# Cheetah: A Framework for Scalable Collective Operations

Project members

## Richard L. Graham
## Pavel Shamis
## Joshua S. Ladd
## Manjunath Gorentla Venkata

Computer Science and Mathematics Division
Oak Ridge National Laboratory
Work supported by US DOE ASCR FASTOS
and Math/CS Institute EASI!

# Need for Efficient Collectives

**What are Collectives ?**

**A global communication operation performed over all processes of a parallel application**

- **A large percentage of High Performance Computing (HPC) application execution time is spent in global communication operations (collectives)**

- **Moving towards exascale systems, the time spent in collectives only increases (unless the applications change)**

- **Collectives are basic building blocks for many parallel programming languages and communication libraries**

# Collectives on Modern HPC Systems

**Modern HPC systems have multiple data paths for communication, the collectives should take advantage of the system architecture to improve the performance of collectives**
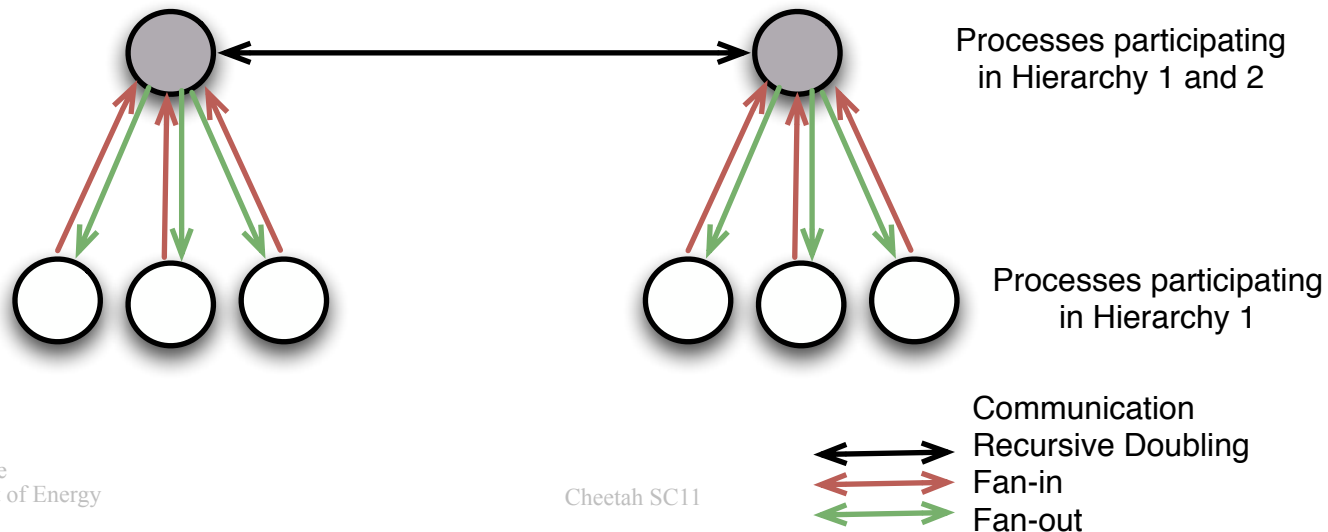
## Previous Approaches

- Collectives optimized by replacing point-to-point communication by shared-memory communication

- Collectives optimized for a particular architecture
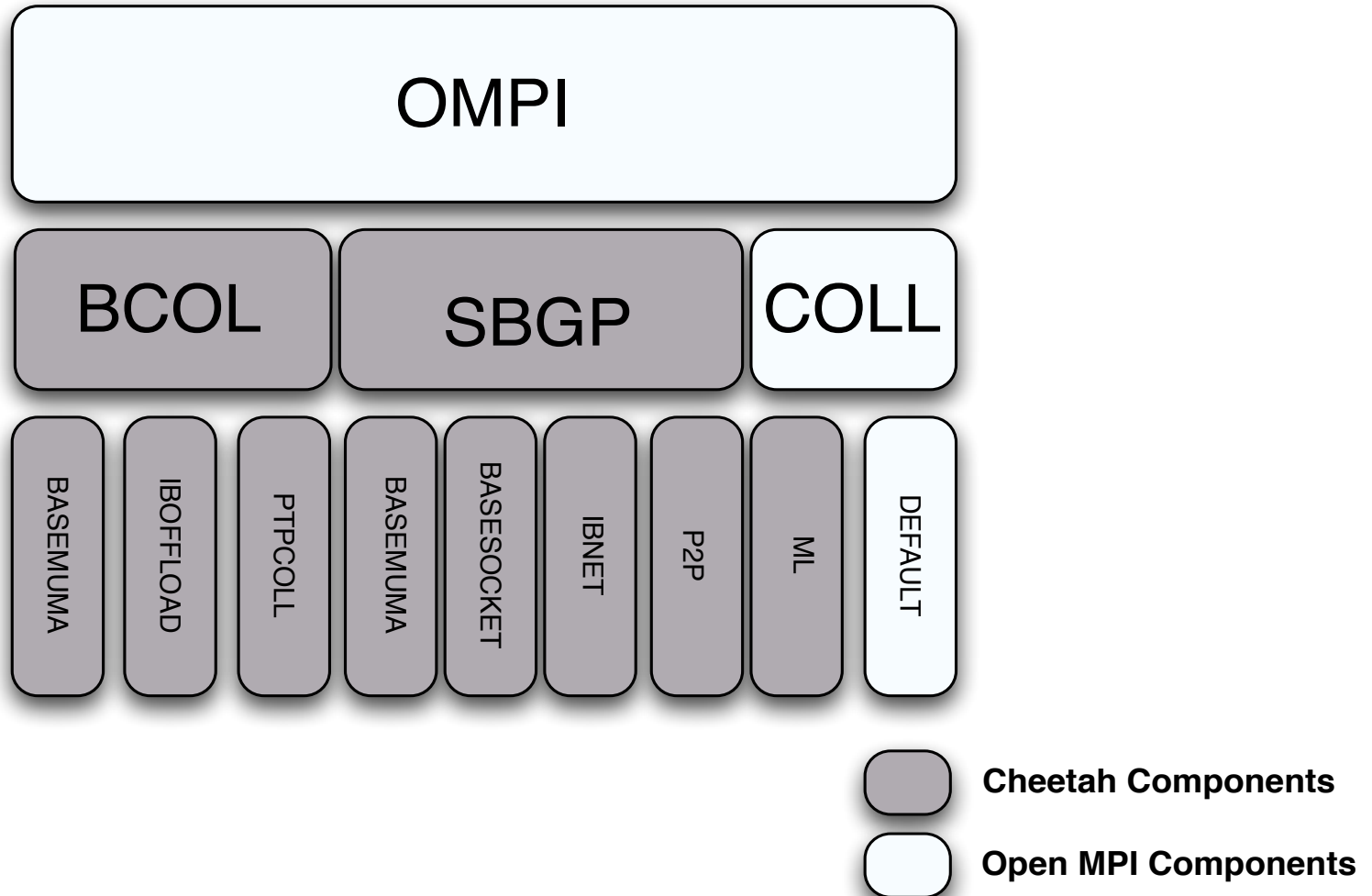
- Collectives with limited hierarchy support

