



Checklist for Selecting and Deploying Scalable Clusters with InfiniBand Fabrics

Lloyd Dickman, CTO

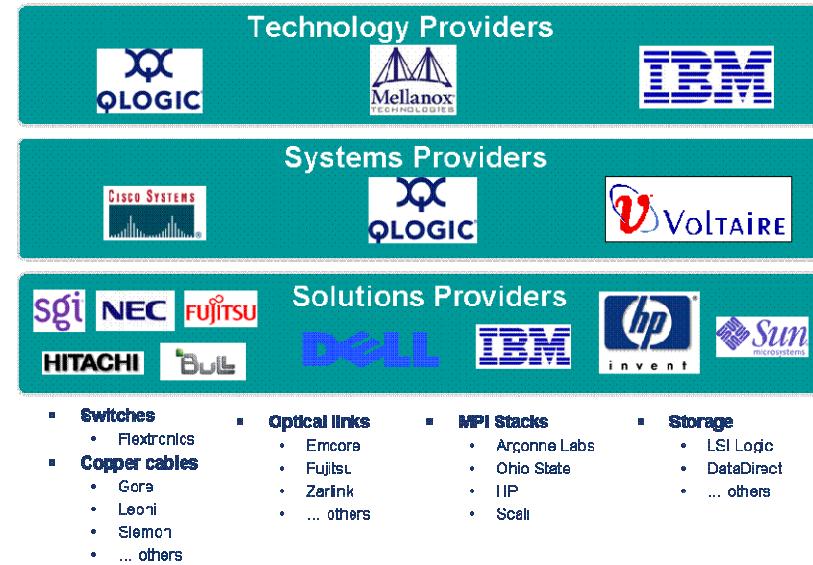
InfiniBand Products
Host Solutions Group
QLogic Corporation

November 13, 2007 @ SC07, Exhibitor Forum

InfiniBand is the Fabric of Choice for HPC

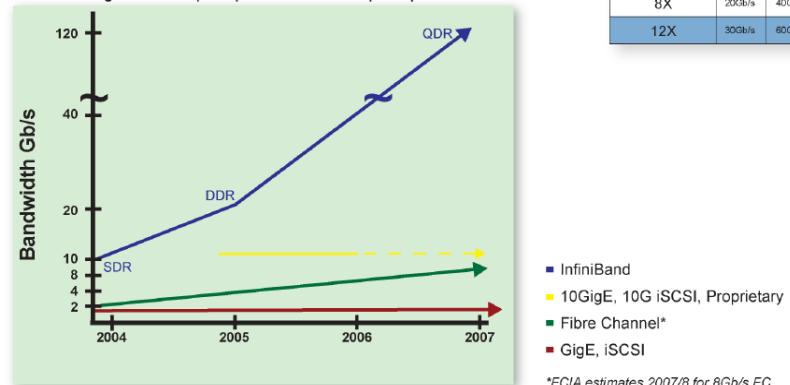


- Large and growing standards-based ecosystem
- Aggressive technology roadmap
- Low-latency, high messaging rate implementations in market
- Evidenced by increasing traction in Top500



InfiniBand Roadmap

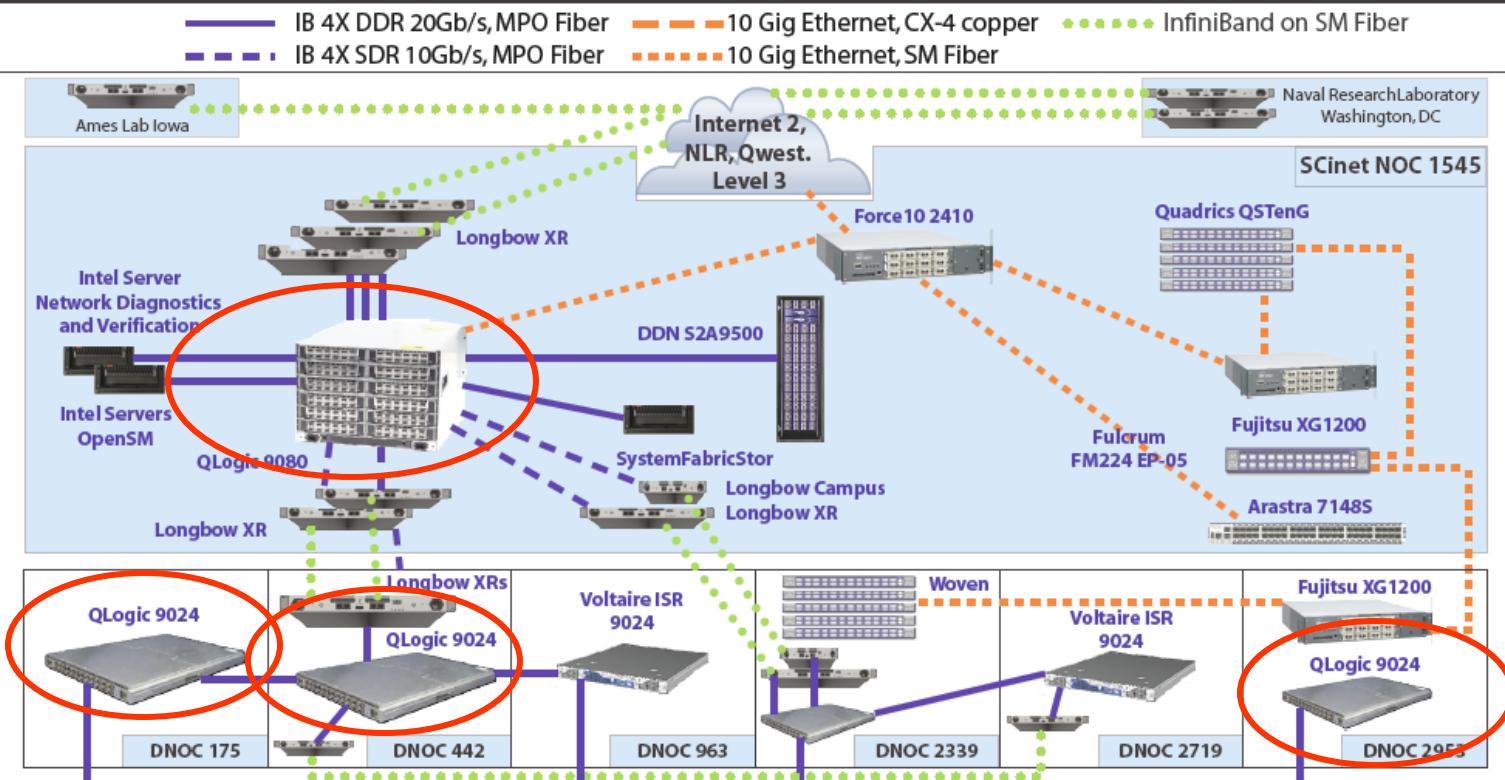
InfiniBand's roadmap outpaces all proprietary and standard based I/O technologies in both pure performance and price/performance



InfiniBand Interoperability



SC07 SCinet OpenFabrics InfiniBand and Low Latency Ethernet Networks



InfiniBand Exhibitors

AMD - 757	Mellanox - 127	Violin/SFW - 158
Ames Lab - 181	Nat Cntr Data Mining - 2623	Voltaire - 1137
HP - 1417	OpenFabrics - 2660	Yotta Yotta - 703
IBM - 1215	QLogic - 261	SC Education Prog - L1
Intel - 820	Sun Microsystems - 514	

iWARP Exhibitors

Arastra - 936	NetEffect - 2753
Chelsio - 2337	OpenFabrics - 2660
Fujitsu - 1429	Woven - 2733
Fulcrum - 2161	

