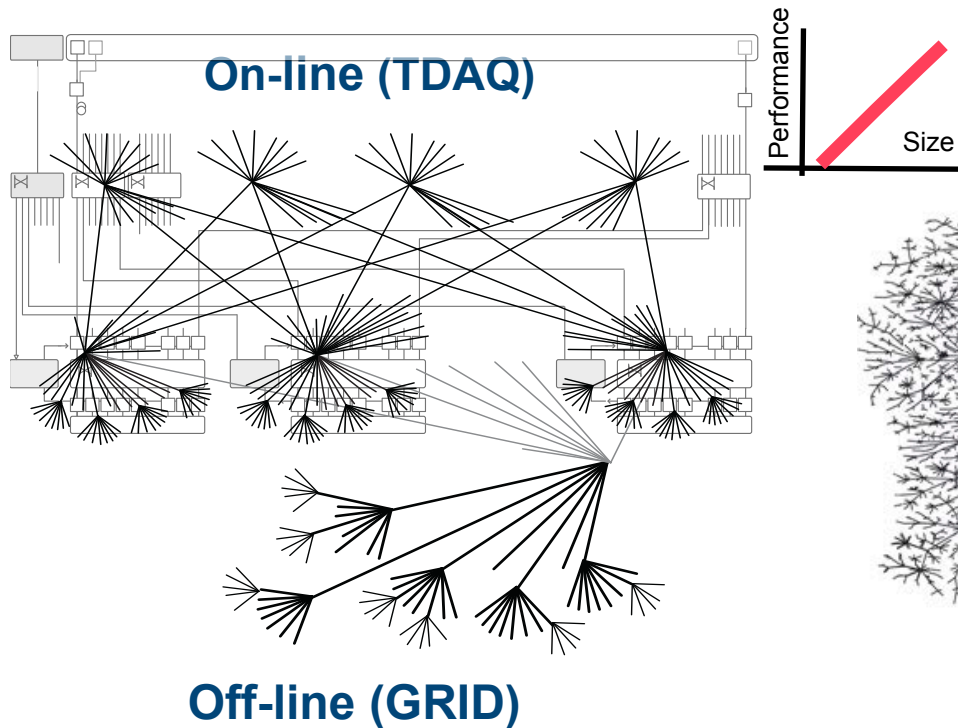
The GRID. A distributed computing infrastructure (~150 kCores, 50 PB disk), uniting resources of HEP institutes around the world to provide seamless access to CPU and storage for the LHC experiments. A common solution for an unprecedented demand (in HEP) of computing power for physics analysis.



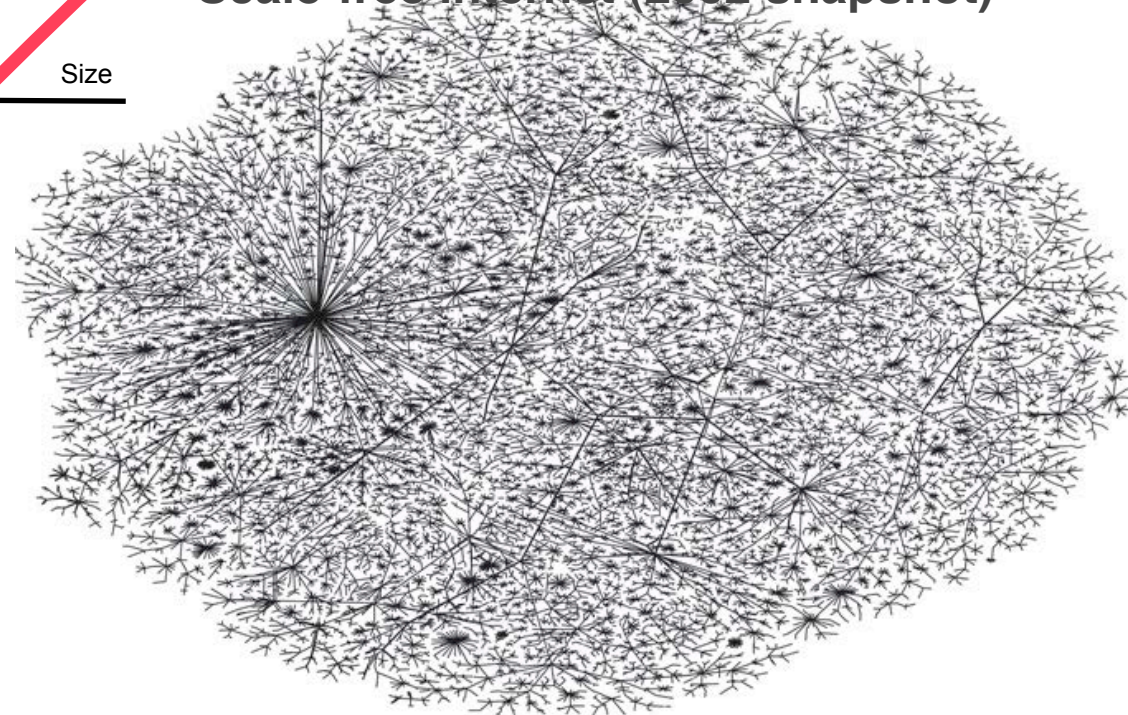
- Tier-0 (CERN) : recording – reconstruction – distribution
- Tier-1 (~10 centres) : storage - reprocessing – analysis
- Tier-2 (~140 centres) : simulation – end-user analysis

On/Off-line **TDAQ (and GRID) systems are, by construction, scale-free systems**; they are capable of operating efficiently, taking advantage of any additional resources that become available or as they change in size or volume of data handled.

Other complex systems. e.g. **the Word Wide Web, show the same behavior.** This is the result of the simple mechanism that allows networks to expand by the addition of new vertices which are attached to existing well-connect vertices.



Scale-free internet (2002 snapshot)