Cheetah SC11

# Cheetah: A Framework for Scalable Collective Operations

- **Cheetah collectives are implemented as a combination of multiple collective primitives**

- **Collective primitive is optimized for a particular data communication path (communication hierarchy)**

- **Collective primitives are progressed asynchronously and independently (when semantics permit)**

**Example: A n-level hierarchial Barrier is a combination of fan-in, fan-out and Barrier collective primitives**



Processes participating in Hierarchy 1 and 2

Processes participating in Hierarchy 1

Communication
Recursive Doubling
Fan-in
Fan-out

OAK RIDGE
National Laboratory

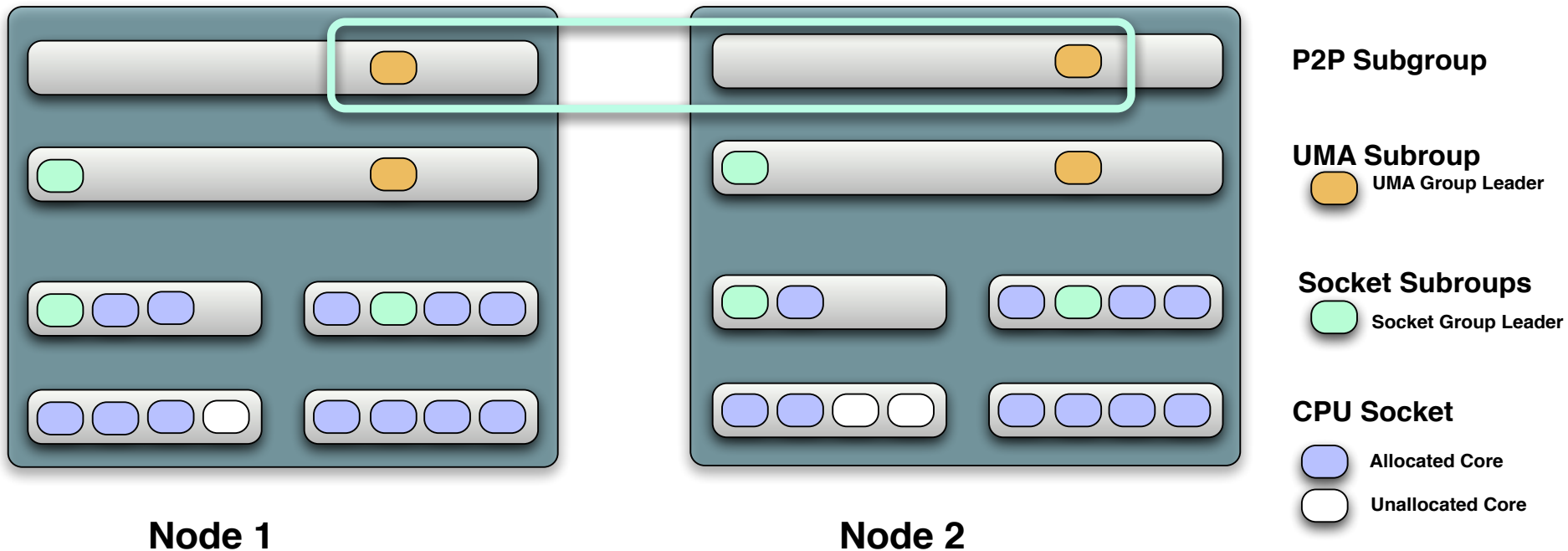# Cheetah is Implemented as a Part of Open MPI

Cheetah SC11

# Cheetah Components and Functionality

- **Base Collectives (BCOL) – Implements basic collective primitives**

- **Subgrouping (SBGP) – Provides rules for grouping the processes**

- **Multilevel (ML) – Coordinates collective primitive execution, manages data and control buffers, and maps MPI semantics to BCOL primitives**

