# Performance Evaluation of Open MPI on Cray XE/XK Systems

**Samuel K. Gutierrez – LANL**

**Nathan T. Hjelm – LANL**

**Manjunath Gorentla Venkata – ORNL**

**Richard L. Graham – ORNL**

**Hot Interconnects 2012**

**Aug 23, 2012**

# A Collaborative Effort

Thursday, August 23, 12

# Outline

**Los Alamos**
NATIONAL LABORATORY
EST.1943

U N C L A S S I F I E D – LA-UR-12-24229

Operated by Los Alamos National Security, LLC for the U.S. Department of Energy's NNSA

**NNSA**

Thursday, August 23, 12

# First Things First – Open MPI Overview

- **Open-Source Implementation of the MPI-2 Standard**

- **Developed and Maintained By**
  - Academia
  - Industry
  - National Laboratories

- **Supports a Range of High-Performance Network APIs**
  - Verbs (Infiniband, RoCE, iWarp)
  - PSM (QLogic/Intel HCAs)
  - MXM (Mellanox HCAs)
  - Portals (Cray SeaStar, Infiniband)
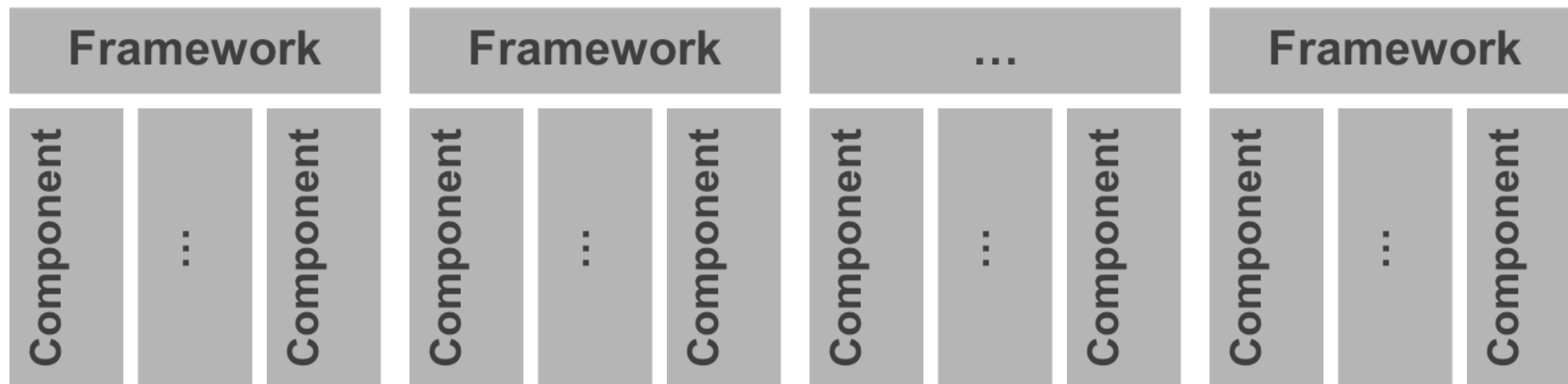  - uGNI (Cray Gemini, Cray Ares)

**Los Alamos**
NATIONAL LABORATORY
EST.1943

Operated by Los Alamos National Security, LLC for the U.S. Department of Energy's NNSA

**U N C L A S S I F I E D – LA-UR-12-24229**

*Slide 4*

NNSA

Thursday, August 23, 12

**User Application**

**MPI API**

**Modular Component Architecture (MCA)**

| Framework | Framework | … | Framework |
|---|---|---|---|
| Component ⋮ Component | Component ⋮ Component | Component ⋮ Component | Component ⋮ Component |

