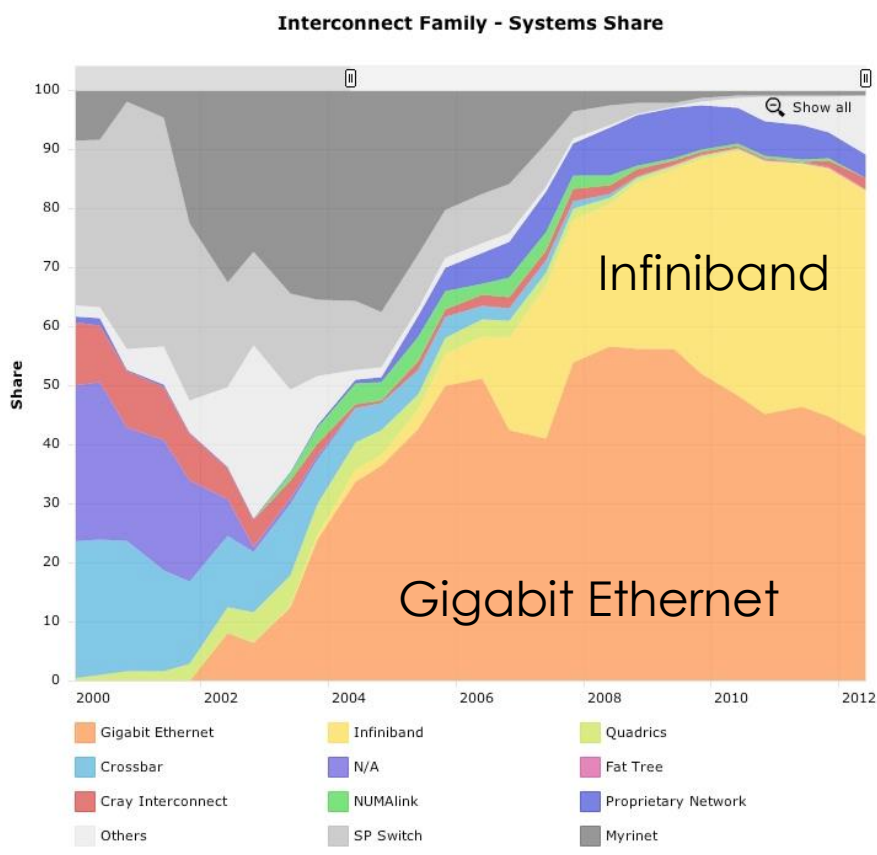# Boosting Event Building Performance Using Infiniband FDR for CMS Upgrade

**Andrew Forrest** – CERN (PH/CMD)

Technology and Instrumentation in Particle Physics Conference
2 June 2014, Amsterdam

# Outline

- Motivations

- A Quick Overview of Infiniband

- CMS Data Acquisition
  - Event Building

- CMS Online Software Framework (XDAQ)

- Integrating Infiniband
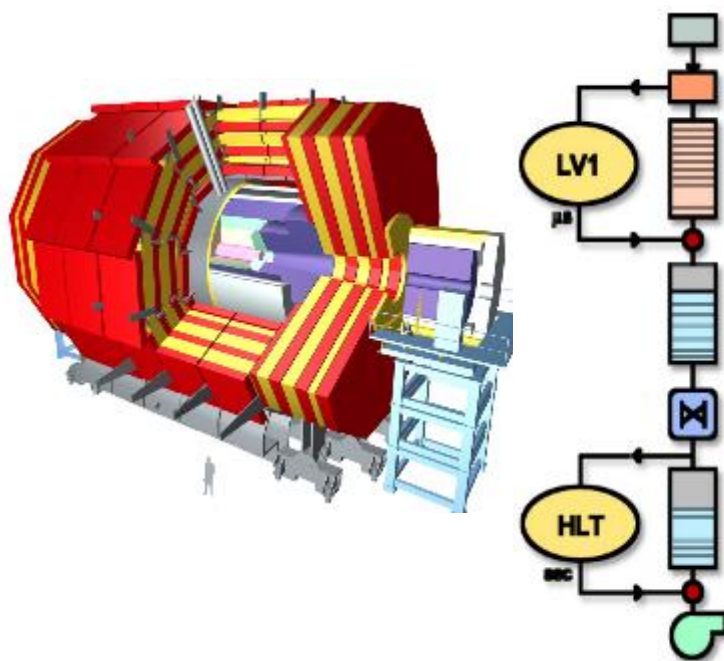
- Preliminary Results

- Conclusions

- DAQ hardware (PC's) from run 1 at end-of-life (>5 years old)

- Cost-effective solution that meets the requirements for run 2

- Opportunity to take advantage of technological advances

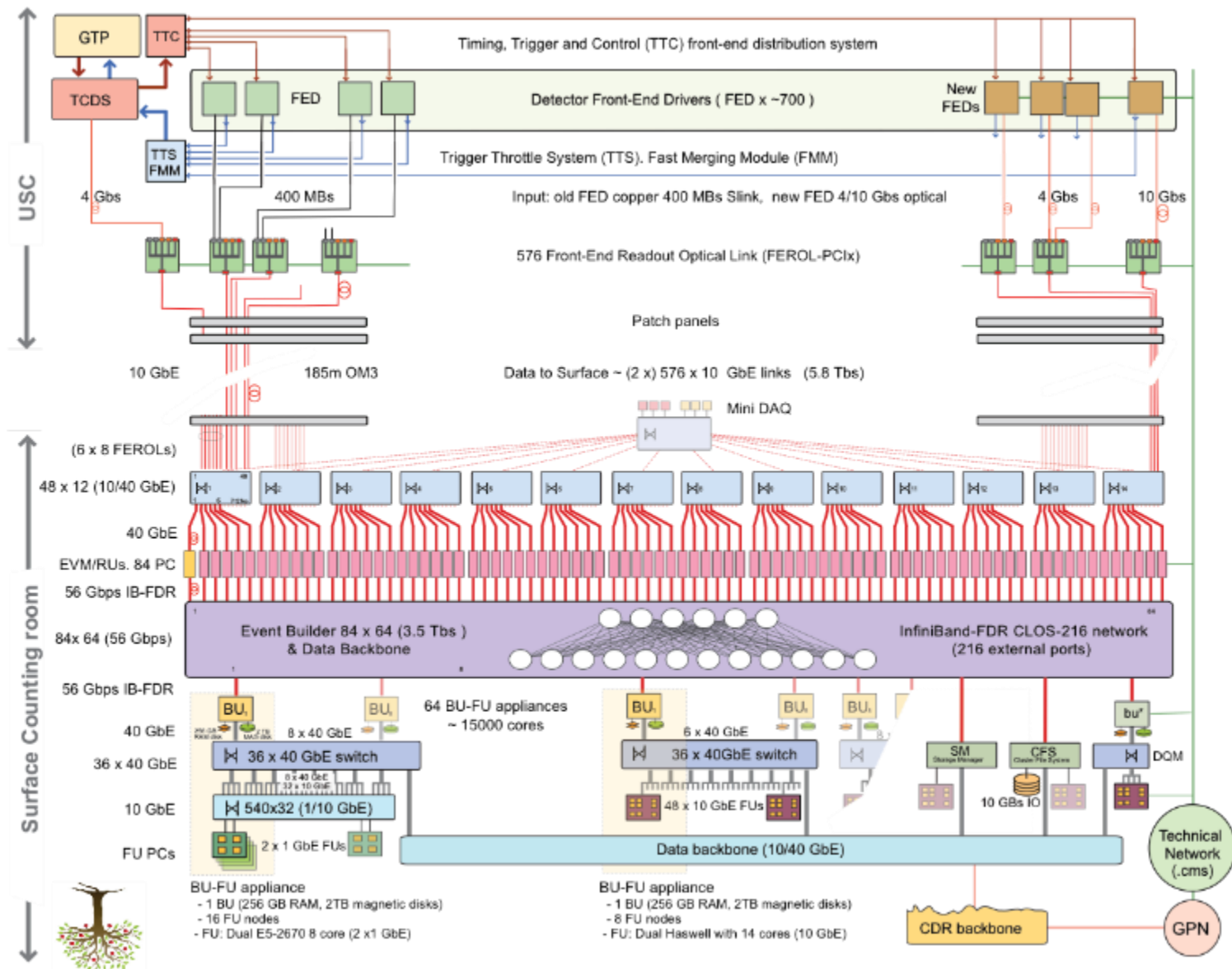**Interconnect Family - Systems Share**



Infiniband

Gigabit Ethernet

| | | |
|---|---|---|
| Gigabit Ethernet | Infiniband | Quadrics |
| Crossbar | N/A | Fat Tree |
| Cray Interconnect | NUMAlink | Proprietary Network |
| Others | SP Switch | Myrinet |

top500.org

# A Quick Overview of Infiniband

- Infiniband is…
  - A switched fabric computer network communications link for high performance
  - Reliable communication
  - Supports message based transfer using send/receive semantic

- Multiple programming methods
  - Infiniband verbs
  - uDAPL (user-level Direct Access Programming Library)
  - IPoIB

