

Overview of InfiniBand Architecture

Dhabaleswar K. (DK) Panda
The Ohio State University
E-mail: panda@cse.ohio-state.edu
<http://www.cse.ohio-state.edu/~panda>

Processing Bottlenecks in Traditional Protocols

- Ex: TCP/IP, UDP/IP
- Generic architecture for all network interfaces
- Host-handles almost all aspects of communication
 - Data buffering (copies on sender and receiver)
 - Data integrity (checksum)
 - Routing aspects (IP routing)
- Signaling between different layers
 - Hardware interrupt whenever a packet arrives or is sent
 - Software signals between different layers to handle protocol processing in different priority levels

Capabilities of High-Performance Networks

- Intelligent Network Interface Cards
- Support entire protocol processing completely in hardware (hardware protocol offload engines)
- Provide a rich communication interface to applications
 - **User-level communication capability**
 - Gets rid of intermediate data buffering requirements
- No software signaling between communication layers
 - All layers are implemented on a **dedicated** hardware unit, and not on a **shared** host CPU

Previous High-Performance Network Stacks

- Virtual Interface Architecture (VIA)
 - Standardized by Intel, Compaq, Microsoft
- Fast Messages (FM)
 - Developed by UIUC
- Myricom GM
 - Proprietary protocol stack from Myricom
- These network stacks set the trend for high-performance communication requirements
 - Hardware offloaded protocol stack
 - Support for fast and secure user-level access to the protocol stack

IB Trade Association

- IB Trade Association was formed with seven industry leaders (Compaq, Dell, HP, IBM, Intel, Microsoft, and Sun)
- **Goal: To design a scalable and high performance communication and I/O architecture by taking an integrated view of computing, networking, and storage technologies**
- Many other industry participated in the effort to define the IB architecture specification
- IB Architecture (Volume 1, Version 1.0) was released to public on Oct 24, 2000
 - Latest version 1.2.1 released January 2008
- <http://www.infinibandta.org>

IB Hardware Acceleration

- Some IB models have multiple hardware accelerators
 - E.g., Mellanox IB adapters
- Protocol Offload Engines
 - Completely implement layers 2-4 in hardware
- Additional hardware supported features also present
 - RDMA, Multicast, QoS, Fault Tolerance, and many more

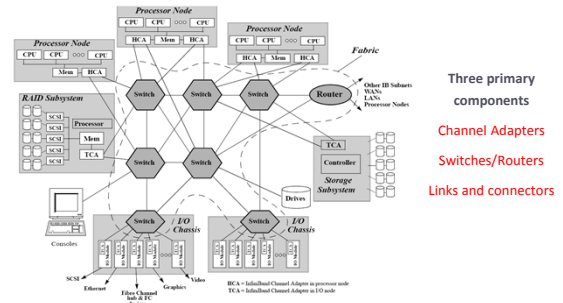
IB Overview

- **InfiniBand**
 - **Architecture and Basic Hardware Components**
 - **Communication Model and Semantics**
 - Communication Model
 - Memory registration and protection
 - Channel and memory semantics
 - **Novel Features**
 - Hardware Protocol Offload
 - Link, network and transport layer features
 - **Management and Services**
 - Subnet Management
 - Hardware support for scalable network management

HPCA '10

7

A Typical IB Network

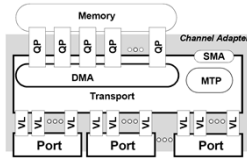


HPCA '10

8

Components: Channel Adapters

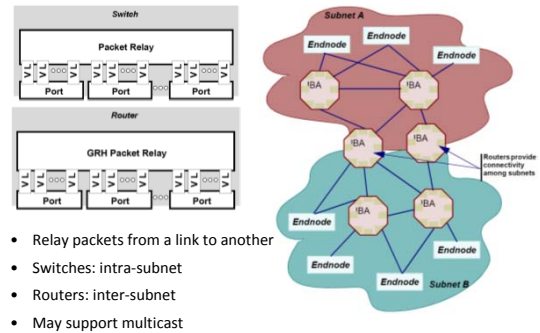
- Used by processing and I/O units to connect to fabric
- Consume & generate IB packets
- Programmable DMA engines with protection features
- May have multiple ports
 - Independent buffering channeled through Virtual Lanes
- Host Channel Adapters (HCAs)



