



**MVAPICH**

MPI, PGAS and Hybrid MPI+PGAS Library

# Boosting the Performance of HPC Applications with MVAPICH and MVAPICH-Plus

A Tutorial at MUG'25

Presented by

**Nathaniel Shineman and Benjamin Michalowicz**

**The MVAPICH Team**

The Ohio State University

<http://mvapich.cse.ohio-state.edu/>

# Overview of the MVAPICH Project

- High Performance open-source MPI Library
- Support for multiple interconnects
  - InfiniBand, Omni-Path, Ethernet/iWARP, RDMA over Converged Ethernet (RoCE), AWS EFA, OPX, Broadcom RoCE, Intel Ethernet, Rockport Networks, Slingshot 10/11
- Support for multiple platforms
  - X86-Xeon, AMD EPYC, OpenPOWER, ARM, Xeon-Phi, GPUs (NVIDIA, Intel, and AMD)
- Started in 2001, first open-source version demonstrated at SC '02
- Supports the latest MPI-4.1 standard
- <http://mvapich.cse.ohio-state.edu>
- Additional optimized versions for different systems/environments:
  - MVAPICH-Plus (Unification of MVAPICH2-X and MVAPICH2-GDR), since 2023
  - MVAPICH2-X (Advanced MPI + PGAS), since 2011
  - MVAPICH2-GDR with support for NVIDIA (since 2014) and AMD (since 2020) GPUs
  - MVAPICH2-MIC with support for Intel Xeon-Phi, since 2014
  - MVAPICH2-Virt with virtualization support, since 2015
  - MVAPICH2-EA with support for Energy-Awareness, since 2015
  - MVAPICH2-Azure for Azure HPC IB instances, since 2019
  - MVAPICH2-X-AWS for AWS HPC+EFA instances, since 2019
- Tools:
  - **OSU MPI Micro-Benchmarks (OMB), since 2003**
  - OSU InfiniBand Network Analysis and Monitoring (INAM), since 2015



- **Used by more than 3,450 organizations in 92 countries (listed under the Users Tab of the MVAPICH page)**
- **More than 1.93 Million downloads from the OSU site directly**
- Empowering many TOP500 clusters (Jun'25 ranking)
  - 21<sup>st</sup>, 10,649,600-core (Sunway TaihuLight) at NSC, Wuxi, China
  - 67<sup>th</sup>, 448, 448 cores (Frontera) at TACC
  - 88<sup>th</sup>, 288,288 cores (Lassen) at LLNL
  - 109<sup>th</sup>, 570,020 cores (Nurion) in South Korea and many others
- Available with software stacks of many vendors and Linux Distros (RedHat, SuSE, OpenHPC, and Spack)
- Partner in the 67<sup>th</sup> ranked TACC Frontera system
- **Empowering Top500 systems for more than 20+ years**

# Architecture of MVAPICH Software Family for HPC and DL/ML

## High Performance Parallel Programming Models

Message Passing Interface  
(MPI)

PGAS  
(UPC, OpenSHMEM, CAF, UPC++)

Hybrid --- MPI + X  
(MPI + PGAS + OpenMP/Cilk)

## High Performance and Scalable Communication Runtime

### Diverse APIs and Mechanisms

Point-to-point Primitives

Collectives Algorithms

Job Startup

Energy-Awareness

Remote Memory Access

I/O and File Systems

Fault Tolerance

Virtualization

Active Messages

Introspection & Analysis

CH4 Netmod

ch4:ucx

ch4:ofi

CH4 Shmmod

posix

ipc

Support for Modern Networking Technology  
(InfiniBand, RoCE, Omni-Path, EFA, Slingshot, etc.)

Transport Protocols

RC

SRD

UD

DC

Modern Features

UMR

ODP

SR-IOV

Multi Rail

Support for Modern Multi-/Many-core Architectures  
(Intel-Xeon, AMD EPYC, ARM, NVIDIA/AMD/Intel GPU)

Transport Mechanisms

POSIX Shared Memory

CMA

XPMM

Modern Features

CXL

NVLink

UALink

\* Upcoming

# Updated User Guide and Resources

- <https://mvapich-docs.readthedocs.io/en/latest/>
- Latest RPM install instructions, CVAR names/aliases/etc.
- Runtime instructions, setting up environment variables, OMB installation and execution, FAQ

# Contents

- CPU Process Mapping with MVAPICH
  - Summary of parallel launchers
- CH4 devices and their impact
  - UCX, OFI, enhanced libfabrics from OSU
- Dynamic Tuning
  - CPU, GPU dynamic tuning
- GPU-Awareness
  - NVIDIA, AMD, Intel GPU support
- Compression (On-The-Fly for Point-to-Point and Collectives)
  - NVIDIA Support
  - AMD Support
- Future Plans

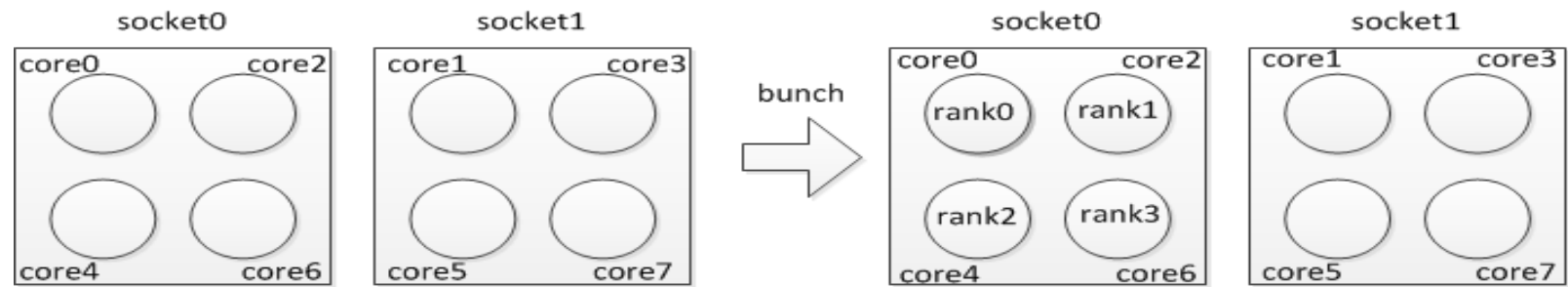
# Common Parallel Launchers

- Hydra – included with MVAPICH/MVAPICH-Plus
  - mpiexec/mpirun binaries
  - Compatible with all commonly used resource managers
    - SLURM, PBS, Flux, etc.
    - Uses PMI1 interface standard
- Srun – included with the Slurm resource manager
  - Mutually exclusive with hydra and other launchers
    - Uses slurm specific PMI2
    - --with-pm=none --with-pmi=slurm

# CPU Process Mapping

- Determines Process-to-Core mapping At Runtime via resource managers and parallel job launchers
- Supported by both Hydra and srun
  - Hydra: “`mpiexec/mpirun ... -bind-to <x> -map-by <y>`” (and use of [hwloc](#))
  - Srun: “`srun --cpu-bind=<x>`”
    - See [slurm documentation](#) for more details
- Common cases
  - “Bunch”
  - “Scatter”
  - “Spread”
  - Custom mappings
- To see process mappings at runtime:
  - Hydra: set `HYDRA_TOPO_DEBUG=1`
  - Srun: set “`SLURM_CPU_BIND_VERBOSE=verbose`” or “`srun --cpu-bind=verbose[,options]`”