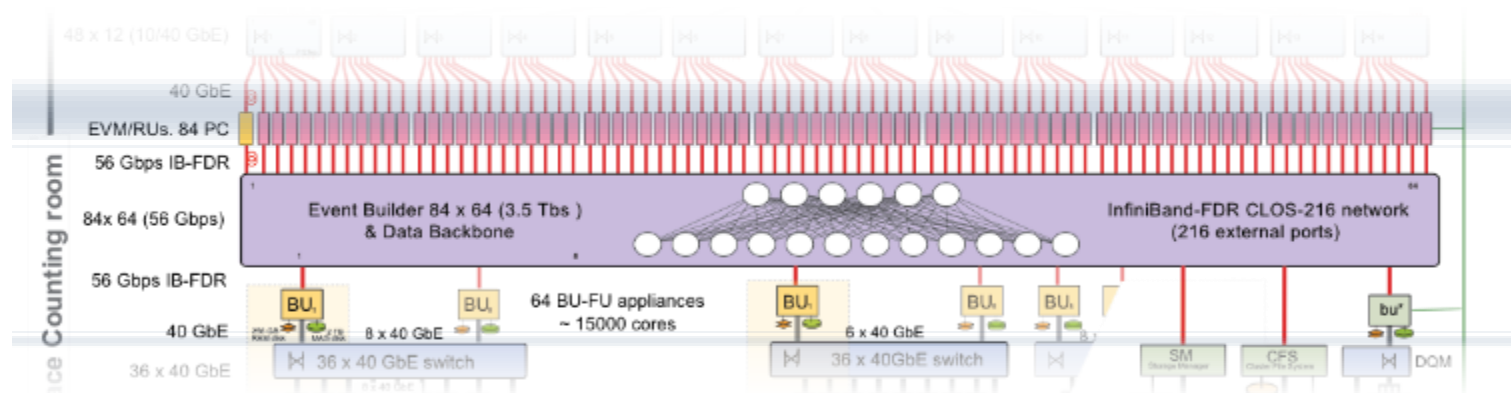- Software available as part of the OFED (OpenFabrics Enterprise Distribution)

# CMS Data Acquisition

| Parameters | |
|---|---|
| Data Sources (FEDs) | ~ 620 |
| Trigger levels | 2 |
| First Level rate | 100 kHz |
| Event size | 1 to 2 MB |
| Readout Throughput | 200 GB/s |
| High Level Trigger | 1 kHz |
| Storage Bandwidth | 2 GB/s |

## Resulting Required Physical Resources

| Resources Type | Run 1 | Run 2 |
|---|---|---|
| RU | ~640 | 84 |
| BU | ~1260 | 64 |
| FU | | ~1260 |
| Throughput Requirement | 100 Gb/s | 200 Gb/s |

Event fragments stored in physically separated memory systems

Readout Units

Input Sources (e.g. Detector Frontend)

Event Manager (EVM)

Builder Networks

Builder Units

Fragments collected into full event by builder units
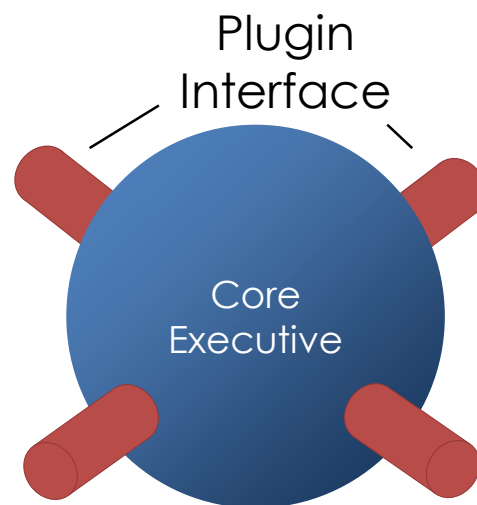
# CMS Online Software Framework (XDAQ)

- The XDAQ is software platform created specifically for the development of distributed data acquisition systems

- Implemented in C++, developed by the CMS DAQ group

- Provides platform independent services, tools for inter-process communication, configuration and control