- **Schedule – Defines the collective primitives that are part of collective operation**

- **Progress Engine – Responsible for starting, progressing, and completing the collective primitives**

OAK RIDGE
National Laboratory

# Grouping the Processes Based on the Communication Hierarchy

## Grouping processes into UMA, Socket, and P2P subgroup



**P2P Subgroup**

**UMA Subroup**
🟠 **UMA Group Leader**

**Socket Subroups**
🟢 **Socket Group Leader**

**CPU Socket**
🟪 **Allocated Core**
⬜ **Unallocated Core**
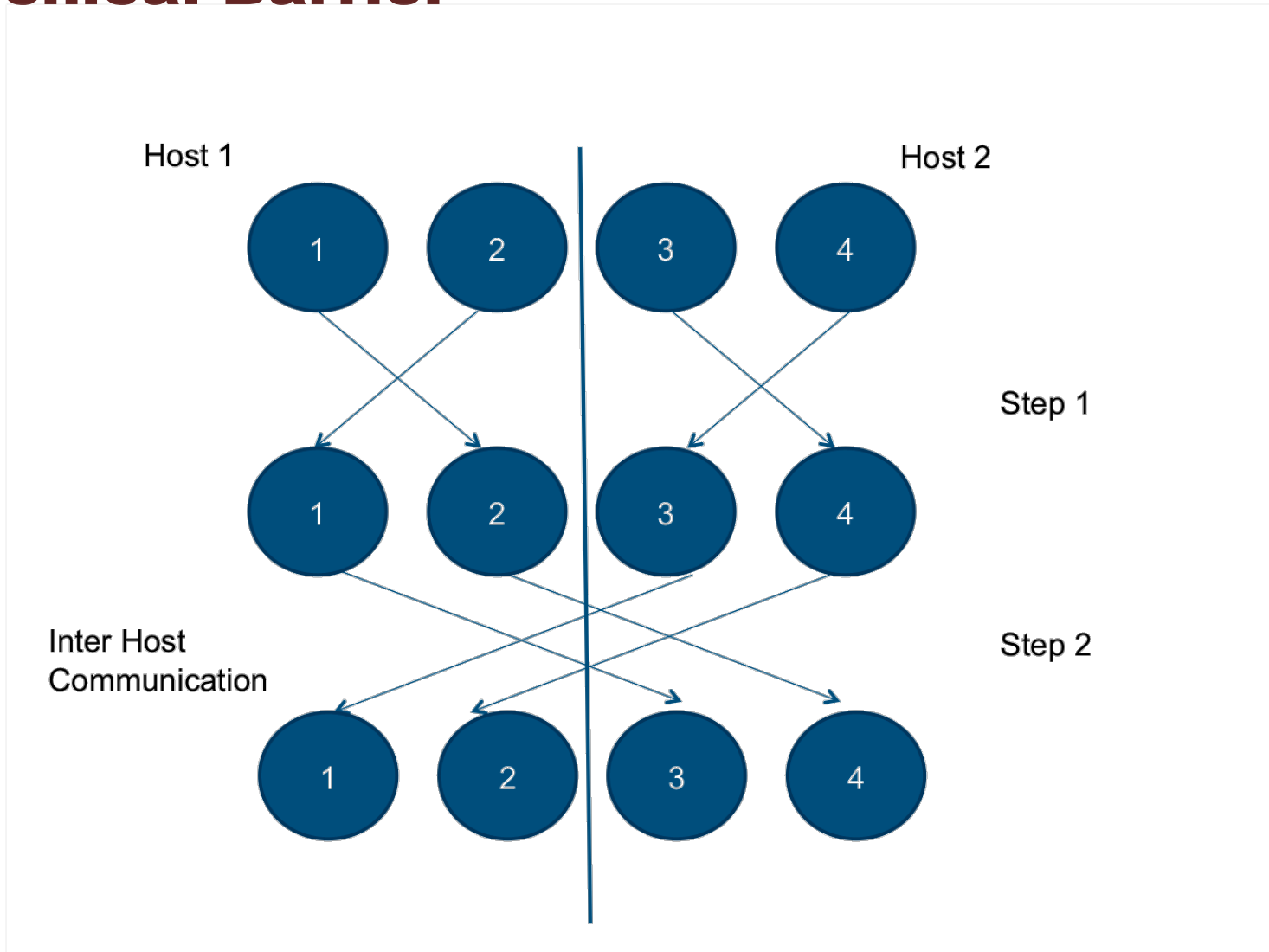
Node 1          Node 2

