

Adventures in Wide Area InfiniBand

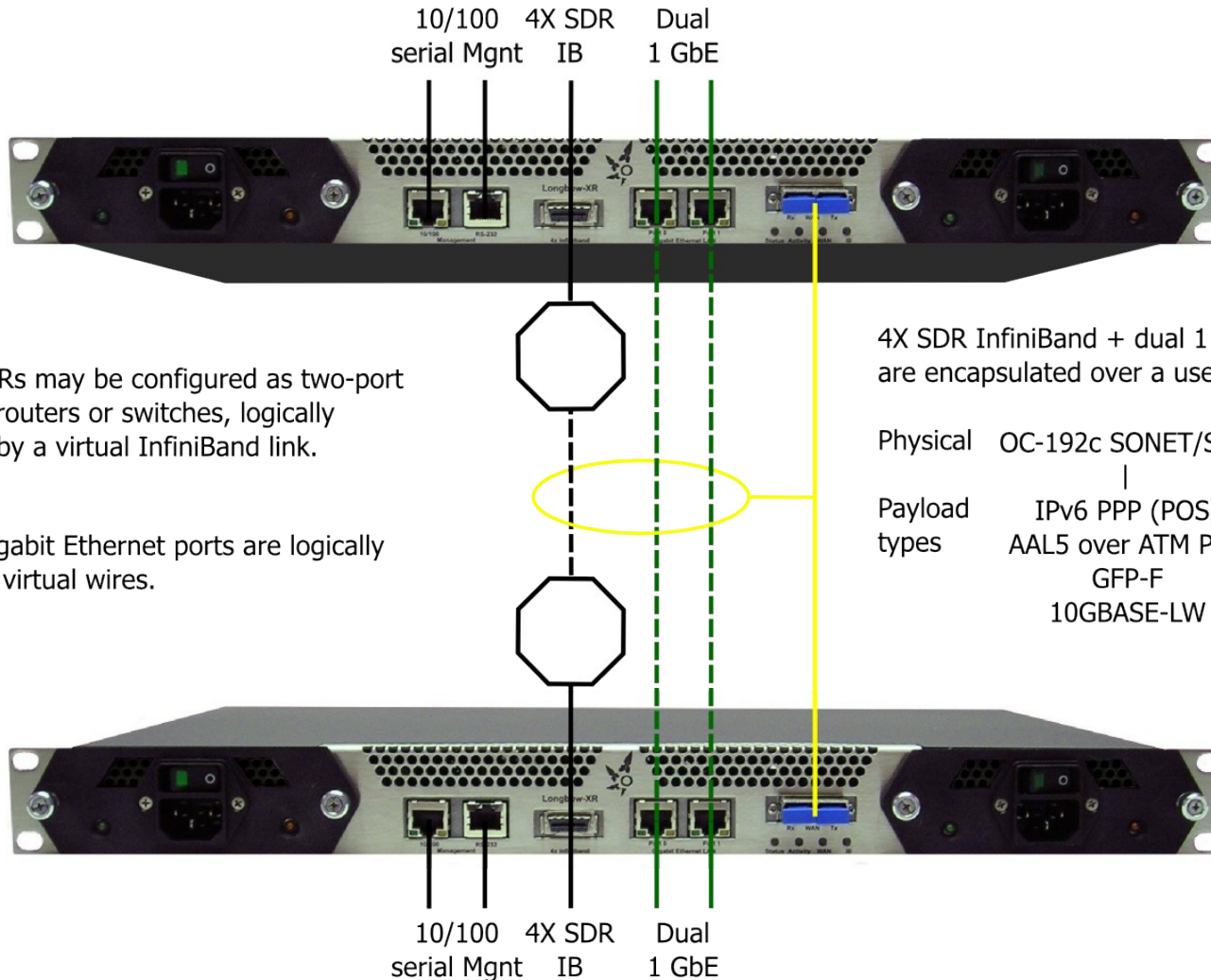


OPEN**FABRICS**
ALLIANCE



David Southwell

Longbow XR in production



Longbow XRs may be configured as two-port InfiniBand routers or switches, logically connected by a virtual InfiniBand link.

The dual 1 Gigabit Ethernet ports are logically connected by virtual wires.

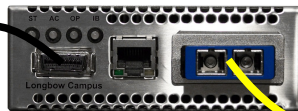
4X SDR InfiniBand + dual 1 Gigabit Ethernet links are encapsulated over a user-configurable WAN:

Physical	OC-192c SONET/SDH	10GBASE-LW/LR
Payload types	IPv6 PPP (POS) AAL5 over ATM PVC GFP-F 10GBASE-LW	IPv6