- Builds upon industrial standards, open protocols and libraries, and is designed according to the object-oriented model
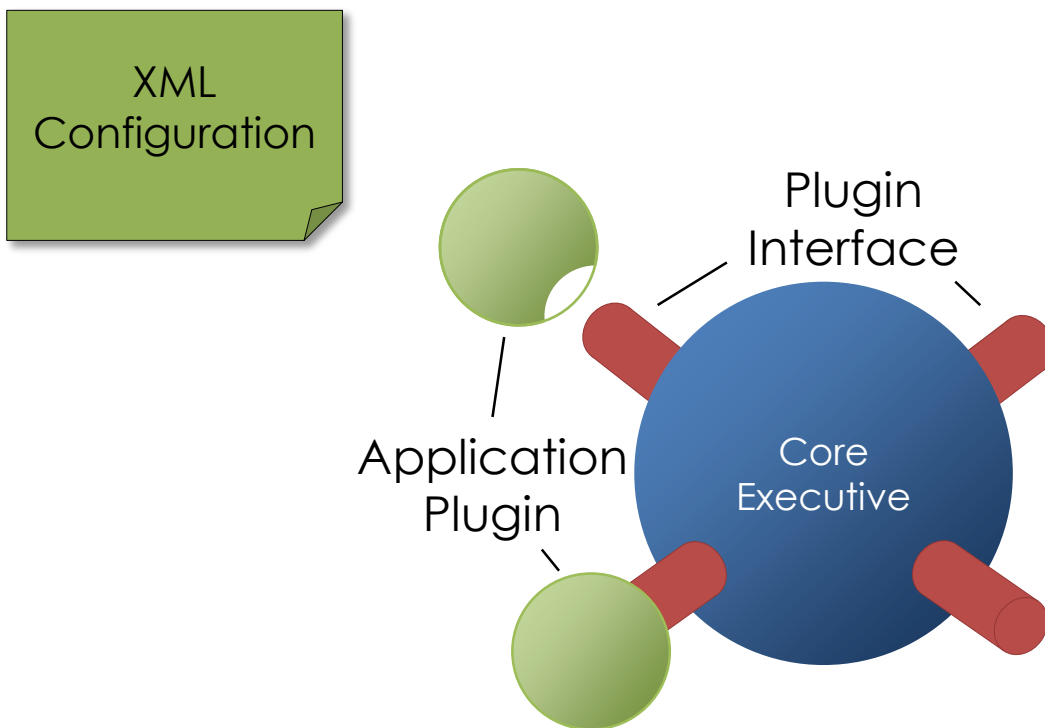
For further information about XDAQ see:
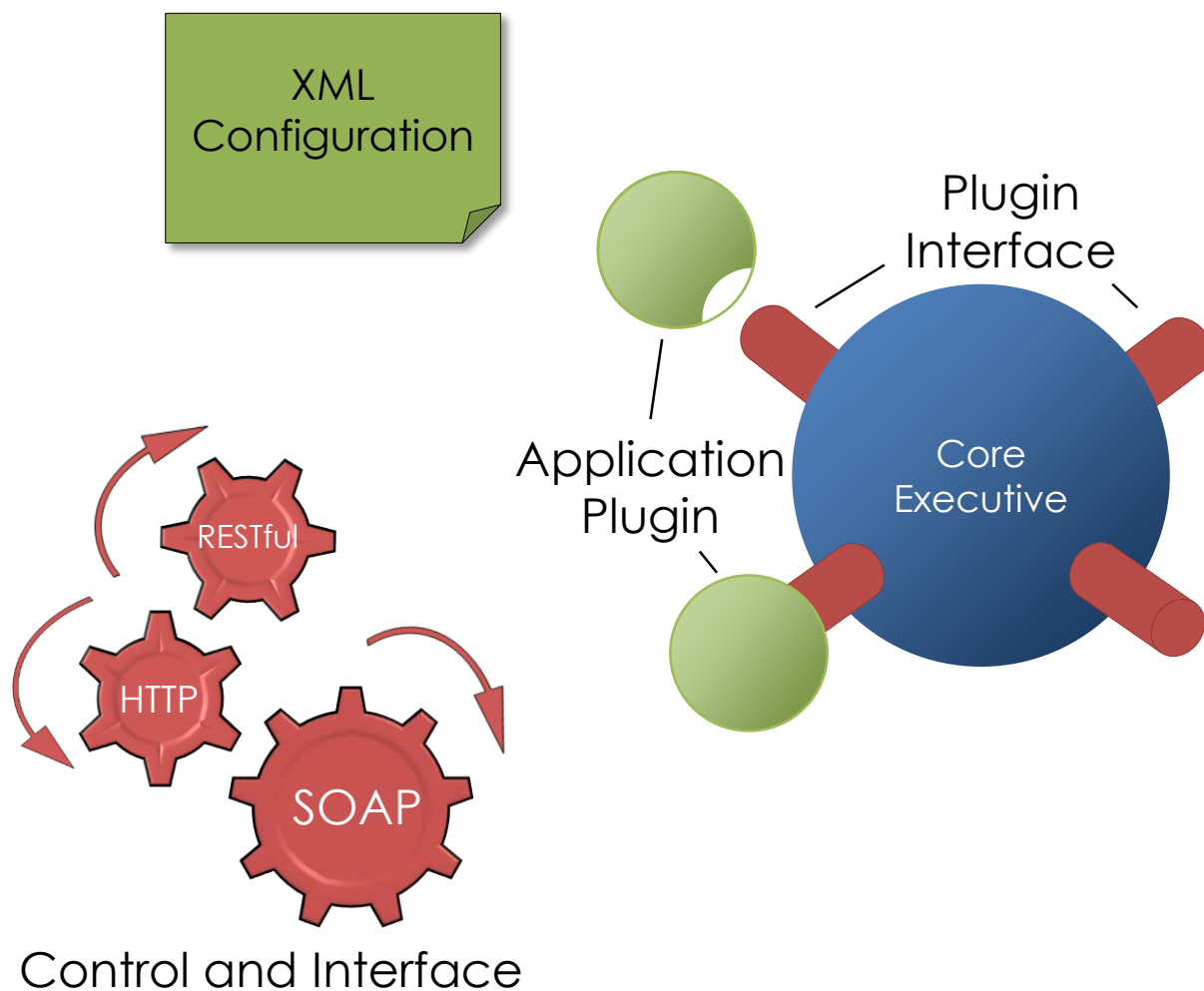
J. Gutleber, S. Murray and L. Orsini, Towards a homogeneous architecture for high-energy physics data acquisition systems published in Computer Physics Communications, vol. 153, issue 2, pp. 155-163, 2003
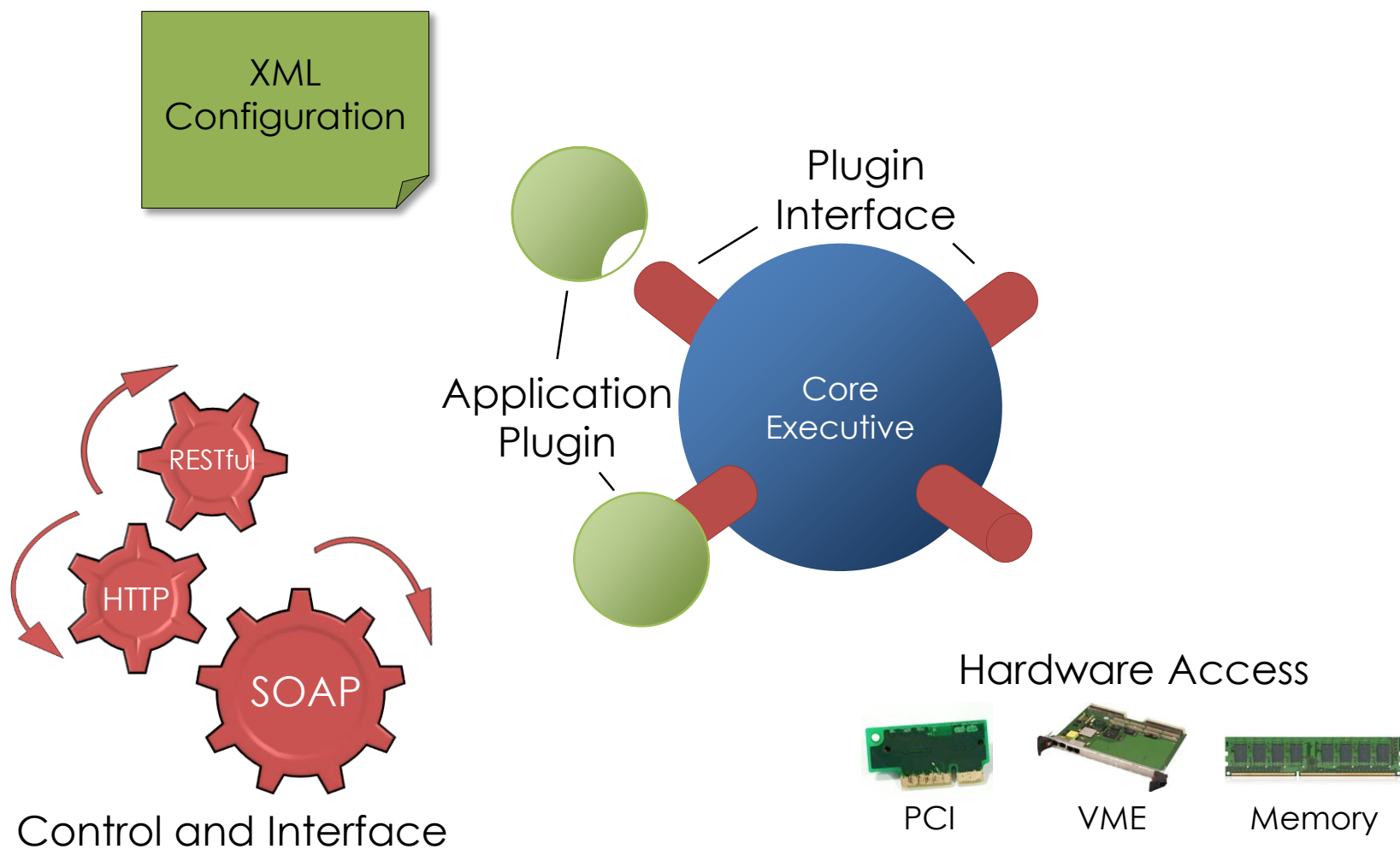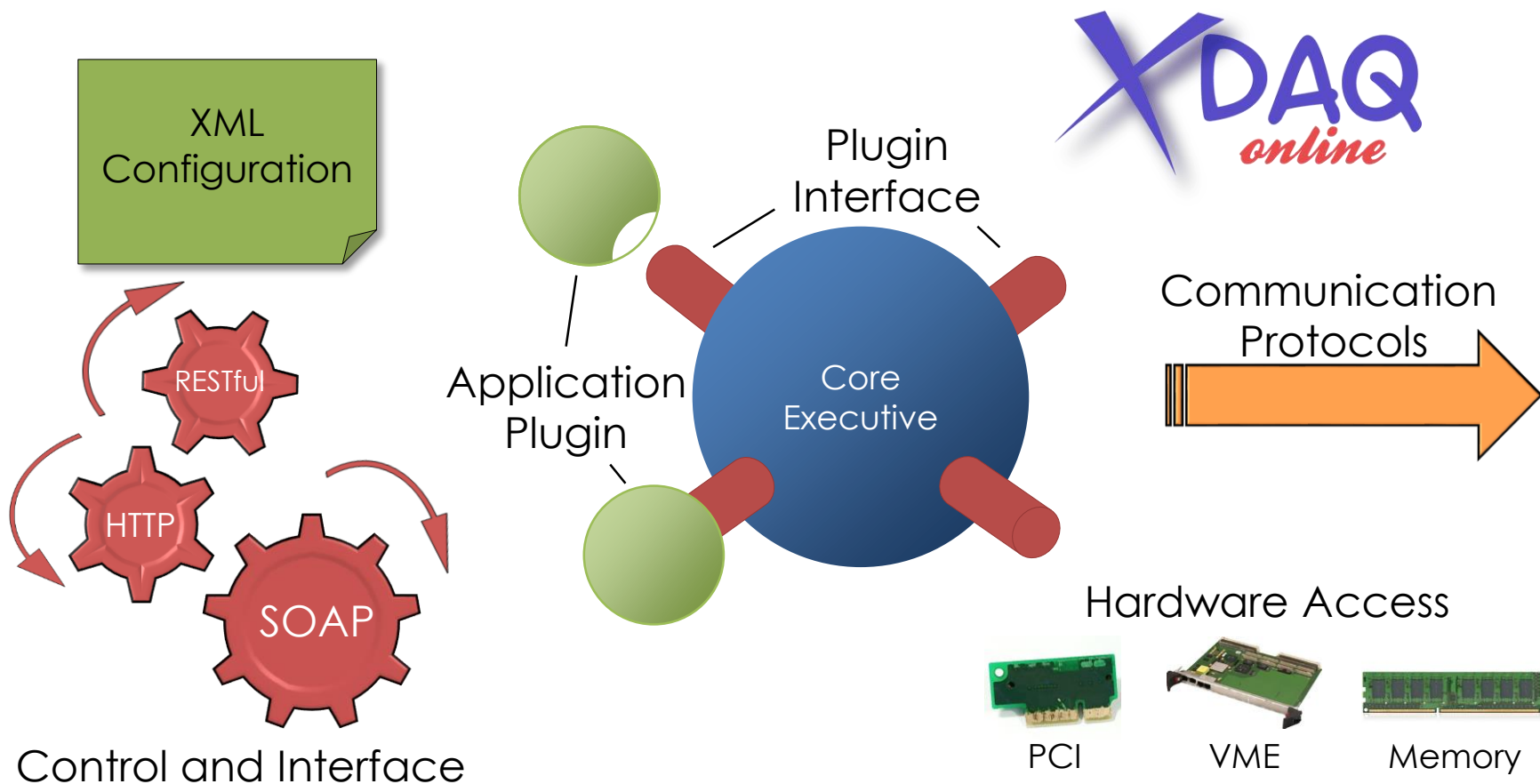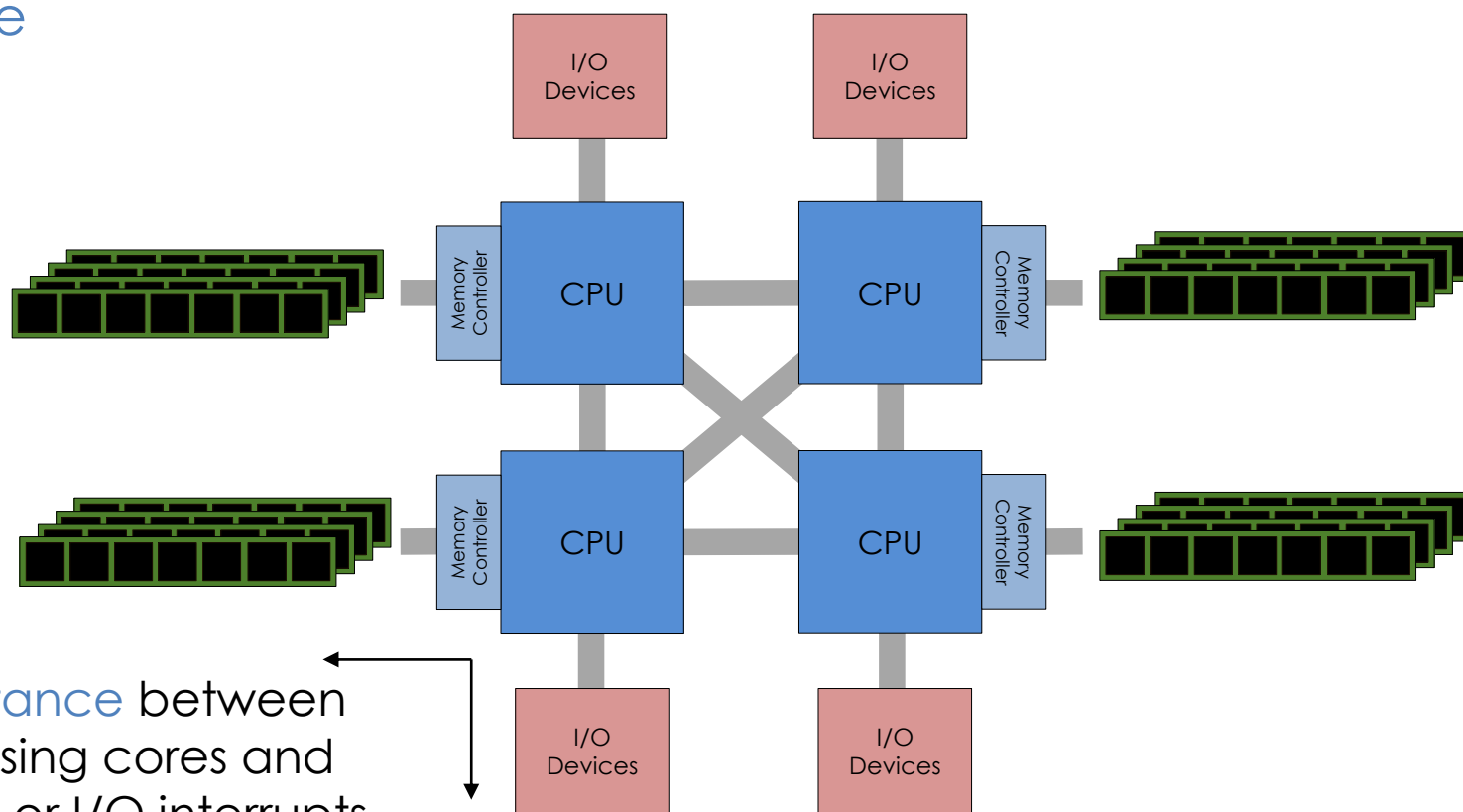http://www.sciencedirect.com/science/article/pii/S0010465503001619