**Los Alamos**
NATIONAL LABORATORY
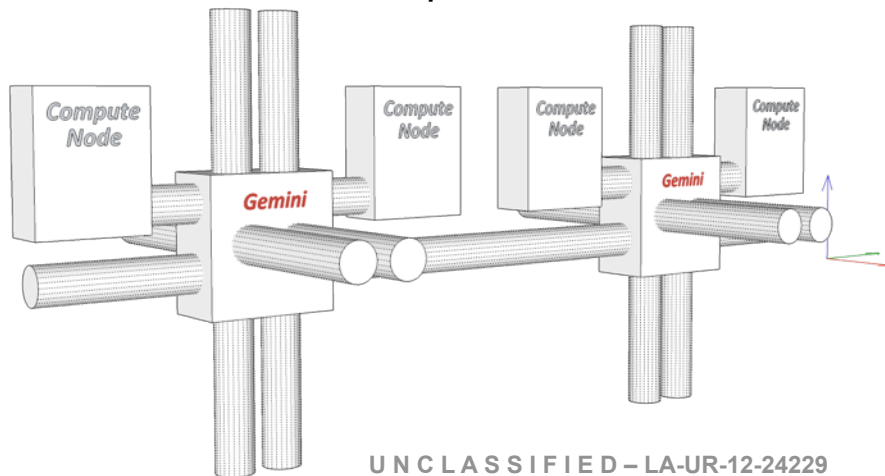EST. 1943

# Open MPI's Plugin Architecture – Main Code Sections[1]

- **Open MPI Layer (OMPI)**
  - MPI API and Support Logic

- **Open Run-Time Environment (ORTE)**
  - Run-Time System

- **Open Portability Access Layer (OPAL)**
  - OS-Specific/Utility Code

**OMPI**

**ORTE**

**OPAL**

**Operating System**

U N C L A S S I F I E D – LA-UR-12-24229

Thursday, August 23, 12

# The Gemini System Interconnect[3] – An Overview

- **Network Used by the Cray XE and XK System Families**
  - Titan, Cielo, Hopper

- **Successor to the Cray SeaStar\* Network Interconnect**

- **3D Torus Network Built of Gemini ASICs**

- **Gemini ASIC**
  - Provides 10 Torus Connections – 2 x (+X, -X, +Z, -Z) – 1 x (+Y, -Y)
  - Provides 2 NICs and a 48-port Router

Operated by Los Alamos National Security, LLC for the U.S. Department of Energy's NNSA

Thursday, August 23, 12

# OB1 PML High-Level Protocol Overview

- **Eager Message Protocol**
  - Uses BTL buffered, inline, and in-place send protocols

- **Remote Get Protocol**
  - 2 Protocol Messages: RGET (ready to send + segment), FIN
  - Available When Registration Cache is Enabled and BTL Implements Get

- **RDMA Pipeline Protocol (Put)**
  - 3 Protocol Messages: RNDV + segment, RDMA, FIN
  - Used When Remote Get protocol is not Available

- **Remote Get Fallback (New)**
  - Essentially a Rendezvous
  - Fallback Initiated by the Receiver During Remote Get Protocol if BTL Get Protocol is not Available

- **Rendezvous (no RDMA)**

Thursday, August 23, 12
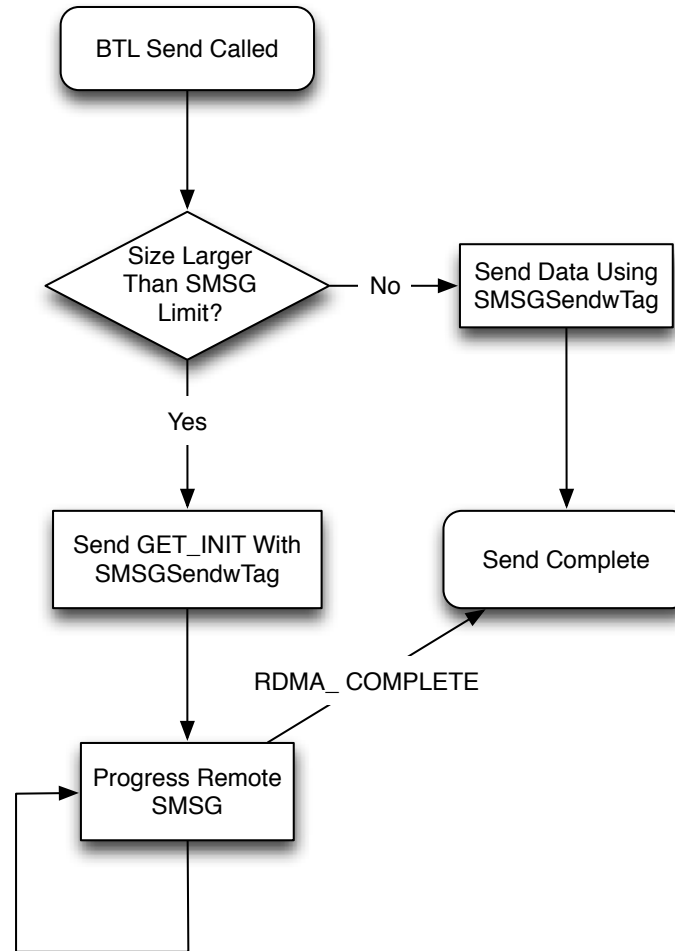
# uGNI BTL Overview