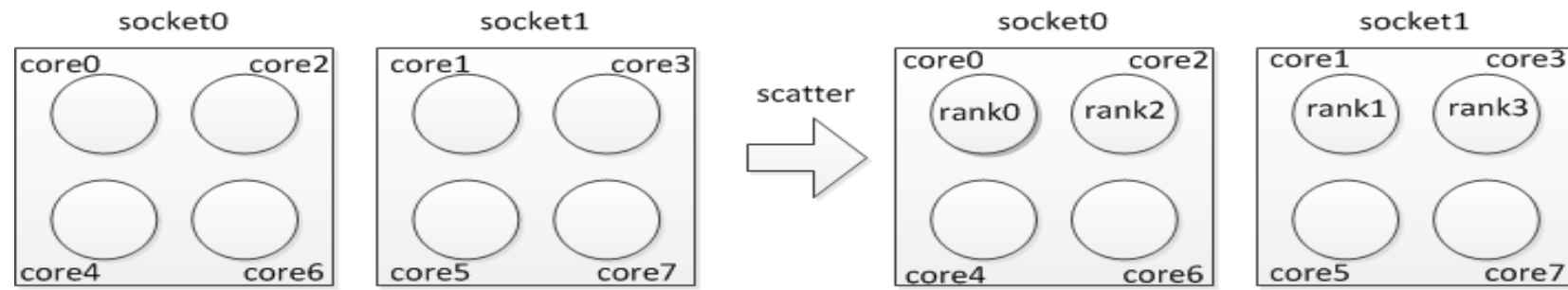
# Bunch Mapping



- Allocates ranks sequentially on cores
  - Optimal for near-neighbor communication
  - Can heavily leverage shared memory, CMA, XPMEM
  - Does not provide optimal resource distribution of memory, network devices, GPU, etc
  - Useful when fully subscribing nodes
- Launcher commands:
  - Hydra: `mpiexec -np <num_tasks> ... -map-by core -bind-to core`
  - Srun: `srun -n <num_tasks> ... --cpu-bind=core --distribution=block:block`



# Scatter Mapping



- Performs a round-robin assignment per-socket
  - Prioritizes network device and memory access over intra-node transfers
  - Can be applied to other map-by resources, ie NUMA or L3 cache
  - Not optimal for applications bound by intra-node communication
  - Useful for sparse allocations with heavy network use and multiple network adapters
- Launcher commands:
  - Hydra: `mpiexec -np <num_tasks> ... -map-by socket -bind-to core`
  - Srun: `srun -n <num_tasks> ... --cpu-bind=core --distribution=cyclic:cyclic`

# Spread Mapping

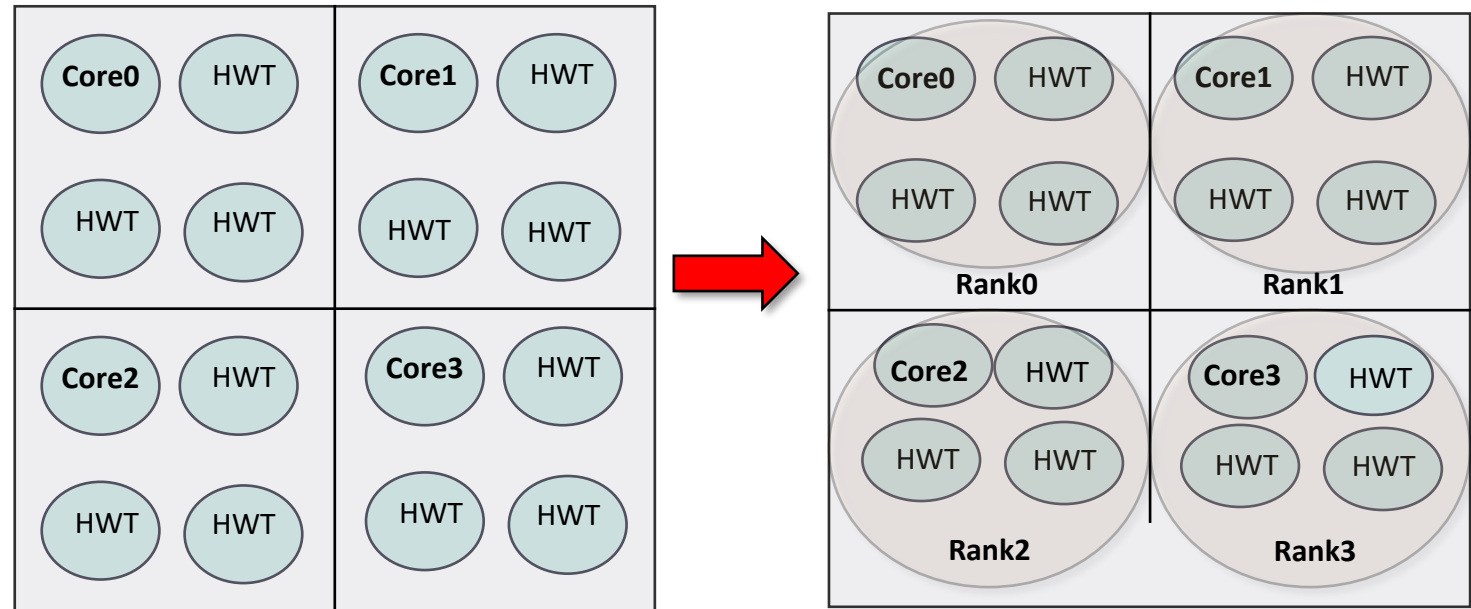
- Handles MPI+Threads workloads

- Works with and without hyperthreading activated for a compute node

- Best practice: divide total cores by requested ppn and map each process to N cores

- Launcher Commands:

- Hydra: `mpiexec --bind-to core --map-by core:<total_cores / ppn>`
- Srun: `srun -n <num_tasks> ... --cpus-per-task=<total_cores / ppn> --ntasks-per-node=<ppn> --cpu-bind=core`