Plugin
Interface

Core
Executive

XML Configuration

Plugin Interface

Application Plugin

Core Executive

XML
Configuration

Plugin
Interface

Core
Executive

Application
Plugin

RESTful

HTTP

SOAP

Control and Interface

XML Configuration

Plugin Interface

Application Plugin

Core Executive

RESTful

HTTP

SOAP

Control and Interface

Hardware Access

PCI          VME          Memory

☐ **Uniform building blocks** - One or more executives per computer contain application and service components



XML Configuration

Plugin Interface

XDAQ online

RESTful

HTTP

SOAP

Application Plugin

Core Executive

Communication Protocols

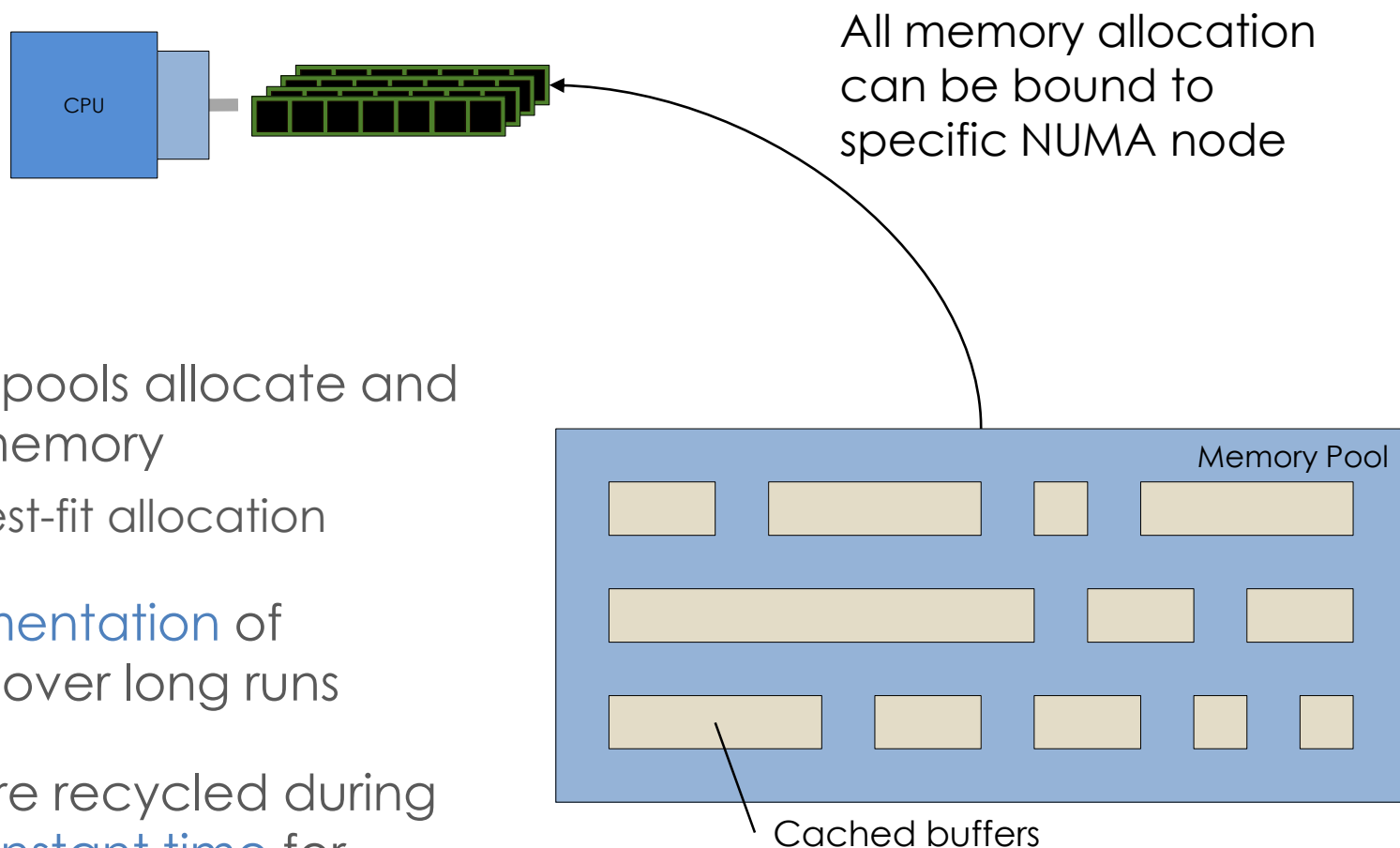Hardware Access

Control and Interface

PCI     VME     Memory

- Non Uniform Memory Access (NUMA)

