



Education

InfiniBand Technology Overview

Dror Goldenberg, Mellanox Technologies

- The material contained in this tutorial is copyrighted by the SNIA.
- Member companies and individuals may use this material in presentations and literature under the following conditions:
 - ◆ Any slide or slides used must be reproduced without modification
 - ◆ The SNIA must be acknowledged as source of any material used in the body of any document containing material from these presentations.
- This presentation is a project of the SNIA Education Committee.

InfiniBand Technology Overview

The InfiniBand architecture brings fabric consolidation to the data center. Storage networking can concurrently run with clustering, communication and management fabrics over the same infrastructure, preserving the behavior of multiple fabrics. The tutorial provides an overview of the InfiniBand architecture including discussion of High Speed – Low Latency, Channel I/O, QoS scheduling, partitioning, high availability and protocol offload. InfiniBand based storage protocols, iSER (iSCSI RDMA Protocol), NFS over RDMA and SCSI RDMA Protocol (SRP), are introduced and compared with alternative storage protocols, such as iSCSI and FCP. The tutorial further enumerates value-add features that the InfiniBand brings to clustered storage, such as atomic operations and end to end data integrity.

Learning Objectives:

- Understand the InfiniBand architecture and feature set.
- Understand the benefits of InfiniBand for networked storage.
- Understand the standard InfiniBand storage protocols.

- Motivation and General Overview
- Protocol Stack Layers
- Storage Protocols over InfiniBand
- Benefits

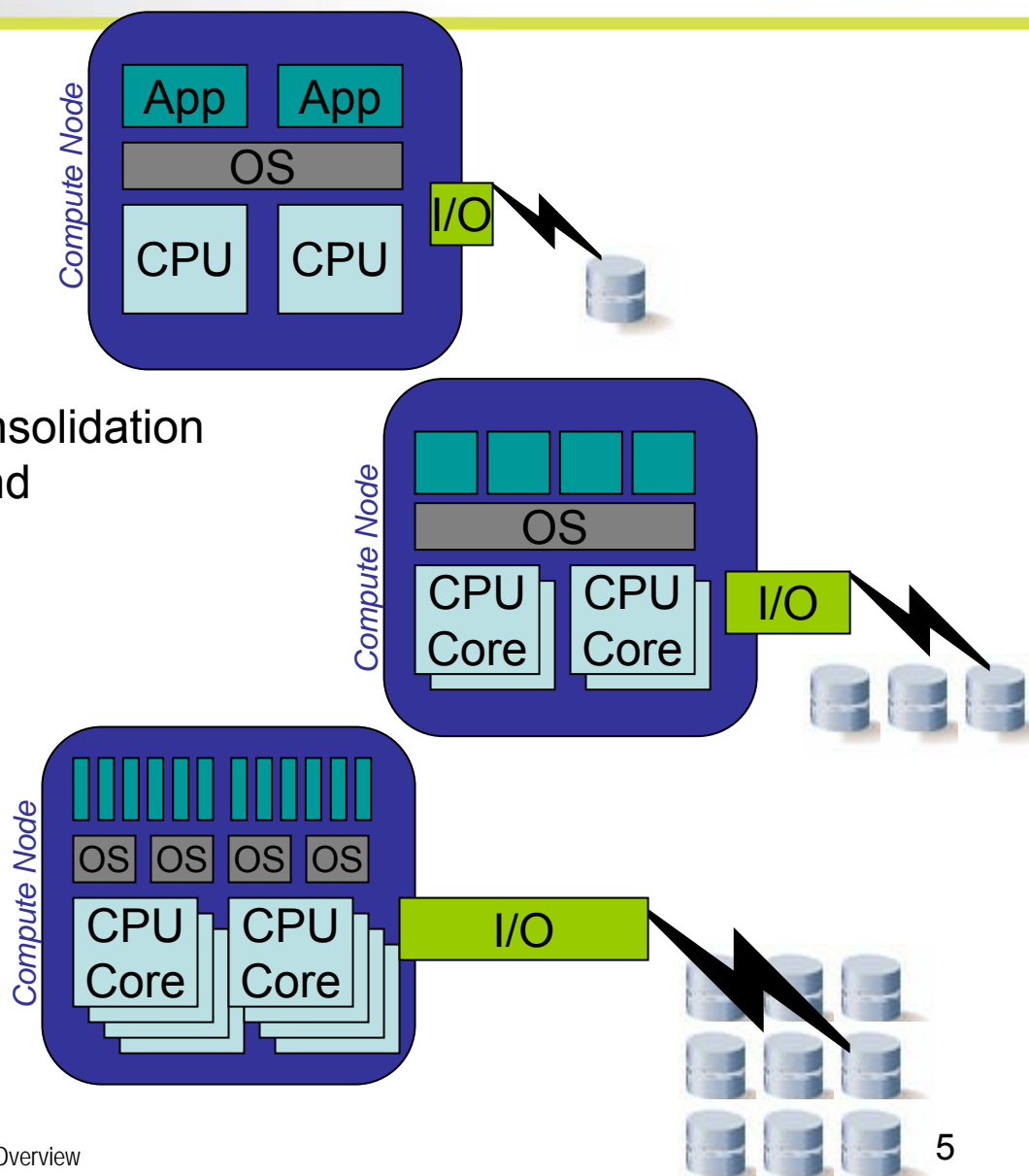
The Need for Better I/O

➤ Datacenter trends

- ◆ Multi-core CPUs
- ◆ Bladed architecture
- ◆ Fabric consolidation
- ◆ Server virtualization & consolidation
- ◆ Increasing storage demand

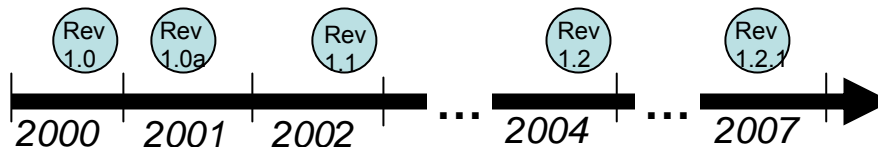
➤ Better I/O is required

- ◆ High capacity
- ◆ Efficient
 - Low latency
 - CPU Offload
- ◆ Scalable
- ◆ Virtualization friendly
- ◆ High availability
- ◆ Performance
- ◆ Low power
- ◆ TCO reduction

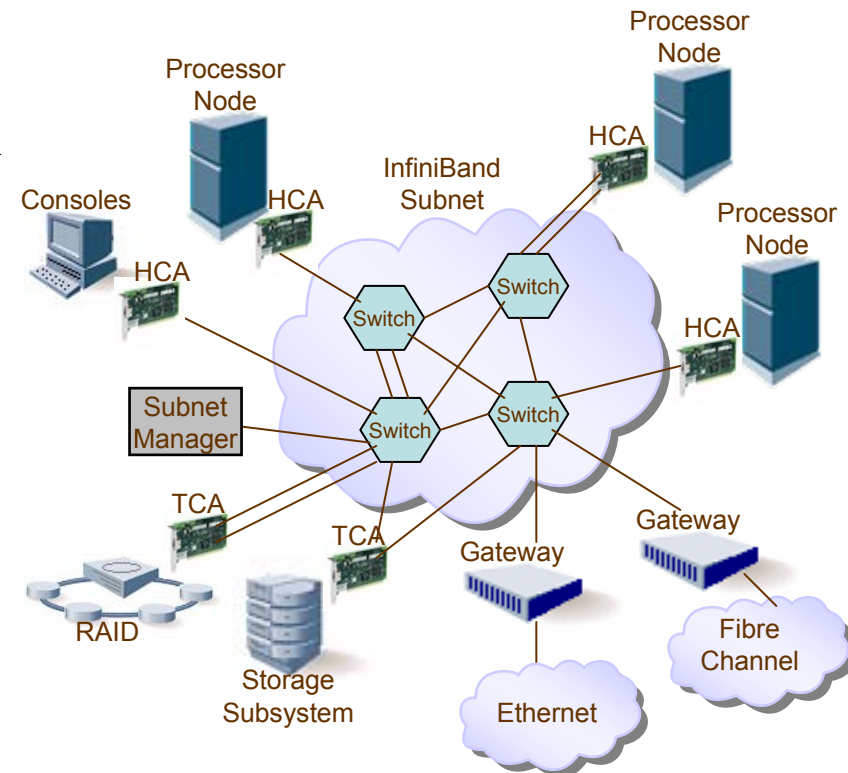


The InfiniBand Architecture

- Industry standard defined by the InfiniBand Trade Association
- Defines System Area Network architecture
 - ◆ Comprehensive specification:
from physical to applications