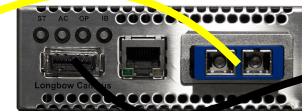
Longbow Campus in production



4X SDR link into
fabric at site A



Up to 10km of
single mode fiber



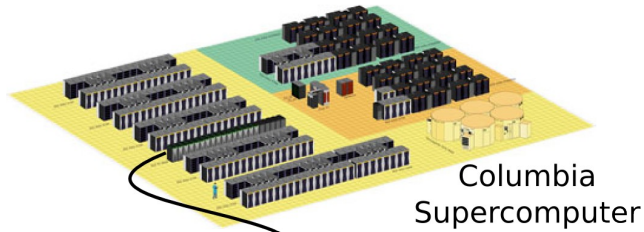
4X SDR link into
fabric at site B

SC|06 NASA/NLR

NASA Ames Research Center - Mountain View, CA

NASA booth (917) Supercomputing 2006 - Tampa, FL

InfiniBand based 3x3 Hyperwall visualization system

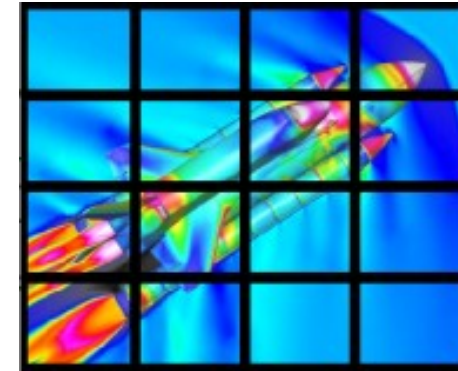


InfiniBand fabric

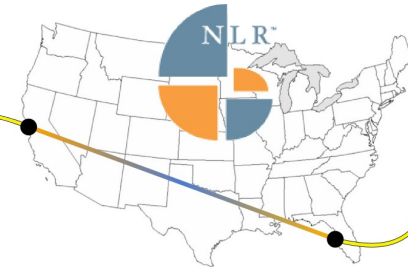
Obsidian Longbow XR



Obsidian Longbow XR



10 Gigabit Ethernet
Cisco 6509 with
IPSec option



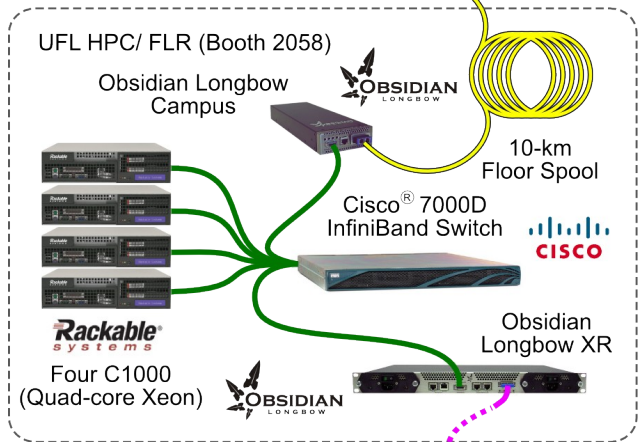
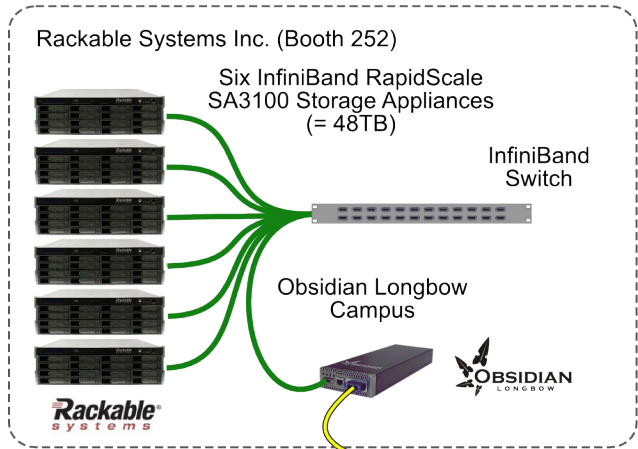
6,600km of 10Gigabit Ethernet layer 3



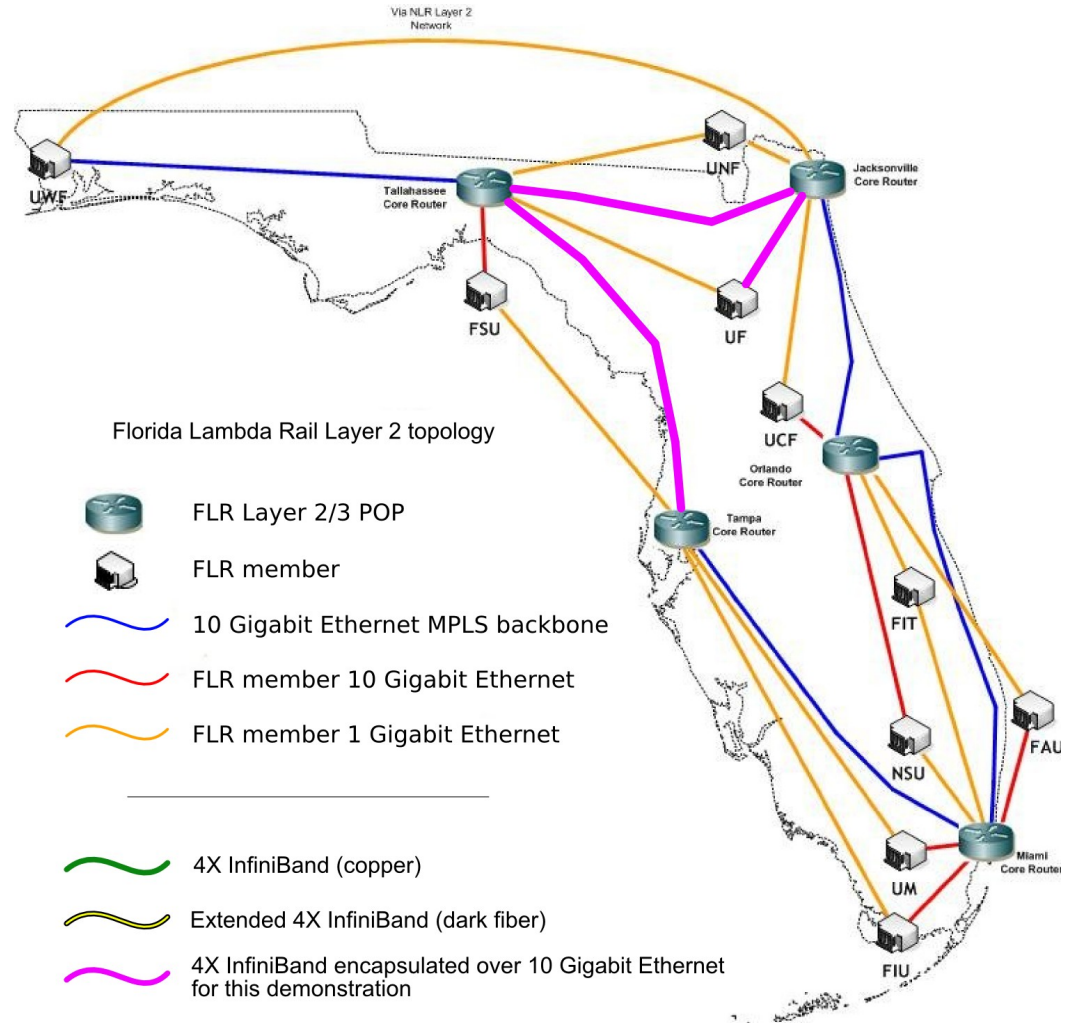
10 Gigabit Ethernet
Cisco 6506 with
IPSec option



