



# Status Report on the LHCb Online System

LHCC Referees Meeting

14 May 2001

Clara Gaspar & Beat Jost

Cern / EP

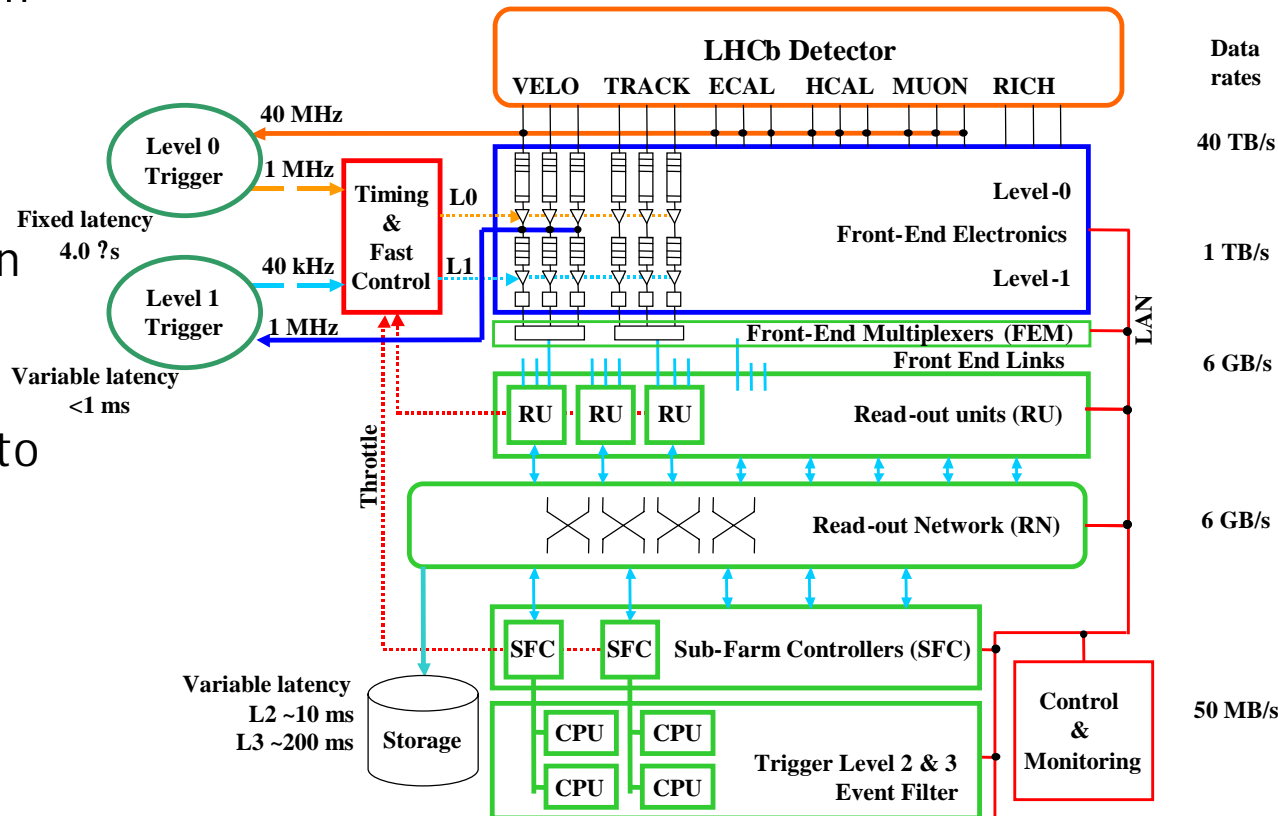
□ The LHCb Online System is composed of

➤ **Dataflow System**

responsible for timing and trigger distribution and the coherent data movement from the front-end Electronics to the data-store.

➤ **Experiment Controls System (ECS)**

controlling and monitoring the operational conditions of the entire experiment



## □ Dataflow System Status

- Timing and Fast Control (TFC)
- Front-End Multiplexing/Readout Unit (FEM/RU)
- Event Building (EB) and switching fabric (Readout Network, RN)
- Event-Filter Farm (EFF)

## □ ECS Status

- SCADA Framework
- Interfacing to Electronics
- Test-Beam DAQ Controls

## □ Summary

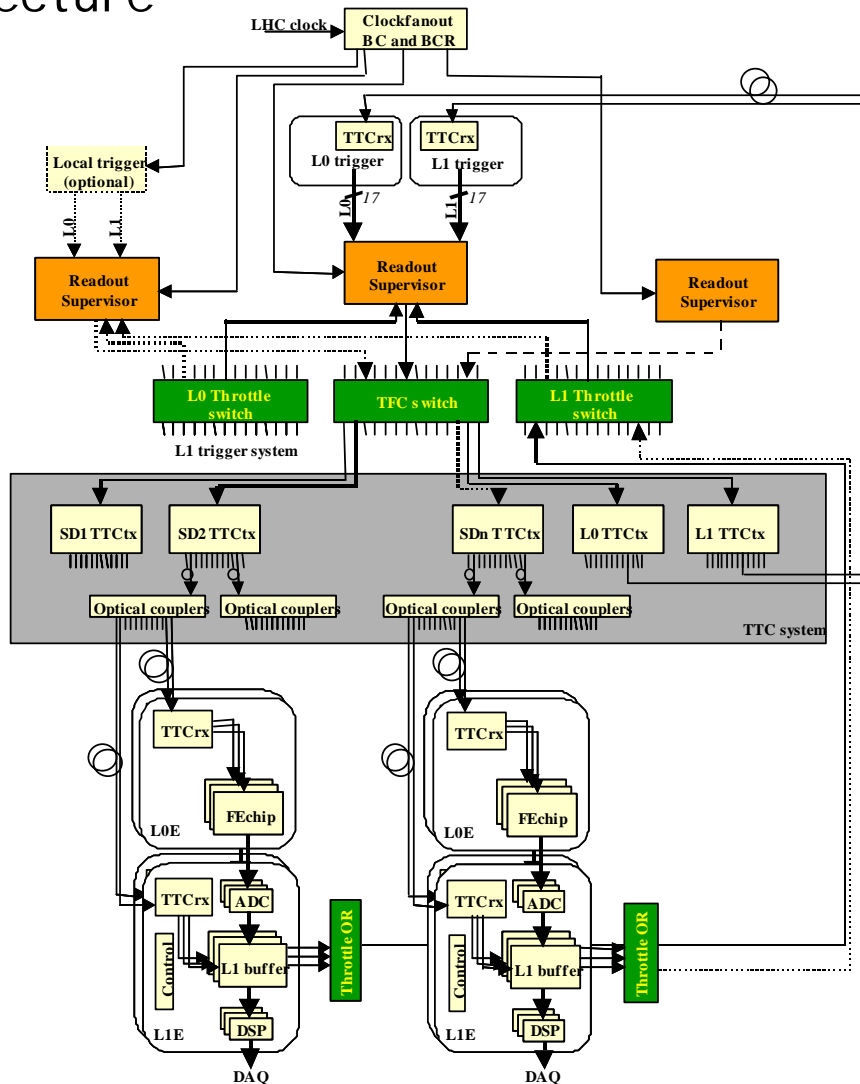
## □ Planning towards TDR



---

# Dataflow System Status

## Architecture



- Two Reviews held with ext. reviewers
  - General architecture and design of switches (Oct 2000)
  - Specifications and Design of Readout Supervisor (RS) (April 2001)
  - Both successful
- Hardware development
  - Switch Prototype (almost) ready
  - RS prototype by October 2001
- Test bed for TFC Components
  - Verified that TTC system can issue broadcasts at Level-0 rate ~1.11MHz (measured rate ~1.6 MHz using TTCvi)
  - Will measure jitter etc. with switch prototype soon.

