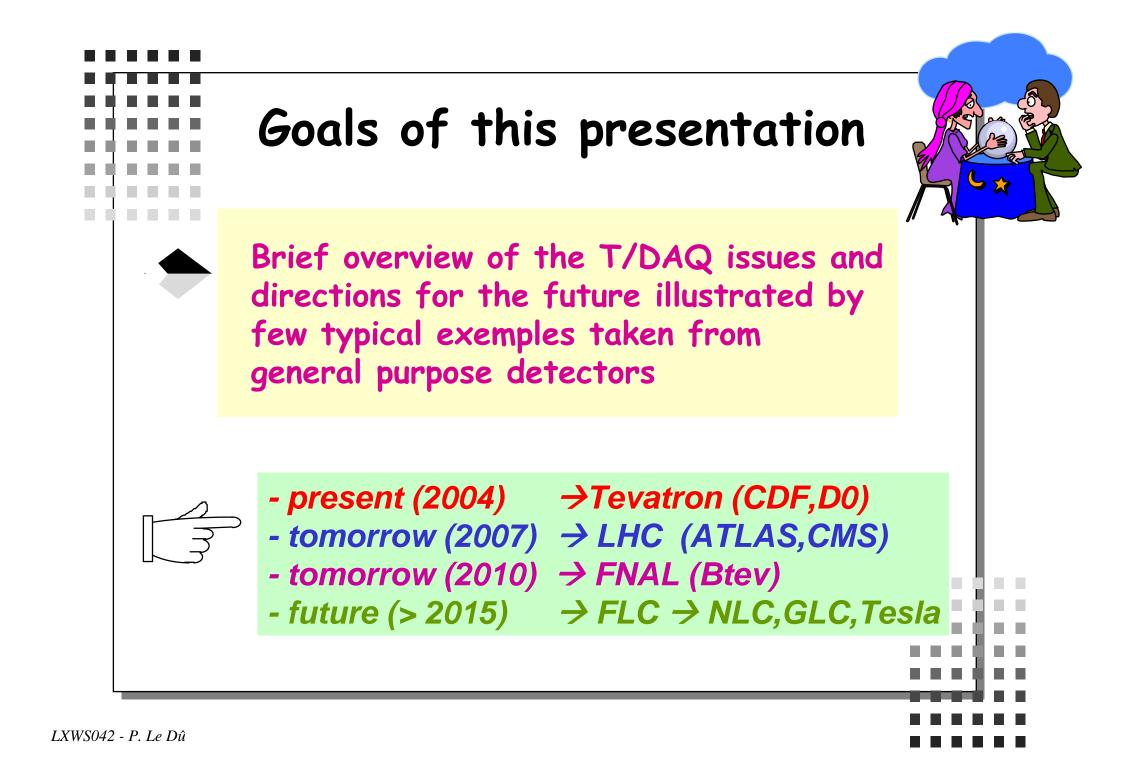
Trigger and Data Acquisition for collider experiments Present and future



- Basic parameters and constraints
- Event selection strategy and evolution of architectures
- Implementation of algorithms : hardware and software
- Boundaries (Trigger levels/DAQ/on-line/off-line)
- Tools,techniques and standards



Credit

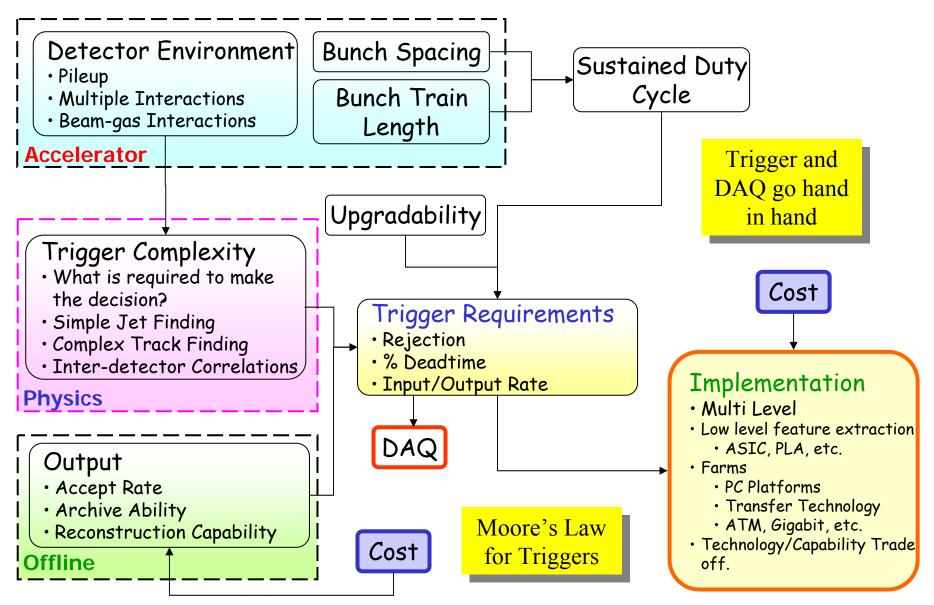
• Trigger/DaQ sessions

Gordon Watts University of Washington, Seattle NSS 2003



M.Mur & S. Anvar (engineers CEA DAPNIA)

Constraints \rightarrow a multiparameters problem

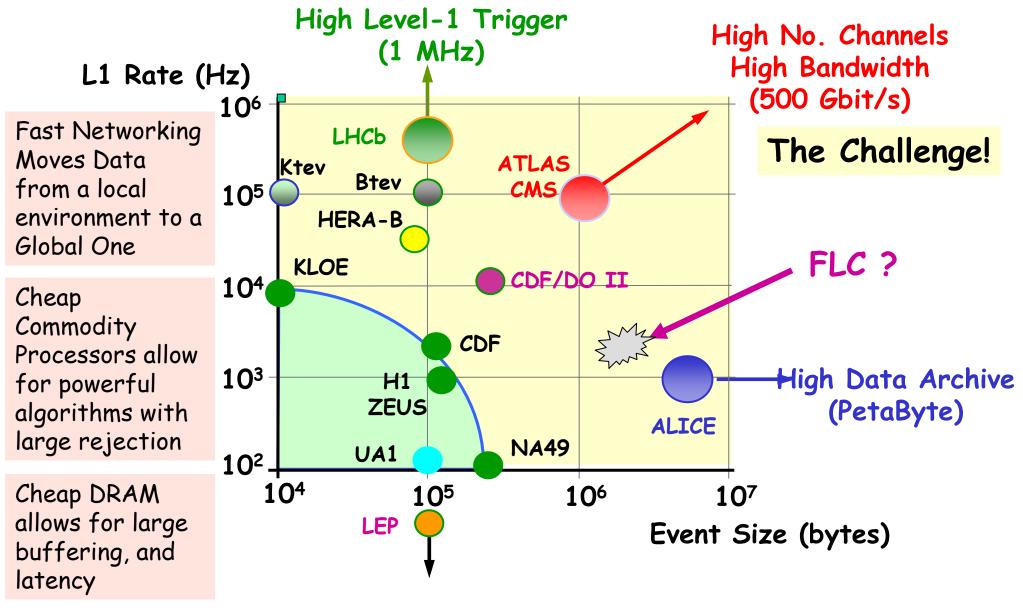


Basic parameters

Exp. <i>Year</i>	Collision rate	Channel count	L1A rate		Event building	Processing. Power		Sociology		
UA's <u>1980</u>	3 µsec	-				5-10 MIPS		150-200		
LEP 1989	10-20 µsec	250 - 500K		-	10 Mbit/sec	100 MIPS		300-500		
BaBar <u>1999</u>	4 ns	150K	2	KHz	z 400 Mbit/s		1000 MIPS		400	
Tevatron	396-132 ns	~ 800 K	10 - 50 KHz		4-10 Gbit/sec	5.10 ⁴ MIPS		500		
LHC 2007	25 ns	200 M*	100 KHz		20-500 Gbit/s	>10 ⁶ MIPS		2000		
NLC 2015	2.5 - 330 ns	800 M*	-		1 Gbit/s	~10⁵ MIPS		5000		
	*	including pixel	s							
			•	Sub-Detector			Tevatron		LHC	

		revation	LIIO	
	Pixel	-	150 M	
	Microstrip	~500 K	~10 M	
-	Fine grain trackers	~100 K	400 K	
	Calorimeters	50 K	200 K	
	Muon	50 K	~1 M	

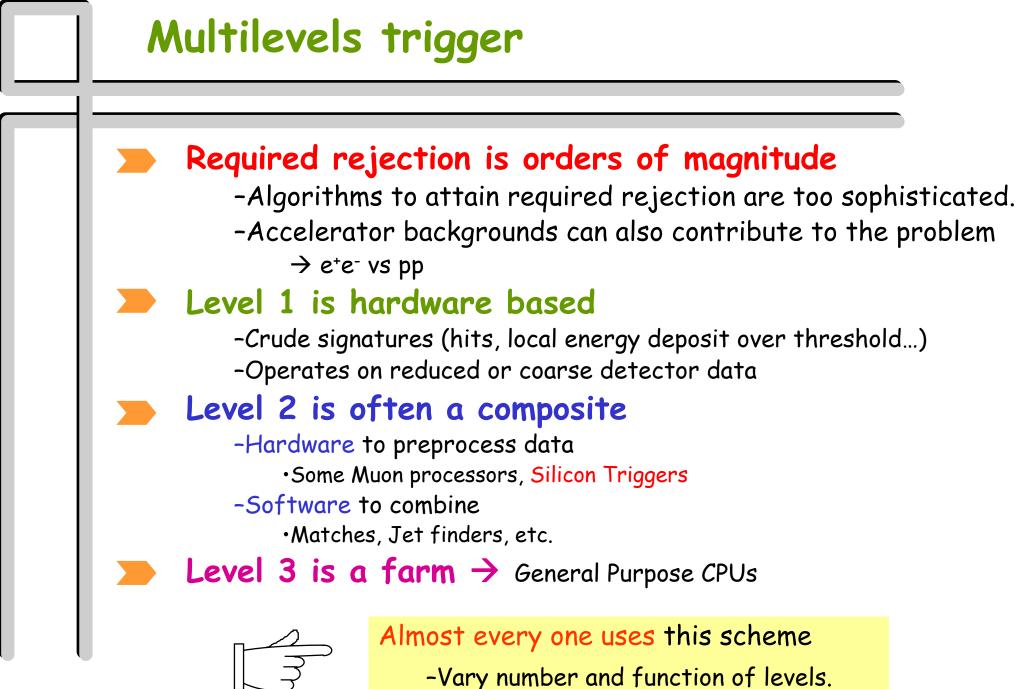
Rate and Data Volume \rightarrow the age of network











$L1 \rightarrow Algorithms in Hardware$

- Cut out simple, high rate backgrounds
 - beam-gas (HERA)
 - z vertex
 - QCD (TeV)
 - jet energy, track matching
- Capabilities are Limited
 - Feature Extraction from single detectors
 - Hadronic or EM energy, tracks, muon stubs
 - Combined at Global Trigger
 - EM Object + Track

Characteristics:

- High speed & Deadtimeless
- Limited ability to modify algorithm
 - But thresholds typically can be modified
- Algorithms Frequently Matched to the detector (and readout) geometry
 - Vertical Design
- Built from Modern
 Components
 - Custom (ASICs)
 - Programmable Logic Arrays (FPGAs, etc.)

The Field Programmable Gate Array

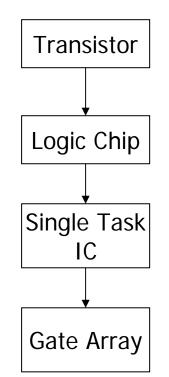
• Board on a chip

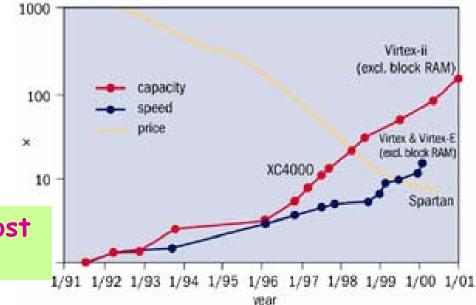
- Revolutionized the way board design is done.
- Logic design is not frozen when the board is laid out.
 - But much faster than running a program in microcode or a microprocessor.
- Can add features at a later time by adding new logic equations.
- Basically a chip of unwired logical gates (flipflops, counters, registers, etc.)