- Architecture supports
 - ◆ Host Channel Adapters (HCA)
 - ◆ Target Channel Adapters (TCA)
 - ◆ Switches
 - ◆ Routers
- Facilitated HW design for
 - ◆ Low latency / high bandwidth
 - ◆ Transport offload



A Comparison of Fabric Technologies

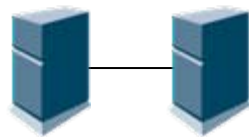
SNIA

Features and Price	Fibre Channel	Standard 10 GbE	InfiniBand
Bandwidth	4Gb/s (4GFC) 8Gb/s (8GFC)	10Gb/s	20Gb/s (4x DDR)
Raw Bandwidth (unidirectional)	400MB/s (4GFC) 800MB/s (8GFC)	1,250 MB/s	2,000 MB/s* (4x DDR) 4,000 MB/s (4x QDR)
Reliable Service	Yes	No	Yes
Fabric Consolidation	Practically no	Practically partial**	Yes
Copper Distance	15m	10GBase-CX4 15m 10GBase-T 100m	Passive SDR 20m/ DDR 10m Active DDR 25m
Optical Distance	100m	10GBase-SR 300m 10GBase-LRM 220m	300m (SDR) 150m (DDR)

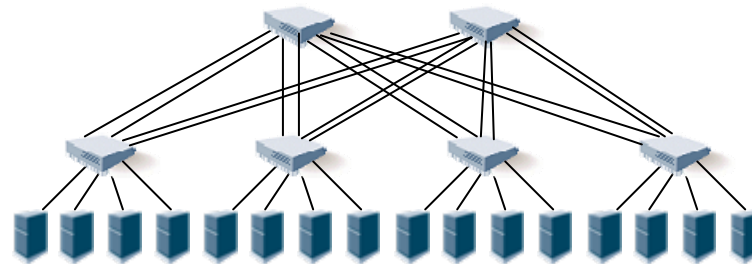
* 1,940 MB/s measured



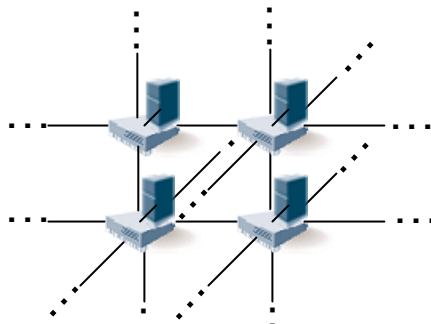
InfiniBand Topologies



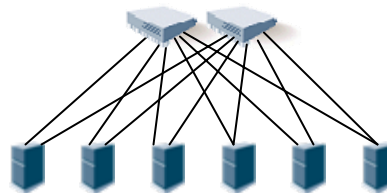
Back to Back



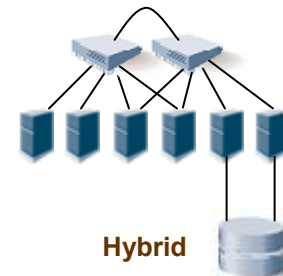
2 Level Fat Tree



3D Torus



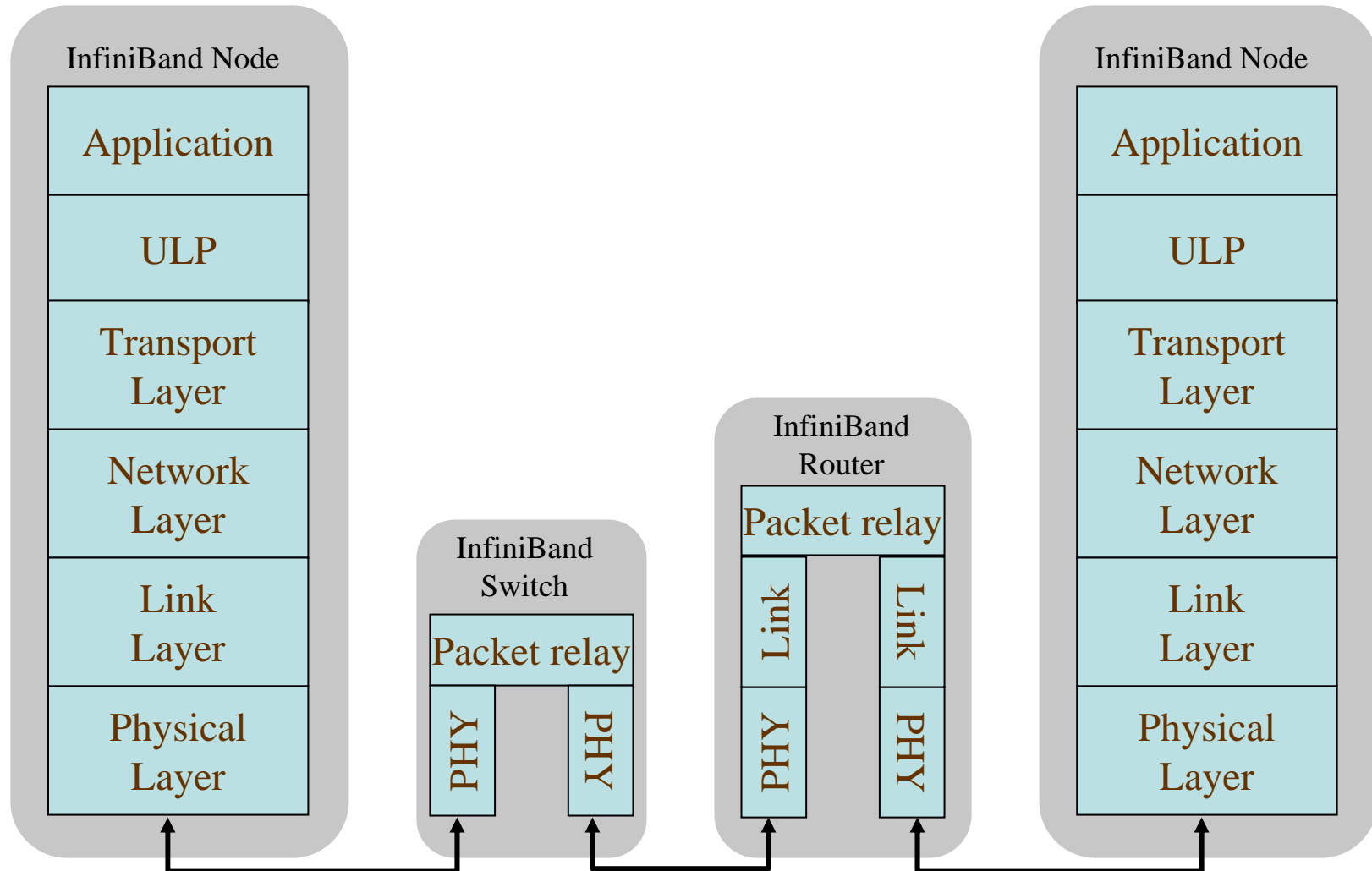
Dual Star



Hybrid

- Example topologies commonly used
- Architecture does not limit topology
- Modular switches are based on fat tree architecture

InfiniBand Protocol Layers



Physical Layer

- Width (1X, 4X, 8X, 12X) including auto-negotiation
- Speed (SDR/DDR/QDR) including auto-negotiation
 - ◆ 4X DDR HCAs are currently shipping
- Power management
 - ◆ Polling / Sleeping
- Connector
 - ◆ Board: MicroGiGaCN
 - ◆ Pluggable: QSFP
- 8/10 encoding
 - ◆ Maintain DC Balance
 - ◆ Limited run length of 0's or 1's
- Control symbols (Kxx.x)
 - ◆ Lane de-skew, auto negotiation, training, clock tolerance, framing

Link Speed (10⁹ bit/sec)

Lane Speed →	SDR (2.5GHz)	DDR (5GHz)	QDR (10GHz)
Link Width ↓			
1X	2.5	5	10
4X	10	20	40
8X	20	40	80
12X	30	60	120

* MicroGiGaCN is a trademark of Fujitsu Components Limited

Physical Layer – Cont'd

Copper Cables*:

Width	Speed	Connector	Min Reach	Type / Power
4X	SDR/DDR	Micro-GiGaCN	20m/10m	Passive
4X	DDR	Micro-GiGaCN	15-25m	Active 0.5-1.5W
12X	SDR/DDR	24pin Micro-GiGaCN	20m/10m	Passive

4X –
MicroGiGaCN→



12X –
24 pair MicroGiGaCN→



Fiber Optics*:

Width	Speed	Connector	Type	Min Reach	Power	Fiber Media
4X	SDR/DDR	Micro-GiGaCN	Media Converter	300m/150m	0.8-1W	12 strand MPO
4X	DDR	Micro-GiGaCN	Optical Cable	100m	1W	12 strand attached

