# Performance Analysis of MPI over InfiniBand on Yellowstone

**Zhengyang Liu**
**Mentor: Dr. John Dennis**

**Collaborators:**
**Prof. Malathi Veeraraghavan (University of Virginia)**
**Prof. Robert D. Russell (University of New Hampshire)**
**Fabrice Mizero (SIParCS)**
**Patrick MacArthur (University of New Hampshire)**

**Aug 1, 2013**

NCAR

UNIVERSITY *of* VIRGINIA

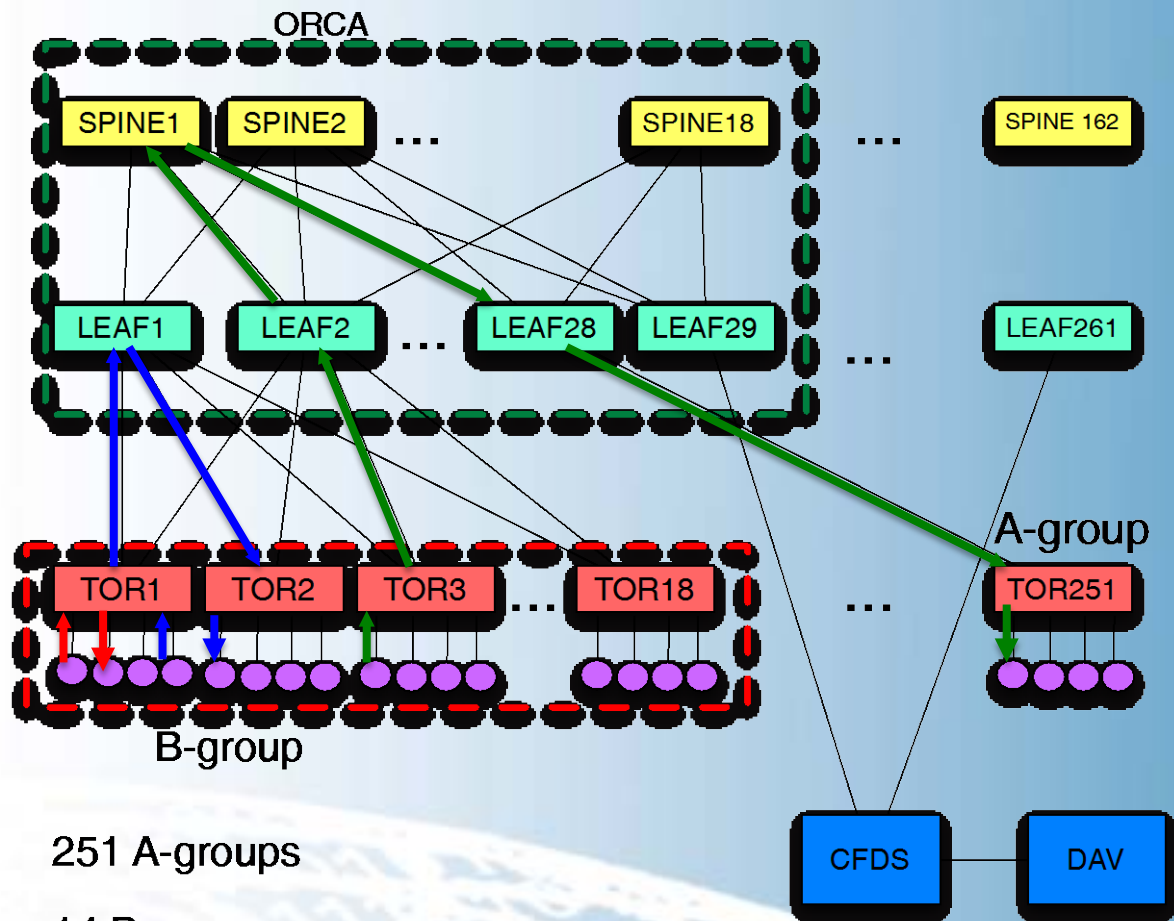Computational & Information Systems Laboratory

# Big Picture

- **Understanding the causes of poor performance of CESM on Yellowstone: a 5-step approach**