□ Currently two candidate solutions

➤ FPGA-based

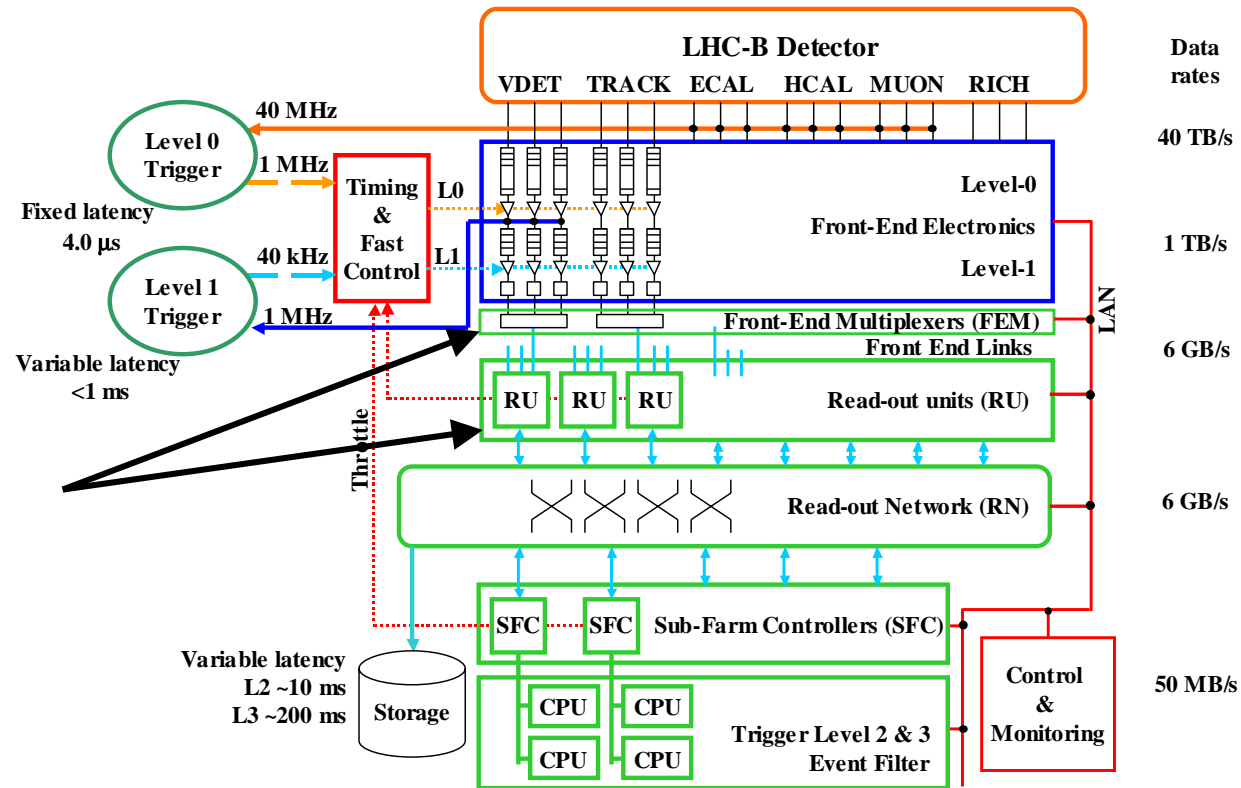
- Conventional technology using standard components
- Prototype exists

➤ Network Processor-based

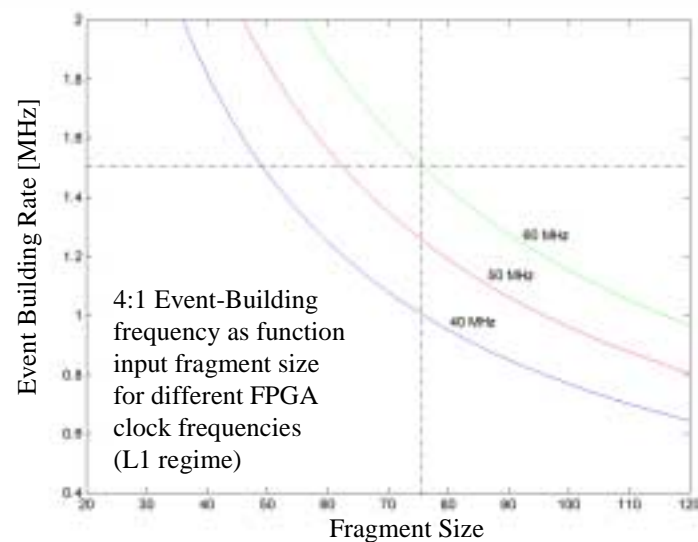
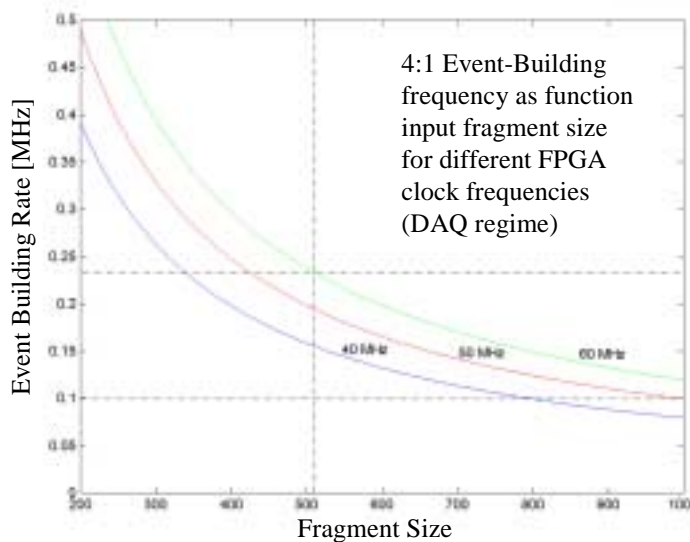
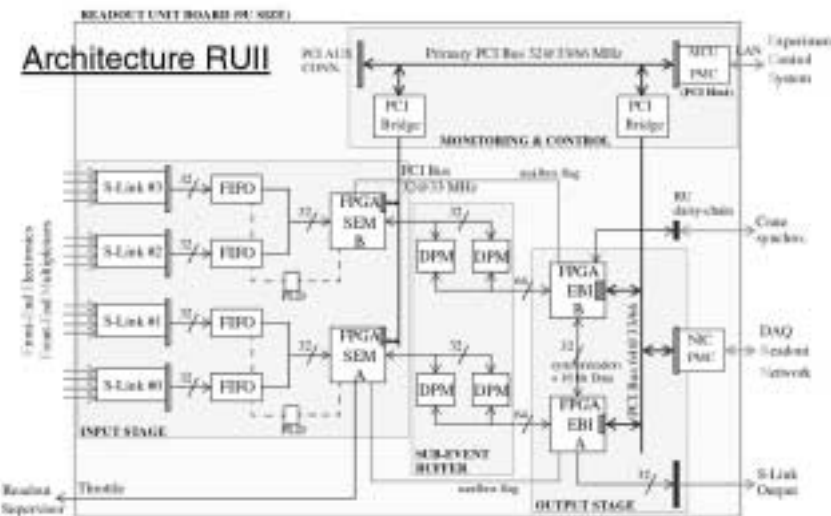
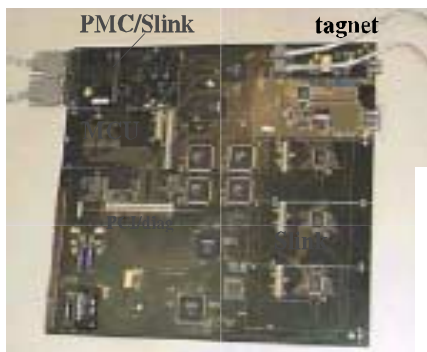
- Fully programmable
- very flexible