# Custom Mappings

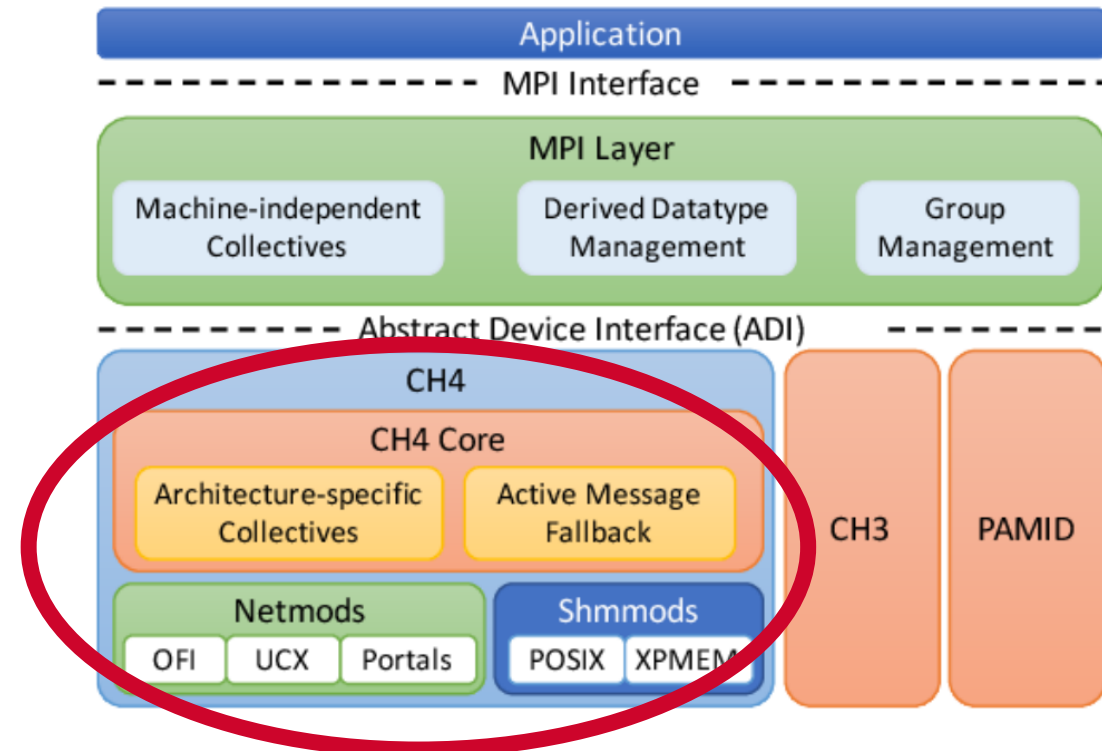
- Not all standard mappings work best for all applications
- Need fine-tuning/architecture/numa-aware process placement
- Hydra custom mappings:
  - Supports binding/mapping to NUMA, socket, core, hwthread, l1/2/3/4/5 cache, pci devices, GPUs, IB HCAs, eth devices, hfi, etc.
- SLURM/srun custom mappings:
  - See [https://slurm.schedmd.com/mc\\_support.html](https://slurm.schedmd.com/mc_support.html) for options
  - Additionally supports binding to [memory regions](#) with “--mem-bind=...”

# Contents

- CPU Process Mapping with MVAPICH
  - Summary of parallel launchers
- CH4 devices and their impact
  - UCX, OFI, enhanced libfabrics from OSU
- Dynamic Tuning
  - CPU, GPU dynamic tuning
- GPU-Awareness
  - NVIDIA, AMD, Intel GPU support
- Compression (On-The-Fly for Point-to-Point and Collectives)
  - NVIDIA Support
  - AMD Support
- Future Plans

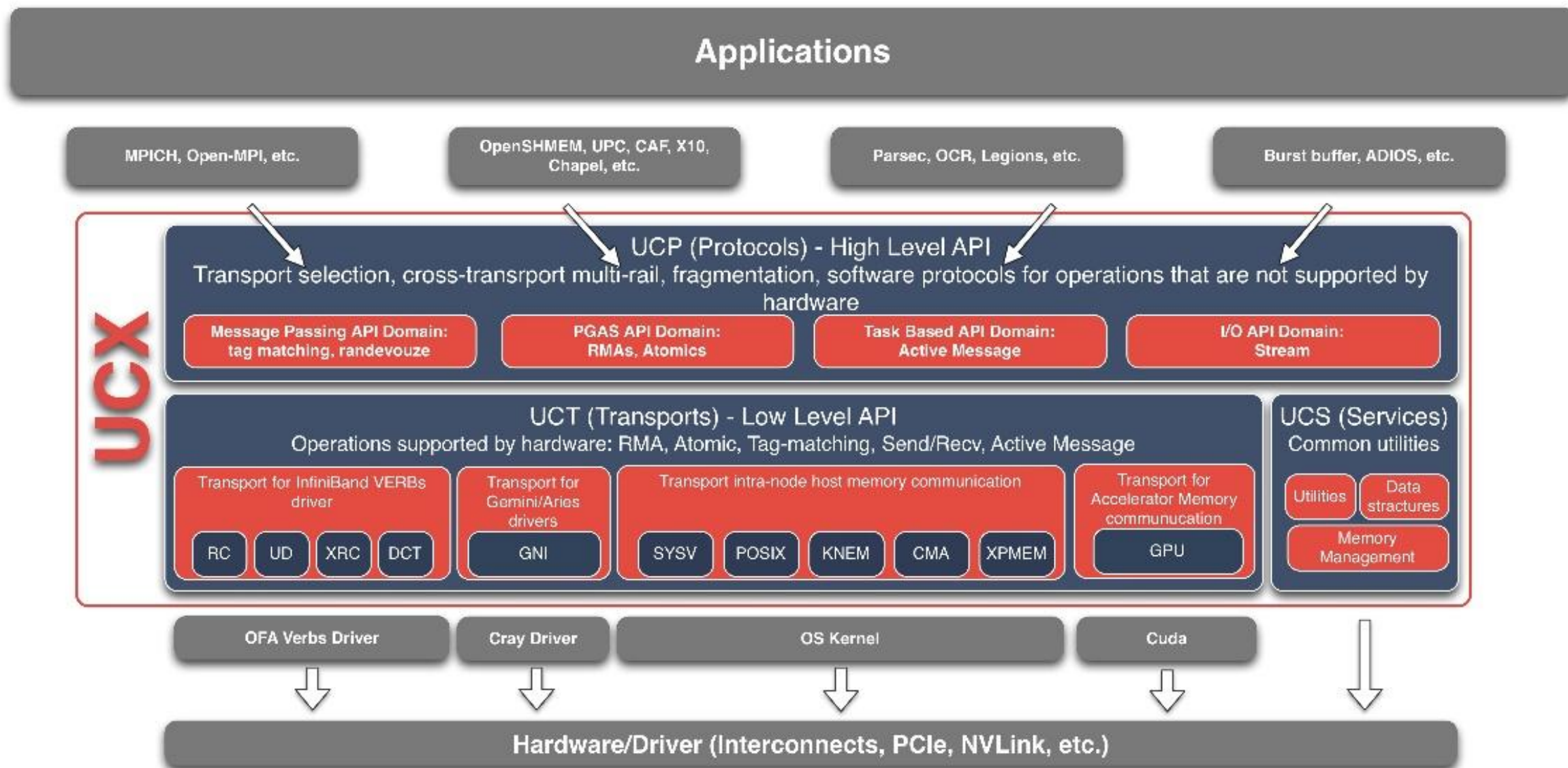
# CH4 Devices – Overview

- CH4 device
  - Replaced CH3 with MPICH 3.4 in 2021
- 2 components to CH4 device
  - Netmods (Network)
  - Shmmods (Shared Memory)
- MVAPICH Netmodes:
  - UCX (Unified Communications X)
    - Recommended for IB/RoCE networks
  - OFI (Open Fabrics Interfaces/Libfabrics)
    - Recommended for all other networks
- Configure with
  - `--with-device=ch4:ofi --with-device=ch4:ucx`



Courtesy: <https://multicore.world/wp-content/uploads/2024/02/multicoreworld24-brightwell.pdf>