4X - MicroGiGaCN
MPO Media Converter →



4X - MicroGiGaCN
Optical Cable →



* currently deployed

➤ Addressing and Switching

- ◆ Local Identifier (LID) addressing
- ◆ Unicast LID - 48K addresses
- ◆ Multicast LID – up to 16K addresses
- ◆ Efficient linear lookup
- ◆ Cut through switching supported
- ◆ Multi-pathing support through LMC

➤ Independent Virtual Lanes

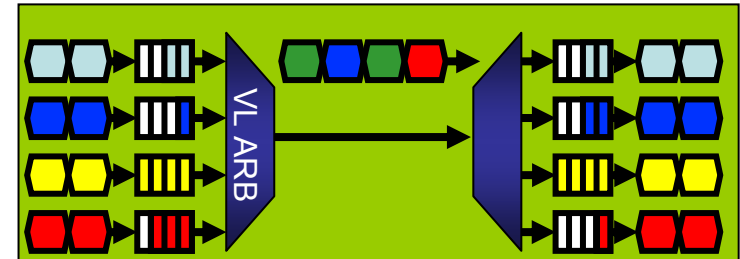
- ◆ Flow control (lossless fabric)
- ◆ Service level
- ◆ VL arbitration for QoS

➤ Congestion control

- ◆ Forward / Backward Explicit Congestion Notification (FECN/BECN)

➤ Data Integrity

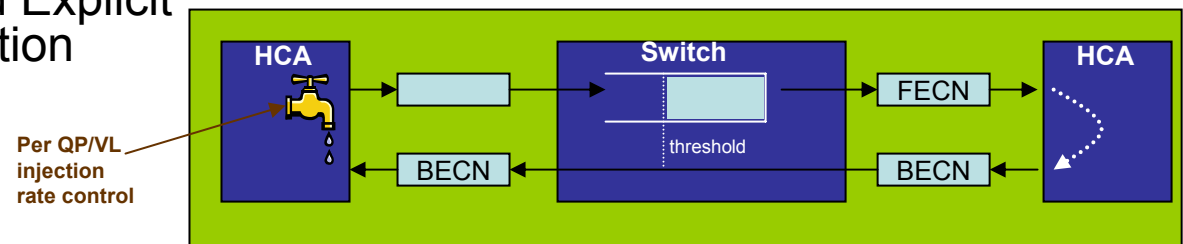
- ◆ Invariant CRC
- ◆ Variant CRC



Independent Virtual Lanes (VLs)



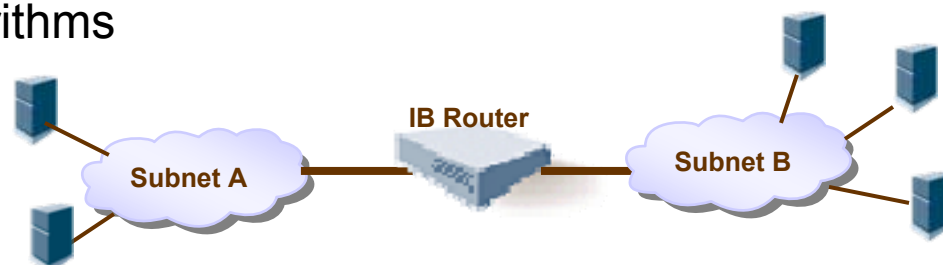
H/L Weighted Round Robin (WRR) VL Arbitration



Efficient FECN/BECN Based Congestion Control

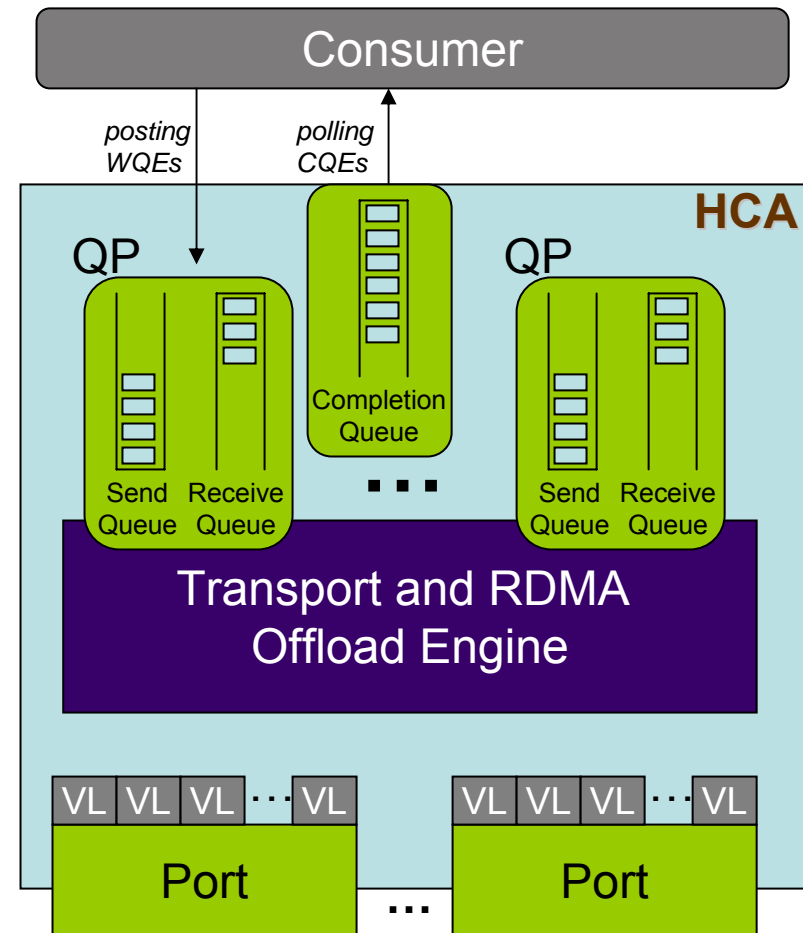
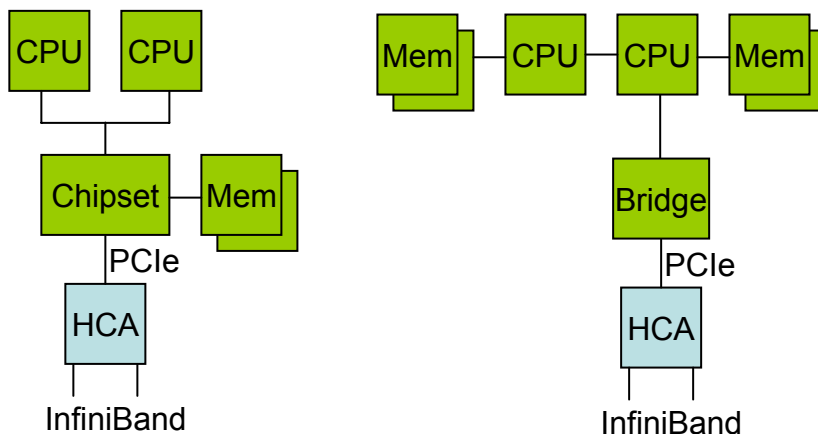
Network Layer

- Global Identifier (GID) addressing
 - ◆ Based on IPv6 addressing scheme
 - ◆ $GID = \{64 \text{ bit GID prefix, } 64 \text{ bit GUID}\}$
 - GUID = Global Unique Identifier (64 bit EUI-64)
 - GUID 0 – assigned by the manufacturer
 - GUID 1..(N-1) – assigned by the Subnet Manager
- Optional for local subnet access
- Used for multicast distribution within end nodes
- Enables routing between IB subnets
 - ◆ Still under definition in IBTA
 - ◆ Will leverage IPv6 routing algorithms



Transport - Host Channel Adapter Model

- Asynchronous interface
 - ◆ Consumer posts work requests
 - ◆ HCA processes
 - ◆ Consumer polls completions
- Transport executed by HCA
- I/O channel exposed to the application
- Transport services
 - ◆ Reliable / Unreliable
 - ◆ Connected / Datagram