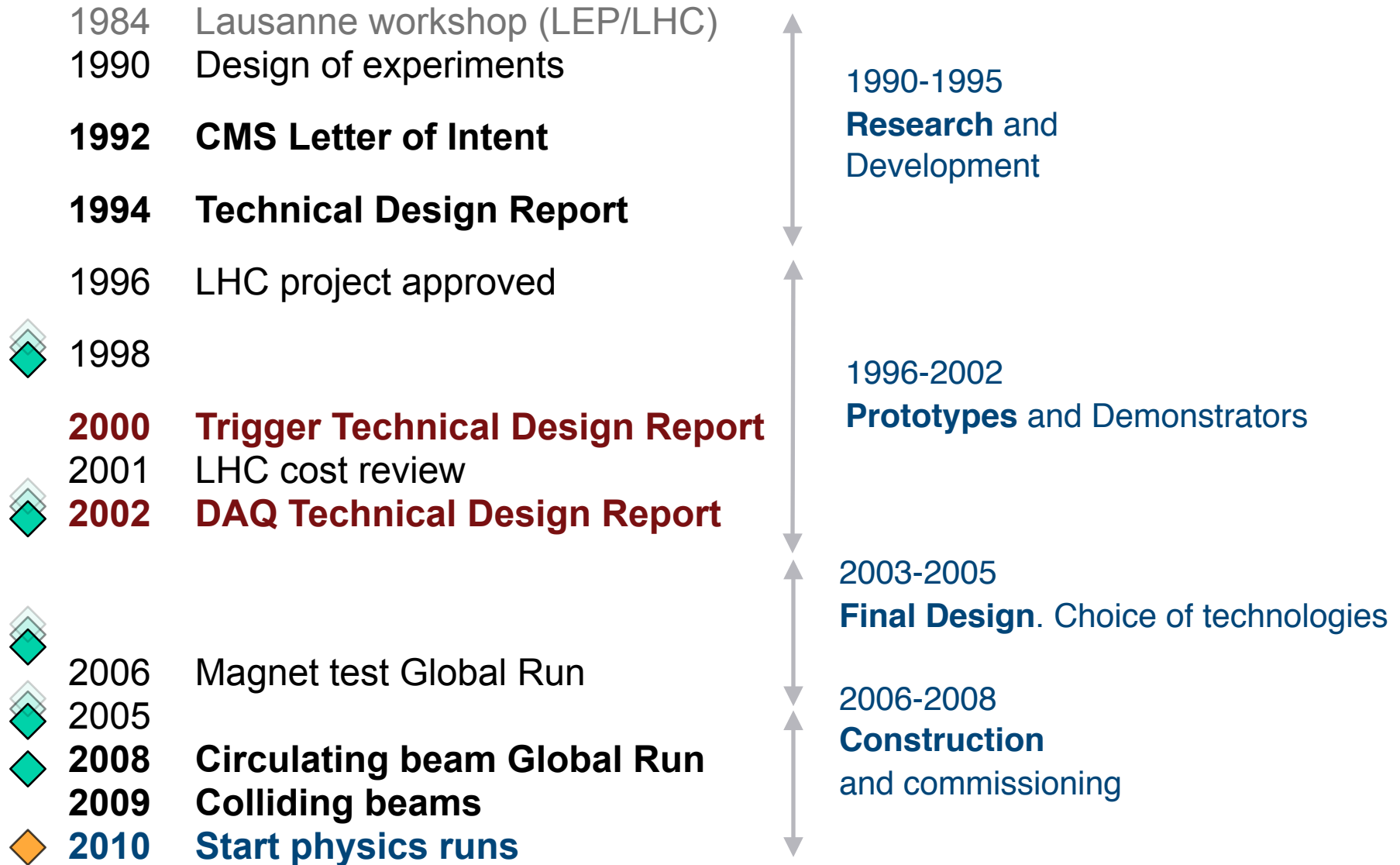


Design issues: Technologies

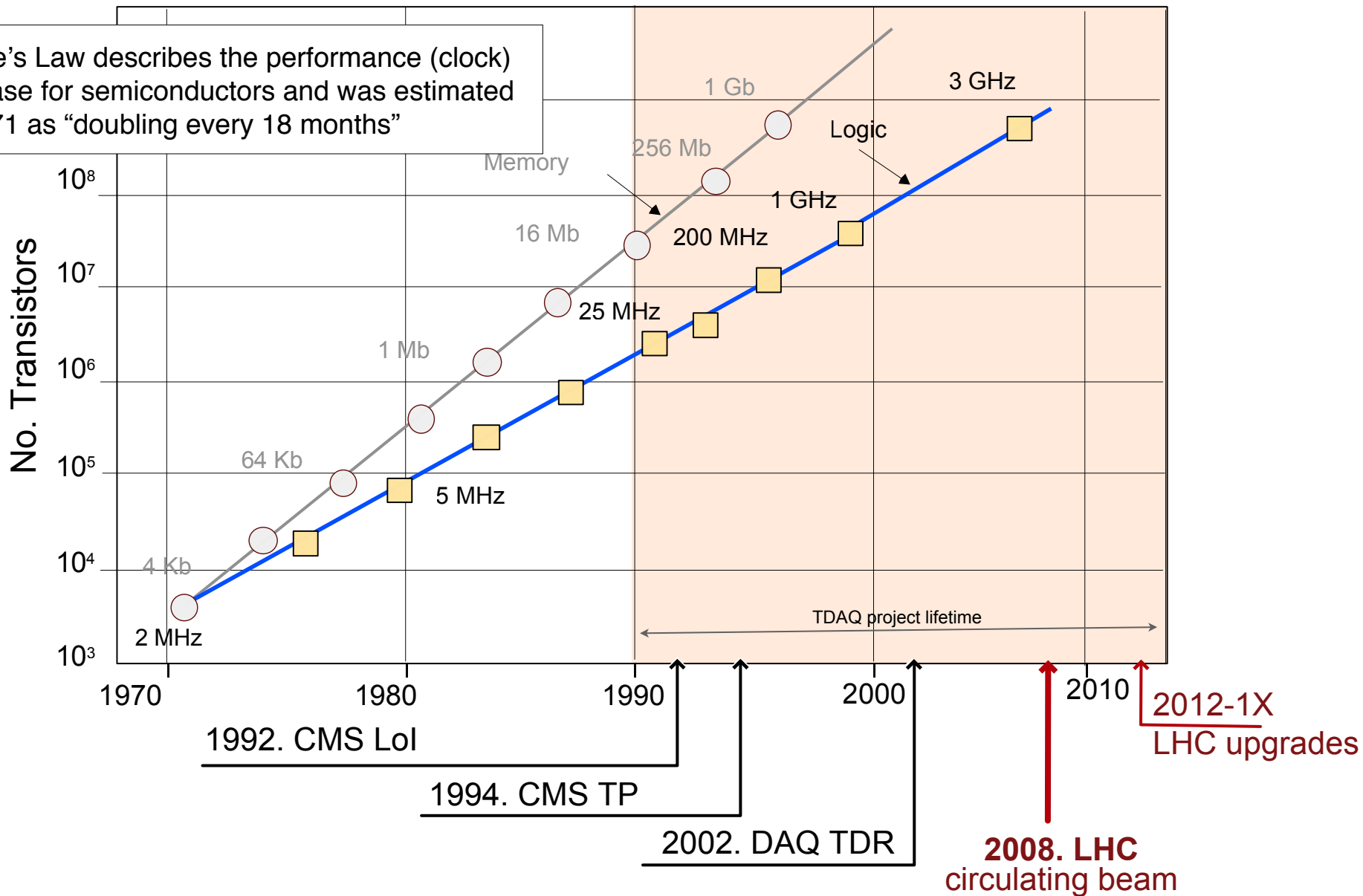
- Project history and information technologies trends
- Predicted and unpredicted evolutions

