

# UCX: An Open Source Framework for HPC Network APIs and Beyond



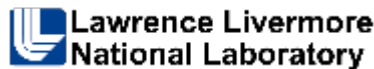
Pavel Shamis (Pasha)  
Principal Research Engineer

# Co-Design Collaboration

## The Next Generation HPC Communication Framework

Collaborative Effort

Industry, National Laboratories and Academia



# Challenges

- Performance Portability (across various interconnects)
  - Collaboration between industry and research institutions
    - ...but mostly industry (because they built the hardware)
- Maintenance
  - Maintaining a network stack is time consuming and expensive
  - Industry have resources and strategic interest for this
- Extendibility
  - MPI+X+Y ?
  - Exascale programming environment is an ongoing debate

# UCX – Unified Communication X Framework

- Unified
  - Network API for multiple network architectures that target HPC programming models and libraries
- Communication
  - How to move data from location in memory A to location in memory B considering multiple types of memories
- Framework
  - A collection of libraries and utilities for HPC network programmers

# History

## MXM

- Developed by Mellanox Technologies
- HPC communication library for InfiniBand devices and shared memory
- Primary focus: MPI, PGAS

## UCCS

- Developed by ORNL, UH, UTK
- Originally based on Open MPI BTL and OPAL layers
- HPC communication library for InfiniBand, Cray Gemini/Aries, and shared memory
- Primary focus: OpenSHMEM, PGAS
- Also supports: MPI

## PAMI

- Developed by IBM on BG/Q, PERCS, IB VERBS
- Network devices and shared memory
- MPI, OpenSHMEM, PGAS, CHARM++, X10
- C++ components
- Aggressive multi-threading with contexts
- Active Messages
- Non-blocking collectives with hw acceleration support

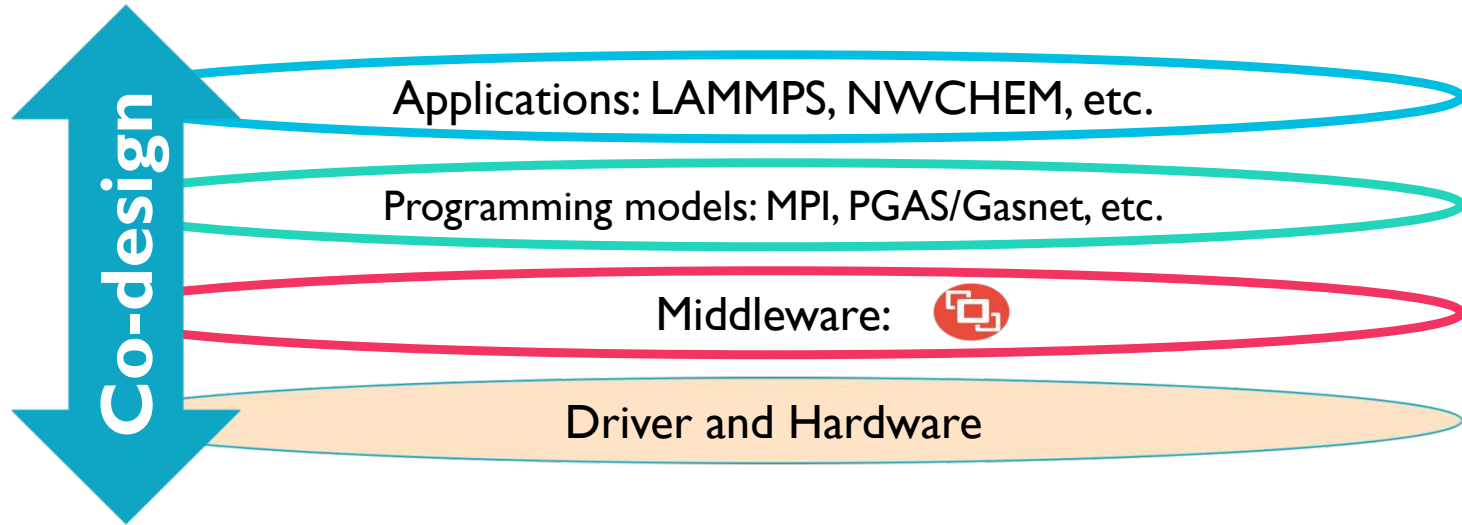
**Decades of community and industry experience in development of HPC software**