# UCX: Overview

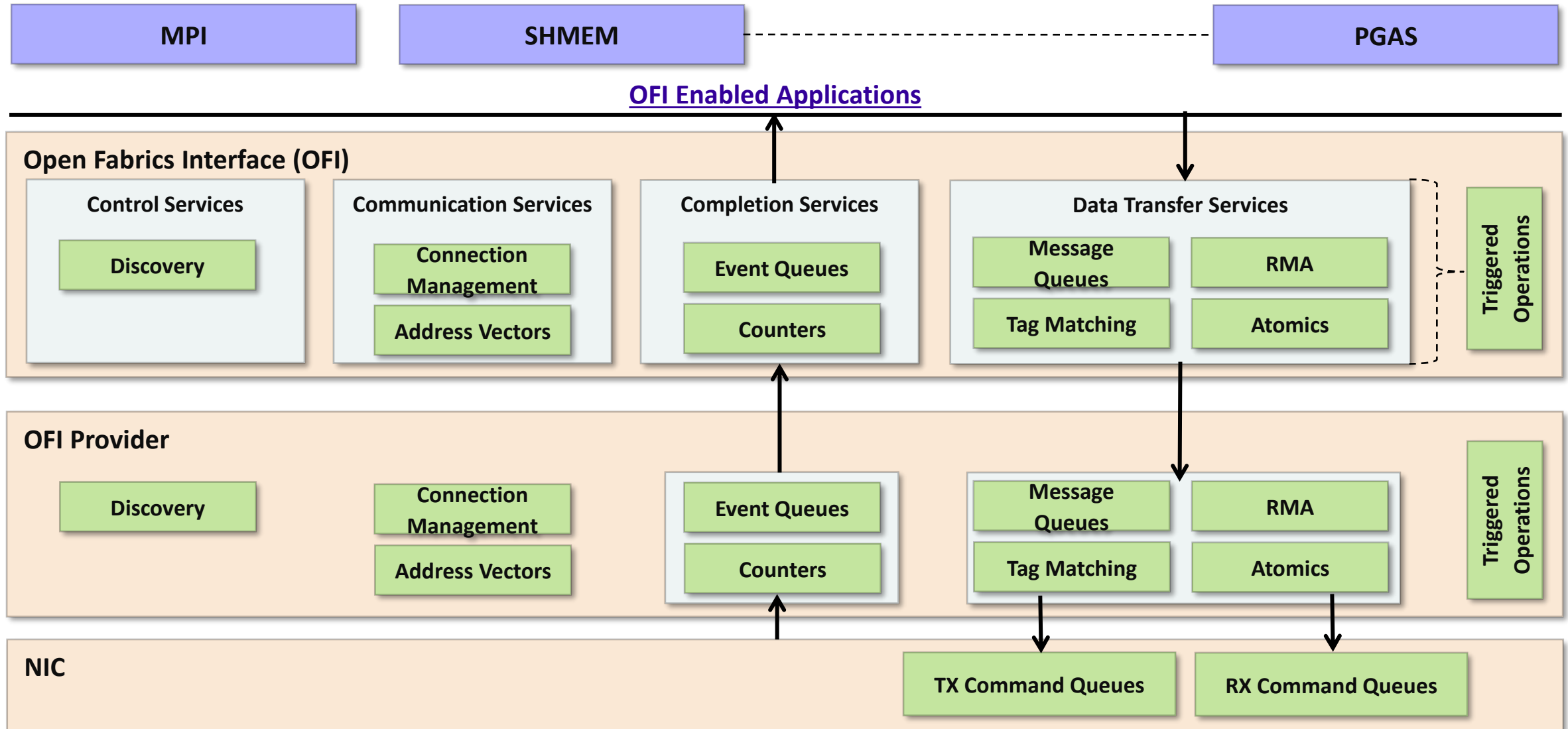


Courtesy: <https://www.openucx.org/>

# Using the UCX Netmod

- Configure UCX path with `--with-ucx=/path/to/ucx` or `embedded` options
- Common UCX environment variables:
  - UCX\_RNDV\_THRESH: transition point for eager/rendezvous message transfer protocols; provide a number in bytes
  - UCX\_MAX\_RNDV/EAGER\_RAILS: number of network devices used per process for eager and rendezvous-based transfers  
(maximum supported is 4 (<https://openucx.readthedocs.io/en/master/faq.html>))
  - UCX\_TLS: comma-separated list of transports used by UCX (rc, ud, shm, cma, xpmem, etc.)
  - UCX\_NET\_DEVICES: sets the device(s) UCX will use for message transports
  - UCX\_LOG\_LEVEL: logging information in various verbiages (trace, info, debug, error, fatal)

# OFI: Overview



Courtesy: <http://www.slideshare.net/seanhefty/ofi-overview?ref=http://ofiwg.github.io/libfabric/>



# Using the OFI Netmod

- Configure OFI path with `--with-libfabric=/path/to/libfabric` or `embedded` options
- Common OFI environment variables:
  - FI\_PROVIDER: forces a specific network provider to be used
    - MVAPICH will attempt to choose the highest performing provider that supports your network
    - Generally not necessary unless a particular network is required for testing
  - MPIR\_CVAR\_CH4\_OFI\_CAPABILITY\_SETS\_DEBUG: set to 1 to view the OFI provider being selected at runtime
  - MVP\_CH4\_OFI\_ENABLE\_HMEM: set to 1 to allow OFI provider to use GPU direct RDMA support
  - Highly provider dependent
  - More details at <https://ofiwg.github.io/libfabric/v2.2.0/man/>

# Cray OFI (Slingshot 11)

- Provided on all HPE/Cray machines equipped with Slingshot 11 interconnects
  - (Frontier, El Capitan, Isambard-AI, Delta-AI, COSMOS, etc.)
- Must be linked when building MVAPICH on HPE/Cray systems
- CXI provider up-streaming in progress
- Linked at compile time
  - Use MVAPICH-Plus ums038 module on frontier
  - Request system specific installation on other sites
  - Link directly to craype libfabric when compiling open source

# Contents

- CPU Process Mapping with MVAPICH
  - Summary of parallel launchers
- CH4 devices and their impact
  - UCX, OFI, enhanced libfabrics from OSU
- Dynamic Tuning
  - CPU, GPU dynamic tuning
- GPU-Awareness
  - NVIDIA, AMD, Intel GPU support
- Compression (On-The-Fly for Point-to-Point and Collectives)
  - NVIDIA Support
  - AMD Support
- Future Plans

# Dynamic Collective Tuning

- New in MVAPICH-Plus 4.1
- Enables on-the-fly collective selection tuning
- Combines collective knowledge with real-time, adaptive, system-specific analysis
- Runs potential algorithms and collects runtime data to determine best collective algorithm at a given time

# Controlling Dynamic Tuning with CVARs

- MVP\_ENABLE\_DYNAMIC\_TUNING – allows dynamic tuning to be activated or deactivated at runtime or within an application
- MVP\_FORCE\_STATIC\_TUNING – force MVAPICH to use static tables
- MVP\_DYNAMIC\_TUNE\_THRESHOLD – when to retest algorithms
- MVP\_DYNAMIC\_TUNE\_TEST\_COUNT – how many tests to perform with each algorithm
- MVP\_DYNAMIC\_TUNE\_LOCAL\_EVAL\_METRIC – how each process reports its "best" result
- MVP\_DYNAMIC\_TUNE\_SELECTION\_METRIC – metric used by the dynamic tuning to select the overall best algorithm

