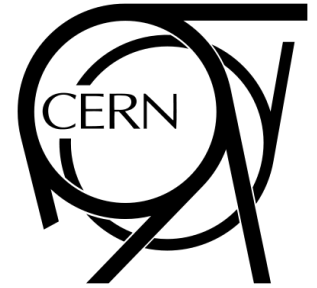




# Online Software



## InfiniBand enabled XDAQ clusters

### CMS DAQ Upgrade Studies



**17<sup>th</sup> November 2011 – CERN - SuperB ETD Meeting**

**Luciano Orsini, Andrea Petrucci – CERN (PH/CMD)**



# Topics

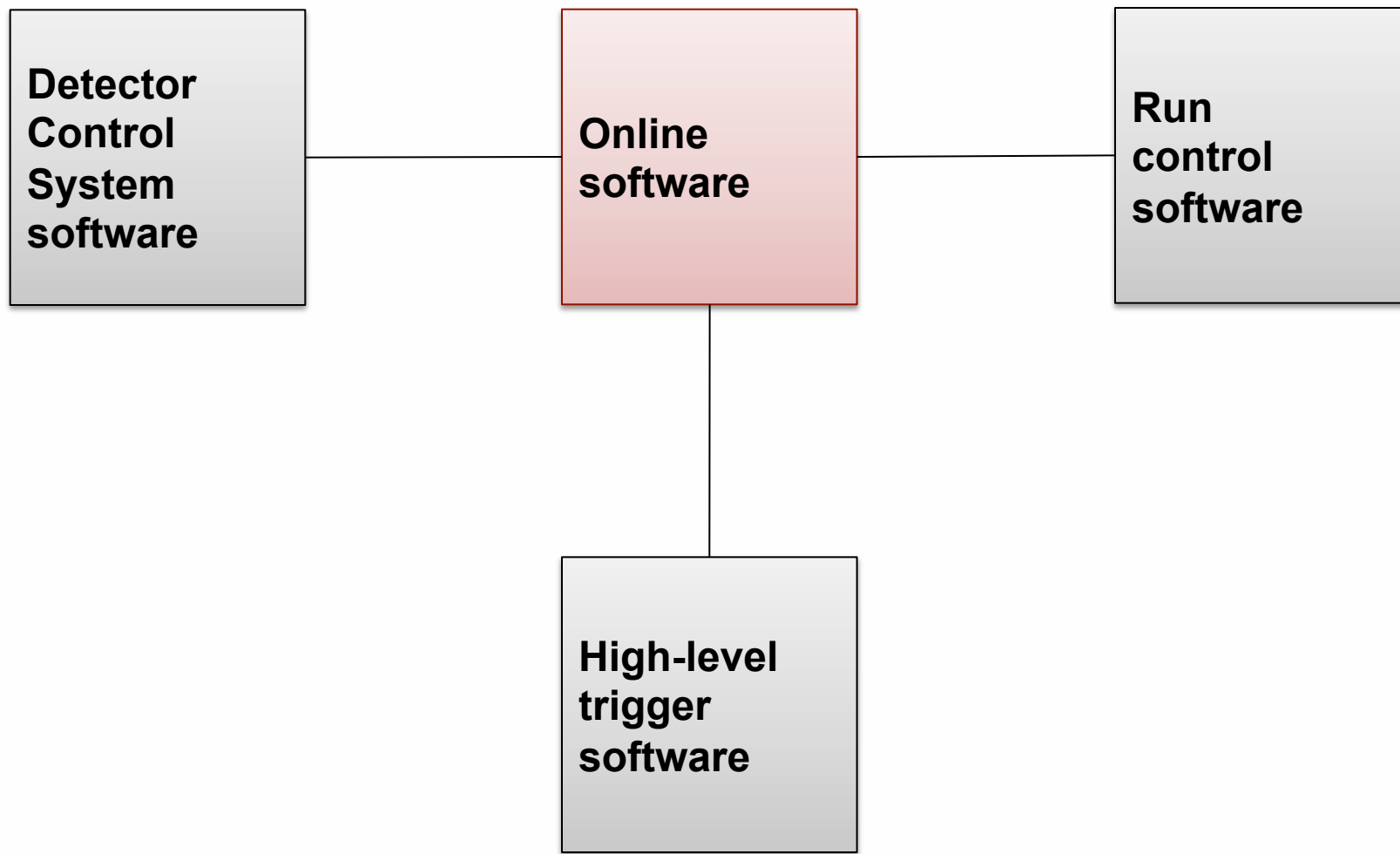
- **XDAQ** (CMS on-line software)
  - Scope
  - **Motivation** and **Requirements**
  - **Architecture**
  - Software Product Line
  - Outlook
  - Configuration Management
- **InfiniBand enabled XDAQ** clusters
  - InfiniBand software architecture
  - Protocol stack overview
  - **XDAQ peer transport** architecture overview
  - Early adopters
  - The **XDAQ uDAPL peer transport**
  - Testbed environment for I/O benchmarks
  - Preliminary measurements
- **Conclusion**



# XDAQ (CMS on-line software)



# Context Diagram





# Motivation

- CMS consists of a **set of sub-projects**
  - Similar to a coordinated set of small experiments
  - Many scenarios: central DAQ, subdetector DAQ, testbeams,
- **Geographically dispersed** participants
- **Autonomous developments**
- High **personnel turnover**
- High **performance** requirements
- Long **lifetime** and need to survive technology generations
- Similar **tasks** to be performed in each sub-detector



# Functional Requirements (TDR)

- **Communication and Interoperability**
  - Transparent use of **communication protocols**
  - Possibility to **add new protocols**
  - **Concurrent** use of **multiple protocols**
- **Device Access**
  - Access to **custom devices**
  - Hardware abstraction layer
- **Configuration, control and monitoring of applications**
  - Inspect and modify simple/complex **parameters**
  - **Allow coordination** of application components
  - **Record** structured **information**
    - **Uniform** logging, **error reporting**, **monitoring**
    - Interface to persistent data stores

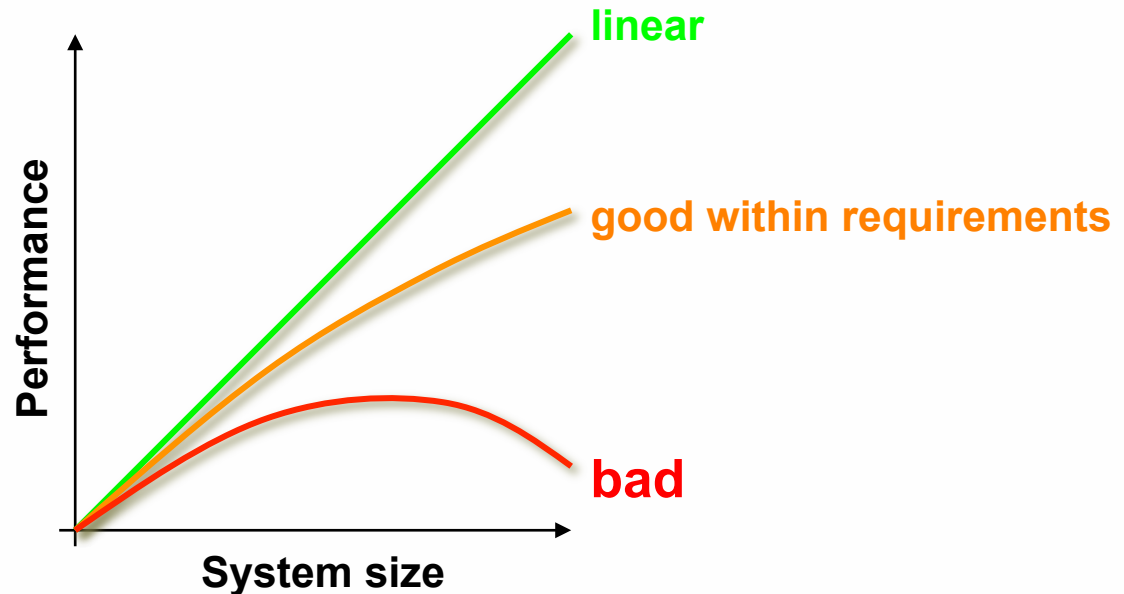


# Non-Functional Requirements (TDR)

- **Maintainability and Portability**
  - Portability across **operating system** and hardware **platforms**
  - Add **new electronics** without functional changes in user software
  - Application **code shall be invariant** with respect to the physical **location** and the **network**
  - Encourage working with **re-usable building blocks**

- **Scalability**

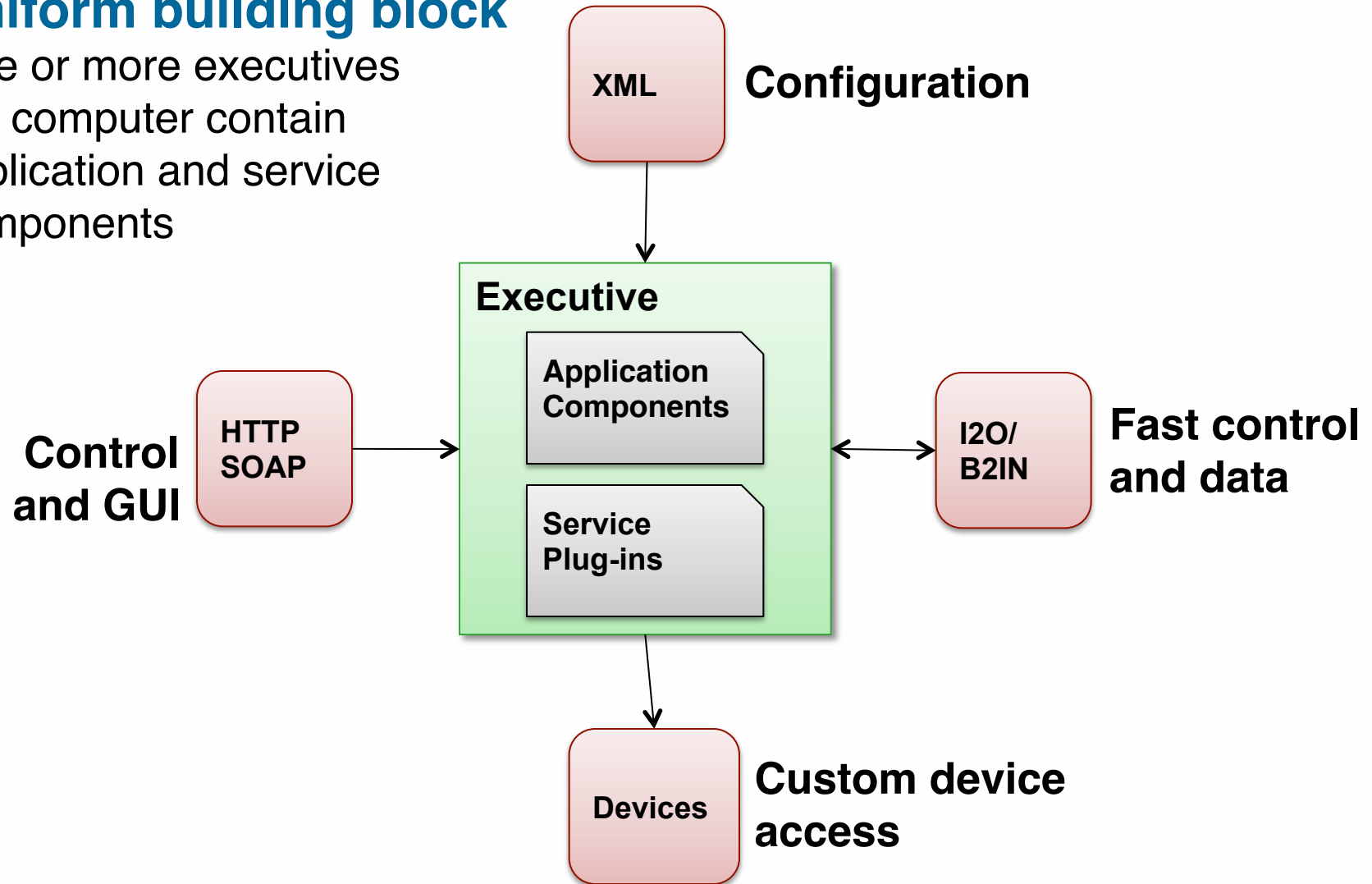
- Operate within requirements if size or volumes change
- Take advantage of **additional resource availability**
- **Overhead** introduced by the software environment must be **constant** for each transmission operation and small with respect to the underlying communication hardware in order not to introduce unpredictable behavior