□ Review for technical assessment end of July 2001

□ Decision on baseline solution end of August



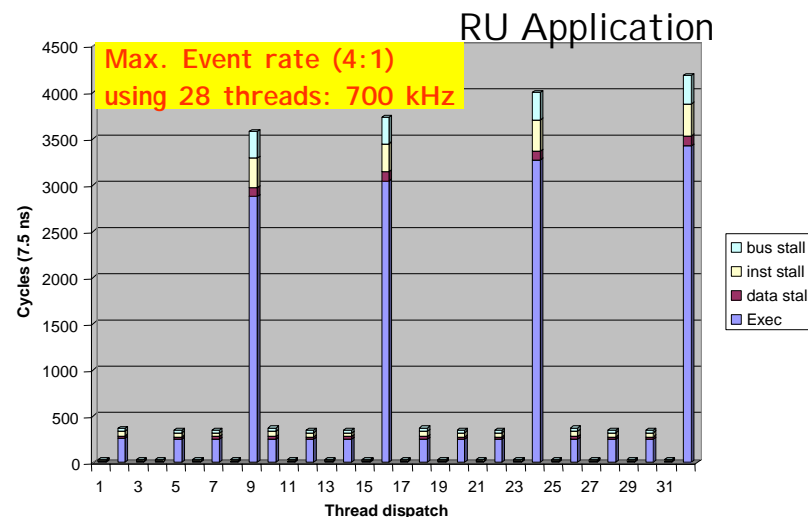
- ❑ New Architecture
  - Better performance, fewer chips
- ❑ Second prototype finished
- ❑ Tests ongoing



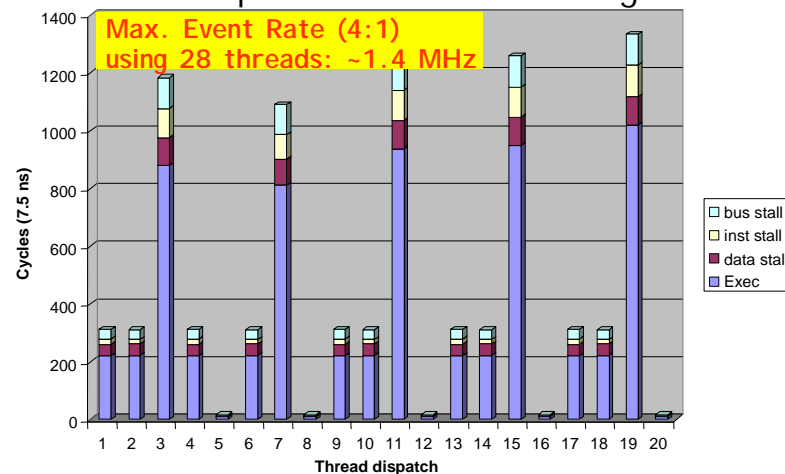
- ❑ Performance simulations show sufficiency for FEM/RU and Level-1 application

- ❑ Network Processors are a new technology gaining very much in momentum in the switch industry. All major chip manufacturers are working on them (IBM, Intel, Siemens, ...)
- ❑ Target market are switch manufacturers using them as input stage of high-speed switches.
- ❑ LHCb software written, debugged and simulated
- ❑ Performance simulated and found sufficient for FEM/RU application and prob. also as L1-RU
- ❑ Hardware reference kit ordered (functionally equivalent to RU but on several boards)
  - Real performance measurements
- ❑ Also some interest in Atlas

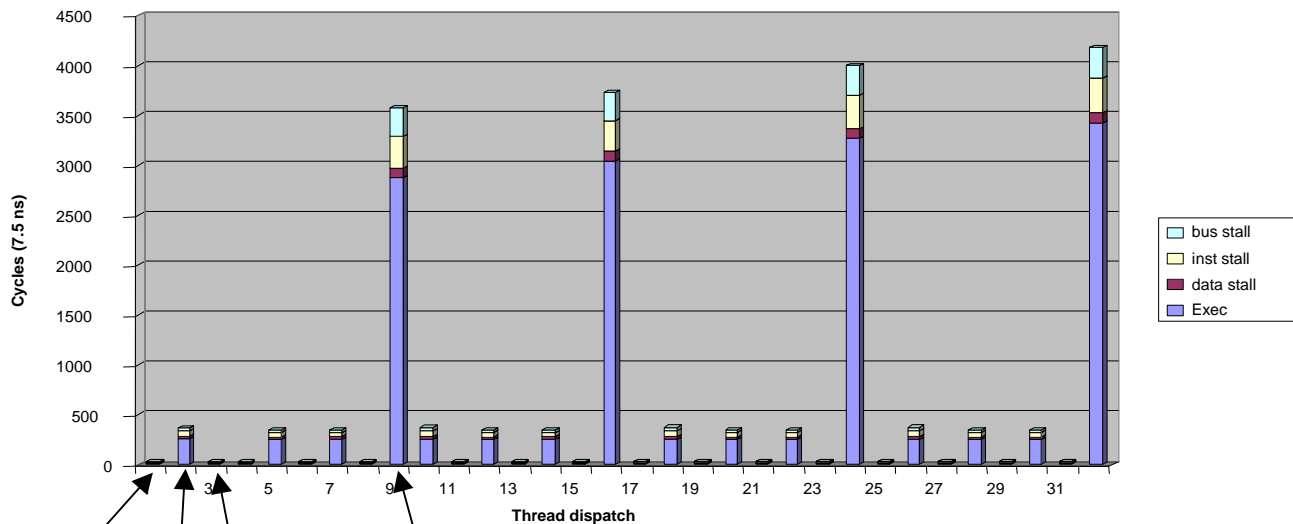
## Simulated Performance (n:1 event-building)