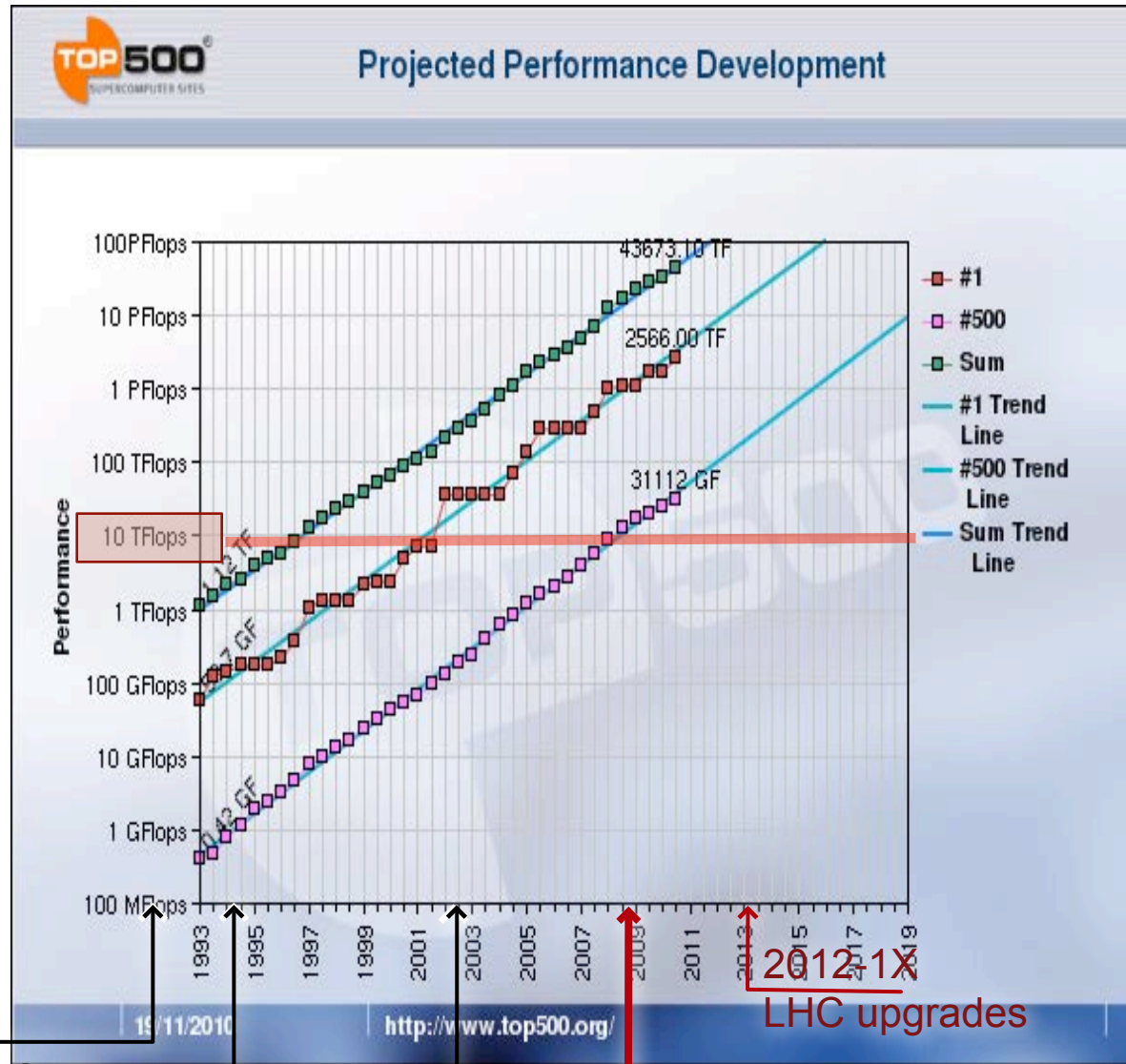


LHC&TDAQ project timeline (the time of a generation)



Moore's Law describes the performance (clock) increase for semiconductors and was estimated in 1971 as "doubling every 18 months"