OAK RIDGE National Laboratory
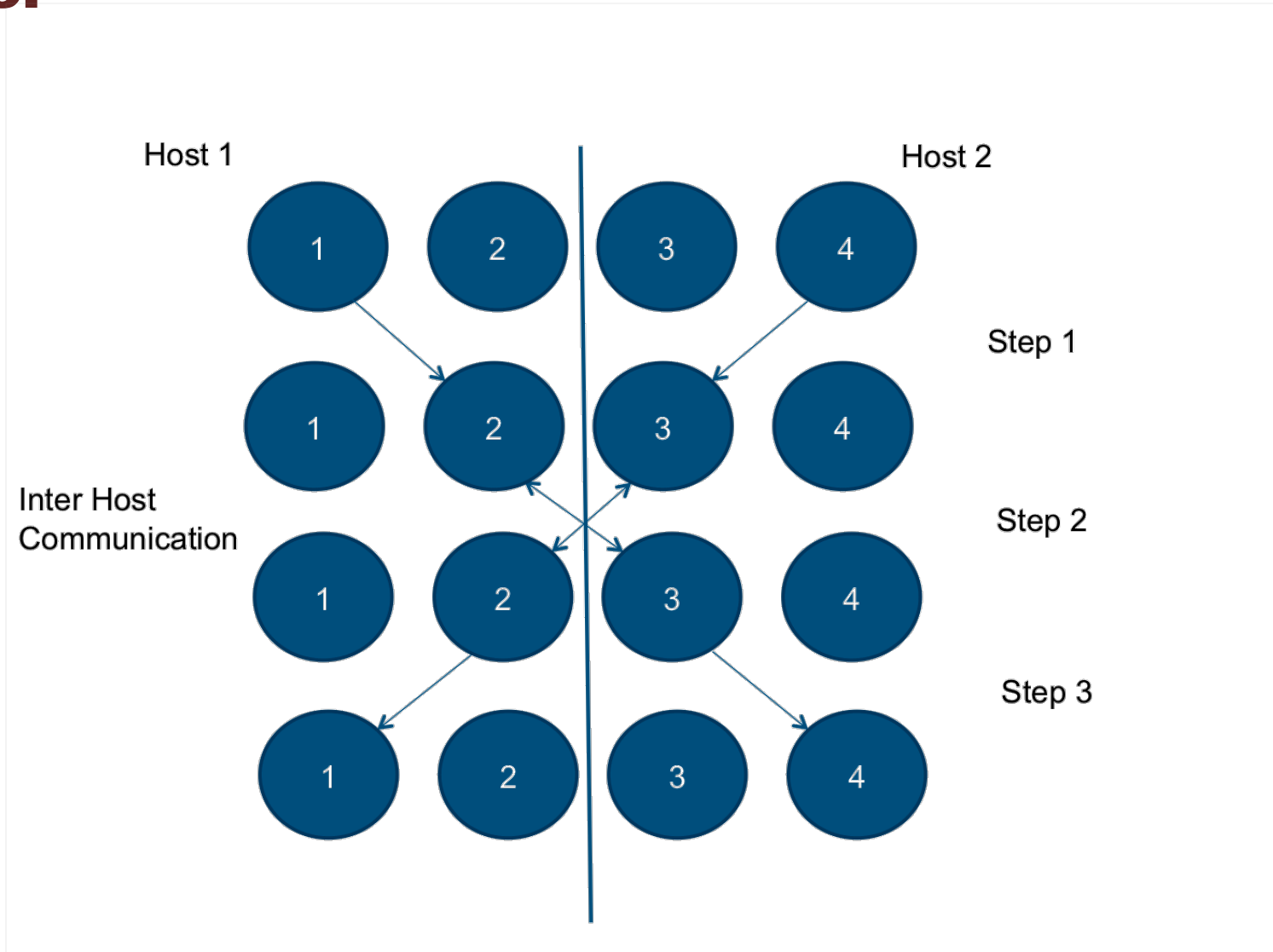
# Hierarchical Collectives : n-level Hierarchical Barrier

- n-level Hierarchical Barrier is a combination of fan-in, fan-out, and recursive k'ing barrier (generalization of recursive doubling) primitives

  – processes participating in the collective operation are grouped into n hierarchies based on communication hierarchy

  – processes in top level hierarchy participate in recursive k`ing Barrier primitive

  – processes in n-1 levels participate in fan-in and fan-out primitives

Cheetah SC11

OAK RIDGE
National Laboratory

# Communication Pattern of a Non-hierarchical Barrier

Cheetah SC11

# Communication Pattern of a Hierarchical Barrier

Cheetah SC11

# Cheetah's Hierarchical Broadcast Algorithms

- **Knownroot Hierarchical Broadcast**
  - the primitives are ordered based on the source of data
  - the primitives are concurrently started after the execution of collective primitive with the source of broadcast
  - uses k-nomial tree for data distribution

- **Unknowroot Hierarchical Broadcast**
  - the primitives are not ordered and started simultaneously
  - the k-nomial tree for data distribution is built dynamically

- **N-ary Hierarchical Broadcast**
  - same as Knownroot algorithm but uses N-ary tree for data distribution

- **Sequential Hierarchical Broadcast**
  - the collective primitives are ordered sequentially
  - there is no concurrent execution