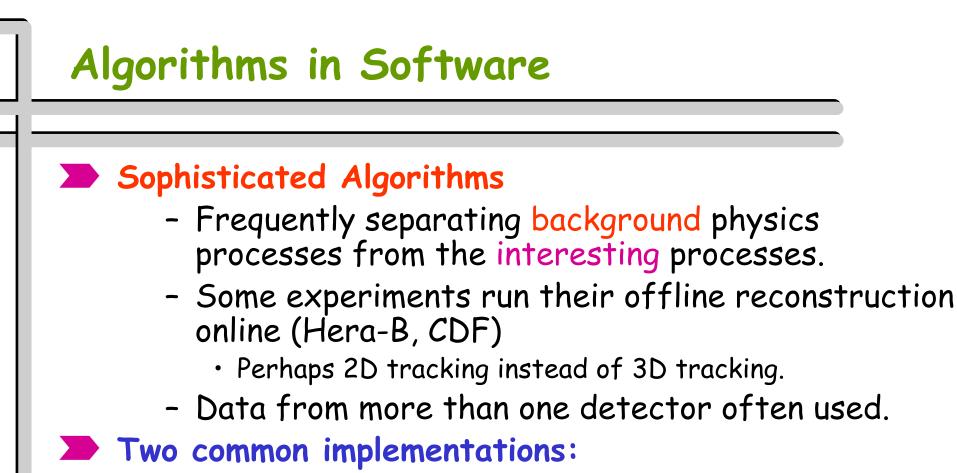
Several Types

Simple Programmable Logic Devices (SPLD) Complex Programmable Logic Devices (CPLD) Field Programmabl Gate Arrays (FPGA) Field Programmable InterConnect (FPIC)

The FPGA allows for the most logic per chip.





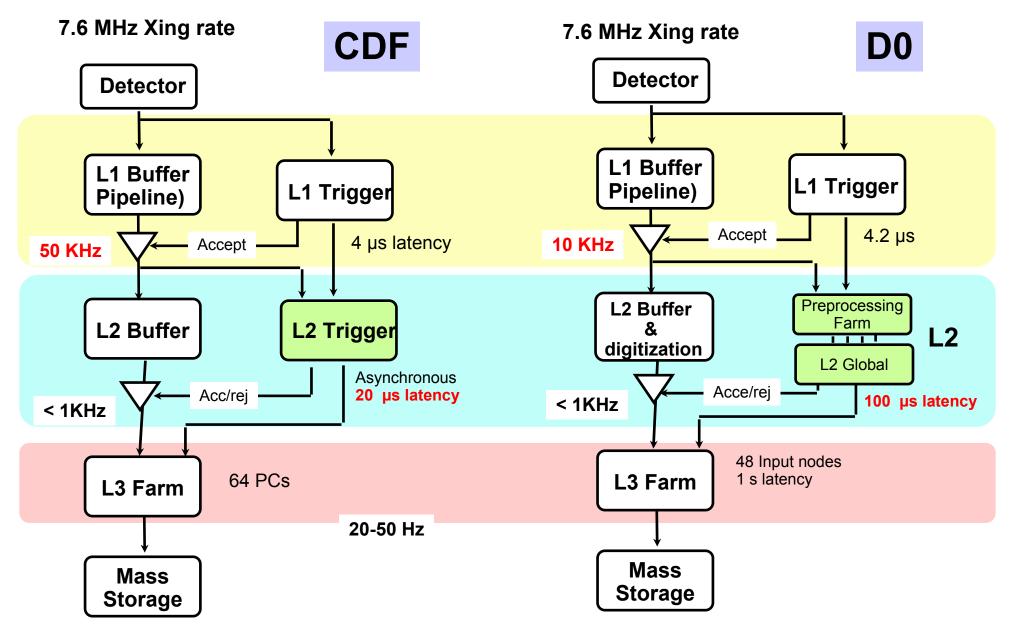


- DSPs tend to live at lower levels of the trigger.
 - Stub finders in muon DØ's muon system.
 - Transputers in ZEUS' L2, DSPs on HERAB L2 & L3
- CPU Farms tend to be at the upper end

Sometimes difficult to classify a system

- Blurring of line between software and hardware.

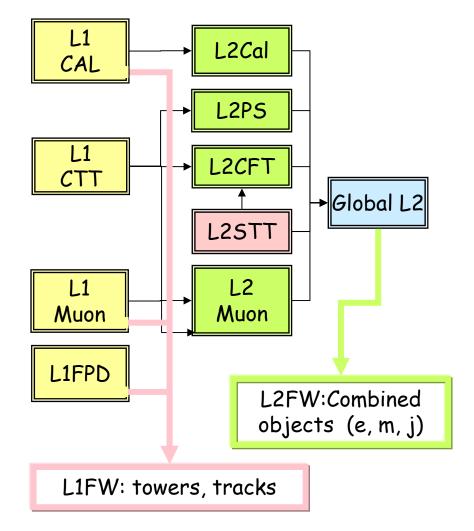
Tevatron architectures (RUN 2)



LXWS042 - P. Le Dû

Hybrids - Level 2 at Tevatron

- Use a combination of both Hardware and General Purpose Processors
- CDF & DØ's Level 2 Trigger
 - Each box is a VME crate
 - Contain specialized cards to land and reformat data from detector front-ends
 - Contain ALPHA processor to run algorithms.



DØ Level 2

Hybrids - SVT

Soth CDF & DØ will trigger on displaced tracks at Level 2

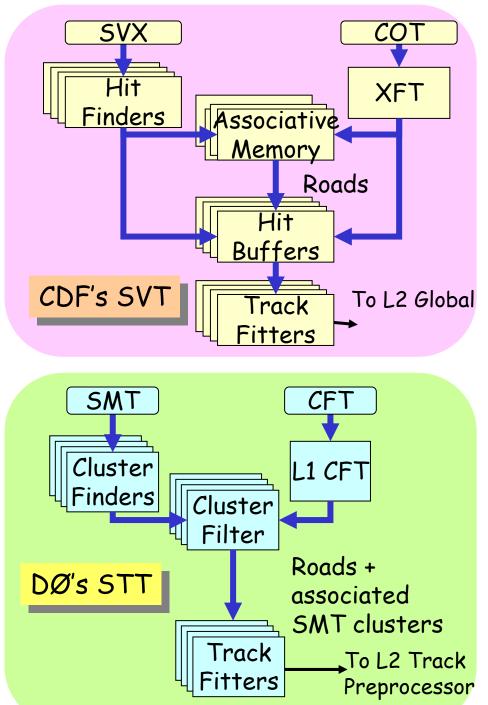
- DØ will use another specialized set of L2 crates.
- Cluster/hit finding
 - CDF: Associative memories @50 KHz
 - DO: FPGA's @ 10 KHz

* Input to Level 2 Global

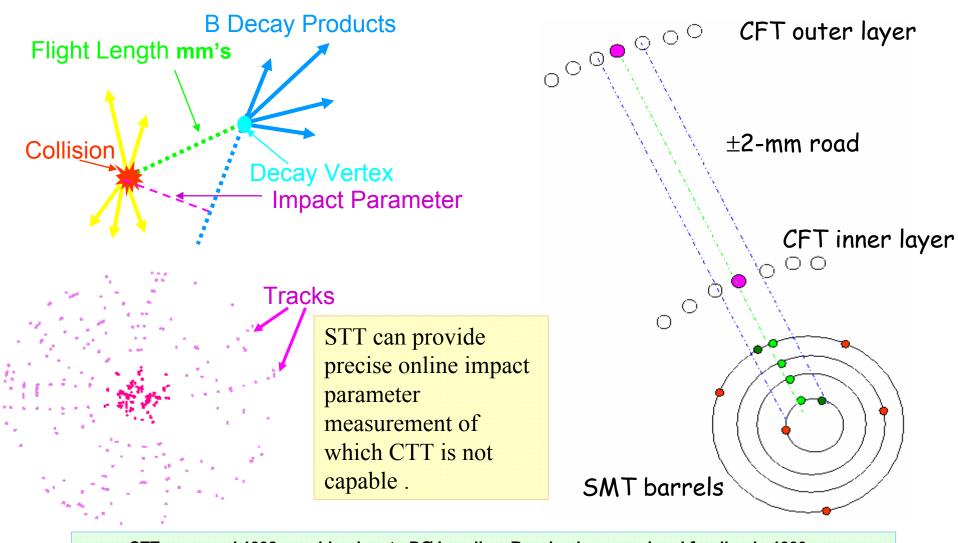
- Select Events with more than 2 displaced tracks
- Resolution is almost as good as offline.
- Refined track parameters

Track Fitters

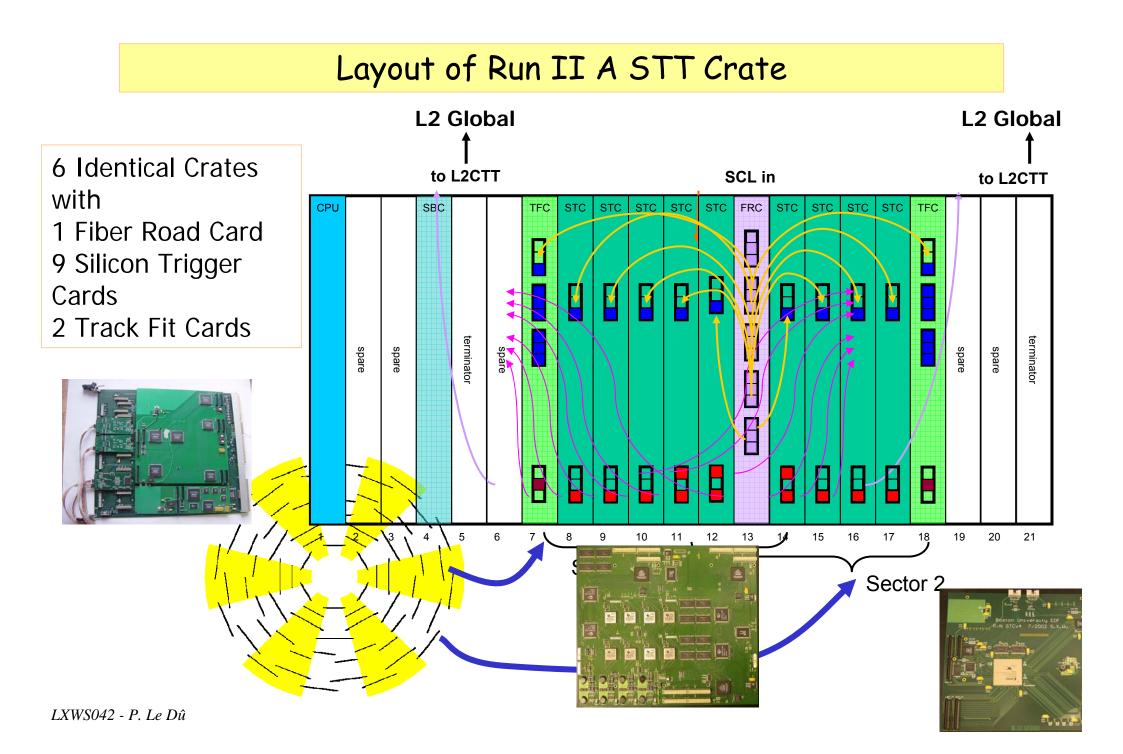
- DSP's Farm.