- **Protocols**
  - Send
    - In-place Send for Small Messages Directly Using Small Message Protocol (SMSG)
    - Buffered Send Using Get for Larger Eager Messages (Eager Get)
  - Get
    - Uses FMA Or BTE
    - Available Only if Source And Destination Segments Are 4-Byte aligned and a Multiple of 4-Bytes in Size
  - Put
    - Uses Fast Memory Access (FMA) or Byte Transport Engine (BTE)
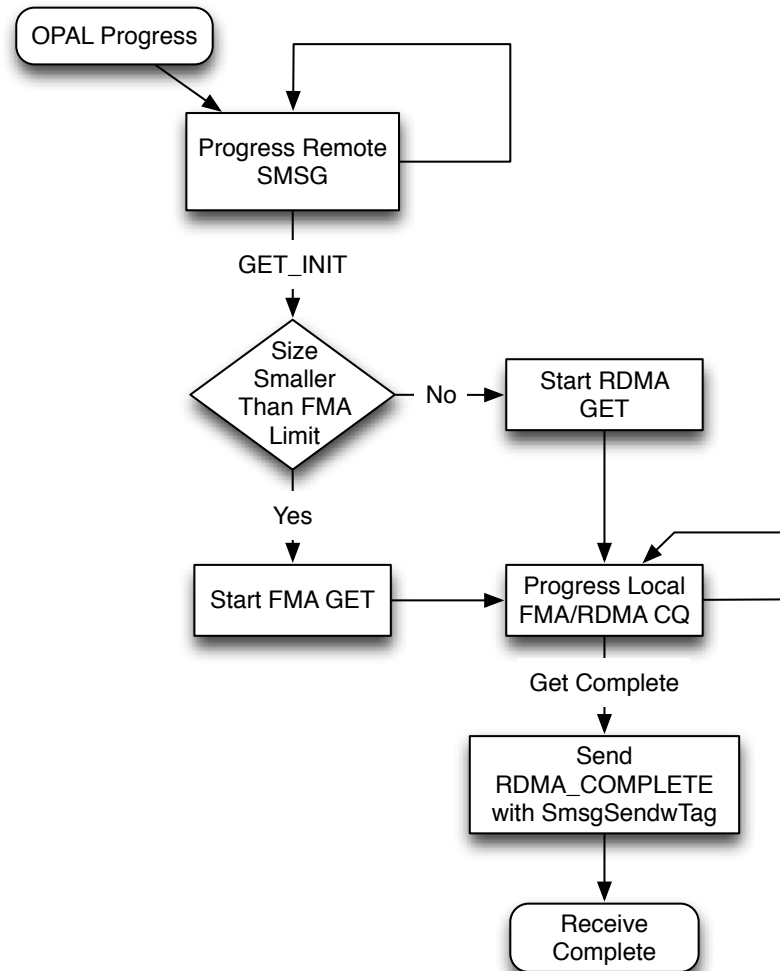    - No Alignment Restrictions

- **Lazy Connection Establishment**
  - Resource Utilization Directly Related to Application Communication Characteristics