1992. CMS LoI

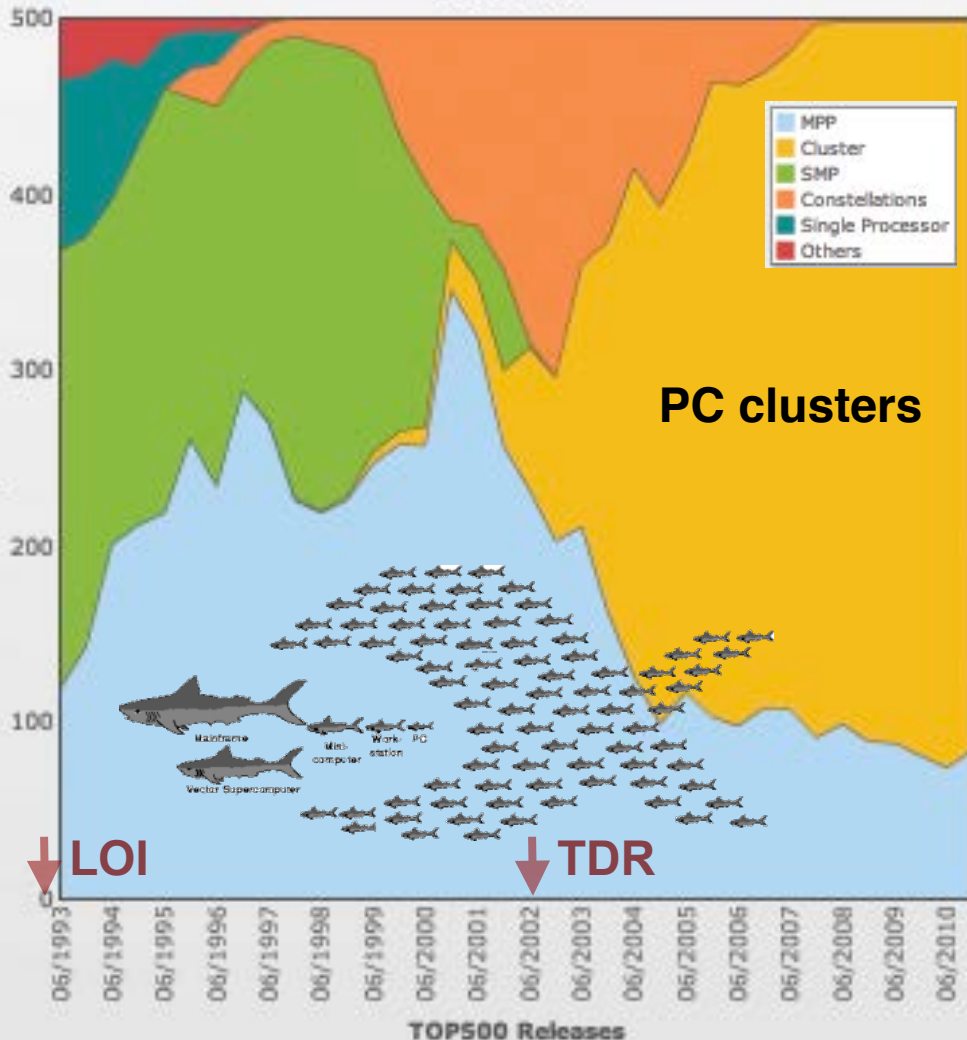
1994. CMS TP

2002. DAQ TDR

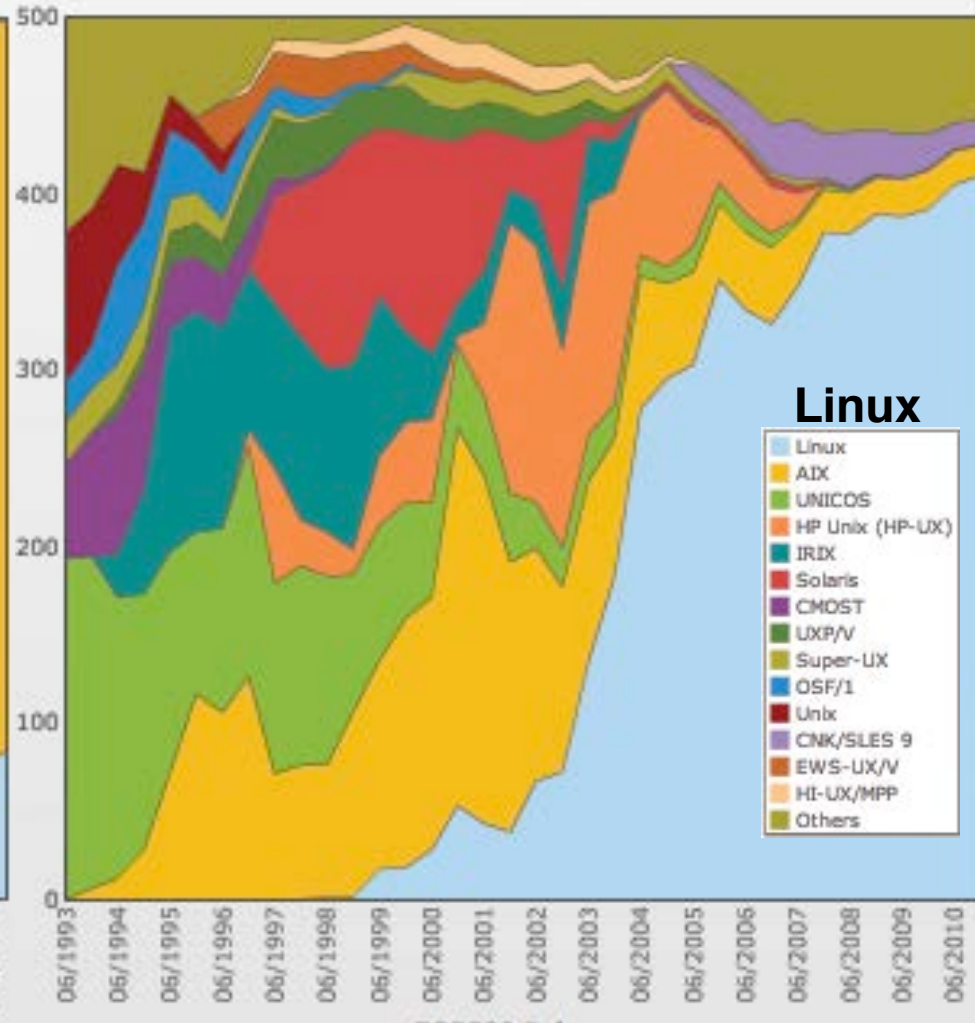
2008. LHC circulating beam

LHC upgrades

Architecture Share Over Time
1993-2010



Operating System Share Over Time
1993-2010





1996. According to Linux Magazine, Digital Domain, a production studio located in Venice, California, produced a large number of visual effects for the film Titanic. During the work on Titanic the facility had approximately **350 SGI CPUs, 200 DEC Alpha CPUs and 5 Tbytes of disk** all connected by a 100 Mbit/s network.

Since 90's:

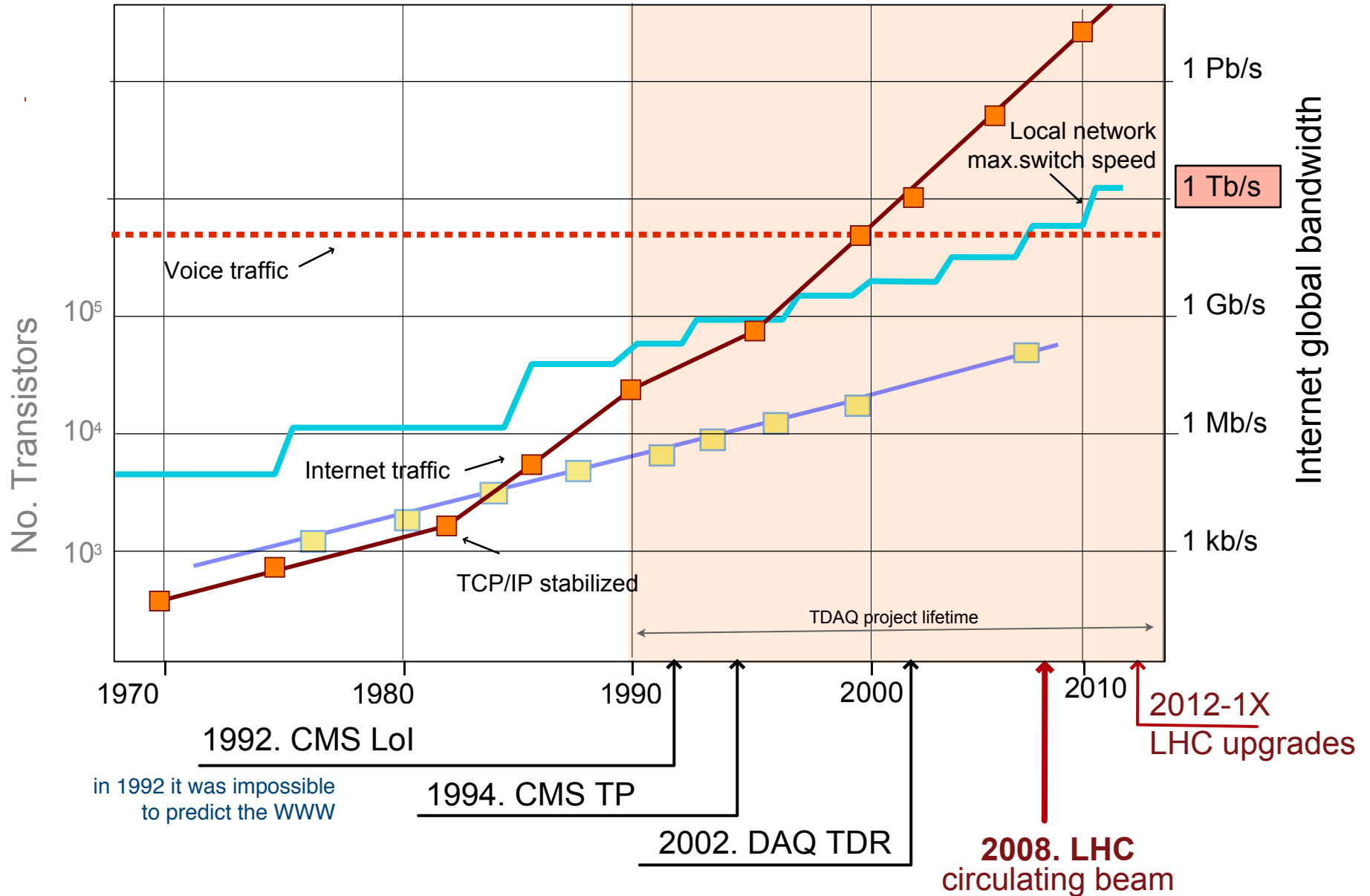
- Large computing power at low cost is made available as **clusters of commodities** (PCs and networks)
- **LINUX** has become the most popular Operating System

CPU estimated in 2002. Total: 4092 s for 15.1 kHz → 271 ms/event. Therefore, a 100 kHz system required about 13 TFLOPs (corresponding to ~30000 CPUs of 2002)