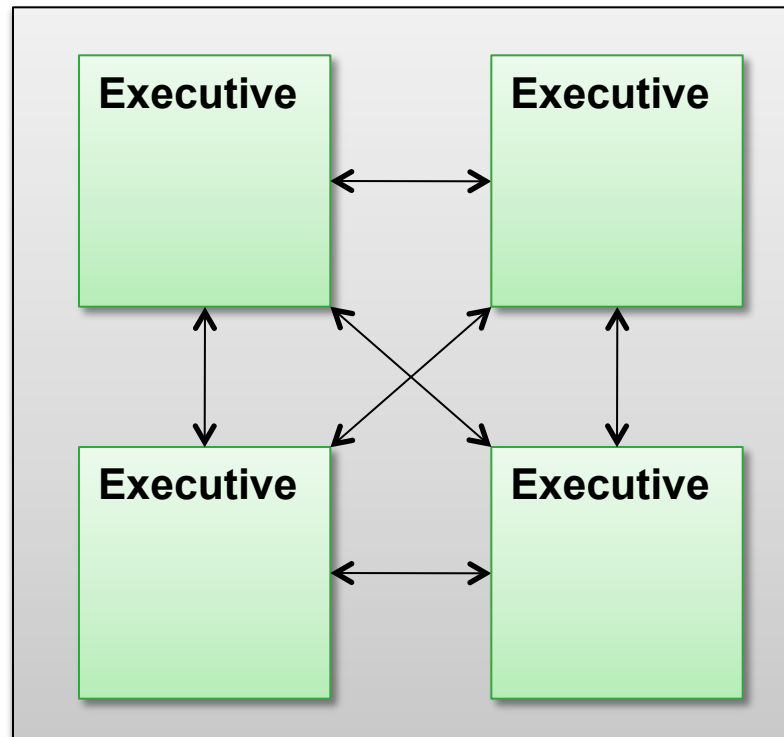
## Uniform building block

One or more executives per computer contain application and service components



## Replicated building blocks

Scalable cluster system architecture



## coretools



Core framework

## powerpack



Reusable applications

## worksuite



CMS specific applications



# Layered View

## Online Software

### Worksuite

Event Builders

Front-end Controllers

External System Interfaces

Detector Specific Applications

### Powerpack

Data Monitoring

Error and Alarming

Job Control

User Interfaces

### Coretools

OS Abstraction

Executive Framework

Hardware Access

Communication Subsystems

Configuration Management Support

### Platforms

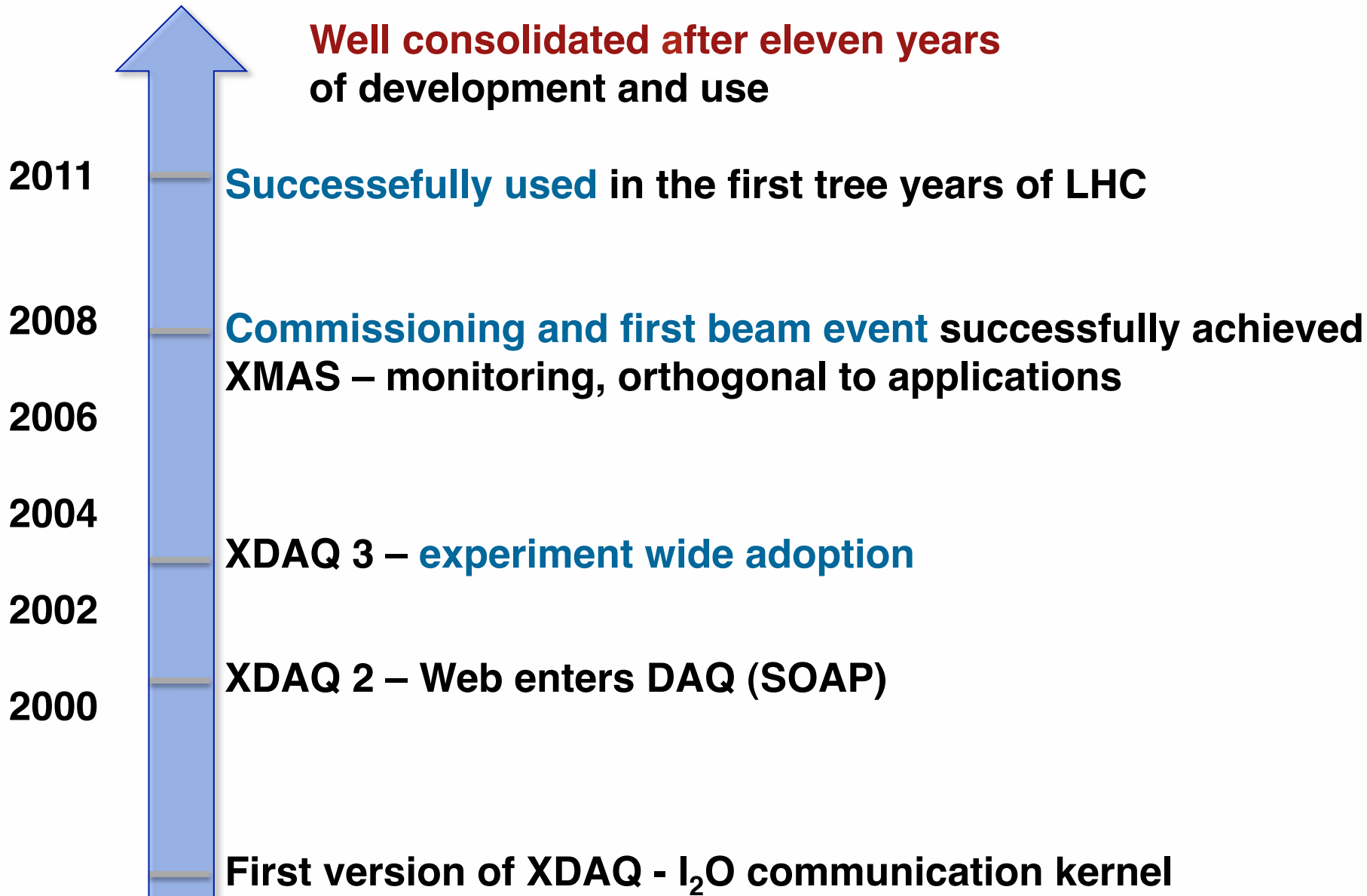
Operating Systems

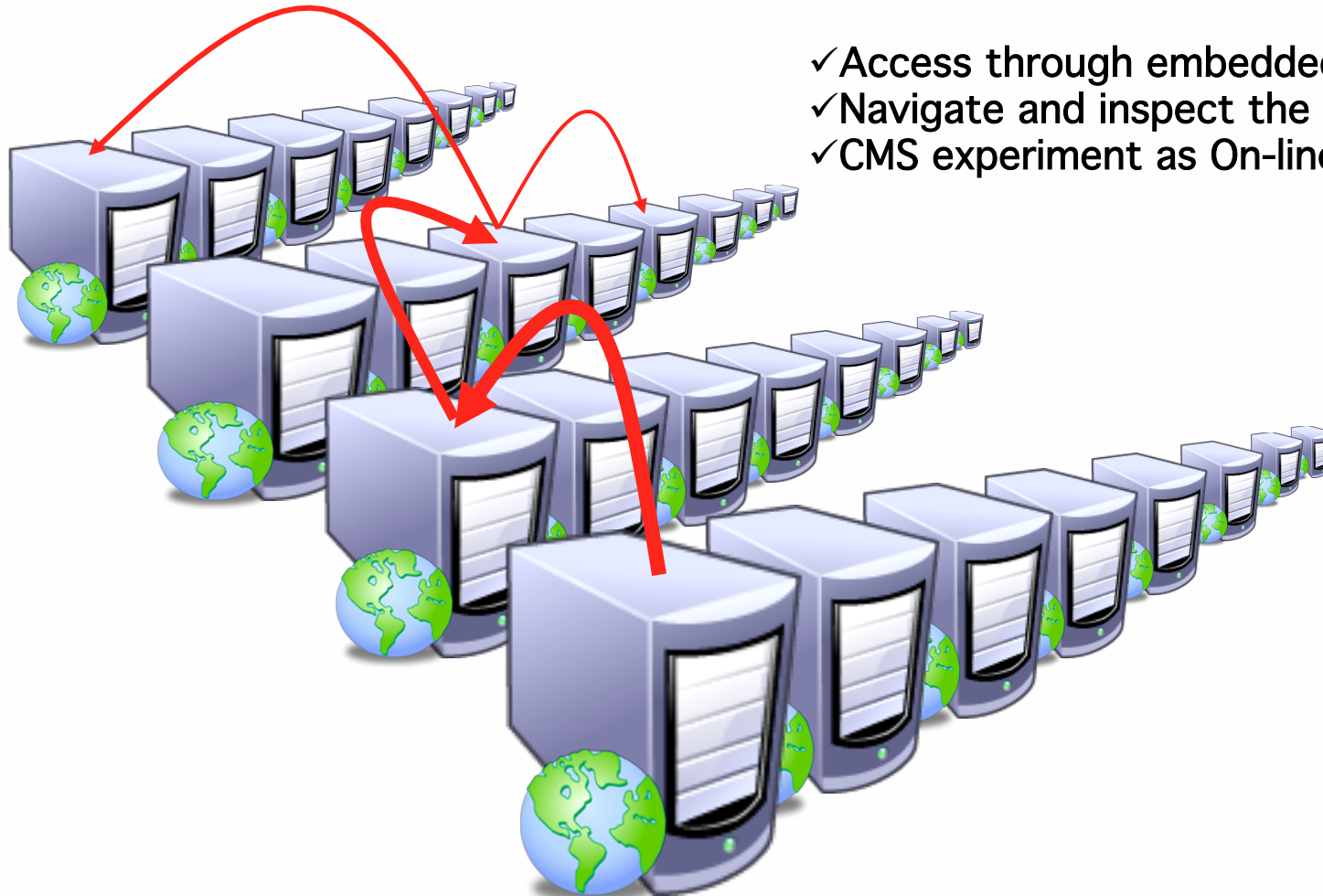
Networking Infrastructures

Hardware Device Interfaces



# Timeline





- ✓ Access through embedded HTTP
- ✓ Navigate and inspect the whole cluster
- ✓ CMS experiment as On-line Wide Web



