# UPC AT SCALE

Yili Zheng

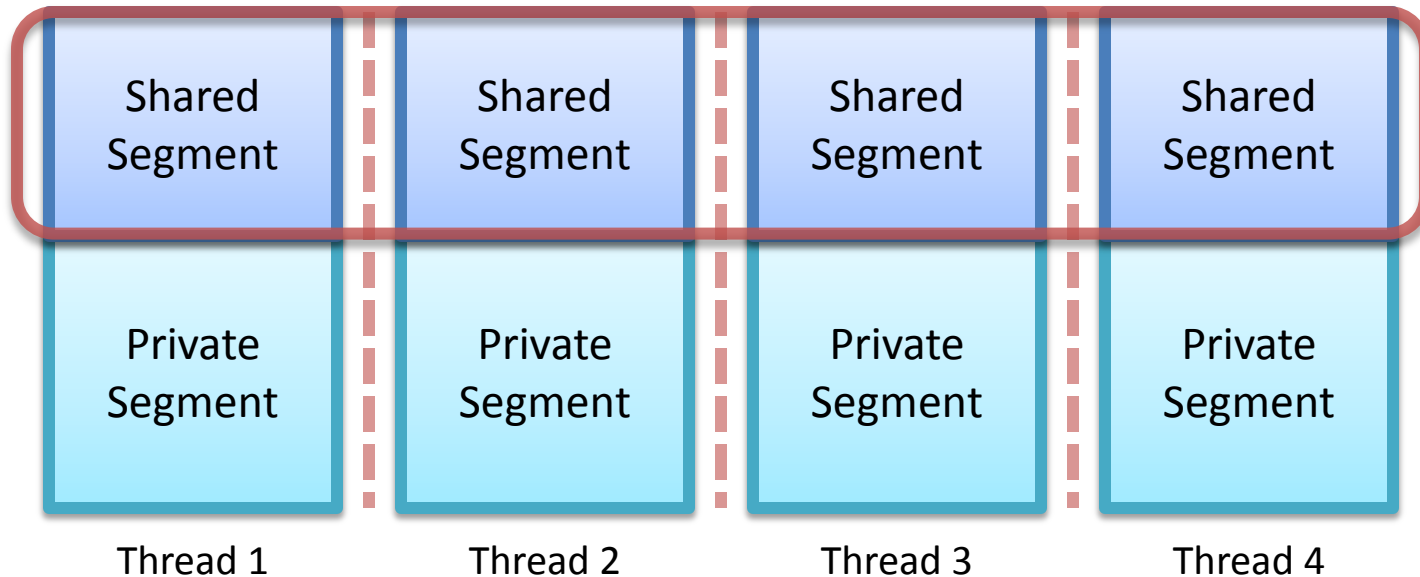Lawrence Berkeley National Laboratory

# Berkeley UPC Team

- Project Lead: Katherine Yelick

- Team members: Filip Blagojevic, Dan Bonachea, Paul Hargrove, Costin Iancu, Seung-Jai Min, Yili Zheng

- Former members: Christian Bell, Wei Chen, Jason Duell, Parry Husbands, Rajesh Nishtala , Mike Welcome

- A joint project of LBNL and UC Berkeley

# Motivation

- Scalable systems have either distributed memory or shared memory without cache coherency
  - Clusters: Ethernet, Infiniband, CRAY XT, IBM BlueGene
  - Hybrid nodes: CPU + GPU or other kinds of accelerators
  - SoC: IBM Cell, Intel Single-chip Cloud Computer (SCC)
- Challenges of Message Passing programming models
  - Difficult data partitioning for irregular applications
  - Memory space starvation due to data replication
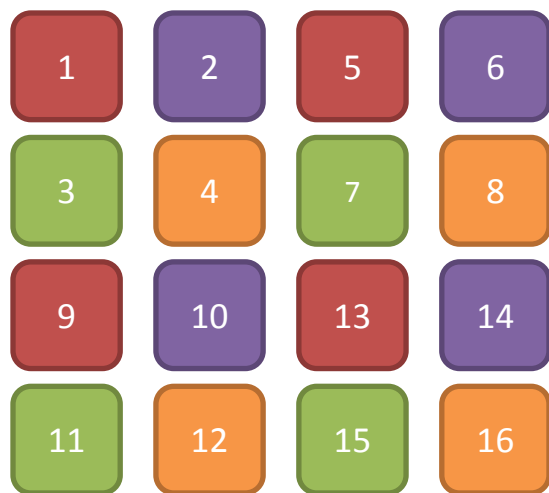  - Performance overheads from two-sided communication semantics

# Partitioned Global Address Space

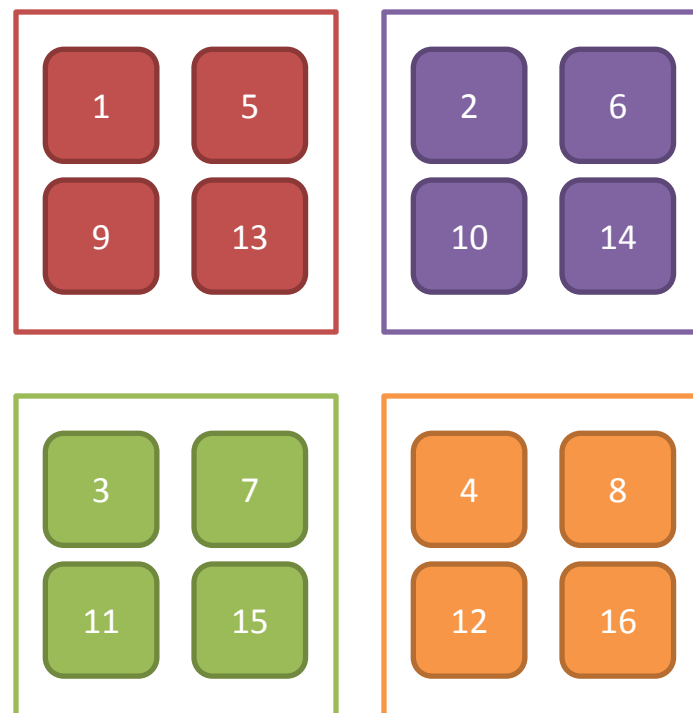| Shared Segment | Shared Segment | Shared Segment | Shared Segment |
|---|---|---|---|
| Private Segment | Private Segment | Private Segment | Private Segment |
| Thread 1 | Thread 2 | Thread 3 | Thread 4 |

- Global data view abstraction for productivity
- Vertical partitions among threads for locality control
- Horizontal partitions between shared and private segments for data placement optimizations
- Friendly to non-cache-coherent architectures