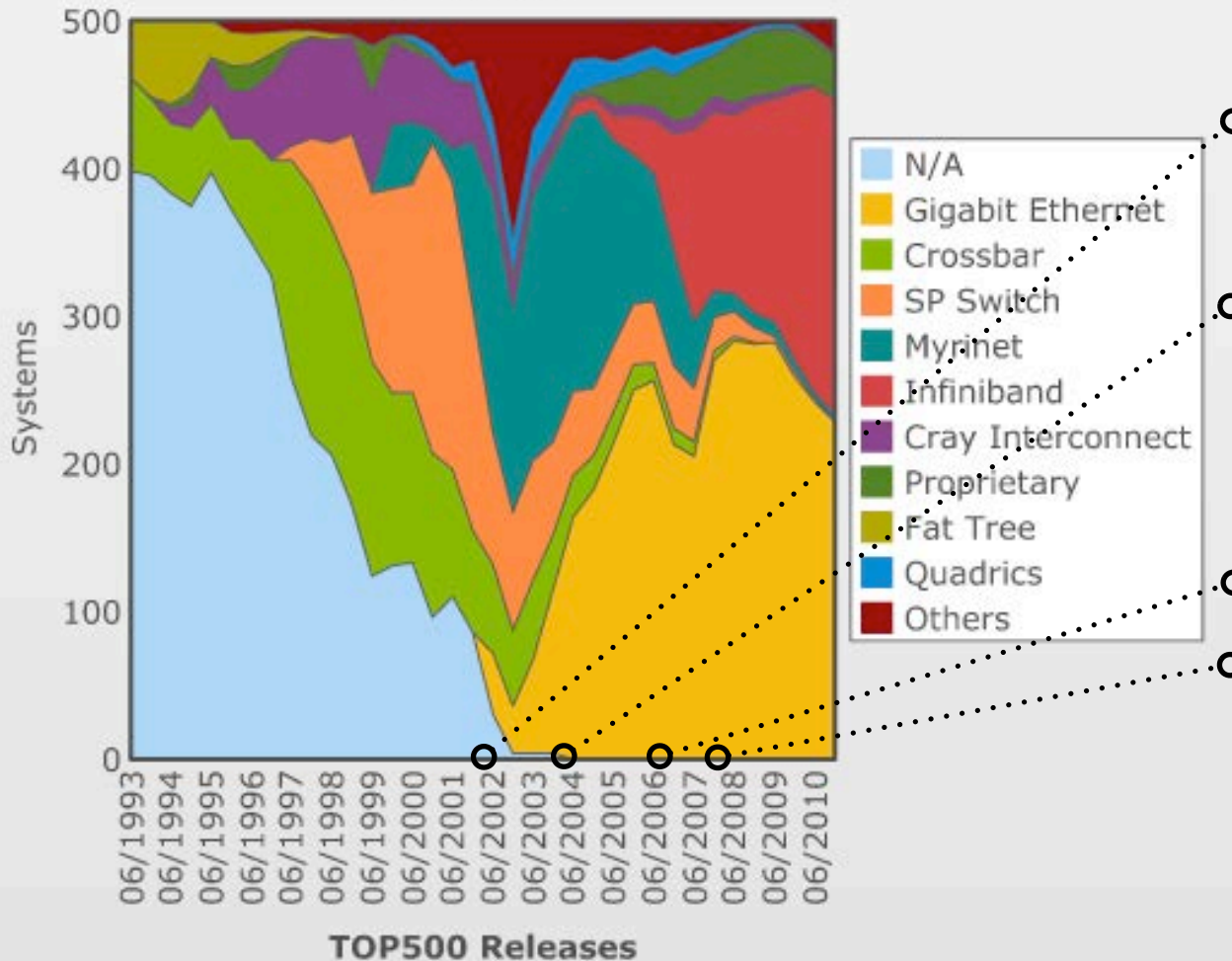
CPU implemented in 2008. The 50% of the HLT system integrated in 2008 consisting of 5000 2.6 GHz CPUs (720 PCs of two quad-core) corresponds to about **10 TFLOPs** in line with the foreseen requirements and in agreement with the Moore law of integrated logic systems (corresponding to a factor 10 in speed every 6 years)



Data communication. Network and Internet traffic trends



**Interconnect Family Share Over Time
1993-2010**



Decision schedule

2002 Data to surface:

- Myrinet used as first layer of readout (FED builder and Data link to surface)

2004 Event builder:

- Gigabit Ethernet routers used for Event builders and DAQ services (controls, mass storage, data link to central Tier0)