Budget

CapEx

+

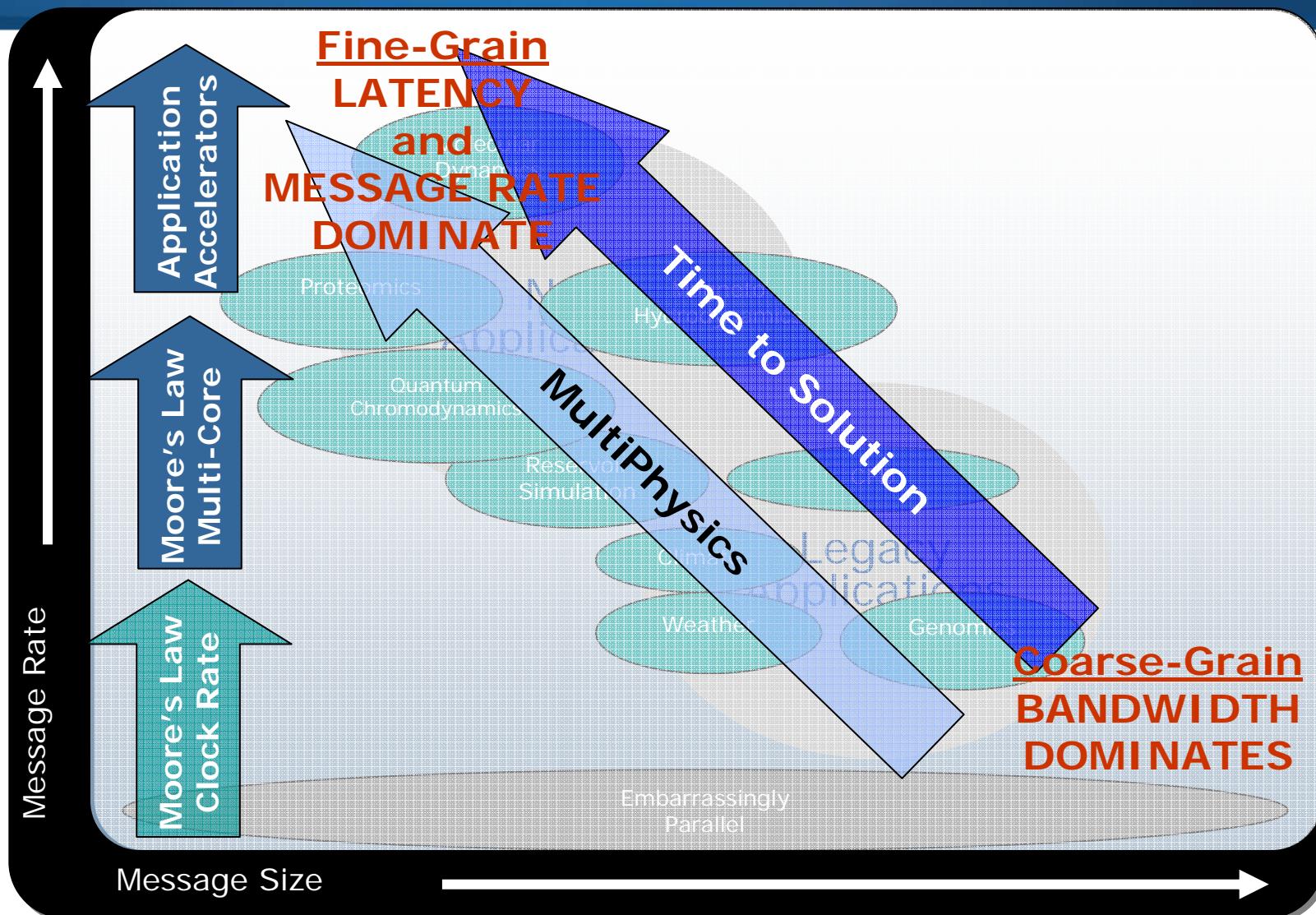
OpEx

Decisions

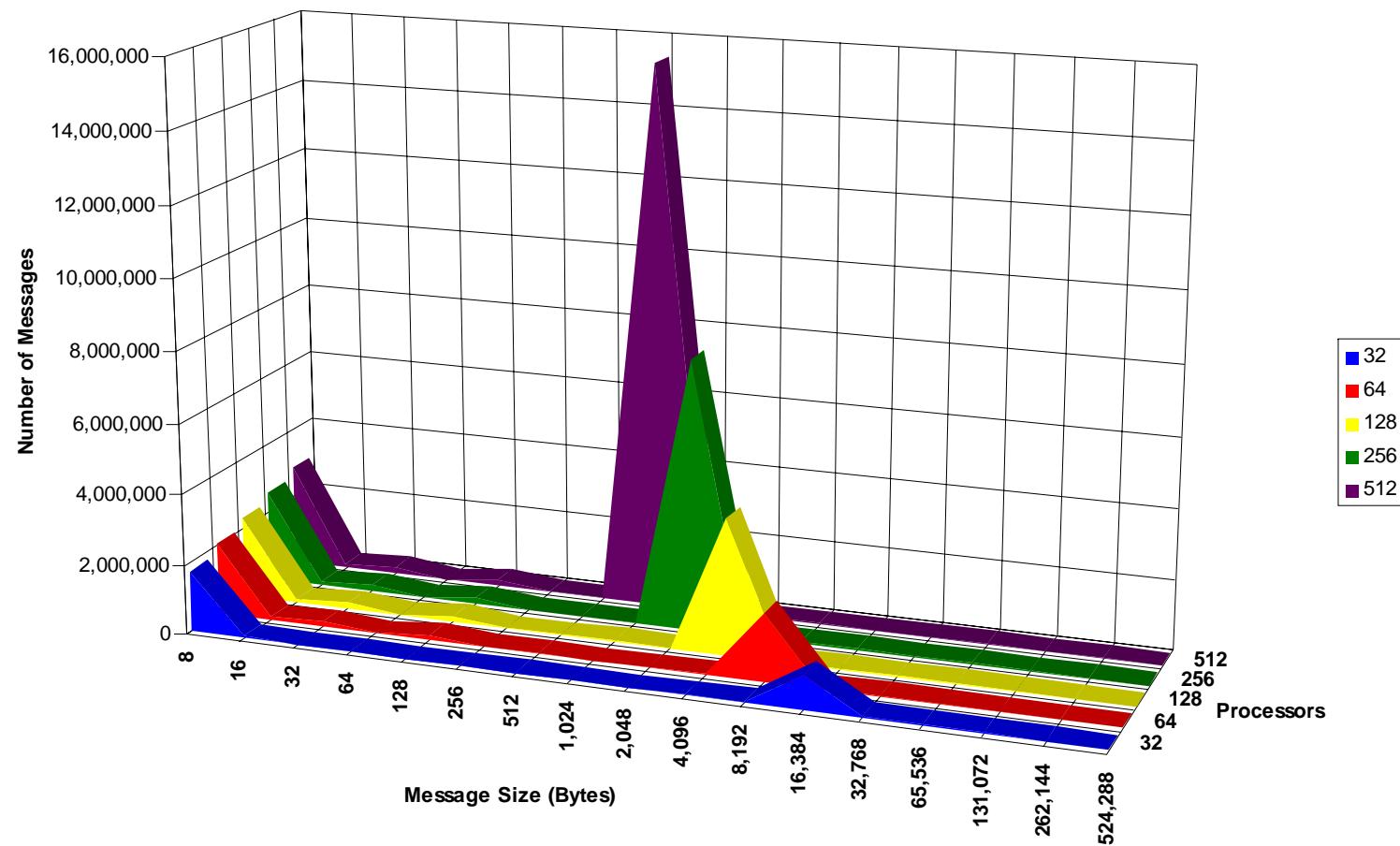
- ✓ Host Architecture
- ✓ Fabric Topology
- ✓ Fabric Convergence / Unification
- ✓ RAS, Power / Cooling
- ✓ Software – host applications, ULPs and fabric

Performance

For HPC, Application Messaging Profiles Matter!