# PGAS Example: Global Matrix Distribution

## Global Matrix View



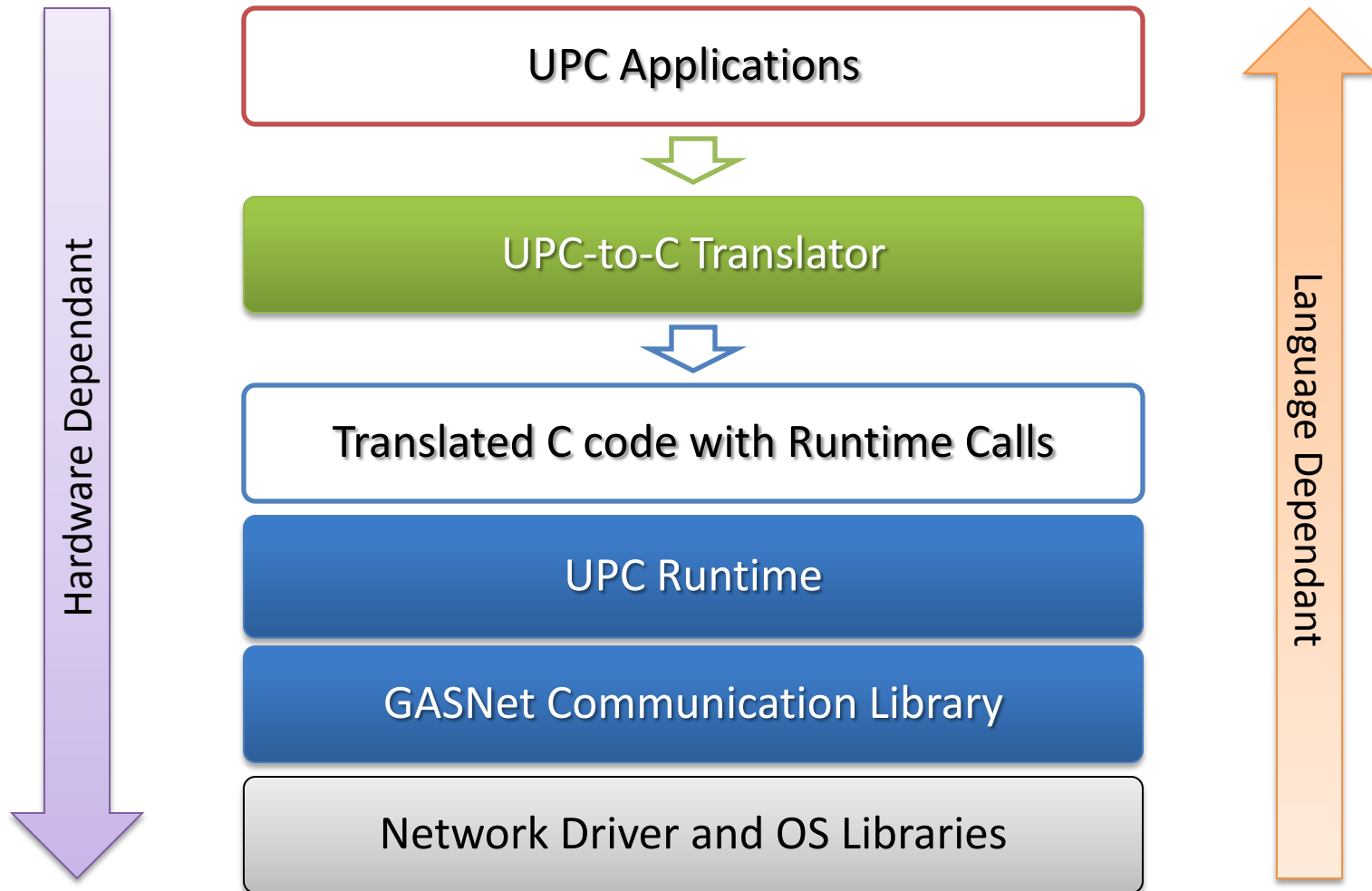## Distributed Matrix Storage

# UPC Overview

- PGAS dialect of ISO C99

- Distributed shared arrays

- Dynamic shared-memory allocation

- One-sided shared-memory communication

- Synchronization: barriers, locks, memory fences

- Collective communication library

- Parallel I/O library

# Key Components for Scalability

- One-sided communication and active messages

- Efficient resource sharing for multi-core systems
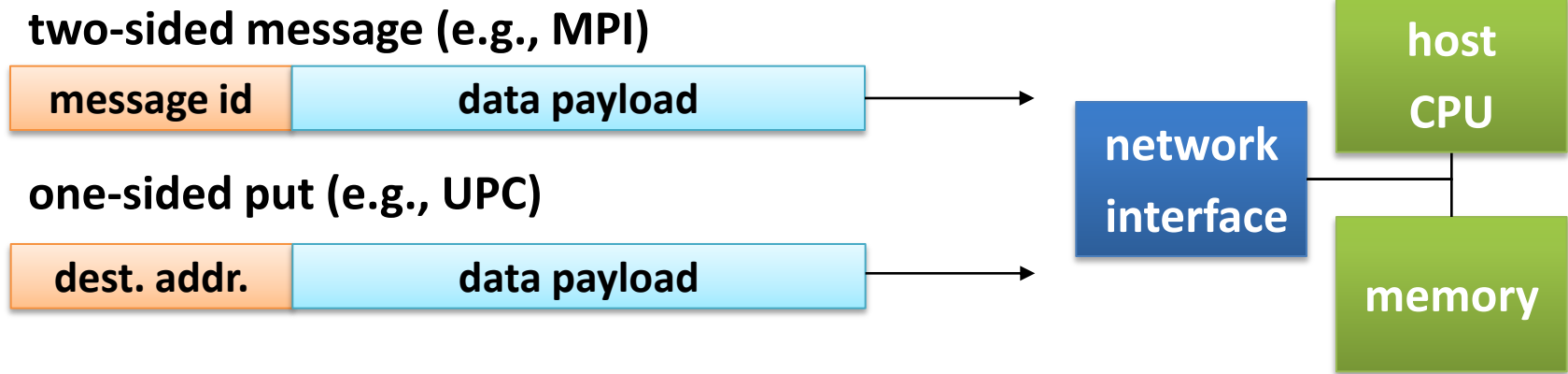
- Non-blocking collective communication

# Berkeley UPC Software Stack

Hardware Dependant

Language Dependant

UPC Applications

UPC-to-C Translator

Translated C code with Runtime Calls

UPC Runtime

GASNet Communication Library

Network Driver and OS Libraries

# Berkeley UPC Features

- Data transfer for complex data types (vector, indexed, stride)
- Non-blocking memory copy
- Point-to-point synchronization
- Remote atomic operations
- Active Messages
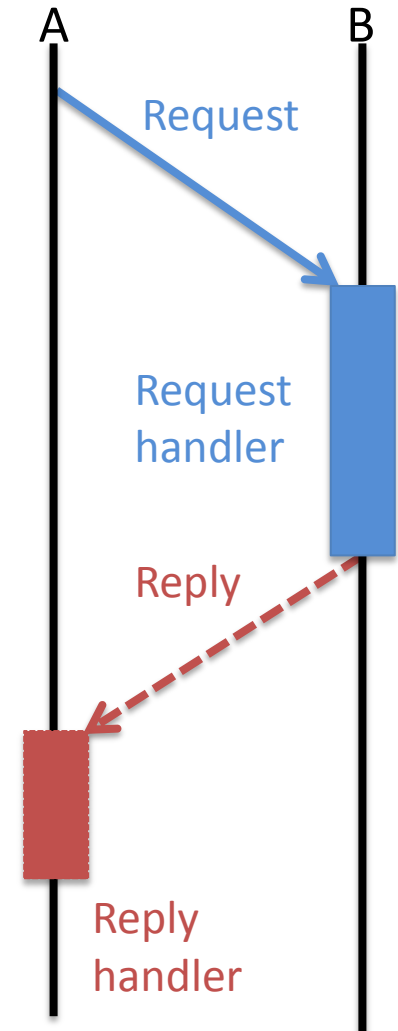- Extension to UPC collectives
- Portable timers