2006 Procurements

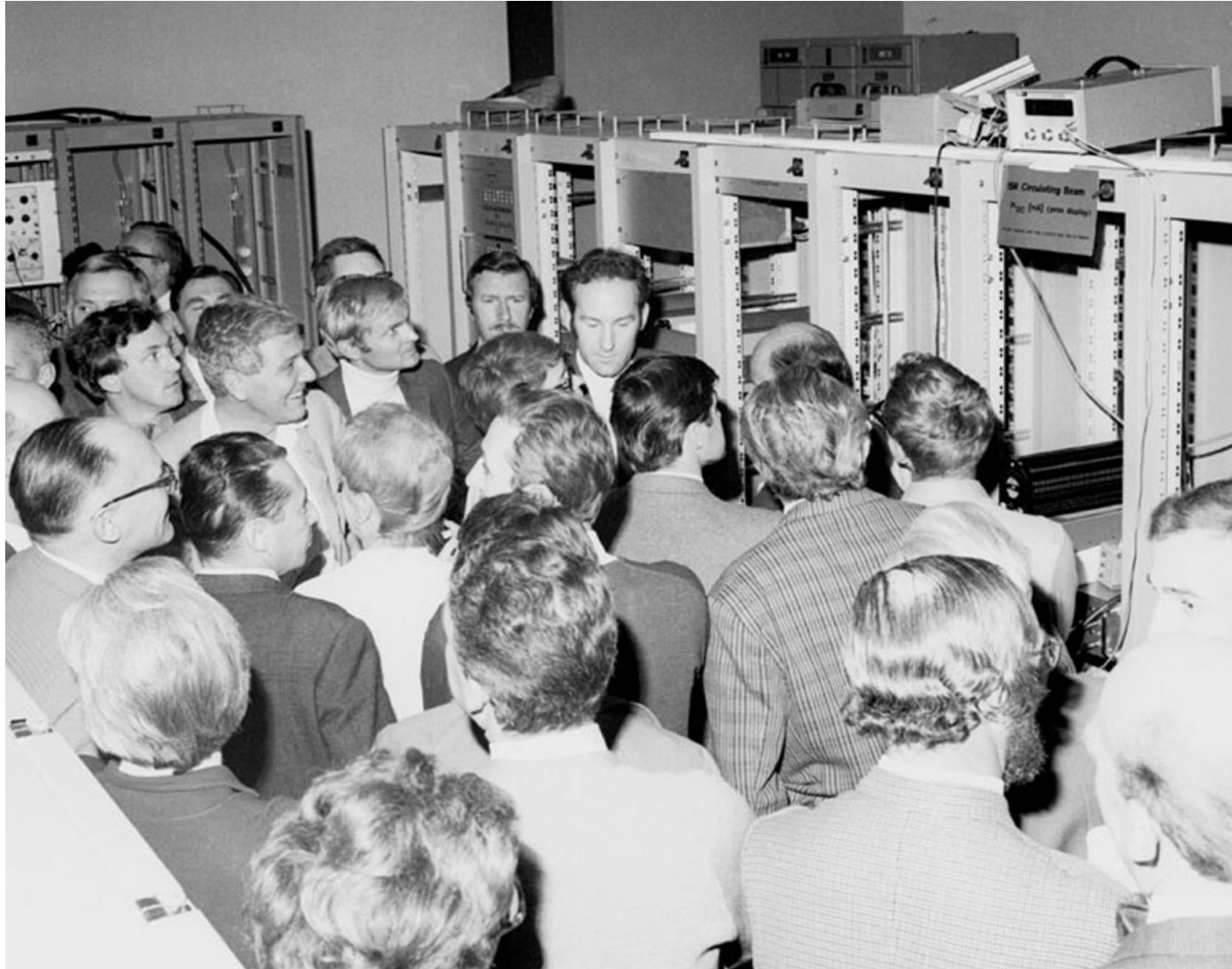
2007 Construction

2008 Commissioning

Unexpected

- Collaborative work
- Network&Computing fusion

ISR. 1970
CR info tools:
Coaxial Cables
Teletype
Telephone



ISR 1970. Voltmeter display, no terminal

ISR. 1970

CR info tools:

Coaxial Cables
Teletype
Telephone

P-aP. 1980

CR info tools:

RS 232
Alpha terminal
Video&Telephone



1980 P-Pbar. A lot of persons in front of one screen

ISR. 1970

CR info tools:
Coaxial Cables
Teletype
Telephone

P-aP. 1980

CR info tools:
RS 232
Alpha terminal
Video&Telephone

LEP. 1990

CR info tools:
RS 232, Ethernet
Graphics terminals
Video&Telephone



1990 LEP. A lot of screens in front of one person

ISR. 1970

CR info tools:

Coaxial Cables
Teletype
Telephone

P-aP. 1980

CR info tools:

RS 232
Alpha terminal
Video&Telephone

LEP. 1990

CR info tools:

RS 232, Ethernet
Graphics terminals
Video&Telephone

LHC 2010

CR info tools:

Wireless
LAN, WAN
Internet, WWW



2010 LHC. The person is on the screen

Cessy: Master&Command control room



Fermilab: Remote Operations Centre



Meyrin: CMS DQM Centre



CR: Any Internet access.....



A general and expandable architecture has been deployed for the **experiments' Run control and monitoring** largely based on the emerging Internet technology developed in the field of **WWW services**



Hard-to-predict in the 90's: The World Wide Web



World Wide Web

The World Wide Web (W3) is a wide-area [hypermedia](#) information system. It is an international initiative aiming to give universal access to a large universe of documents.

Everything there is online about W3 is linked directly or indirectly to this document, including an [executive summary](#) of the project, [Mailing lists](#), [Policy](#), November's [W3 news](#), [Frequently Asked Questions](#).

What's out there?
Pointers to the world's online information, [subjects](#), [W3 servers](#), etc.

Help
on the browser you are using

Software Products
A list of W3 project components and their current state. (e.g. [Line Mode](#), [X11 Viola](#), [NeXTStep](#), [Servers](#), [Tools](#), [Mailrobot](#), [Library](#))

Technical
Details of protocols, formats, program internals etc

Bibliography
Paper documentation on W3 and references.

People
A list of some people involved in the project.

History
A summary of the history of the project.

1992

Since the start of the exploitation of large accelerator laboratories around the world, the design and operation of High Energy Physics experiments have required an ever increasing number of participating institutions and collaborators. From tens of institutions and hundreds participants during the Collider and LEP period up to **hundreds of institutions and thousands scientists** in today LHC experiments.

At the end of 80's with the digitalization of information and the growing support of information infrastructures (computer centers and Internet), a tool was needed to improve the collaboration between physicists and other researchers in the high energy physics community.

The **World Wide Web** originally was intended for this purpose, however fusing together networking, document/information management and interface design it has become in few years the most popular instrument to provide seamless access to any kind of information that is stored in many millions of different geographical locations. In addition, it stimulated the expansion of network infrastructures and the development of new software and hardware services based on common standards (TCP/IP, HTML, SOAP, XML,.... GRID, CLOUD,...)

2011

CERN European Organization for Nuclear Physics

Structure: EN, IT, BE, TE, HR, FP, GS, DGU, DGS

Physics Experiments & Research Library & Archives

BBC Mobile News Sport Weather

NEWS EUROPE

Home UK Africa Asia-Pac Europe Latin America Mid-East South Asia US & Canada

4 May 2011 Last updated at 15:33 GMT

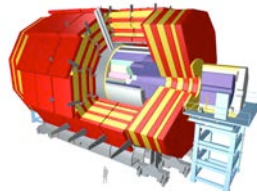
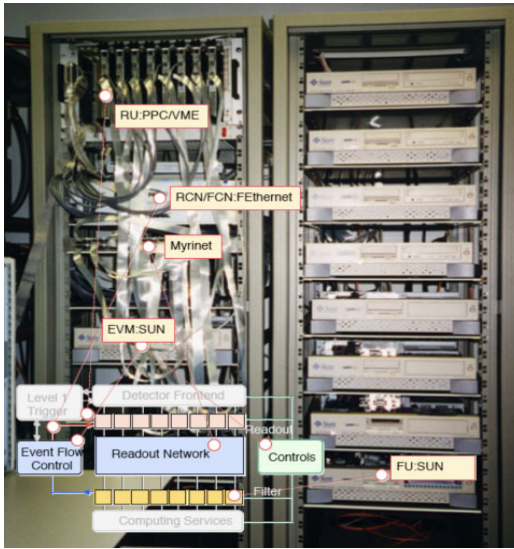
EU moves to tighten border controls

The EU Commission says reimposing border checks in the Schengen zone may be necessary when faced with extraordinary flows of migrants. Migrants set sights on France. Q&A: Schengen Agreement

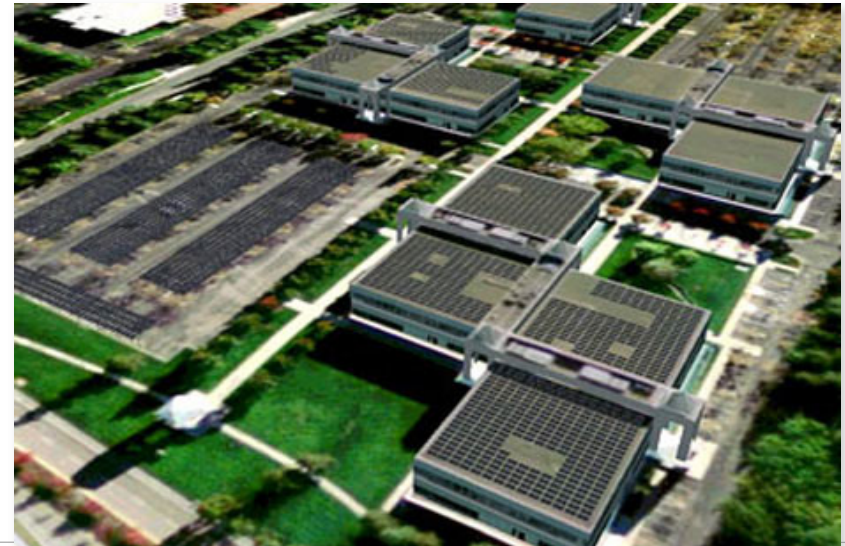
Attack after Turkey PM poll rally

A policeman is killed during an attack in northern Turkey, staged within days of a visit by Prime Minister Recep Tayyip Erdogan, Turkish media say.

