# The ATLAS Read-Out System
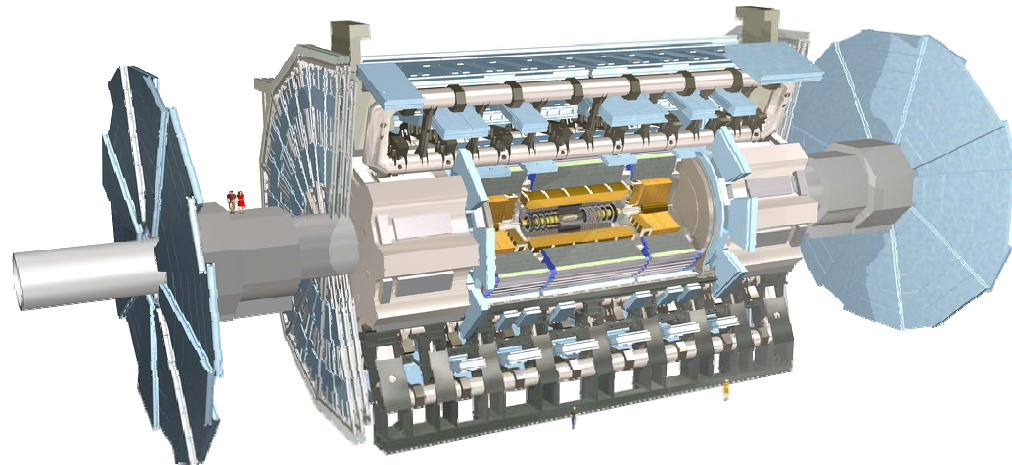## Performance with first data and perspective for the future
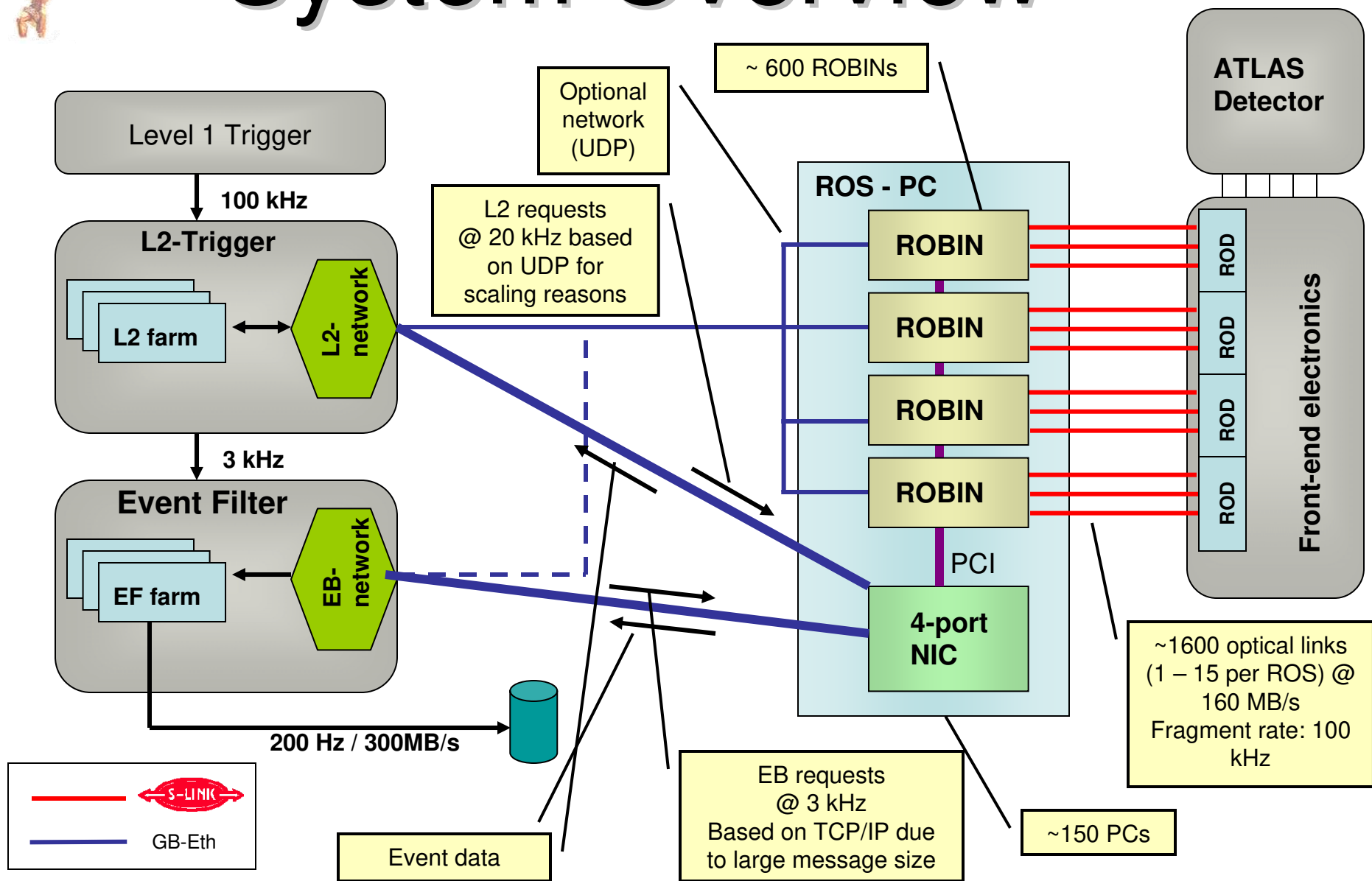
A. Misiejuk[1], G. Crone[2], D. Della Volpe[3], B. Gorini[4], B. Green[1], **M. Joos**[4], G. Kieft[5], K. Kordas[8], A. Kugel[6], N. Schroer[6], P. Teixeira-Dias[1], L. Tremblet[4], J. Vermeulen[5], F. Wickens[7,] P. Werner[4]