- XDAQ provides utilities to take advantage



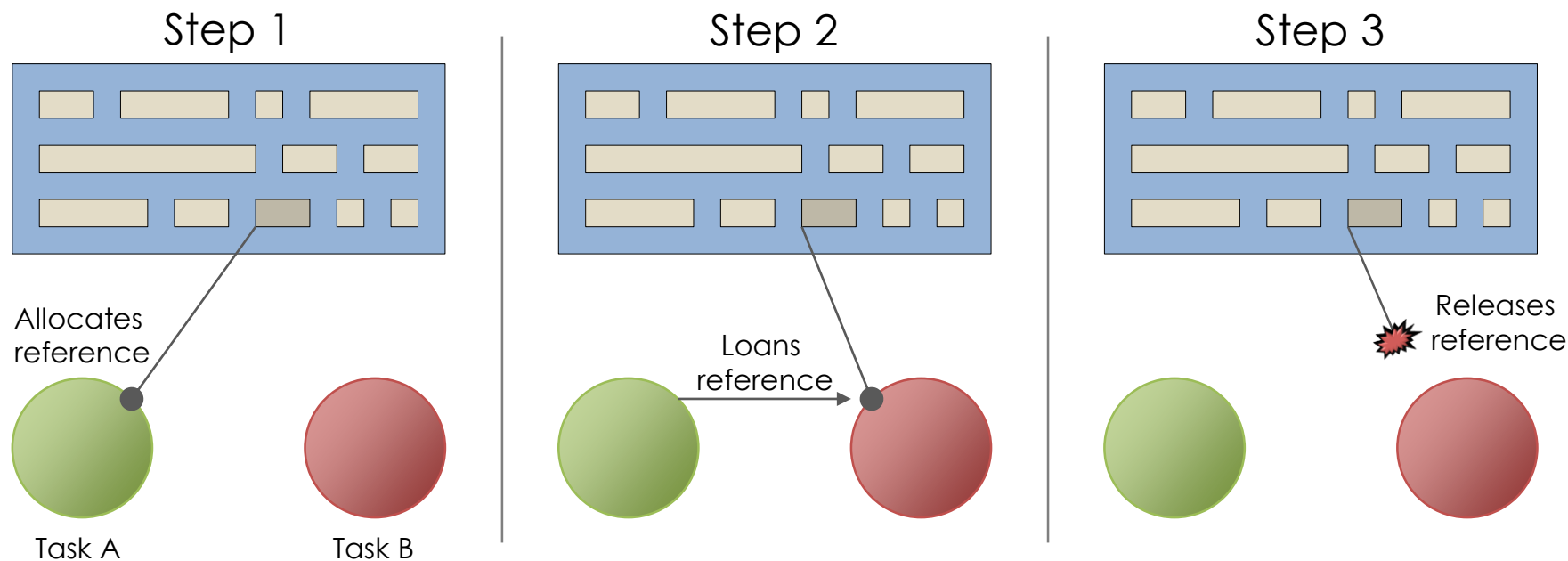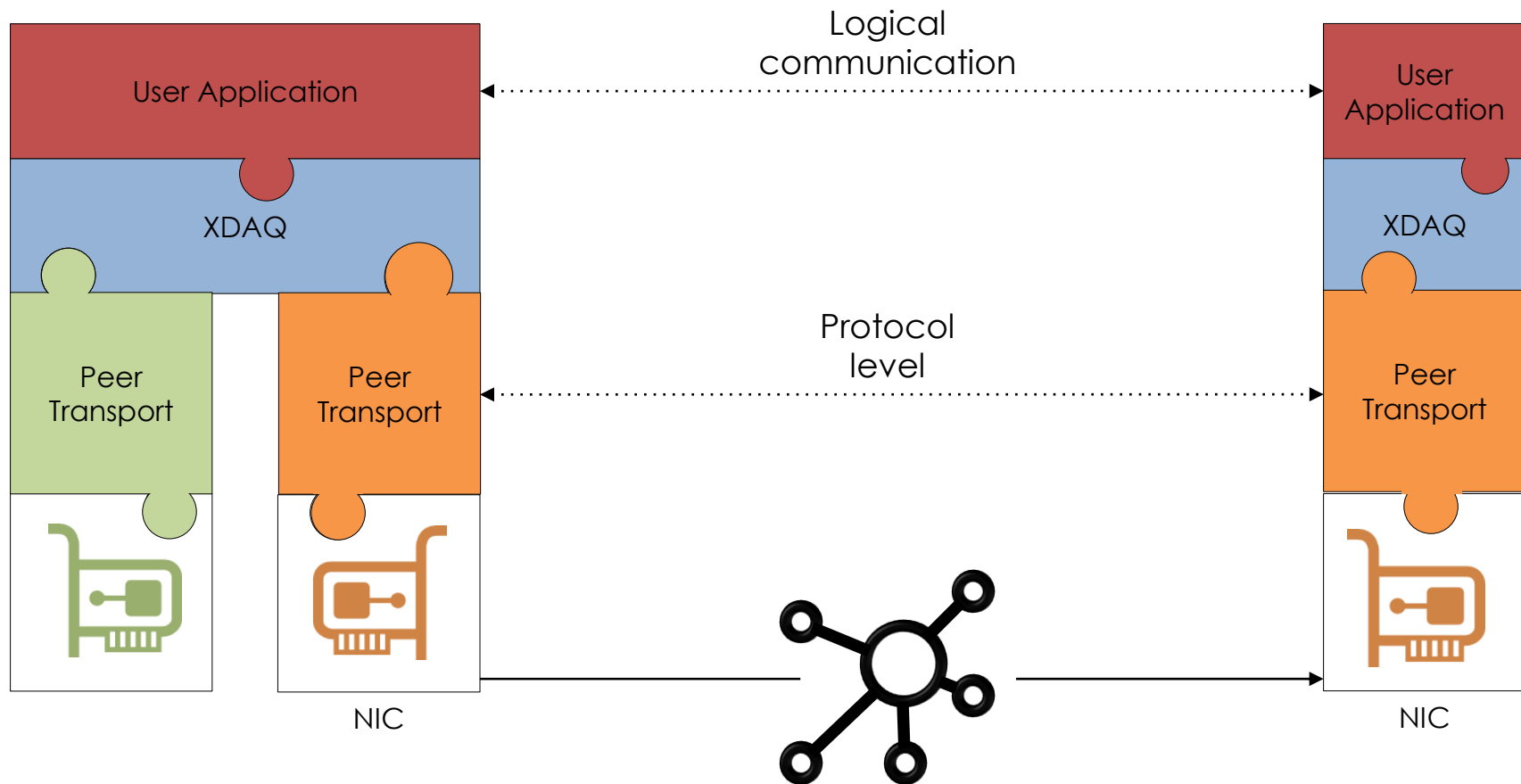The distance between processing cores and memory or I/O interrupts varies for each core

Workloop
Thread

CPU

Assignment
of work

Application
Thread

- ◻ Workloops can be bound to run on specific CPU cores by configuration

- ◻ Work assigned by application

- ◻ Workloops provide easy use of threads

All memory allocation can be bound to specific NUMA node

- ▣ Memory pools allocate and cache memory
  - ▣ $\log_2$ best-fit allocation

- ▣ No fragmentation of memory over long runs

- ▣ Buffers are recycled during runs – constant time for retrieval

Memory Pool

Cached buffers

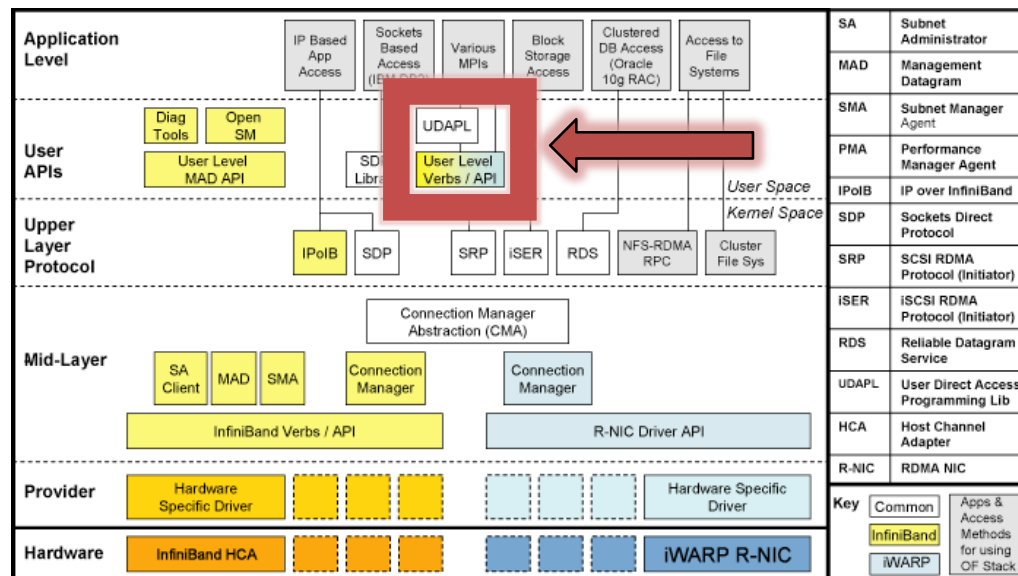- Buffer loaning allows zero-copy of data between software layers and processes

- User application is network and protocol independent

- Routing defined by XDAQ configuration

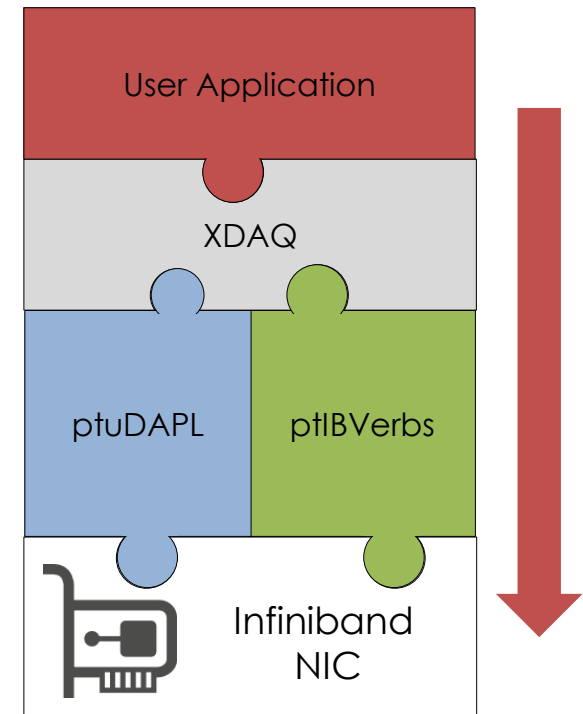- Connections setup through Peer to Peer model