  – Experimental execution and data collection

  – CESM trace analysis

  – IBMgtSim: routing study
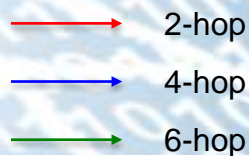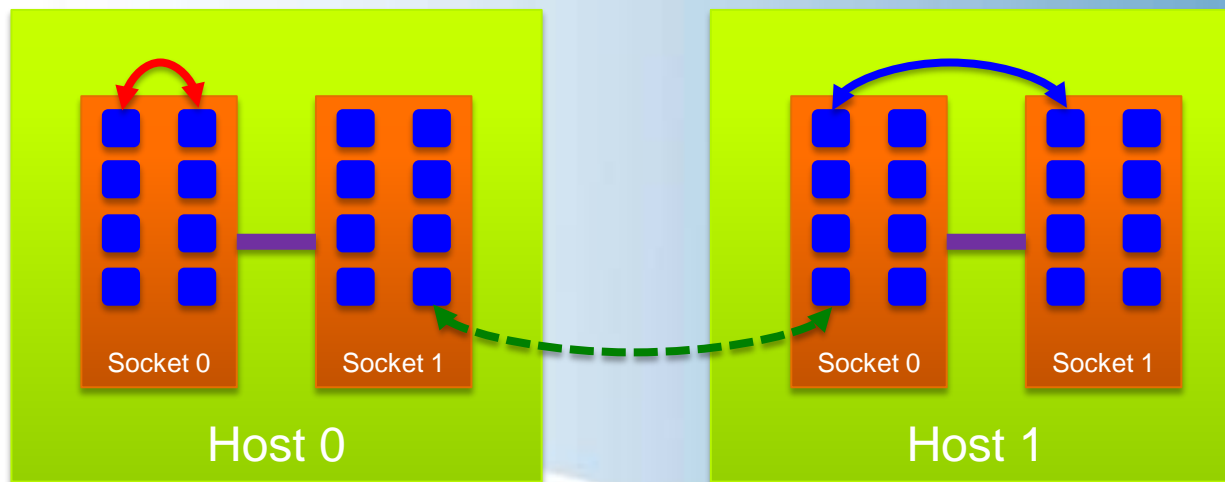
  – Network simulation

  – Integrated simulation

NCAR    UNIVERSITY of VIRGINIA

# Yellowstone network



ORCA

| SPINE1 | SPINE2 | ... | SPINE18 | ... | SPINE 162 |

| LEAF1 | LEAF2 | ... | LEAF28 | LEAF29 | ... | LEAF261 |

A-group

| TOR1 | TOR2 | TOR3 | ... | TOR18 | ... | TOR251 |

B-group

251 A-groups

14 B-groups

9 ORCAs

CFDS — DAV

→ 2-hop
→ 4-hop
→ 6-hop

*Credit: Dr. John Dennis

3

# Communication Patterns



Legend:
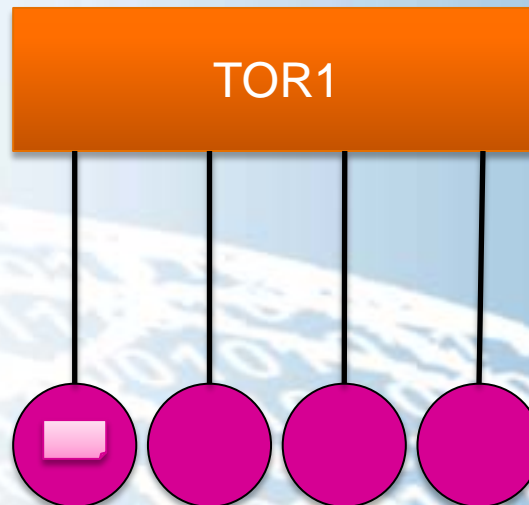- Intra-socket (via shared cache/memory)
- Inter-socket (via shared memory over QuickPath Interconnect)
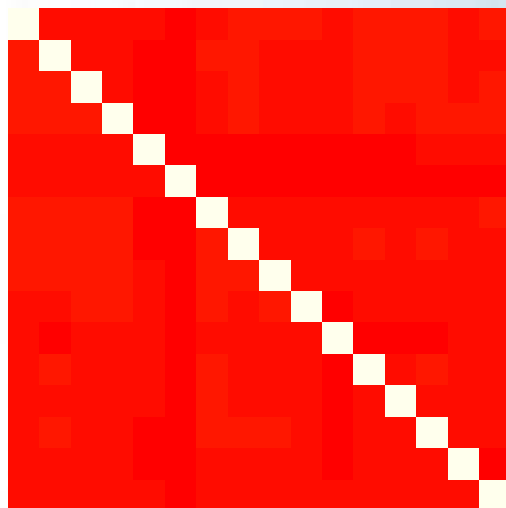- Inter-node (via InfiniBand)