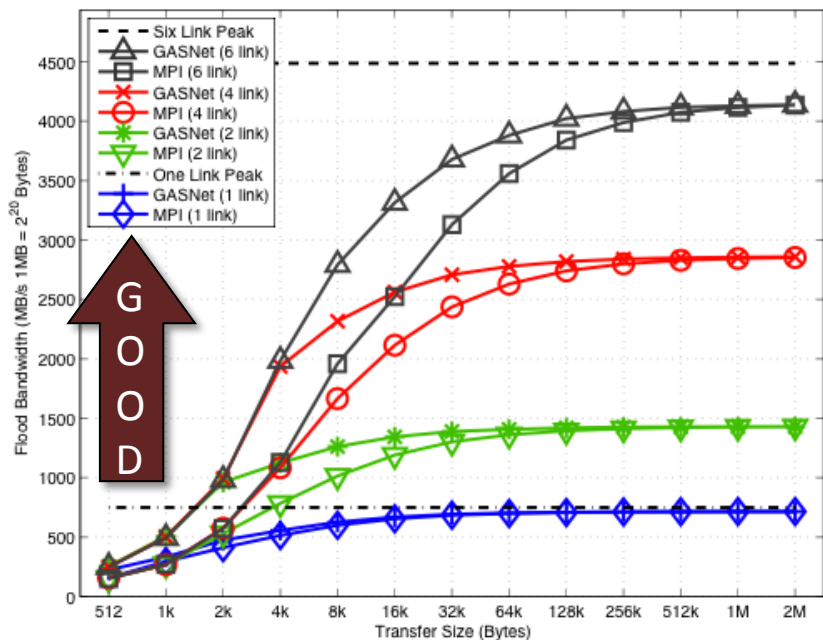# One-Sided vs. Two-Sided Messaging

**two-sided message (e.g., MPI)**

| message id | data payload |
|---|---|

**one-sided put (e.g., UPC)**

| dest. addr. | data payload |
|---|---|

**network interface**

**host CPU**

**memory**

- Two-sided messaging
  - Message does not contain information about the final destination; need to look it up on the target node
  - Point-to-point synchronization implied with all transfers
- One-sided messaging
  - Message contains information about the final destination
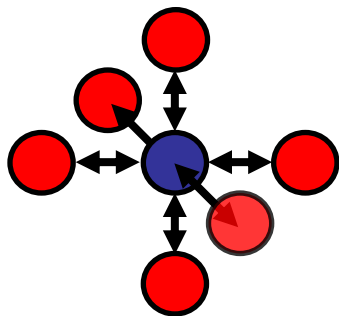  - Decouple synchronization from data movement

# Active Messages

- Active messages = Data + Action
- Key enabling technology for both one-sided and two-sided communications
  - Software implementation of Put/Get
  - Eager and Rendezvous protocols
- Remote Procedural Calls
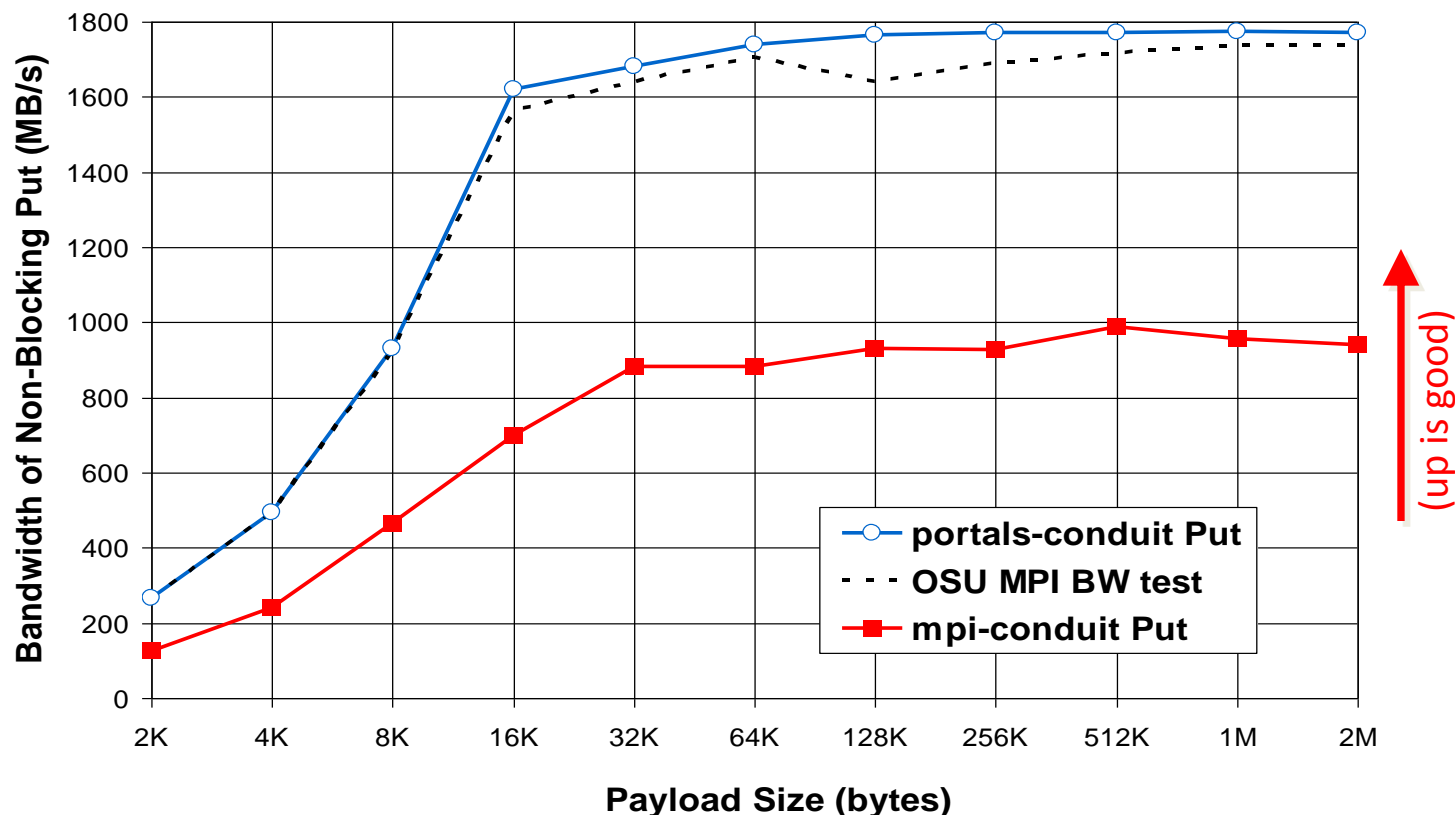  - Facilitate "owner-computes"
  - Spawn asynchronous tasks