[1]Royal Holloway University of London, [2]University College London, [3]Universita & INFN, Napoli, [4]CERN, [5]Nikhef, Amsterdam, [6]Ruprecht-Karls-Universitaet Heidelberg, [7]Rutherford Appleton Laboratory, [8]University Bern

# System Overview



Level 1 Trigger

100 kHz

**L2-Trigger**

L2 farm

L2-network

3 kHz

**Event Filter**

EF farm

EB-network

200 Hz / 300MB/s

Optional network (UDP)

~ 600 ROBINs

L2 requests @ 20 kHz based on UDP for scaling reasons

**ROS - PC**

**ROBIN**

**ROBIN**

**ROBIN**

**ROBIN**

PCI

**4-port NIC**

**ATLAS Detector**

ROD

ROD

ROD

ROD

ROD

**Front-end electronics**

~1600 optical links (1 – 15 per ROS) @ 160 MB/s
Fragment rate: 100 kHz

EB requests @ 3 kHz
Based on TCP/IP due to large message size

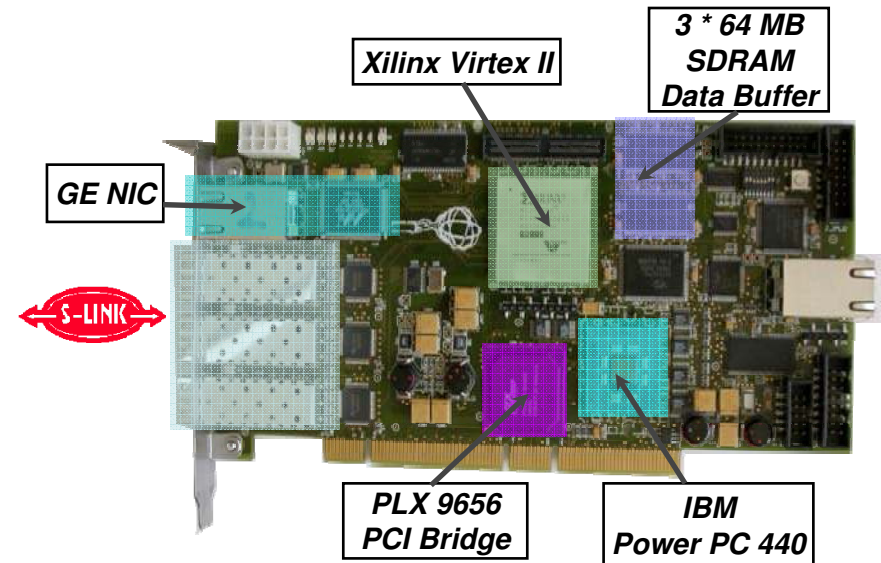~150 PCs

Event data

S-LINK

GB-Eth

# Building blocks

## The ReadOut System (ROS) PC

- Houses 1 to 5 ROBIN cards (typically 4 cards)
- Configures and controls the ROBINs
- Reads data from the ROBINs and provides it to the Second-Level Trigger and to the Event Builder
- Receives clear requests for event fragments and forwards them to the ROBINs
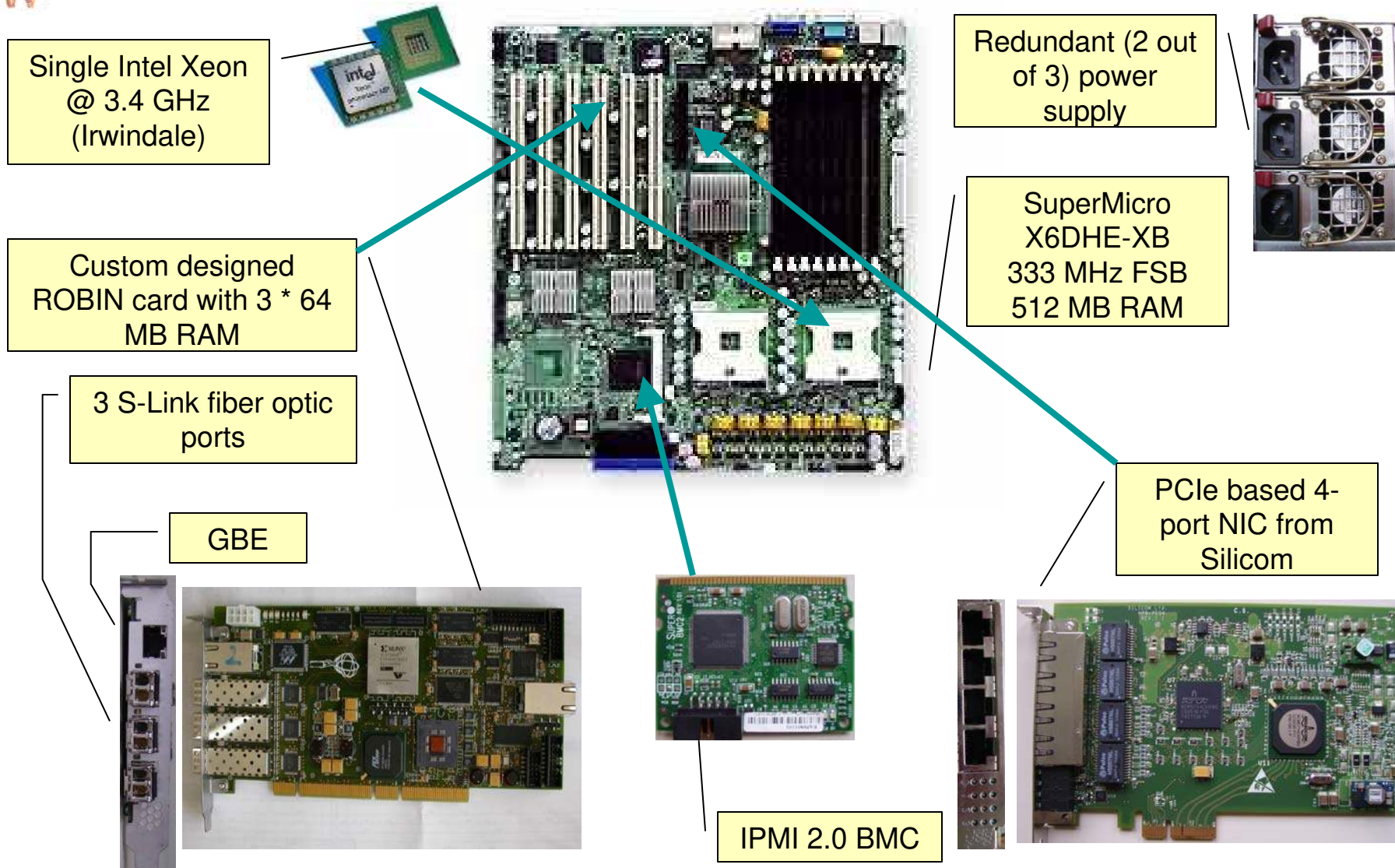- Interfaces to the operational and physics monitoring systems