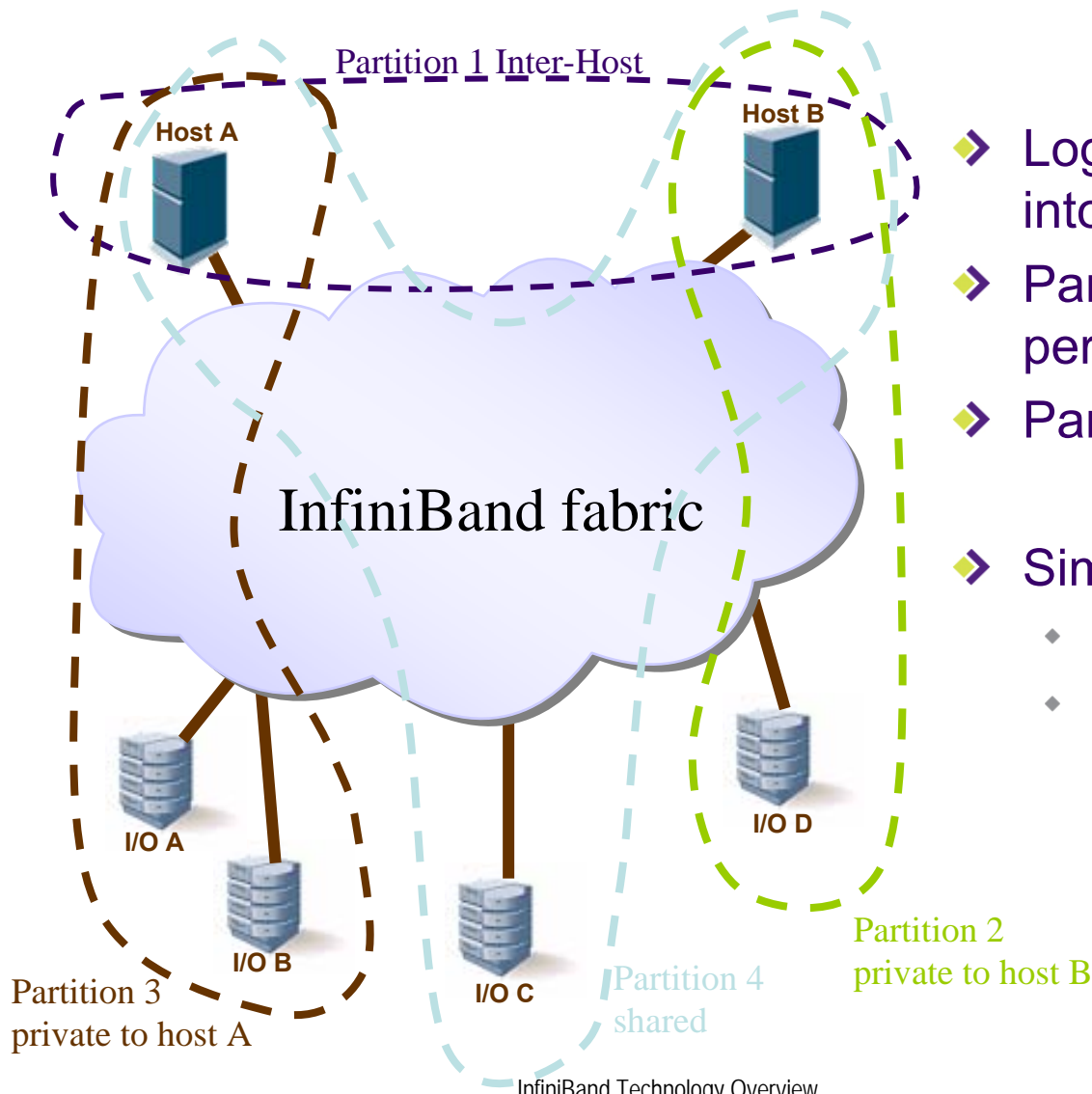


➤ Queue Pair (QP) – transport endpoint

- ◆ Asynchronous interface
 - Send Queue, Receive Queue, Completion Queue
- ◆ Full transport offload
 - Segmentation, reassembly, timers, retransmission, etc
- ◆ Operations supported
 - Send/Receive – messaging semantics
 - RDMA Read/Write – enable zero copy operations
 - Atomics – remote Compare & Swap, Fetch & Add
 - Memory management - Bind/Fast Register/Invalidate

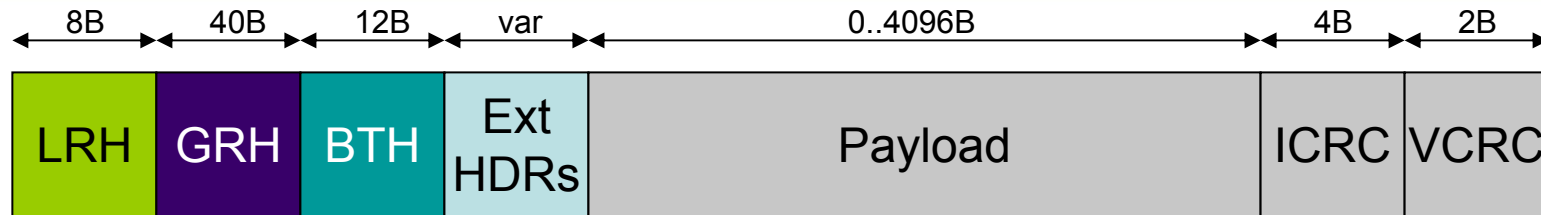
➤ Kernel bypass

- ◆ Enables low latency and CPU offload
- ◆ Enabled through QPs, Completion Queues (CQs), Protection Domains (PD), Memory Regions (MRs)



- Logically divide the fabric into isolated domains
- Partial and full membership per partition
- Partition filtering at switches
- Similar to
 - ◆ FC Zoning
 - ◆ 802.1Q VLANs

InfiniBand Packet Format



InfiniBand Data Packet

VL	LVer	SL	rsvd	LNH	DLID
rsvd	Len			SLID	

LRH

Opcode	SM	Pad	TVer	Partition Key
rsvd	Destination QP			
A	rsvd	PSN		

BTH

IPVer	TClass	Flow Label	
Payload Len		Next Header	Hop Lim
SGID[127:96]			
SGID[95:64]			
SGID[63:32]			
SGID[31:0]			
DGID[127:96]			
DGID[95:64]			
DGID[63:32]			
DGID[31:0]			

GRH (Optional)

Extended headers:

- Reliable Datagram ETH (4B)
- Datagram ETH (8B)
- RDMA ETH (16B)
- Atomic ETH (28B)
- ACK ETH (4B)
- Atomic ACK ETH (8B)
- Immediate Data ETH (4B)
- Invalidate ETH (4B)