Application Messaging Patterns Change with Scale



Source: PathScale measurements, averages per node, DL_POLY_2 AWE Medium case, 19dec2005

QLogic InfiniBand Provides Superior Application Scaling



Comparison	Benchmark	QLE7140 SDR	QLE7240 DDR	QLE7280 DDR	vs Competitive
Point-to-Point Scaling	Point-to-Point Bandwidth	OSU Bandwidth (peak)	0.9 GB/s	1.5 GB/s	1.9 GB/s +24%
	Point-to-Point Latency	OSU Latency (0 Byte)	1.6 μ s	1.5 μ s	1.4 μ s Tie
	Point-to-Point Message Rate (non-coalesced)	OSU Message Rate @ 4 processes/node	9 M/s	12 M/s	13 M/s 2.6X
	Scalable Latency	HPCC Random Ring Latency @ 32 cores	1.5 μ s	1.4 μ s	1.3 μ s 1.8X
	Scalable Message Rate (non-coalesced)	HPCC MPI Random Access @ 32 cores	.12 GUPs	.12 GUPs	.13 GUPs 7.6X

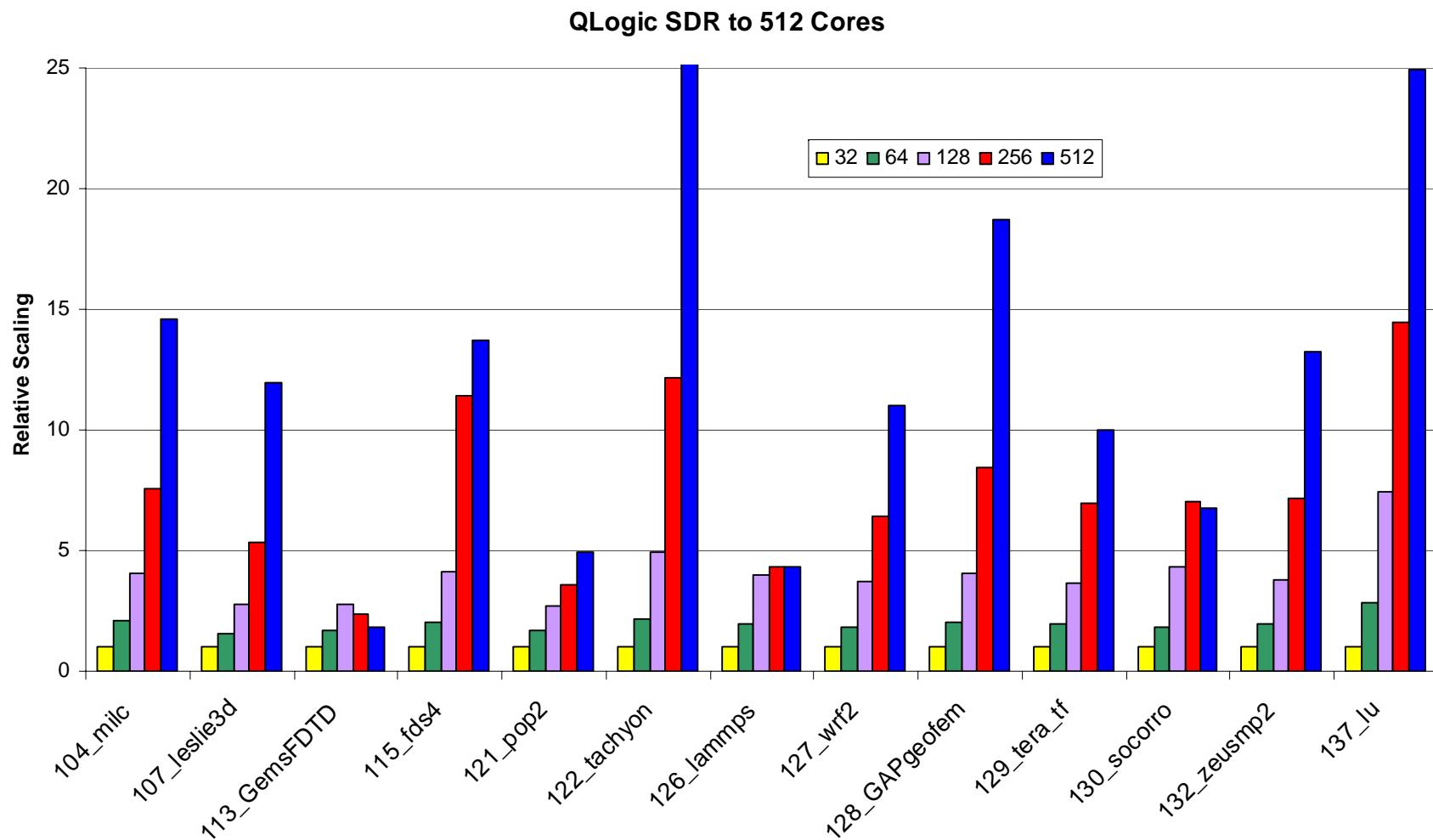
Point-to-point microbenchmarks are no longer sufficient predictors of performance at scale.

Scaling results are for 8 nodes / 32 cores.

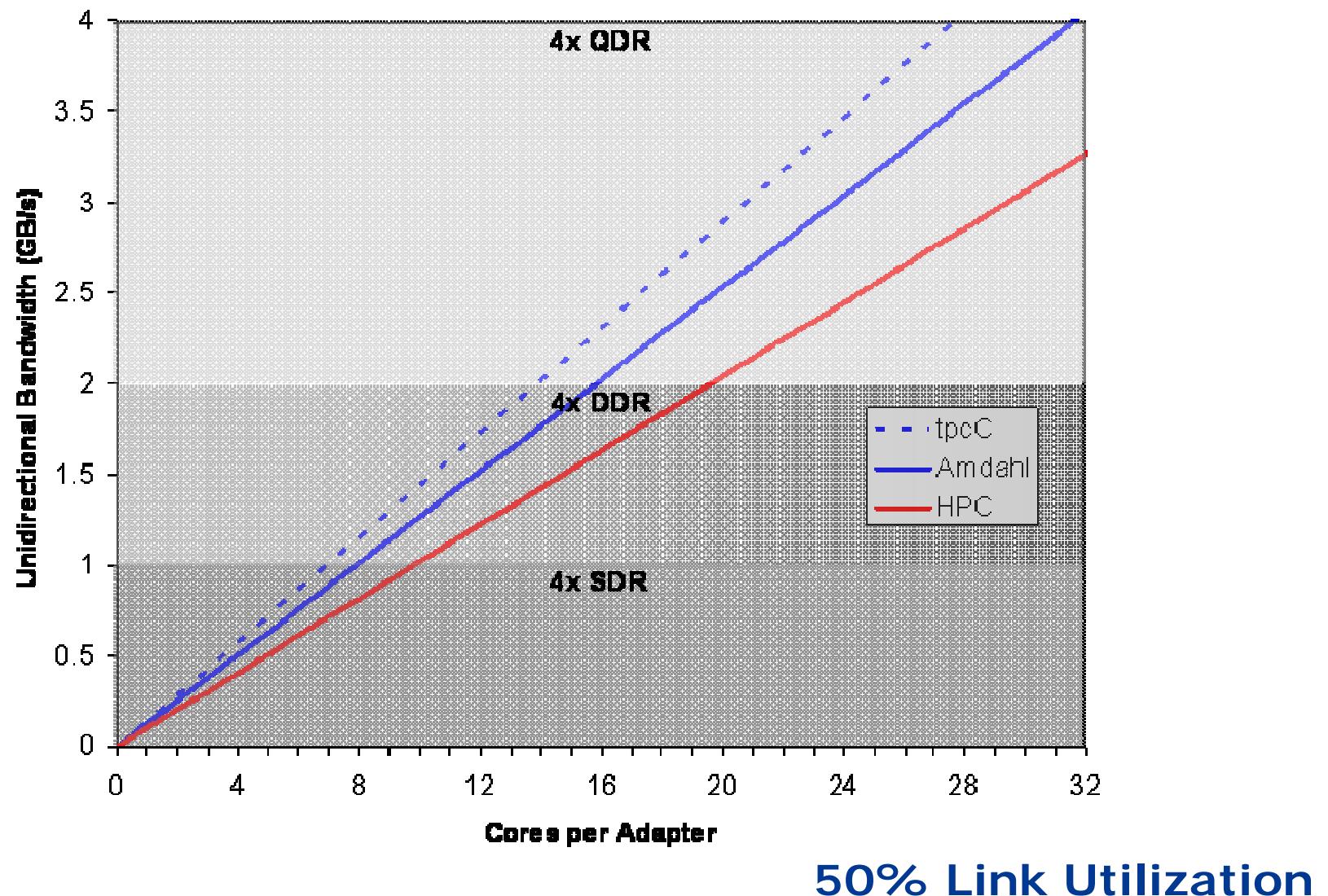
Latency results include a single switch crossing.