Cheetah SC11

OAK RIDGE
National Laboratory

# Experimental Setup

- **Hardware :**

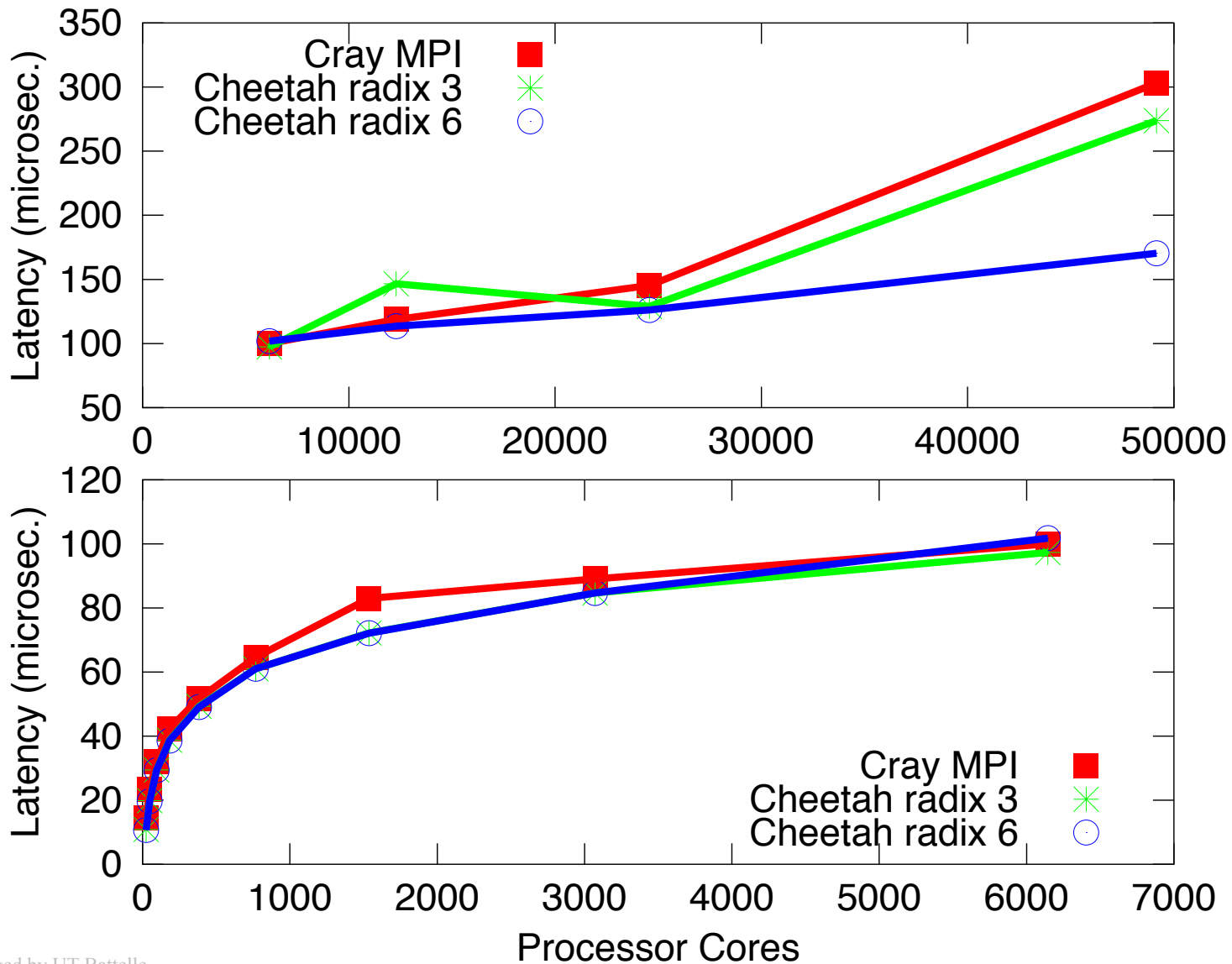  **Jaguar (a Cray XT 5 supercomputer)**
  - 18,688 Compute Nodes
  - 2.6 GHz AMD Opteron (Istanbul)
  - SeaStar 2+ Routers connected in a 3D torus topology

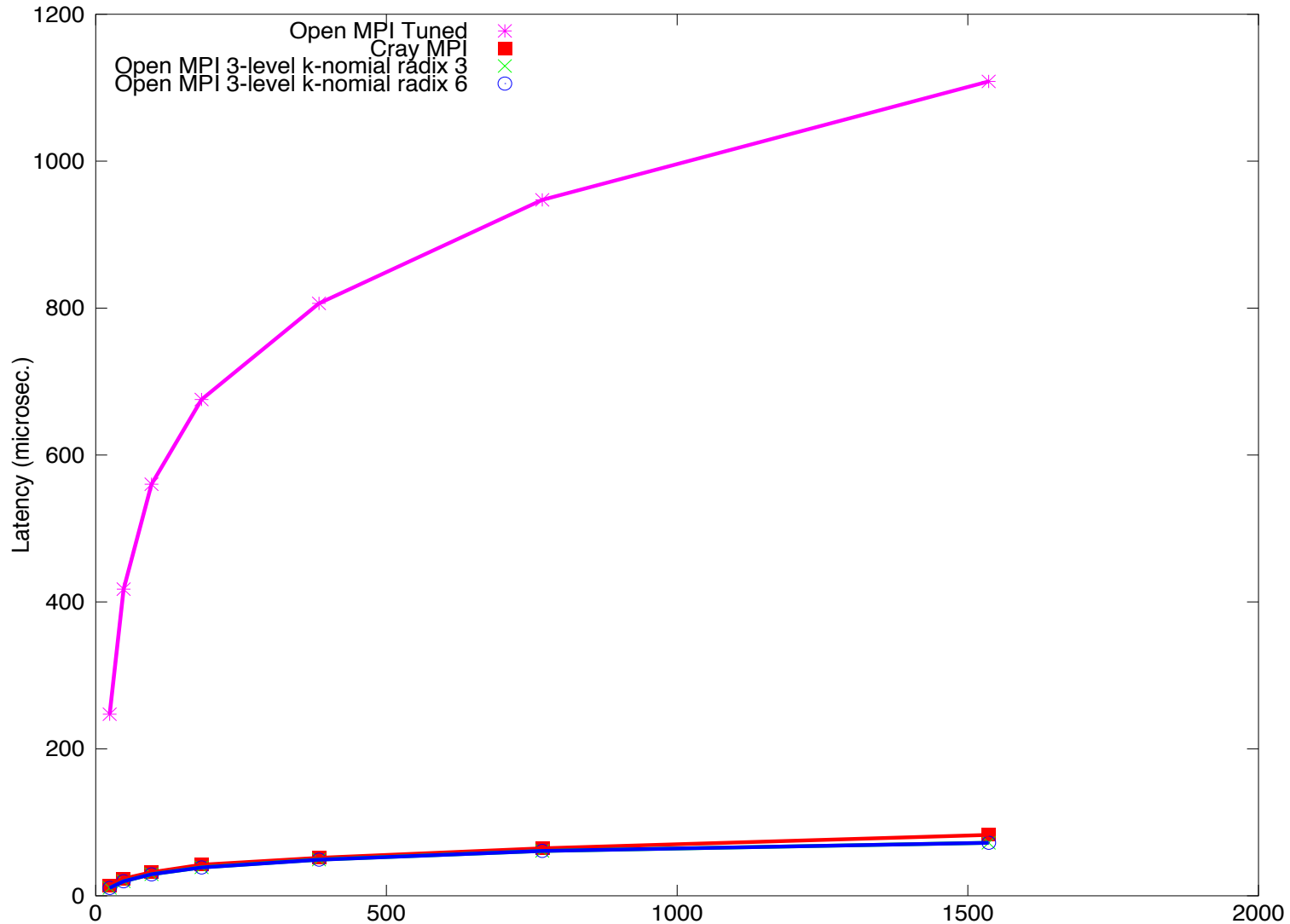- **Benchmarks :**
  - Point-to-Point : OSU Latency and Bandwidth
  - Collectives :
    - Barrier in a tight loop
    - Broadcast in a tight loop