HPCA '10

9

Components: Switches and Routers



HPCA '10

10

Components: Links & Repeaters

- Network Links
 - Copper, Optical, Printed Circuit wiring on Back Plane
 - Not directly addressable
- Traditional adapters built for copper cabling
 - Restricted by cable length (signal integrity)
 - For example, QDR copper cables are restricted to 7m
- Intel Connects: Optical cables with Copper-to-optical conversion hubs (acquired by Emcore)
 - Up to 100m length
 - 550 picoseconds copper-to-optical conversion latency
- Available from other vendors (Luxtera)
- Repeaters (Vol. 2 of InfiniBand specification)



(Courtesy Intel)

HPCA '10

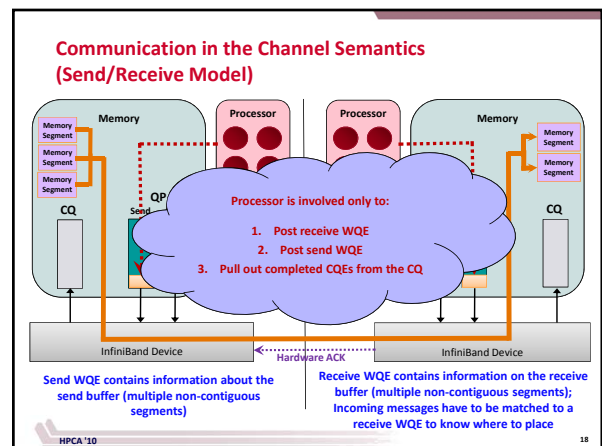
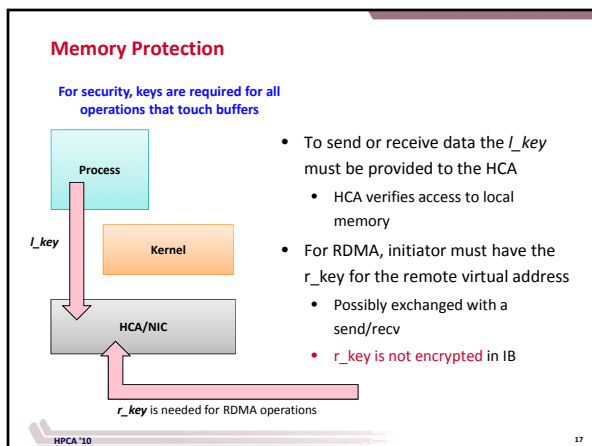
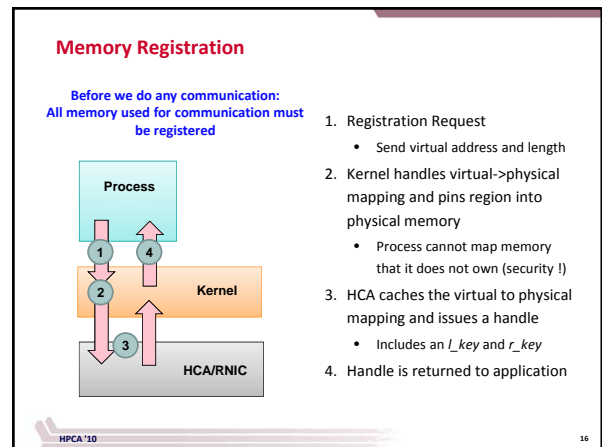
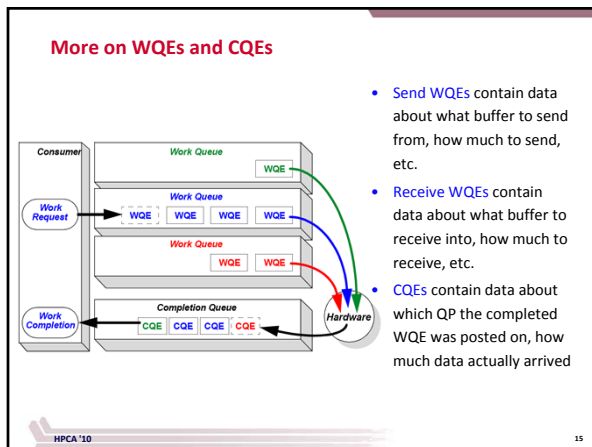
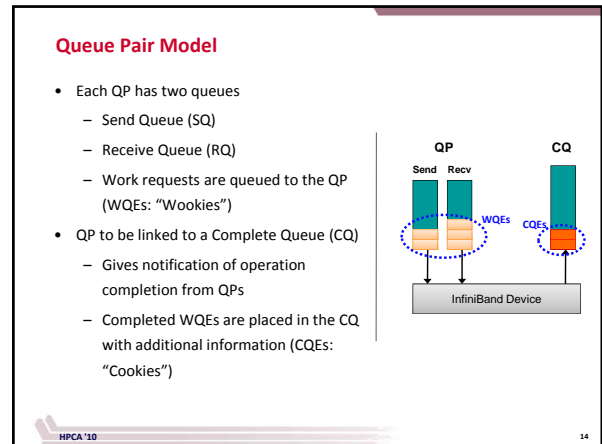
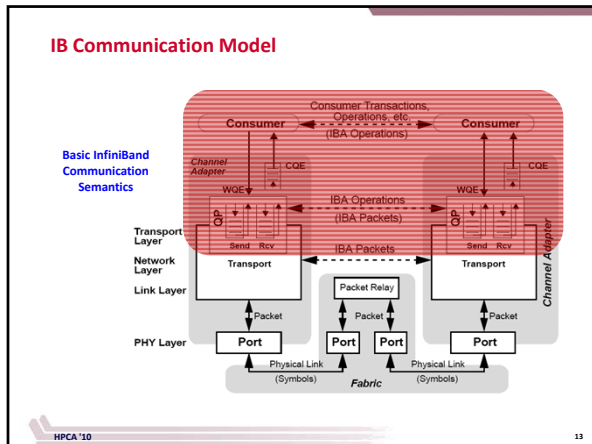
11