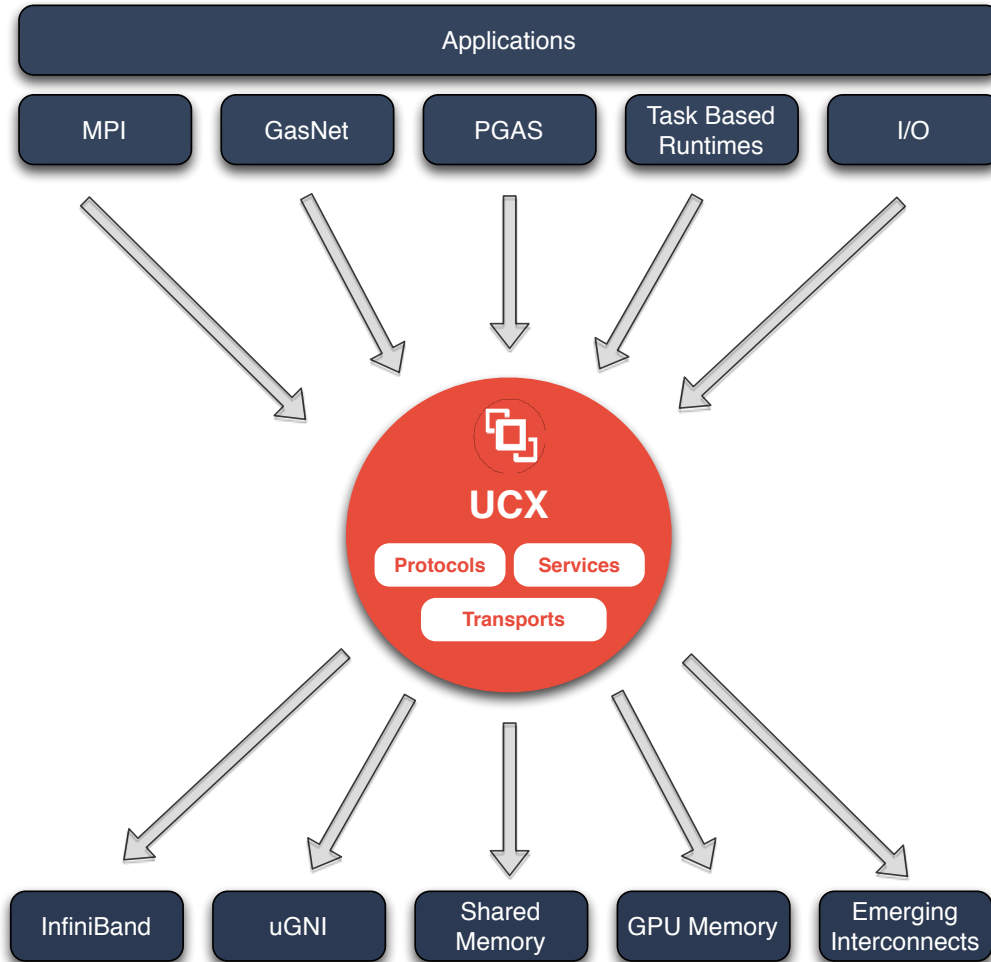
# What we are doing differently...