Thursday, August 23, 12

# uGNI BTL Eager Get Protocol Details (Send)

Thursday, August 23, 12

# uGNI BTL Eager Get Protocol Details (Receive)



OPAL Progress

Progress Remote SMSG

GET_INIT

Size Smaller Than FMA Limit

No → Start RDMA GET

Yes

Start FMA GET

Progress Local FMA/RDMA CQ

Get Complete

Send RDMA_COMPLETE with SmsgSendwTag

Receive Complete

Operated by Los Alamos National Security, LLC for the U.S. Department of Energy's NNSA

Thursday, August 23, 12

# Vader BTL Overview

- **MPICH Nemesis-like Design**
  - Lock-Free Message Queues
  - "Fast Boxes" – I.e. Per-Peer Receive Queues for Short Messages

- **Copy Backend Changes Based on Message Size**
  - E.g. *bcopy [a,b)* - *memcpy* Otherwise
  - User Tunable with *Good* Defaults

- **Cross-Process Memory Mapping Allows for RDMA-Like Semantics**
  - Copy-In/Copy-Out (CICO) Avoided
  - No Backing Store Required
  - Heavy Use of Registration Cache to Amortize Attach Latency
  - Exposes Both Put and Get Interfaces to PML Layer

- **XPMEM Support Requires Kernel Patch and User-Level Library**
  - Already Available and Leveraged by Cray's Native MPI Implementation

**Los Alamos**
NATIONAL LABORATORY
EST.1943

U N C L A S S I F I E D – LA-UR-12-24229

*Slide 12*

Operated by Los Alamos National Security, LLC for the U.S. Department of Energy's NNSA

**NNSA**

Thursday, August 23, 12

# Test Environment

- **Testing Platforms**
  - **Cielito** - 1088 Core XE6
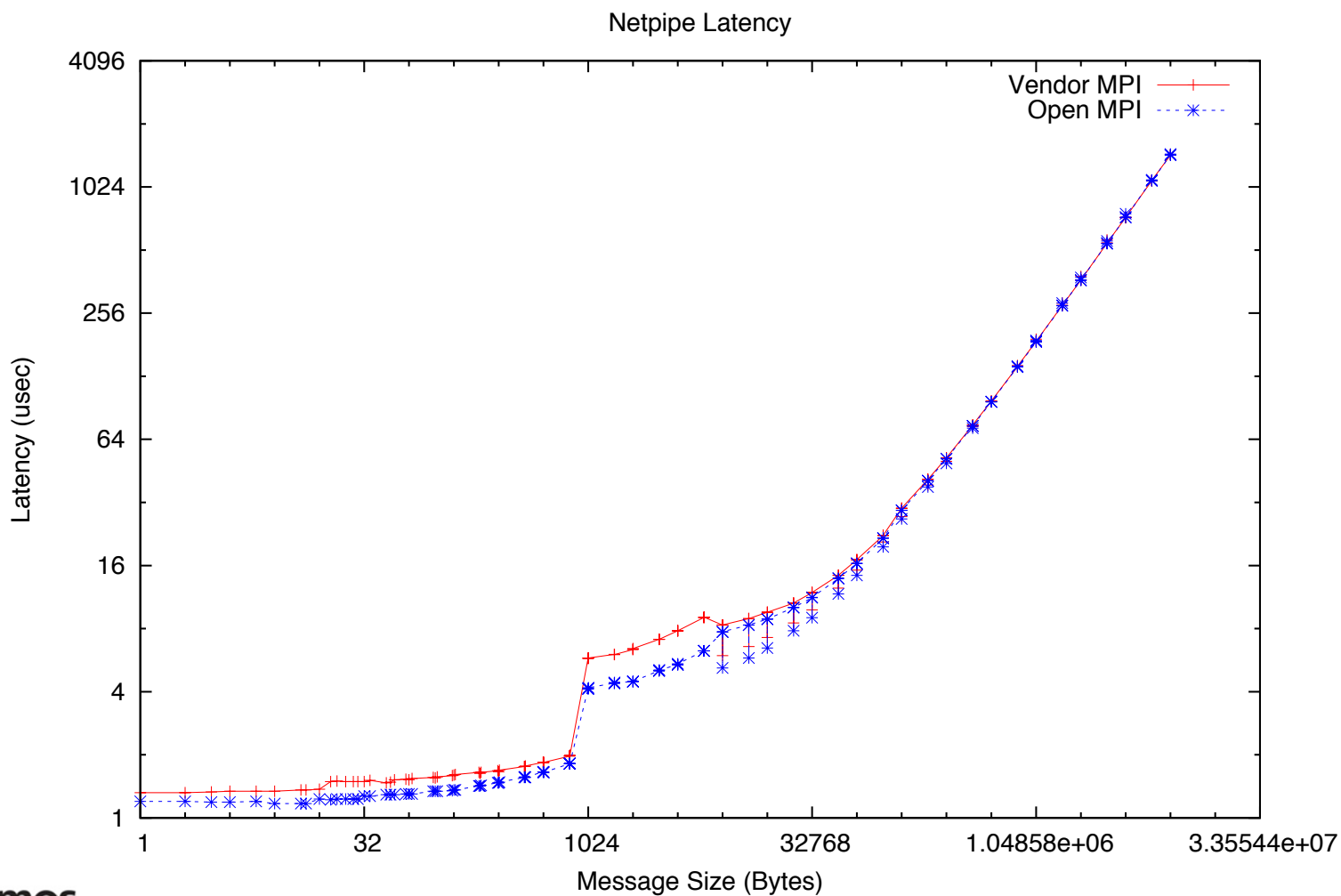  - **Cielo** - 142,304 Core XE6