## The ROBIN PCI card

- Receives event fragments from sub-detector specific front-end electronics (RODs) via 3 optical links
- Buffers events during the decision latency of the Second-Level Trigger and the time required for building of events accepted by L2
- The optical link is based on the S-LINK interface: 32 bit @ 40MHz = 160 MB/s

Xilinx Virtex II

3 * 64 MB SDRAM Data Buffer

GE NIC

S-LINK

PLX 9656 PCI Bridge

IBM Power PC 440

# Hardware Components

Single Intel Xeon @ 3.4 GHz (Irwindale)

Redundant (2 out of 3) power supply

Custom designed ROBIN card with 3 * 64 MB RAM

SuperMicro X6DHE-XB 333 MHz FSB 512 MB RAM

3 S-Link fiber optic ports

GBE

PCIe based 4-port NIC from Silicom

IPMI 2.0 BMC

# Hardware Components



68 ROS PC (liquid Argon sub detector)

15/07/2006

# Operational Monitoring
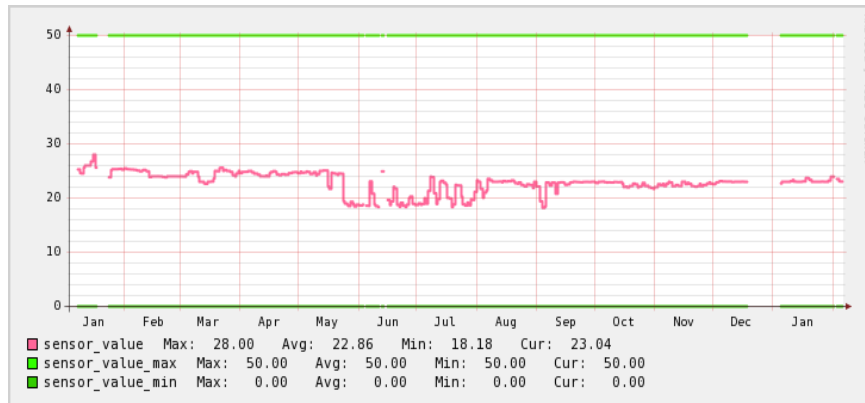
**System status**



Service Status Details For Host 'pc-til-ros-lbc-01'

- Lowest level: IPMI & ssh
- Server level: Nagios
- User level: Web browser

Features:
- History charts
- Automatic E-mail notification in case of problems

**System temperature**



**Fan speed**

# Hardware reliability

| Type of component | Number of installed units | Number of broken units | Failures per year [%] | Failures in 2008 |
|---|---|---|---|---|
| PC Motherboard | 150 | 3 | 0.77 | 1 |
| CPU | 150 | 1 | 0.26 | 0 |
| Memory DIMM | 300 | 3 | 0.39 | 1 |
| Power Supply module | 450 | 4 | 0.34 | 1 |
| IPMI BMC | 150 | 4 | 1.03 | 0 |
| CPU ventilator | 150 | 2 | 0.51 | 0 |
| chassis ventilator | 450 | 1 | 0.09 | 0 |
| 4-port NIC | 150 | 1 | 0.26 | 1 |
| ROBIN cards - broken | 614 | 23 | 1.45 | 1 |
| ROBIN cards – intermittent errors (Firmware issues) | 614 | 36 | 2.26 | 15 |

Average age of the hardware: 2.6 years

# System integration issues

- The individual PSUs of the PCs generate a significant inrush current peak when power is restored after a power cut
- This may cause breakers in the rack to trip
- First solution: Power staggering barrettes
  - But recent questions on long-term reliability
- So now plan second solution
  - Thermistors



Inrush current of ROS PCs



Power staggering barrette

# ROS software Architecture

*The application retrieves data fragments from the ROBINs, combines them in a unique fragment and sends it to L2/EB*

Multi-threaded C++ program running under Linux (SLC4)