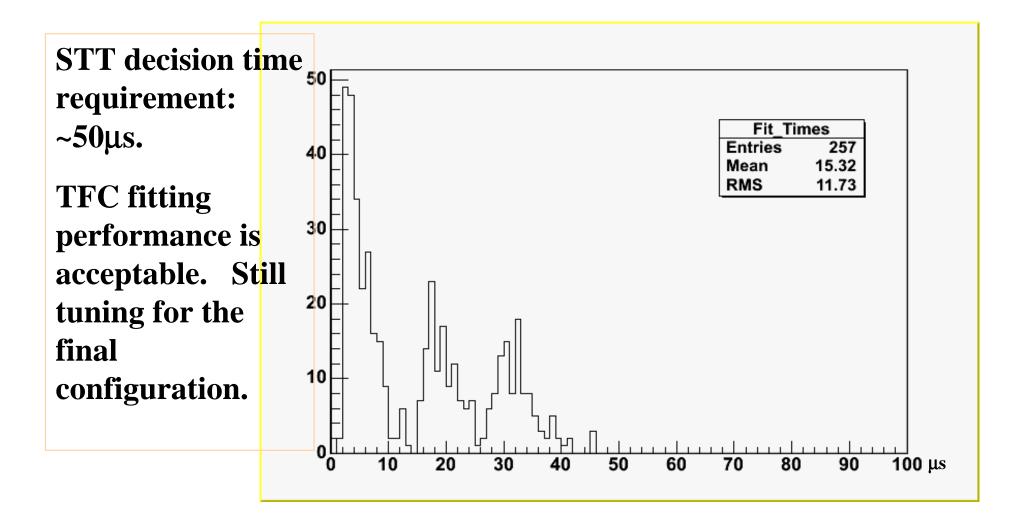
Silicon Track Trigger(STT)



STT proposed 1998 as addendum to DØ baseline. Received approval and funding in 1999



STT Performance I

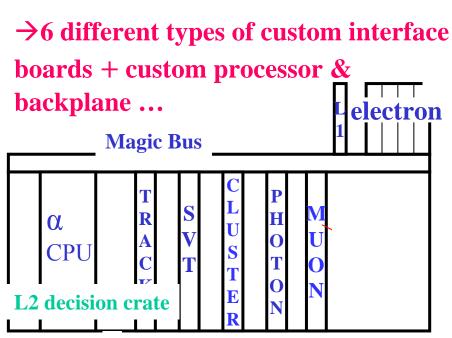


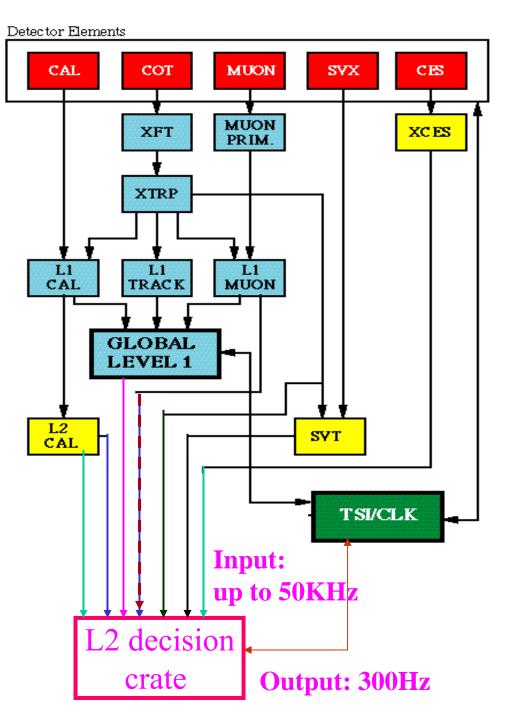
CDF L2 legacy decision Crate

Technical requirement: need a **FAST** way to collect/process many inputs...

→With the technology available back then (1990s), had to design custom (alpha) processor & backplane (magicbus) ...

→had to deal with the fact that each data input was implemented in a different way ...



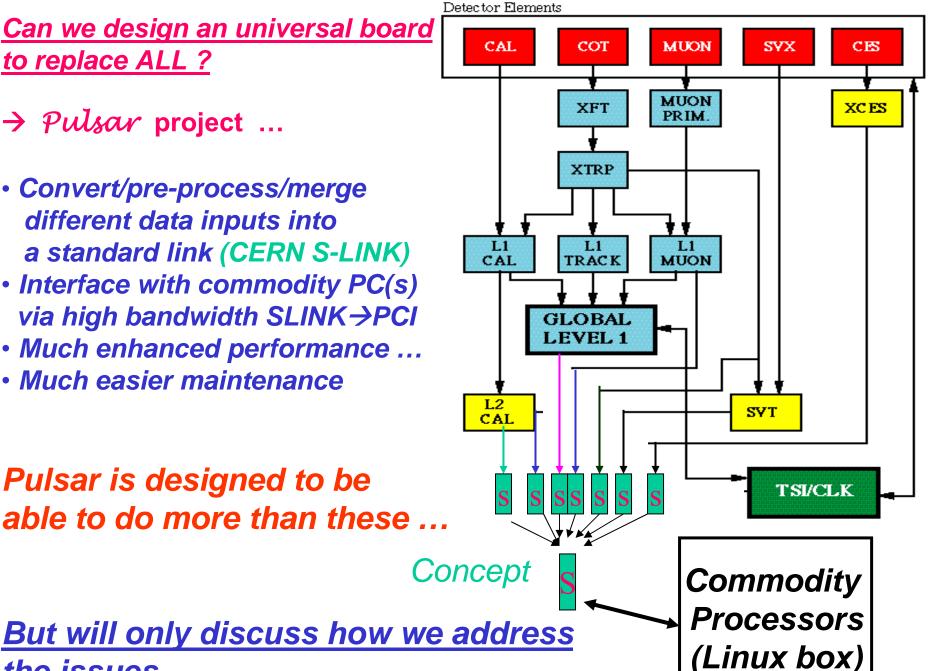


Can we design an universal board to replace ALL ?

- \rightarrow Pulsar project ...
- Convert/pre-process/merge different data inputs into a standard link (CERN S-LINK)
- Interface with commodity PC(s) via high bandwidth SLINK→PCI
- Much enhanced performance ...
- Much easier maintenance

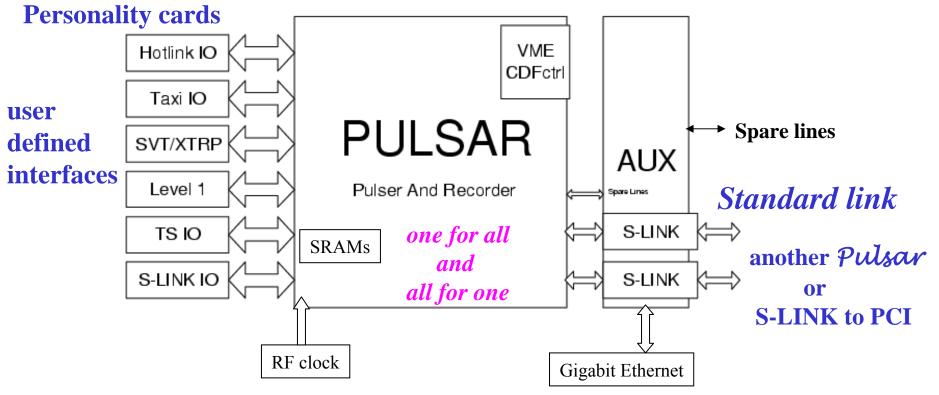
the₄issues

Pulsar is designed to be able to do more than these ...

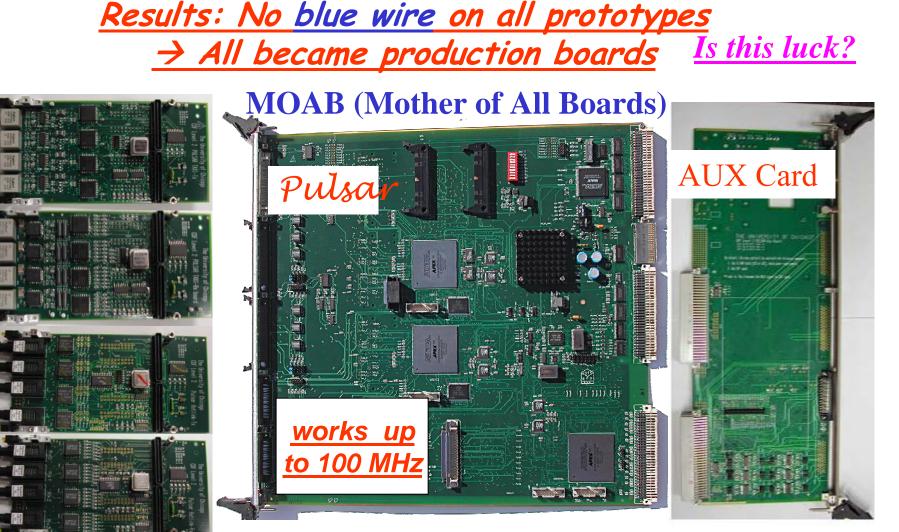


Pulsar is designed to be: <u>Modular, universal & flexible, fully self-testable (board</u> <u>&system level)</u>

Has ALL interfaces L2 decision crate has



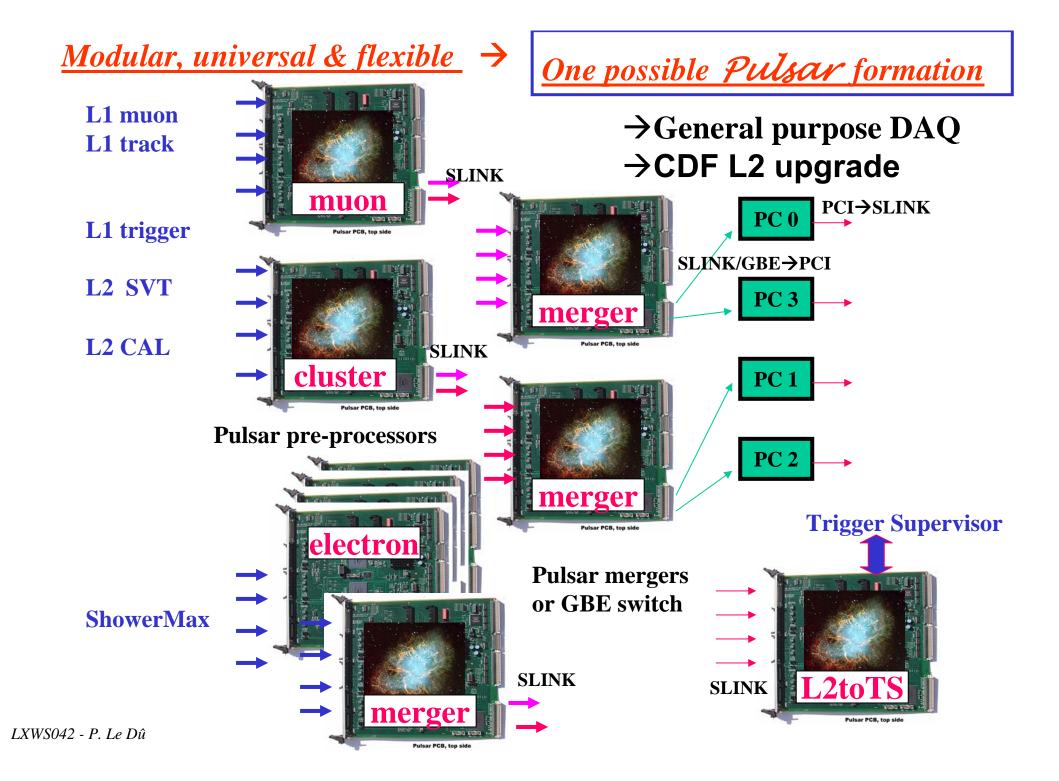
All interfaces are bi-directional (Tx & Rx) Lego-style, open design ...

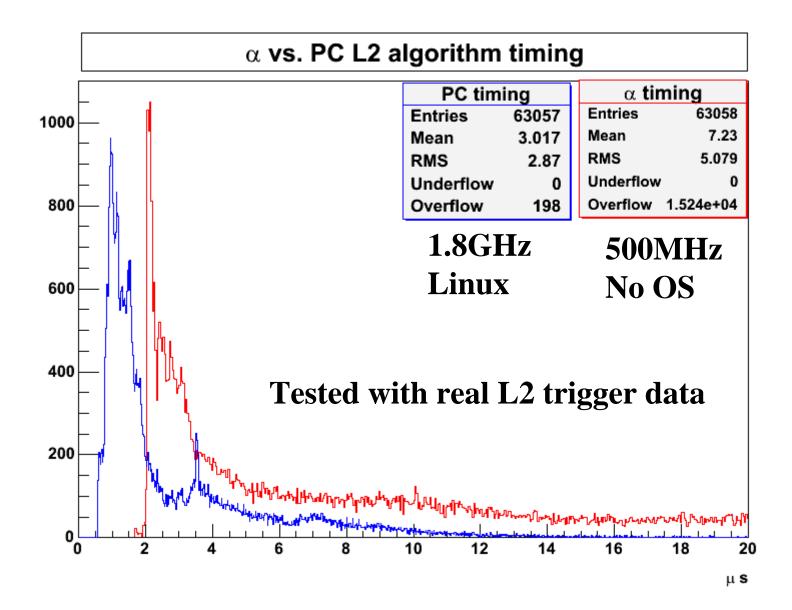


4 types of custom mezzanine LXWS042 - P. Cards (Tx/Rx)

- Design&
 Verifications:
 - ~ 3 people in
 - ~ 9 months

- Initial checkout of ALL interfaces:
 - ~ 2 people in
 - ~ 1 month





Uniformity & modularity & flexibility

Lego-style, general purpose design, backward & forward compatible. Many applications within & outside CDF: Plan to replace/upgrade > 10 different types of CDF trigger board Compatible with S-LINK standard \rightarrow commodity processors ... Knowledge gained transferable to *and from* LHC community...