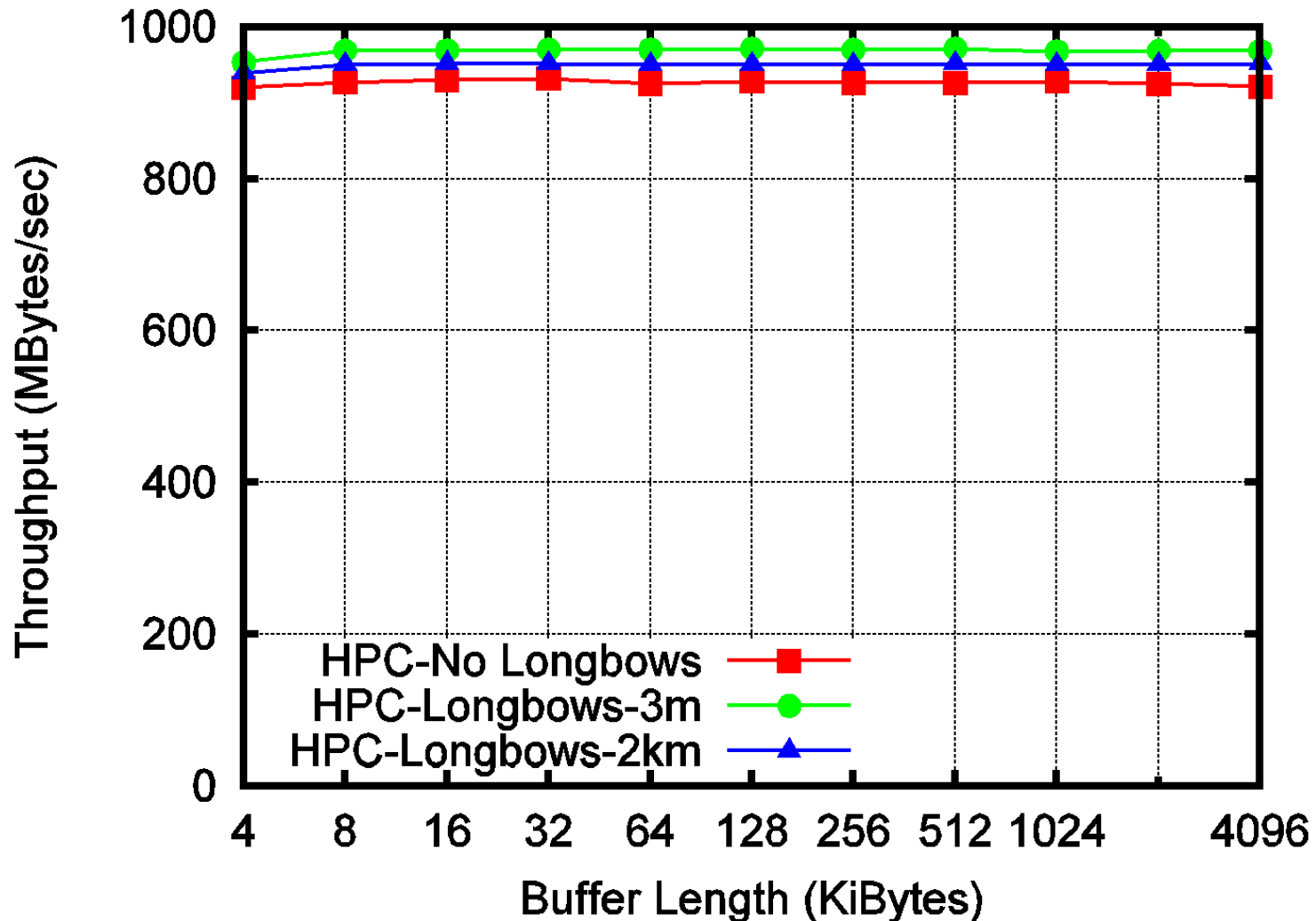
SC|06 UFLHPC/ FLR/ Rackable



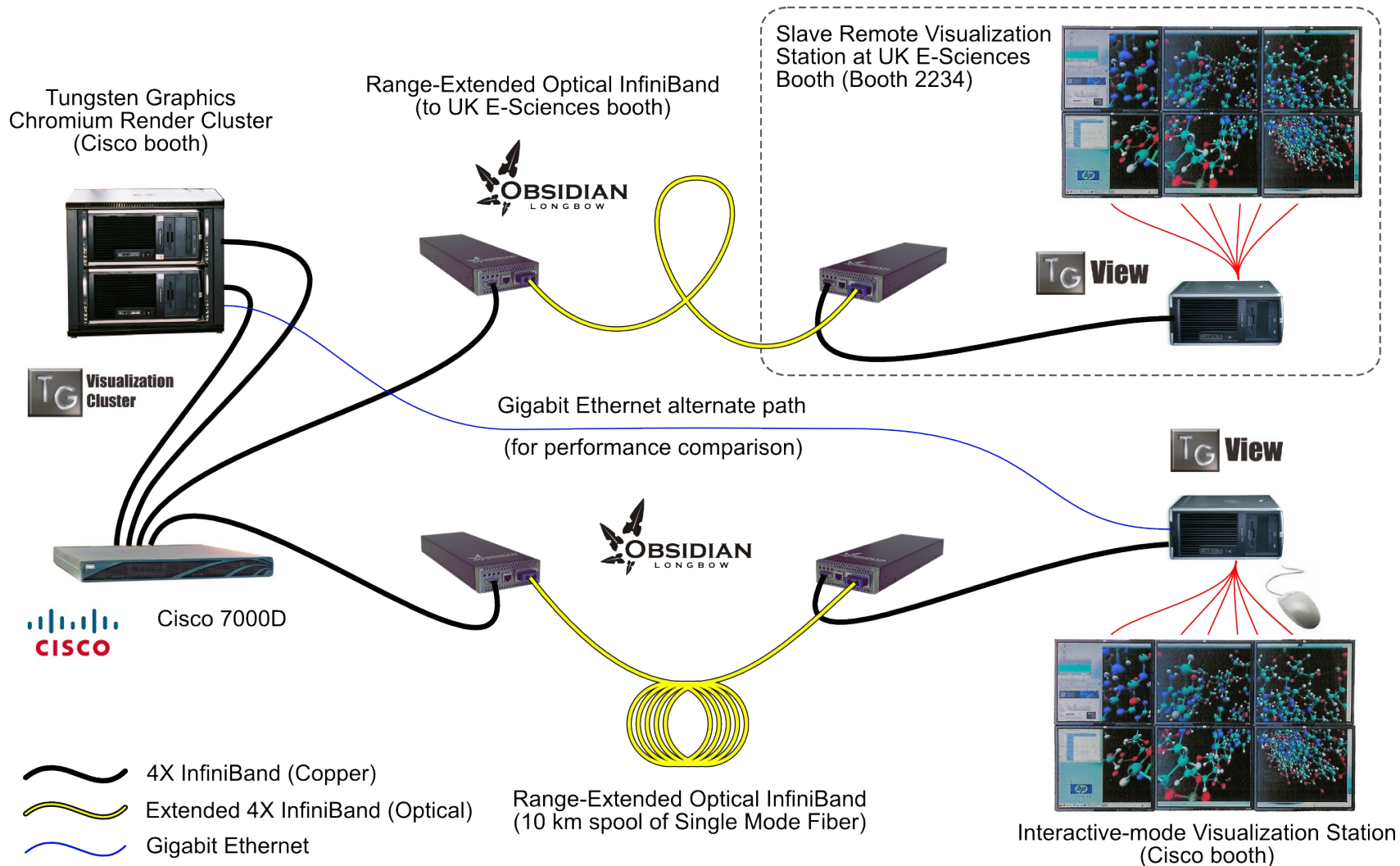
1100km of 10 Gigabit Ethernet to UFL HPC and more Rackable InfiniBand storage...



(UFL HPC Campus)