Results use native MPIS: InfiniPath MPI for QLogic

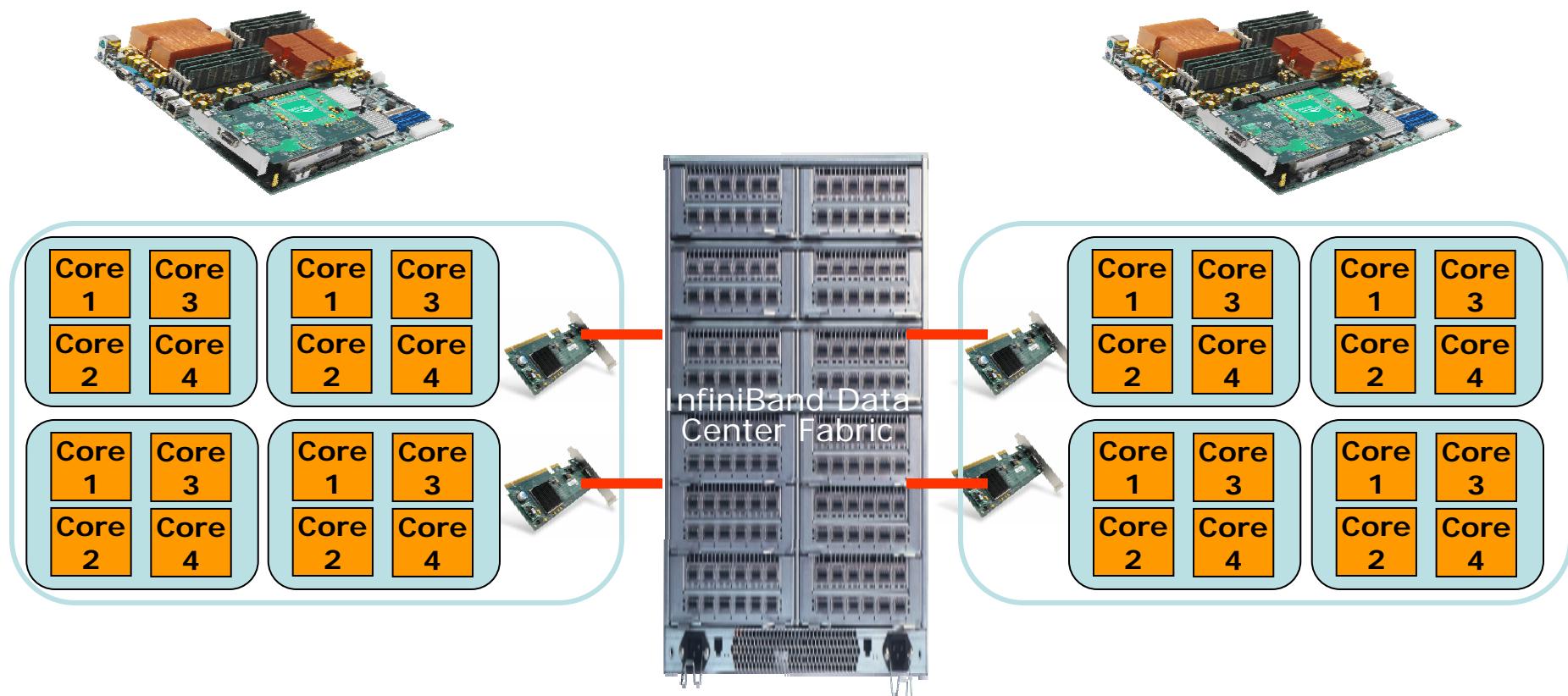
SPEC MPI2007 Results Confirm Scalability



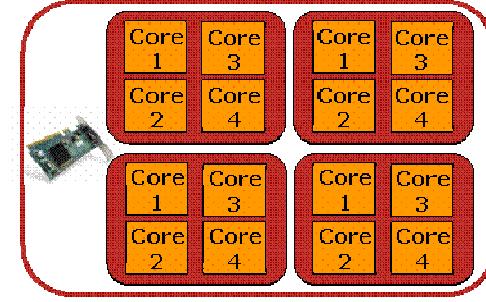
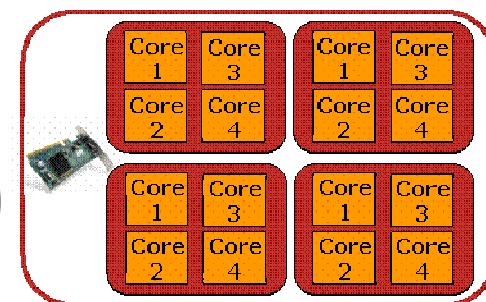
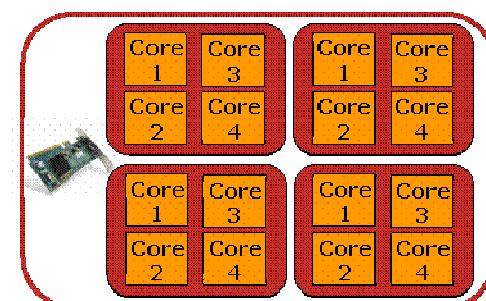
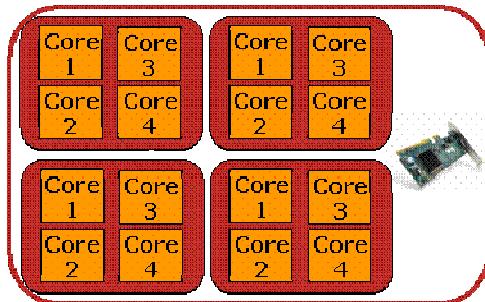
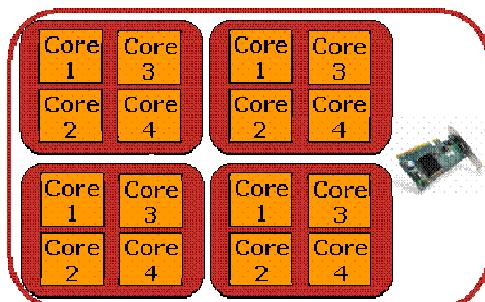
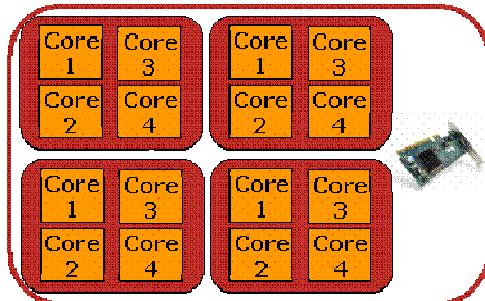
Bandwidth Must Account for ALL Traffic