2008 The CMS HLT center on CESSY and hundreds Off-line GRID computing centres 10^5 cores



2008 One of Google data center 10^6 cores

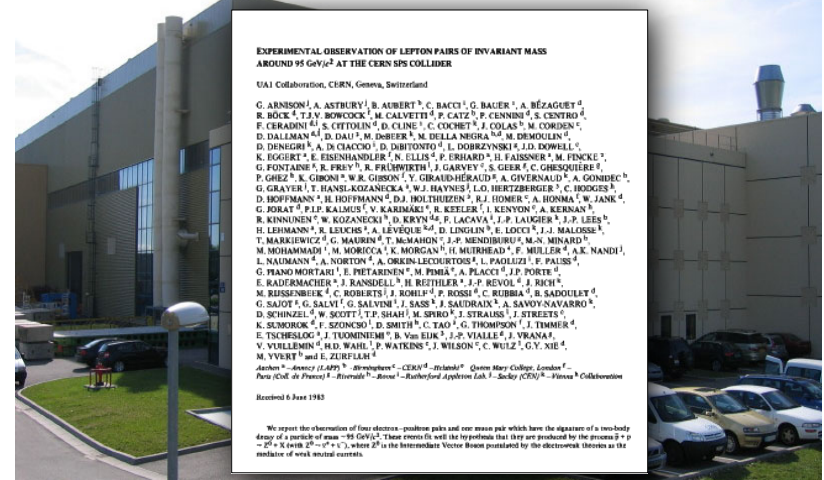
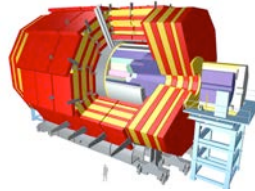
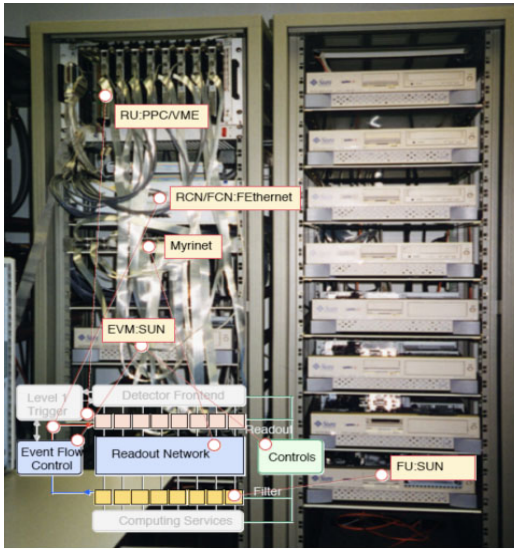




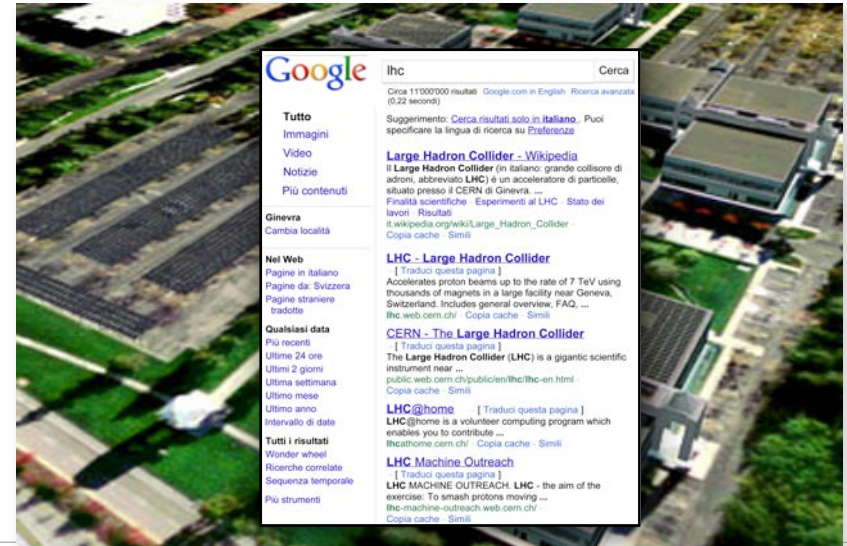
Hard-to-predict in the 90's (II): the same model elsewhere



2008 The CMS HLT center on CESSY and hundreds Off-line GRID computing centres 10^5 cores



2008 One of Google data center 10^6 cores



In operation ...

<http://cmsonline.cern.ch/daqStatusSCX/aDAQmon/DAQstatusGre.jpg>

<http://cmsonline.cern.ch/daqStatusSCX/aDAQmon/LayoutB.jpg>


<http://cmsonline.cern.ch/daqStatusSCX/aDAQmon/wDAQ/aDAQmonScreenDumps1.html>



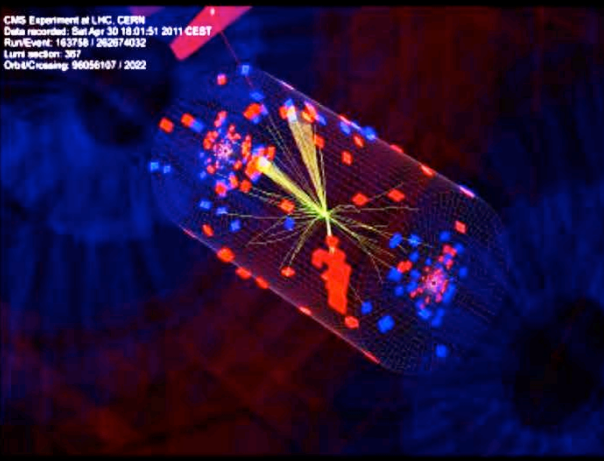
Data acquisition in operation



26/04/11 Tue 12:34: 9.01Z: 1164: 12/04/11 Tue 21:10


 30/04/11 PROTON PHYSICS DAQ state Run Number **Lv1 rate** Ev. <Size> kB DeadTime(AB) Acc. Hz (%) **HLT <CPU>**
 Sat 18:04:30 STABLE BEAMS Running 163758 30.929 kHz 401.0 [224.8] 0.673 % 30991.1 (100.0%) 21.96 %