Design & verification methodology

simulation & simulation: single/multi-board/trace & cross talk analysis ... → no single design or layout error (blue wire) on all prototypes

Testability & commissioning strategy

Board & system level self-testability fully integrated in the design, Suitable to develop and tune an upgrade system in stand-alone mode Minimize impact on running experiment during commissioning phase:

Maintenance & firmware version control

Built-in maintenance capability; Common firmware library (VHDL), CVS ...

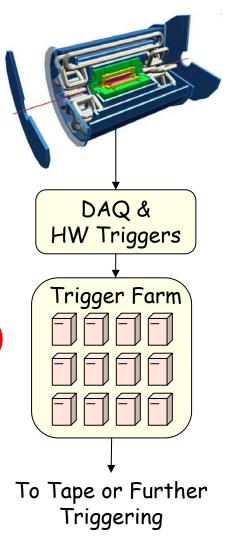
Documenting the work

http://hep.uchicago.edu/~thliu/projects/Pulsar/

Farm Triggers

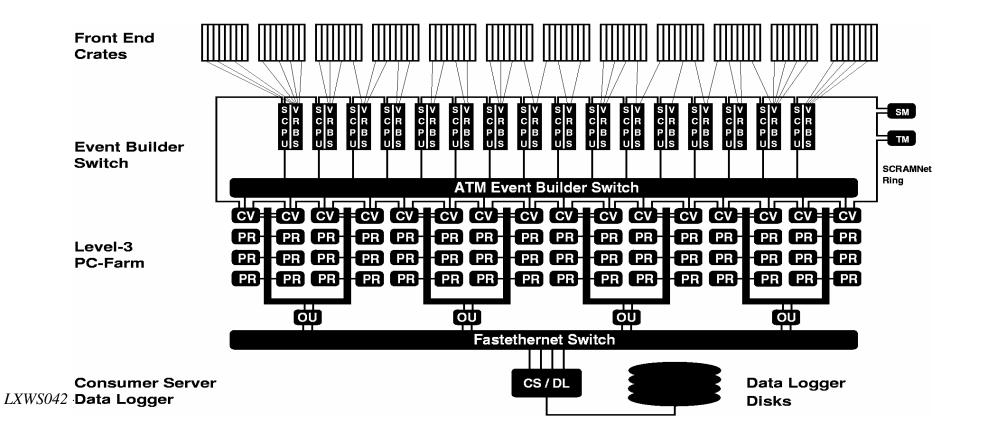
Industry has revolutionized how we do farm processing

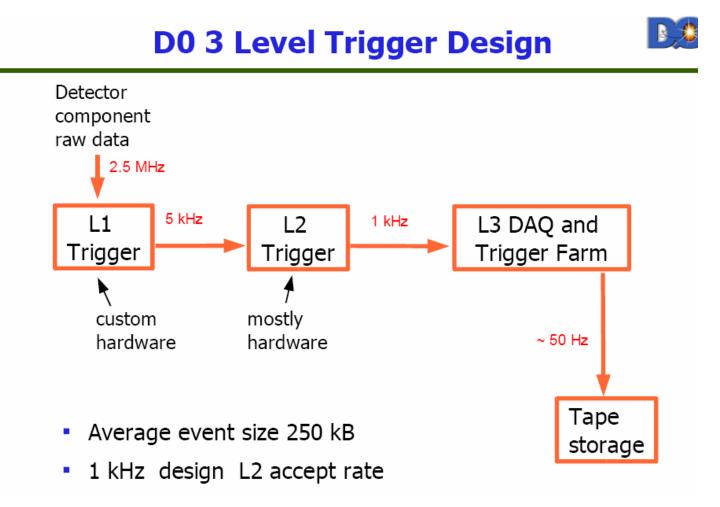
- Wide scale adoption by HEP.
- Online and Offline.
- One CPU, one event
 - Not massively parallel
- Data must be fed to farm
 - As many ways as experiments (push/pull)
 - Flow Management & Event Building
- > Farm Control: an issue !
 - Distributed system can have 100's of nodes that must be kept in sync.



CDF

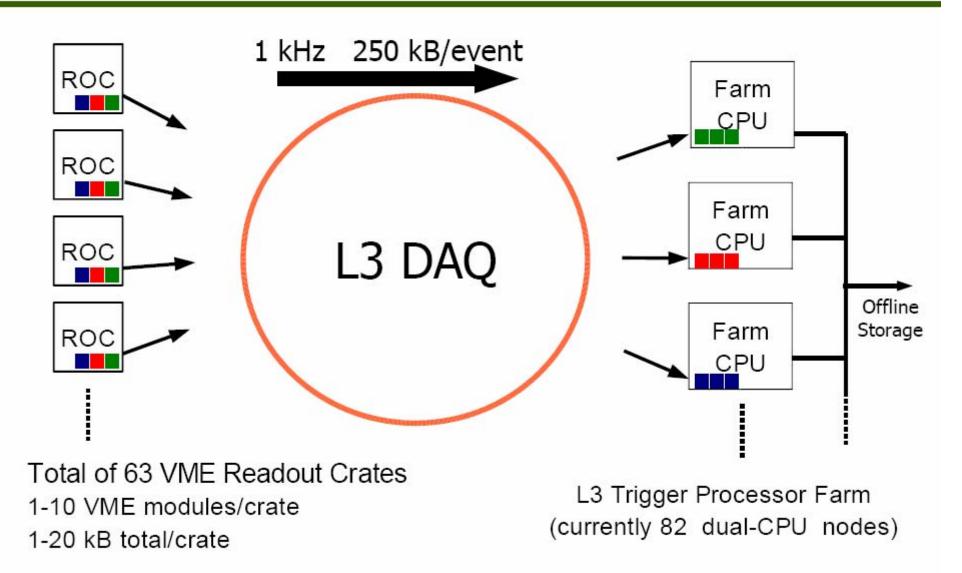
- Data Flows when Manager has Destination
 - Event Builder Machines
- One of first to use large switch in DAQ.
- Dataflow over ATM switch
 - Traffic Shaping
- Backpressure provided by Manager.





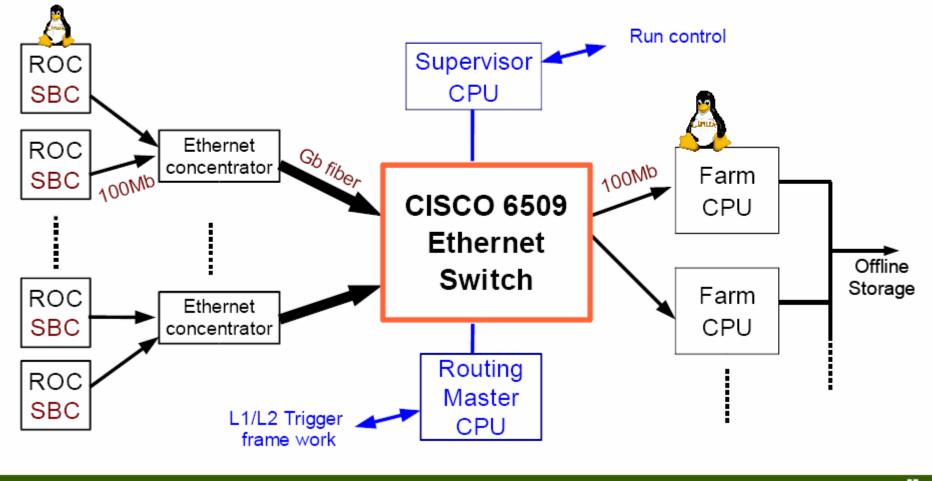
L3 DAQ Requirements





Commodity Ethernet L3 DAQ Design

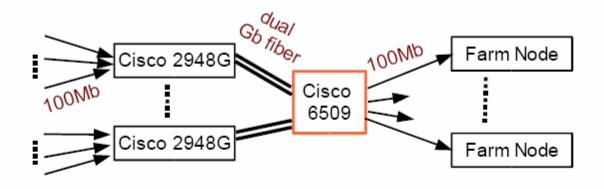
VME Single Board
 Ethernet TCP/IP
 Linux farm nodes
 Computers (SBC), Linux



7

Cisco 6509 Ethernet Switch





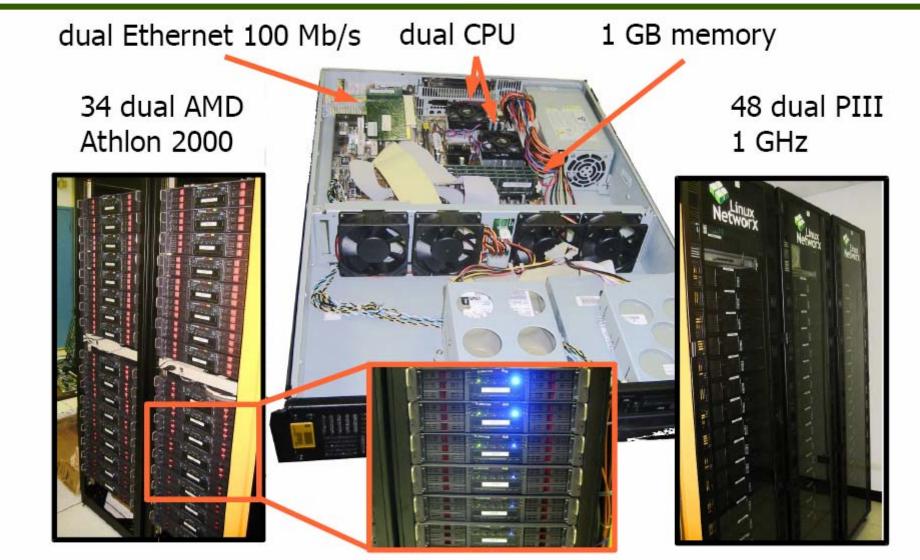
- Backplane easily handles
 250 MB/s design rate
- Using:
 - One 16 x 1 Gb/s module
 - Two 48 x 100 Mb/s modules with 1 MB output buffer per port
- Room for expansion







The L3 Trigger Farm



(Farm typically ~50% busy, <100 ms/event on average)



Error Recovery

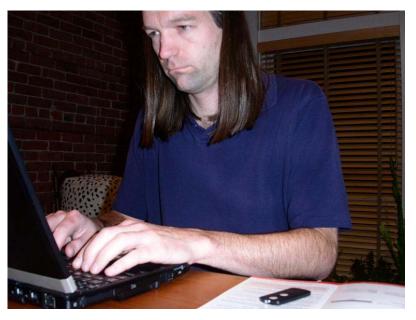


Potential source of inefficiency: occasional human errors by DAQ shifters



The Human Factor

- System are complex
 - 1000's of interacting computers
 - Complex Software with 100's of thousands of lines of code
- We get the steady state behavior right.
- What about the shifter who does a DAQ system reset 3 times in a row in a panic of confusion?
- The expert playing from home?



Anonymous messing up DAQ system from home

The Microsoft Problem?