Host-to-Fabric Rails



Fabric Topology Considerations



- Multi-rail
- Fat Tree / FBB
- Oversubscription
- Asymmetric Topology

QLogic InfiniBand Fabric Products



Core Fabric Switches

Multi-Protocol Fabric Director Modules

INFINIBAND™
Trade Association

IEEE
802
FCIA

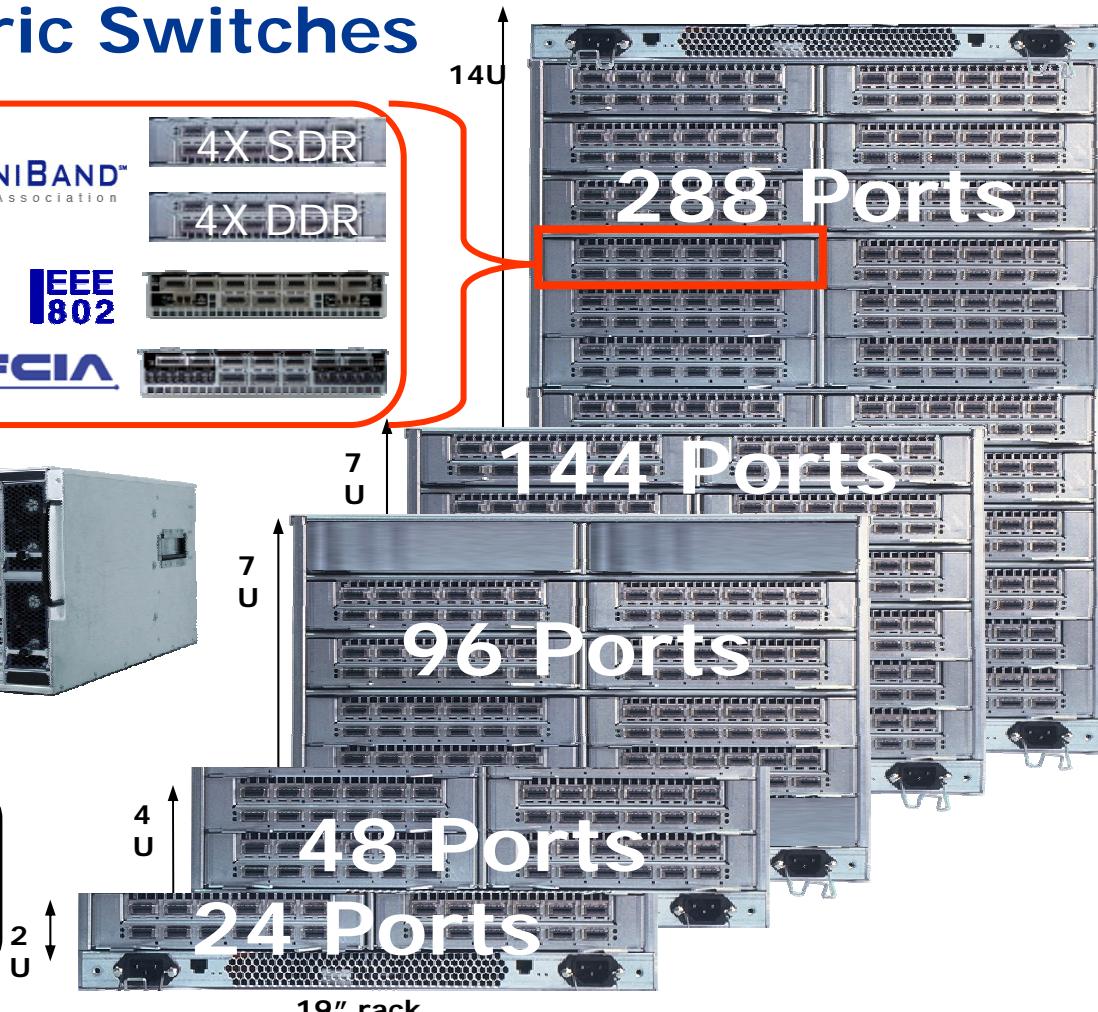
SDR and DDR Spines

Edge Fabric Switch

24 Ports

1 U

2 U



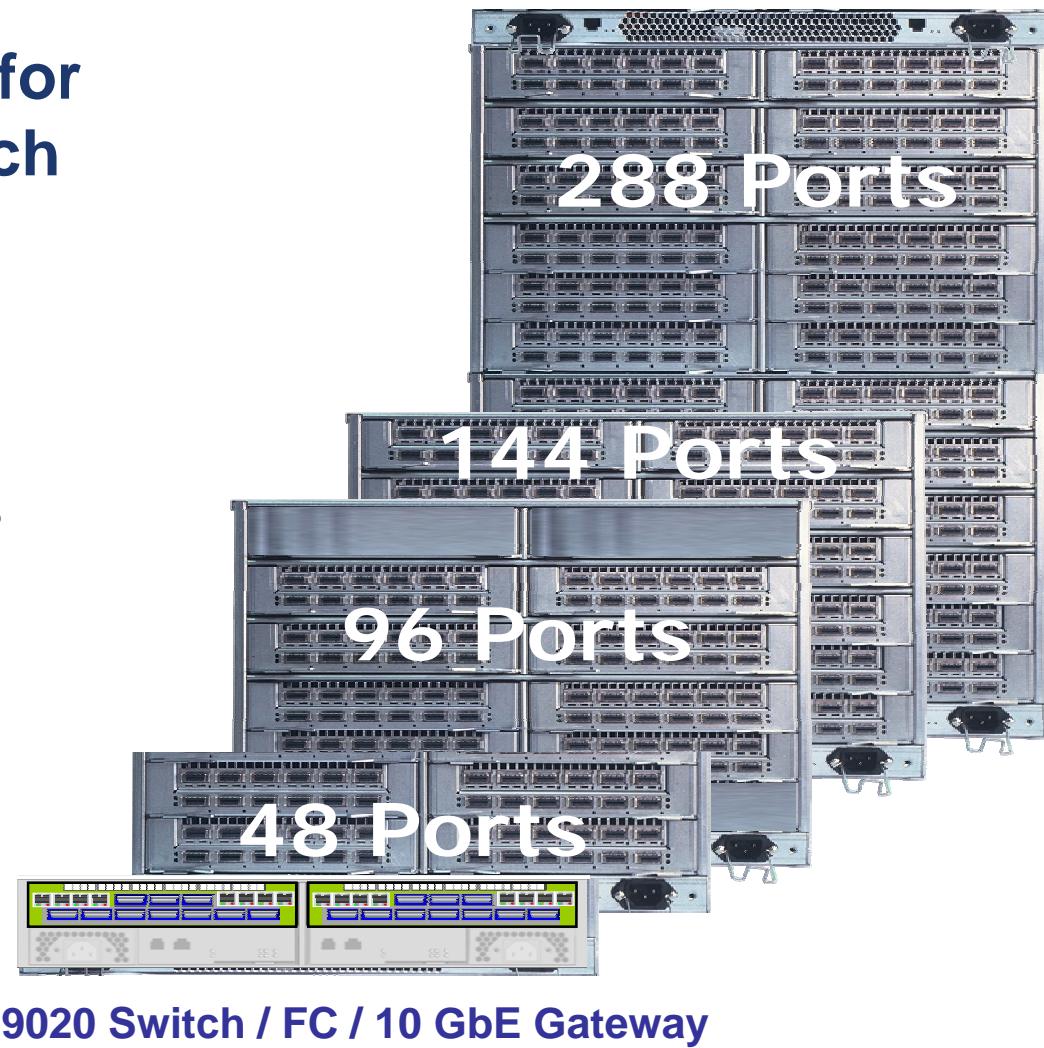
Production quality DDR since April 2006,

Over 50,000 ports to date with extensive OEM qualifications

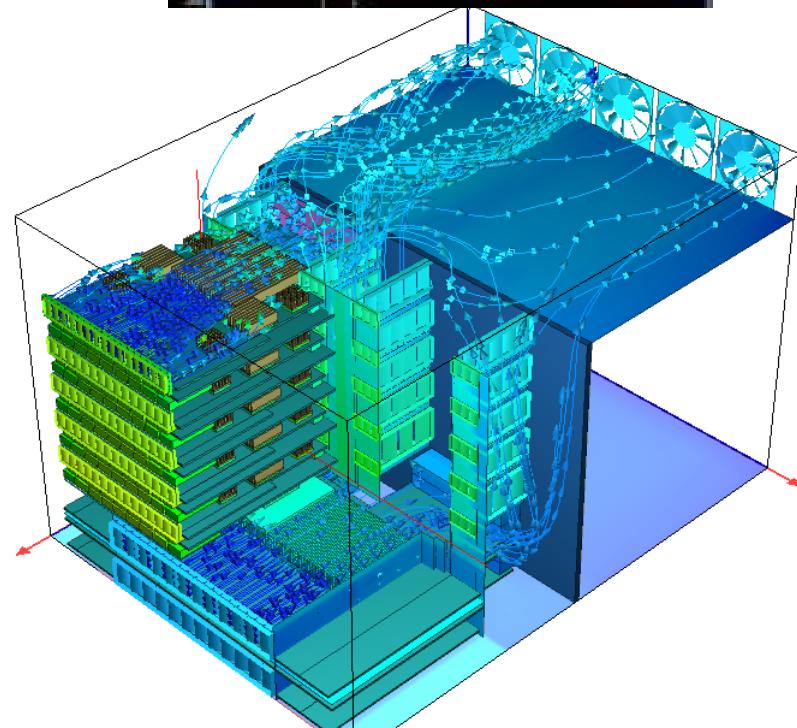
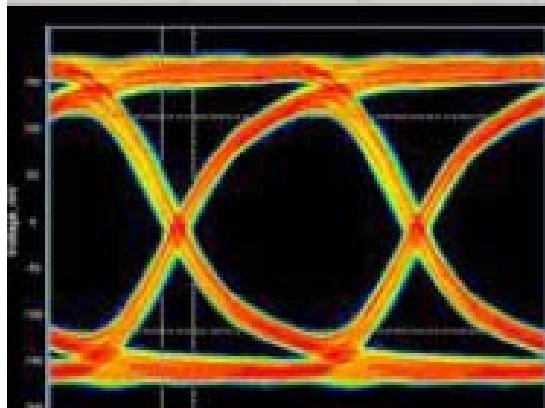
Matching Fabric to Node Count