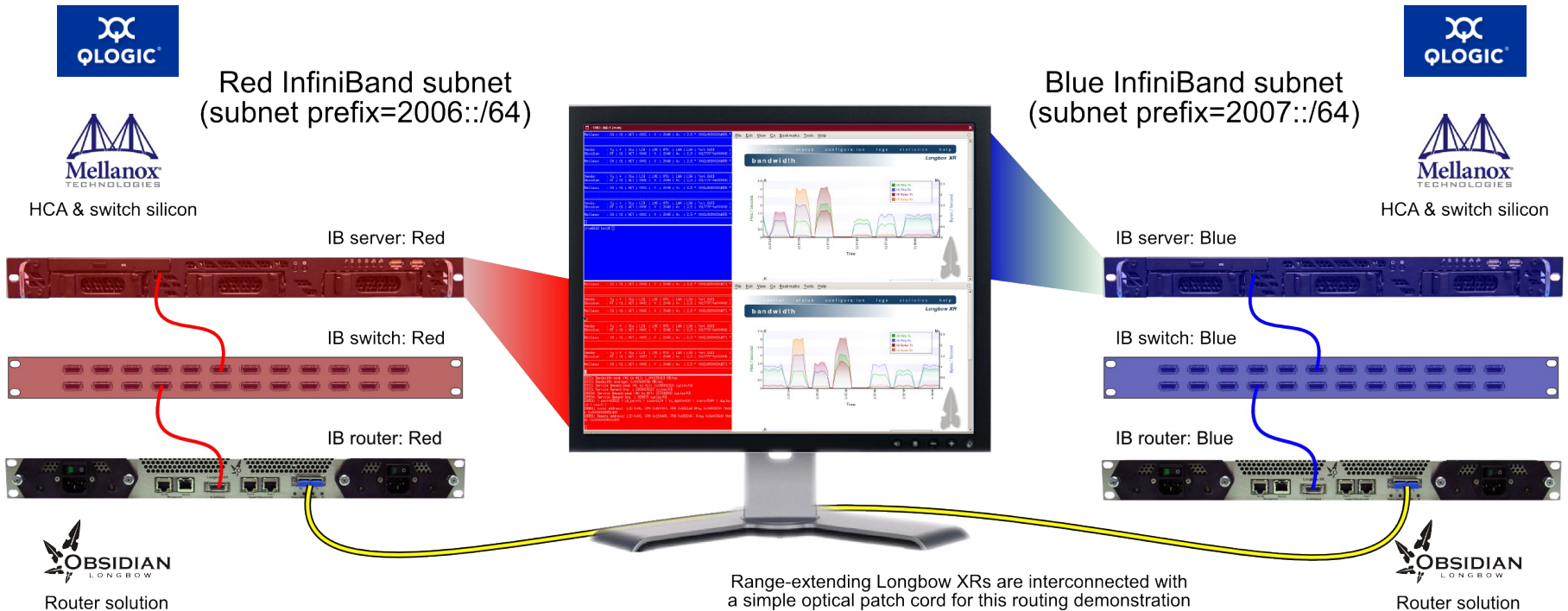


SC|06 Cisco/ TungstenGraphics

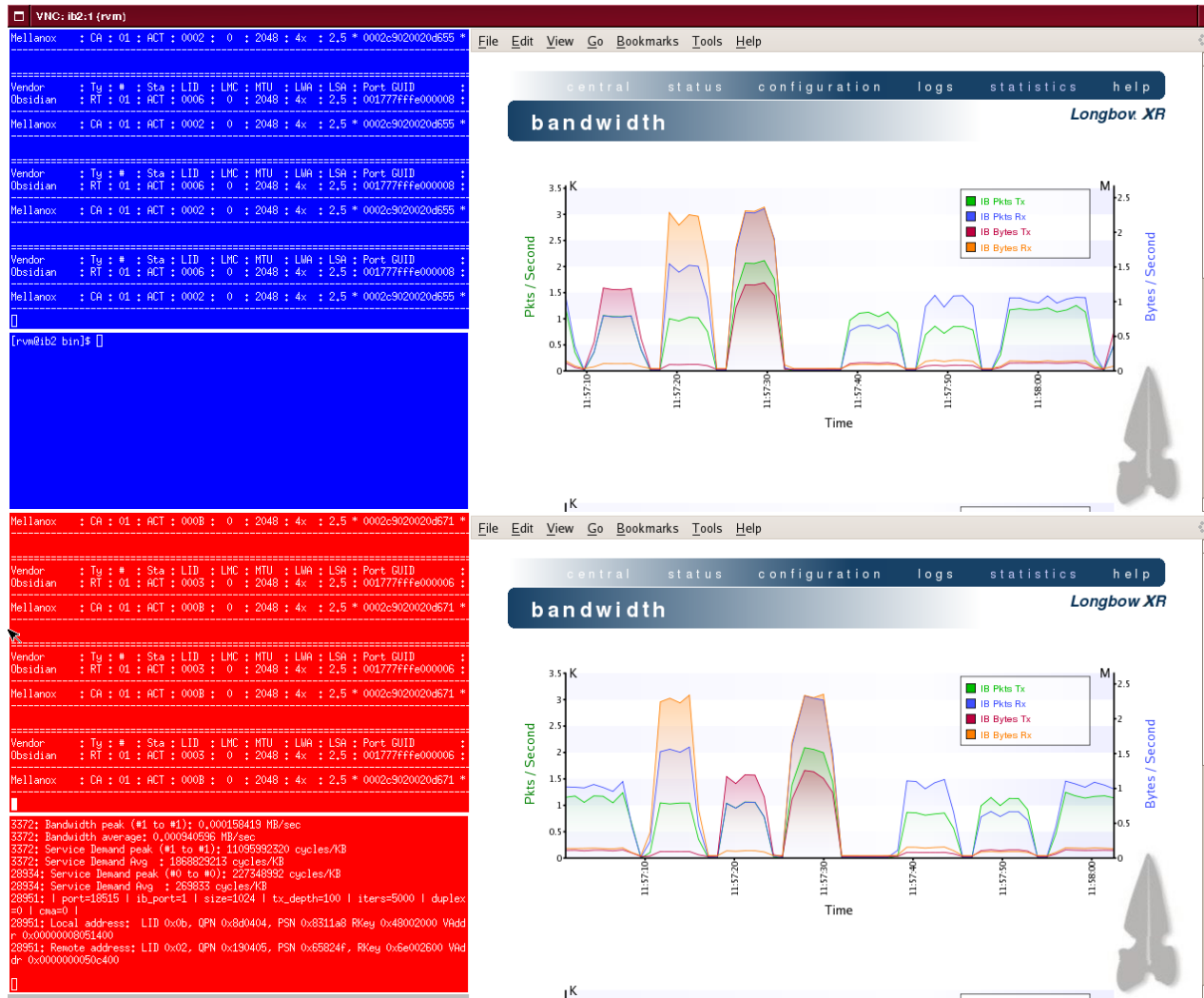




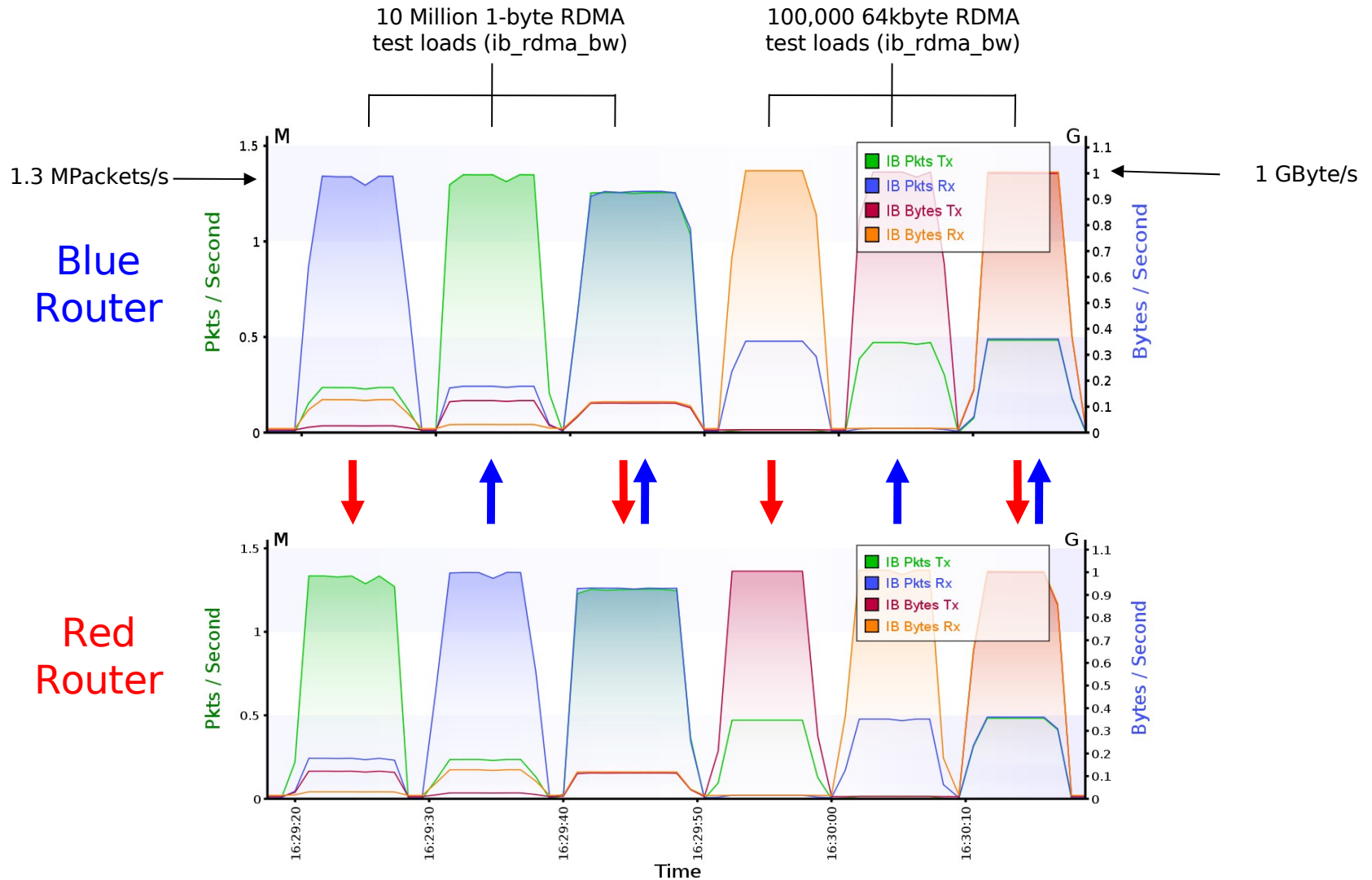
SC|06 Mellanox/ QLogic



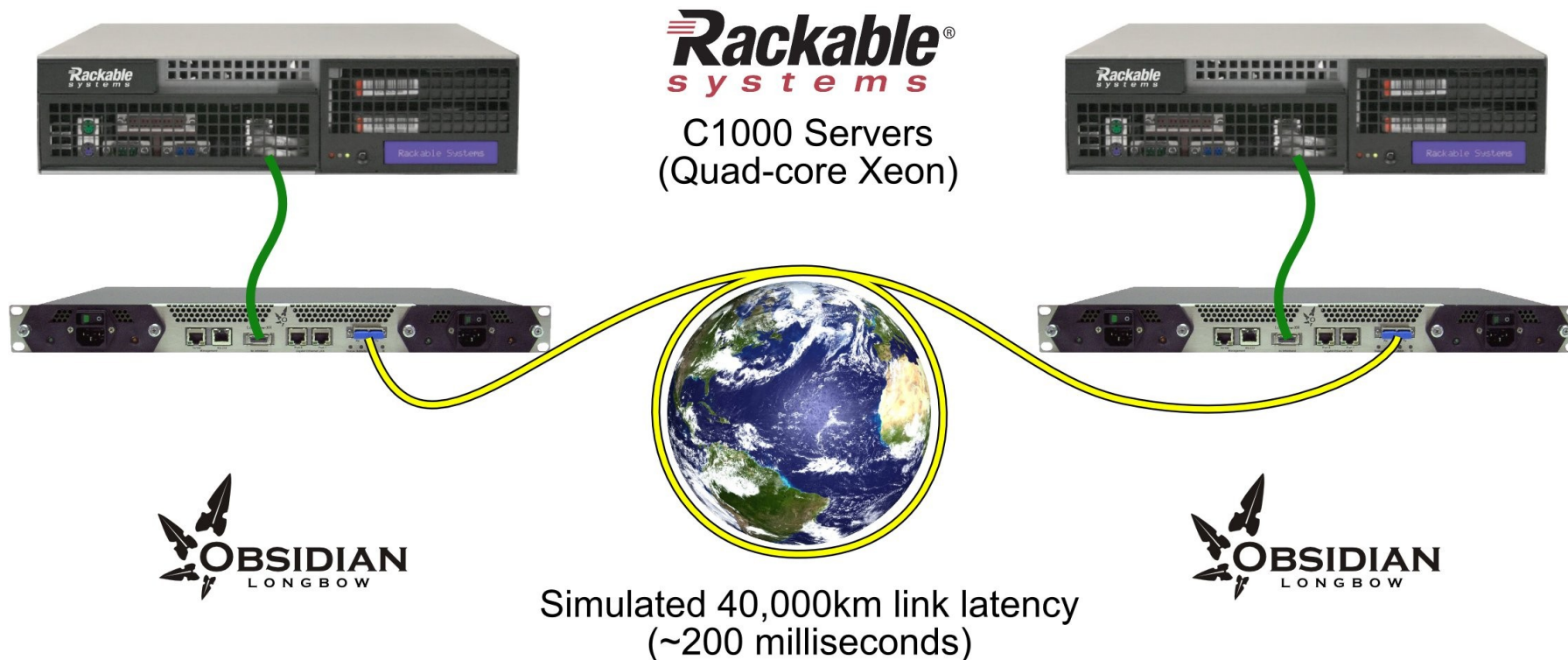
SC|06 Mellanox/ QLogic



SC|06 Mellanox/ QLogic



AFCEA West 2007



Copper 4X SDR InfiniBand



Routed 4X SDR InfiniBand encapsulated over OC-192c ATM



LD JCTD-Concept Of Operations

Advanced Search and Visualization
Advanced data search-and-retrieval to access, integrate, and visualize heterogeneous distributed media, systems, and sites

Better storage and Caching
Integrated, coherent very large-scale (petabytes – 10^{15} to exabytes – 10^{18}) data storage architecture

