

#### Experience with Infiniband for CMS Event Building

Andrea Petrucci – CERN (PH/CMD)

Istituto Nazionale di Fisica Nucleare (INFN), CNAF Seminar 19 September 2014, Sala Garisenda, Via Ranzani n. 13/2, Bologna, Italy

### Outline



- CMS Experiment in LHC run 2
  Event builder changes
- Feasibility studies
- Infiniband Verbs (ptIBVerbs)
- DAPL (ptUDAPL)
- IB Verbs vs DAPL
- Performance tuning
- Preliminary Results
- Conclusion

## CMS Experiment in LHC run 2

### CMS DAQ Requirements for LHC Run 2





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

### Event builder changes







### CMS DAQ System for LHC Run 2



## Feasibility studies

### Networking technologies



Our feasibility studies are focused in two network technologies

#### Ethernet

- 10/40 Gigabit Ethernet (different vendors)
- iWARP (RDMA) TCP/IP full offload (Chelsio T4 Unified Wire Adapters)
- performance measurements using TCP/IP and DAPL (Direct Access Programming Library- OpenFabrics)

#### Infiniband

- 4x quad data rate (QDR)
  - 40 Gb/s 8B/10B encoding -32 Gb/s data rate
- 4x fourteen data rate (FDR)
  - 56 Gb/s 64B/66B encoding 54.54 Gb/s data rate
- performance measurements using DAPL (Direct Access Programming Library- OpenFabrics), IB verbs (Infiniband verbs) and IPoIB (IP over InfiniBand)



9

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

<u>http://www.openfabrics.org/</u>



....

Curbon et

Application		IP Based B	ockets Based Various	Block	Clustered DB Access	Access to	>	SA	Administrator	ALLIANCE
Level		App Access (IBI	ccess MPIs M DB2)	Storage Access	(Oracle 10g RAC)	File Systems		MAD	Management Datagram	
	Diag Open Tools SM		UDAPL			-+		SMA	Subnet Manager Agent	
User APIs	User Level	SDP	User Level					РМА	Performance Manager Agent	
	MAD API						ser Space	IPolB	IP over InfiniBand	
Upper							emel Space	SDP	Sockets Direct Protocol	
Protocol		IPolB SDP	SRP	iser R	DS RPC	File	Sys	SRP	SCSI RDMA Protocol (Initiator)	
			Connection Mar	ager				ISER	iSCSI RDMA Protocol (Initiator)	
Mid-Layer	Abstraction (CMA)					RDS	Reliable Datagram Service			
	Client MAD SM	Mana Mana	ager	Manager				UDAPL	User Direct Access Programming Lib	
InfiniBand Verbs / API				R-NIC Driver API			HCA	Host Channel Adapter		
								R-NIC	RDMA NIC	
Provider	Hardware Specific Driver				Н	ardware \$ Drive	Specific r	Key Co	Apps & Access	
Hardware	InfiniBand HCA				iV	VARP I	R-NIC	Inf i	IniBand Methods for using WARP OF Stack	

Comparison of the Stacks (Infiniband vs TCP/IP)



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



### Infiniband Peer Transports



- 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



### Event builder software



The ptuDAPL and ptIBVerbs are pluggable components to transparently access the IB network in XDAQ – CMS online framework



## Infiniband verbs (ptIBVerbs)

### IB verbs specification



- RDMA read and write requires additional communication to alert remote processes of completed transactions. Therefore,
- We are using the 'Send' and 'Receive' operations with 'Reliable Connections'

Operation	UD	UC	RC	RD
Send	Х	Х	Х	Х
Receive	Х	Х	Х	Х
RDMA Write		Х	Х	Х
RDMA Read			х	Х
Atomic: Fetch and Add/Cmp and Swap			х	Х
Max Message Size	MTU	2GB	2GB	2GB