- Optimal Fabric Size for today's 24-port switch chips
 - 1 tier – 24 ports
 - 2 tier – 288 ports
 - 3 tier – 3,456 ports



Advanced Design Capabilities



■ Signal Integrity

- Internal and external bit lanes capable of 10Gb/sec
- Significant signal integrity challenge
- Leveraging established QLogic center for signal integrity analysis at headquarters

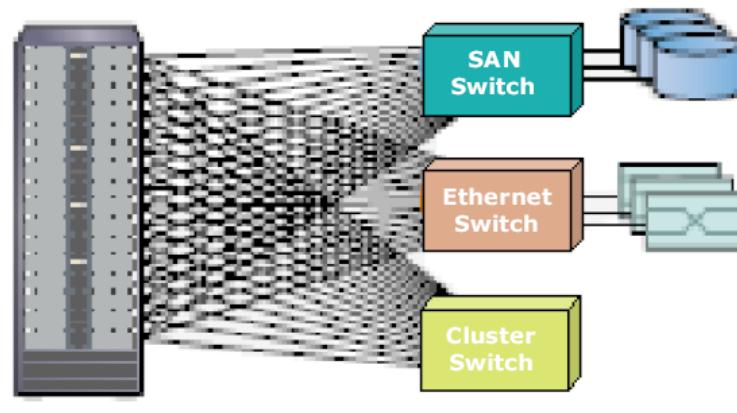
■ Thermal Analysis

- Chip power nearly triple previous generation
- Leveraging QLogic CFD capabilities

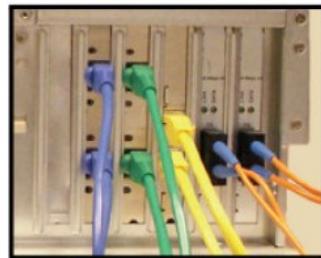
Unified InfiniBand Fabric Reduces Cost



Typical Cluster Interconnect

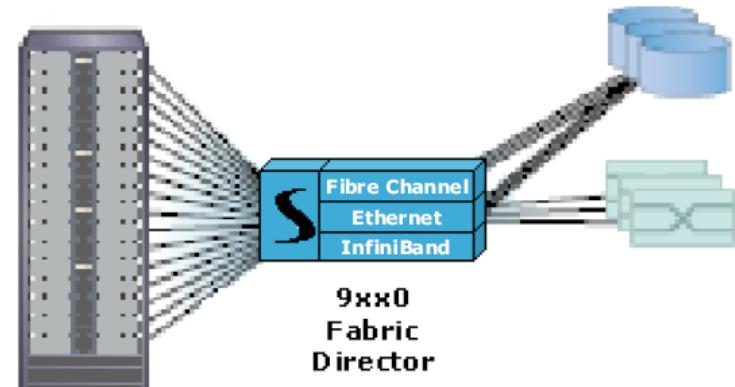


Server Cluster

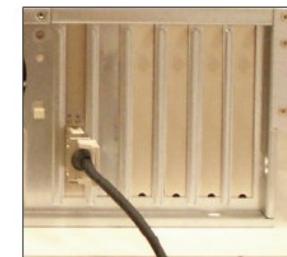


Communications
Computing
Management
Storage

Unified Fabric Interconnect



Server Cluster



“One Wire”



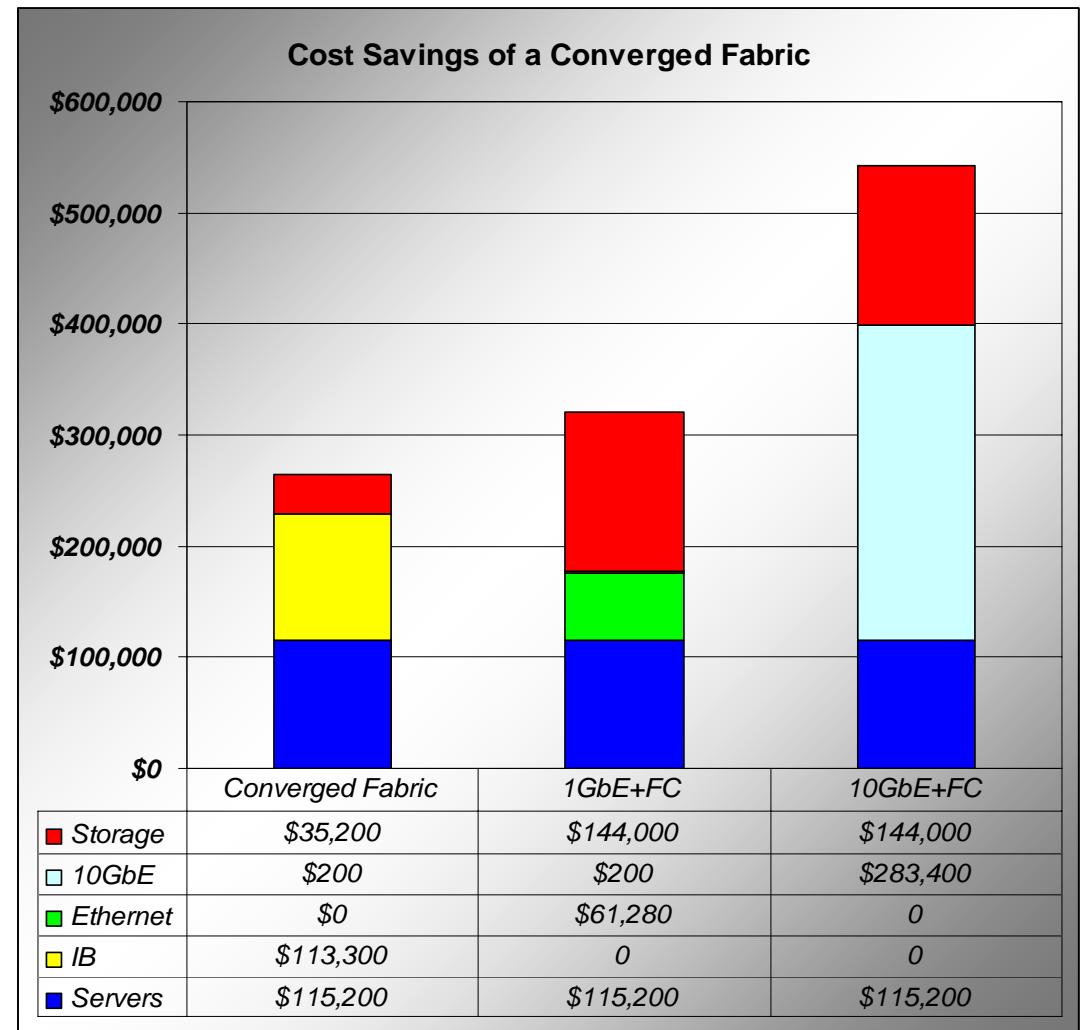
**Unified Fabric
Deployments
Increasing**

Communications
Computing
Management
Storage

Cost and Power Advantages of Converged Fabric



- **CF cost and power savings**
 - FC HBAs and switches
 - Ethernet NICs and switches
 - Cables
- **Additional cost savings**
 - Simpler management: single network
- **Comparisons for a 64-node cluster**
 - CF advantage increases for larger clusters
 - 30-38% cost savings on a 400/800-node FSC cluster
- **Power savings from CF**
 - ~30% vs. GigE