# Latency Benchmark: mpi_pingpong

- **Approximate one-way latency by measuring round-trip latency**
- **Results represent ideal latencies between nodes**
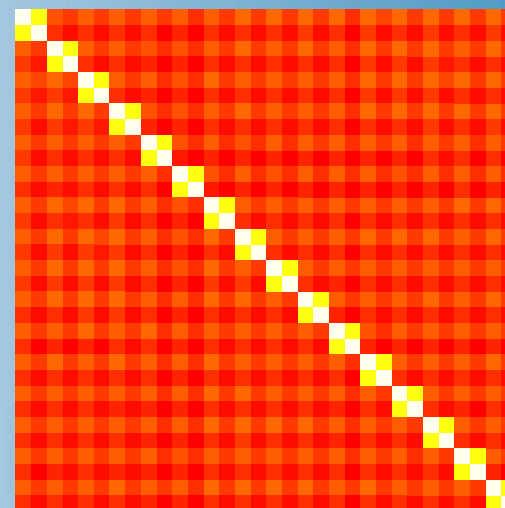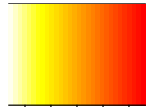
# Jellystone Results
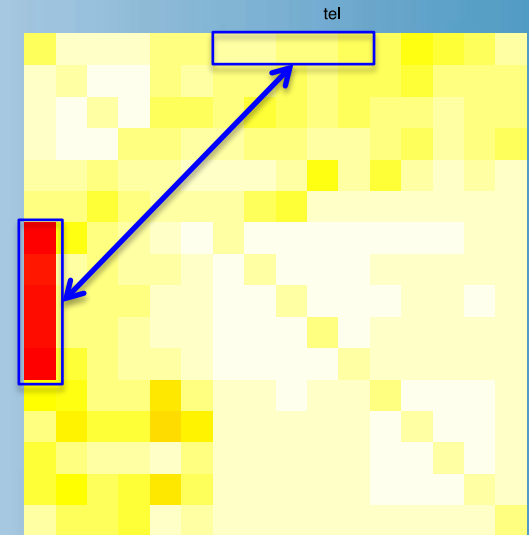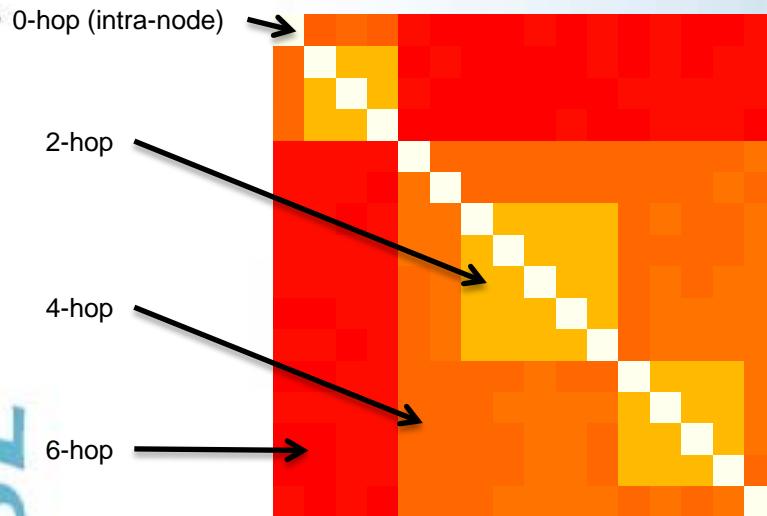
Color Key

0.8  1.2  1.6
Value

Color Key

ue

0-byte messages; all hosts connected to the same TOR; no competing processes

6

# Yellowstone Results

tmap by Node (Median)

0-hop (intra-node)