➤ Hop by hop

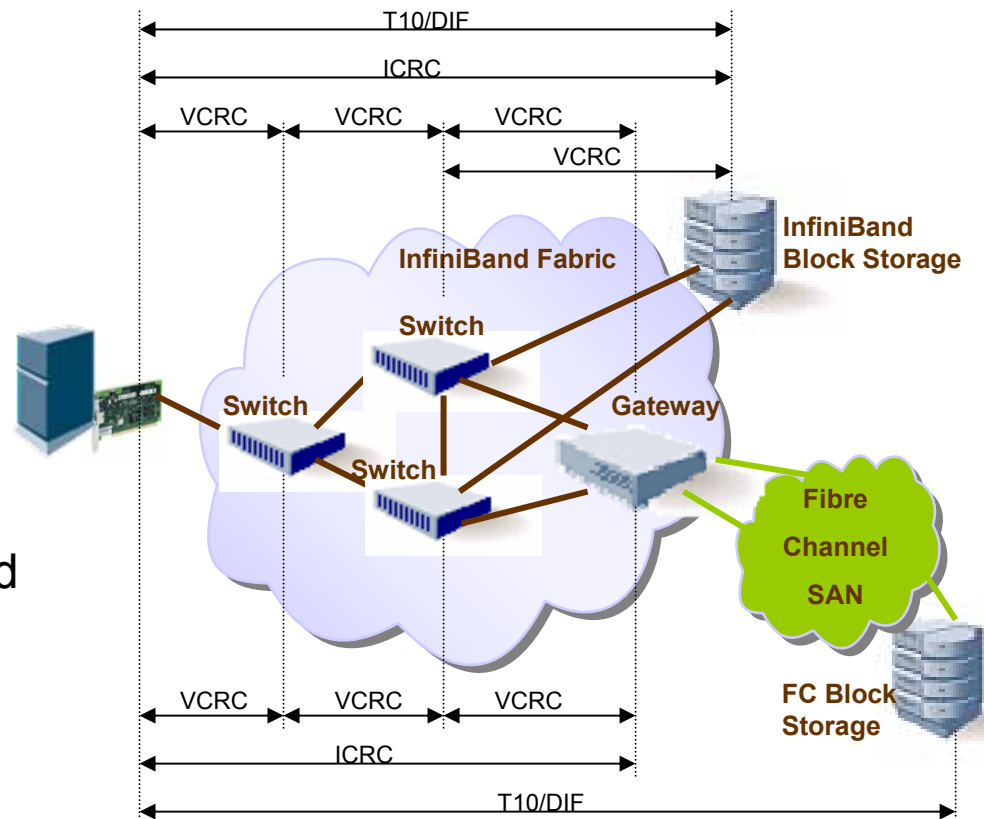
- ♦ VCRC – 16 bit CRC
- ♦ CRC16 0x100B

➤ End to end

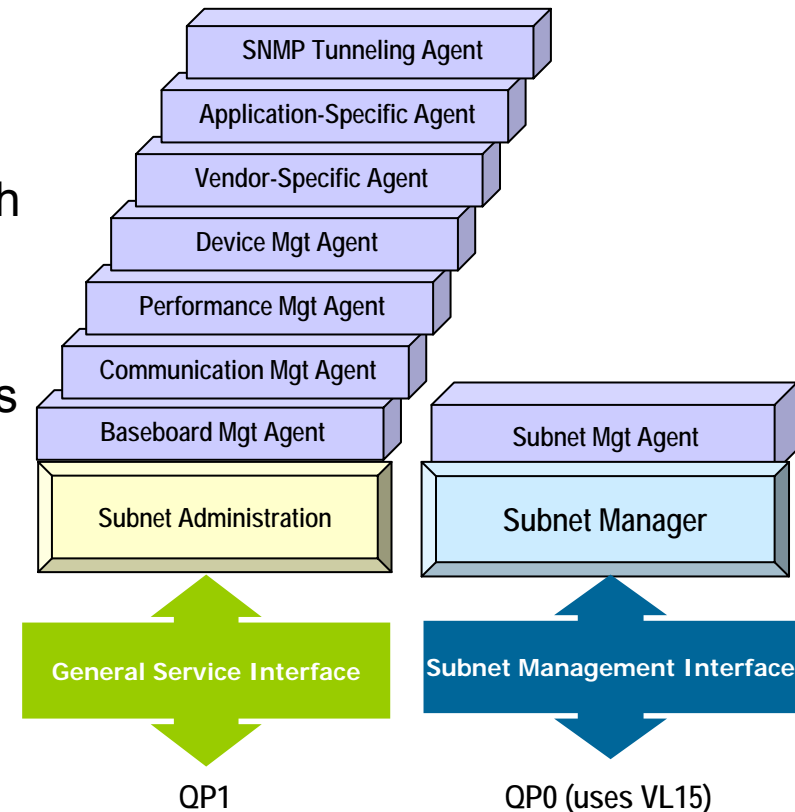
- ♦ ICRC – 32 bit CRC
- ♦ CRC32 0x04C11DB7
- ♦ Same CRC as Ethernet

➤ Application level

- ♦ T10/DIF Logical Block Guard
 - Per block CRC
- ♦ 16 bit CRC 0x8BB7



- **Subnet Manager (SM)**
 - ◆ Configures/Administers fabric topology
 - ◆ Implemented at an end-node or a switch
 - ◆ Active/Passive model when more than one SM is present
 - ◆ Talks with SM Agents in nodes/switches
- **Subnet Administration**
 - ◆ Provides path records
 - ◆ QoS management
- **Communication Management**
 - ◆ Connection establishment processing



Upper Layer Protocols

- ULPs connect InfiniBand to common interfaces
- Supported on mainstream operating systems

- **Clustering**

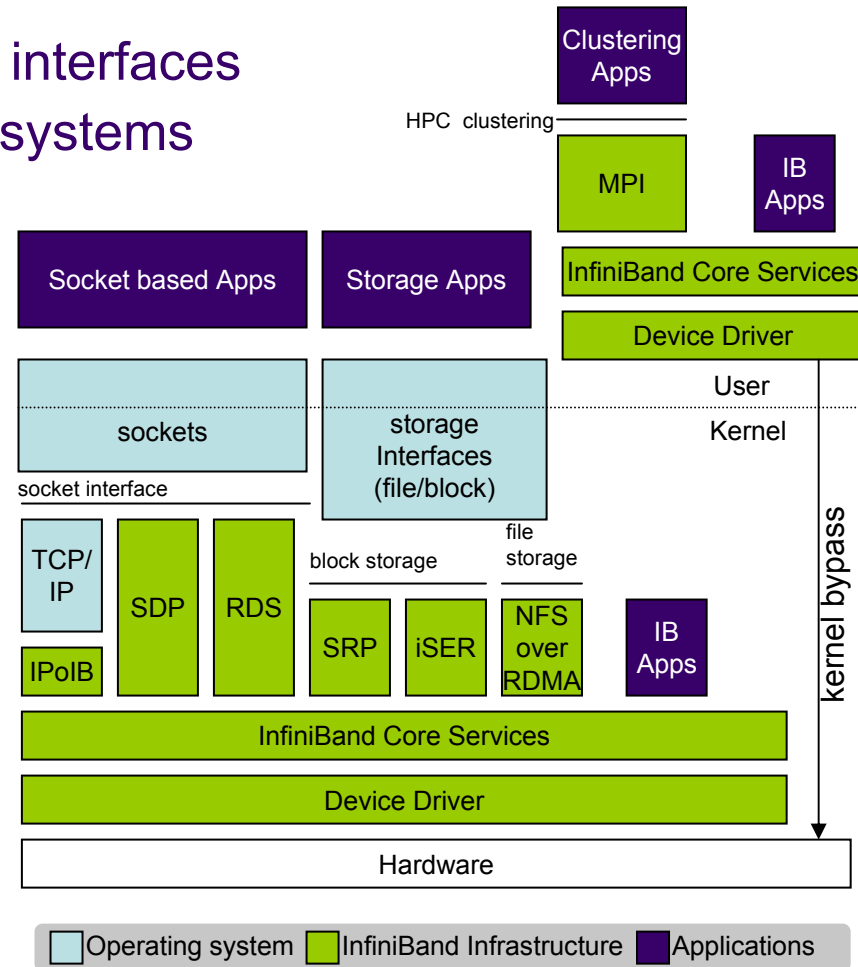
- ♦ MPI (Message Passing Interface)
- ♦ RDS (Reliable Datagram Socket)

- **Network**

- ♦ IPoIB (IP over InfiniBand)
- ♦ SDP (Socket Direct Protocol)

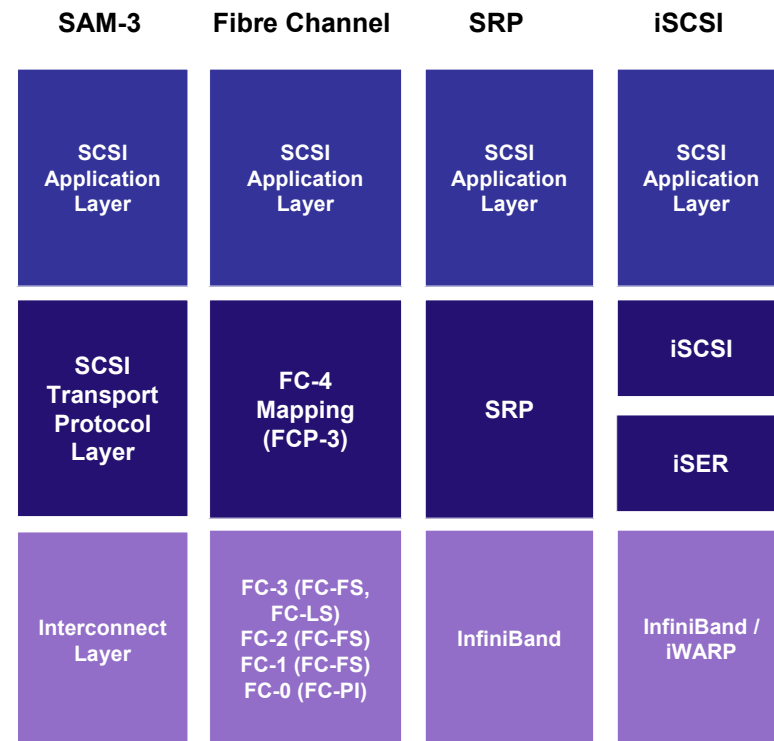
- **Storage**

- ♦ SRP (SCSI RDMA Protocol)
- ♦ iSER (iSCSI Extensions for RDMA)
- ♦ NFSoRDMA (NFS over RDMA)



InfiniBand Block Storage Protocols

- **SRP - SCSI RDMA Protocol**
 - ◆ Defined by T10
- **iSER – iSCSI Extensions for RDMA**
 - ◆ Defined by IETF IP Storage WG
 - ◆ InfiniBand specifics (e.g. CM) defined by IBTA
 - ◆ Leverages iSCSI management infrastructure
- **Protocol offload**
 - ◆ Use IB Reliable Connected
 - ◆ RDMA for zero copy data transfer