Cheetah SC11

OAK RIDGE
National Laboratory

# Cheetah's Barrier Collective Outperforms the Cray MPI Barrier by 78%



Managed by UT-Battelle
for the U.S. Department of Energy

Cheetah SC11

# Cheetah's Barrier Collective Outperforms the Open MPI default Barrier



Managed by UT-Battelle
for the U.S. Department of Energy

Cheetah SC11

# Cheetah's N-ary Broadcast Outperforms the Cray MPI Broadcast by 52% for 1 byte message



Cray MPI
Cheetah three-level known k-nomial
Cheetah three-level known n-ary
Cheetah three-level unknown k-nomial
Cheetah three-level sequential bcast
Cheetah two-level known k-nomial
Cheetah two-level known n-ary
Cheetah two-level unknown k-nomial
Cheetah two-level sequential bcast

# Cheetah's N-ary Broadcast Outperforms the Cray MPI by 67% for 2 KB message



Legend:
- Cray MPI
- Cheetah three-level known k-nomial
- Cheetah three-level known n-ary
- Cheetah three-level unknown k-nomial
- Cheetah three-level sequential bcast
- Cheetah two-level known k-nomial
- Cheetah two-level known n-ary
- Cheetah two-level unknown k-nomial
- Cheetah two-level sequential bcast

Y-axis: Latency (microsec.)
X-axis: Processor cores
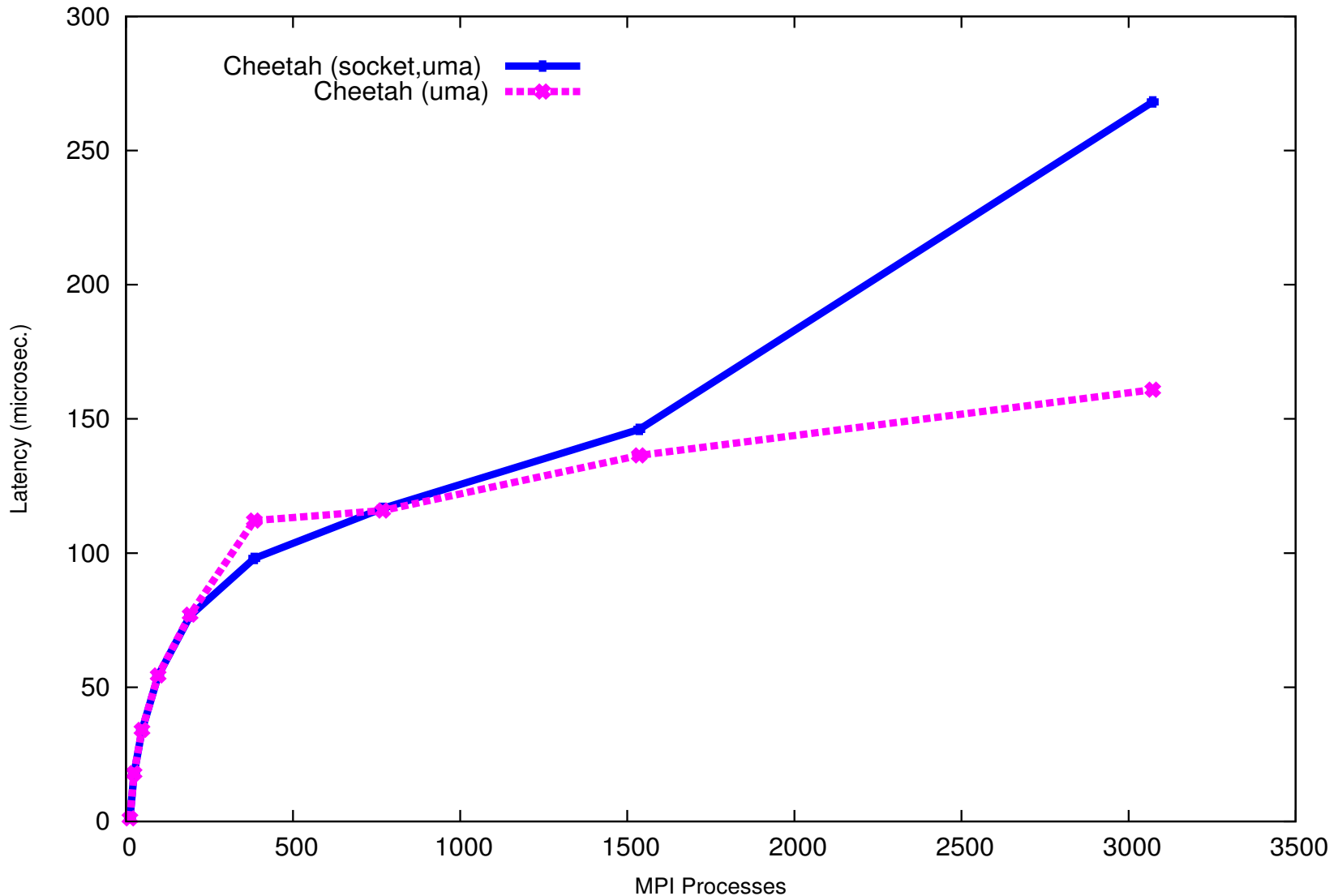
Cheetah SC11