2-hop

4-hop

6-hop

tel

0-byte messages

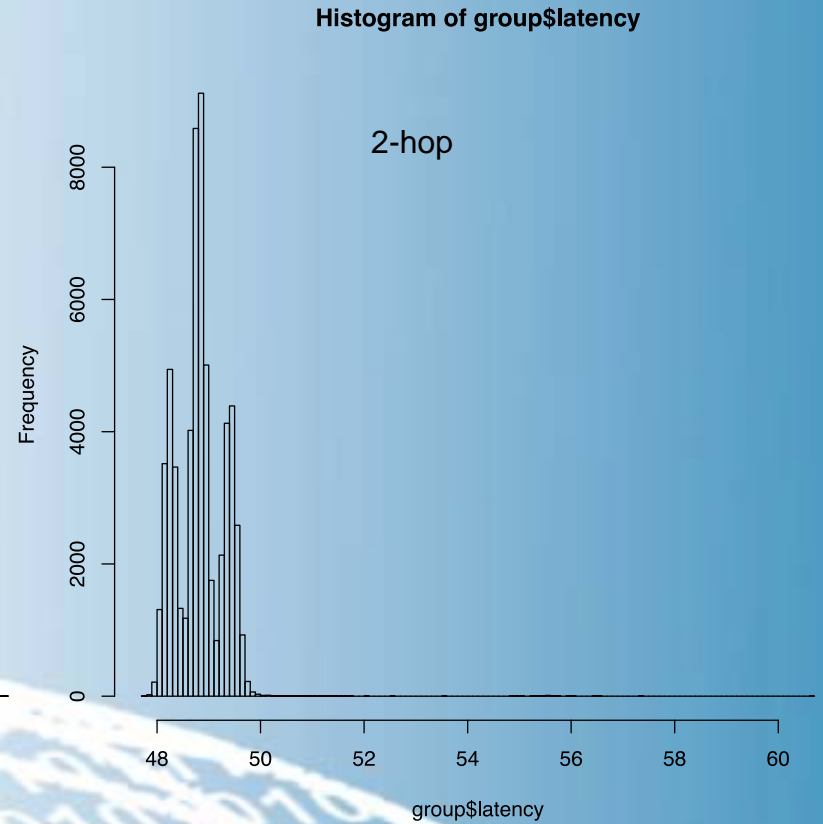256 KB messages
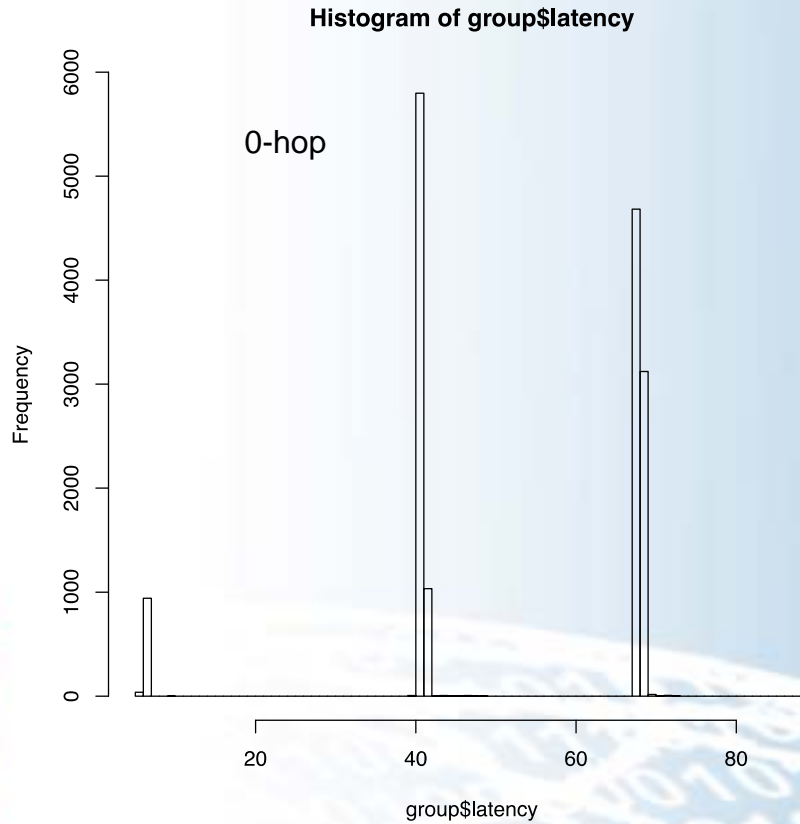
NCAR · NSF · UNIVERSITY of VIRGINIA

# Latency vs. # of Hops

- **Experiment:**

  – mpi_pingpong on 1024 cores

  – 1,048,576 communication pairs*

  – 256 KB messages

Unit: μs

| # of Hops | # of pairs | Min. | Avg. | Max. |
|-----------|-----------|-------|-------|--------|
| 0 | 15,680 | 5.95 | 52.20 | 88.29 |
| 2 | 59,904 | 47.72 | 48.83 | 60.64 |
| 4 | 588,736 | 49.55 | 53.30 | 114.10 |
| 6 | 332,016 | 52.75 | 56.92 | 159.30 |