Self Healing Systems? (too smart for their own good)



D0 L3 DAQ built from commodity hardware

- 63 VME sources to 82 node processor farm
- 250kB events at 1kHz, without incurring dead time
- Ethernet TCP/IP communication

One year to design and build

- Hardware availability
- Open software tools
- On-hand expertise VME, Linux, TCP/IP
- Prior DAQ experience

Stable performance

- Stable, reliable, low maintenance and smooth operation
- Meets current and future needs of D0



NA60 readout system \rightarrow ALICE prototype

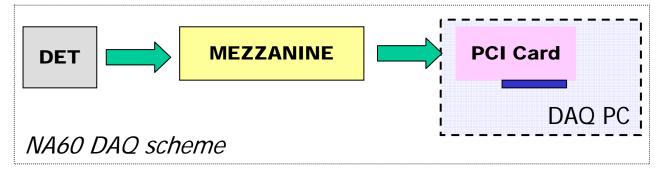
► PCI based system:

Good performances / low cost

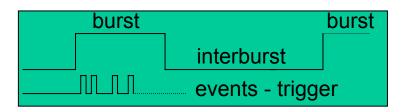
Readout of several different detectors

General purpose PCI readout card

Detector specific mezzanine



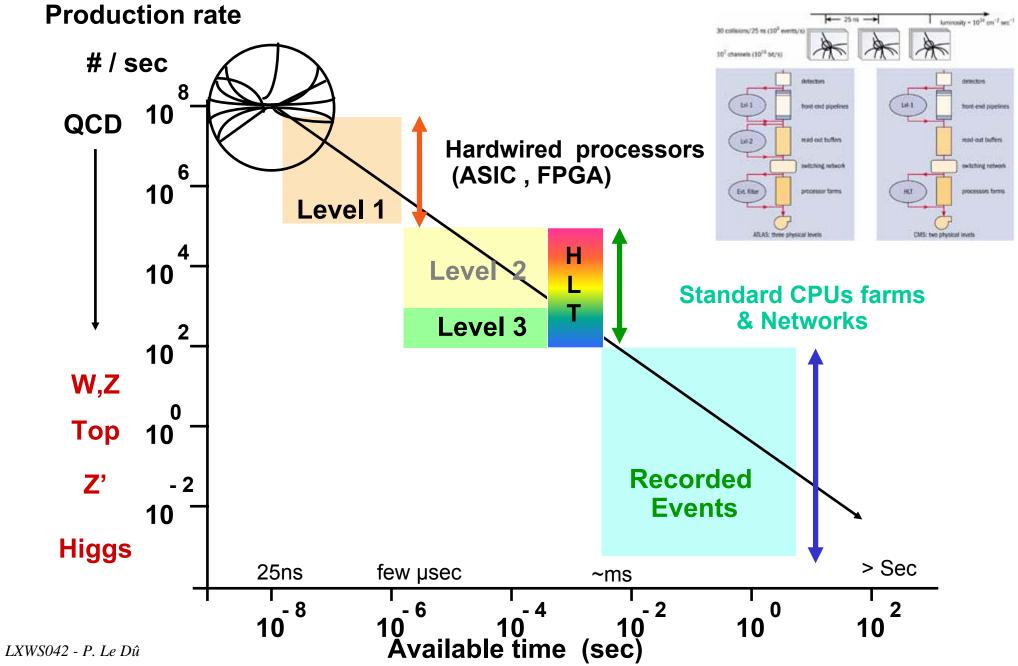
- NA60 readout is spill-buffered
 - DAQ handshake with readout system



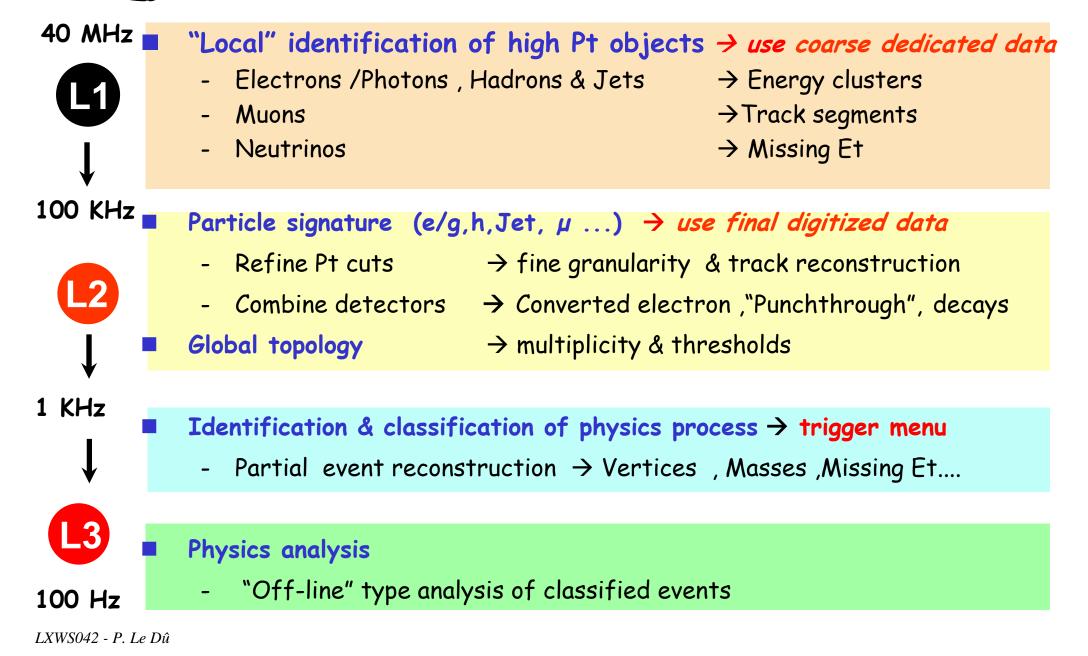


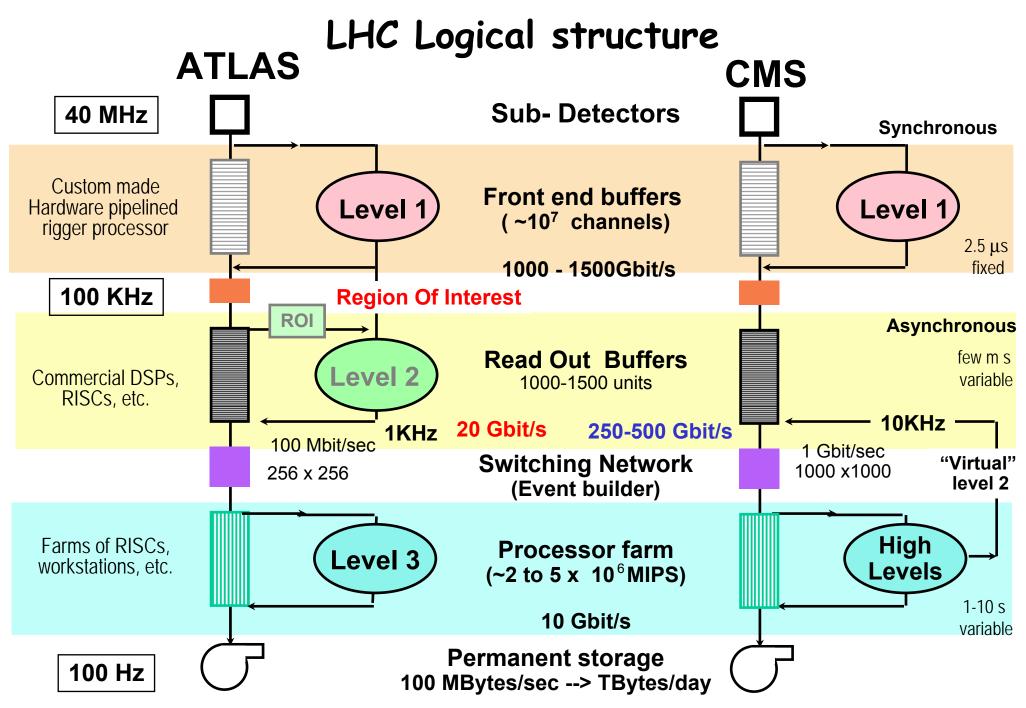


LHC Multilevels Selection scheme



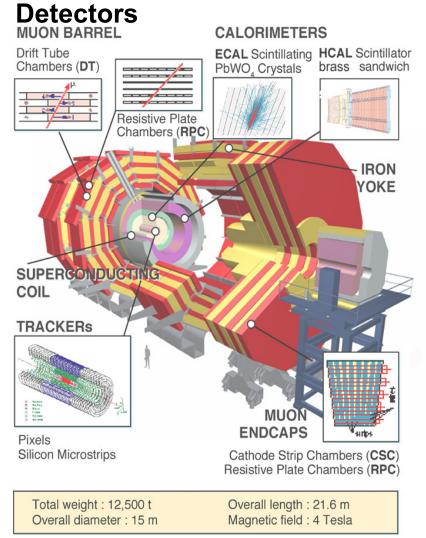
Strategy for event selection @ LHC





LXWS042 - P. Le Dû

Requirements and design parameters



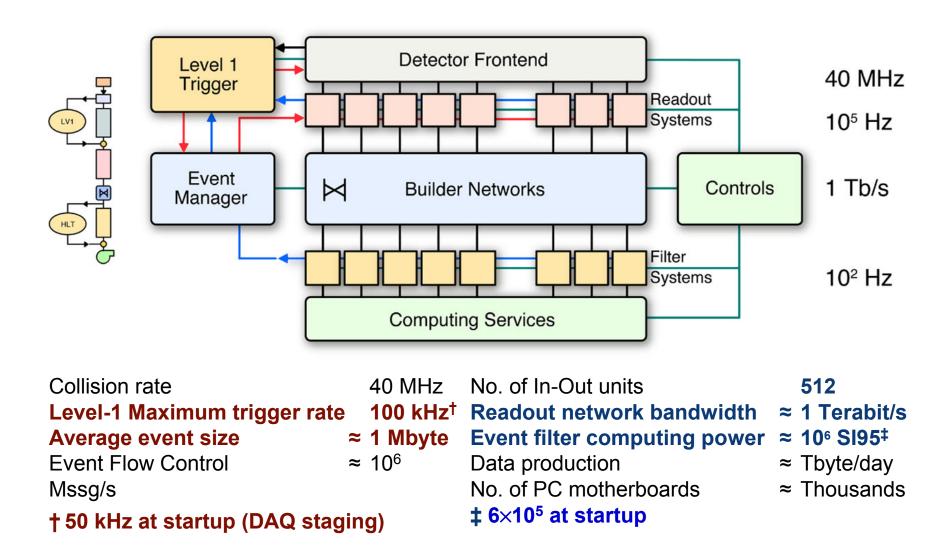
Detector Channels Control Ev. Data

Pixel	6000000	1 GB		70 (kB)
Tracker	10000000	1 GB		300
Preshower	145000	10 MB		110
ECAL	85000	10 MB		100
HCAL	14000	100 kB		80
Muon DT	200000	10 MB		10
Muon RPC	200000	10 MB		5
Muon CSC	400000	10 MB		80
Trigger		1 GB		20
			~	800

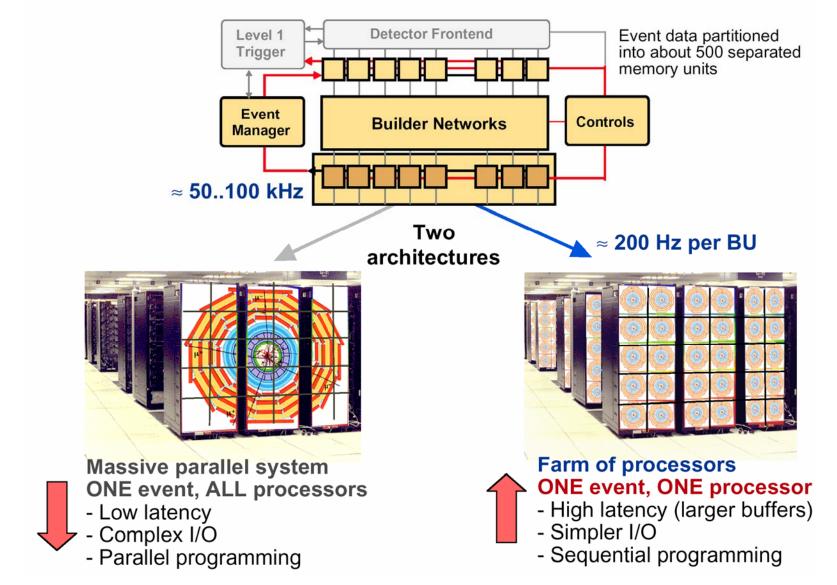
Event size Max LV1 Trigger **Online rejection** System dead time