## Optimized for small fragments







Input Handling  
Fragment #1  
Event #1

Input Handling  
Fragment #2  
Event #1

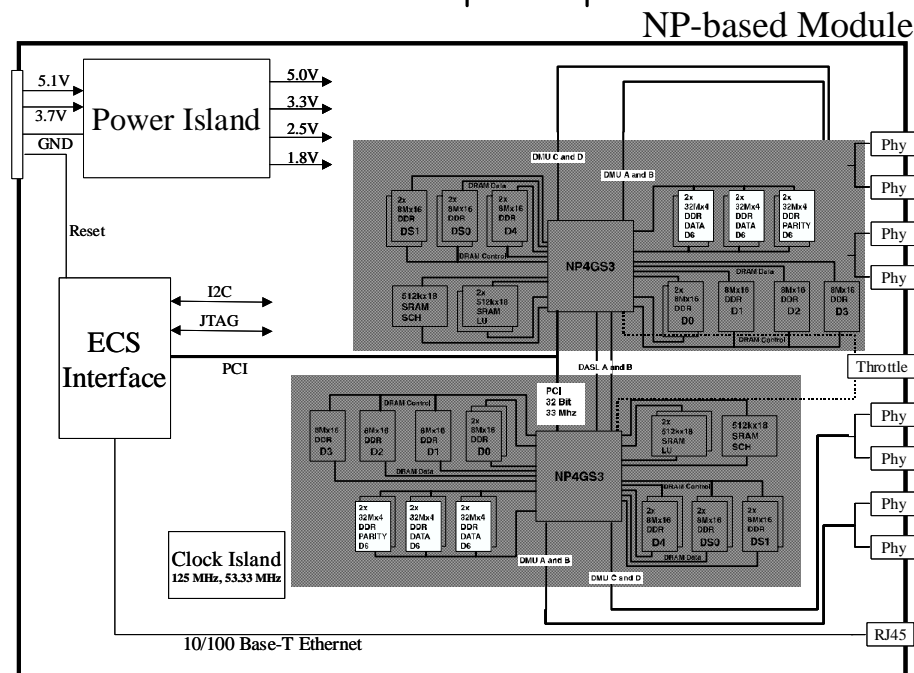
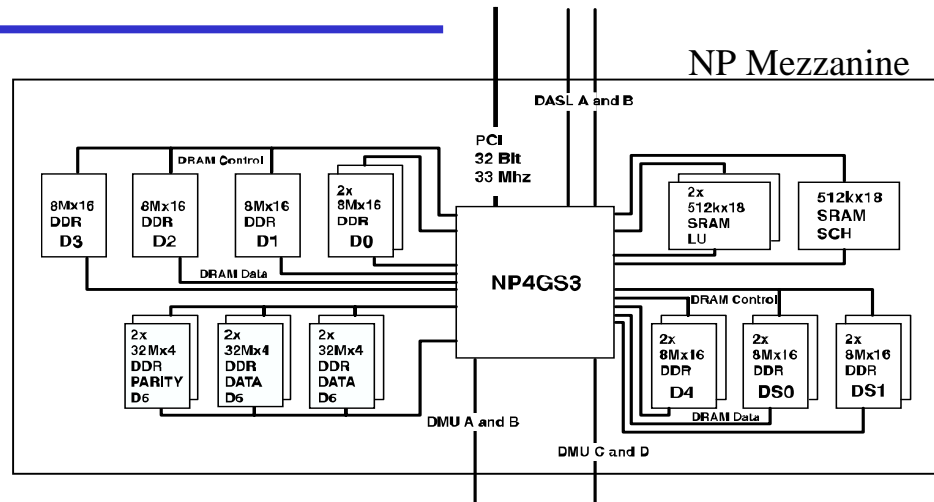
Book-keeping  
Fragment #1  
Event #1

Event-Building  
For Event #1

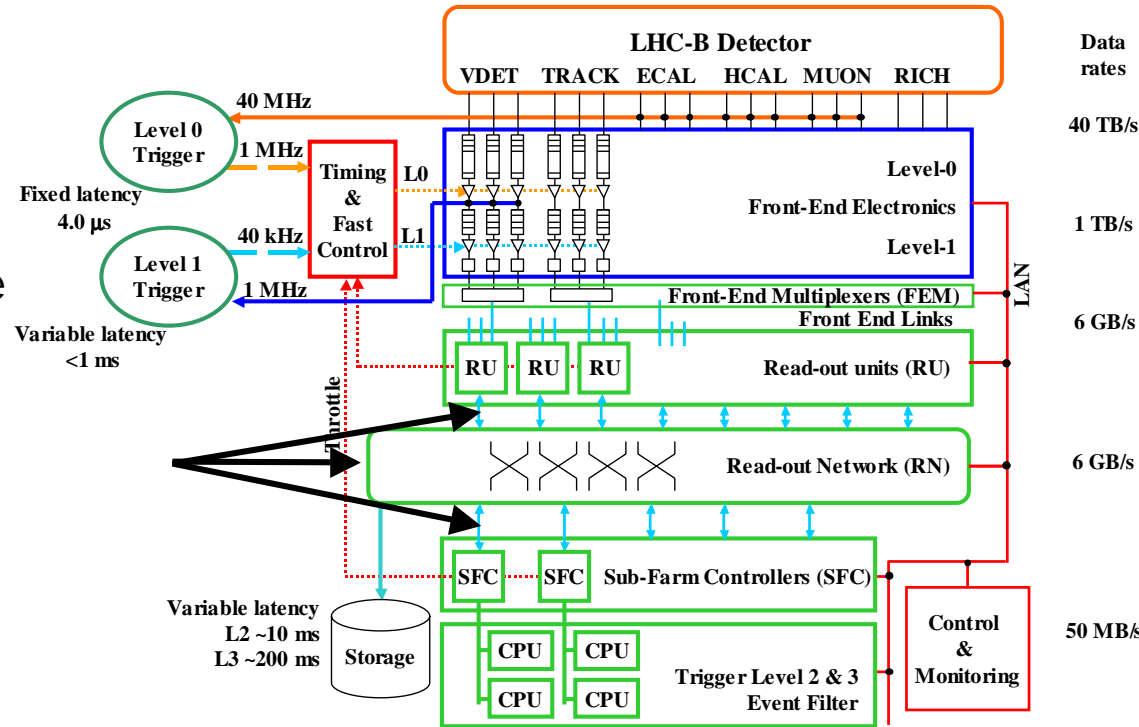
All to be verified with  
real hardware!