*routing data available for 996,336 pairs

8

Histogram of group$latency — 0-hop
Histogram of group$latency — 2-hop

**Histogram of group$latency**

4-hop

**Histogram of group$latency**
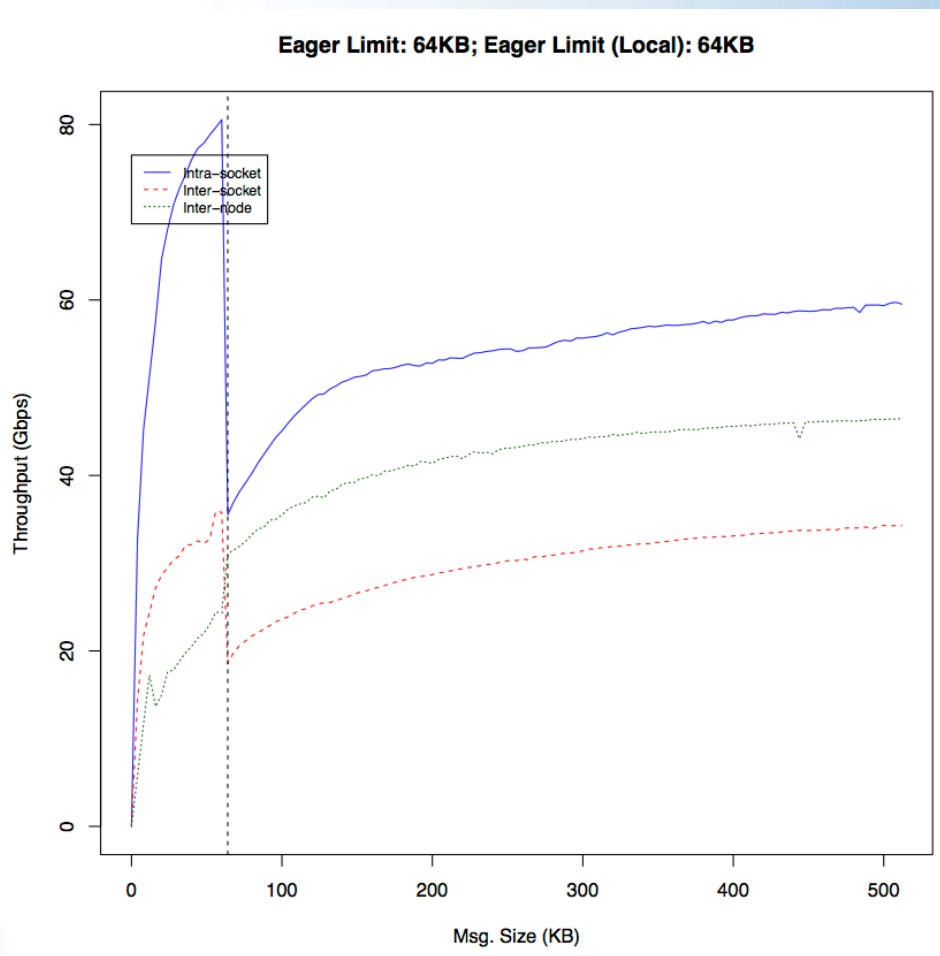
6-hop

# Bandwidth Benchmark: mpi_bw

- **Measures throughput between two MPI ranks**

- **3 communication patterns:**
  - Intra-socket
  - Inter-socket
  - Inter-node

- **2 communication protocols:**
  - Eager protocol
  - Rendezvous protocol

# Communication Protocols

- **Rendezvous Protocol: buffer negotiation before sending**

- **Eager Protocol: send directly without confirming available buffer space**

- **InfiniBand: Eager protocol uses SEND/RECV verbs (two-sided communication); Rendezvous protocol uses WRITE/READ verbs (one-sided communication)**

- **Eager Limit: threshold below which Eager protocol is used**

# Jellystone Results



Eager Limit: 64KB; Eager Limit (Local): 64KB

- **Intra-node throughput decreases when msg. size > eager limit**

- **Inter-node throughput increases when msg. size > eager limit**

- **Inter-node communication faster than inter-socket communication: RDMA vs shared memory**

# Summary

- **Identified contention through mpi_pingpong benchmarks**
- **Studied effect of different communication patterns/protocols on throughput**

# Future Work

- **Analyses of larger data sets**
  - \> 500 million data points
  - Analysis needs to be parallelized
- **Study interaction between MPI and InfiniBand**
  - Open-source MPI implementations
  - Network sniffing

# Thank You

Zhengyang Liu
zl4ef@virginia.edu