A      B

Request

Request handler

Reply

Reply handler

# GASNet Bandwidth on BlueGene/P



* Kumar et. al showed the maximum achievable bandwidth for DCMF transfers is 748 MB/s per link so we use this as our peak bandwidth
See "The deep computing messaging framework: generalized scalable message passing on the blue gene/P supercomputer", Kumar et al. ICS08
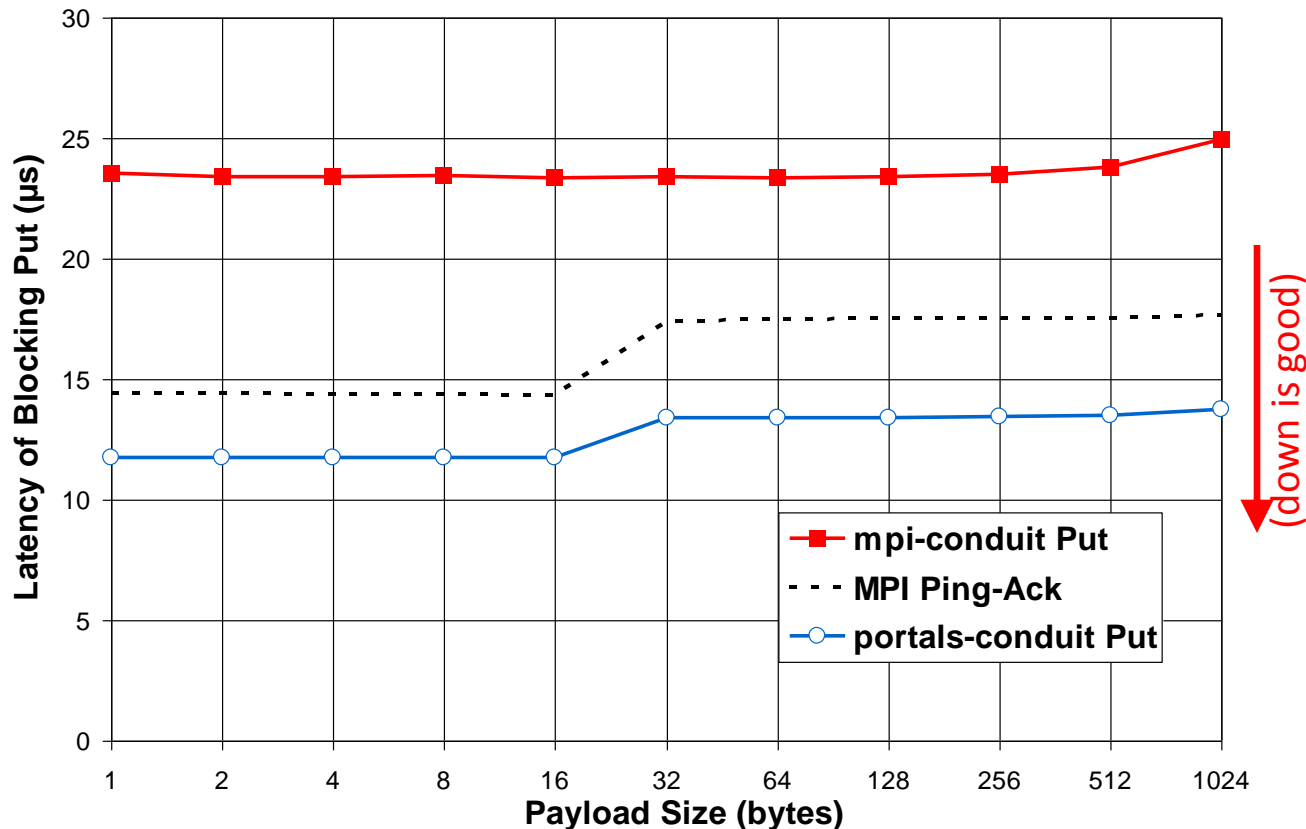
- Torus network
  - Each node has six 850MB/s* bidirectional links
  - Vary number of links from 1 to 6
- Consecutive non-blocking puts on the links (round-robin)
- Similar bandwidth for large-size messages
- GASNet outperforms MPI for mid-size messages
  - Lower software overhead
  - More overlapping

See "Scaling Communication Intensive Applications on BlueGene/P Using One-Sided Communication and Overlap",  Rajesh Nishtala, Paul Hargrove, Dan Bonachea, and Katherine Yelick, *IPDPS 2009*

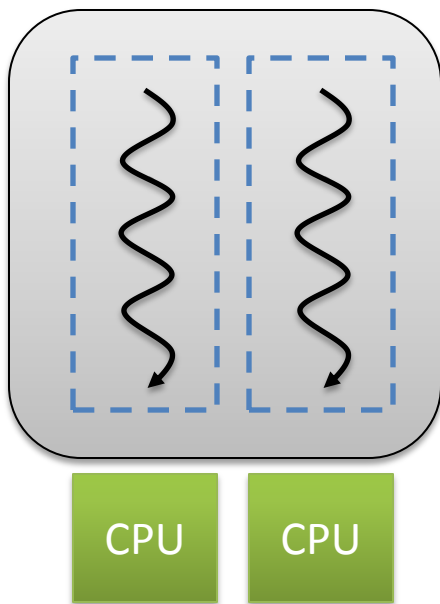# GASNet Bandwidth on Cray XT4



Slide source: Porting GASNet to Portals: Partitioned Global Address Space (PGAS) Language Support for the Cray XT, Dan Bonachea, Paul Hargrove, Michael Welcome, Katherine Yelick, CUG 2009
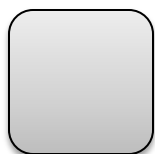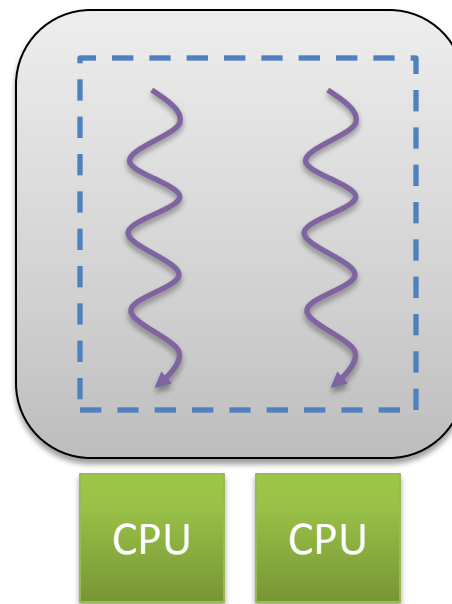
# GASNet Latency on Cray XT4



Slide source: Porting GASNet to Portals: Partitioned Global Address Space (PGAS) Language Support for the Cray XT, Dan Bonachea, Paul Hargrove, Michael Welcome, Katherine Yelick, CUG 2009