~ 1 Mbyte 100 kHz 99.999% ~ 0%

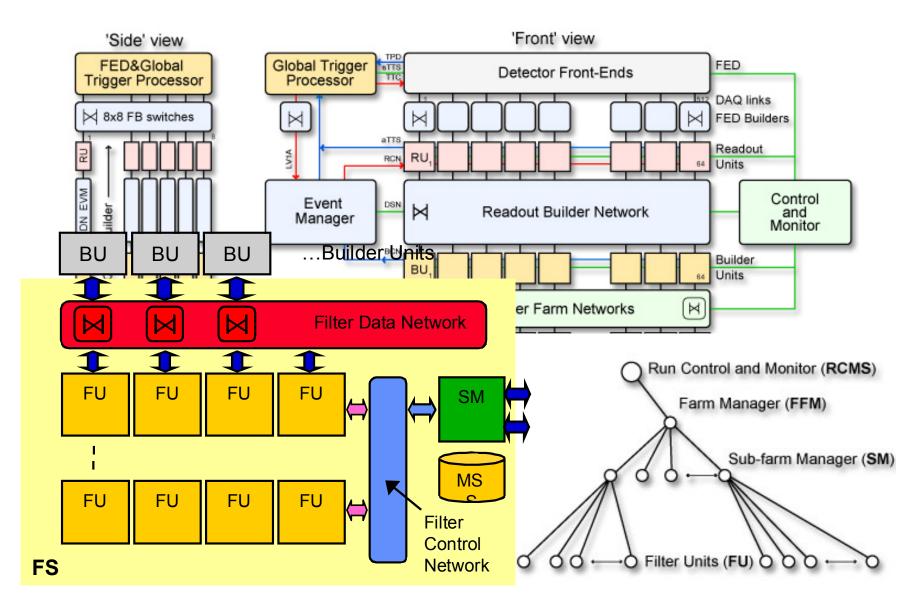
CMS DAQ Baseline



Filter Farm

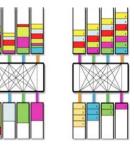


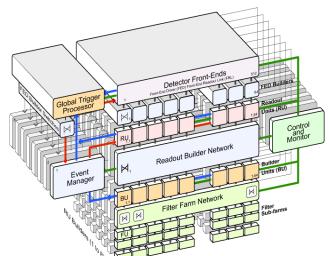
Filter Farm



LXWS042 - P. Le Dû

Event Builder Technology





DEMONSTRATOR results: Myrinet: baseline technology: • Data to surface FED and Readout Builders

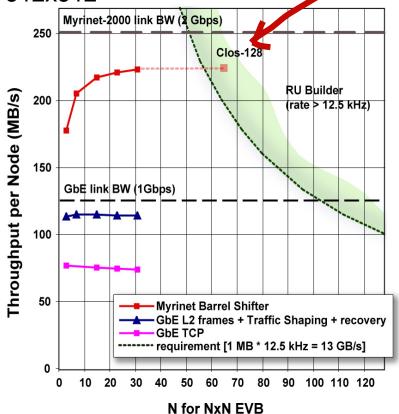
- Feasible with 'today's' (2003) technology

GbEthernet: backup technology: • first results up to 31x31+BM. Raw packet and TCP/IP

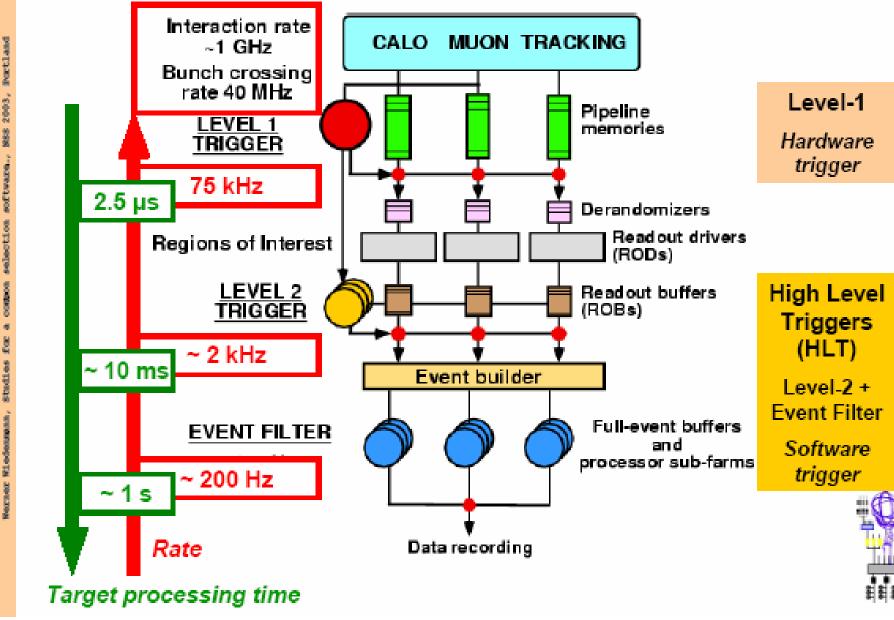
- for nominal conditions (16 kB and RMS = 8 kB): 75 MB/s (60% of wire) speed
- it scales, not clear if scaling persists to larger N or higher switch load
- longer term: with good NIC, fast PC, jumbo frames etc might be possible to reach 100 MB/s per node and twice that for two-rail configuration
- option to track

Two-stage Event Builder Architecture:

- Allows scaling and staging.
 RU-builder: 16 kB fragments at 12.5 kHz easier than 2 kB at 100 kHz
- EVB switch size 8x8 and 64x64 rather than 512x512



The ATLAS Trigger



Regions of Interest concept (ROI)

> Bandwidth/Physics Compromise

- Farm usually reads out entire detector
- HW often look at single detector

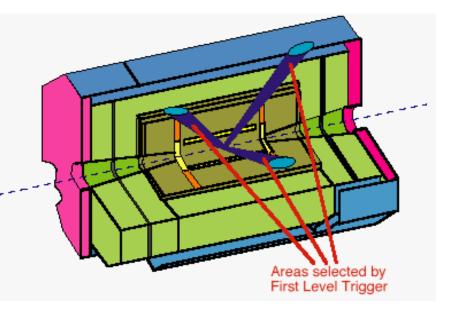
> ROI sits in the middle

- CP Farm CPU requests bits of detector
 - Uses previous trigger info to decide what regions of the detector it is interested in.
- Once event passes ROI trigger, complete readout is requested and triggering commences.

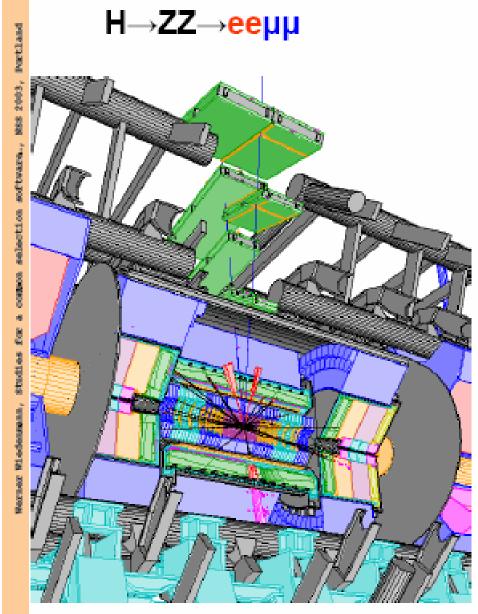
> Flexible, but not without problems.

- Pipelines in front ends must be very deep
- Farm decisions happen out-of-eventorder
 - Pipelines must be constructed appropriately.





Region of Interest (Rol) Concept



- Level-1
 - Uses only coarse calorimeter and muon spectrometer data
 - Local signatures dominate selection
 - No matching of different detectors
- Rol
 - The <u>Region of Interest</u> is the <u>geometrical location</u> of a LVL1 signature (identified high p_T object)
 - Allows access to local full granularity data of each relevant detector
 - <Rol/Level-1 accept> ~ 1.6
- Level-2
 - Seeded with Rol
 - Matching of full detector data within Rol
 - Uses ~ 2% of the full event data for decision

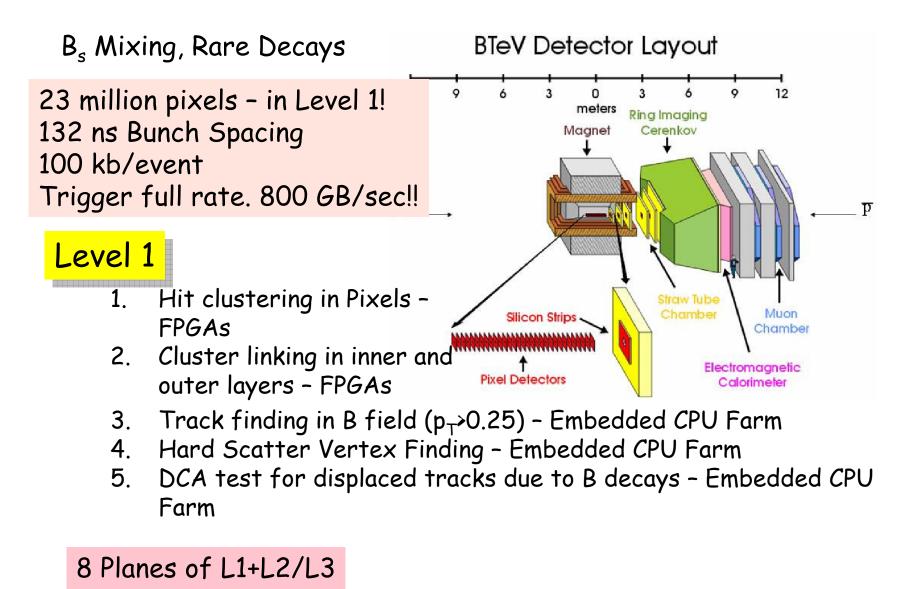
High Level Triggers

- Level-2 (75 kHz→2 kHz, 10 ms target processing time)
 - Selection software runs in Level-2 Processing Unit (L2PU) in multiple worker threads (each thread processes one event)
 - Feature extraction in Rol with specialized algorithms that are optimized for speed and combine information of all sub-detectors sequentially
 - Event Filter (2 kHz→200 Hz, 1-2 s target processing time)
 - Independent Processing Tasks (PT) run selection software on Event Filter (EF) farm nodes
 - Full event reconstruction (seeded by Level-2 result) with offlinetype algorithms which have access to full calibration data
 - Software based on ATLAS offline environment ATHENA/GAUDI
- HLT Selection Software Framework requirements
 - Common to Level-2 and EF
 - Possibility to move algorithms from Level-2 to EF for optimization
 - Development and evaluation in offline environment

Need offline interfaces in Level-2: Steering Controller

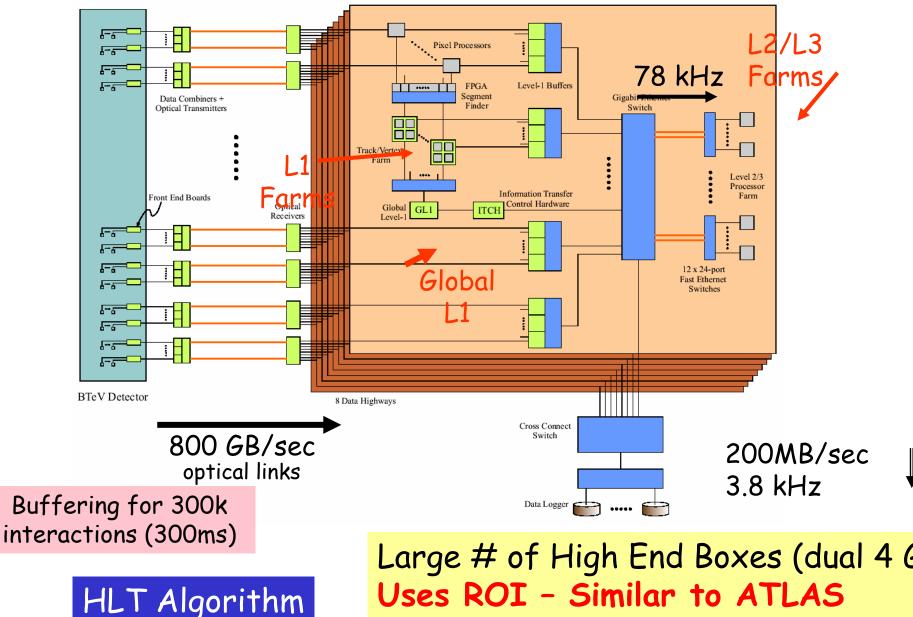


BTeV



 $LXWS042 - P. Le D\hat{u}$ (round robin)