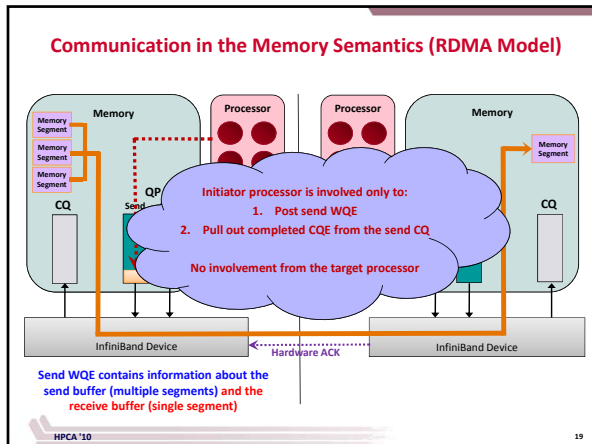
IB Overview

- **InfiniBand**
 - **Architecture and Basic Hardware Components**
 - **Communication Model and Semantics**
 - **Communication Model**
 - **Memory registration and protection**
 - **Channel and memory semantics**
 - **Novel Features**
 - Hardware Protocol Offload
 - Link, network and transport layer features
 - **Management and Services**
 - Subnet Management
 - Hardware support for scalable network management