- ❑ 4:1 Multiplexing/Event-Building
- ❑ Input handling takes almost no time
- ❑ Book-keeping takes ~400 cycles à 7.5 ns per fragment
- ❑ Event-Building takes ~3800 cycles for 4 Fragments (incl. book-keeping)
- ❑ Time per fragment incl. Book-keeping:  
 $\sim(4 \cdot 400 + 3400) / 4 =$   
 $\sim 1250 \text{ cycles} = 9.4 \mu\text{s}$  or  $\sim 100 \text{ kHz}$
- ❑ There are 28 threads running in parallel, hence 2.8 MHz of fragments can be treated. At 4 fragments per event this leads to  $\sim 700 \text{ kHz}$  of events
- ❑ For small fragments this rate is  $\sim 5.6 \text{ MHz}$  of Fragments or  $\sim 1.4 \text{ MHz}$  of Events (4:1)

- ❑ Possible implementation scenarios being studied. One example could be a mezzanine card containing the NP+memories together with a simple mother board.
- ❑ Very generic board. Can be used as
  - FEM/RU (n:1 Multiplexer/EB)
  - 4x4 switch
  - 4:4 event-builder (input stage of SFC)
- ❑ One single module for all applications in LHCb and only one network technology (Gb Ethernet)
- ❑ Negotiations started on external design/production of boards



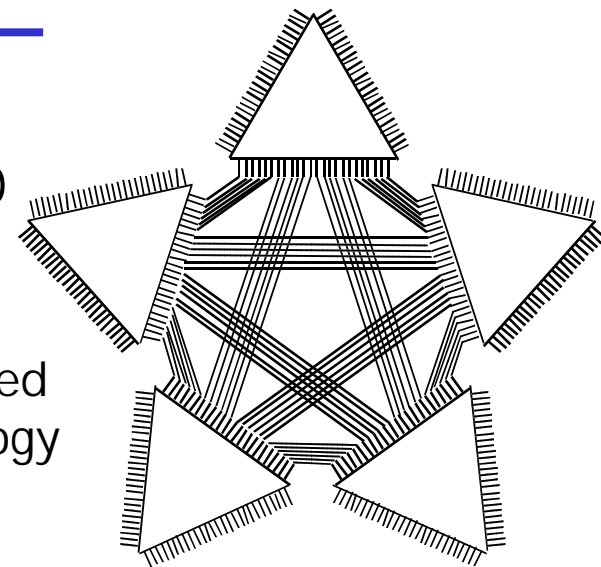
- ❑ No deviation from Full Readout protocol (Push-through, i.e. all data are transferred to farm CPU)
- ❑ Baseline technology is xGb Ethernet (limit to 1 Gb/s at RU/SFC)
- ❑ We have setup a 2x2-port development test bed but also have access sometimes to a 32x32 test bed of CMS
- ❑ Restarting simulation effort to confirm scalability
- ❑ Looking into alternative topologies to minimize number of ports (cost)
- ❑ Studied smart NIC (Network Interface Card)
  - > sufficient performance up-to ~90 kHz L1-rate



- ❑ Relies on the fact that traffic is very unidirectional
- ❑ Working assumption: 4x4 switch as building block
- ❑ Mixing source/destination on same module
- ❑ Still full bandwidth available for event-building

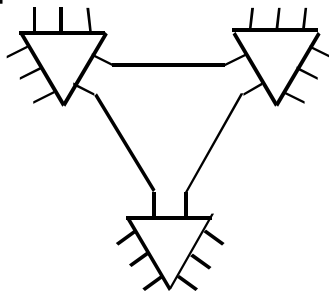
'Final' network:  
270 ports with 90  
4x4 modules