- UCX consolidates multiple industry and academic efforts
  - Mellanox MXM, IBM PAMI, ORNL/UTK/UH UCCS, etc.
- Supported and maintained by industry
  - IBM, Mellanox, NVIDIA, Pathscale, ARM

# What we are doing differently...

- Co-design effort between national laboratories, academia, and industry







# A Collaboration Efforts



- Mellanox co-designs network API and contributes MXM technology
  - Infrastructure, transport, shared memory, protocols, integration with OpenMPI/SHMEM, MPICH
- ORNL & LANL co-designs network API and contributes UCCS project
  - InfiniBand optimizations, Cray devices, shared memory
- ARM co-designs the network API and contributes optimizations for ARM eco-system
- NVIDIA co-designs high-quality support for GPU devices
  - GPUDirect, GDR copy, etc.
- IBM co-designs network API and contributes ideas and concepts from PAMI
- UH/UTK focus on integration with their research platforms

# Licensing

- Open Source
  - BSD 3 Clause license
  - Contributor License Agreement – BSD 3 based

# UCX Framework Mission

- Collaboration between industry, laboratories, and academia
- Create open-source production grade communication framework for HPC applications
- Enable the highest performance through co-design of software-hardware interfaces
- Unify industry - national laboratories - academia efforts

## API

Exposes broad semantics that target data centric and HPC programming models and applications

## Performance oriented

Optimization for low-software overheads in communication path allows near native-level performance

## Production quality

Developed, maintained, tested, and used by industry and researcher community

## Community driven

Collaboration between industry, laboratories, and academia

## Research

The framework concepts and ideas are driven by research in academia, laboratories, and industry

## Cross platform

Support for Infiniband, Cray, various shared memory (x86-64 and Power), GPUs

Co-design of Exascale Network APIs

# Architecture

# UCX Framework

## UC-P for Protocols

High-level API uses UCT framework to construct protocols commonly found in applications

### Functionality:

Multi-rail, device selection, pending queue, rendezvous, tag-matching, software-atomics, etc.

## UC-T for Transport

Low-level API that expose basic network operations supported by underlying hardware. Reliable, out-of-order delivery.

### Functionality:

Setup and instantiation of communication operations.

## UC-S for Services

This framework provides basic infrastructure for component based programming, memory management, and useful system utilities

### Functionality:

Platform abstractions, data structures, debug facilities.

# A High-level Overview

## Applications

MPICH, Open-MPI, etc.

OpenSHMEM, UPC, CAF, X10,  
Chapel, etc.

Parsec, OCR, Legions, etc.

Burst buffer, ADIOS, etc.

UCX

### UC-P (Protocols) - High Level API

Transport selection, cross-transport multi-rail, fragmentation, operations not supported by hardware