**ROS Application**

PowerPC processor in ROBIN runs C program booted from FLASH memory

**Data from detector**

**ROBIN**          **ROBIN**

Request handlers

Event store    PPC          Event store    PPC

**Data requests from HLT**

**Clear requests from DFM**

FPGA          FPGA

Request receiver

ROS request queue

**PCI bus**

ROBIN request queues

- ⬭ = Linux process
- ⬭ = DAQ thread
- ⬭ = Control thread
- ○ = Scheduler

**Data to HLT**

→ Data fragment
→ Control message

# System performance in 2008

- The ROS PCs were powered almost permanently
- Most of the time they were used for data taking with a trigger on cosmic events
- ROS data output statistics for August to October:
  - ~900 TB of data
- All detectors were commissioned at full S-Link (ROD-ROS) speed
- During selected periods high-rate tests were performed with pre-loaded data (ROBIN not involved)
  - 136 ROS PCs delivered 5.3 GB/s to HLT (40 MB/s/ROS)

A continuous run of more than 4 days

Output of ROS to HLT during high rate test; 40 MB/s

# Performance of the standard ROS

## Test Setup

**GbEth.**

**Requester PCs, emulate HLT nodes**

### Lab measurement

Differences wrt deployed system:

- All network traffic based on TCP/IP
- ROBINs generate data internally



- Canonical fragment size in ATLAS: ~ 1 kByte (256 words)
- Meets the original requirements (~20 kHz L2 rate) for small fragment sizes
- Too slow for large fragments
- Bottleneck seems to be the ROS CPU (ROBINs & NICs can sustain much higher rates)

# Performance of the optimized ROS

The original configuration of the OS & drivers turned out to be inefficient.
Finally we obtained the best performance by:

- Turning hyperthreading off
- Using a uni-processor kernel
- Tuning the interrupt coalescence of the network driver
- Changing the SELinux configuration



3 KHz EB
3 ROLs per L2 request

2 GBE links
Optimized system

100 kHz L1 rate
TCP/IP only

- L2 rate (optimized)
- L2 rate (default)
- Bandwidth (optimized)
- Bandwidth (default)

- Performance is now OK for all fragment sizes
- However would like more headroom

# Higher ROS performance - motivation

- ## ATLAS (upgrade) phases
  - Phase 0 (until 2013, luminosity: up to $1*10^{34}$ cm$^{-2}$ s$^{-1}$)
    - Need more ROS performance to:
      - have headroom for ROS PCs with high L2 request rates
      - compensate for higher rates due to modified thresholds of the L2 trigger
      - allow for additional bandwidth-demanding types of triggers. E.g.:
        - » Inner detector full scan for b-physics
        - » Calorimeter full scan for missing $E_T$
  - Phase 1 (2013 – 2017, luminosity: up to $3*10^{34}$ cm$^{-2}$ s$^{-1}$)
    - Higher data rates due to increased luminosity
    - Still use (current) ROS PCs & ROBINs
    - Requires more network bandwidth (switches, ROBINs & ROS)
  - Phase 2 (from 2018, luminosity: up to $10*10^{34}$ cm$^{-2}$ s$^{-1}$)
    - Much higher data rates
    - Replace ROS system

# Higher ROS performance - options
## (for phase 0 & 1)

- The main bottleneck of the ROS is the network interface to the HLT
  - Only 2 GBE links per ROS
  - Network protocol (mix of UDP and TCP/IP) handled by the CPU of the ROS PC
- 3 approaches to solve the network limitation

| | | |
|---|---|---|
| Install smart NICs (to offload CPU from the TCP/IP protocol) | Replace the motherboard, CPU and memory of the ROS PCs with faster hardware and connect additional GBE lines from the ROS PCs to the HLT network | Connect the ROBINs directly to the HLT network |

GBE port of ROBIN (UDP only)

SuperMicro X7DB8-X MB with 2 * 2.66 GHz quad core Xeon and RAM @ 667 MHz

Chelsio S320E 2-port smart NIC

This may be one of the last MBs with >3 64bit PCI slots -> Development of PCIe based ROBIN started

# Impact of the Smart NIC

Laboratory measurement
- TCP/IP only

- No significant gain in performance
- Driver for smart NIC not yet tuned for best performance (waiting for optimized driver from manufacturer)
- More work required but overall potential seems to be low

3 KHz EB
3 ROLs per L2 request

2 GBE links
Optimized system

100 kHz L1 rate
TCP/IP only



**External L2 request rate (kHz)** / **Total transfer bandwidth (MB/s)** vs **Event fragment size (words)**

Legend:
- L2 rate smart NIC
- L2 rate default NIC
- Bandwidth smart NIC
- Bandwidth default NIC

# Impact of faster MB, CPU & RAM



Test Setup

fast ROS PC

NIC | PCI | ROBIN ROBIN ROBIN ROBIN

GbEth.