➤ 30% fewer  
Modules compared  
to **Banyan** topology

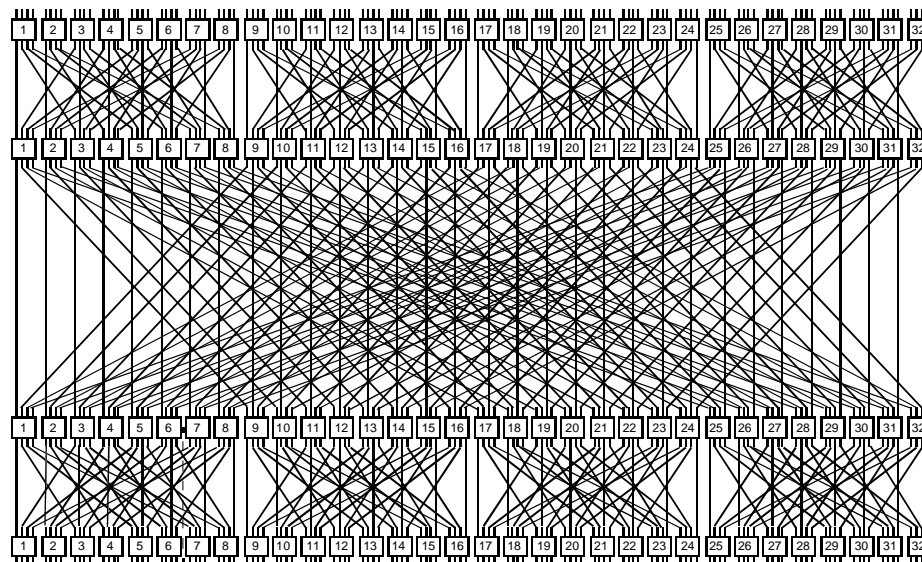


## Basic Building Block

➤ 3 4x4 switches  
to make an 18  
port-switch

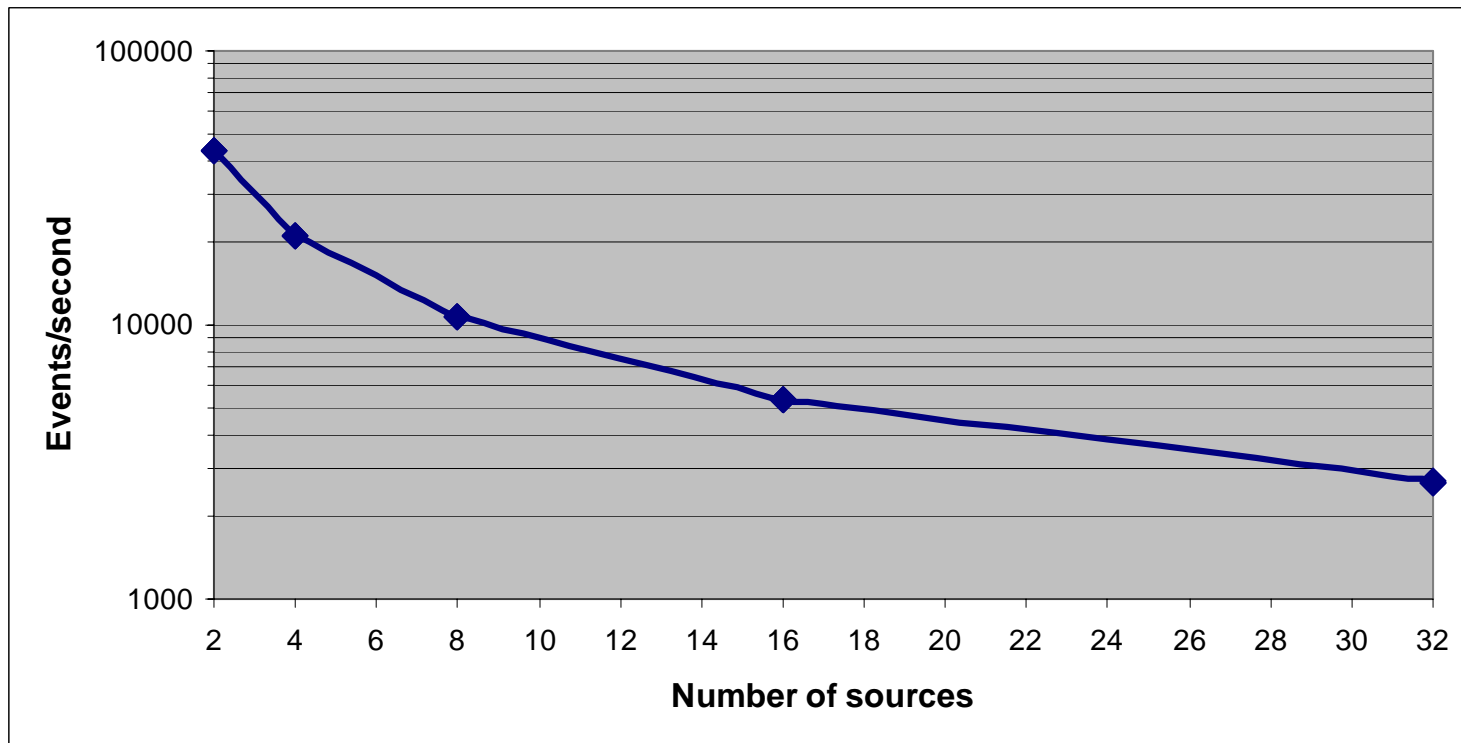


128 X 128 complete connexion based on 32 X 32 sub-switches



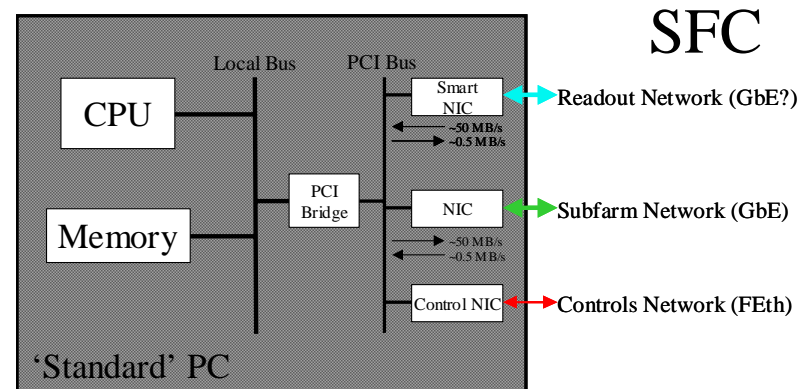
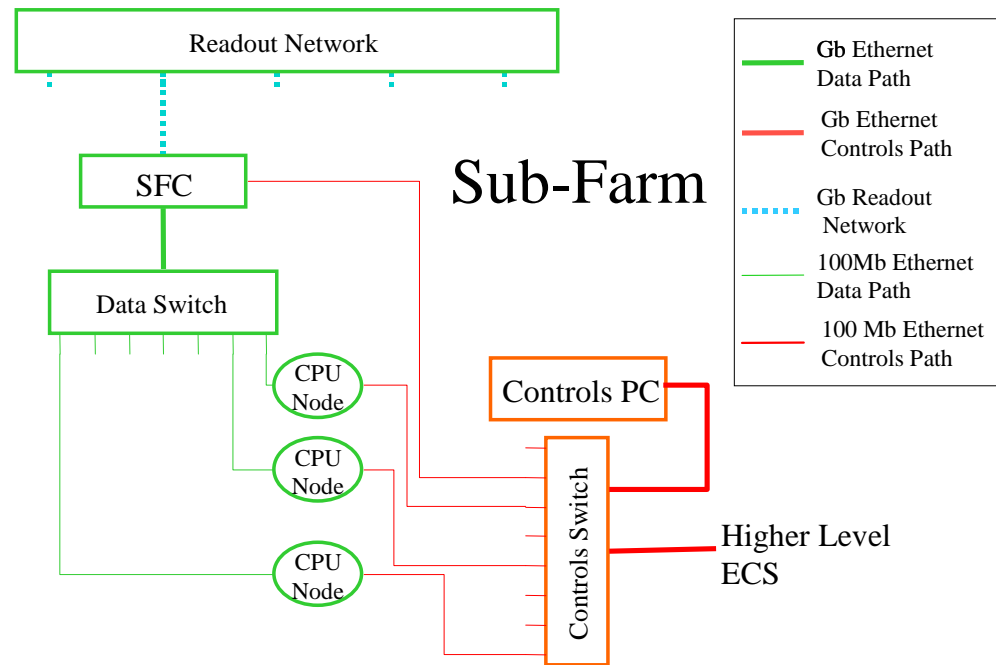
# Smart NIC performance

- ❑ Measurements emulating the sending of data fragments from several sources and event-building at the destination using two smart NICs.
- ❑ Results show that ~90 kHz from 100 sources can be sustained (Code not yet optimized)



# Sub-Farm and SFC Architecture

- ❑ So far very limited effort invested. Ideas...
- ❑ Each sub-farm is an autonomous entity
- ❑ It will have a strictly separated controls and data path, reflected throughout (individual CPU nodes, SFC)
- ❑ Typically 10-20 CPUs per SFC and per controls PC
- ❑ Possibly the controls PCs will have another hierarchical level
- ❑ **Should scale very nicely...**