# Dynamic Tuning Use Cases

- Improves microbenchmark performance on unknown systems
  - Match iteration count with MVP\_DYNAMIC\_TUNE\_THRESHOLD for best results
  - Set `-x 200 -i 1000` for best performance on OMB
- Respond to application specific performance
- Handle real-time network constraints
- Feedback encouraged

# Contents

- CPU Process Mapping with MVAPICH
  - Summary of parallel launchers
- CH4 devices and their impact
  - UCX, OFI, enhanced libfabrics from OSU
- ⑩ Dynamic Collective Tuning
  - ✂ CPU, GPU dynamic tuning
- GPU-Awareness
  - NVIDIA, AMD, Intel GPU support
- Compression (On-The-Fly for Point-to-Point and Collectives)
  - NVIDIA Support
  - AMD Support
- Future Plans

# MPI + CUDA - Naive

- Data movement in applications with standard MPI and CUDA interfaces

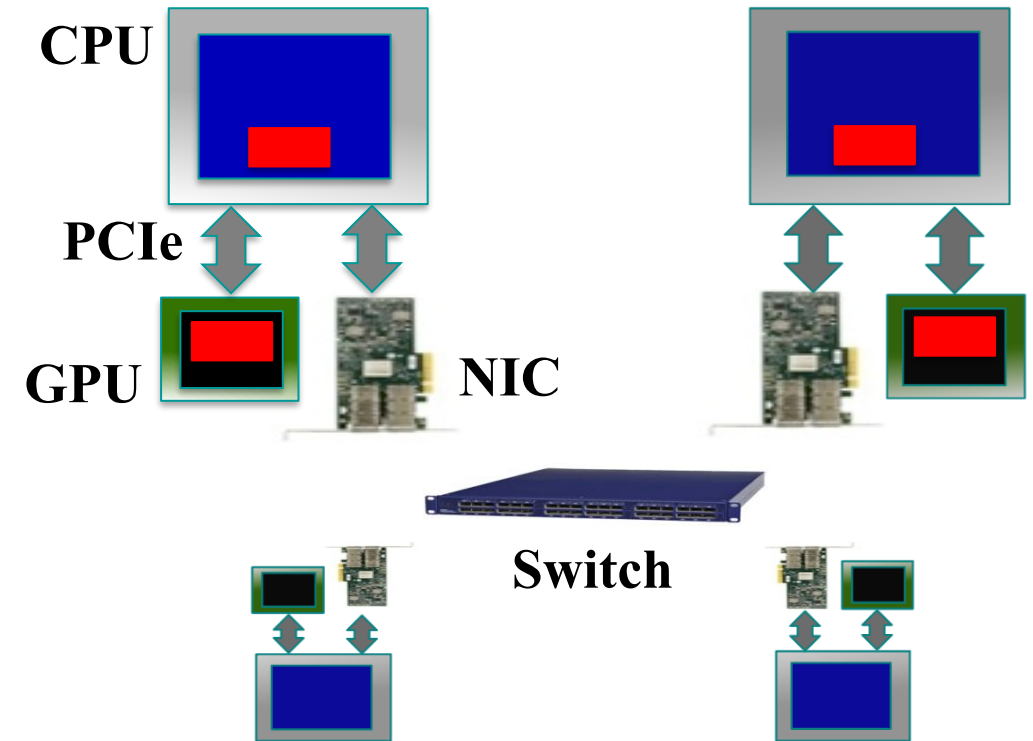
## At Sender:

```
cudaMemcpy(s_hostbuf, s_devbuf, . . .);  
MPI_Send(s_hostbuf, size, . . .);
```

## At Receiver:

```
MPI_Recv(r_hostbuf, size, . . .);  
cudaMemcpy(r_devbuf, r_hostbuf, . . .);
```

*High Productivity and Low Performance*





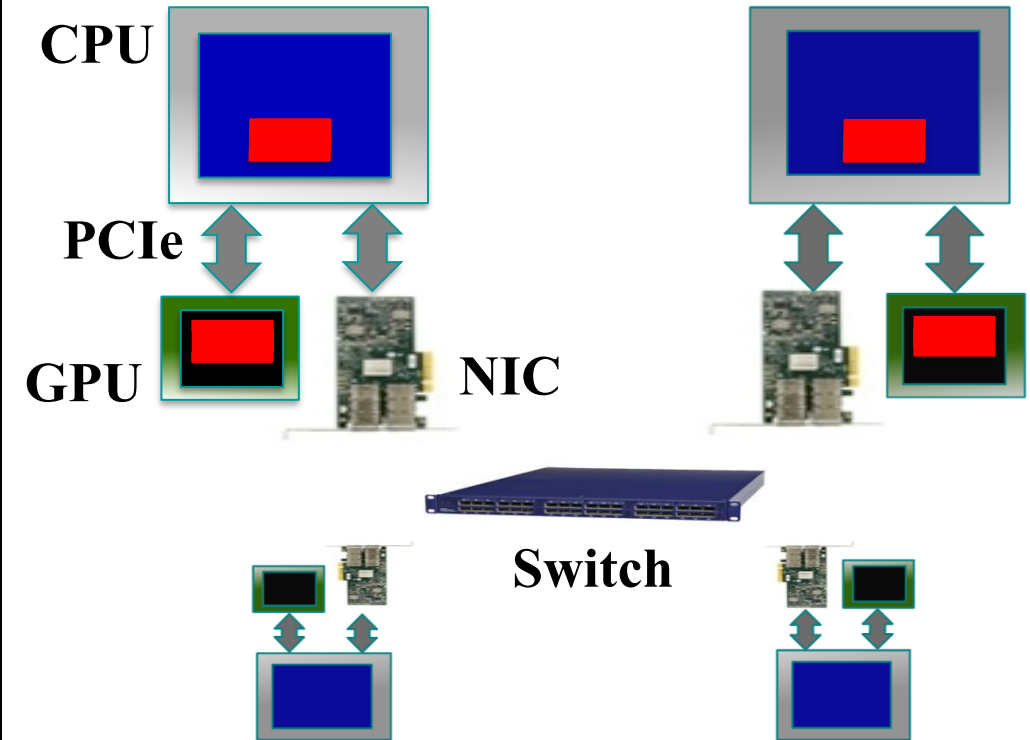
# MPI + CUDA - Advanced

- Pipelining at user level with non-blocking MPI and CUDA interfaces

## At Sender:

```
for (j = 0; j < pipeline_len; j++)  
    cudaMemcpyAsync(s_hostbuf + j * blk, s_devbuf + j *  
        blksize, ...);  
for (j = 0; j < pipeline_len; j++) {  
    while (result != cudaSuccess) {  
        result = cudaStreamQuery(...);  
        if(j > 0) MPI_Test(...);  
    }  
    MPI_Isend(s_hostbuf + j * block_size, blksize . . .);  
}  
MPI_Waitall();
```