- **Microbenchmarks**
  - **NetPIPE –** Measure Lat/BW Benchmark
  - **AMG2006 –** Algebraic Multi-grid Solver
  - **LAMMPS –** Classical Molecular Dynamics Code
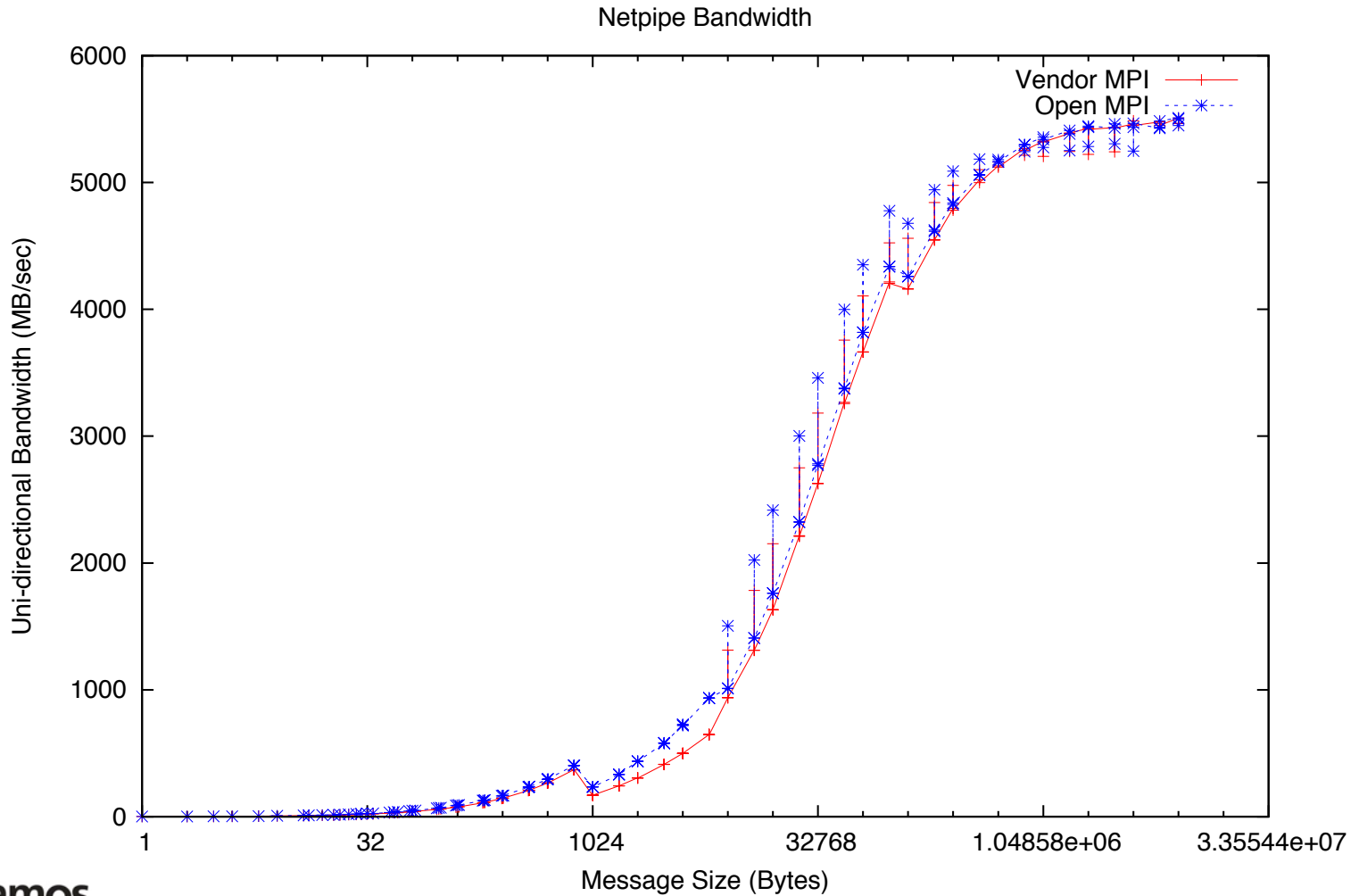  - All Microbenchmarks Were Run on Live Production System