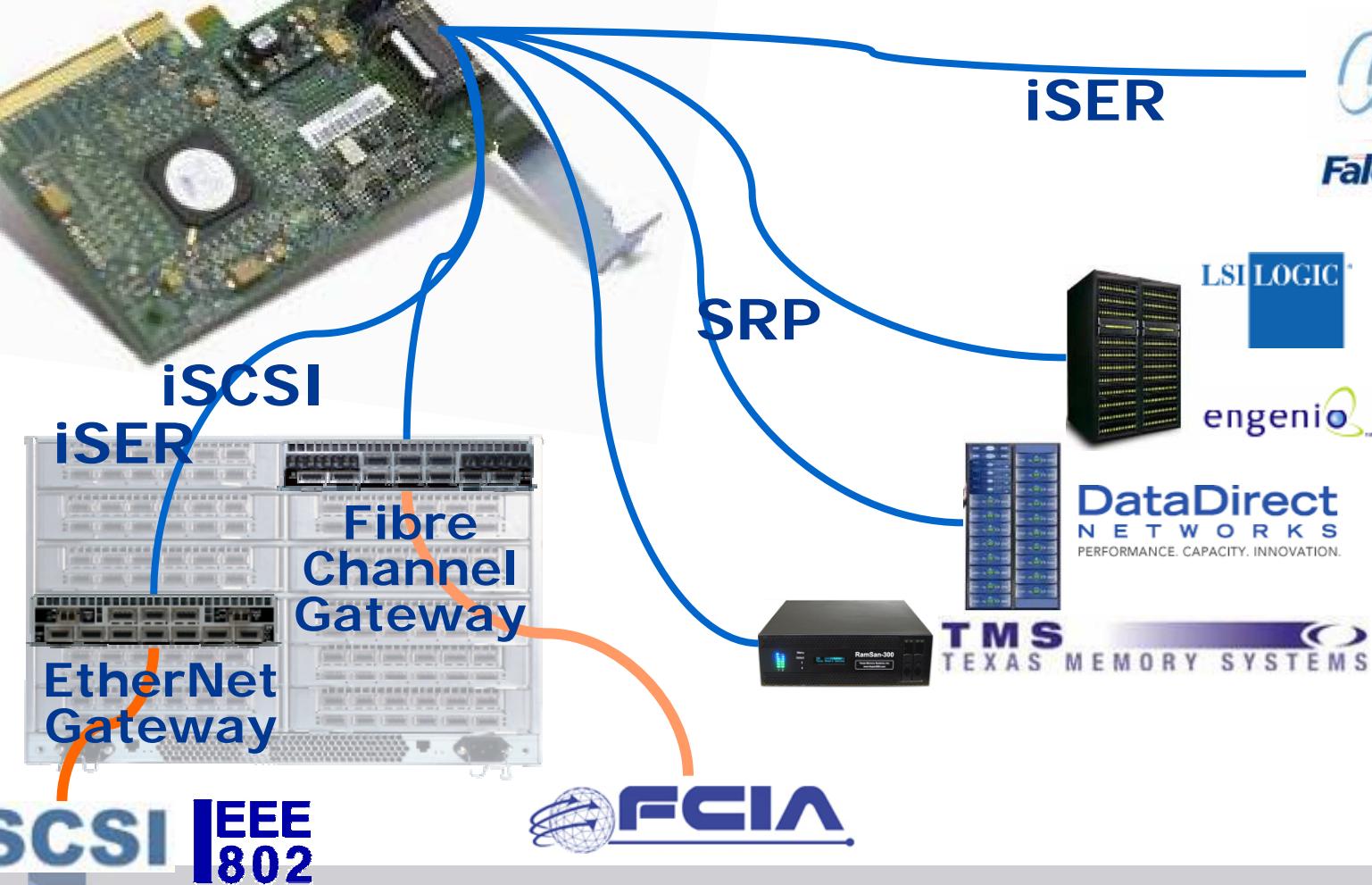


Source: HP BoF, Converged Fabrics, SC07

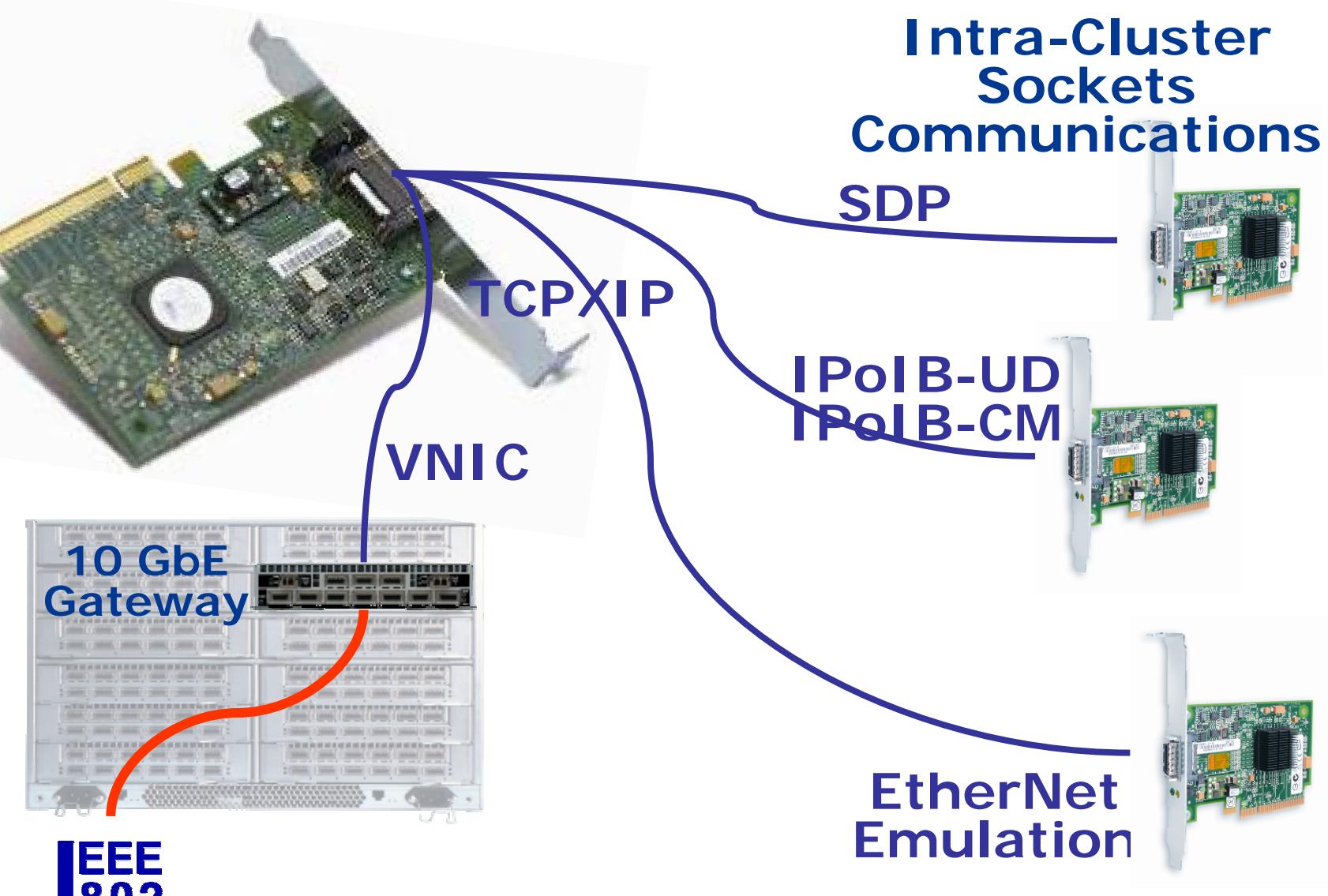
Connectivity to ALL Storage



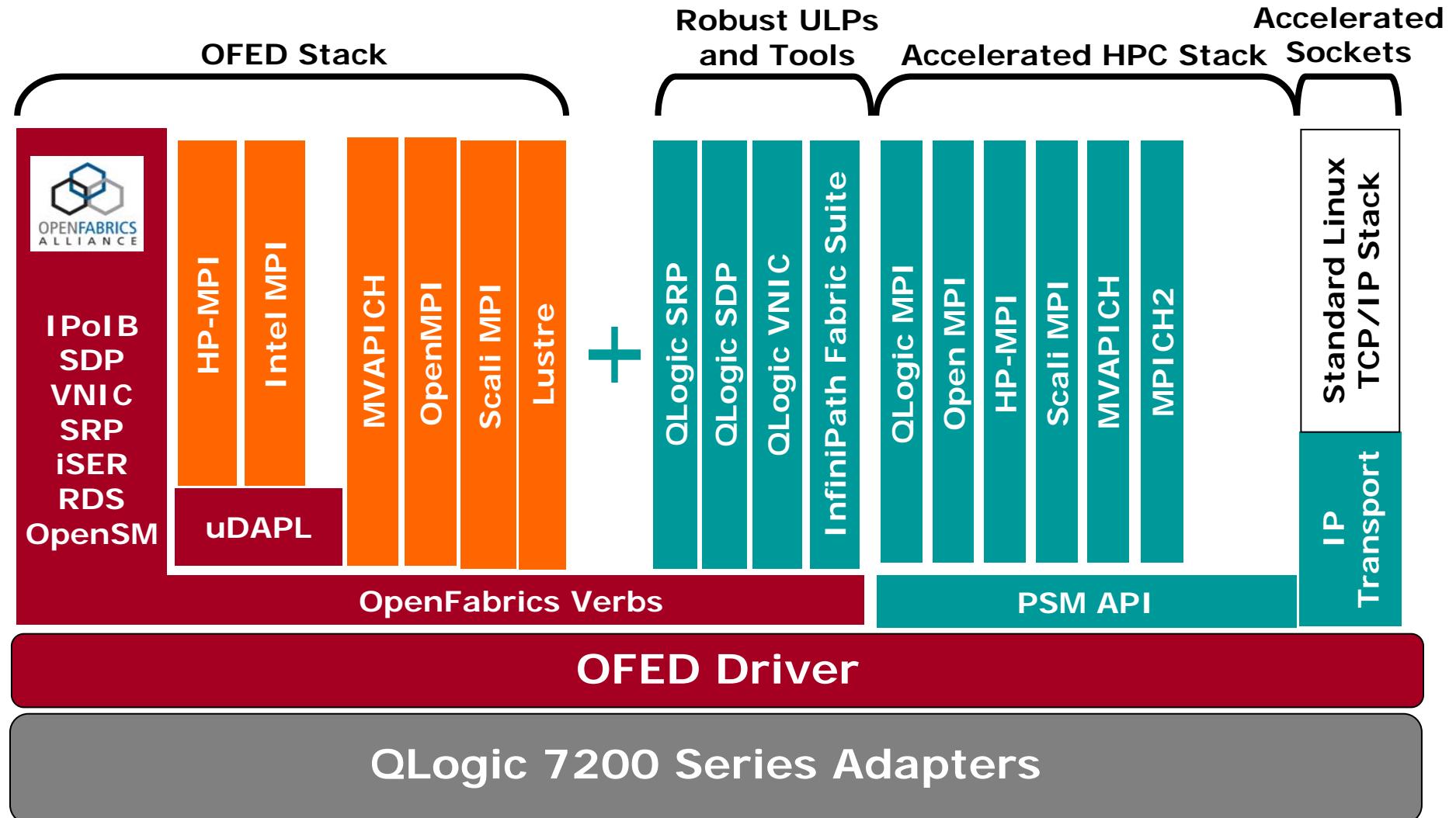
Direct Connect
InfiniBand
Storage



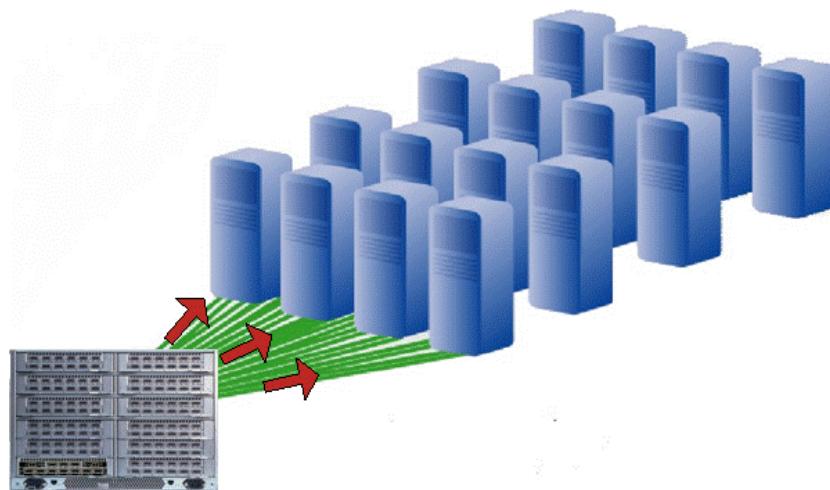
Connectivity to ALL Sockets Applications and Ethernet Networks



OFED+ Provides OpenFabrics Benefits Plus Optional Value Added Capabilities



Large Fabrics Challenges



- **Scale of InfiniBand deployments growing rapidly**
 - 1000s of Nodes, soon 10,000s of Nodes
- **Fabric is the cluster**
 - 1 wire vision being deployed now
 - Requires high degree of fabric stability
- **Fabric Diagnosis and Resiliency**
 - Many moving parts at large scale
- **Demands of End Nodes on SA increasing**

InfiniBand Fabric Suite

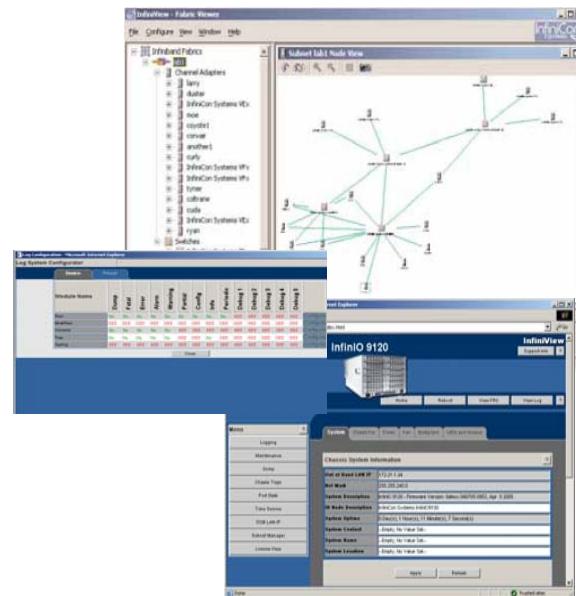


```
Fast Fabric IB Host Setup Menu
Host List: /etc/sysconfig/iba/hosts

0) Edit Config and Select/Edit Hosts Files      [ Skip ]
1) Verify Hosts via Ethernet ping             [ Skip ]
2) Verify rsh/rcp Configured                  [ Skip ]
3) Setup Password-less ssh/scp              [Perform]
4) Copy /etc/hosts to all hosts               [ Skip ]
5) Show uname -a for all hosts                [ Skip ]
6) Configure IPoIB IP Address                 [ Skip ]
7) Reboot Hosts                               [ Skip ]
8) Refresh ssh Known Hosts                   [ Skip ]
9) Run a command on all hosts                [ Skip ]