---

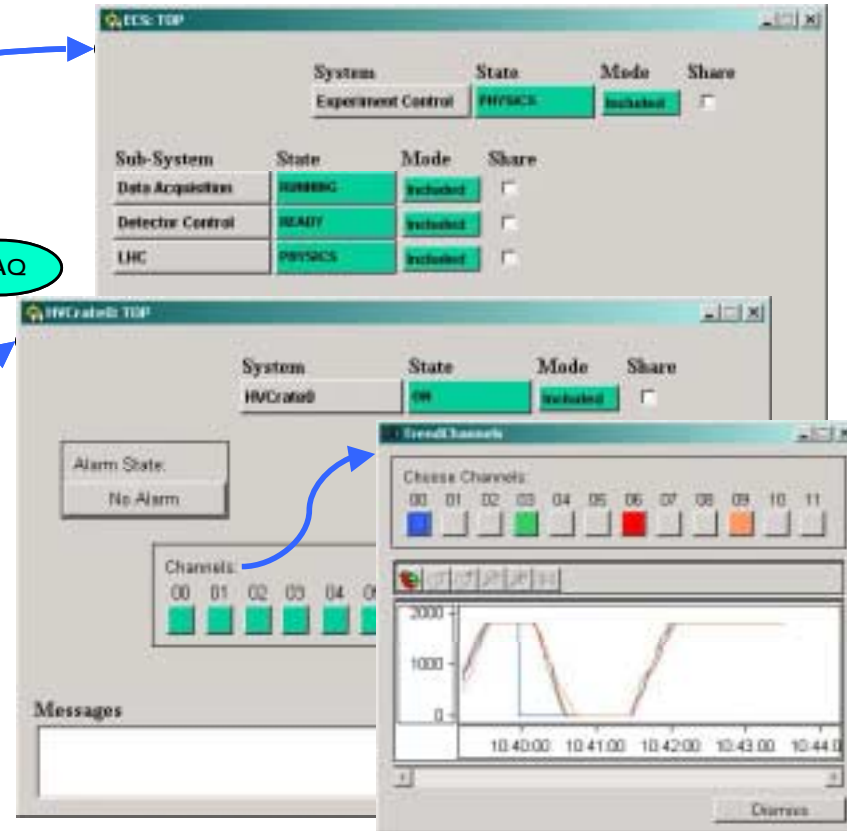
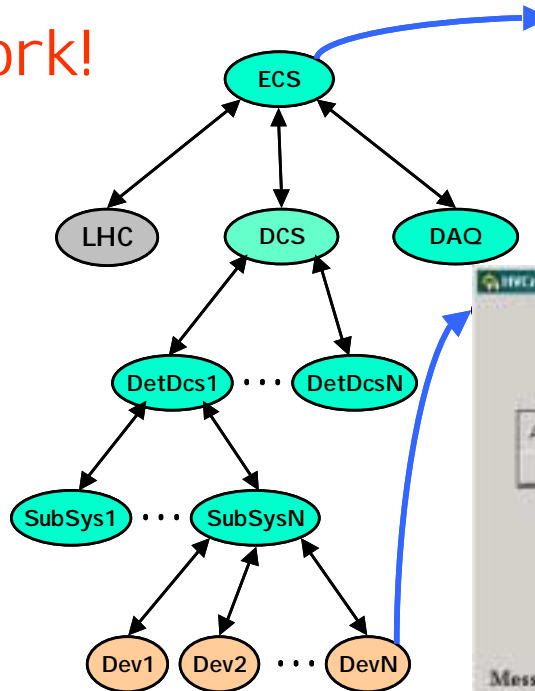
# Experiment Control System (ECS) Status

- ❑ LHCb will have an integrated Experiment Controls System, governing the operational state of the all components of the entire experiment (Detector, DAQ, environment,...)

## ❑ AIM: No Duplication of work!

- ❑ The ECS is based on a hierarchical **architecture**, that is supported and enforced by a common software **framework** (first version available)

- ❑ Components are modeled as **Devices**, devices are, e.g.
  - HV/LV Power Supplies
  - RUs
  - L2/3 CPU node
  - L2/3 Algorithm





Basic software infrastructure is a commercial SCADA (Supervisory Controls And DAq) system (PVSS II), common to all LHC experiments and others...

- Joint COntrols Project (JCOP)
  - SCADA (acquisition, licensing)
  - Architecture Working Group
  - Controls Framework
    - ↳ Including Common Components for HV, LV, etc.
  - Data Interchange Working Group
  - Gas Working Group
  - GIF
  - Rack Control

- LHCb Specific Activities
  - Test Beam Upgrade to PVSS-based run-controls
  - Interface to electronics
  - Eventually, interfacing to "Slow Controls" devices

# Example: CPU Monitoring using SCADA

**Vision\_1: pcomon.pnl**

File Panel ?

Image Name	PID	CPU	CPU Time	Memory
System	2	1	00:03:33	200 K
smss.exe	20	0	00:00:00	20 K
winlogon.exe	34	0	00:00:00	36 K
services.exe	40	0	00:00:04	1768 K
lsass.exe	43	0	00:00:00	636 K
spoolss.exe	68	0	00:00:01	480 K
testloop.exe	69	100	00:00:16	40 K
RpcSs.exe	83	0	00:00:00	664 K
inetd32.exe	89	0	00:00:00	76 K
rtvscan.exe	92	0	00:24:33	3216 K
lprserv.exe	98	0	00:00:05	436 K
msiexec.exe	106	0	00:00:00	308 K
pstores.exe	111	0	00:00:00	68 K
MSTask.exe	114	0	00:00:00	132 K

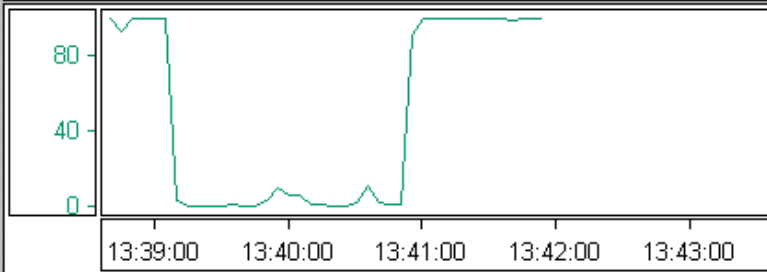
  

Processes	CPU Usage	Memory Usage
39	100 %	60424 K

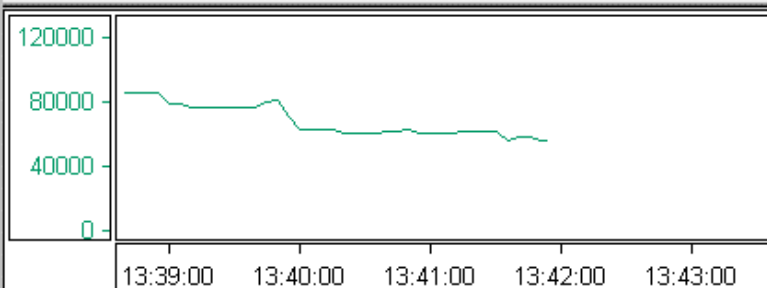
**Performance**

PCEPDELP01

CPU Time (%)



Memory Usage (Kb)

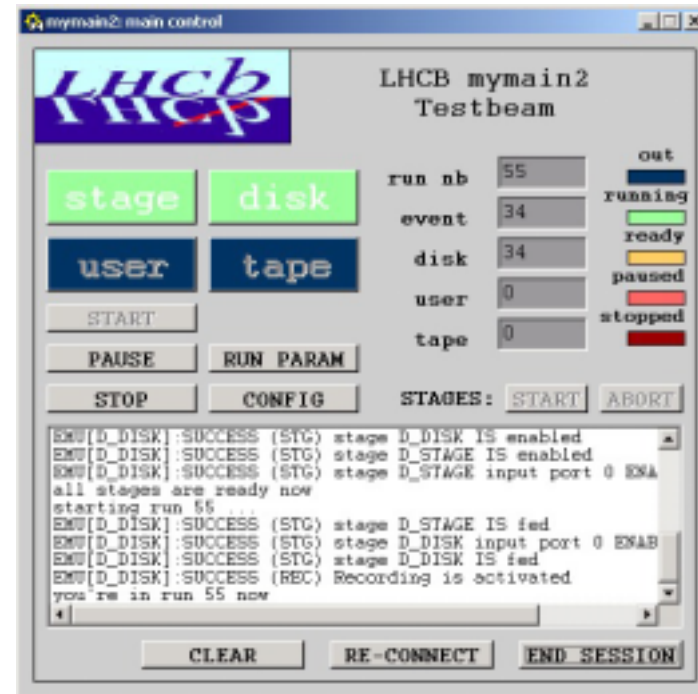


# Test-Beam Controls upgrade

Replace original run-control by PVSS-based run control

- Prove of concept of general idea of integrated controls system (run control & slow control)