# Outlook

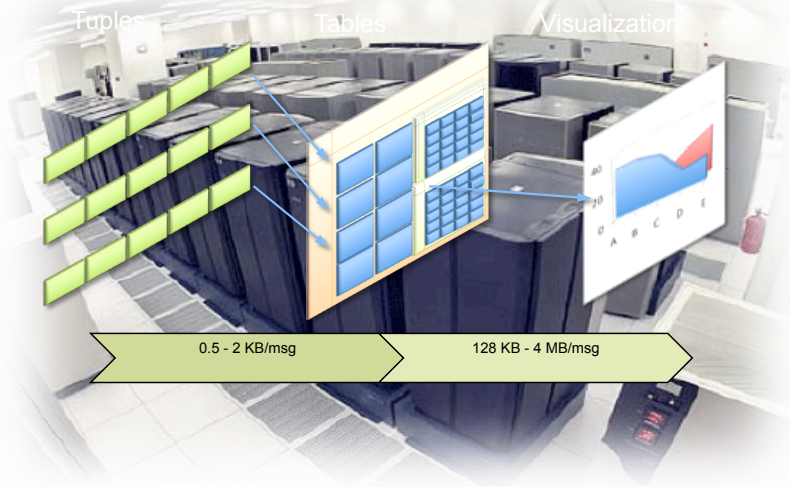
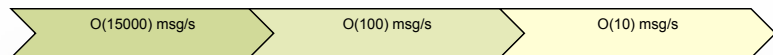
- Integrated web technologies
- Monitoring
- Errors and Alarms
- Reusable Event Builder

RU8 Version 3.4 Halted

Summary	Standard monitoring
throughput 0.000000e+00	class RU
average 0.000000e+00	instance 8
rate 0.000000e+00	hostname http://rubu9.cmsdaqpreseries:1972
rms 0.000000e+00	deltaT 1.010108e+00
	deltaN 0
	deltaSumOfSquares 0.000000e+00
	deltaSumOfSizes 0
	nbSuperFragments 0
	stateName Halted
	I2O_RU_DATA_READY_Payload 0
	I2O_RU_DATA_READY_LogicalCount0 0
	I2O_RU_DATA_READY_I2oCount 0
	I2O_RU_READOUT_Payload 0
	I2O_RU_READOUT_LogicalCount 0
	I2O_RU_READOUT_I2oCount 0

Standard configuration

- nbEvtIdsInBuilder 4096
- blockFIFOCapacity 16384



XDAQ online

Reset Up Re-arm Start 12/01/2008 9:22 End 12/01/2008 9:22 UTC Retrieve

- b32dev(310/0)
  - RU Builder Slice 1(309/0)
    - BUS(0/0)
    - RU(309/0)
    - EVM-0(0/0)
    - FED Builder(0/0)
    - JOBS(0/0)
    - Other(1/0)

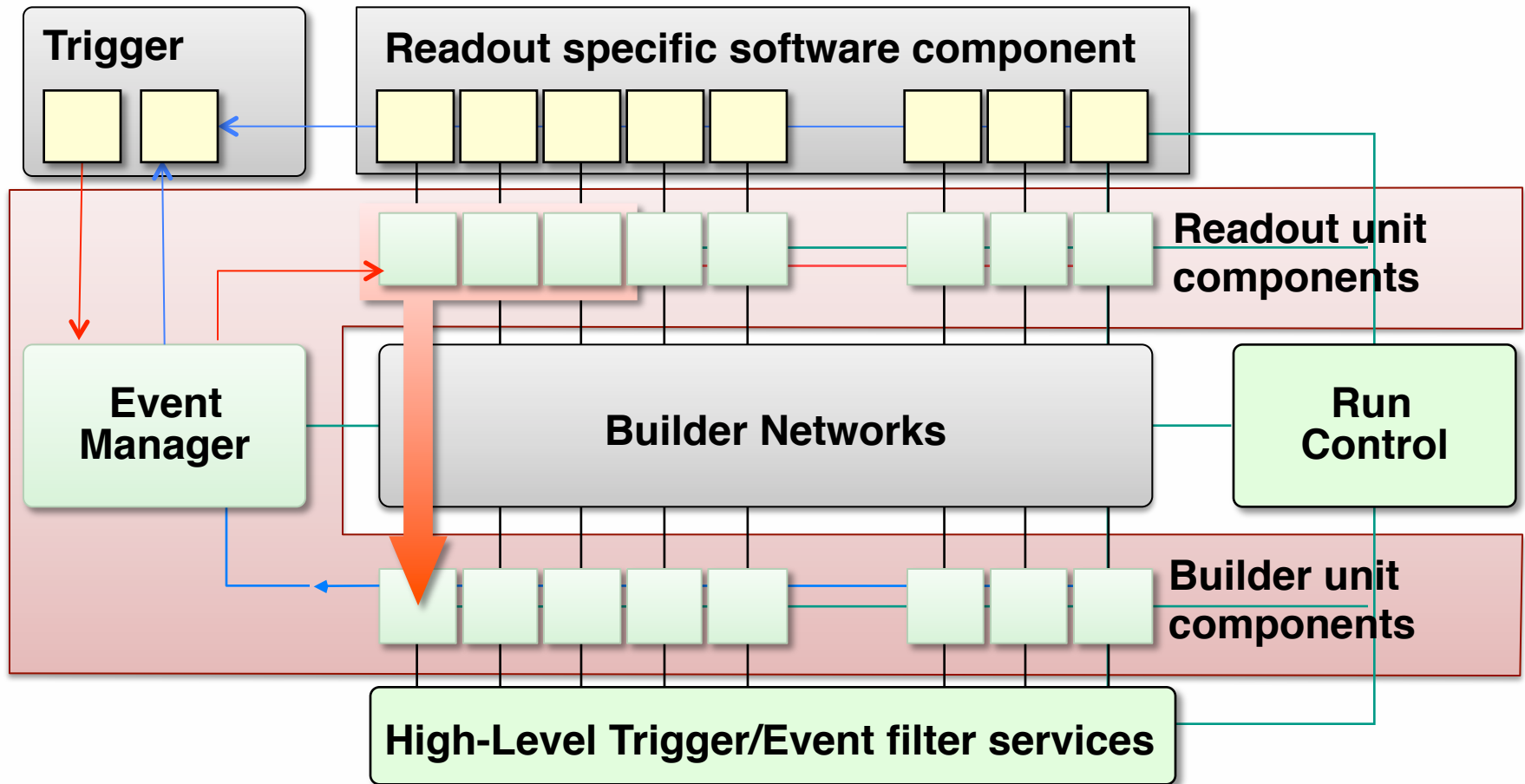
RU Builder Slice 1	FED Builder	JOBS
Active/Ack	Active/Ack	Active/Ack
Fatal 0/0	Fatal 0/0	Fatal 0/0
Error 309/0	Error 0/0	Error 0/0
Warn 0/0	Warn 0/0	Warn 0/0

Other
Active/Ack
Fatal 1/0
Error 0/0
Warn 0/0



# Generic Event Builder

- Configurable in size and network technology
- Customization at boundaries through pluggable components





- Identification
  - definition of packages
  - versioning
- Traceability
  - Status accounting
  - Documented change history
  - Tickets
- Planning
  - Milestones
  - Priorities, People
- Release
  - Source and binary releases
  - Parallel releases
  - Upgrades





# Configuration Management Tools



IT Department



WEB Services

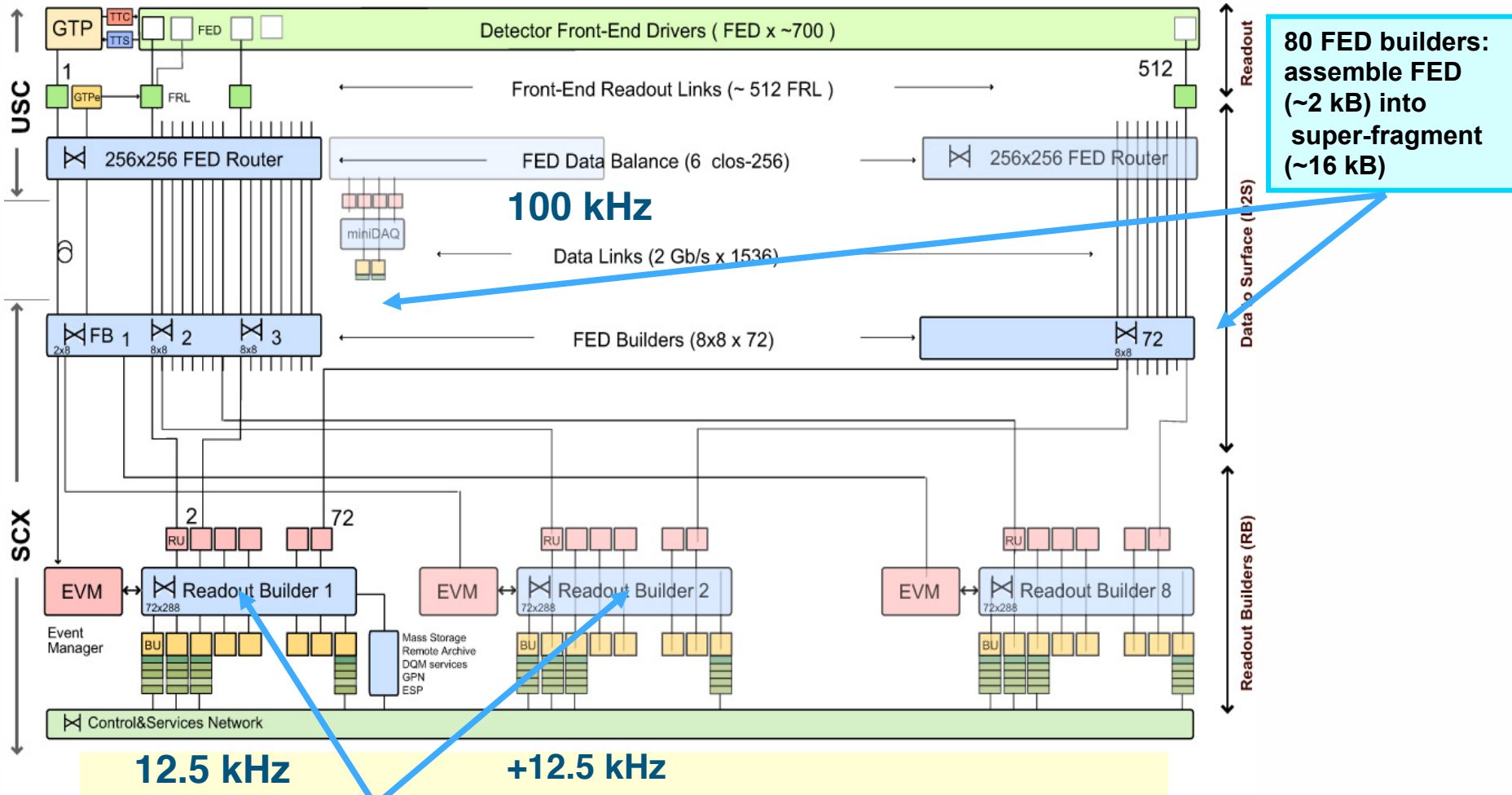


E-Groups





# CMS DAQ System





# CMS DAQ 2011

- Full readout of all sub-detectors
- 64 FedBuilders, Myrinet (~1000 applications)
- 8 DAQ Slices, 100 GB/s event builder, Force10 (~1600 applications)
  - Slice with 80 RUs x 125 BUFUs – TCP/IP over GE
- Event Filter (~8200 instances):
  - 720 8-core PCs (Intel Xeon E5430- 2.66 GHz, 16 GB):
    - ~5000 instances of CMSSW-HLT
    - ~100ms/evt
  - 288 8-core PCs (Intel Xeon X5650 2.67 GHz, 24 GB)
    - ~3200 instances of CMSSW-HLT
    - ~81ms/evt
- Storage Manager (16 instances): 2 GB/s, 250 TB buffer.
- CMS on-line farm ~2500 PCs