<<Similar at receiver>>



*Low Productivity and High Performance*

# GPU-Aware (CUDA-Aware) MPI Library: MVAPICH-Plus

- Standard MPI interfaces used for unified data movement
- Takes advantage of Unified Virtual Addressing ( $\geq$  CUDA 4.0)
- Overlaps data movement from GPU with RDMA transfers

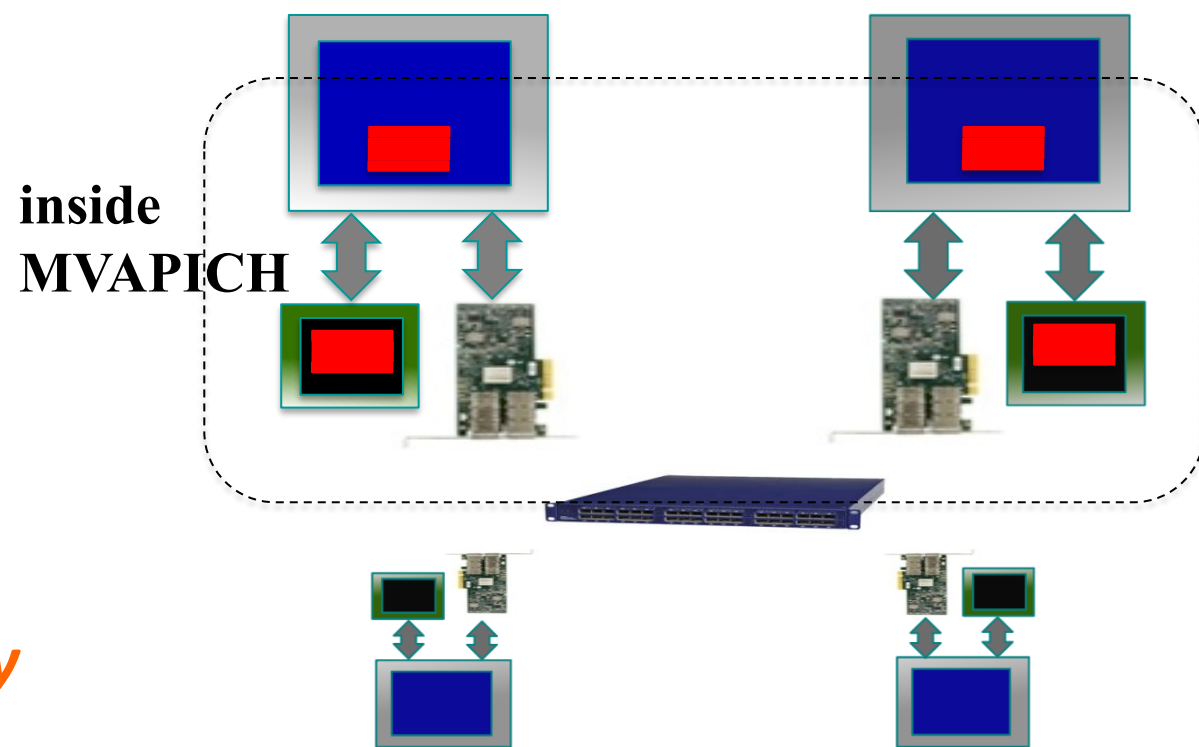
## At Sender:

```
MPI_Send(s_devbuf, size, ...);
```

## At Receiver:

```
MPI_Recv(r_devbuf, size, ...);
```

*High Performance and High Productivity*



# GPU Support in MVAPICH

- Recommended in MVAPICH-Plus
  - Basic support only in MVAPICH, uses base MPICH implementation
- Supports NVIDIA, AMD, and Intel GPUs
- Includes optimized designs for pt2pt and collective operations
- Enhanced GPU kernels for certain collectives
- Utilizes GDRCOPY, GPU IPC, and kernel based transfers

# Tuning GDRCOPY Designs in MVAPICH-Plus

Parameter	Significance	Default	Notes
MVP_ENABLE_GPU_GDRCOPY	• Enable / Disable GDRCOPY-based designs	1 (Enabled)	• Always enable • Enables NVIDIA GDRCOPY and AMD LargeBAR support
MVP_GDRCOPY_MAX_SIZE_H2D	• Maximum message size for using GDRCOPY from Host to Device	32768 Bytes	• Tune for your system • AMD and newer NVIDIA systems can use higher values
MVP_GDRCOPY_MAX_SIZE_D2H	• Maximum message size for using GDRCOPY from Device to Host	2048 Bytes	• Tune for your system • Typically lower than H2D for NVIDIA systems • Set equal to D2H for AMD

- Refer to **CVAR** section of MVAPICH-Plus user guide for more information
- <https://mvapich-docs.readthedocs.io/en/mvapich-plus/cvar.html>

# Tuning GPU IPC Designs in MVAPICH-Plus

Parameter	Significance	Default	Notes
MVP_CH4_IPC_GPU_P2P_THRESHOLD	• Minimum message size to use GPU IPC protocol	16K Bytes	• Tune for your system • Typically at or around the GDRCOPY H2D limit
MVP_CH4_IPC_GPU_READ_WRITE_PROTOCOL	• Select the IPC protocol used for GPU transfers (read/write)	read	• Generally keep on read
MVP_ALLREDCE_IPC_MESSAGE_SIZE_THRESHOLD	• Minimum message size to use IPC aware Allreduce	16K Bytes	• Tune for your system • IPC aware collectives have some registration overhead
MVP_REDCCE_IPC_MESSAGE_SIZE_THRESHOLD	• Minimum message size to use IPC aware Reduce	16K Bytes	• Tune for your system • IPC aware collectives have some registration overhead

- Refer to **CVAR** section of MVAPICH-Plus user guide for more information
- <https://mvapich-docs.readthedocs.io/en/mvapich-plus/cvar.html>

# Enhancements for APUs

- Supported by MVAPICH-Plus
  - Provides enhanced support for MI300A APUs
- Supports buffers allocated on CPU or GPU
- Enables flexible GPU transfers using both CPU and GPU techniques

# Tuning Unified APU Designs in MVAPICH-Plus

Parameter	Significance	Default	Notes
MVP_GPU_UNIFIED_DIRECT_COPY_MAX	<ul style="list-style-type: none"><li>Maximum number of bytes for which MVAPICH-Plus will use direct CPU memcpy support</li></ul>	4M	<ul style="list-style-type: none"><li>Tune for your system</li><li>Only effective on MI300A unified APU builds</li><li>Higher message sizes will use kernel driven methods</li></ul>

- Refer to **CVAR** section of MVAPICH-Plus user guide for more information
- <https://mvapich-docs.readthedocs.io/en/mvapich-plus/cvar.html>