# Execution Models on Multi-core – Process vs. Thread

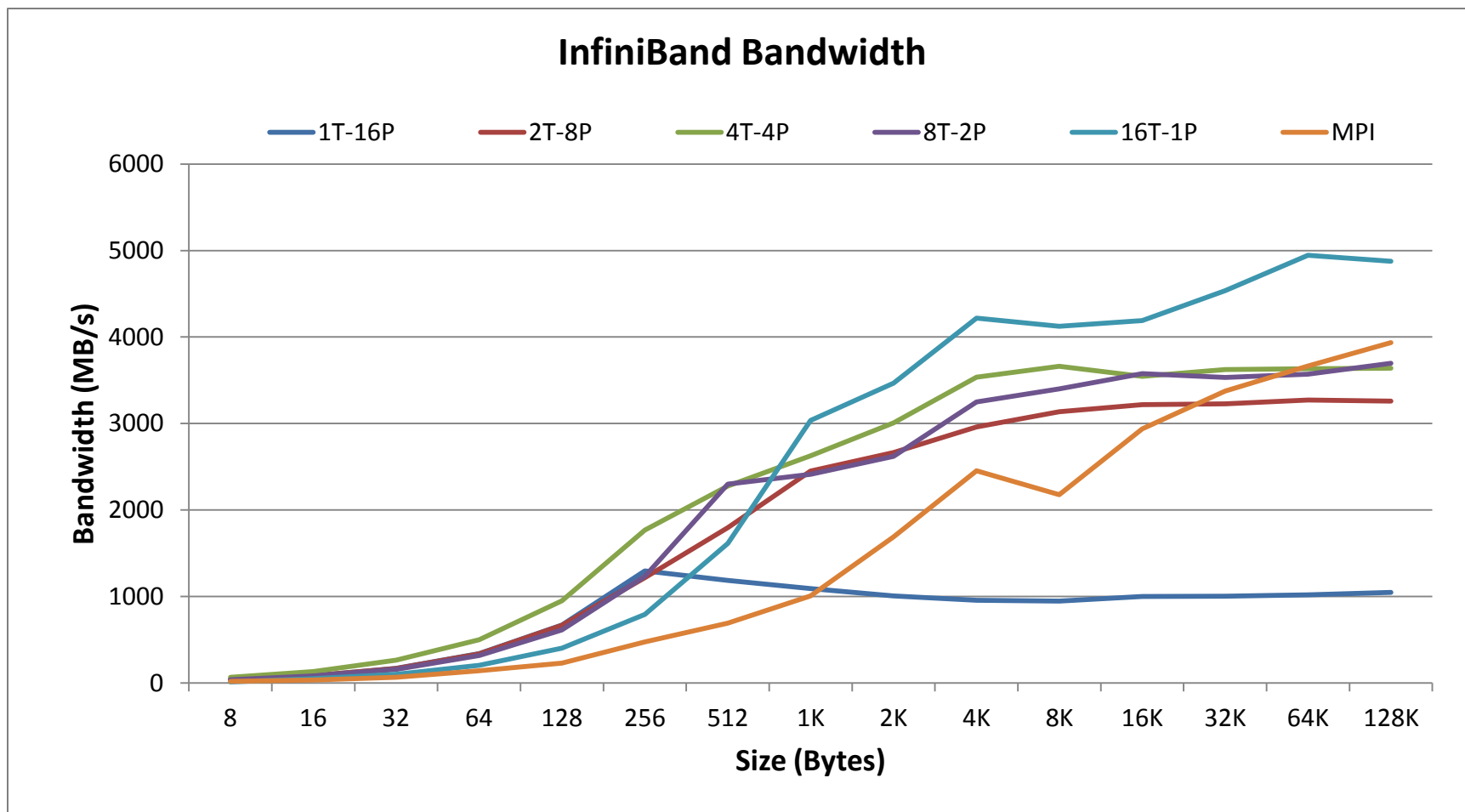Map UPC threads to Processes

Map UPC threads to Pthreads

CPU    CPU

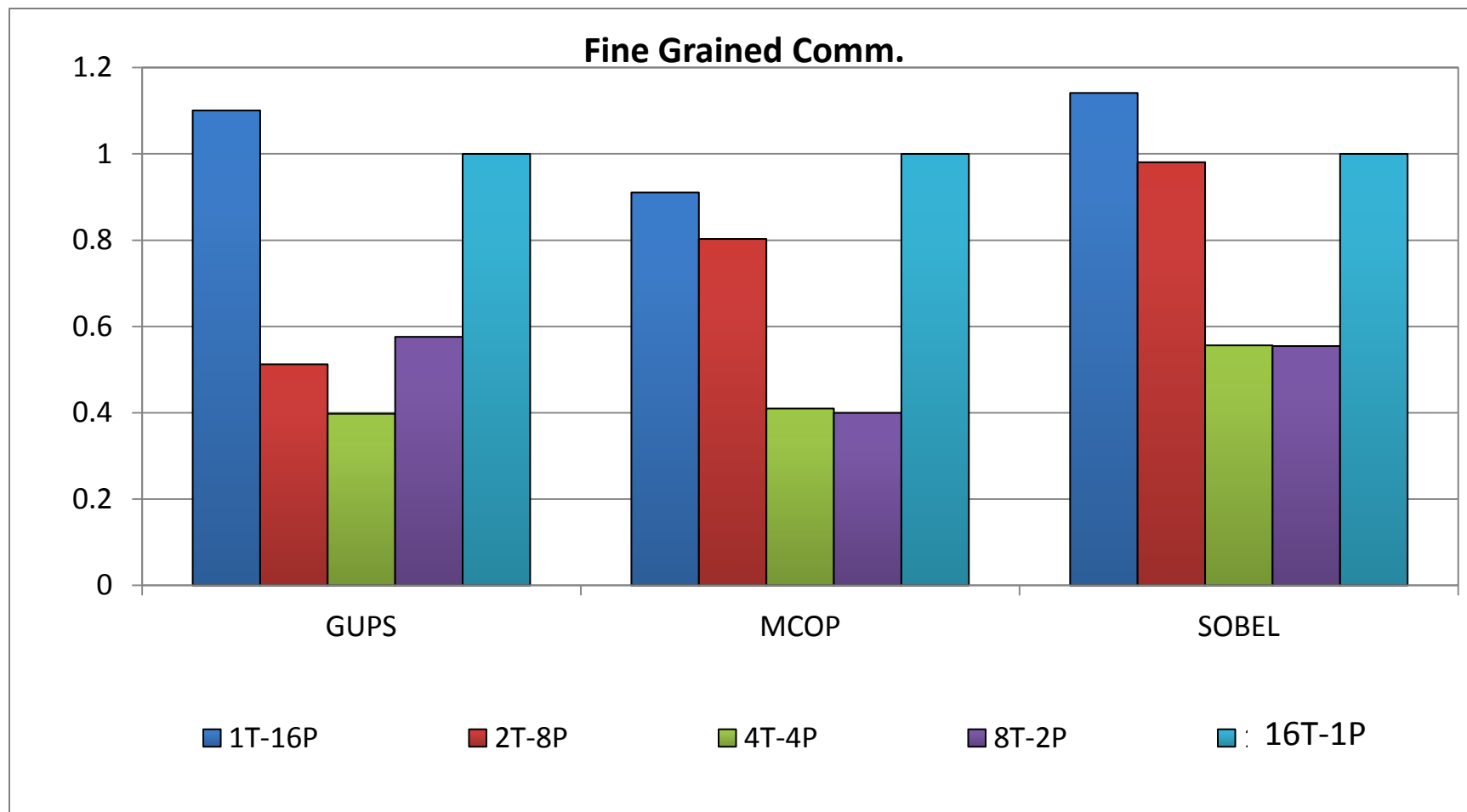CPU    CPU

Physical Shared-memory

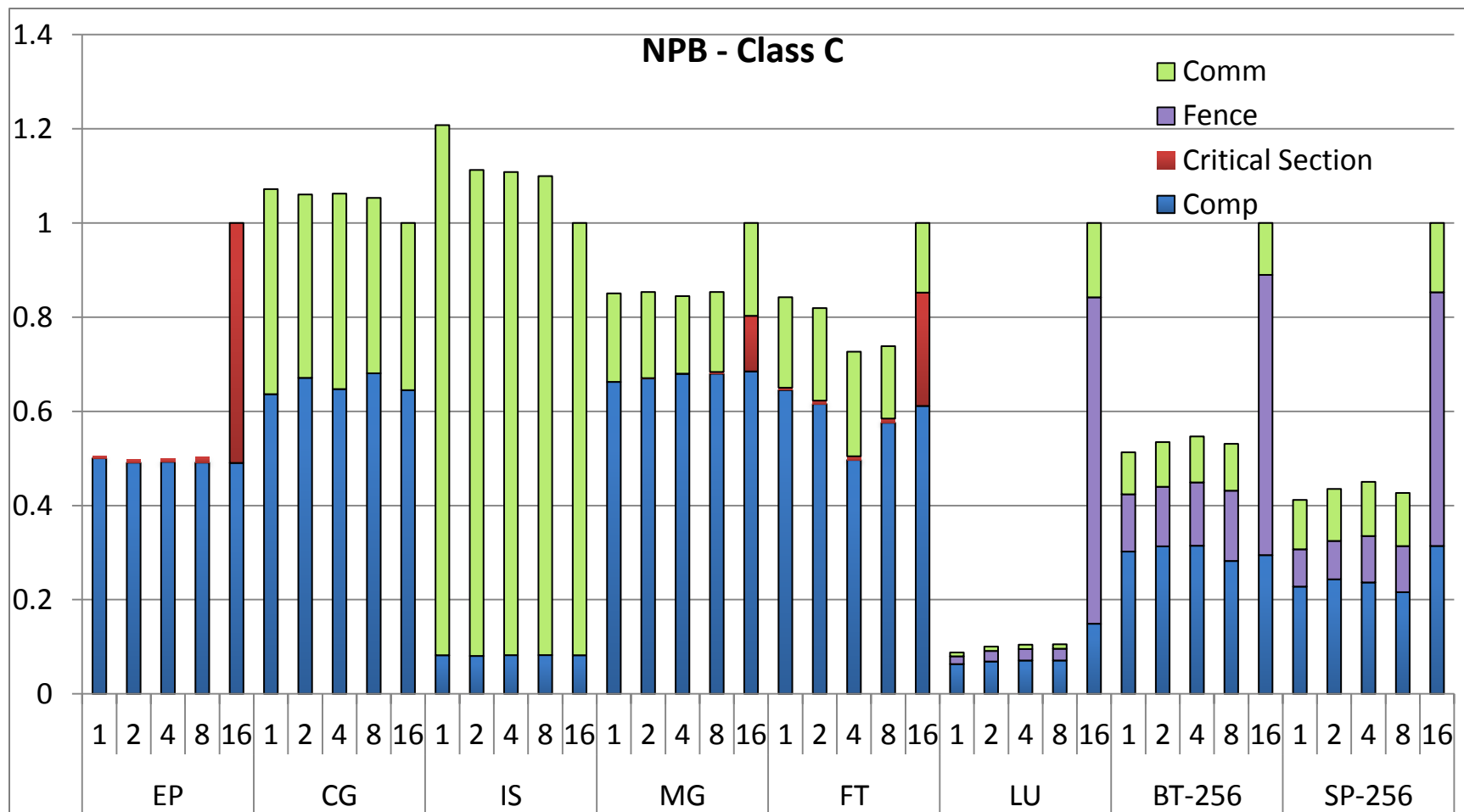Virtual Address Space