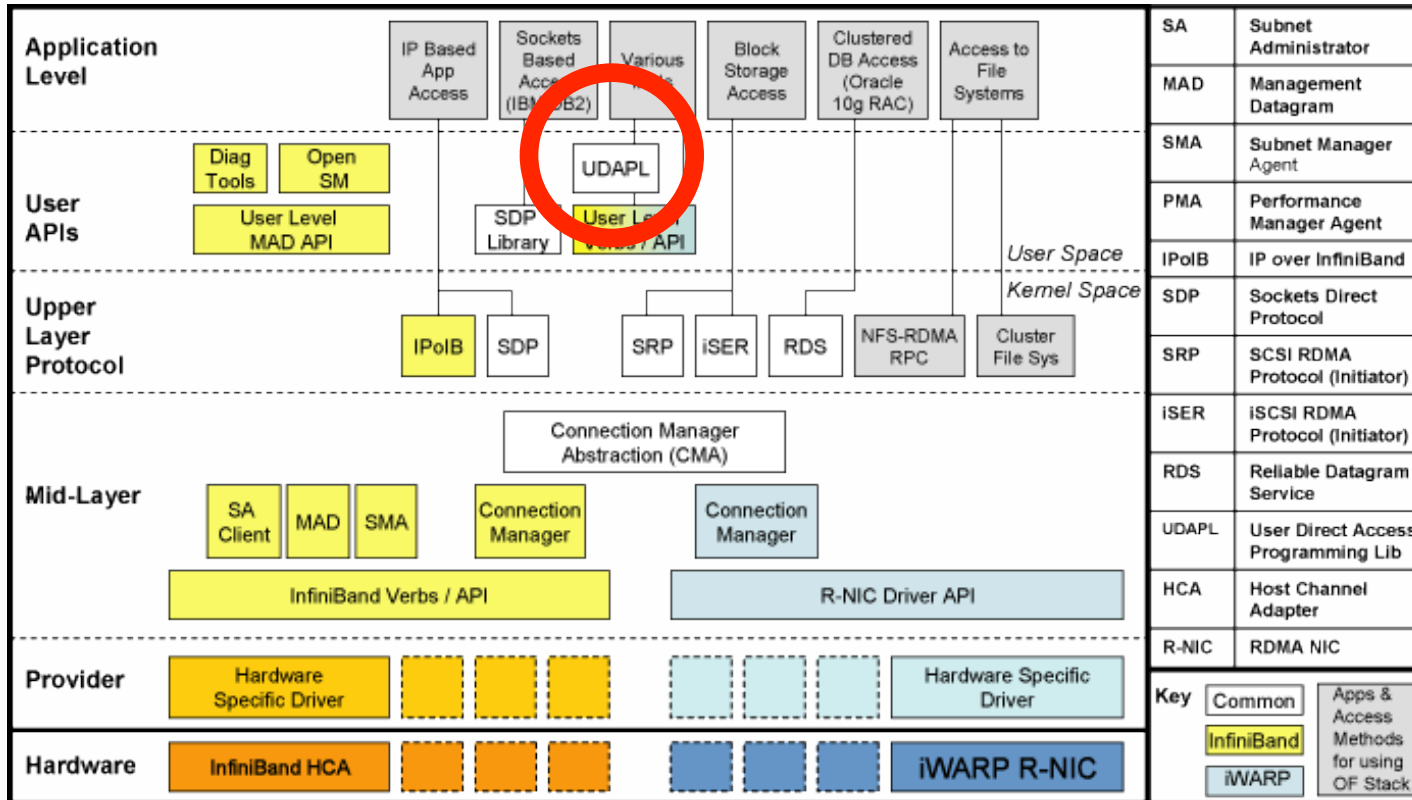
# InfiniBand enabled XDAQ clusters



# The OFED Stack (source: OpenFabrics Alliance)

A unified, cross-platform, transport-independent software stack for RDMA and kernel bypass

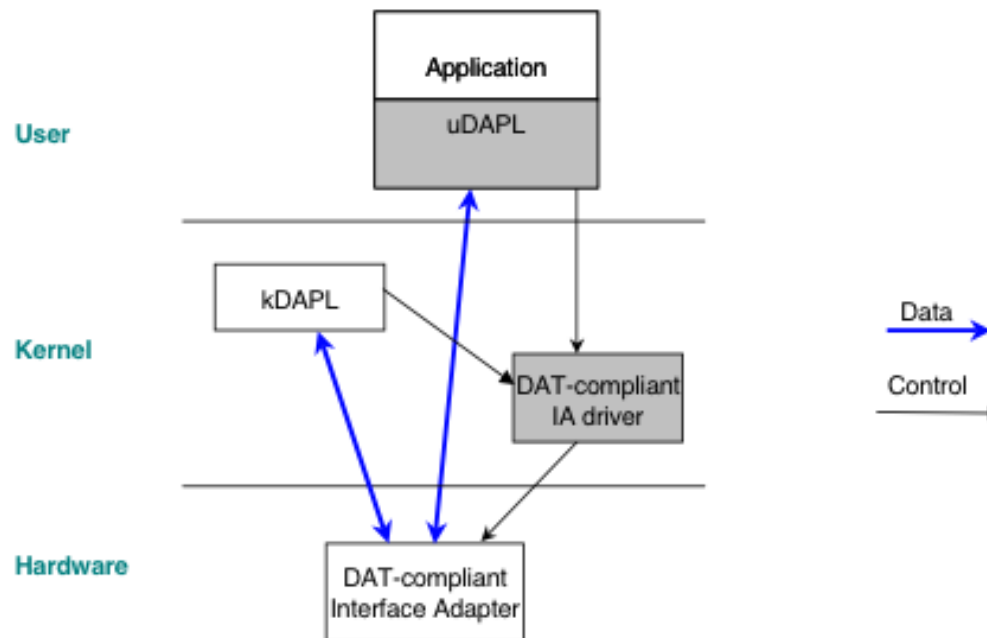
□ <http://www.openfabrics.org/>





# DAT Model (source DAT Collaborative)

- Developed by DAT collaborative
  - ❑ <http://www.datcollaborative.org/>
- Transport and platform (OS) independent
- Define user (**uDAPL**) and kernel (kDAPL) APIs
- DAT supports **reliable** connection
- Data Transfer Operations **send, receive**, rdma\_read, rdma\_write
- uDAPL Version 2.x, January, 2007





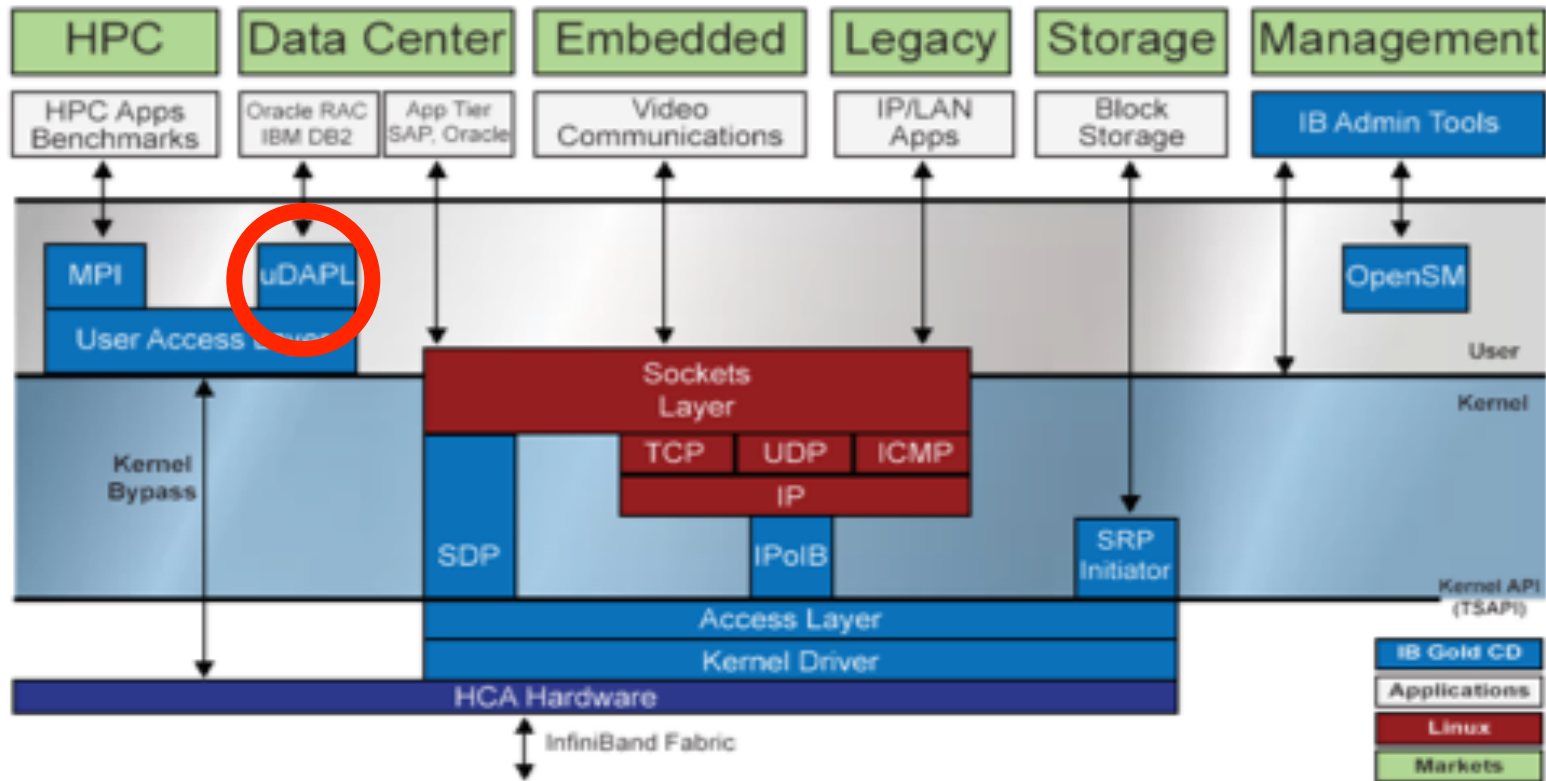
# Other APIs

- IPoIB (IP over Infiniband)
- MPI (Message Passing Interface)
- SDP (Socket Direct Protocol)
- RDS (Reliable Datagram Sockets)