Message Passing API Domain:  
tag matching, rendezvous

PGAS API Domain:  
RMAs, Atomics

Task Based API Domain:  
Active Messages

I/O API Domain:  
Stream

### UC-T (Hardware Transports) - Low Level API

RMA, Atomic, Tag-matching, Send/Recv, Active Message

Transport for InfiniBand VERBS driver

RC

UD

XRC

DCT

Transport for Gemini/Aries drivers

GNI

Transport for intra-node host memory communication

SYSV

POSIX

KNEM

CMA

XPMM

Transport for Accelerator Memory communication

GPU

### UC-S (Services)

Common utilities

Utilities

Data structures

Memory Management

OFA Verbs Driver

Cray Driver

OS Kernel

Cuda

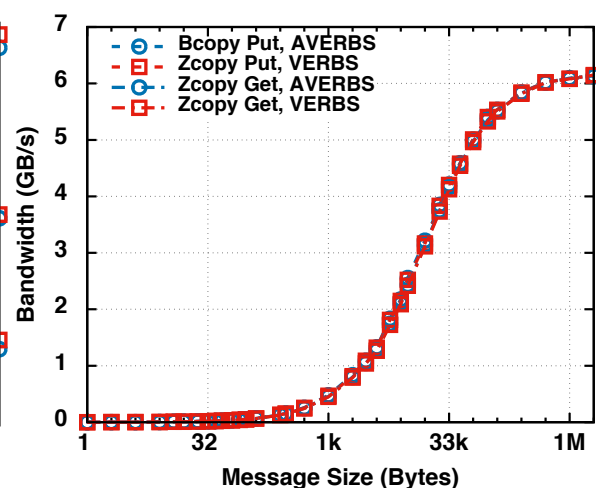
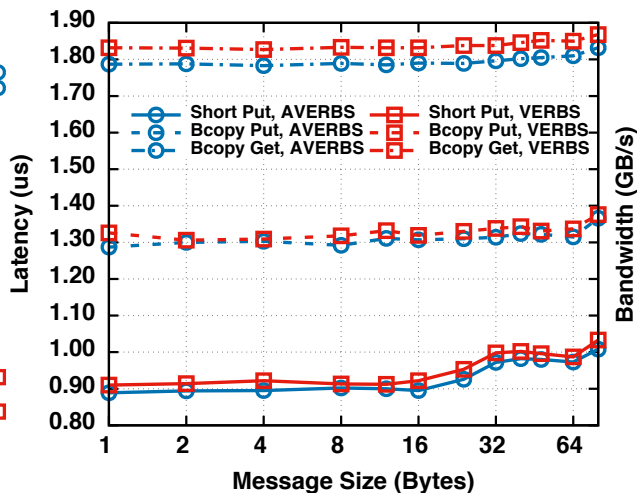
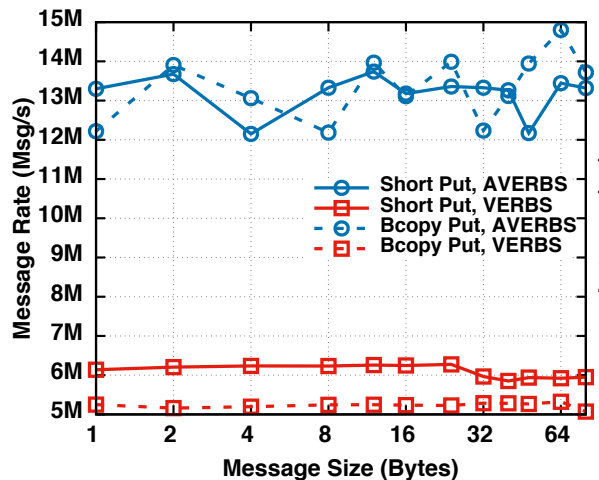
Hardware

# UCP API (DRAFT) Snippet

(<https://github.com/openucx/ucx/blob/master/src/ucp/api/ucp.h>)

- **ucs\_status\_t ucp\_put(ucp\_ep\_h ep, const void \*buffer, size\_t length, uint64\_t remote\_addr, ucp\_rkey\_h rkey)**  
*Blocking remote memory put operation.*
- **ucs\_status\_t ucp\_put\_nbi(ucp\_ep\_h ep, const void \*buffer, size\_t length, uint64\_t remote\_addr, ucp\_rkey\_h rkey)**  
*Non-blocking implicit remote memory put operation.*
- **ucs\_status\_t ucp\_get(ucp\_ep\_h ep, void \*buffer, size\_t length, uint64\_t remote\_addr, ucp\_rkey\_h rkey)**  
*Blocking remote memory get operation.*
- **ucs\_status\_t ucp\_get\_nbi(ucp\_ep\_h ep, void \*buffer, size\_t length, uint64\_t remote\_addr, ucp\_rkey\_h rkey)**  
*Non-blocking implicit remote memory get operation.*
- **ucs\_status\_t ucp\_atomic\_add32(ucp\_ep\_h ep, uint32\_t add, uint64\_t remote\_addr, ucp\_rkey\_h rkey)**  
*Blocking atomic add operation for 32 bit integers.*
- **ucs\_status\_t ucp\_atomic\_add64(ucp\_ep\_h ep, uint64\_t add, uint64\_t remote\_addr, ucp\_rkey\_h rkey)**  
*Blocking atomic add operation for 64 bit integers.*
- **ucs\_status\_t ucp\_atomic\_fadd32(ucp\_ep\_h ep, uint32\_t add, uint64\_t remote\_addr, ucp\_rkey\_h rkey, uint32\_t \*result)**  
*Blocking atomic fetch and add operation for 32 bit integers.*
- **ucs\_status\_t ucp\_atomic\_fadd64(ucp\_ep\_h ep, uint64\_t add, uint64\_t remote\_addr, ucp\_rkey\_h rkey, uint64\_t \*result)**  
*Blocking atomic fetch and add operation for 64 bit integers.*
- **ucs\_status\_ptr\_t ucp\_tag\_send\_nb(ucp\_ep\_h ep, const void \*buffer, size\_t count, ucp\_datatype\_t datatype, ucp\_tag\_t tag, ucp\_send\_callback\_t cb)**  
*Non-blocking tagged-send operations.*
- **ucs\_status\_ptr\_t ucp\_tag\_recv\_nb(ucp\_worker\_h worker, void \*buffer, size\_t count, ucp\_datatype\_t datatype, ucp\_tag\_t tag, ucp\_tag\_t tag\_mask, ucp\_tag\_recv\_callback\_t cb)**  
*Non-blocking tagged-receive operation.*