# Point-to-Point Performance – Process vs. Thread



InfiniBand Bandwidth

# Application Performance – Process vs. Thread



Fine Grained Comm.
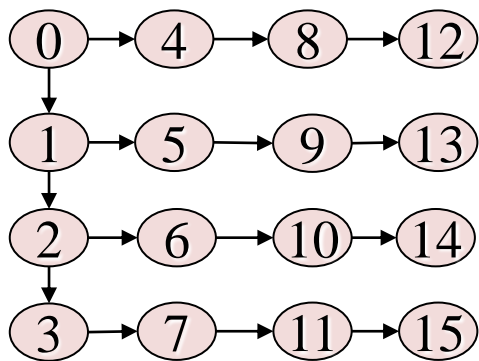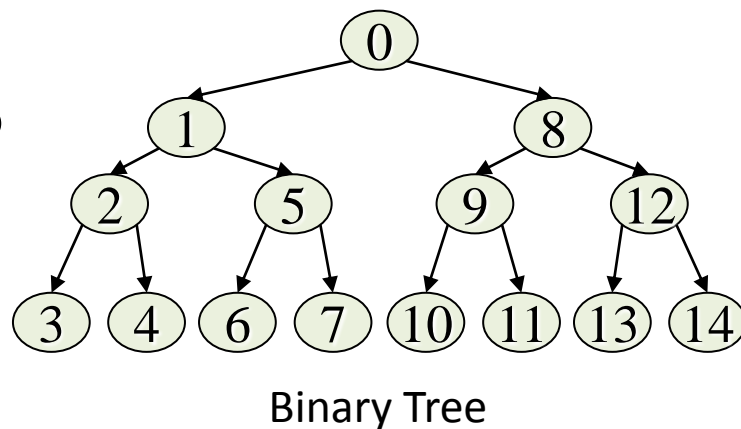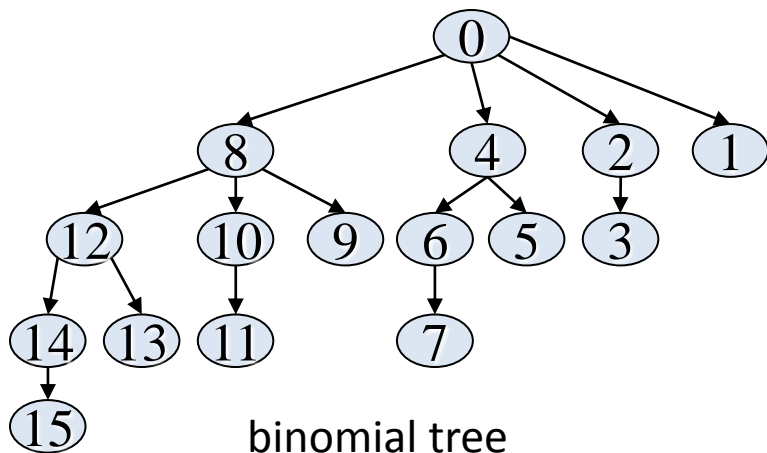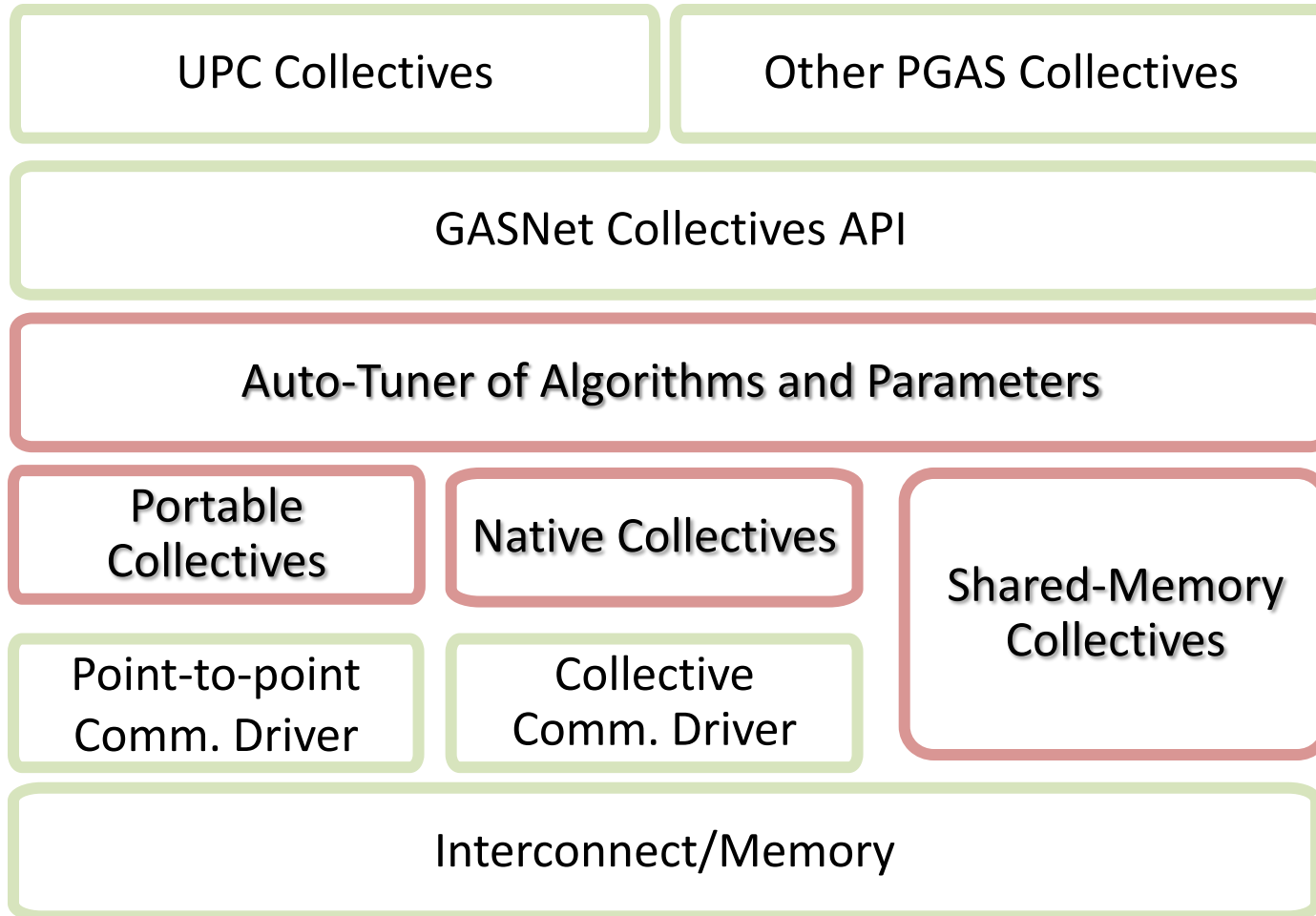
# NAS Parallel Benchmarks – Process vs. Thread