**Requester PCs, emulate HLT nodes**

3 KHz EB
3 ROLs per L2 request

100 kHz L1 rate
TCP/IP only

- L2 rate fast PC (3 GBE links)
- L2 rate default PC (2 GBE links)
- Bandwidth fast PC (3 GBE links)
- Bandwidth default PC (2 GBE links)

Event fragment size (words)

2 GBE links
TCP/IP only

2 "neighboring" cores

♦ CPU affinity masks

# of used CPU cores

**Best performance with 2 (of 8) cores
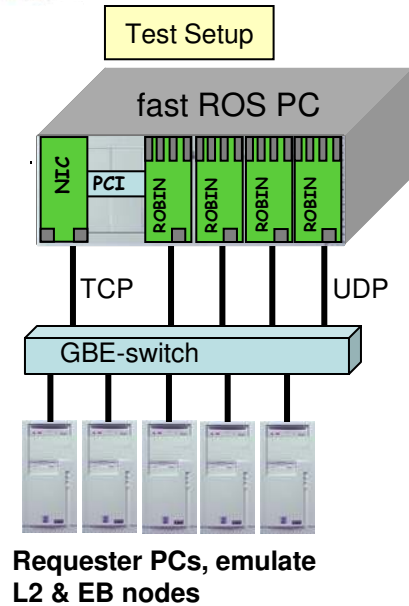-> effect of (extensive) use of mutexes?**

- L2 request rate (almost) fragment size independent
- L2 request rate increases by 50% to 150%
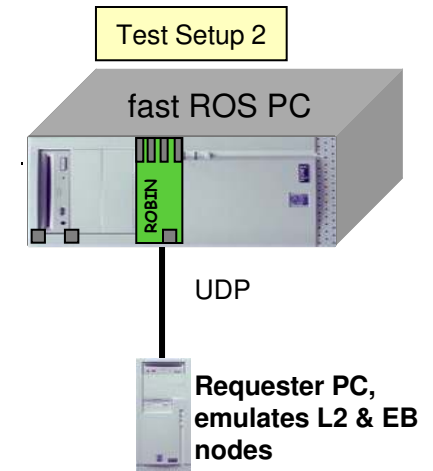- (expensive) ROBIN cards can be reused

# Impact of Read-Out via PC & ROBIN cards

Test Setup

fast ROS PC

NIC | PCI | ROBIN | ROBIN | ROBIN | ROBIN

TCP                    UDP

GBE-switch

**Requester PCs, emulate L2 & EB nodes**

Test Setup 2

fast ROS PC

ROBIN

UDP

**Requester PC, emulates L2 & EB nodes**

Full system test with (small HLT) farm
Only preliminary results so far

- Functionality OK
- current test system limits performance

Simplified set-up for tests at the ROBIN level

- This ROS configuration has the potential to deliver more performance than the ROS with the faster motherboard & CPU
- Further optimization of the system (software) required

# Summary and Conclusions

- Since its installation in 2006/2007 the ROS system has worked very reliably
- The ROS in its current configuration meets the requirements that were specified in the ATLAS Technical Design Report
- Several alternatives exist for the further improvement of the performance
    - More detailed tests have to be carried out in the deployed system to better understand the relative advantages and disadvantages of these alternatives
- The development of a PCIe based ROBIN has been started
    - Because motherboards with at least 4 64-bit PCI slots become difficult to find
    - Faster PPC CPU will also improve ROBIN performance
- Based on today's understanding of the ATLAS TDAQ (HLT rejection factor and algorithms) as well as the planned upgrades of ATLAS and LHC the current ROS architecture fulfills the requirements of phase 0 & 1 of ATLAS