OAK RIDGE National Laboratory

# Cheetah's N-ary Broadcast Outperforms the Cray MPI by 8% for 16 MB message

Cheetah SC11

# Cheetah's Collective Performance with Different Levels of Hierarchy

Managed by UT-Battelle
for the U.S. Department of Energy

Cheetah SC11

OAK
RIDGE
National Laboratory

# Cheetah's two-level hierarchical Barrier outperforms three-level Barrier
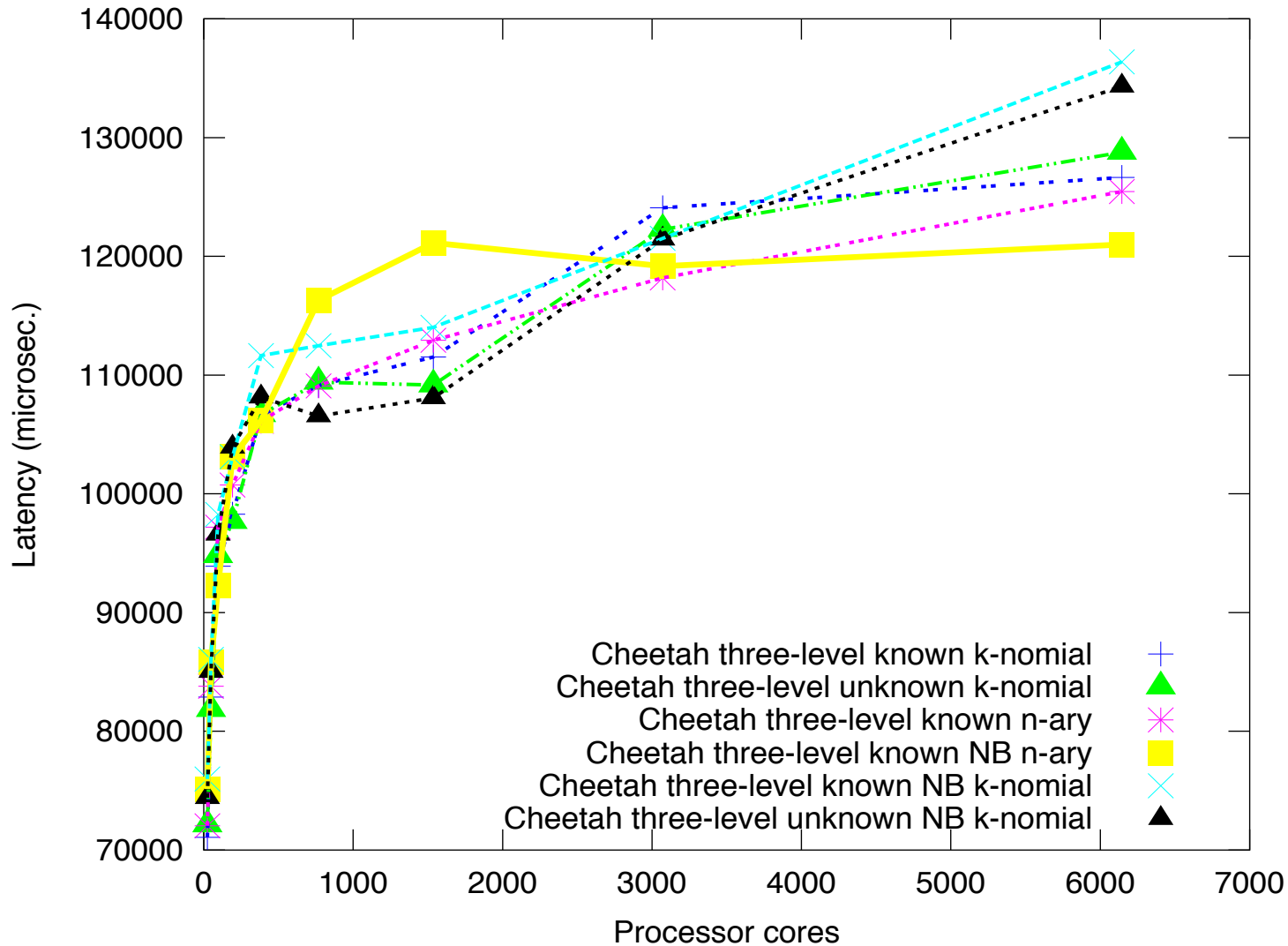
# Cheetah's three-level hierarchical Broadcast outperforms two-level Broadcast