# Collective Communication for PGAS

- Communication patterns similar to MPI: broadcast, reduce, gather, scatter and alltoall

- Global address space enables one-sided collectives

- Flexible synchronization modes provide more communication and computation overlapping opportunities

# Collective Communication Topologies



binomial tree

Binary Tree

Fork Tree

Radix 2 Dissemination

# GASNet Module Organization

UPC Collectives

Other PGAS Collectives

GASNet Collectives API

Auto-Tuner of Algorithms and Parameters

Portable Collectives

Native Collectives

Shared-Memory Collectives

Point-to-point Comm. Driver

Collective Comm. Driver

Interconnect/Memory

# Auto-tuning Collective Communication
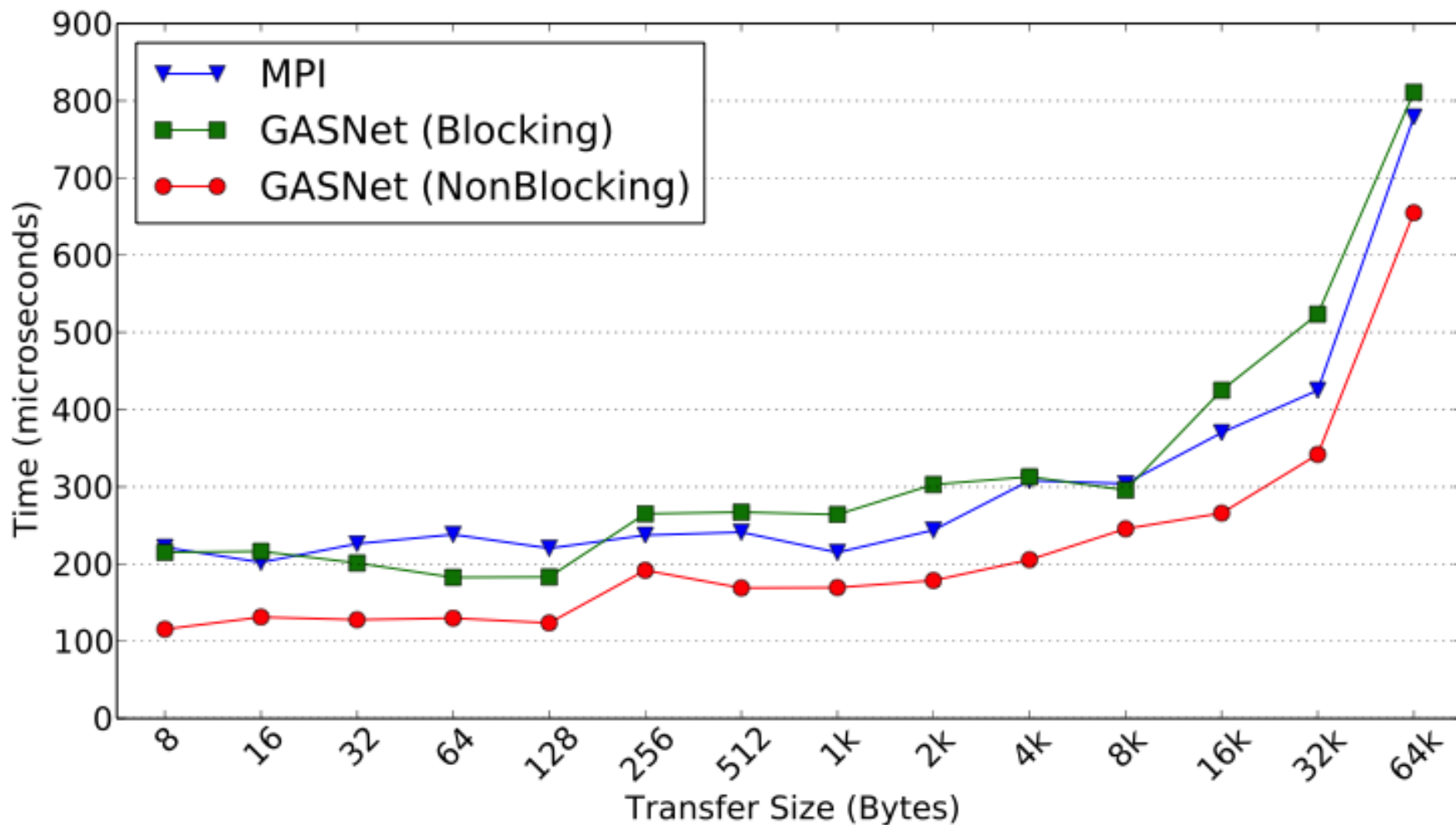
## Offline tuning

- Optimize for platform common characteristics

- Minimize runtime tuning overhead

## Online tuning

- Optimize for application runtime characteristics
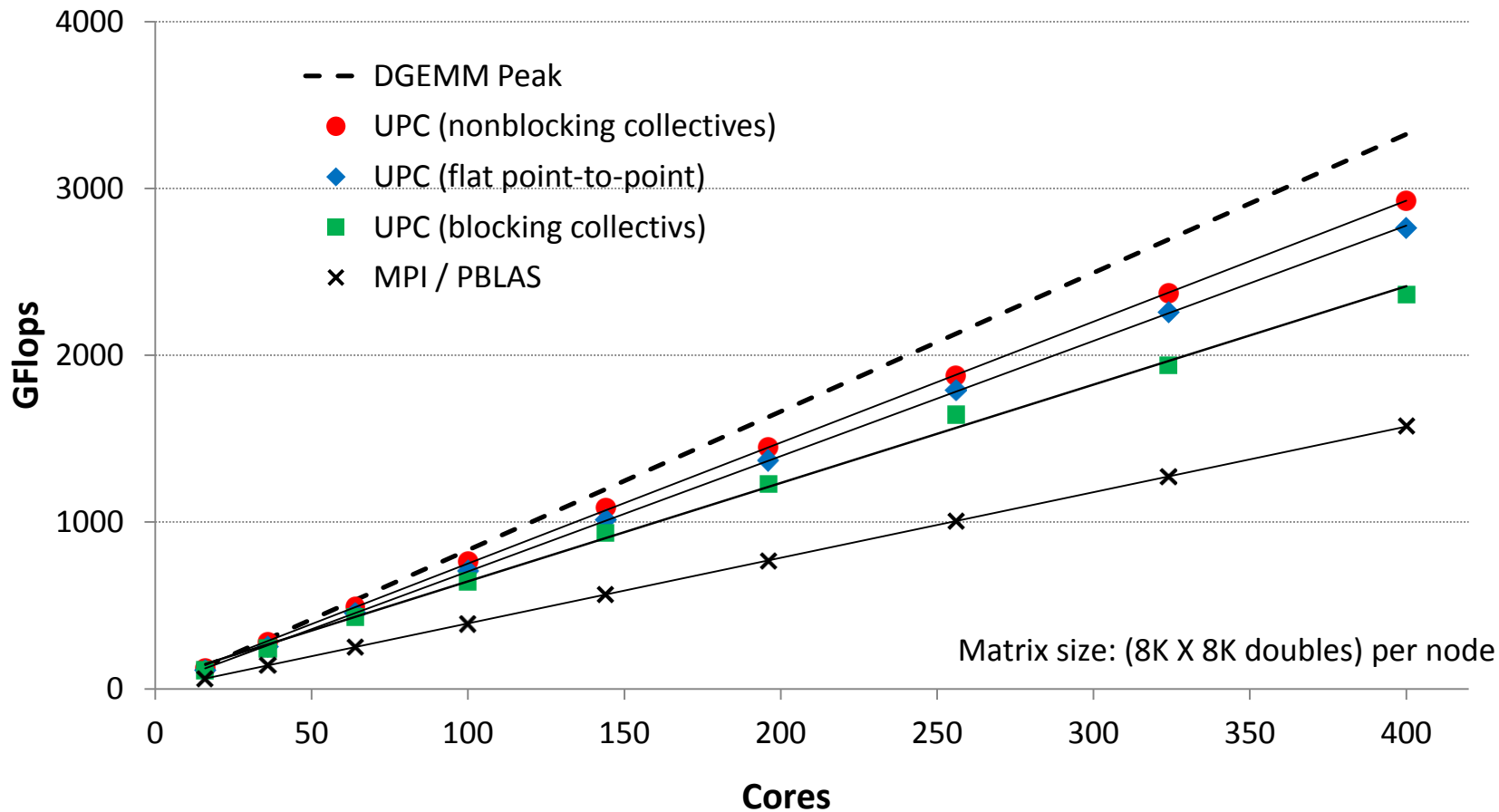
- Refine offline tuning results

| Performance Influencing Factors | Performance Tuning Space |
|---|---|
| Hardware<br>▪ CPU<br>▪ Memory system<br>▪ Interconnect<br>Software<br>▪ Application<br>▪ System software<br>Execution<br>▪ Process/thread layout<br>▪ Input data set<br>▪ System workload | Algorithm selection<br>▪ Eager vs. rendezvous<br>▪ Put vs. get<br>▪ Collection of well-known algorithms<br>Communication topology<br>▪ Tree type<br>▪ Tree fan-out<br>Implementation-specific parameters<br>▪ Pipelining depth<br>▪ Dissemination radix |

# Broadcast Performance



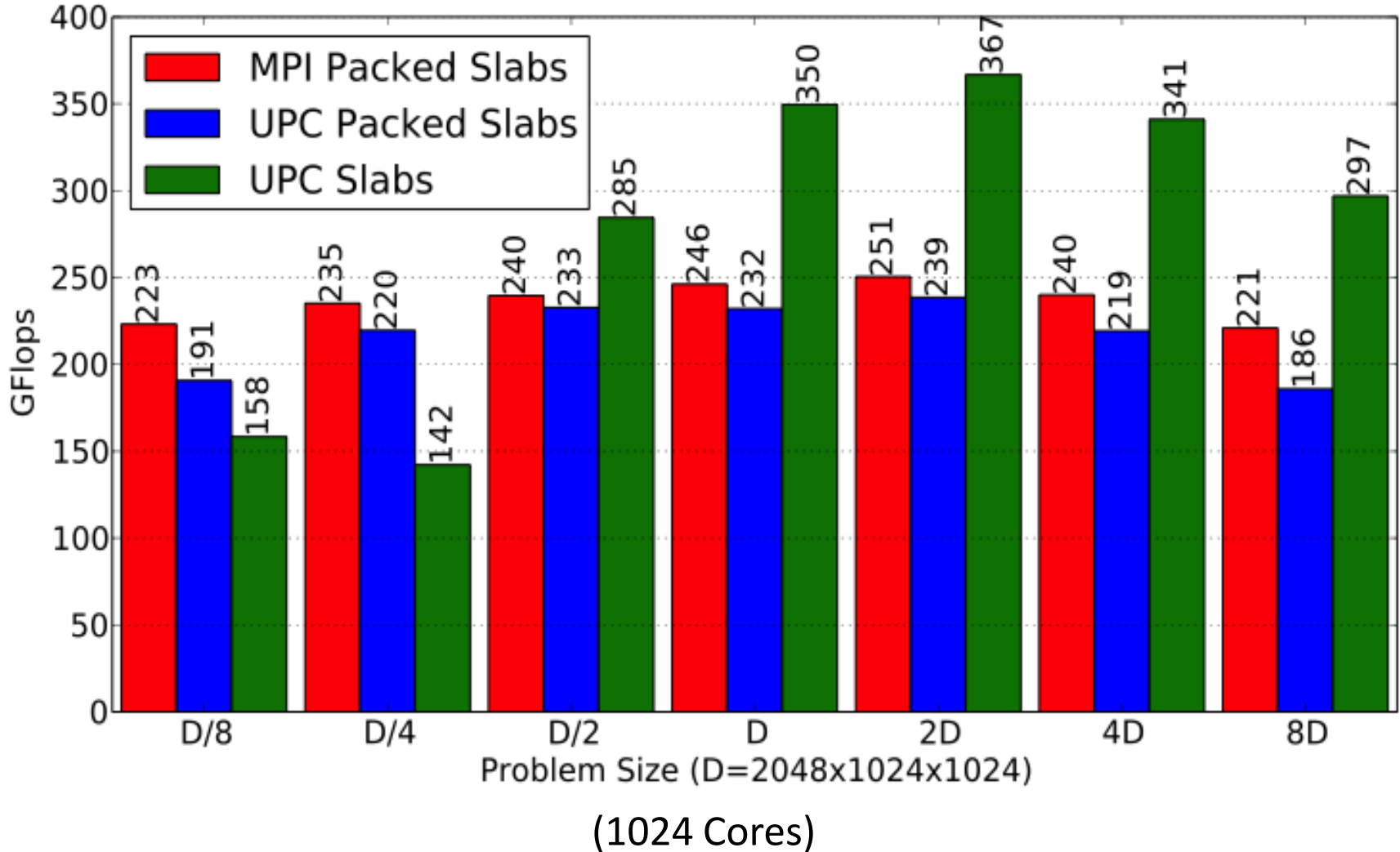Cray XT4 Nonblocking Broadcast (1024 Cores)

# Matrix-Multiplication on Cray XT4



Matrix size: (8K X 8K doubles) per node

# Choleskey Factorization on Sun Constellation (Infiniband)



2048 cores on Ranger
Matrix size: 240K

GFlops chart:
- UPC team collectives: 4097
- Hand-coded UPC: 3757
- Naïve UPC (get-based): 3118

# FFT Performance on Cray XT4



(1024 Cores)

# FFT Performance on BlueGene/P



MPI FFT of HPC Challenge as of July 09 is ~4.5 Tflops on 128k Cores.

# Summary

- PGAS provides programming convenience similar to shared-memory models

- UPC has demonstrated good performance comparable to MPI at large scale.

- Interoperable with other programming models and languages including MPI, FORTRAN and C++

- Growing UPC community with actively developed and maintained software implementations
  - Berkeley UPC and GASNet: http://upc.lbl.gov
  - Other UPC compilers: Cray UPC, GNU UPC, HP UPC and IBM UPC
  - Tools: TotalView and Parallel Performance Wizard (PPW)