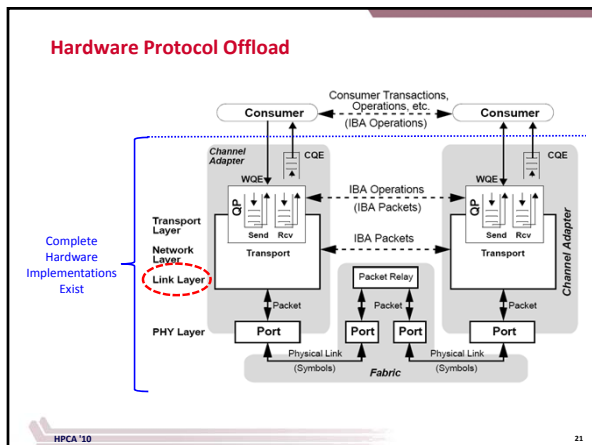
HPCA '10

12





- ### IB Overview
- **InfiniBand**
 - Architecture and Basic Hardware Components
 - Communication Model and Semantics
 - Communication Model
 - Memory registration and protection
 - Channel and memory semantics
 - **Novel Features**
 - **Hardware Protocol Offload**
 - Link, network and transport layer features
 - Management and Services
 - Subnet Management
 - Hardware support for scalable network management
- HPCA '10 20



- ### Link Layer Capabilities
- **CRC-based Data Integrity**
 - **Buffering and Flow Control**
 - Virtual Lanes, Service Levels and QoS
 - Switching and Multicast
 - IB WAN Capability
- HPCA '10 22

- ### CRC-based Data Integrity
- Two forms of CRC to achieve both early error detection and end-to-end reliability
 - Invariant CRC (ICRC) covers fields that do not change per link (per network hop)
 - E.g., routing headers (if there are no routers), transport headers, data payload
 - 32-bit CRC (compatible with Ethernet CRC)
 - End-to-end reliability (does not include I/O bus)
 - Variant CRC (VCRC) covers everything
 - 16-bit CRC
 - Erroneous packets do not have to reach the destination
 - Early error detection
- HPCA '10 23

- ### Buffering and Flow Control
- IB provides an absolute credit-based flow-control
 - Receiver guarantees that it has enough space allotted for N blocks of data
 - Occasional update of available credits by the receiver
 - Has no relation to the number of messages, but only to the total amount of data being sent
 - One 1MB message is equivalent to 1024 1KB messages (except for rounding off at message boundaries)
- HPCA '10 24

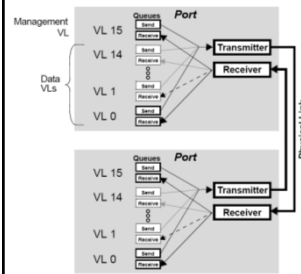
Link Layer Capabilities

- CRC-based Data Integrity
- Buffering and Flow Control
- **Virtual Lanes, Service Levels and QoS**
- **Switching and Multicast**
- **IB WAN Capability**

HPCA '10

25

Virtual Lanes



- Multiple virtual links within same physical link
 - Between 2 and 16
- Separate buffers and flow control
 - Avoids Head-of-Line Blocking
- VL15: reserved for management
- Each port supports one or more data VL

HPCA '10

26

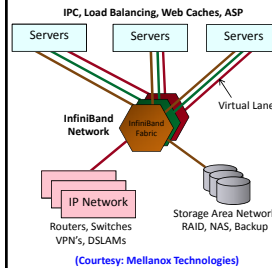
Service Levels and QoS

- Service Level (SL):
 - Packets may operate at one of 16 different SLs
 - Meaning not defined by IB
- SL to VL mapping:
 - SL determines which VL on the next link is to be used
 - Each port (switches, routers, end nodes) has a SL to VL mapping table configured by the subnet management
- Partitions:
 - Fabric administration (through Subnet Manager) may assign specific SLs to different partitions to isolate traffic flows

HPCA '10

27

Traffic Segregation Benefits



- InfiniBand Virtual Lanes allow the multiplexing of multiple independent logical traffic flows on the same physical link
- Providing the benefits of independent, separate networks while eliminating the cost and difficulties associated with maintaining two or more networks

HPCA '10

28

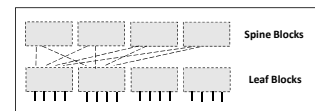
Switching (Layer-2 Routing) and Multicast

- Each port has one or more associated LIDs (Local Identifiers)
 - Switches look up which port to forward a packet to based on its destination LID (DLID)
 - This information is maintained at the switch
- For multicast packets, the switch needs to maintain multiple output ports to forward the packet to
 - Packet is replicated to each appropriate output port
 - Ensures at-most once delivery & loop-free forwarding
 - There is an interface for a group management protocol
 - Create, join/leave, prune, delete group