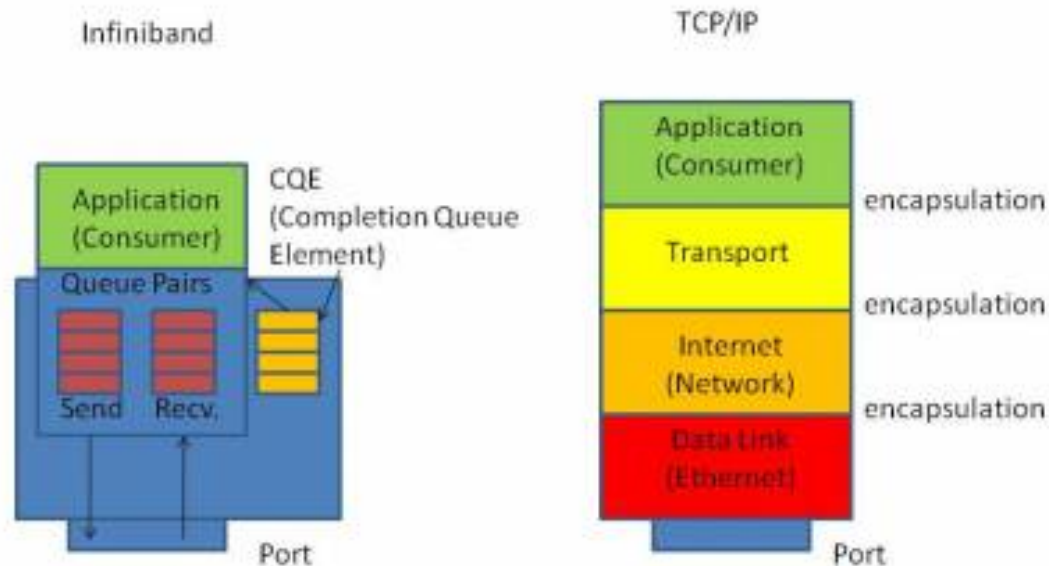


# Protocol Stack View



# Comparison of the Stacks

The protocol is defined as a very thin set of **zero copy** functions when compared to thicker protocol implementations such as TCP/IP



- I<sub>2</sub>O messaging on InfiniBand transparent to XDAQ applications



31	3	24, 23	2	16, 15	1	8, 7	0	0
MessageSize				MessageFlags		VersionOffset		
Function (FFH)		InitiatorAddress			TargetAddress			
InitiatorContext								
TransactionContext								
OrganizationID				XFunctionCode				
PrivatePayload = function parameters								
PrivatePayload								



# XDAQ IB Early Adopters

- **XDAQ evaluation at GSI** for FAIR (Facility for Antiproton and Ion Research) by J. **Adamczewski-Musch**, **H.G. Essel**, **S. Linev** at GSI Helmholtzzentrum für Schwerionenforschung GmbH, Darmstadt
  - ❑ <https://www.gsi.de/documents/DOC-2006-Oct-34-1.pdf>
- Implemented **very first prototype** for XDAQ peer transport ptDAPL for **I<sub>2</sub>O**
- **Many Thanks** to GSI people
  - Their **interest** in XDAQ technology
  - Extremely useful **contribution** for enabling XDAQ with Infiniband
  - **Shorten** learning curve for new technology
  - **Reference** measurements for **comparison**



# LHCb IB Collaboration

- Infiniband Event-Builder Architecture Test-beds for Full Rate Data Acquisition in LHCb – CHEP 2010
  - ❑ <http://cdsweb.cern.ch/record/1302037/files/LHCb-TALK-2010-151.pdf>
- **Many thanks** to LHCb people (Niko Neufeld, Jean Christophe Garnier)
  - Provide access to their network cluster with IB installation ready to use
  - Reference measurements for comparison



# LHCb IB Environment (2010)

- 8 nodes
- Qlogic 12300
  - QDR 4X
  - 32 Gb/s
- Processor type Intel Xeon E5520
- Processors x cores x clock (GHz) 2 x 4 x 2.27
- RAM (GiB) 3
- HCA qle7340 4x QDR
- SLC5/64 bit - Kernel 2.6.18
- CMSOS release 11
  - GEVB (Generic Event Builder)
  - PTuDAPL (XDAQ Peer Transport uDAPL)



- Reverse engineering of GSI ptDAPL
  - Identified overhead in buffer management (**virtual mapping** for all packets in use for all send/recv operations)
  - Implementation based on non blocking TCP peer transport not effective for use of uDAPL API
  - Not optimized workloop usage
  - Based on DAT Spec 1.x
- Full re-factoring into new ptuDAPL
  - Use of smart memory pool based on uDAPL memory region allocator (random access to memory with no intermediate management by using **cookies**)
  - Profiting for inherent **non blocking and queuing of uDAPL** API for minimizing latency (removed intermediate output fifo)
  - Several optimizations in workloops, wrapper classes etc.
  - All I/O operations centered on dedicated uDAPL memory pool
  - Based on DAT Spec 2.x

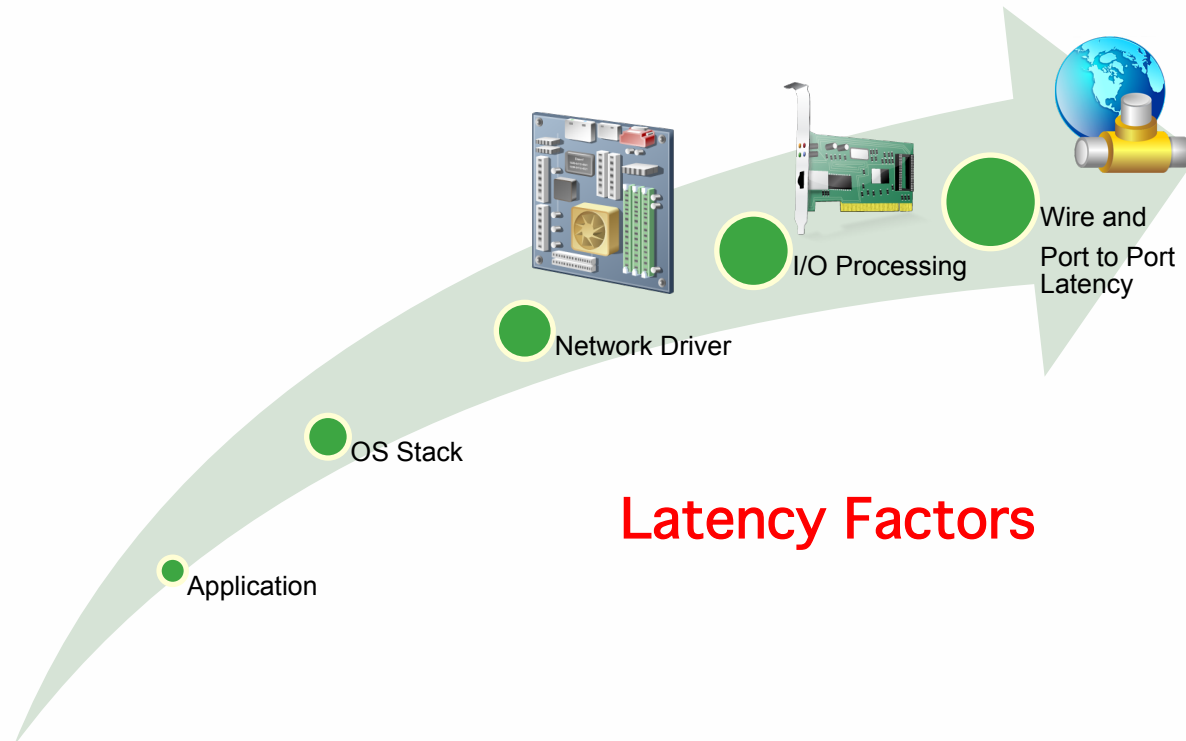


# Reference XDAQ Applications

- **Roundtrip**
  - Used to measure latency
- **MStreamIO**
  - Used to measure throughput
- **Event Builders**
  - **GEVB** Generic Event Builder
  - **RUBuilder** Official CMS event builder



- Simple XDAQ application to compute the One-way delay
- Time packet to travel from a specific source to a specific destination and back again
- One-way latency is measured by timing a round-trip message and dividing the obtained result by two

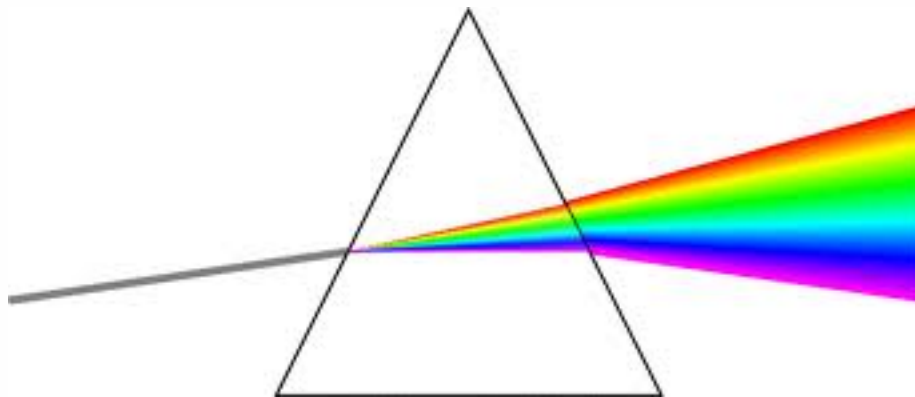


## Latency Factors

- **Unidirectional** throughput (bandwidth) is measured using a unidirectional send of N messages. Time sampling is done at the receiver side and starts with the first incoming message
- **Bidirectional** throughput measures the case where both of the two hosts are busy sending data to each other. Note that the theoretical maximum here is twice the link speed for full duplex media



- **RUBuilder**
  - Currently used in central DAQ
  - Dataflow protocol **5 hops** (allocate, confirm, send, cache and release)
- **GEVB**
  - Same functionalities and interfaces as RUBuilder
  - Dataflow control **3 hops** (allocate, ship and cache)
  - It was used for sub-detector data acquisition testbeams before TDR and adoption of RUBuilder





# Preliminary Measurements



# Test Planning

## UDAPL

- **RTT, I/O**
- **Throughput** 1x1, 1x2, 1x4
- **GEVB** 1x1, 1x2, 1x4, 2x1, 4x1, 2x2, 3x3
- **RUBuilder** 1x1, 1x2, 1x4, 2x1, 4x1, 2x2, 3x3

## ATCP(10GE)

- **RTT**
- **Throughput** 1x1, 1x2, 1x4
- **GEVB** 1x1, 1x2, 1x4, 2x1, 4x1, 2x2, 3x3
- **RUBuilder** 1x1, 1x2, 1x4, 2x1, 4x1, 2x2, 3x3

## IPoIB

- **RTT,I/O**
- **Throughput** 1x1, 1x2, 1x4
- **GEVB** 1x1, 1x2, 1x4, 2x1, 4x1, 2x2, 3x3
- **RUBuilder** 1x1, 1x2, 1x4, 2x1, 4x1, 2x2, 3x3