# Contents

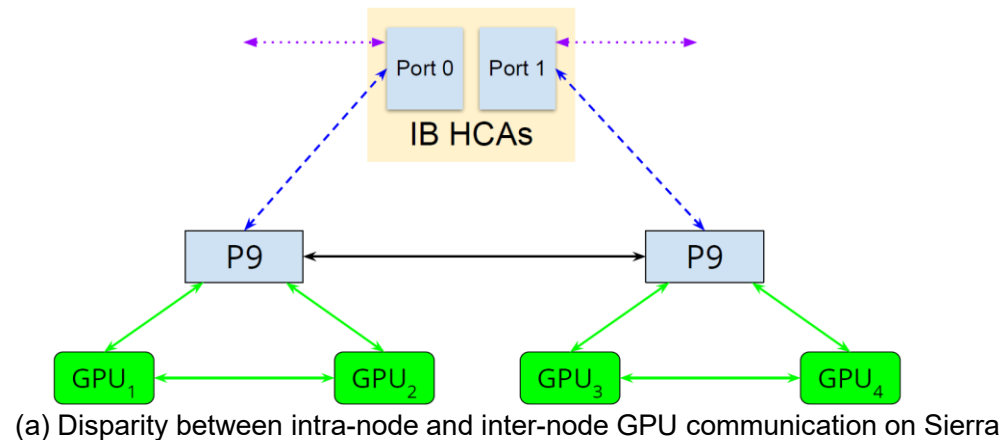
- CPU Process Mapping with MVAPICH
  - Summary of parallel launchers
- CH4 devices and their impact
  - UCX, OFI, enhanced libfabrics from OSU
- Dynamic Tuning
  - CPU, GPU dynamic tuning
- GPU-Awareness
  - NVIDIA, AMD, Intel GPU support
- Compression (On-The-Fly for Point-to-Point and Collectives)
  - NVIDIA Support
  - AMD Support
- Future Plans



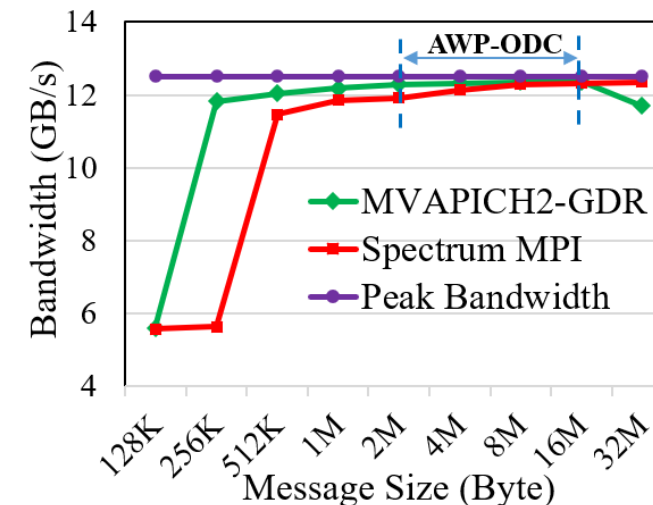
# MVAPICH-Plus: “On-the-fly” Compression – Motivation

- For HPC and data science applications on modern GPU clusters
  - With larger problem sizes, applications exchange **orders of magnitude more data** on the network
  - Leads to significant **increase in communication times** for these applications on larger scale (AWP-ODC)
  - On modern HPC systems, there is **disparity** between intra-node and inter-node GPU communication bandwidths that prevents efficient scaling of applications on larger GPU systems
  - CUDA-Aware MPI libraries **saturate the bandwidth** of IB network
  - **Compression** can reduce the data size and lower the pressure on network with limited bandwidth

→ 3-lane NVLink (75 GB/s)   ↔ X-Bus (64 GB/s)   ↔ 8-lane PCIe Gen4 (16 GB/s)   ↔ Infiniband EDR (12.5 GB/s)



OpenPOWER supercomputer [1]



[1] K. S. Khorassani, C.-H. Chu, H. Subramoni, and D. K. Panda, “Performance Evaluation of MPI Libraries on GPU-enabled OpenPOWER Architectures: Early Experiences”, in International Workshop on Open-POWER for HPC (IWOPH 19) at the 2019 ISC High Performance Conference, 2018.

# MVAPICH-Plus Compression Support

- ZFP
  - Lossy compression
  - Requires MVAPICH-Plus enhanced ZFP library
  - Bundled with MVAPICH-Plus on request
- Upcoming support for lossless compression algorithms
- Support for pt2pt and some collectives
  - Alltoall, Allgather, Allreduce, Reduce\_scatter

# Tuning Pt2pt Compression Designs in MVAPICH-Plus

Parameter	Significance	Default	Notes
MVP_ENABLE_PT2PT_GPU_COMPRESSION	<ul style="list-style-type: none"><li>• Enable/disable compression for pt2pt GPU transfers</li></ul>	0 (disabled)	<ul style="list-style-type: none"><li>• Enable only when application can tolerate compression related data loss</li></ul>
MVP_PT2PT_GPU_COMPRESSION_THRESHOLD	<ul style="list-style-type: none"><li>• Minimum message size to use pt2pt compression</li></ul>	1M Bytes	<ul style="list-style-type: none"><li>• Tune for your system/application</li></ul>

- Refer to **CVAR** section of MVAPICH-Plus user guide for more information
- <https://mvapich-docs.readthedocs.io/en/mvapich-plus/cvar.html>

# Tuning Collective Compression Designs in MVAPICH-Plus

Parameter	Significance	Default	Notes
MVP_ENABLE_COMPRESSION	<ul style="list-style-type: none"><li>• Enable/disable compression enabled collectives</li></ul>	0 (disabled)	<ul style="list-style-type: none"><li>• Enable only when application can tolerate compression related data loss</li></ul>
MVP_<coll>_GPU_COMPRESSION_THRESHOLD	<ul style="list-style-type: none"><li>• Minimum message size for using compression enabled collective</li></ul>	4M Bytes	<ul style="list-style-type: none"><li>• Supports alltoall, allgather, allreduce, reduce_scatter for coll</li><li>• Only effective when compression enabled algorithm is in use</li></ul>
MVP_<coll>_INTRA_ALGORITHM	<ul style="list-style-type: none"><li>• Set to osu_gpu_compression to utilize compression designs</li></ul>	auto/dynamic	<ul style="list-style-type: none"><li>• Compression designs will only be enabled when the osu_gpu_compression algorithm is forced</li></ul>

- Refer to **CVAR** section of MVAPICH-Plus user guide for more information
- <https://mvapich-docs.readthedocs.io/en/mvapich-plus/cvar.html>

# Contents

- CPU Process Mapping with MVAPICH
  - Summary of parallel launchers
- CH4 devices and their impact
  - UCX, OFI, enhanced libfabrics from OSU
- GPU-Awareness
  - NVIDIA, AMD, Intel GPU support
- Dynamic Tuning
  - CPU, GPU dynamic tuning
- Compression (On-The-Fly for Point-to-Point and Collectives)
  - NVIDIA Support
  - AMD Support
- Future Plans

# MVAPICH-Plus 4.2 Upcoming Release

- Incorporation of new research from MVAPICH team
  - Enhanced intra-node performance (CPU shared memory)
    - Alltoall, Allreduce, Reduce
  - Enhanced GPU support for AMD, Intel, NVIDIA GPUs
  - Unified Compression Layer
    - Support for lossless compression algorithms
- SHARP collective support

# OSU Microbenchmarks (OMB) Version 8.0

- Total refactor of MPI benchmarks for easier contributions
  - Reduced files, less duplicated code, easier patches
- Enhanced runtimes to facilitate easier MPI library testing
  - Similar to IMB all-in-one tests
  - Run multiple tests from one MPI launch command
    - ``mpiexec -np <np> ... ./omb_coll [benchmark1[,benchmark2,...]]>`
- Will remain backwards compatible