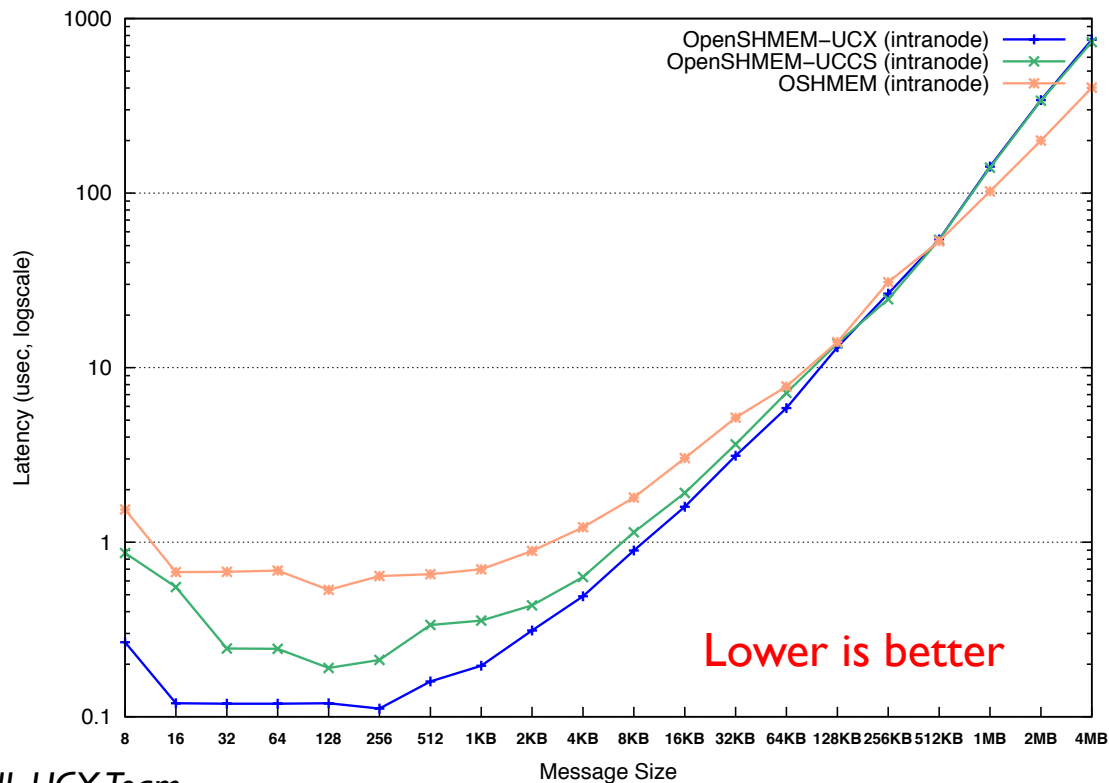
# Preliminary Evaluation ( UCT )