• Reliable Datagrams is not currently implemented in Mellanox HCA

## Infiniband Verbs Concepts



- **Device** The physical Host Channel Adapter (HCA)
- **Context** A user 'context' within which a user can interact with the HCA (1 per device per application)
- **Protection Domain** Memory is registered to a Protection Domain, which in turn registers the memory with the device
- Memory Region A Memory Region represents a buffer that is registered to a Protection Domain
- Queue Pair The IBV equivalent to a socket which internally consists of a send queue and a receive queue
- Work Request Send and receive requests are represented as Work Requests
- **Completion Queue** Queue Pairs report completion of Work Requests through Completion Queues (there can be many Queue Pairs for a single Completion Queue)
- Asynchronous Events Queue Errors and other asynchronous events are reported through a special asynchronous event queue

### Software Stack





### Infiniband Architecture Principles





### Peer Transport Memory Management





### Parallelism



- 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



### Infiniband Connections



- Infiniband supports "reliable connections", but they are really reliable, connectionless datagram messaging services
- Queue Pairs store routing information for the remote Queue Pair, which is used to route all outgoing packets
- A connected pair of Queue Pair's effectively have their internal send and receive queues connected to each other respectively
- The relationship between Queue Pairs is 1 to 1 (no multicast)



## ptIBV Connection Setup





Connector/Acceptor based on ACE (The ADAPTIVE Communication Environment)

### Events and Error Handling

•

٠



• Two types of notification, work completions and asynchronous event, which are sent through separate queues



### The Life of a Queue Pair





## DAPL (ptUDAPL)

- Developed by DAT collaborative
  - <u>http://www.datcollaborative.org/</u>
- Transport and platform (OS) independent
- User-Level Direct Access Transport APIs (UDAPL)
- Kernel-Level Direct Access Transport APIs (kDAPL)
- DAT supports reliable connection
- Data Transfer Operations send, receive, rdma\_read, rdma\_write
- UDAPL Version 2.x, January, 2007







The ptuDAPL is a pluggable component to transparently access the DAT library in XDAQ – CMS online framework

- All I/O operations centered on dedicated memory pool based on uDAPL memory region allocator
- Profiting for inherent non blocking and queuing of uDAPL API for minimizing latency
- Full zero-copy between CMS online applications and DAPL driver
- Based on DAT Spec 2.x



## IB verbs vs DAPL



ptIBV	ptuDAPL		
Complex API	Simple API		
High Performance	High performance		
High stability	Average stability		
High flexibility	Moderate flexibility		
Well documented	Poorly documented		
Open connection method	Connection method		
Core Infiniband Specification	3 <sup>rd</sup> Party		

## Performance Tuning

## Non-Uniform Memory Access (NUMA)



The distance between processing cores and memory or I/O interrupts varies within Non-Uniform Memory Access (NUMA) architectures



Non-Uniform Memory Access (NUMA)

### Performance Factors









Memory affinity



#### TCP Custom Kernel Settings

### Readout Unit Machine Affinities





#### **Affinities**

- 40 GbE I/O Interrupts
- 40 GbE Socket reading
- Infiniband I/O Interrupts
- DAQ threads
- DAQ Memory allocation
- □ Cores available to operating system

## Preliminary Results

### Test Setup



- 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-to-N Setup



- 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 Building Setup



- 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



### **Event Building**





### Larger Scale Tests



ptIBVerbs used for larger scale tests

- up to 48x48 using Infiniband CLOS network
- Preliminary N-to-N tests



### ptIBVerbs Performance – N-to-N



#### Message Size v Throughput



Experience with Infiniband for CMS Event Building, INFN CNAF, 19 September 2014, Sala Garisenda, Via Ranzani n. 13/2, Bologna, Italy, Andrea Petrucci (CERN) 40

### DAQ1 v DAQ2 RU Performance







### 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 ?



- Thank you for listening
- Are there any questions?



Backup Slides

### Links and References



- CMS Online Software Support
  - Trac <u>https://svnweb.cern.ch/trac/cmsos</u>
  - Mailing list (e-group) <u>cms-os-users@cern.ch</u>
  - User's Support Guide <u>https://edms.cern.ch/file/1001791/</u>
  - Documentation available on EDMS
    - https://edms.cern.ch/nav/CMS-00000001/CMS-000009924

#### Links to additional documentation and tools

- Book "Version Control with Subversion" <u>http://svnbook.red-bean.com/</u>
- Subversion <u>http://subversion.tigris.org/</u>
- Trac <u>http://trac.edgewall.org/</u>