# OMB Public GitHub

- NOWLAB Github: <https://github.com/OSU-Nowlab>
  - Currently features many open-source developments from MVAPICH team
  - Increase opportunities for community engagement
- Coming with v8.0 release



# Funding Acknowledgments

## *Funding Support by*



## *Equipment Support by*



# Acknowledgments to all the Heroes (Past/Current Students and Staffs)

## Current Students (Under/Graduate)

- N. Alnaasan (Ph.D.)
- Q. Anthony (Ph.D.)
- C.-C. Chen (Ph.D.)
- T. Chen (Ph.D.)
- N. Contini (Ph.D.)
- S. Gumaste (Ph.D.)
- J. Hatef (Ph.D.)
- G. Kuncham (Ph.D.)
- S. Lee (Ph.D.)
- B. Michalowicz (Ph.D.)
- J. Oswal (Ph.D.)
- T. Tran (Ph.D.)
- L. Xu (P.h.D.)
- S. Xu (Ph.D.)
- J. Yao (Ph.D.)
- S. Zhang (Ph.D.)
- S. Mohammad (M.S.)
- B. Lampe (B.S.)
- N. Klein (B.S.)

## Current Research Specialist

- R. Motlagh

## Current Software Engineers

- N. Shineman
- M. Lieber

## Past Research Scientists

- K. Hamidouche
- S. Sur
- X. Lu
- M. Abduljabbar
- A. Shafi

## Past Students

- A. Awan (Ph.D.)
- A. Augustine (M.S.)
- P. Balaji (Ph.D.)
- M. Bayatpour (Ph.D.)
- R. Biswas (M.S.)
- S. Bhagvat (M.S.)
- A. Bhat (M.S.)
- D. Buntinas (Ph.D.)
- L. Chai (Ph.D.)
- B. Chandrasekharan (M.S.)
- S. Chakraborty (Ph.D.)
- N. Dandapanthula (M.S.)
- V. Dhanraj (M.S.)
- C.-H. Chu (Ph.D.)
- T. Gangadharappa (M.S.)
- K. Gopalakrishnan (M.S.)
- R. Gulhane (M.S.)
- J. Hashmi (Ph.D.)
- M. Han (M.S.)
- W. Huang (Ph.D.)
- A. Jain (Ph.D.)
- J. Jani (M.S.)
- W. Jiang (M.S.)
- J. Jose (Ph.D.)
- M. Kedia (M.S.)
- K. S. Khorassani (Ph.D.)
- S. Kini (M.S.)
- M. Koop (Ph.D.)
- P. Kousha (Ph.D.)
- K. Kulkarni (M.S.)
- R. Kumar (M.S.)
- S. Krishnamoorthy (M.S.)
- K. Kandalla (Ph.D.)
- M. Li (Ph.D.)
- P. Lai (M.S.)
- J. Liu (Ph.D.)
- M. Luo (Ph.D.)
- A. Mamidala (Ph.D.)
- G. Marsh (M.S.)
- V. Meshram (M.S.)
- A. Moody (M.S.)
- S. Naravula (Ph.D.)
- R. Noronha (Ph.D.)
- X. Ouyang (Ph.D.)
- S. Pai (M.S.)
- S. Potluri (Ph.D.)
- J. Queiser (M.S.)
- K. Raj (M.S.)
- R. Rajachandrasekar (Ph.D.)
- B. Ramesh (Ph.D.)
- D. Shankar (Ph.D.)
- G. Santhanaraman (Ph.D.)
- N. Sarkauskas (B.S. and M.S.)
- V. Sathu (M.S.)
- N. Senthil Kumar (M.S.)
- A. Singh (Ph.D.)
- J. Sridhar (M.S.)
- S. Srivastava (M.S.)
- H. Subramoni (Ph.D.)
- S. Sur (Ph.D.)
- K. K. Suresh (Ph.D.)
- K. Vaidyanathan (Ph.D.)
- A. Vishnu (Ph.D.)
- J. Wu (Ph.D.)
- W. Yu (Ph.D.)
- J. Zhang (Ph.D.)
- Q. Zhou (Ph.D.)
- N. Chmura (B.S.)

## Past Faculty

- H. Subramoni

## Past Senior Research Associate

- J. Hashmi

## Past Programmers

- A. Reifsteck
- D. Bureddy
- J. Perkins
- B. Seeds
- A. Gupta
- N. Pavuk

## Past Research Specialist

- M. Arnold
- J. Smith

## Past Post-Docs

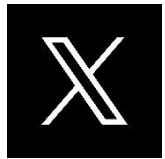
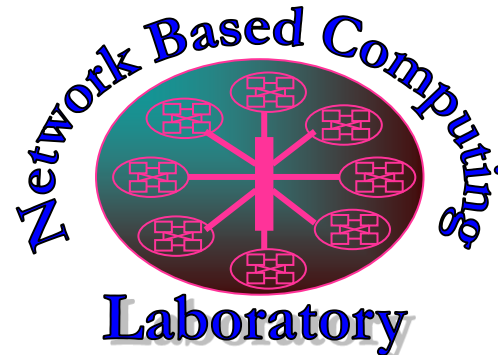
- D. Banerjee
- X. Besson
- M. S. Ghazimirsaeed
- H.-W. Jin
- J. Lin
- M. Luo
- E. Mancini
- K. Manian
- S. Marcarelli
- A. Ruhela
- J. Vienne
- H. Wang

# Interested in more? Check out our Tutorials at IEEE HOTI 2025!

- “High-Performance and Smart Network Technologies for HPC and AI” -- <https://hoti.org/tutorials-smart-network.html>
  - Friday, 22 August 2025, 13:00 – 16:30 Pacific Time (16:00 – 19:30 Eastern Time)
  - Discussion on Networking Technologies across IB, RoCE, Slingshot, EFA, NVLink/Switch, AMD InfinityFabric/xGMI, DPUs/SmartNICs
  - Case Studies and Results using MVAPICH/MVAPICH-Plus and OMB Features shown here!!
- “Principles and Practice of Scalable and Distributed Deep Neural Networks Training and Inference” -- <https://hoti.org/tutorials-dl-training.html>
  - Friday, 22 August 2025, 8:30 – 12:00 Pacific Time (11:30 – 15:00 Eastern Time)
  - Discussion on Parallelizing Training and Inference for Deep Neural Networks and Large Language Models on heterogeneous CPU/GPU/DPU architectures on modern HPC clusters
- HOTI Registration is free!! <https://hoti.org/register.html>

# Thank You!

[shineman.5@osu.edu](mailto:shineman.5@osu.edu), [michalowicz.2@osu.edu](mailto:michalowicz.2@osu.edu)



*Follow us on*

<https://x.com/mvapich>

Network-Based Computing Laboratory

<http://nowlab.cse.ohio-state.edu/>



The High-Performance MPI/PGAS Project

<http://mvapich.cse.ohio-state.edu/>



The High-Performance Big Data Project

<http://hibd.cse.ohio-state.edu/>



The High-Performance Deep Learning Project

<http://hidl.cse.ohio-state.edu/>