BTeV Trigger Design

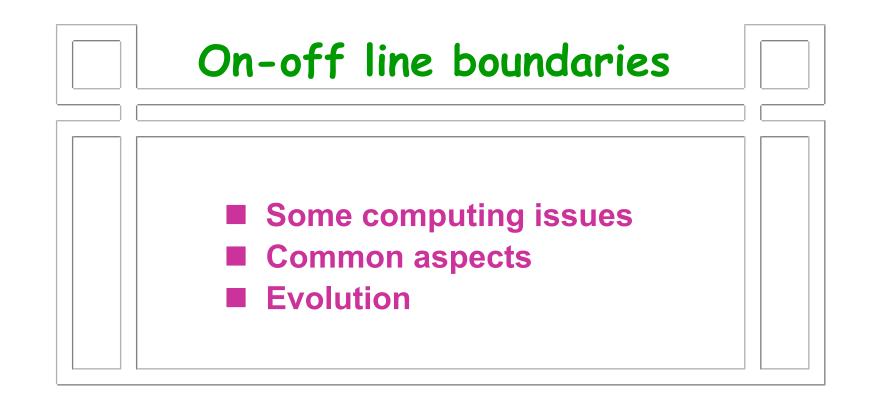


LXWS042 - P. Le Dû

Large # of High End Boxes (dual 4 GHz) Uses ROI - Similar to ATLAS But L2/L3 are interleaved on same box

Evolution \rightarrow 2005-2010

- LHC (ATLAS & CMS) \rightarrow Two levels trigger
 - L1 = physics objects (e/g,jet,m..) using dedicated data
 - L2 + L3 = High Levels « software » Triggers using « digitized data »
- Complex algorithms like displaced vertices are moving downstream
 - New L1 Calorimeter trigger (sign electron, jets and taus)
 - CDF/DO : L2 vertex trigger
 - LHCb/Btev : LO/L1 b trigger
- Use as much as possible comodity products (HLT)
 - No more « Physic » busses \rightarrow VME,PCI ..
 - Off the shelf technology
 - Processor farms
 - Networks switches (ATM, GbE)
 - Commonly OS and high level languages



Detectors will produce a huge data volume



-Few Tbytes/year ! (according to the present models) -Is it reasonable (and usefull!) -Can we reduce it ? \rightarrow 100

A personal, critical .. non conformist view !



About On-Off line boundaries

> Detectors are becoming more stable and less faulty

- High efficiency, Low failure rate
- Powerfull "on-line" diagnostics and error recovery (expert systems)

On-line computing is increasing and not doing only "data collection": More complex analysis is moving on-line

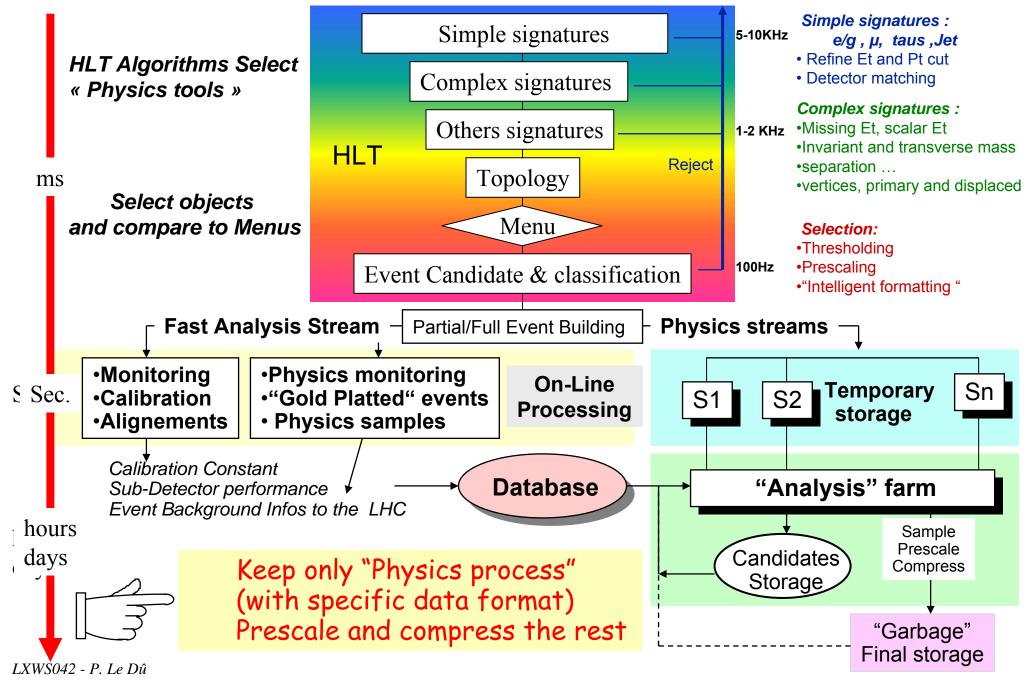
- » "Off-line" type algorithms early in the selection chain (b tag ..)
- » Selection of "data streams" --> Important role of the "Filter"
- » Precise alignment needed for triggering
- » Detector calibration using Physics process available
- » On-line calibration and correction of data possible

> Common aspect \rightarrow

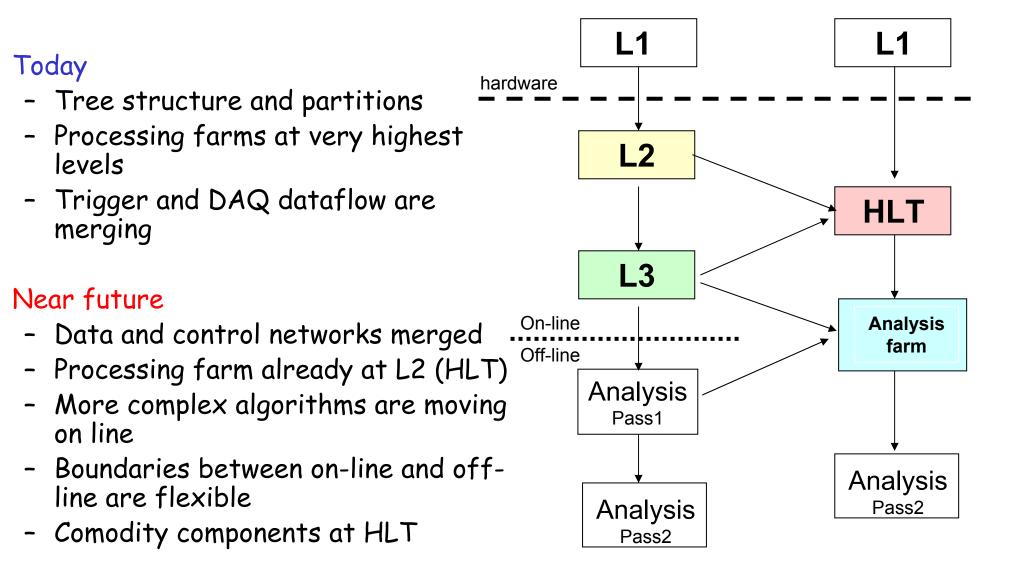
- > Algorithms, Processing farms, Databases...
- > use similar hw/sw components (PC farms..)



Trigger & Event Analysis common strategy

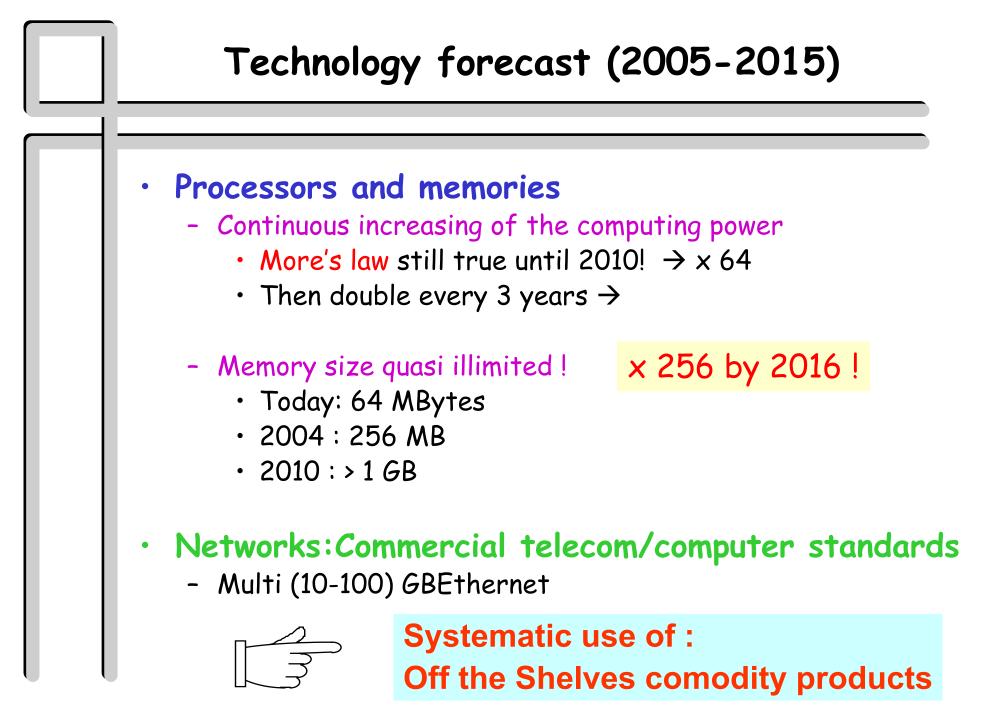


Summary of T/DAQ architecture evolution

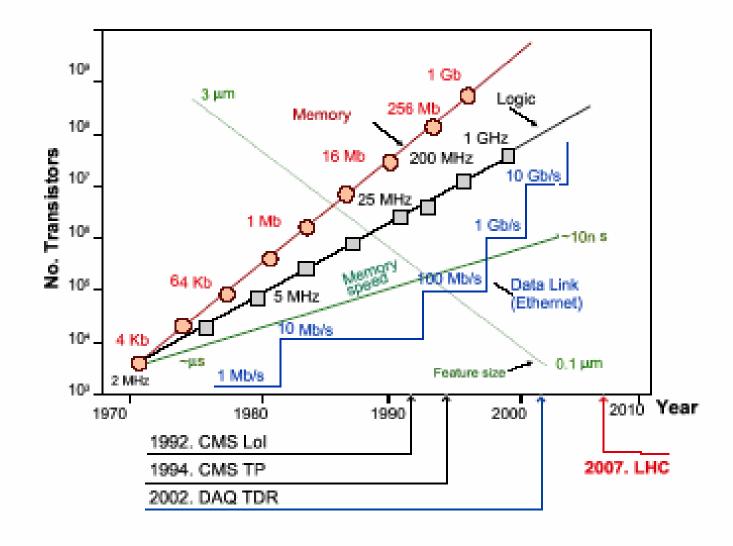


•

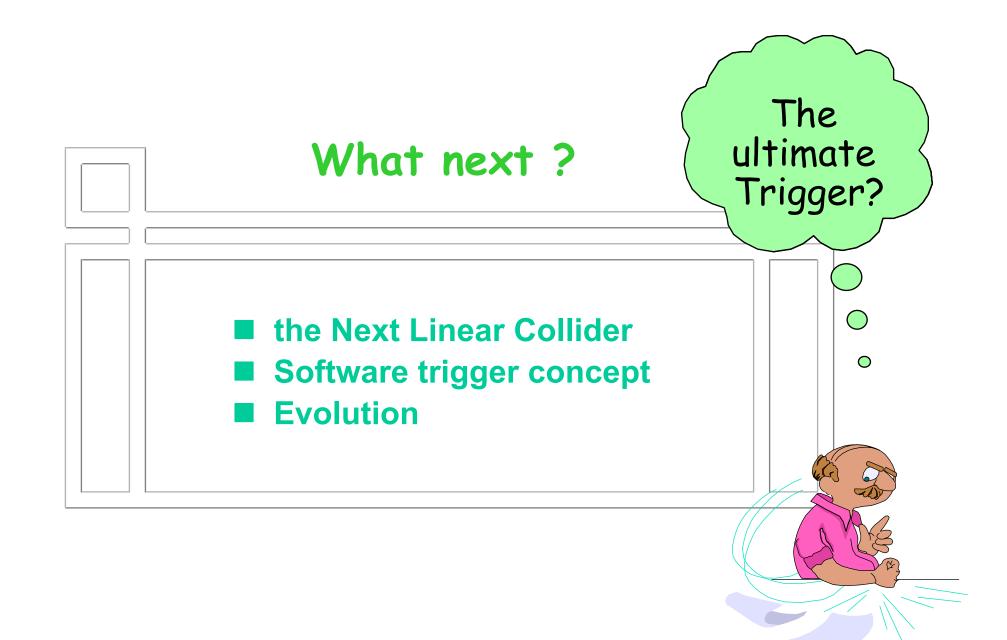
•



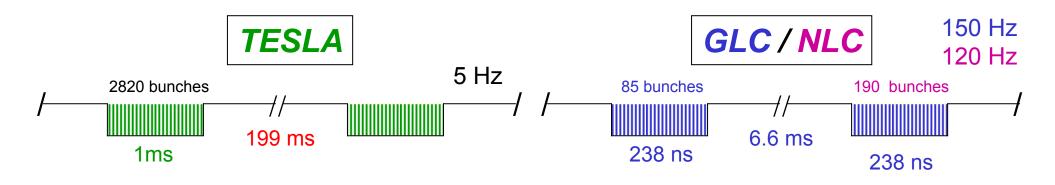
Evolution of technologies since 30 years



