

Jahanzeb Maqbool Hashmi

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

OBJECTIVE

Over 11 years of experience in the area of High Performance Computing (HPC) with strong focus on application-aware hardware software co-design, end-to-end system architecture, performance engineering and benchmarking, research and development of large-scale distributed communication software for HPC and AI workloads.

EDUCATION

Ph.D. in Computer Science and Engineering, [The Ohio State University](#), Columbus, Ohio, USA 2015 – 2020

Thesis: Designing High Performance Shared-Address-Space and Adaptive Communication Middlewares for Next-Generation HPC Systems

Advisor: [Dhabaleswar K. \(DK\) Panda](#)

Skills Learned: R&D of communication libraries (MPI, PGAS) on parallel architectures (CPUs, GPUs, Networks), architecture-aware algorithm design, performance engineering and optimizations on large-scale systems, optimal workload mapping on heterogeneous architectures

M.S. in Computer Engineering, [Ajou University](#), Suwon, South Korea 2012 – 2014

Thesis: Exploring Performance and Efficiency of Multicore ARM Cluster for High Performance Computing

Skills Learned: Performance benchmarking, Application evaluations on parallel architectures

B.S. Information Technology, [National University of Science and Technology](#), Islamabad, Pakistan 2007 – 2011

Thesis: Implementation and Evaluation of Scientific Simulations on HPC Architectures

Skills Learned: Parallel Programming (MPI, Open MP, CUDA), Porting serial codes to Parallel

PROFESSIONAL EXPERIENCE

• **Senior High Performance Compute Architect**, [NVIDIA Corporation](#), USA. March 2021 – Present

- Working with hardware, software, and product teams to design NVIDIA's next-generation HPC and AI systems as part of the GPU Architecture group
- Identifying scaling limiters of important workloads, addressing limiters through hardware software co-design, and proposing system designs for future-generation systems
- Leading HPC performance projections working group to evaluate our current and next-generation architectures for key HPC and AI workloads. These insights influence various architectural design choices of our datacenter products and help management define future roadmap
- Leading HPC competitive analysis working group to analyze – from chip-level analysis to whole datacenter architecture – and project application performance on competitive products, and provide data to help our product teams in product positioning

• **Senior Research Associate**, [The Ohio State University](#), USA. June 2020 – March 2021

- Worked on the design and development of high-performance MPI library for next-generation HPC and Cloud systems with multi-core CPUs (AMD Rome, Intel Xeon, IBM POWER9, ARM A64FX) and many-core GPUs (NVIDIA, AMD)
- Led the design and development of a generalized hierarchical MPI collective communications framework for modern CPU and GPU systems
- Led the design and development of MVAPICH2-GDR, a high-performance GPU-aware MPI library, for NVIDIA and AMD based multi-GPU systems
- Mentored Ph.D. and Masters students' thesis on diverse research areas such as communication over multi-core CPUs and many-core GPUs, novel designs for message-passing communication over high-speed interconnects.

• **Graduate Research Associate**, [The Ohio State University](#), USA. August 2016 – May 2020

- Designed an adaptive and topology aware algorithm for mapping of MPI processes to hardware cores by capturing the communication-patterns of AI and HPC applications
- Worked collaboratively on efficient parallelization of large-scale distributed DNN training (data and model parallel) on CPU and GPU systems
- Designed and developed a truly zero-copy XPMEM-based inter-process (IPC) communication with shared address space principle targeting manycore architectures
- Designed a novel algorithm to cache data layouts to help mitigate the performance costs of layout translation of MPI derived datatypes. This work was nominated for the Best Paper Award at IPDPS'19

- Worked on PGAS libraries e.g., OpenSHMEM, UPC++ and task-based programming models e.g., Kokkos with MPI backend
- **Department Fellow**, [The Ohio State University](#), USA. Aug 2015 - July 2016
 - I was awarded the prestigious Department Fellowship by Ohio State University. Under the fellowship, I made research contributions to *PGAS over MPI* project at Network Based Computing Laboratory.
 - I designed, developed, and analyzed an MPI based communication conduit for UPC++ asynchronous PGAS programming model.
 - I implemented an MPI+UPC++ hybrid version of popular LULESH application and demonstrated the performance benefits against pure MPI or UPC++ implementations.

TECHNICAL SKILLS

- Hardware Software co-design for HPC and AI
- End-to-end system design and optimization for TCO
- Performance Engineering and benchmarking
- Parallel Programming Systems – MPI, OpenMP, CUDA, HIP, PGAS
- Deep Learning – Tensorflow, PyTorch, Horovod
- Programming Languages – C, C++, Java, Bash, Python
- Tools – Git, LaTeX, NVIDIA Nsight Suite, GDB, PerfAPI, mpiP, Valgrind

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