a) Copy a file to all hosts                  [ Skip ]
b) View ibtest result files                 [ Skip ]

P) Perform the selected actions
N) Select None

X) Return to Previous Menu (or ESC)
```



Lowering InfiniBand OpEx

- **Automates Cluster Install, Verification and Administration**
 - For both Hosts and Switches
- **Rich set of tools and capabilities for ongoing Fabric Administration**
 - Easy to use TUI and CLI
 - Easy integration into site specific tools
- **Rapidly Pinpoints fabric errors**
 - Extensive verification and diagnostic capabilities
- **Industry Leading Performance and Scalability**
 - 1024 node fabric initialization in <20 seconds
 - Fabric Change in < 12 seconds
 - Trap Based event notification & subscription supported
 - Load balances LID(s) across ISL(s)
- **Highly Resilient**
 - No single point of failure
 - Automatic Fabric Reconfiguration
- **Complete IBTA Support**
 - Optional SA queries (Trace Route, Link records, etc)
- **Flexible Deployment Options**
 - Embedded Fabric Management
 - Host Fabric Management

Please Visit Us to Learn More About Complete InfiniBand – Booth #261



Cluster Administration using InfiniBand Fabric Suite

Advanced Technology Demonstrations

HPCtrack Partnership Program

9020 InfiniBand Switch and Multi-Protocol Gateway

Fluent Application running multiple data sets on a converged fabric with Fibre Channel storage via a gateway

InfiniBand DDR Adapter for highly scaling applications





Thank You