LXWS042 - P. Le Dû

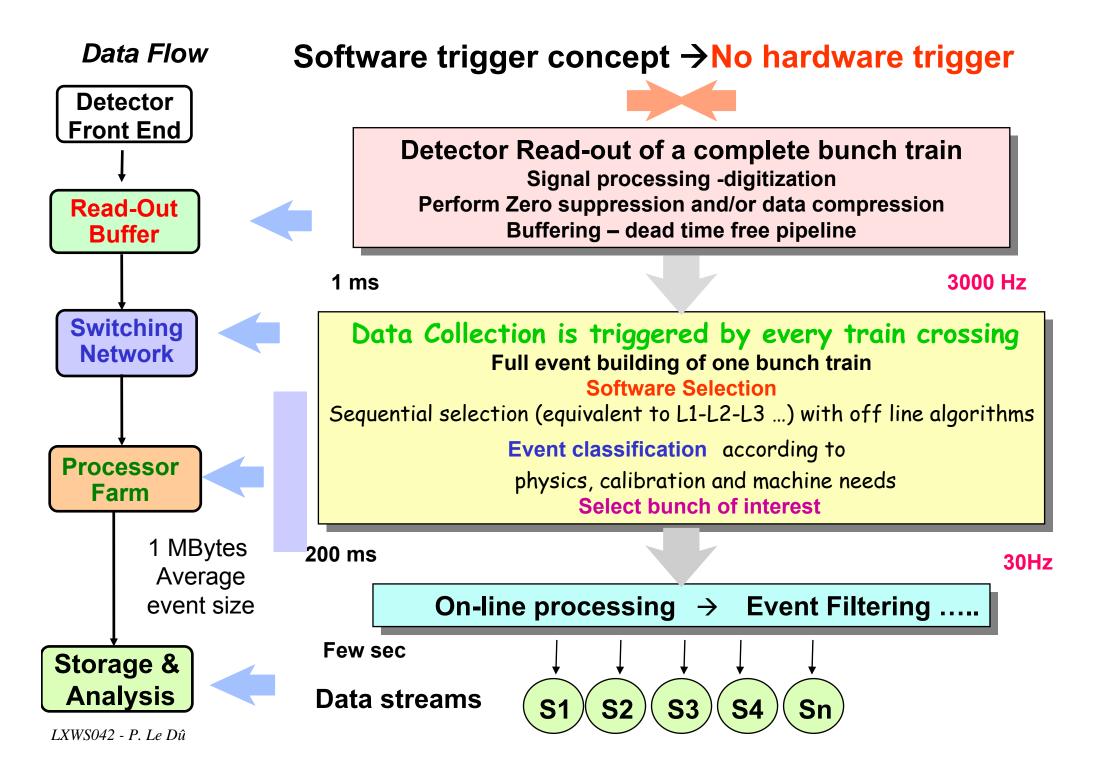


LC beam structure



- Relatively long time between bunch trains : 199 ms
- Rather long time between bunches: 337 ns
- Rather long bunch trains (same order as detector read-out time: 1 ms

- Relatively long time between bunch trains (same order as read-out time): 6.6 ms
- Very short time between bunches: 2.8 ns/1.4 ns
- Rather short pulses : 238 ns



Advantages \rightarrow all

Flexible

- fully programmable
- unforeseen backgrounds and physics rates easily accomodated
- Machine people can adjust the beam using background events
- > Easy maintenance and cost effective
 - Commodity products : Off the shelf technology (memory, switches, procsessors)
 - Commonly OS and high level languages
 - on-line computing ressources usable for « off-line »

> Scalable :

- modular system

$_{ m b}$ Looks like the ' ultimate trigger '

 \rightarrow satisfy everybody : no loss and fully programmable

LXWS042 - P. Le Dû

Consequences on detector concept

Constraints on detector read-out technology

- TESLA: Read 1ms continuously

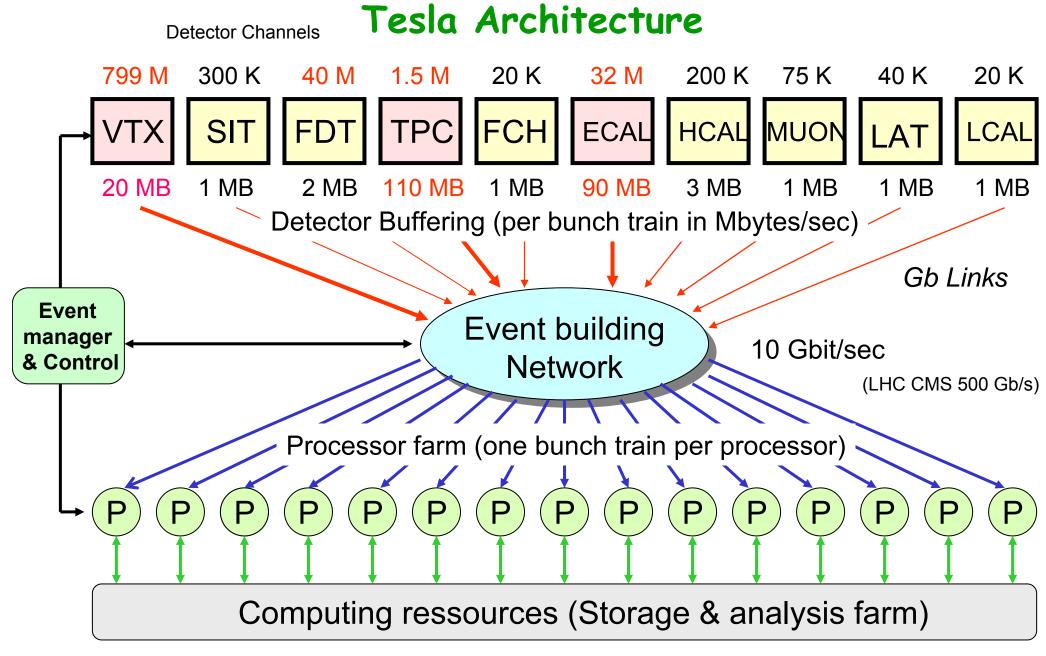
- VTX: digitizing during pulse to keep VTX occupancy small
- TPC : no active gating

- JLC/NLC :

- 7 ms pulse separation
 - detector read out in 5 ms
- 3 ns bunch separation
 - off line bunch tagging

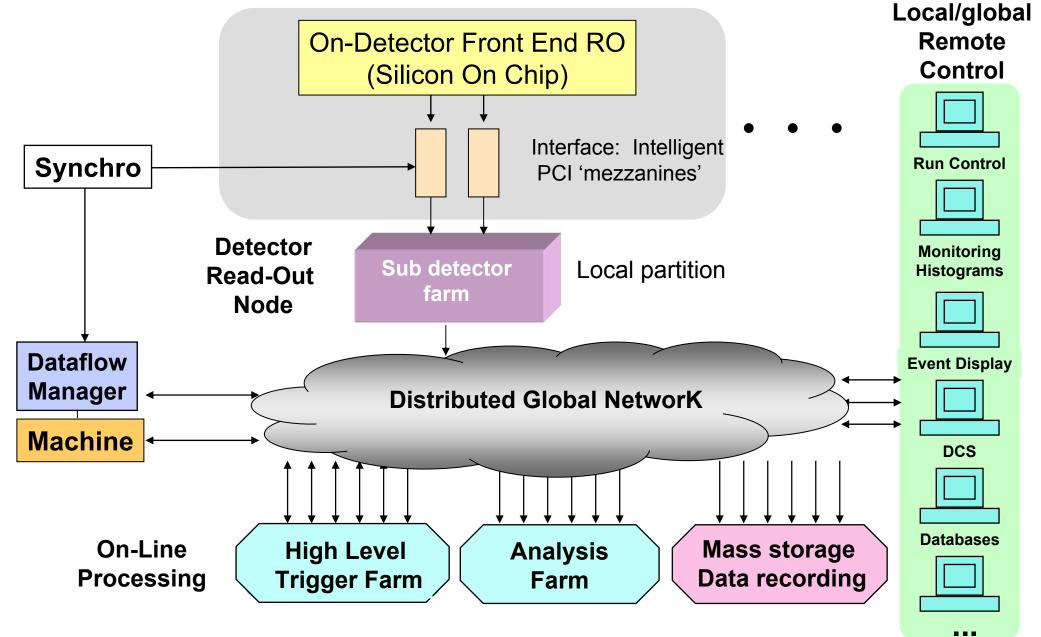
Efficient/cheap read-out of million of front end channels should be developped

- silicon detectors (VTX and SiWcalorimeters)



30 Mbytes/sec \rightarrow 300TBytes/year

FLC 'today' Network model





* Evolution of standards \rightarrow no more HEP!

- HEP: NIM (60s), CAMAC (70s), FASTBUS (80s),
- Commercial OTS : VME (90s), PCI (2000) \rightarrow CPCI?
 - High Speed Networking
 - Programmable Logic Arrays (FPGAs)
 - Commodity CPUs

* Looking ahead \rightarrow where is the market moving now?

- No wide parallel data buses in crates
- Backplanes used for power distribution, serial I/O, special functions
- Small networked devices & embeded systems
- High speeGb/s fiber & copper serial data link
- Wireless devices and very-local area network



Summary

- Higher level trigger decisions are migrating to the lower levels
 Software Migration is following functional migration
 - Correlations that used to be done at Level 2 or Level 3 in are now done at Level 1.
 - More complex trigger (impact parameter!) decisions at earlier times (HLT) → Less bandwith out of detector??

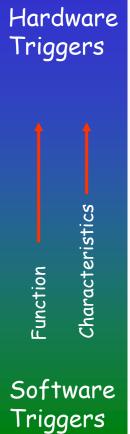
> Boundaries

- L2 and L3 are merging into High Levels Triggers
- DAQ and trigger data flow are merging
- On-line and off-line boundaries are flexible

> Recent Developments in Electronics

- Line between software and hardware is blurring
- Complex Algorithms in hardware (FPGAs)
- Possible to have logic designs change after board layout
- Fully commercial components for high levels.

> "Software trigger" possible for the future LC



FLC issues for the next steps

NLC/GLC scheme ?

- Particularities need initiate discussions... Soon!
- Interface with machine
 - Which infos are needed?
- Gaine experience with Model Driven Technologies
 - Follow XML implementation experience (CMS)
- Modelling using 'simple tools'
 - Quantitative evaluation of the size of the system
- What about GRID ???
- ✤ Define 'boundaries' → functional block diagram
- Integration of GDN & DAQ ?
- Trigger criteria & algorithms