# Integrating Infiniband

- Evaluation of two programming libraries
  - uDAPL
  - verbs

- Both from OpenFabric distribution

- uDAPL has connection support

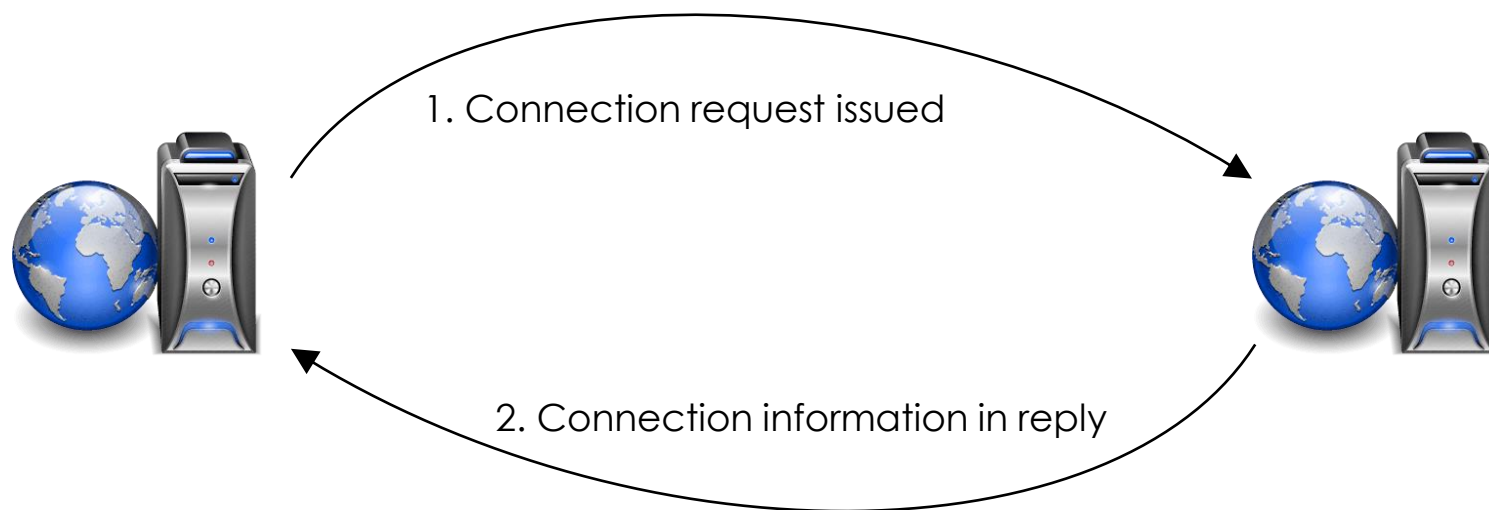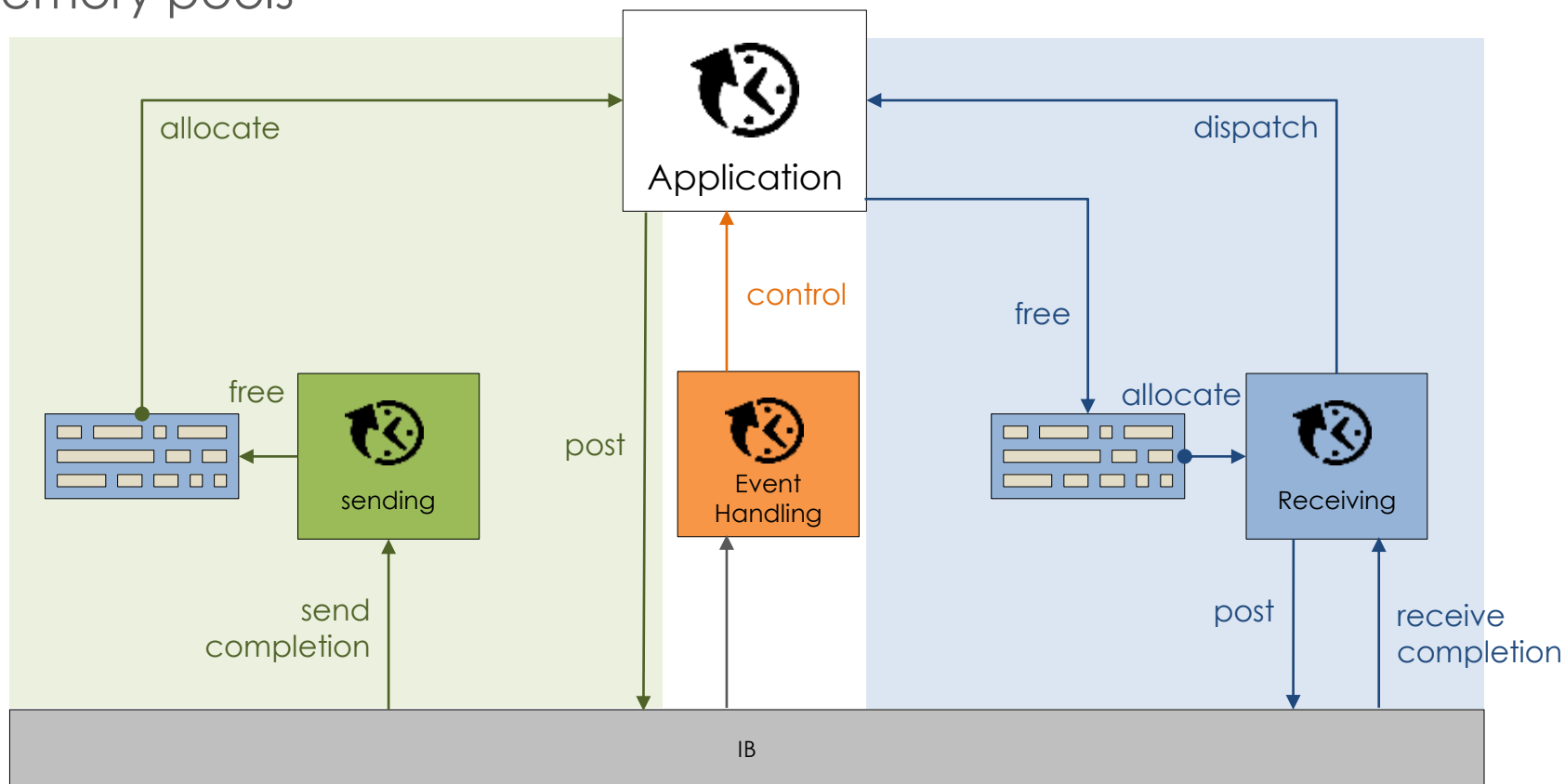- verbs is lower level in the software stack

- Two new XDAQ Peer Transports
  - uDAPL -> ptuDAPL
  - verbs   -> ptIBVerbs

- Full integration into XDAQ framework

- Event based API

- Send/receive with reliable connections

- Buffer loaning – zero-copy

- Memory pools automatically register memory with the NIC
  - Translation of virtual to physical addresses
  - Pinning memory to avoid swapping

- uDAPL provides IP address based connections

- verbs leaves the question of connecting peers open…

- For ptIBVerbs, a custom connection mechanism was implemented based upon IPoIB

1. Connection request issued

2. Connection information in reply

- Work is distributed across several XDAQ workloops

- Workloops are bound to run on one or more cores

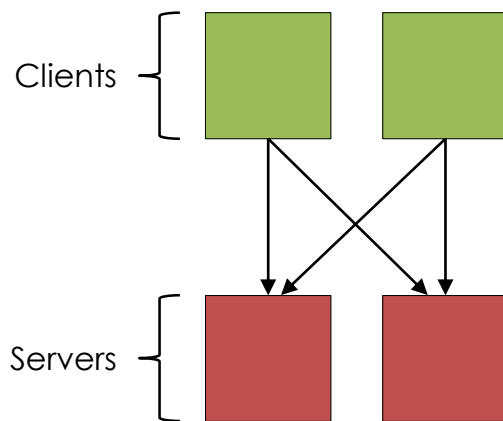- Sending/receiving operations separated and use different memory pools
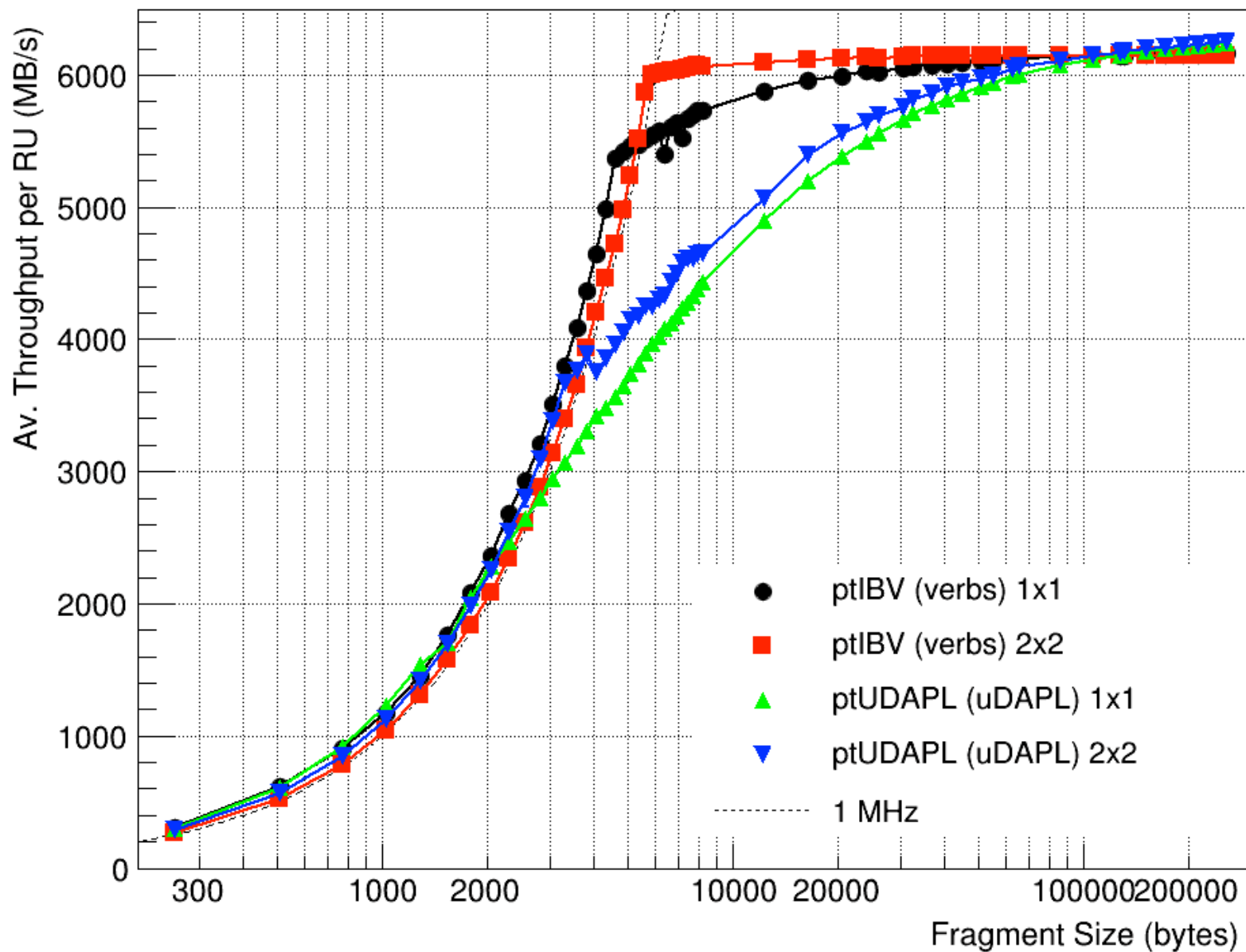
# Preliminary Results