HPCA '10

29

Destination-based Switching/Routing



An Example IB Switch Block Diagram (Mellanox 144-Port)

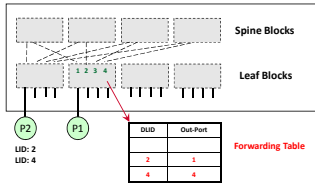
Switching: IB supports Virtual Cut Through (VCT)
 Routing: Unspecified by IB SPEC
 Up*/Down*, Shift are popular routing engines supported by OFED

- Fat-Tree is a popular topology for IB Clusters
- Different over-subscription ratio may be used

HPCA '10

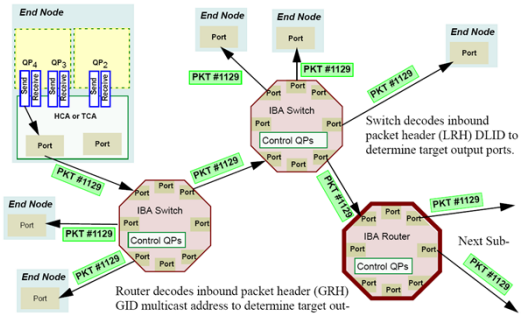
30

IB Switching/Routing: An Example



- Someone has to setup these tables and give every port an LID
 - “Subnet Manager” does this work (more discussion on this later)
- Different routing algorithms may give different paths

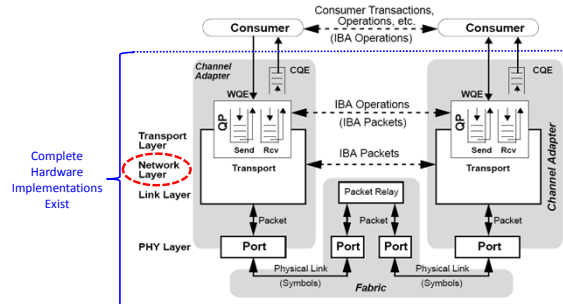
IB Multicast Example



IB WAN Capability

- Getting increased attention for:
 - Remote Storage, Remote Visualization
 - Cluster Aggregation (Cluster-of-clusters)
- IB-Optical switches by multiple vendors
 - Obsidian Research Corporation: www.obsidianresearch.com
 - Network Equipment Technology (NET): www.net.com
 - Layer-1 changes from copper to optical; everything else stays the same
 - Low-latency copper-optical-copper conversion
- Large link-level buffers for flow-control
 - Data messages do not have to wait for round-trip hops
 - Important in the wide-area network

Hardware Protocol Offload

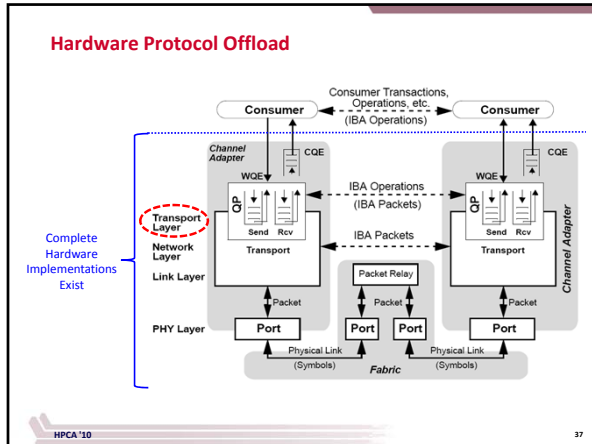


IB Network Layer Capabilities

- Most capabilities are similar to that of the link layer, but as applied to IB routers
 - Routers can send packets across subnets (subnet are management domains, not administrative domains)
 - Subnet management packets are consumed by routers, not forwarded to the next subnet
- Several additional features as well
 - E.g., routing and flow labels

Routing and Flow Labels

- Routing follows the IPv6 packet format
 - Easy interoperability with Wide-area translations
 - Link layer might still need to be translated to the appropriate layer-2 protocol (e.g., Ethernet, SONET)
- Flow Labels allow routers to specify which packets belong to the same connection
 - Switches can optimize communication by sending packets with the same label in order
 - Flow labels can change in the router, but packets belonging to one label will always do so



