# Infiniband Event-Builder Architecture Test-beds for Full Rate Data Acquisition in LHCb

Enrico Bonaccorsi, Juan Manuel Caicedo Carvajal, Jean-Christophe Garnier, Guoming Liu, Niko Neufeld and Rainer Schwemmer

CERN

Computing in High Energy and Nuclear Physics, 2010

# Introduction

- Upgrade of the experiment involving a multi-terabit readout
- Interesting technologies
  - 10 Gigabit Ethernet : successor of Gigabit Ethernet
  - InfiniBand : Challenger
- Is InfiniBand an interesting alternative to 10 Gigabit Ethernet ?
- What software implementation for the Event-Builder ?

# Outline

# Outline

# Outline

# Outline

# Protocols

# Figures

|                       | Current | Upgrade |
|-----------------------|---------|---------|
| Event size            | 35 kB   | 100 kB  |
| Read-out rate         | 1 MHz   | 30 MHz  |
| Sources               | 313     | 1000    |
| Sinks                 | 1500    | > 3000  |
| Event-rate to storage | 2 kHz   | 10 kHz  |

# Aims

- Study the two main multi-gigabit technologies
- Study different network topologies
  - ▶ Cf. PS10-3-338 from Guoming Liu
- Study various software stacks
  - ▶ Push and pull protocols
  - ▶ Non-oriented connections protocols: Unreliable (UDP) and Realiable (RDS) datagrams
  - ▶ Oriented connections protocols: TCP, SDP
- Sources are FPGA, simple protocols are more interesting

# Outline

# Overview

- Distributed computing
- Low latency
- Very promising bandwidth
- Data rate (and not signal rate):

|      | **SDR**    | **DDR**    | **QDR**    | **FDR**     | **EDR**     |
|------|------------|------------|------------|-------------|-------------|
| **1X**  | 2 Gbit/s   | 4 Gbit/s   | 8 Gbit/s   | 14 Gbit/s   | 25 Gbit/s   |
| **4X**  | 8 Gbit/s   | 16 Gbit/s  | 32 Gbit/s  | 56 Gbit/s   | 100 Gbit/s  |
| **12X** | 24 Gbit/s  | 48 Gbit/s  | 96 Gbit/s  | 168 Gbit/s  | 300 Gbit/s  |

- SDR, DDR and QDR use 8B/10B encoding
- Aggregation of links in units of 4
- Only a few vendors
- Huge software stack

# Software stack: OpenFabrics Enterprise Distribution



OpenFabrics Software Stack

# IPoIB

- Standard IP encapsulation over InfiniBand fabrics
- Relies on 2 modes of InfiniBand
  - Unreliable datagram: Max MTU = link MTU = 4096 B
  - Connected mode: Max MTU = $2^{31} B$
- Implementation using the libc socket
- No changes to your current source code
- No RDMA

IP Based
App
Access

IPoIB

# Sockets Direct Protocol (SDP)

- Defines a standard wire protocol over IB fabrics
- Local IP assignments and IP resolution using IPoIB
- Supports only stream sockets (*SOCK_STREAM*)
- 2 ways to use
  - Minimal reimplementation
  - LD_PRELOAD=libsdp.so
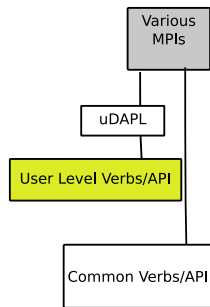- Minor changes to your current source code
- Good use of RDMA

```
Sockets
Based
Access
```

```
SDP Lib
```

```
SDP
```

# Reliable Datagram Sockets (RDS)

- Developed by Oracle and added to the OFED
- Not restricted to InfiniBand
- 1 socket to talk to N destinations
- Included in Linux kernel 2.6.30
- New source code
- Minimal use of RDMA

clustered
DB Access

RDS

# Verbs, MPIs and standard limits

- Allows full use of RDMA
- List of standard verbs, syntax is vendor specific
- About 1200 pages of verbs out of the 1700 page InfiniBand Architecture Specification
- If used correctly, maximum performances
- Each vendor brings its own non standard optimized MPI implementation

Various MPIs

uDAPL

User Level Verbs/API

Common Verbs/API

# Outline

# Configuration



Server  Server  Server  Server

qlogic 12300
- QDR 4x
- 32 Gb/s

Server  Server  Server  Server

# Server configuration

| Processor type | Intel Xeon E5520 |
|---|---|
| Processors x cores x clock (GHz) | 2 x 4 x 2.27 |
| RAM (GiB) | 3 |
| HCA | qle7340 4x QDR |
| kernel | 2.6.18 |

| net.core.rmem_default | 16777216 |
|---|---|
| net.core.wmem_default | 16777216 |
| net.core.rmem_max | 16777216 |
| net.core.wmem_max | 16777216 |
| net.core.netdev_max_backlog | 250000 |

# Implementation

- Node synchronization via Precise Time Protocol (ptp)
- From scratch, plain C/C++ for Linux, focusing on network performance and troubleshooting
- Portability InfiniBand/10GbE
  - TCP/IP stack done at kernel/API level, except RAW Ethernet/IP
  - SDP only related to InfiniBand
  - RDS fully portable
  - RDMA: InfiniBand verb syntax used by some 10GbE suppliers

# Preliminary performance results: PUSH over UDP

# Preliminary performance results

# First experience feedback

- A few issues experienced at the beginning
  - Performance drop over time
  - Negotiation problems
  - Congestion management to be investigated
  - Far from line rate
- Lack of online literature
- Always go to support, they are willing to work with you
- IPoIB not optimized according to the QLogic support
- MPIs are the best !
  - Dedicated implementations should allow maximum performance

# Summary

- Push protocol over unreliable datagrams established
  - Optimization possible: aggregation of several IB interfaces over one IB link
- Future work
  - Implement PUSH over RDS and open MPI
  - Implement PULL over RDS, TCP, SDP and open MPI
    - TCP, SDP and MPIs might be too complex for FPGAs
    - Still interesting to know their performance