Cheetah SC11

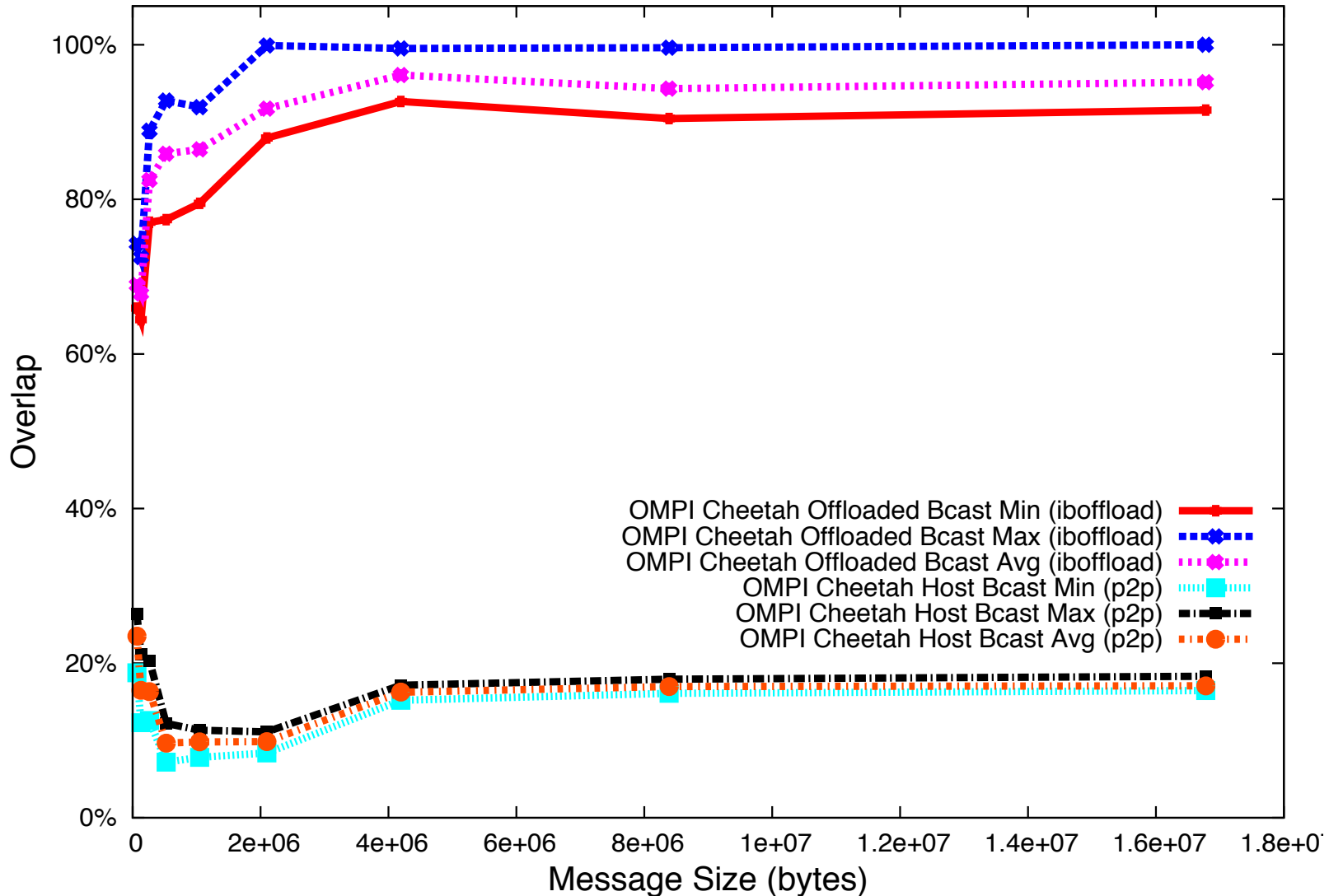# Cheetah's Non-blocking Collective Performance and Overlap

OAK RIDGE
National Laboratory

# Cheetah's Non-blocking Collectives Performance is Comparable to Blocking Collectives



Managed by UT-Battelle
for the U.S. Department of Energy

Cheetah SC11

# Cheetah's offloaded Broadcast Collective provides 96% Overlap



OMPI Cheetah Offloaded Bcast Min (iboffload)
OMPI Cheetah Offloaded Bcast Max (iboffload)
OMPI Cheetah Offloaded Bcast Avg (iboffload)
OMPI Cheetah Host Bcast Min (p2p)
OMPI Cheetah Host Bcast Max (p2p)
OMPI Cheetah Host Bcast Avg (p2p)

# Summary

- **Cheetah's medium message Broadcast outperforms the Cray MPI Broadcast by 67%**

- **Cheetah's Barrier outperforms the Cray MPI's Barrier by 78%**

- **No one-fixed hierarchy configuration suits all collectives**

- **Cheetah's non-blocking collective performance is similar to its blocking collective performance, and offloaded Cheetah collective provides 96% overlap**

- **The key to the performance and scalability of the collective operations**
  - **Collectives customized for communication hierarchies**
  - **Concurrent execution of sub-operations**
  - **Scalable resource usage techniques**
  - **Asynchronous semantics and progress**

OAK RIDGE
National Laboratory

# Contact

## Richard L. Graham

**Application Performance Tools**
**Computer Science and Mathematics Division**
**(865) 356-3469**
**rlgraham@ornl.gov**