IB Transport Services

Service Type	Connection Oriented	Acknowledged	Transport
Reliable Connection	Yes	Yes	IBA
Unreliable Connection	Yes	No	IBA
Reliable Datagram	No	Yes	IBA
Unreliable Datagram	No	No	IBA
RAW Datagram	No	No	Raw

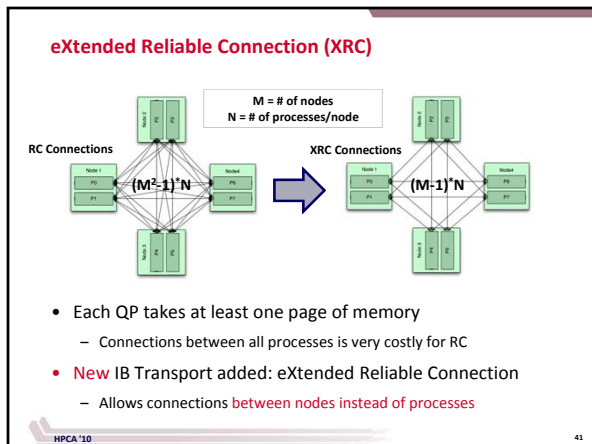
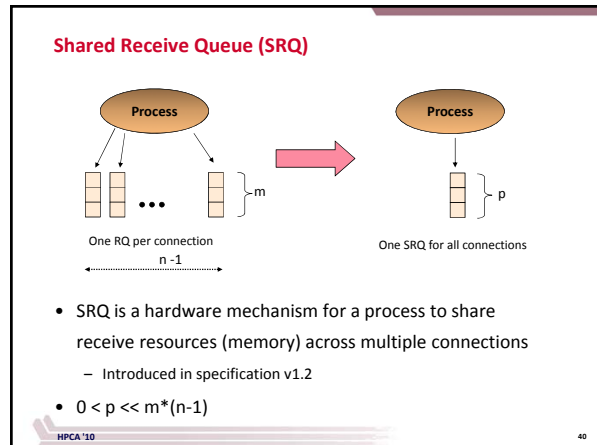
- Each transport service can have zero or more QPs associated with it
 - E.g., you can have four QPs based on RC and one QP based on UD

HPCA '10 38

Trade-offs in Different Transport Types

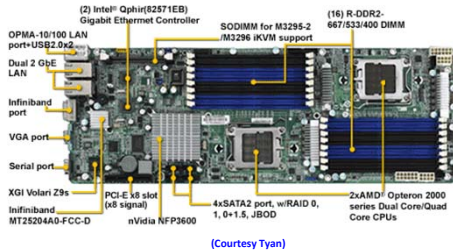
Attribute	Reliable Connection	Reliable Datagram	Unreliable Datagram	Unreliable Connection	Raw Datagram (both IPv6 & ethernet)
Scalability (M processes on N Processor nodes communicating with all processes on all nodes)	M*N QPs required on each processor node, per CA	M QPs required on each processor node, per CA	M QPs required on each processor node, per CA	M*N QPs required on each processor node, per CA	1 QP required on each end node, per CA
Corrupt data detected	Yes	Yes	Yes	No	No
Data delivery guarantee	Data delivered exactly once	Yes, per connection	Yes, packets from any one source QP are ordered to multiple destination QPs	No	No
Data order guaranteed	Yes, per connection	Yes, packets from any one source QP are ordered to multiple destination QPs	No	Unordered and duplicate packets are detected.	No
Data loss detected	Yes	Yes	No	Yes	No
Error recovery	Reliable: Errors are detected at both the requestor and the responder. The requestor can transparently recover from errors (retransmission, alternate path, etc.) without any involvement of the client application. QP processing is halted only if the destination is inoperable or all fabric paths between the channel adapters have failed.	Unreliable: Packets with errors are detected at both the requestor and the responder. The requestor is not informed.	Unreliable: Packets with errors, including sequence errors, are detected and may be logged by the responder. The requestor is not informed.	Unreliable: Packets with errors are not delivered. The requestor and responder are not informed of dropped packets.	Unreliable: Packets with errors are not delivered. The requestor and responder are not informed of dropped packets.