- Pavel Shamis, et al. "UCX: An Open Source Framework for HPC Network APIs and Beyond," HOT Interconnects 2015 - Santa Clara, California, US, August 2015
  - Two HP ProLiant DL380p Gen8 servers
  - Mellanox SX6036 switch, Single-port Mellanox Connect-IB FDR (10.10.5056)
  - Mellanox OFED 2.4-1.0.4. (VERBS)
  - Prototype implementation of Accelerated VERBS (AVERBS)



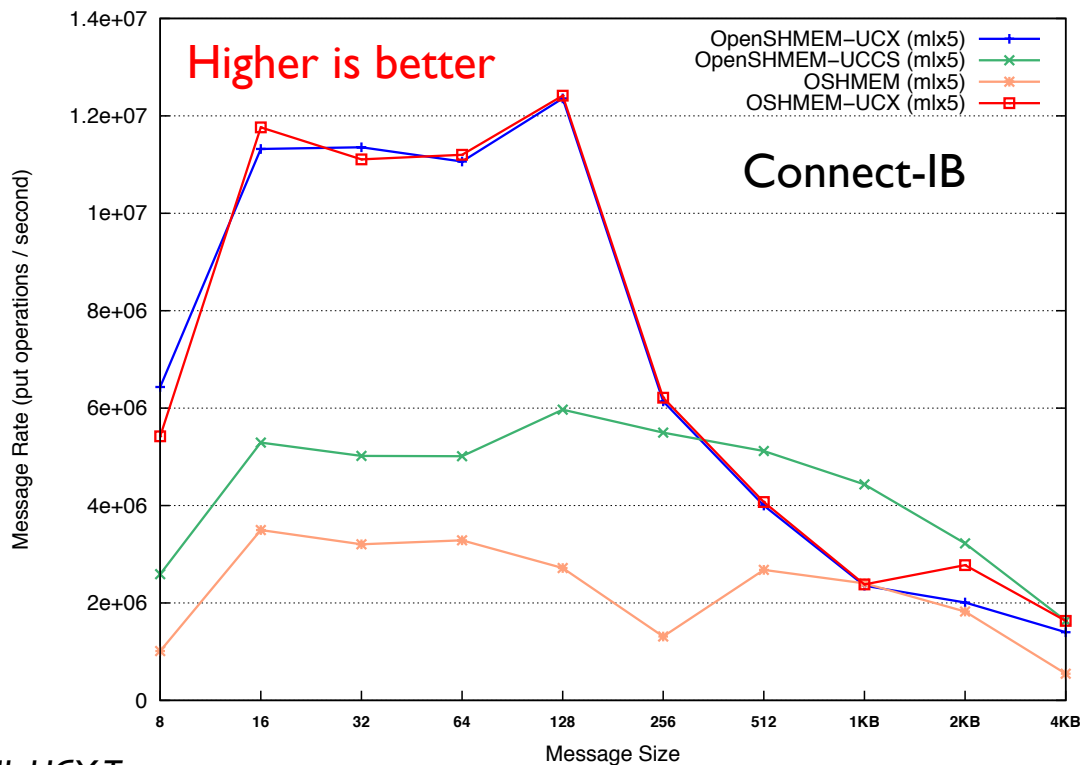
# OpenSHMEM and OSHMEM (OpenMPI) Put Latency (shared memory)