## Variable sizes

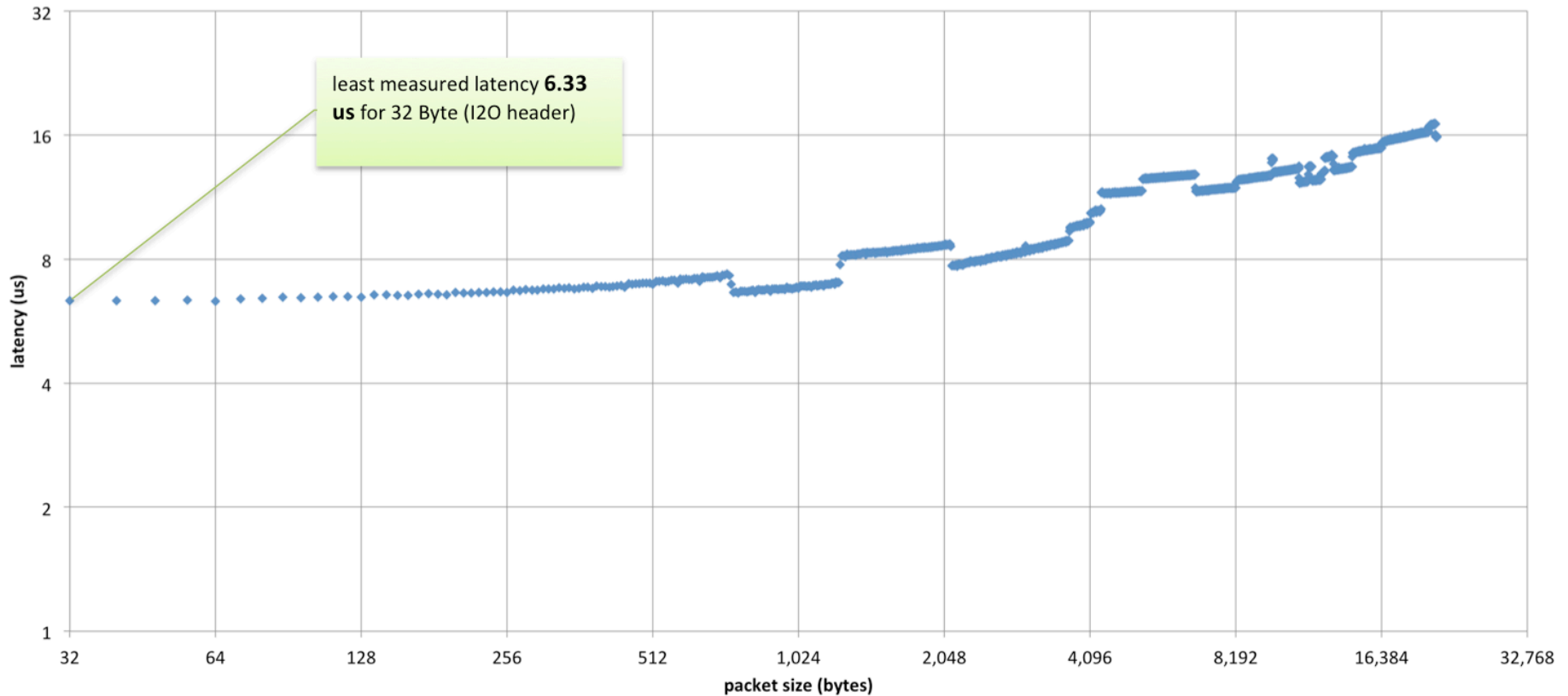
- **GEVB** both uDAPL/IPoIB
- **RUBuilder** uDAPL/IPoIB



# RTT

## Roundtrip latency

XDAQ uDAPL peer transport preliminary measurements on QLogic 4x QDR - CMS DAQ upgrade studies  
September 2011 by L.Orsini, A. Petrucci



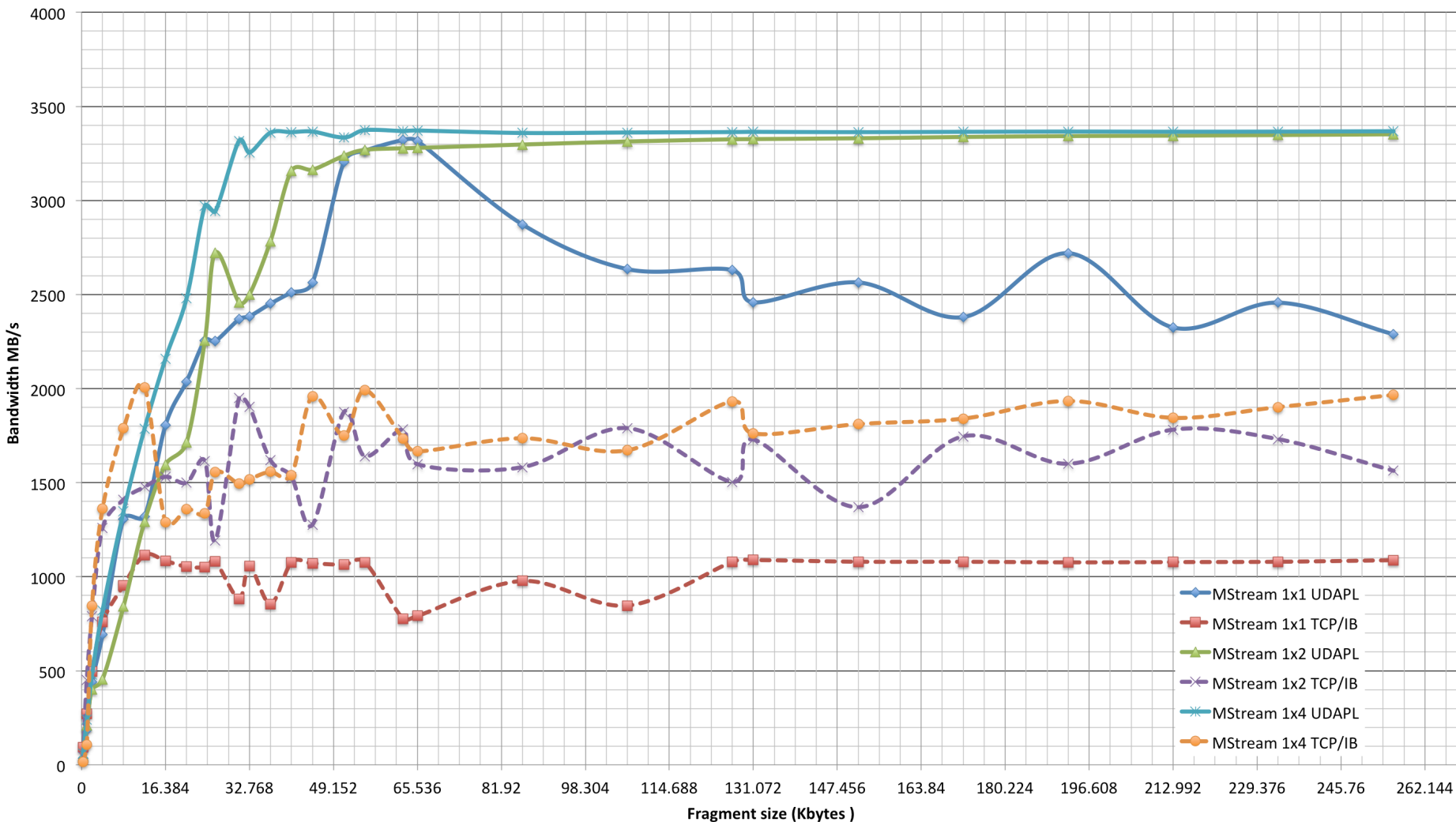


# Throughput

## Stream I/O 1 to N (UDAPL vs TCP/IB)

XDAQ uDAPL peer transport preliminary measurements on QLogic 4x QDR - CMS DAQ upgrade studies

September 2011 by L.Orsini, A. Petrucci

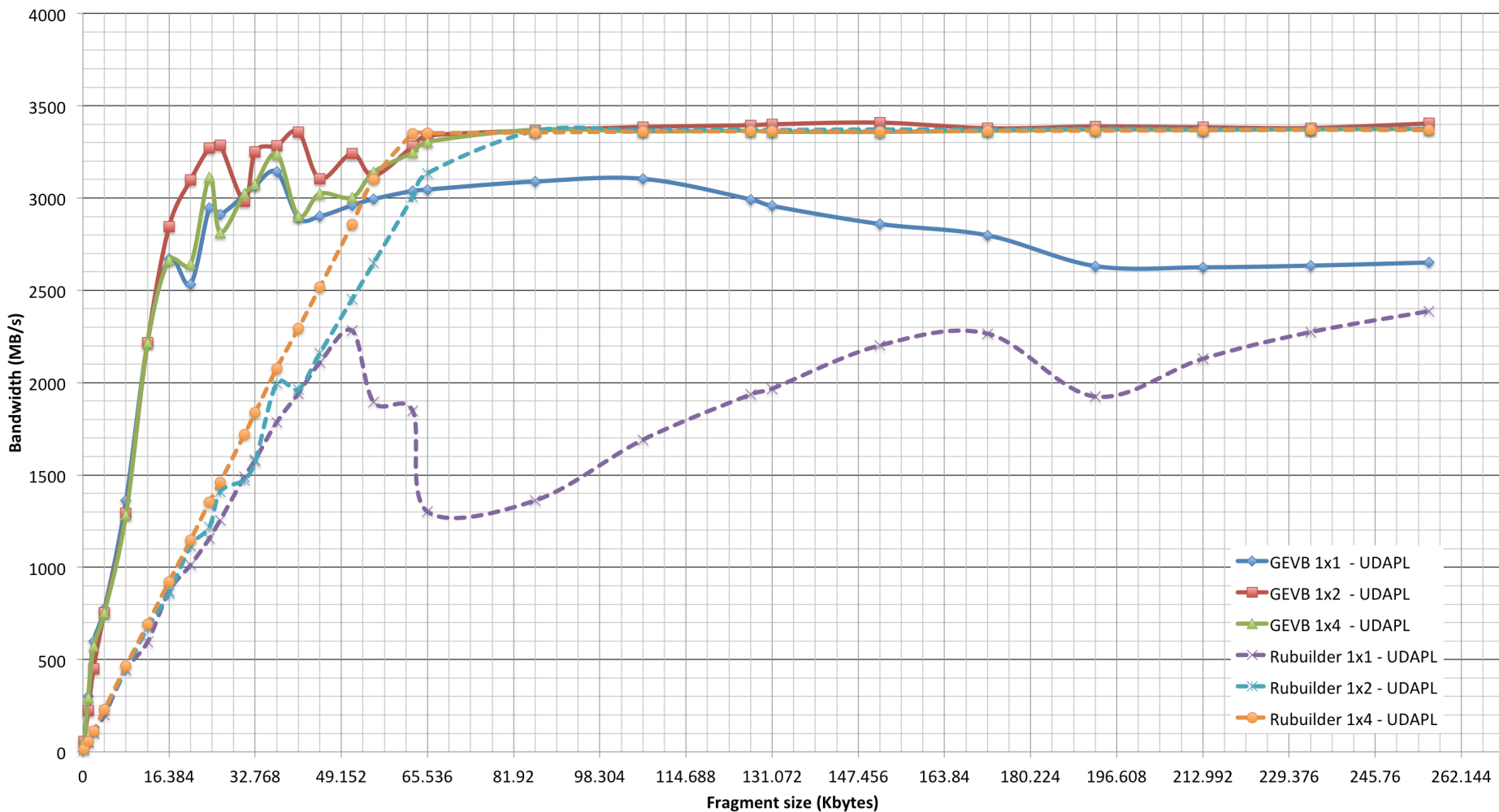




# EVB-1xN ptuDAPL

## Bandwidth 1xN - GEVB vs Rubuilder (UDAPL)

XDAQ uDAPL peer transport preliminary measurements on QLogic 4x QDR - CMS DAQ upgrade studies  
September 2011 by L. Orsini, A. Petrucci