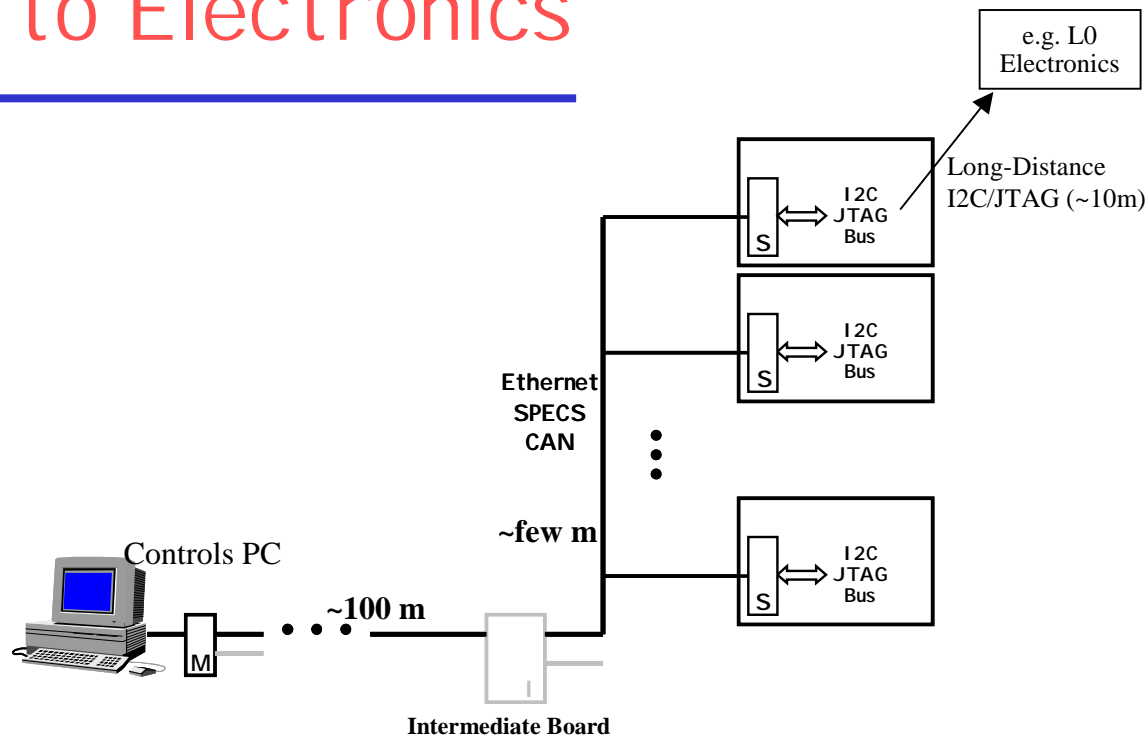
New Test Beam Run-Control Panel



Add other ('slow controls') devices when and where needed

- Data from the PS
- High/Low Voltage
- Gas
- ...

- ❑ LHCb standardized on 3 interfaces to Electronics (board-based)
  - Credit-Card PCs for electronics boards in the barracks
  - SPECS (evolution of SPAC) and CAN-ELMB (Atlas) for areas with radiation (SEU), but no total-dose problem, **provided they are SEU-hard (Sep '01)**.
- ❑ For front-end chips usually I<sup>2</sup>C or JTAG is used for controls and monitoring.
- ❑ I<sup>2</sup>C and JTAG will be generated over long-distance (~10 m) from SPECS/CAN-ELMB



- S (slave) can be
- CC-PC (not in cavern)
  - SPECS
  - CAN-ELMB (Atlas)

## ❑ TFC Test-bed set up

- Verified that TTC system can issue Level-1 decisions at required rate (1.11 MHz) in test-bed
- TFC switch prototype almost ready
- Readout Supervisor designed, prototype in October '01
- Design reviews done

## ❑ FEM/RU

- Second prototype of FPGA-based RU ready and being simulated/tested
- Investigating Network Processors (-> uniform module for entire LHCb DAQ)

## ❑ Event-Building

- Looking at alternative topologies for Readout Network
- Restarting simulation effort of RN
- Measured performance of "smart" NICs, ok.
- Test-bed set up for development and performance measurements

## ❑ First ideas on CPU farm

## ❑ Overall review of Dataflow subsystem planned for September

## ❑ ECS

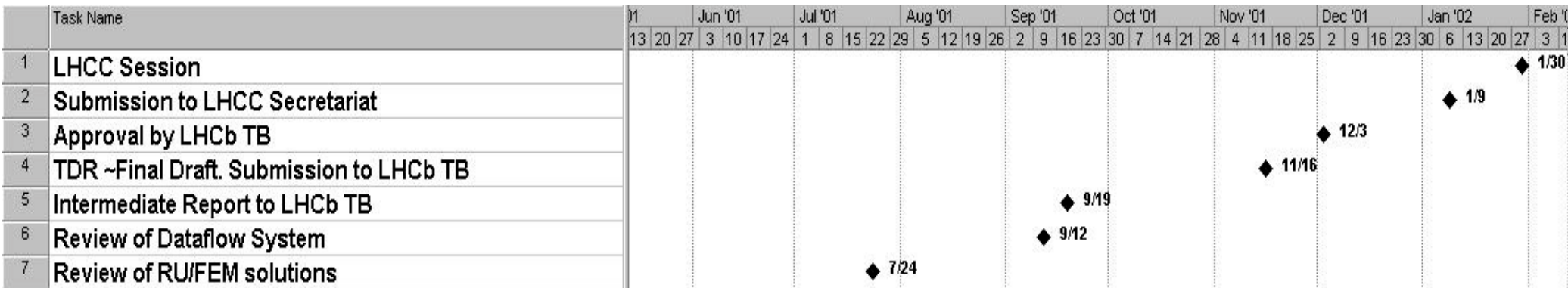
- Integration of Dataflow controls and Slow controls in a single Experiment Controls System
- Hierarchical architecture defined, first version of software framework and tools available
- Hardware interfaces to (board-level) electronics standardized
- Proof of concepts in test-beam



# Planning towards the TDR



# TDR Planning



- Working backwards from the LHCC session in January 2002
- Make use only of the foreseen LHCb weeks (no special TBs, etc.)
- TDR Format:
  - 'Terse' main text with supporting notes (à la TP)