CMS Experiment at LHC, CERN
 Date recorded: Sat Apr 30 18:01:51 2011 CEST
 Run# 163758 / 262574032
 Lumi section: 357
 Orbit Crossing: 90956187 / 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/Run2011/5e32/v8.3/HLT/V4
 #LS Random ON
 PreScaleIndex Physics ON
 Tracker HV ON CalibCyc ON
 Pixel HV ON
 Physics DECLARED
 #Lv1(GT) **266915432** 55 FEDCRC
 Lv1 Rate **30.929 kHz**
 Pending Lv1 FBI occ. %
 #Frag. in RU Max
 Min Min
 BnW (MB/s) **1.2E+4** EvSize (kB) **401**
 Events in BU
 Pending Req. Rceiv.-Disc.
 #P> 53179 P.M-m
 #Running FUs 53183 A.M-m
 100.00% <FU-CPU>
A **318.8 Hz** 100
 BnW MB/s Disks usage <SM-CPU>
 EventRate Hz Free space TB
 Stored **50002951**
 Time to fill disk 3 of srv-c2c06-14 > week
 TIER1 TRANSFER ON

Beam setup & DCS states history

LHC mode: PROTON PHYSICS, STABLE BEAMS

Run# 163758 history time window (2.6 H)





Technologies and TDAQ: summary



- Continued rapid processor performance growth following Moore's law
- Network bandwidth grew at rate even higher than that of Moore's Law
- Commodity products everywhere
- Open software model (Linux) became a standard
- Internet explosion included software and hardware services
- Global network services allow scale-free decentralisation of data analysis and storage

However

High technology products limited life time (3-5 years)

Long-term maintenance issues and equipment replacement plans

The LHC program – machine, experiments, computing – has taken 25 years to complete, from conception to operation. It is expected to operate for as long with some major upgrades in the performance of the machine.

25 years is several generations of technology changes which the TDAQ design has to take into account for its maintenance plans and for the next machine upgrade program. Even more use of emerging computing and communication Internet services will (have to) be made



SLHC upgrade



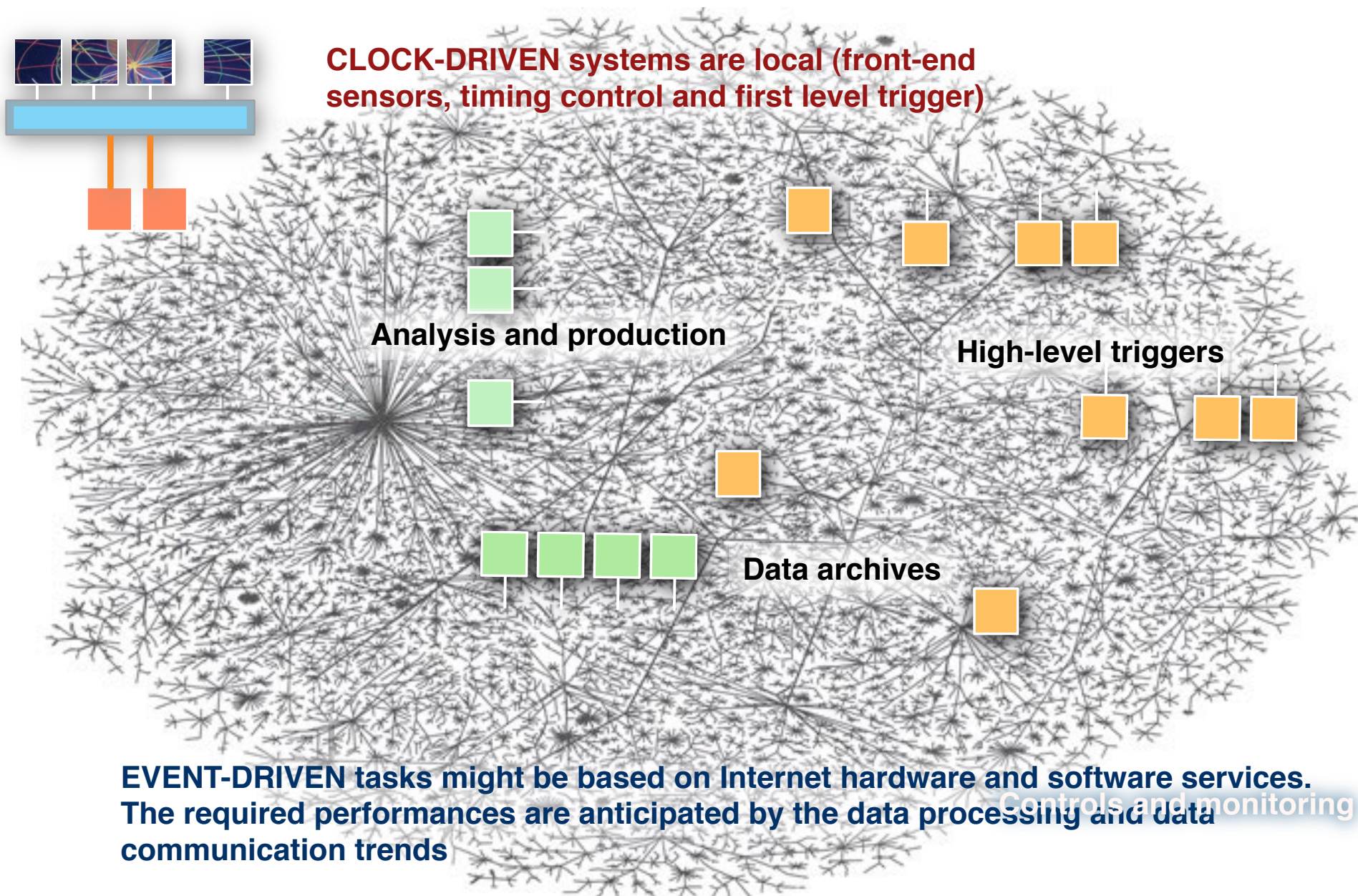
The current TDAQ design should be considered complete, but its implementation is not final. It is expected to change with time, accelerator and experiment conditions to provide the maximum possible flexibility to execute a physics selection on-line.

Luminosity increase (2018...) will require

- New front-end electronics and readout links
- Higher **level-1 selection power** (to maintain 100 kHz max. output)
- **Readout network (>10 Tb/s)** with an order of magnitude higher
- **More on-line computing power/mass-storage**

The upgrade program will include:

- All very front-end systems and selection logic will **still be based on custom design**. However **new telecommunication technologies** (e.g. TCA etc.) can be employed to interconnect data concentrators, level-1 logic modules and to interface the detector readout with commercial standards.
- **Level1 trigger**. Wider scope, new architectures including arrays of logic nodes. Use standards for data communication and central logic
- **Front-End digitizers**, new rad hard data links and a new timing and trigger distribution system (distribute event type, HLT destination etc.).
- **Data to Surface links (10 Tb/s)** has to be replaced (2005 proprietary technology life time and 10 time the speed). Likely with standards e.g. 1000x10Gb/s data links (not yet a Moore law for data links)
- Event data fragments will be tagged with trigger type and HLT destination. Event builder and High Level Trigger will be **embedded in an single data network** (real-time internet clusters/grid like?) which includes local/central data archives and off-line





Conclusion



Computing power evolved as expected -- if not faster -- as needed by experiments at the LHC.

Digital information technology as well as the Internet have generated the drive for the development of higher bandwidth networks, along with the expansion of world-wide infrastructures to interconnect computing and data routing centers.

The computing and communications challenges that were posed by experiments in high energy physics (HEP) have not only presented themselves as high-end applications of the most advanced technologies. They have also been a source of inspiration for the development of new ones.

Even more importantly, HEP promises to maintain its dual role of client/motivator.