

Jahanzeb Maqbool Hashmi

<https://jahanzeb-hashmi.github.io>

SUMMARY

15+ years of experience in performance engineering for large-scale distributed systems, with deep expertise in GPU/accelerator co-design, collective communication libraries (NCCL, MPI), and workload-driven performance modeling for next-generation HPC/AI hardware.

RELEVANT TECHNICAL HIGHLIGHTS

- **Performance modeling and projections:** Model speed-of-light compute and communication performance on NVIDIA systems by establishing roofline targets, identifying scaling limiters, and innovating architectural improvements to increase performance, power efficiency, and TCO.
- **HW/SW Co-design:** Influence next-generation of NVIDIA GPUs/CPU/networks by providing comparative analysis of various design trade-offs across NVIDIA and competitive platforms.
- **Performance engineering:** Extensive use of NVIDIA internal and public tools to analyze performance bottlenecks of key HPC/AI applications from chip-level microarchitecture up through full datacenter-scale systems.
- **GPU-aware collective communication:** Led design of high-performance GPU-aware MPI library (MVAPICH2-GDR) for NVIDIA and AMD multi-GPU systems with focus on zero-copy point-to-point and collective communication, architecture-aware hierarchical collectives, and RDMA communication.

PROFESSIONAL EXPERIENCE

- **Senior High Performance Compute Architect, NVIDIA Corporation, USA.** March 2021 – Present
 - Model speed-of-light system performance (compute and communication) to establish GPU utilization roofline targets for next-generation accelerator hardware.
 - Identify scaling limiters and co-design applications and communication libraries to achieve maximum architecture-allowed performance and efficiency.
 - Lead the HPC performance projections working group evaluating current and next-gen accelerator architectures for key HPC and AI/LLM workloads — insights directly inform datacenter product architecture and roadmap decisions.
 - Lead the HPC competitive analysis working group, projecting application and system performance — from kernel/chip-level analysis to full-scale datacenter architecture — on internal and external accelerator products.
- **Senior Research Associate, The Ohio State University, USA.** June 2020 – March 2021
 - Designed and developed a high-performance MPI library for next-generation HPC/cloud systems spanning multi-core CPUs and GPUs (NVIDIA, AMD).
 - Led design and development of a generalized hierarchical MPI collective communications framework optimized for modern CPU+GPU node architectures.
 - Led design and development of MVAPICH2-GDR, a GPU-aware MPI library for NVIDIA/AMD multi-GPU systems, from architecture through production deployment.
 - Mentored Ph.D. and M.S. students on communication runtime design for multi-core CPUs, many-core GPUs, and high-speed interconnects.
- **Graduate Research Associate, The Ohio State University, USA.** August 2015 – May 2020
 - Designed an adaptive, topology-aware algorithm mapping MPI processes to hardware cores based on communication patterns of AI and HPC workloads.
 - Collaborated on efficient parallelization of large-scale distributed DNN training (data- and model-parallel) across CPU and GPU systems.
 - Designed and developed a zero-copy XPMEM-based inter-process communication layer exploiting shared address space for manycore/multi-GPU architectures.
 - Designed a novel data-layout caching algorithm to mitigate performance costs of MPI derived-datatype layout translation (Best Paper Award nominee, IPDPS '19).
 - Worked on PGAS libraries (OpenSHMEM, UPC++) and task-based programming models (Kokkos) with an MPI communication backend.

TECHNICAL SKILLS

- **GPU & Accelerator Systems:** CUDA/HIP, NCCL collective communication, GPU-aware MPI/SHMEM, NVIDIA Nsight Suite, multi-GPU/multi-node performance profiling
- **Performance Engineering:** Roofline modeling, FLOPs/bandwidth/interconnect throughput analysis, performance engineering
- **System Design:** Hardware/software co-design for HPC and AI, end-to-end system design and optimization for perf, power, and TCO
- **Programming Languages:** C, C++, Python, Java

EDUCATION

Ph.D. in Computer Science and Engineering (HPC), The Ohio State University, Columbus, Ohio, USA 2015 – 2020
M.S. in Computer Engineering, Ajou University, Suwon, South Korea 2012 – 2014
B.S. Information Technology, National University of Science and Technology, Islamabad, Pakistan 2007 – 2011

SELECT PUBLICATIONS

For complete list of publications, please refer to my [Google Scholar](#).

1. B. Ramesh, **J. Hashmi**, S. Xu, A. Shafi, M. Ghazimirsaeed, M. Bayatpour, H. Subramoni, and D. K. Panda. “Towards Architecture-aware Hierarchical Communication Trees on Modern HPC Systems”, in proceeding of *28th IEEE International Conference on High Performance Computing, Data, Analytics and Data Science*, Dec. 2021. [[Best Paper Finalist](#)]
2. **J. Hashmi**, C. Chu, S. Chakraborty, M. Bayatpour, H. Subramoni, and DK Panda. “FALCON-X: Zero-copy MPI Derived Datatype Processing on Modern CPU and GPU Architectures“, *Journal of Parallel and Distributed Computing (JPDC)*, Volume 144, October 2020, Pages 1-13, doi.org/10.1016/j.jpdc.2020.05.008
3. **J. Hashmi**, S. Xu, B. Ramesh, M. Bayatpour, H. Subramoni, and D. K. Panda. “Machine-agnostic and Communication-aware Designs for MPI on Emerging Architectures”, in proceeding of *34th IEEE International Parallel and Distributed Processing Symposium (IPDPS '20)* , May 2020
4. **J. Hashmi**, S. Chakraborty, M. Bayatpour, H. Subramoni, D. K. Panda. “FALCON: Efficient Designs for Zero-copy MPI Datatype Processing on Emerging Architectures”, in proceeding of *33rd IEEE International Parallel and Distributed Processing Symposium (IPDPS '19)* , May 2019 [[Best Paper Finalist](#)]
5. **J. Hashmi**, S. Chakraborty, M. Bayatpour, H. Subramoni, D. K. Panda. “Designing Efficient Shared Address Space Reduction Collectives for Multi-/Many-cores”, in proceeding of *32nd IEEE International Parallel and Distributed Processing Symposium (IPDPS '18)* , May 2018
6. **J. Hashmi**, S. Chakraborty, M. Bayatpour, H. Subramoni, D. K. Panda. “Design and Characterization of Shared Address Space MPI Collectives on Modern Architectures”, in proceeding of *The 19th Annual IEEE/ACM International Symposium in Cluster, Cloud, and Grid Computing (CCGRID '19)* , May 2019
7. S. Chakraborty, M. Bayatpour, **J. Hashmi**, H. Subramoni, D. K. Panda. “Cooperative Rendezvous Protocols for Improved Performance and Overlap”, in proceeding of *IEEE/ACM International Conference for High Performance Computing, Networking, Storage, and Analysis (SC '18)* , Nov 2018 [[Best Paper Finalist](#)]
8. M. Bayatpour, **J. Hashmi**, S. Chakraborty, H. Subramoni, P. Kousha, D. K. Panda. “SALaR: Scalable and Adaptive Designs for Large Message Reduction Collectives”, in proceeding of *IEEE International Conference on Cluster Computing (CLUSTER 2018)* , Sep 2018 [[Best Paper Award](#)]
9. A. Awan, K. Hamidouche, **J. Hashmi**, and D. K. Panda. “S-Caffe: Co-designing MPI Runtimes and Caffe for Scalable Deep Learning on Modern GPU Clusters”, in proceeding of *22nd ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP '17)* , February 2017