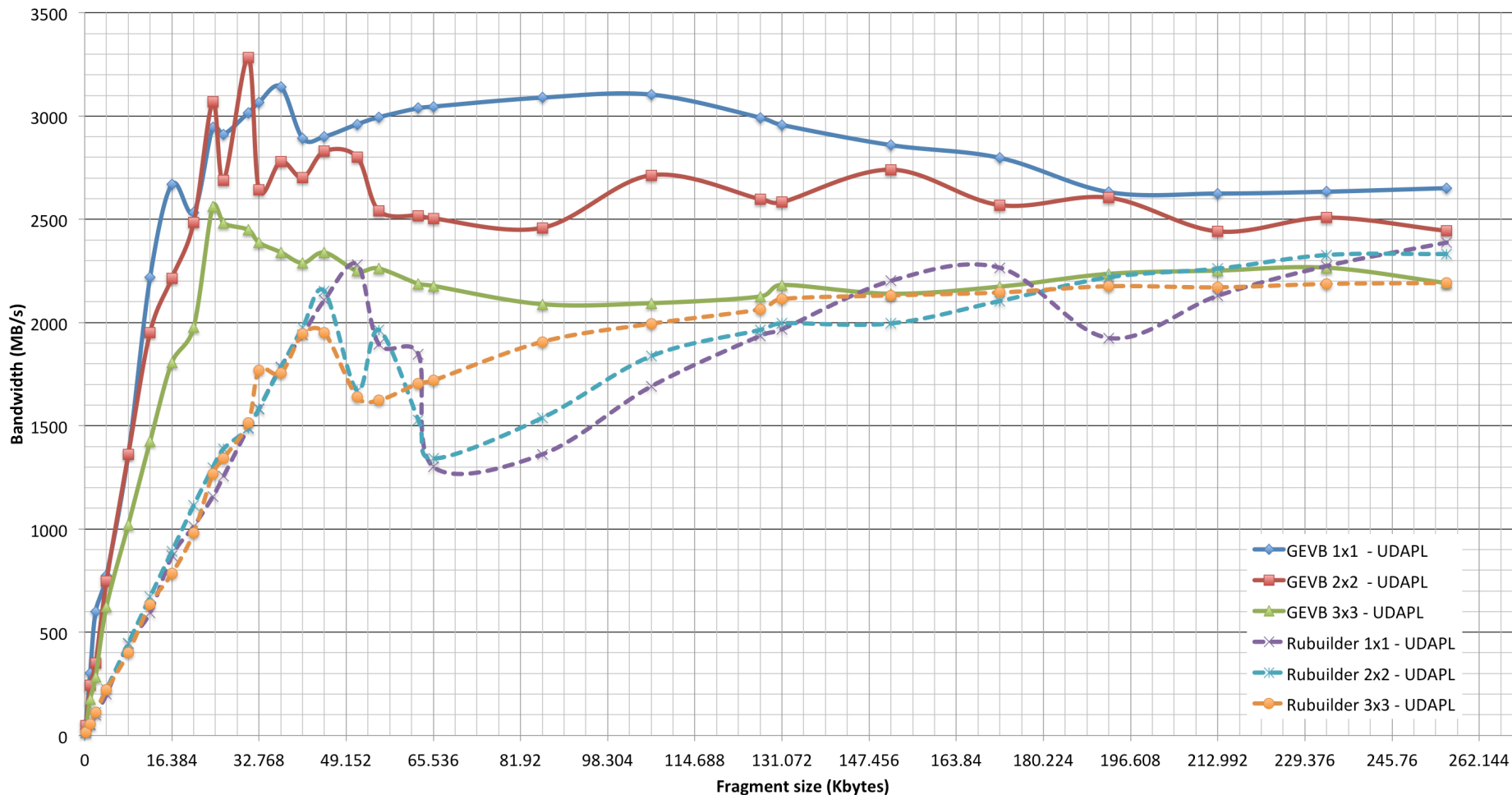


# EVB-NxN ptuDAPL

## Bandwidth NxN - GEVB vs Rubuilder (UDAPL)

XDAQ uDAPL peer transport preliminary measurements on QLogic 4x QDR - CMS DAQ upgrade studies

September 2011 by L.Orsini, A. Petrucci

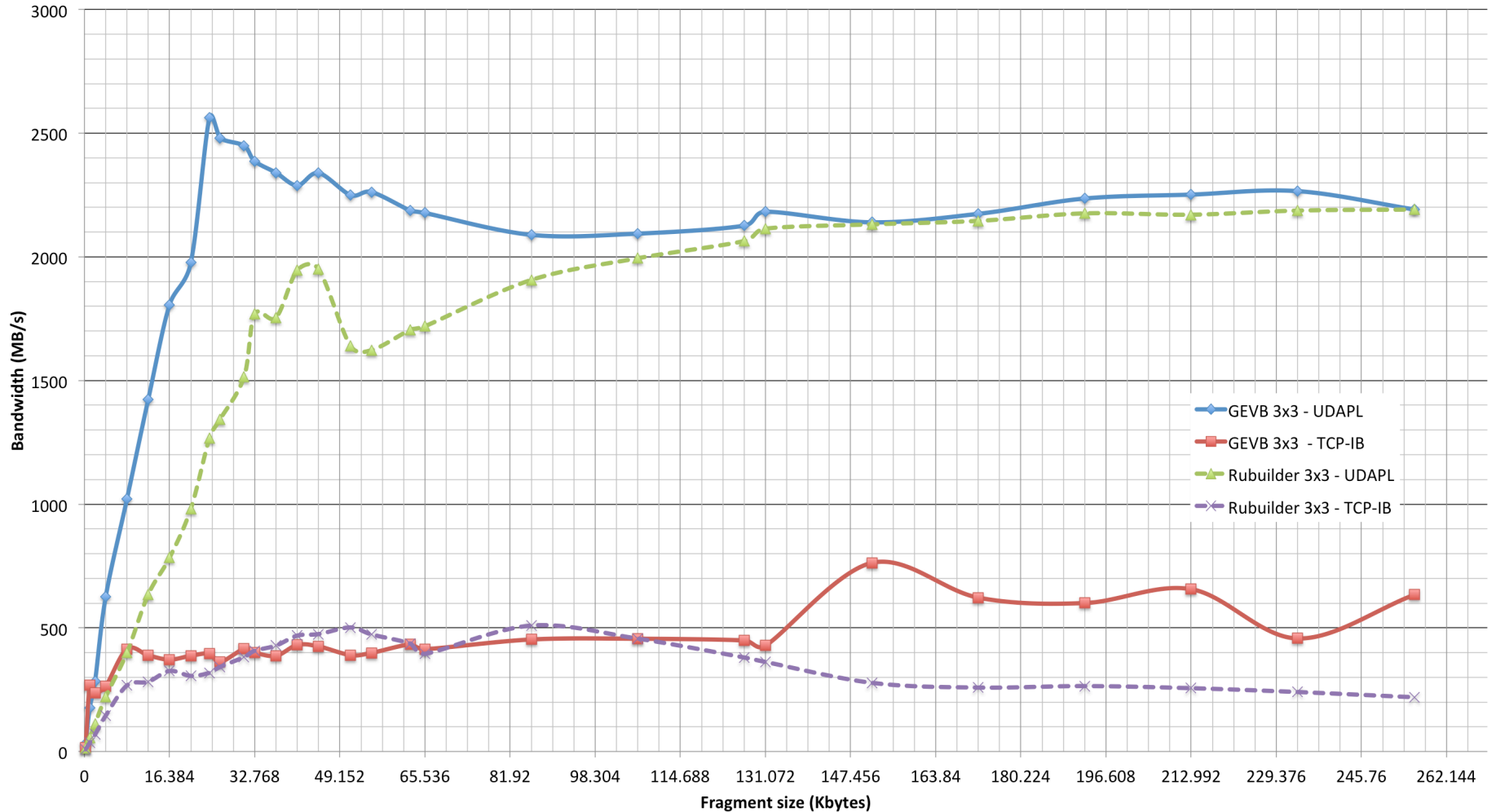




# EVB-3x3 ptuDAPL vs ptaTCP

## Bandwidth 3x3 - GEVB vs Rubuilder (UDAPL and TCP/IB)

XDAQ uDAPL peer transport preliminary measurements on QLogic 4x QDR - CMS DAQ upgrade studies  
September 2011 by L.Orsini, A. Petrucci





# Preliminary results

- Use rubuilder **off-the-shelf** with both peer transport
- Use default configuration for both DAPL and TCP
- Latency pt/**u**DAPL (RTT) 6.33 us
- pt/**u**DAPL **over 3 GB/S** with packet size greater than **26 Kbytes**
- pt/atcp reached **2.4 GB/s** with packet size greater than **172 Kbytes**
- Largest size of event builder **3x3** (6 nodes + 1 node for event manager)
  - maximum bandwidth **2.5 GB/s at 24 Kbytes with pt/uDAPL**
  - maximum bandwidth **762 MB/s at 150 Kbytes with pt/atcp**



# Conclusion

- Building a DAQ system as a process of **assembly re-usable components** in a predetermined way **rather than** a **programming** task.
- Achieved **a uniform DAQ product-line** for all CMS data acquisition application scenarios ranging from single CPU setups to the final systems comprising thousands of nodes.

*<http://svnweb.cern.ch/trac/cmsos>*



# Backup

## I/O 1 to 1 (MBs)

XDAQ uDAPL peer transport preliminary measurements on QLogic 4x QDR - CMS DAQ upgrade studies - September 2011 by L.Orsini, A. Petrucci