SRP - Data Transfer Operations

➤ Send/Receive

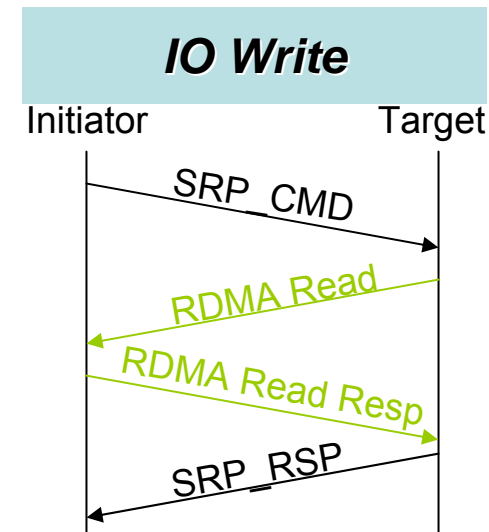
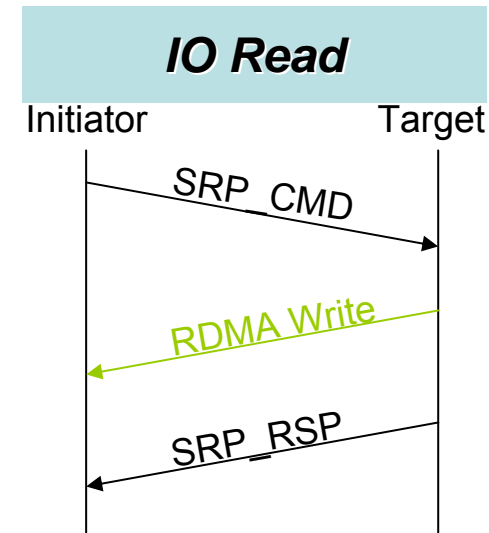
- ◆ Commands
- ◆ Responses
- ◆ Task management

➤ RDMA – Zero Copy Path

- ◆ Data-In
- ◆ Data-Out

➤ iSER uses the same principles

- ◆ Immediate/Unsolicited data allowed through Send/Receive



Data Transfer Summary

	SRP	iSER	iSCSI	FCP
Request	SRP_CMD (SEND)	SCSI-Command (SEND)	SCSI-Command	FCP_CMND
Response	SRP_RSP (SEND)	SCSI-Response (SEND)	SCSI-Response (or piggybacked on Data-In PDU)	FCP_RSP
Data-In Delivery	RDMA Write	RDMA Write	Data-In	FCP_DATA
Data-Out Delivery	RDMA Read RDMA Read Resp.	RDMA Read RDMA Read Resp.	R2T Data-Out	FCP_XFER_RDY FCP_DATA
Unsolicited Data-Out Delivery		Part of SCSI-Command (SEND) Data-Out (SEND)	Part of SCSI- Command Data-Out	FCP_DATA
Task Management	SRP_TSK_MGMT (SEND)	Task Management Function Request/ Response (SEND)	Task Management Function Request/ Response	FCP_CMND

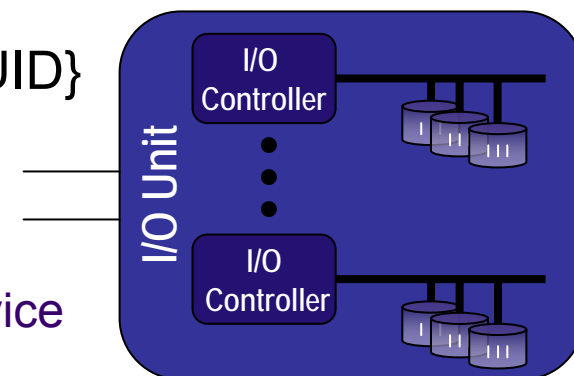
➤ Discovery methods

- ◆ Persistent Information {Node_GUID:IOC_GUID}
- ◆ Subnet Administrator (Identify all ports with CapabilityMask.IsDM)
- ◆ Configuration Manager (CFM)
 - Locate the Device Administrator through Service Record
- ◆ Boot Manager
- ◆ Boot Information Service

➤ Identifiers

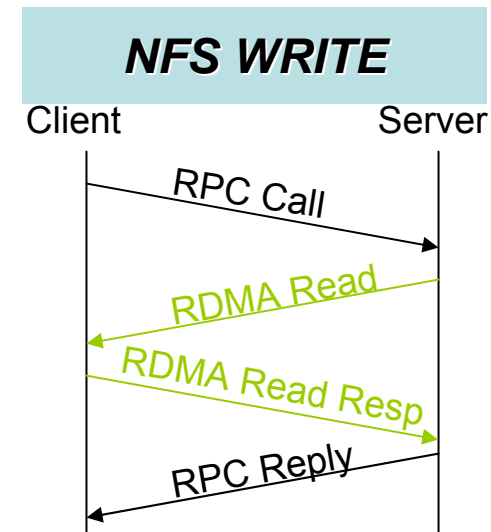
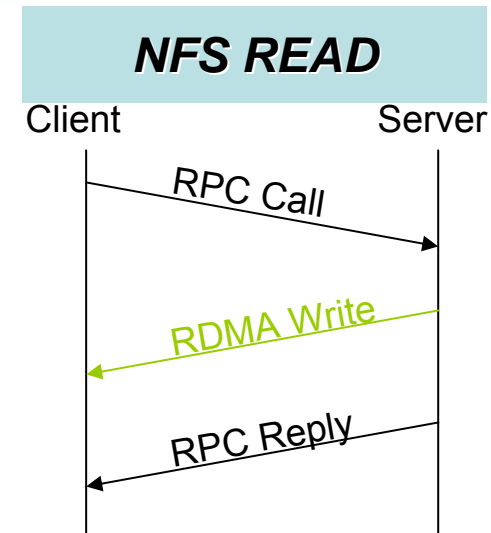
- ◆ Per LUN WWN (through INQUIRY VPD)
- ◆ SRP Target Port ID
{IdentifierExt[63:0], IOC GUID[63:0]}
- ◆ Service Name – SRP.T10.{PortID ASCII}
- ◆ Service ID – Locally assigned by the IOC/IOU

InfiniBand I/O Model

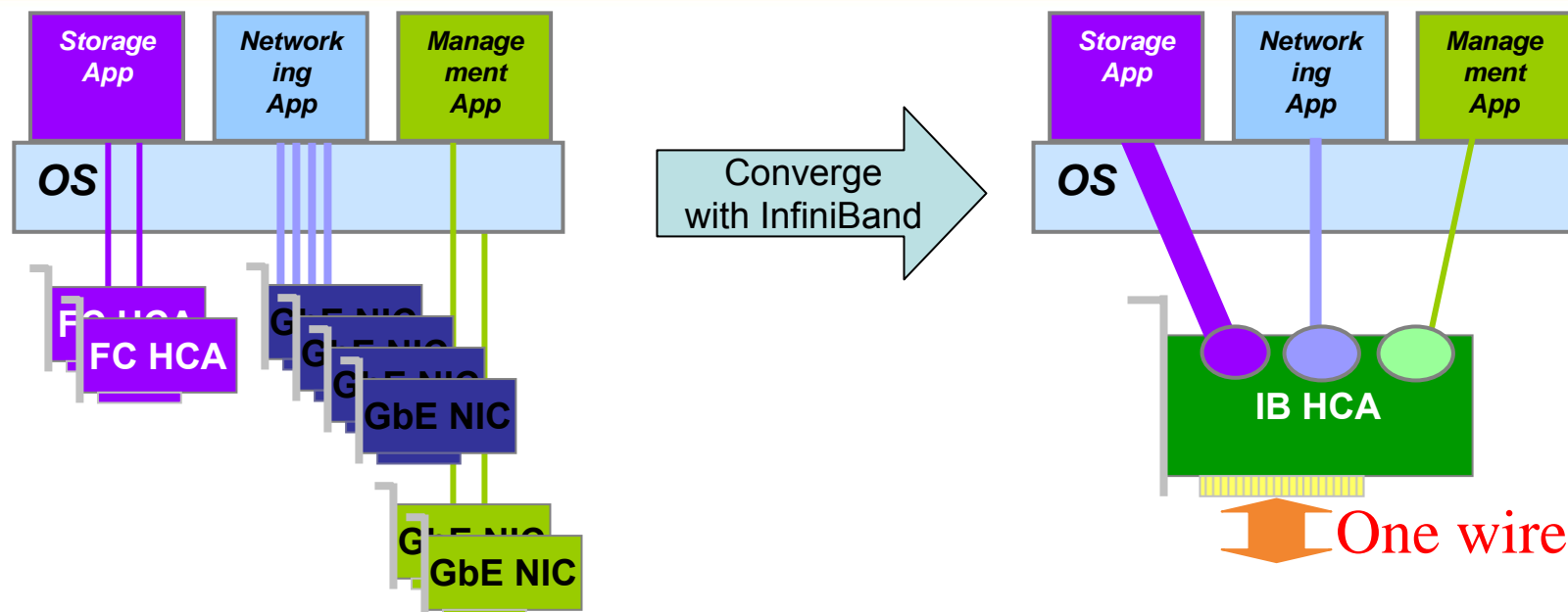


- Leverages all iSCSI infrastructure
 - ◆ Using IP over InfiniBand
- Same iSCSI mechanisms for discovery (RFC 3721)
 - ◆ Static Configuration {IP, port, target name}
 - ◆ Send Targets {IP, port}
 - ◆ SLP
 - ◆ iSNS
- Same target naming (RFC 3721/3980)
 - ◆ iSCSI Qualified Names (iqn.)
 - ◆ IEEE EUI64 (eui.)
 - ◆ T11 Network Address Authority (naa.)