- **Launcher**
  - orterun

Thursday, August 23, 12

# NetPIPE Latency on XE6 (on ASIC)



Netpipe Latency

Los Alamos
NATIONAL LABORATORY
EST.1943

Operated by Los Alamos National Security, LLC for the U.S. Department of Energy's NNSA

NNSA

Thursday, August 23, 12

# NetPIPE Bandwidth on XE6 (on ASIC)



Netpipe Bandwidth

Thursday, August 23, 12

# Microbenchmark – LAMMPS



LAMMPS ASC Benchmark for Scaled-size Lennard-Jones Liquid

Thursday, August 23, 12

# Microbenchmark – AMG2006



AMG2006 ASC Benchmark (3D 7-Point Laplace Problem on a Cube)

Operated by Los Alamos National Security, LLC for the U.S. Department of Energy's NNSA

Thursday, August 23, 12

# Conclusion and Ongoing/Future Work

- **Conclusion**
  - Bandwidth, Latency, and Scalability Similar to Vendor MPI Implementation

- **Stabilization/Optimization**
  - Improve Launch Scalability (Over a Minute to Launch 131072 MPI Tasks)
  - Investigating New Protocols (Shared Message Queue-- MSGQ)
  - Reduce Memory Requirements

- **Improved Collective Performance Using uGNI Atomics**

- **Work with Friendly Testers**

- **Prepare for General Release in Open MPI 1.7.0**

Los Alamos
NATIONAL LABORATORY
EST.1943

Operated by Los Alamos National Security, LLC for the U.S. Department of Energy's NNSA

NNSA

Thursday, August 23, 12

# Thanks!

Thursday, August 23, 12

# Questions?

- **Questions?**

- **Comments?**

Thursday, August 23, 12

# References

- **[1] <u>Open MPI</u>. 13 Feb. 2012 <open-mpi.org>.**

- **[2] R. Alverson, et al., "The Gemini System Interconnect," in High Performance Interconnects (HOTI), 2010 IEEE 18th Annual Symposium on, Aug. 2010, pp. 83 –87.**

**Los Alamos**
NATIONAL LABORATORY
EST.1943

**U N C L A S S I F I E D – LA-UR-12-24229**

*Slide 21*

Operated by Los Alamos National Security, LLC for the U.S. Department of Energy's NNSA

Thursday, August 23, 12