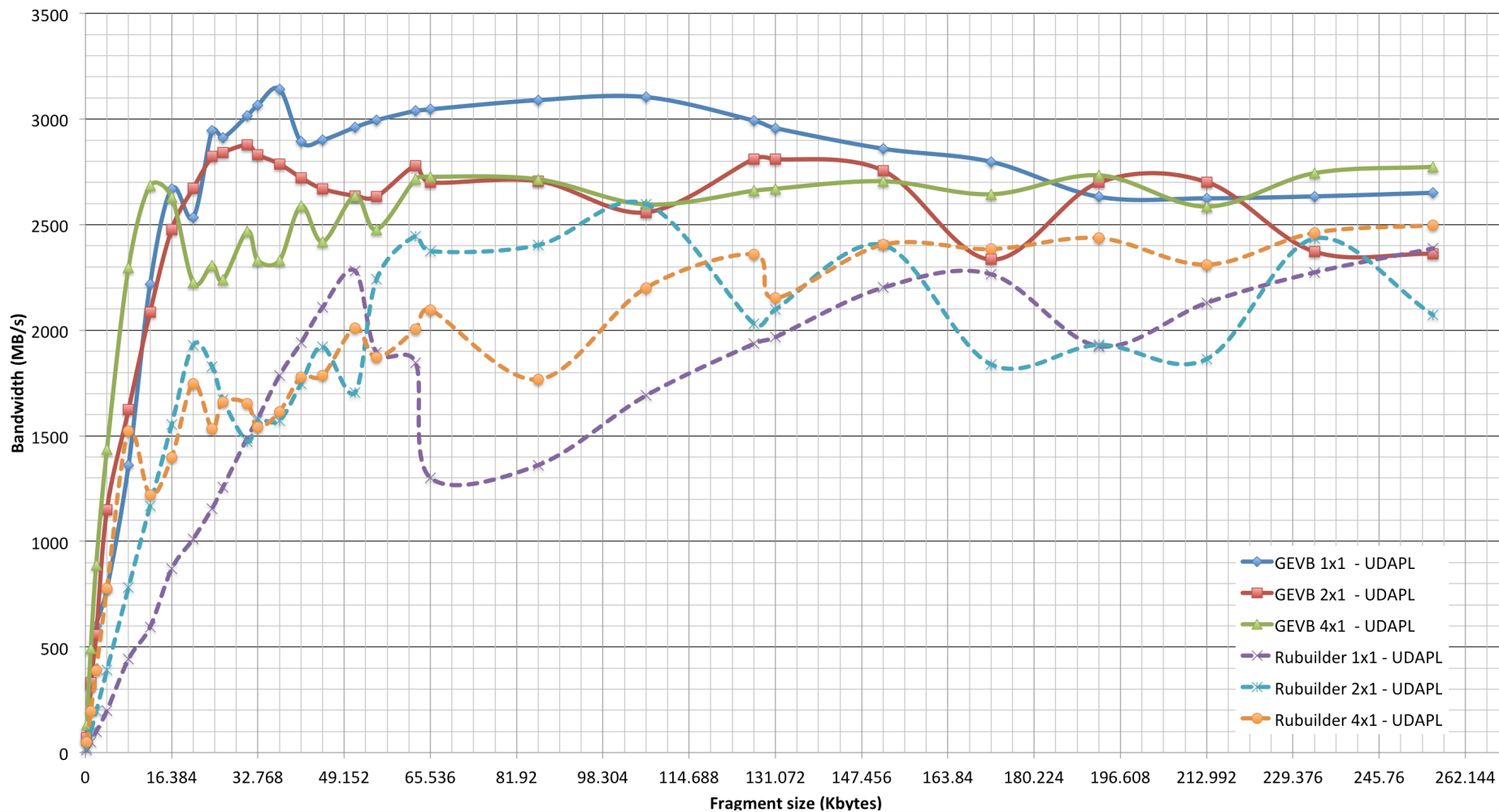


# EVB-Nx1 ptuDAPL

## Bandwidth Nx1 - GEVB vs Rubuilder (UDAPL)

XDAQ uDAPL peer transport preliminary measurements on QLogic 4x QDR - CMS DAQ upgrade studies

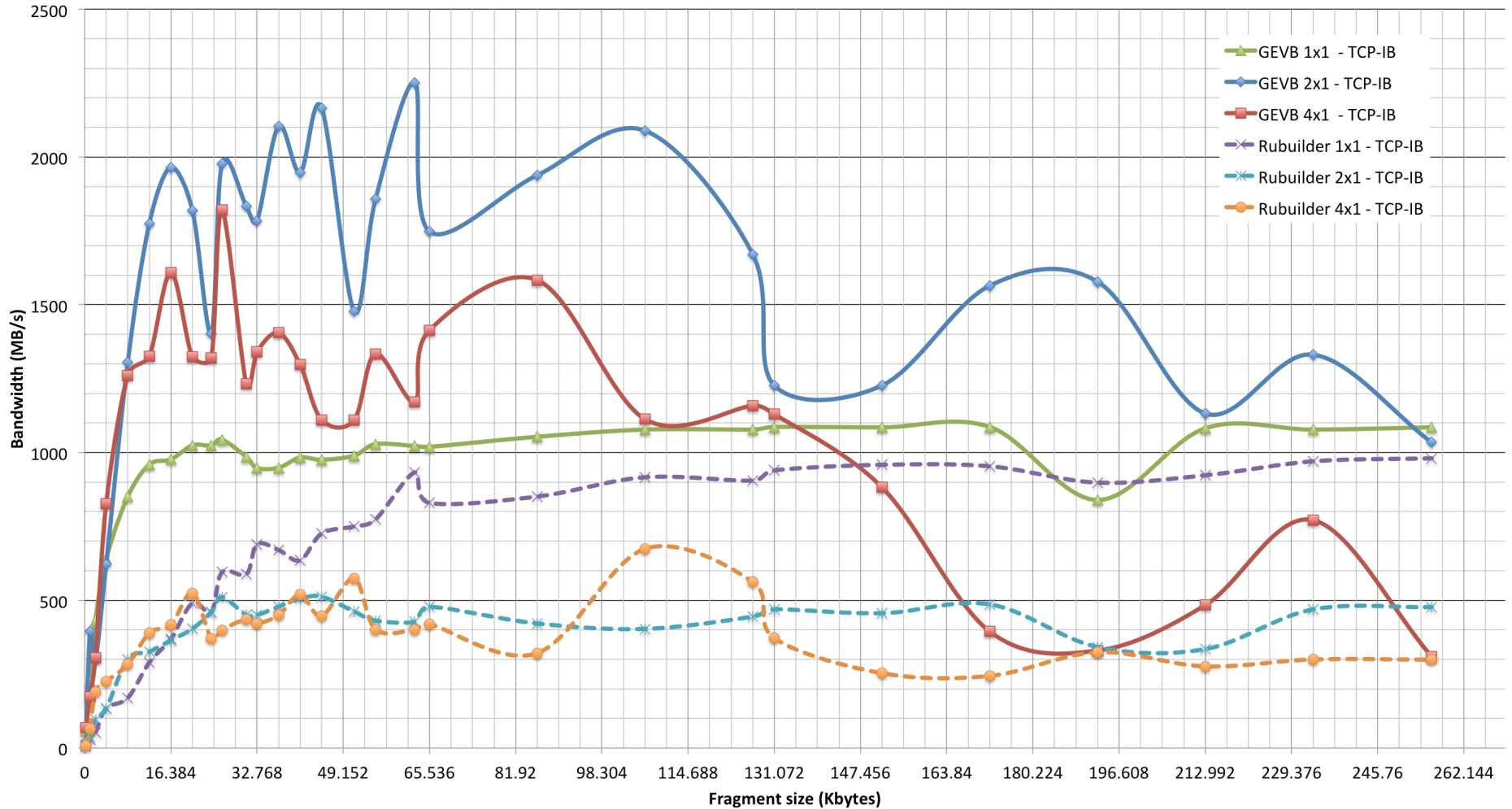
September 2011 by L.Orsini, A. Petrucci





# EVB-Nx1 ptaTCP

Bandwidth NxN - GEVB vs Rubuilder (TCP/IB)  
XDAQ uDAPL peer transport preliminary measurements on QLogic 4x QDR - CMS DAQ upgrade studies  
September 2011 by L.Orsini, A. Petrucci



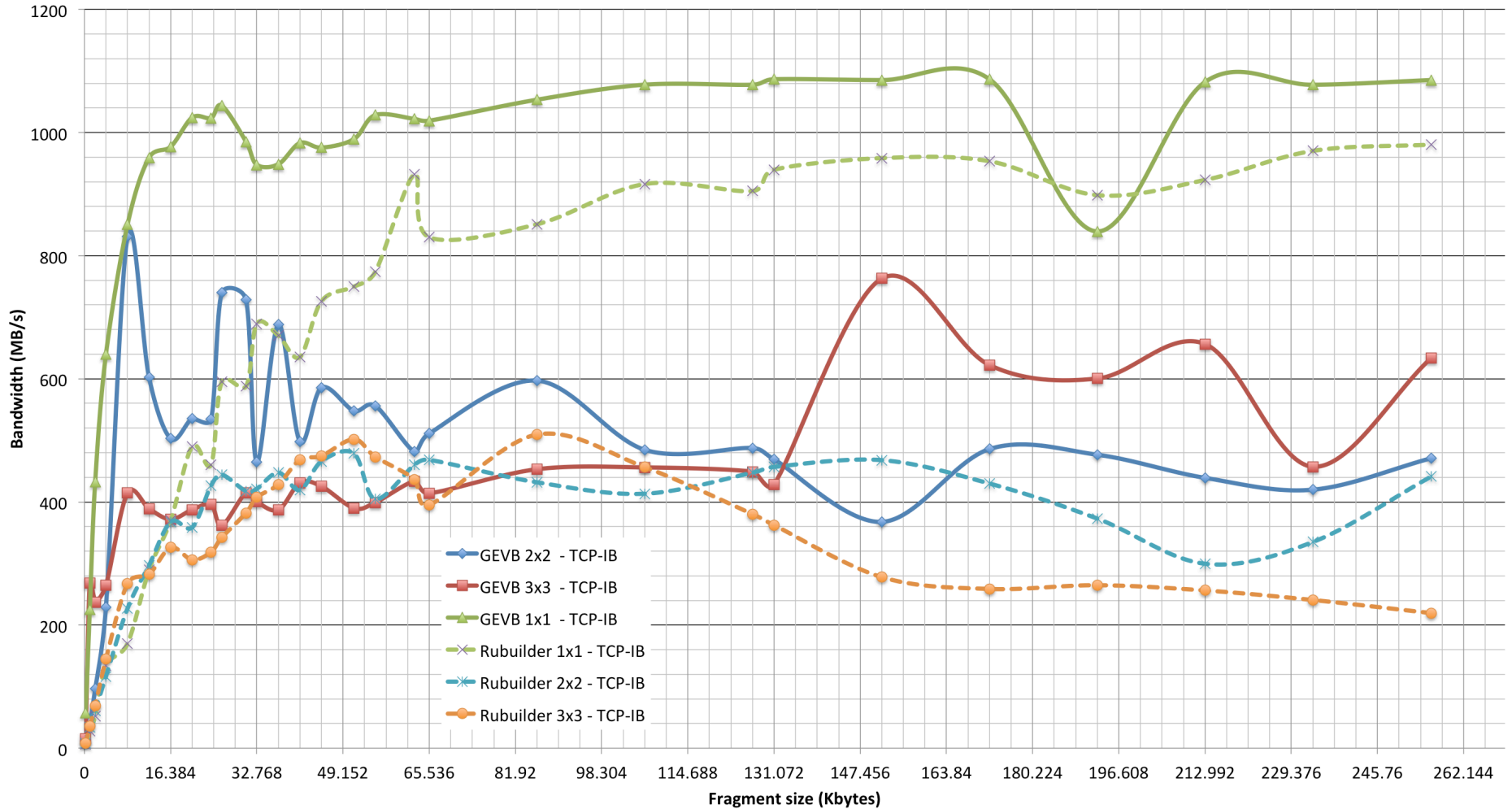




# EVB-NxN ptaTCP

## Bandwidth NxN - GEVB vs Rubuilder (TCP/IB)

XDAQ uDAPL peer transport preliminary measurements on QLogic 4x QDR - CMS DAQ upgrade studies  
September 2011 by L.Orsini, A. Petrucci





# EVB-1xN pt**a**TCP

## Bandwidth 1xN - MstreamIO vs Rubuilder (TCP/IB)

XDAQ uDAPL peer transport preliminary measurements on QLogic 4x QDR - CMS DAQ upgrade studies  
September 2011 by L.Orsini, A. Petrucci