- Defined by IETF
 - ◆ ONC-RPC extensions for RDMA
 - ◆ NFS mapping
- RPC Call/Reply
 - ◆ Send/Receive – if small
 - ◆ Via RDMA Read chunk list - if big
- Data transfer
 - ◆ RDMA Read/Write – described by chunk list in XDR message
 - ◆ Send – inline in XDR message
- Uses InfiniBand Reliable Connected QP
 - ◆ Uses IP extensions to CM
 - ◆ Connection based on IP address and TCP port
 - ◆ Zero copy data transfers



I/O Consolidation



- ❖ Slower I/O
- ❖ Different service needs – different fabrics
- ❖ No flexibility
- ❖ More ports to manage
- ❖ More power
- ❖ More space
- ❖ Higher TCO

- ❖ High bandwidth pipe for capacity provisioning
- ❖ Dedicated I/O channels enable convergence
 - ◆ For Networking, Storage, Management
 - ◆ Application compatibility
 - ◆ QoS - differentiates different traffic types
 - ◆ Partitions – logical fabrics, isolation
- ❖ Gateways - Share remote Fibre Channel and Eth ports
 - ◆ Design based on average load across multiple servers
 - ◆ Scale incrementally – add Ethernet/FC/Server blades
 - ◆ Scale independently

- Multi-port HCAs
 - ◆ Covers link failure
- Redundant fabric topologies
 - ◆ Covers link failure
- Link layer multi-pathing (LMC)
- Automatic Path Migration (APM)
- ULP High Availability
 - ◆ Application level multi-pathing (SRP/iSER)
 - ◆ Teaming/Bonding (IPoIB)
 - ◆ Covers HCA failure and link failure

Performance Metrics

➤ IB Verbs

- ◆ Latency
 - RDMA Write 0.99us
 - RDMA Read 1.87us (roundtrip)
- ◆ Bandwidth
 - 1.5-1.9GB/s (unidirectional)
 - 3.0-3.4GB/s (bidirectional)
 - Depends on PCIe (2.5-5GT/s)

➤ Clustering (MPI)

- ◆ Latency 1.2us
- ◆ Message rate 30M msg/sec

➤ Block Storage (SRP)

- ◆ Bandwidth (1MB I/O, no RAID)
 - I/O Read 1.4GB/s
 - I/O Write 1.2GB/s

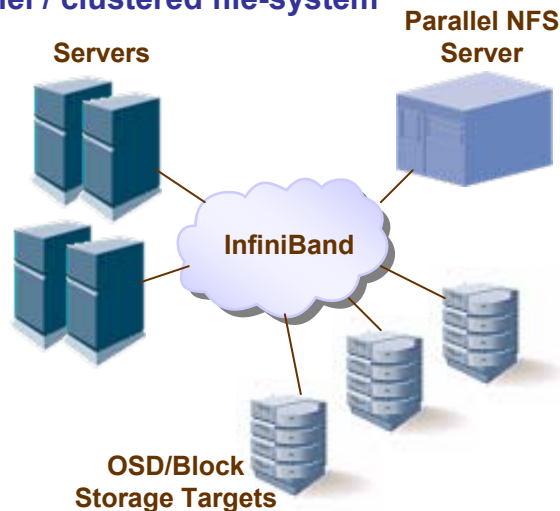
➤ File Storage (NFSoverRDMA)

- ◆ Read 1.3GB/s
- ◆ Write 0.59GB/s

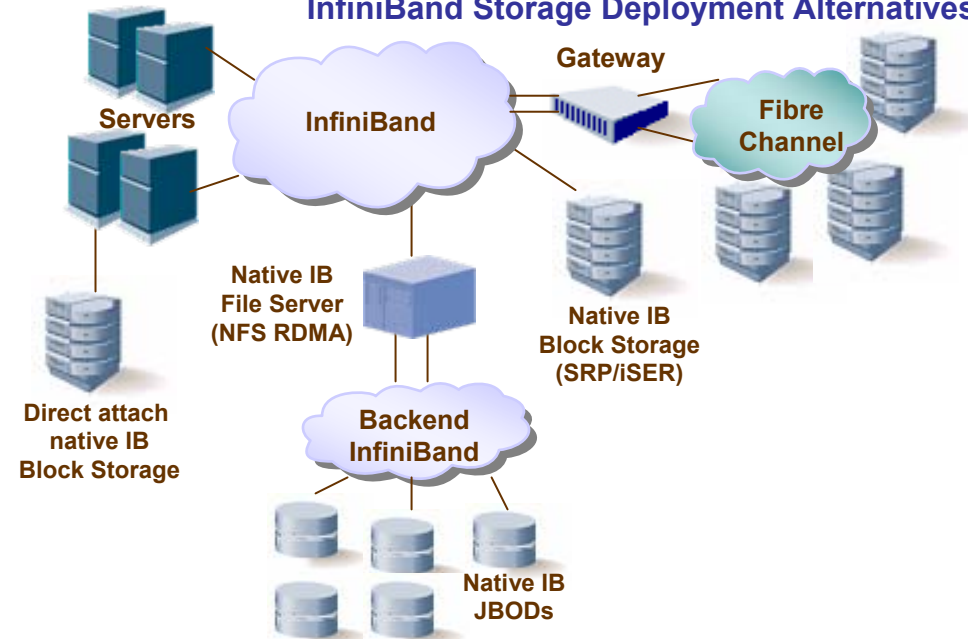
InfiniBand Storage Opportunities & Benefits

- Clustering port can connect to storage
- High Bandwidth Fabric
- Fabric consolidation (QoS, partitioning)
- Efficiency – full offload and zero copy
- Gateways
 - ◆ One wire out of the server
 - ◆ Shared remote FC ports - scalability

Parallel / clustered file-system



InfiniBand Storage Deployment Alternatives



➤ Clustered/Parallel storage, Backend fabric benefits:

- ◆ Combined with clustering infrastructure
- ◆ Efficient object/block transfer
- ◆ Atomic operations
- ◆ Ultra low latency
- ◆ High bandwidth

- Datacenter developments require better I/O
 - ◆ Increasing compute power per host
 - ◆ Server virtualization
 - ◆ Increasing storage demand
- InfiniBand I/O is a great fit for the datacenter
 - ◆ Layered implementation
 - ◆ Brings fabric consolidation
 - ◆ Enables efficient SAN, Network, IPC and Management traffic
 - ◆ Price/Performance
 - ◆ Gateways provide scalable connectivity to existing fabrics
- Existing storage opportunities with InfiniBand

- Please send any questions or comments on this presentation to SNIA: tracknetworking@snia.org