HPCA '10 39



- ### IB Hardware Products
- Many IB vendors: Mellanox, Voltaire and Qlogic
 - Aligned with many server vendors: Intel, IBM, SUN, Dell
 - And many integrators: Appro, Advanced Clustering, Microway
 - Broadly two kinds of adapters
 - Offloading (Mellanox) and Onloading (Qlogic)
 - Adapters with different interfaces:
 - Dual port 4X with PCI-X (64 bit/133 MHz), PCIe x8, PCIe 2.0 and HT
 - MemFree Adapter
 - No memory on HCA → Uses System memory (through PCIe)
 - Good for LOM designs (Tyan S2935, Supermicro 6015T-INFB)
 - Different speeds
 - SDR (8 Gbps), DDR (16 Gbps) and QDR (32 Gbps)
 - Some 12X SDR adapters exist as well (24 Gbps each way)
- HPCA '10 42

Tyan Thunder S2935 Board



(Courtesy Tyan)

Similar boards from Supermicro with LOM features are also available

HPCA '10

43

IB Hardware Products (contd.)

- Customized adapters to work with IB switches
 - Cray XD1 (formerly by Octigabay), Cray CX1
- Switches:
 - 4X SDR and DDR (8-288 ports); 12X SDR (small sizes)
 - 3456-port “Magnum” switch from SUN → used at TACC
 - 72-port “nano magnum”
 - 36-port Mellanox InfiniScale IV QDR switch silicon in early 2008
 - Up to 648-port QDR switch by SUN
 - New IB switch silicon from Qlogic introduced at SC '08
 - Up to 846-port QDR switch by Qlogic
- Switch Routers with Gateways
 - IB-to-FC; IB-to-IP

HPCA '10

44

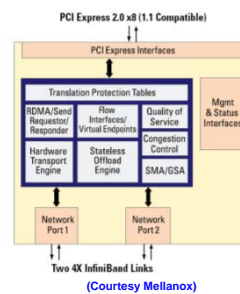
IB Software Products

- Low-level software stacks
 - VAPI (Verbs-Level API) from Mellanox
 - Modified and customized VAPI from other vendors
 - New initiative: Open Fabrics (formerly OpenIB)
 - <http://www.openfabrics.org>
 - Open-source code available with Linux distributions
 - Initially IB; later extended to incorporate iWARP
- High-level software stacks
 - MPI, SDP, IPoIB, SRP, iSER, DAPL, NFS, PVFS on various stacks (primarily VAPI and OpenFabrics)

HPCA '10

45

Mellanox ConnectX Architecture



- Early adapter supporting IB/10GE convergence
 - Support for VPI and IBoE
- Includes other features as well
 - Hardware support for Virtualization
 - Quality of Service
 - Stateless Offloads

HPCA '10

46

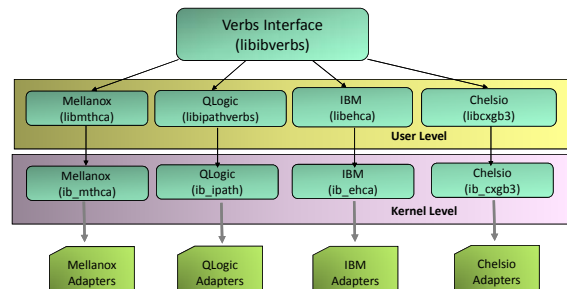
OpenFabrics

- Open source organization (formerly OpenIB)
 - www.openfabrics.org
- Incorporates both IB and iWARP in a unified manner
 - Support for Linux and Windows
 - Design of complete stack with ‘best of breed’ components
 - Gen1
 - Gen2 (current focus)
- Users can download the entire stack and run
 - Latest release is OFED 1.4.3
 - OFED 1.5 is underway

HPCA '10

47

OpenFabrics Stack with Unified Verbs Interface



HPCA '10

48