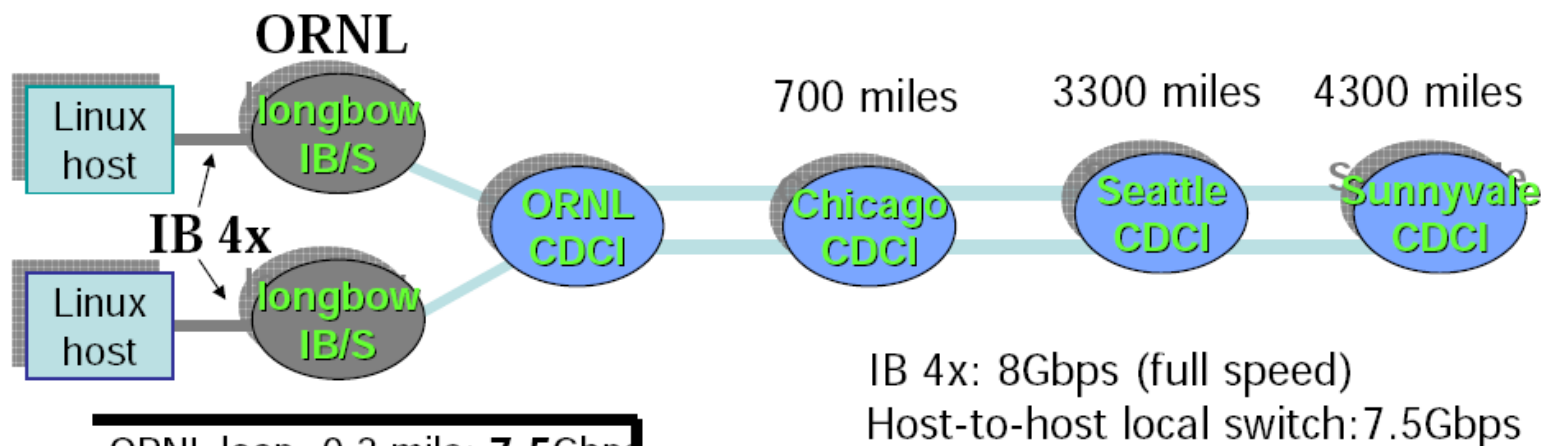
Moving Data to Users

Bigger Pipes
Expanded wideband backbone (10 Gb/sec threshold; 40 Gb/sec objective) linking very large data stores on top of emerging GIG

Infiniband Over SONET

- Infiniband is effective data transport protocol for storage networks (few miles):
- TCP is not easily extended and not optimal for such data transfers

Question: Is IB effective over wide-area? - Yes



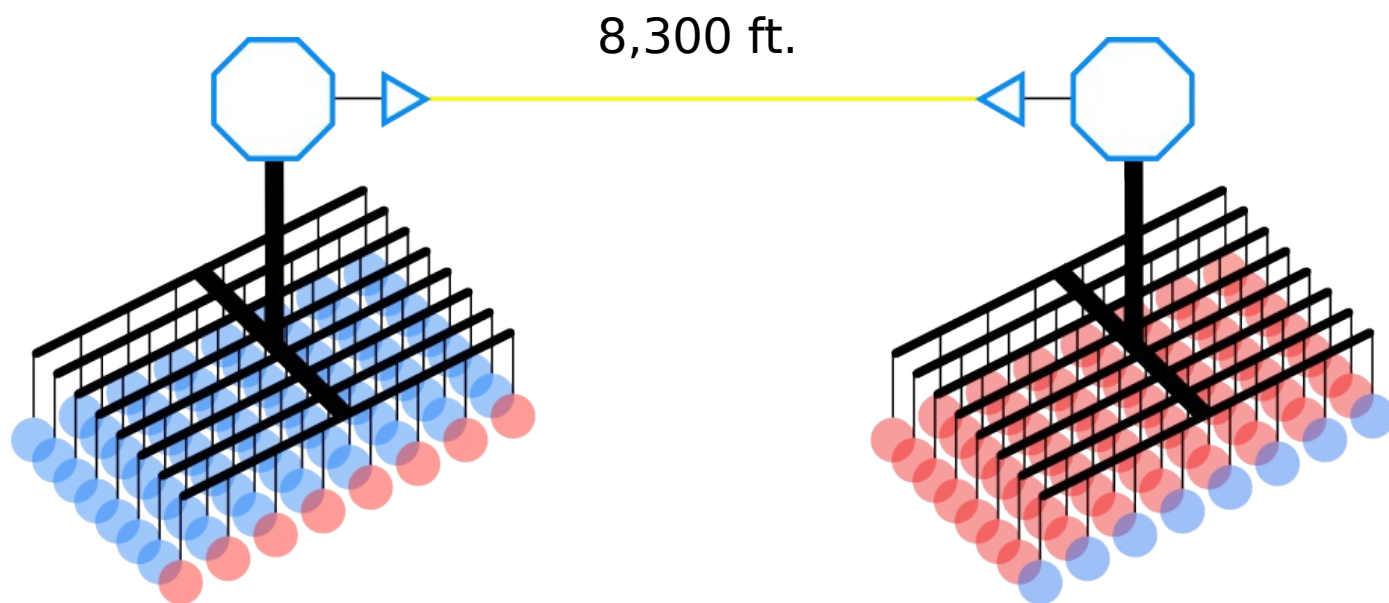
ORNL loop -0.2 mile: **7.5Gbps**

ORNL-Chicago loop – 1400 miles: **7.46Gbps**

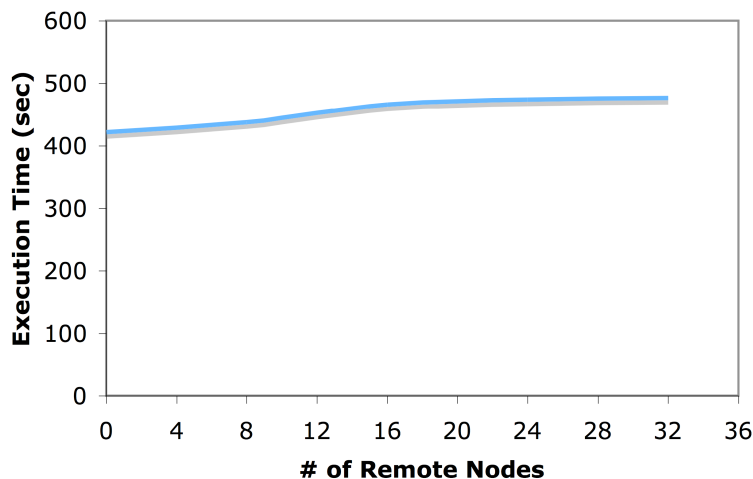
ORNL- Chicago - Seattle loop – 6600 miles: **7.23Gbps**