**Many thanks to the following individuals
for their contributions to this tutorial.**

SNIA Education Committee

**Bill Lee
Ron Emerick
Walter Dey**

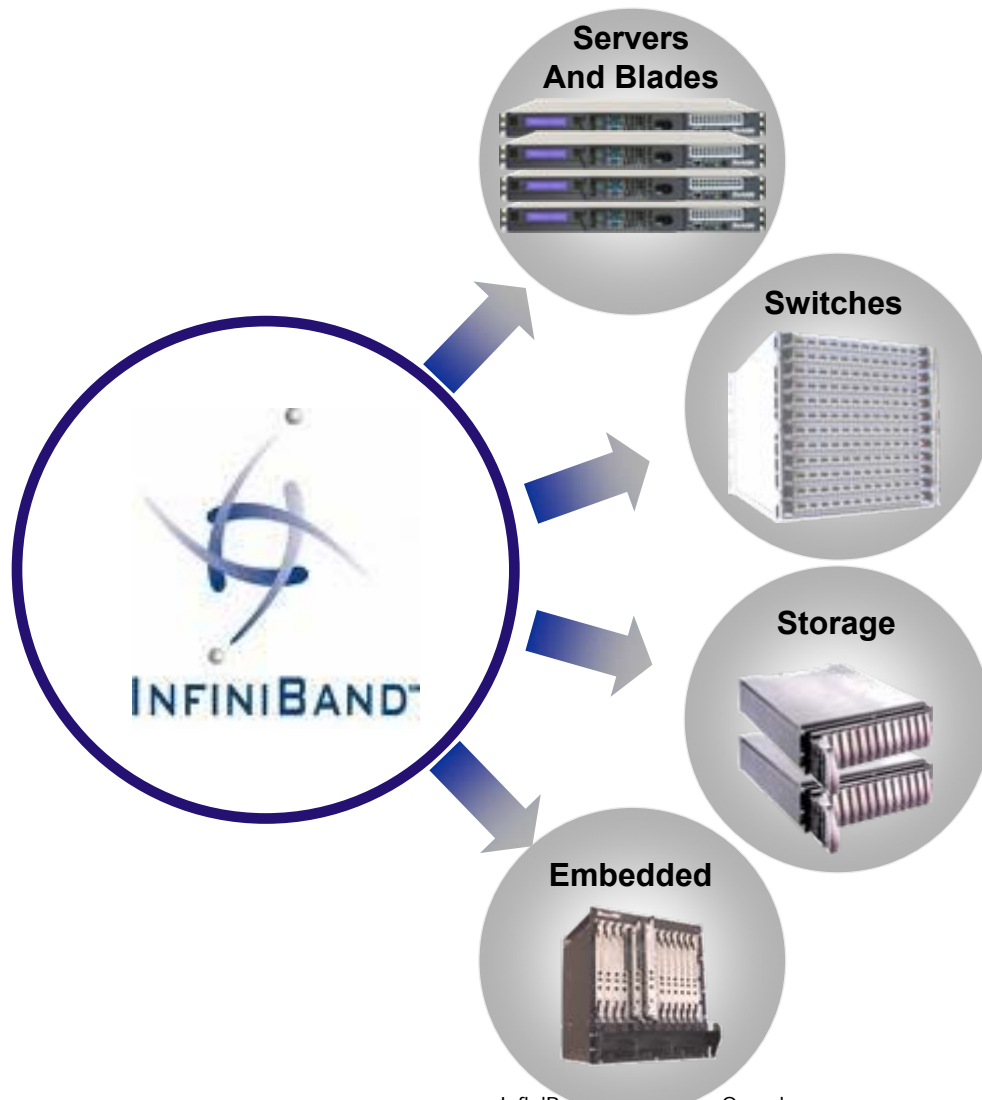
**Howard Goldstein
Sujal Das**



Backup

Education

Interconnect: A Competitive Advantage



End-Users

Enterprise Data Centers

- Clustered Database
- eCommerce and Retail
- Financial
- Supply Chain Management
- Web Services

High-Performance Computing

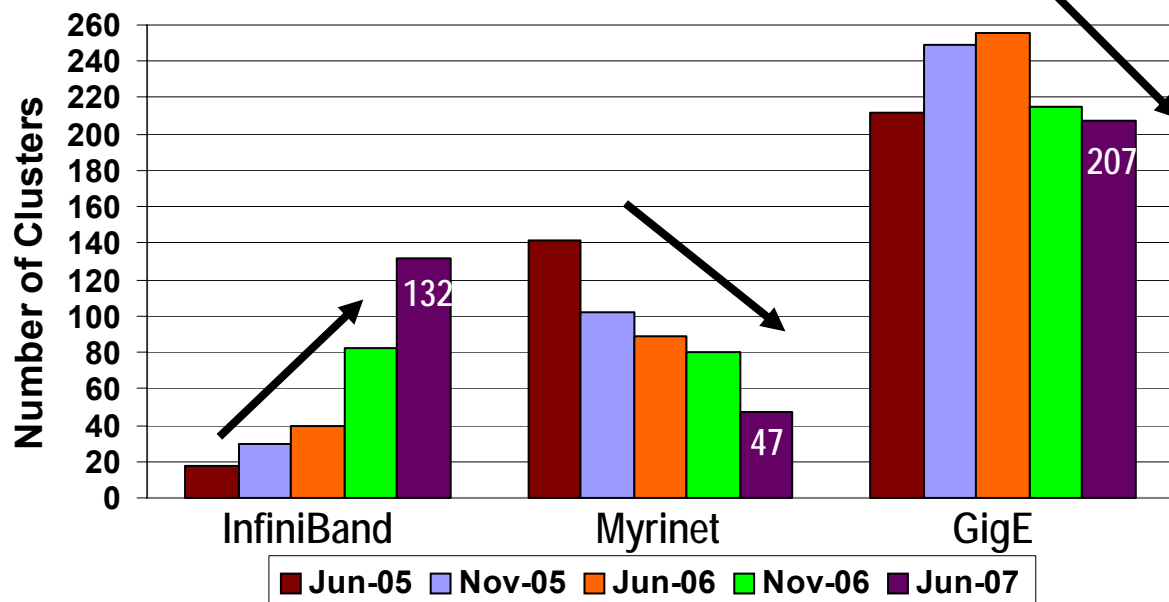
- Biosciences and Geosciences
- Computer Automated Engineering
- Digital Content Creation
- Electronic Design Automation
- Government and Defense

Embedded

- Communications
- Computing and Storage Aggregation
- Industrial
- Medical
- Military

Interconnect Trends – Top500

Top500 Interconnect Trends



Growth rate from Nov 06 to June 07 (6 months)

- ♦ InfiniBand: **+61%**
- ♦ Myrinet: **-41%**
- ♦ GigE: **-4%**

Growth rate from June 06 to June 07 (year)

- ♦ InfiniBand: **+230%**
- ♦ Myrinet: **-47%**
- ♦ GigE: **-19%**

61% growth for InfiniBand from Nov 2006, 230% growth from June 2006

Source: <http://www.top500.org/list/2007/06/>

The TOP500 project was started in 1993 to provide a reliable basis for tracking and detecting trends in high-performance computing.

➤ Data Centers

- ◆ Clustered database, data warehousing, shorter backups, I/O consolidation, power savings, virtualization, SOA, XTP

➤ Financial

- ◆ Real-time risk assessment, grid computing and I/O consolidation

➤ Electronic Design Automation (EDA) and Computer Automated Design (CAD)

- ◆ File system I/O is the bottleneck to shorter job run times

➤ High Performance Computing

- ◆ High throughput I/O to handle expanding datasets

➤ Graphics and Video Editing

- ◆ HD file sizes exploding, shorter backups, real-time production

- InfiniBand software is developed under OpenFabrics Open source Alliance

<http://www.openfabrics.org/index.html>



- InfiniBand standard is developed by the InfiniBand® Trade Association

<http://www.infinibandta.org/home>



Reference

- InfiniBand Architecture Specification Volume 1-2
Release 1.2
 - ◆ www.infinibandta.org
- IP over InfiniBand
 - ◆ RFCs 4391, 4392, 4390, 4755 (www.ietf.org)
- NFS Direct Data Placement
 - ◆ <http://www.ietf.org/html.charters/nfsv4-charter.html>
- iSCSI Extensions for RDMA Specification
 - ◆ <http://www.ietf.org/html.charters/ips-charter.html>
- SCSI RDMA Protocol, DIF
 - ◆ www.t10.org

- APM - Automatic Path Migration
- BECN - Backward Explicit Congestion Notification
- BTH - Base Transport Header
- CFM - Configuration Manager
- CQ - Completion Queue
- CQE - Completion Queue Element
- CRC - Cyclic Redundancy Check
- DDR - Double Data Rate
- DIF - Data Integrity Field
- FC - Fibre Channel
- FECN - Forward Explicit Congestion Notification
- GbE - Gigabit Ethernet
- GID - Global IDentifier
- GRH - Global Routing Header
- GUID - Globally Unique IDentifier
- HCA - Host Channel Adapter
- IB - InfiniBand
- IBTA - InfiniBand Trade Association
- ICRC - Invariant CRC
- IPoIB - Internet Protocol Over InfiniBand
- IPv6 - Internet Protocol Version 6
- iSER - iSCSI Extensions for RDMA
- LID - Local IDentifier
- LMC - Link Mask Control
- LRH - Local Routing Header
- LUN - Logical Unit Number
- MPI - Message Passing Interface
- MR - Memory Region
- NFSoRDMA - NFS over RDMA
- OSD - Object based Storage Device
- OS - Operating System
- PCIe - PCI Express
- PD - Protection Domain
- QDR - Quadruple Data Rate
- QoS - Quality of Service
- QP - Queue Pair
- RDMA - Remote DMA
- RDS - Reliable Datagram Socket
- RPC - Remote Procedure Call
- SAN - Storage Area Network
- SDP - Sockets Direct Protocol
- SDR - Single Data Rate
- SL - Service Level
- SM - Subnet Manager
- SRP - SCSI RDMA Protocol
- TCA - Target Channel Adapter
- ULP - Upper Layer Protocol
- VCRC - Variant CRC
- VL - Virtual Lane
- WQE - Work Queue Element
- WRR - Weighted Round Robin