- Small scale with 4 nodes on 1 switch
    - 1x1 and 2x2 (RU x BU) tests for…

- N-to-N and event building tests

- Each node has…
    - Dual socket Intel Xeon E5-2670 8-core processors @ 2.6 GHz
    - 16 GB RAM per socket (NUMA)
    - Mellanox Connect X-3 VPI Infiniband FDR network card
    - OFED v 2.x
    - Scientific Linux (CERN) 6

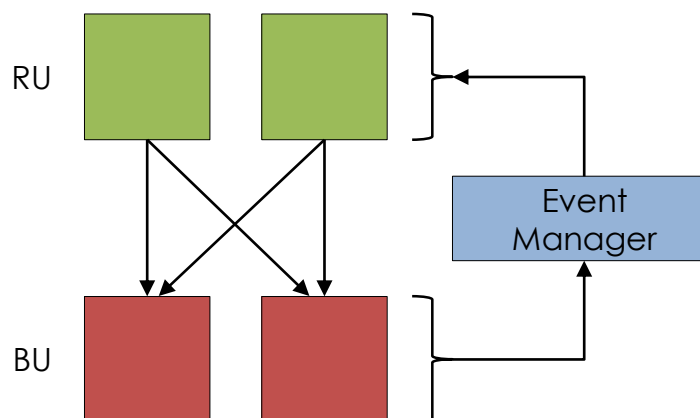|     | SDR | DDR | QDR | FDR-10 | FDR | EDR |
| --- | --- | --- | --- | --- | --- | --- |
| 1X | 2 | 4 | 8 | 9.67 | 13.64 | 25 |
| 4X | 8 | 16 | 32 | 38.79 | 54.54 | 100 |
| 12X | 24 | 48 | 96 | 116.36 | 163.64 | 300 |

Effective unidirectional theoretical throughput in Gb/s

- N clients each send to N servers for each 'message'

- The measurement is the rate of receiving in the receivers

- No additional processing

- Fixed sized messages, round robin dispatching in the senders

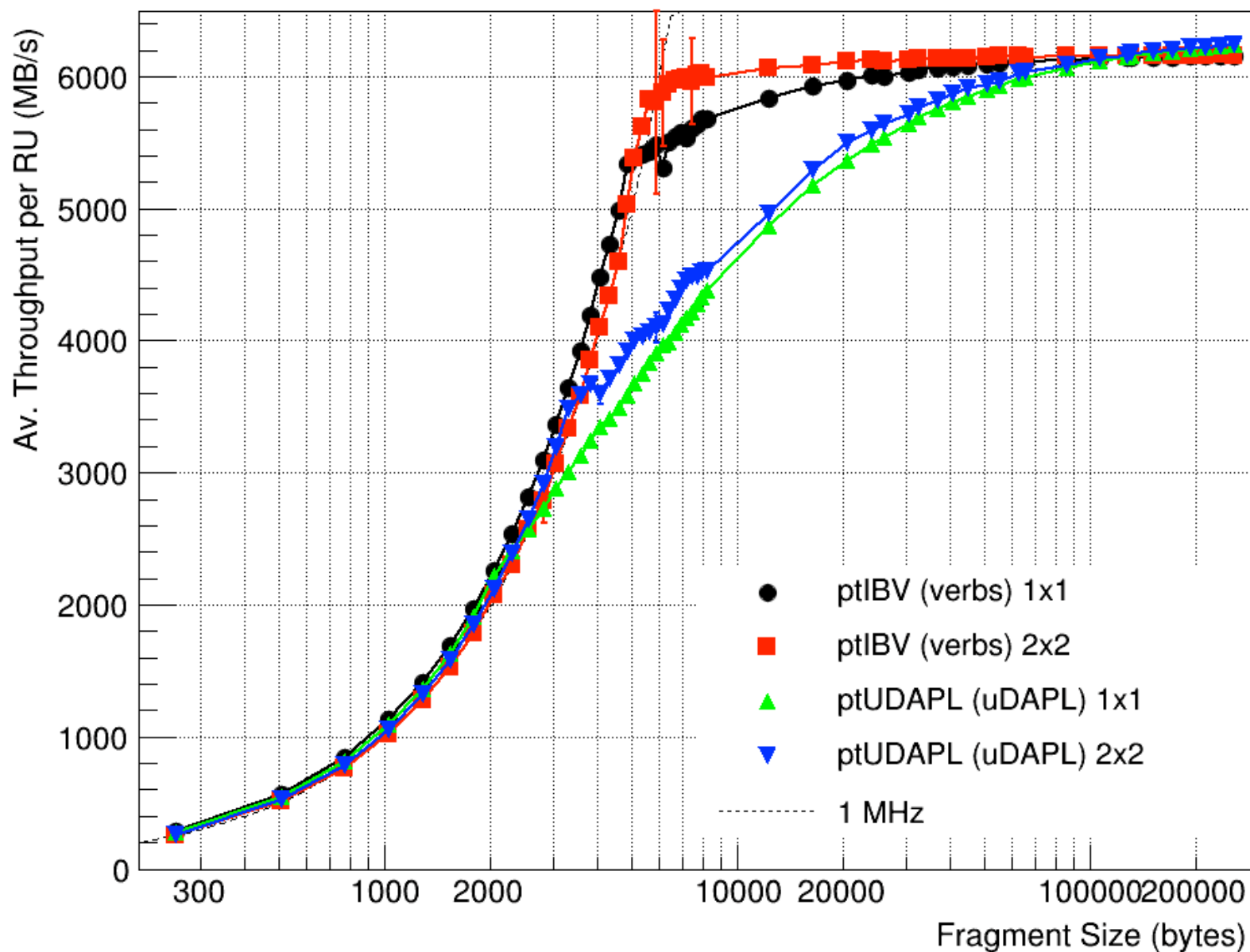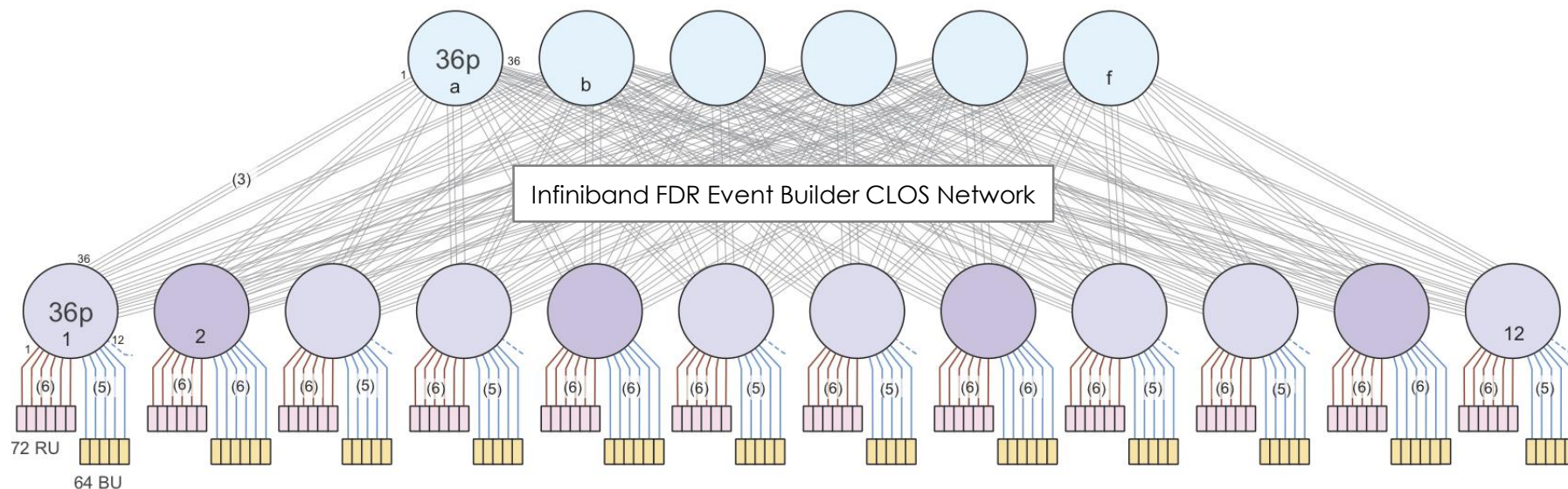- Test to show the performance for unidirectional throughput