ORNL – Chicago – Seattle - Sunnyvale loop – 8600 miles: **7.20Gbps**

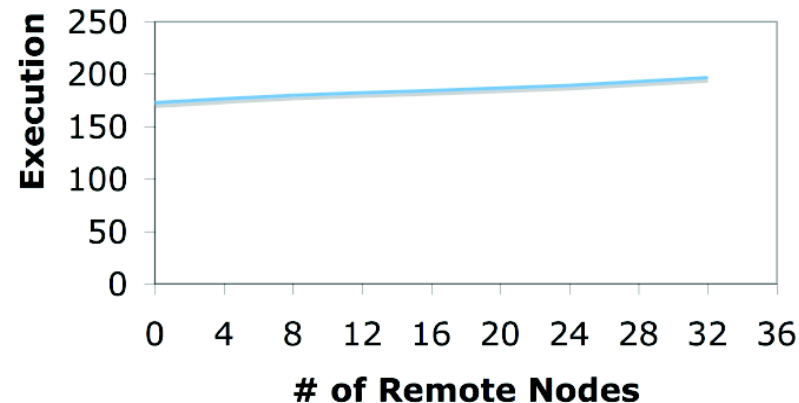
IB 4x: 8Gbps (full speed)
Host-to-host local switch: 7.5Gbps



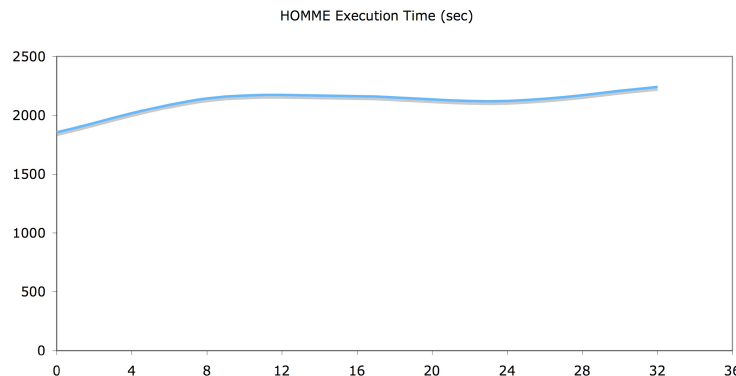
Σ # active cores = 64 (constant)



MIMD Lattice Computation (MILC)

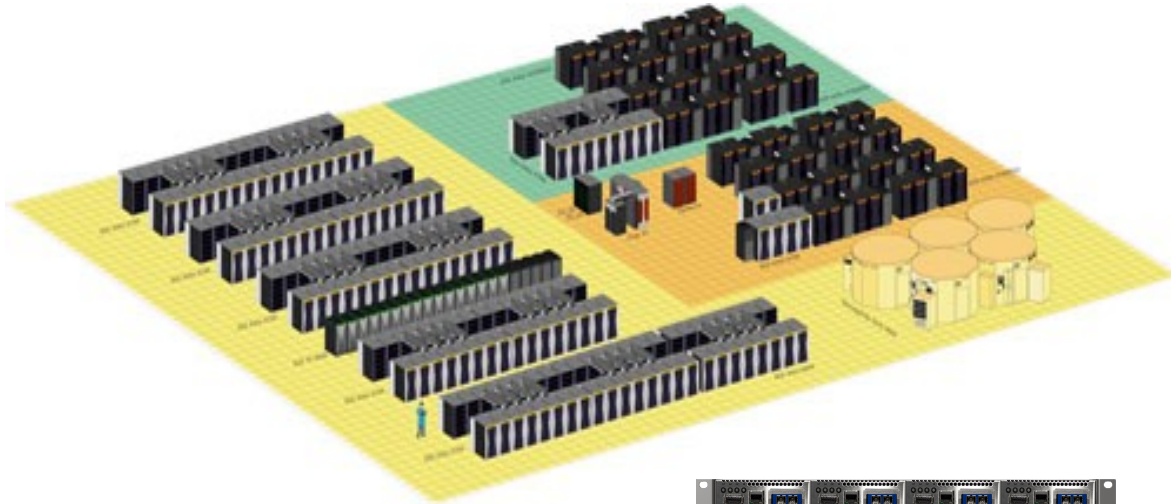


Weather Research and Forecasting (WRF) Code



High Order Methods Modelling Environment (HOMME)

NASA Columbia Expansion



64Gbits/s @ 11.5 μ s latency
(2km of fiber = 10 μ s)



Closing...

New InfiniBand application domains are being validated:

- Efficient WAN transport – reliable RDMA without TCP/IP tuning issues
- Data center replication/ disaster recovery
- High fidelity remote visualization
- Low latency optical networking
- InfiniBand storage area networks
- Cluster clustering – Campus Area Grid
- Supercomputer expansion – building annexation
- Streaming media transport
- Streaming big science transport

...but hurdles remain...

Closing...

We see a strong need for :

- OpenFabrics support for InfiniBand routing for subnet isolation
- Efficient IPoIB (to facilitate selling fat WAN links dedicated to InfiniBand)
- Parallel file system support to tolerate much higher link latencies

Thank You



OPEN**FABRICS**
A L L I A N C E