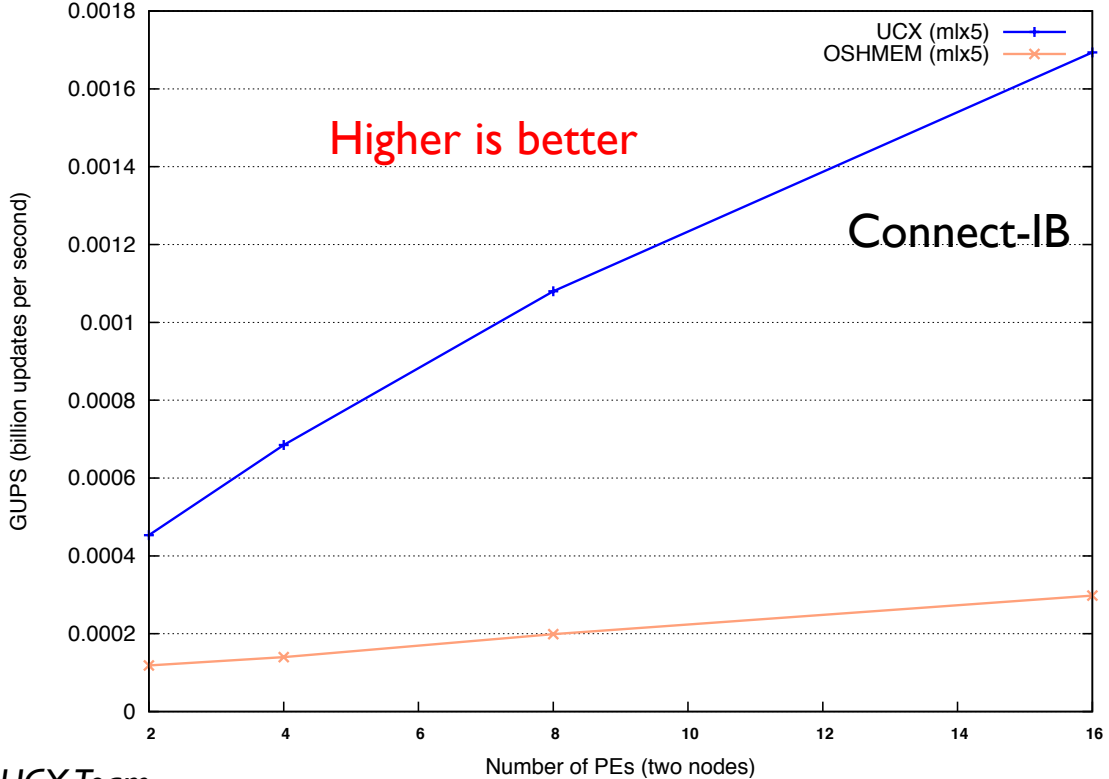
Slide courtesy of ORNL UCX Team



# OpenSHMEM and OSHMEM (OpenMPI) Put Injection Rate



# OpenSHMEM and OSHMEM (OpenMPI) GUPs Benchmark

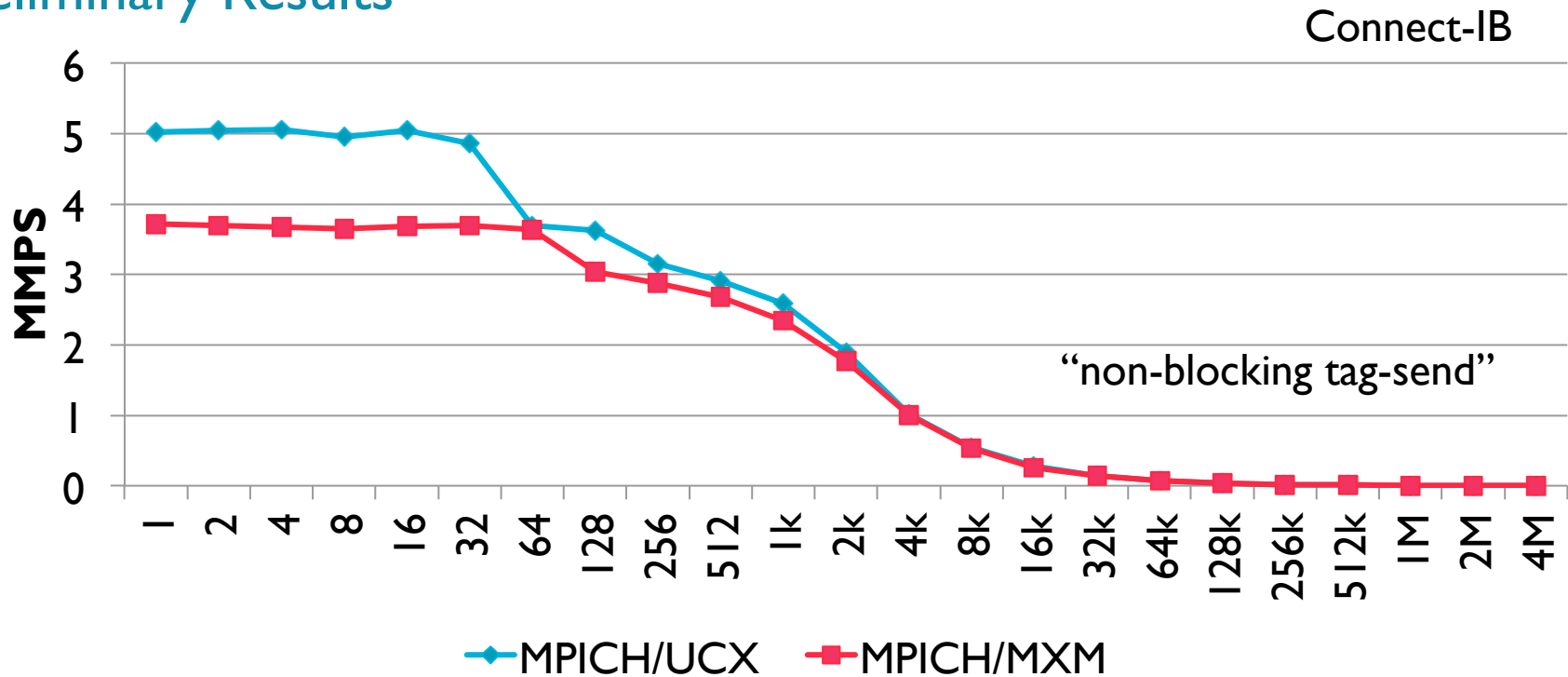


Slide courtesy of ORNL UCX Team



# MPICH - Message rate

## Preliminary Results



# Where is UCX being used?

- Upcoming release of Open MPI 2.0 (MPI and OpenSHMEM APIs)
- Upcoming release of MPICH
- OpenSHMEM reference implementation by UH and ORNL
- PARSEC – runtime used on Scientific Linear Libraries

# What Next ?

- UCX Consortium !
  - <http://www.csm.ornl.gov/newsite/>
- UCX Specification
  - Early draft is available online:  
<http://www.openucx.org/early-draft-of-ucx-specification-is-here/>
- Production releases
  - MPICH, Open MPI, Open SHMEM(s), Gasnet, and more...
- Support for more networks and applications and libraries
- UCX Hackathon 2016 !
  - Will be announced on the mailing list and website



Fork me on GitHub

<https://github.com/orgs/openucx>

WEB: [www.openucx.org](http://www.openucx.org)

Contact: [info@openucx.org](mailto:info@openucx.org)

Mailing List:

<https://elist.ornl.gov/mailman/listinfo/ucx-group>  
[ucx-group@elist.ornl.gov](mailto:ucx-group@elist.ornl.gov)

# Questions ?



## Unified Communication - X Framework



: [www.openucx.org](http://www.openucx.org)

Contact: [info@openucx.org](mailto:info@openucx.org)

WE B: <https://github.com/orgs/openucx>

Mailing List:

<https://elist.ornl.gov/mailman/listinfo/ucx-group>  
[ucx-group@elist.ornl.gov](mailto:ucx-group@elist.ornl.gov)