- Event fragments are generated in the RU's

- Fully built event are dropped in the BU's

- The measurement is the rate of receiving in the BU's
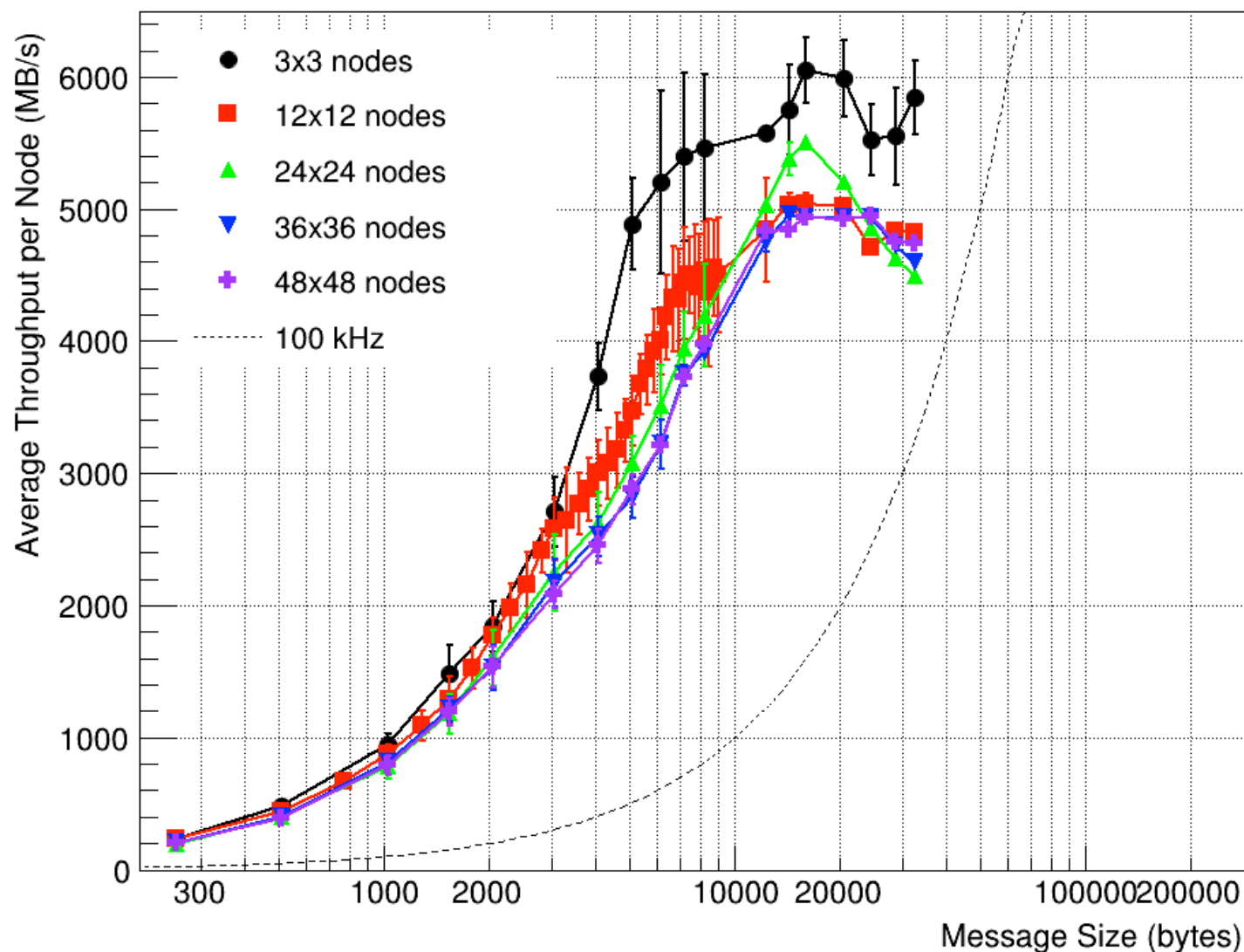
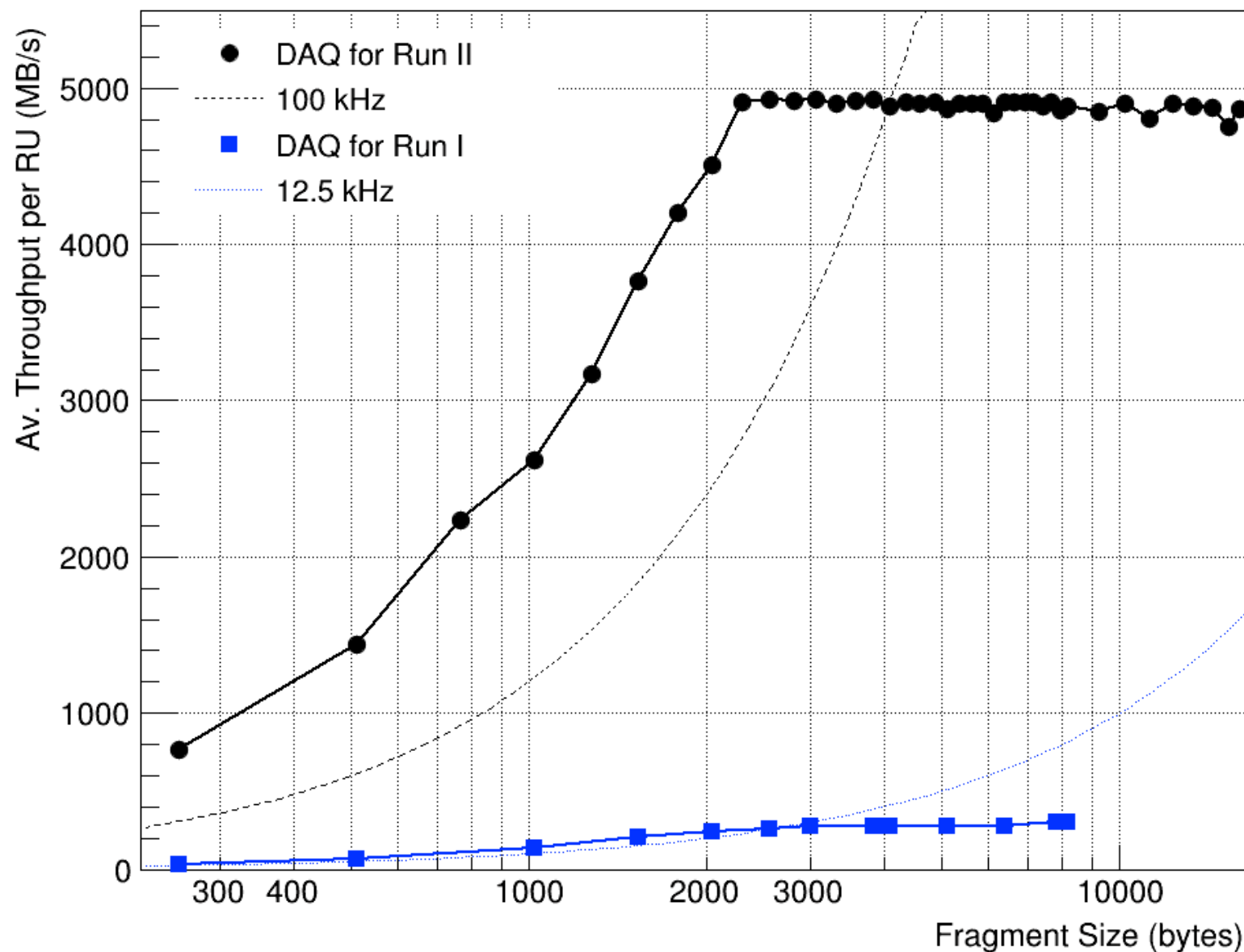- Additional control messages with Event Manager

- ☐ ptIBVerbs used for larger scale tests
  - ☐ up to 48x48 using Infiniband CLOS network

- ☐ Preliminary N-to-N tests

Infiniband FDR Event Builder CLOS Network

Message Size v Throughput

# Conclusions

- Infiniband works well with event building applications

- and the CMS Online Software framework (XDAQ)

- CMS DAQ will be using ptIBV for data flow in LHC run 2

- Performance compared to DAQ 1 allows for an order of magnitude of reduction in physical resources for event building

- In the future…
  - Full DAQ2 tests

- Questions?

# Additional Materials