

Shaoyu Wang

🏠 Luoyu Road 1037, Wuhan, Hubei 430074, China

✉️ wsyy0619@gmail.com • 📞 (+86) 136-9813-2563 • 🌐 shawllewyw

EDUCATION	Huazhong University of Science and Technology (HUST) ▪ B.E. in Computer Science and Technology, GPA: 3.72 / 4.0	Wuhan, China Sep 2020 – Jun 2024 (Expected)
RESEARCH INTEREST	My current focus is improving AI system performance across the whole computing stack from large scale training and serving to hardware/software co-design. I am also enthusiastic in computer architecture , specifically general-purpose GPU (GPGPU) architecture and Domain-Specific Accelerators (DSA).	
PUBLICATIONS	Visual Exploratory Analysis for Designing Large-Scale Network-on-Chip Architectures: A Domain Expert-Led Design Study Shaoyu Wang*, Hang Yan*, Katherine E. Isaacs, Yifan Sun <i>IEEE Transactions on Visualization and Computer Graphics (TVCG)</i>	
RESEARCH EXPERIENCE	HPC-AI Lab at National University of Singapore (NUS) Research Intern, Advised by Prof. Yang You Speculation-based sequence batching for Large Language Model serving (in progress for ICML2024) ▪ Created an asynchronous serving system that incorporates an optimized dual-stage pipeline of prediction and serving processes, improving system throughput. ▪ Developed a dynamic programming based sequence batching algorithm, utilizing pre-serving profiling to improve efficiency and performance of batch computing. ▪ Established a mechanism to detect and re-process sequences with prediction failure, ensuring reliable and consistent serving outcomes.	Jul 2023 – Present
	Scalable Architecture Lab at College of William & Mary Research Intern, Advised by Prof. Yifan Sun and Prof. Katherine E. Issacs Wafer-Scale GPU architecture performance analysis and optimization ▪ Proposed a load-balanced Network-on-Chip (NoC) routing policy, improving NoC bandwidth utilization. ▪ Introduced a metrics collector within the MGPUSim simulator to gather network traffic data. Visualization of mesh Network-on-Chip (NoC) traffic ▪ Integrated with an architecture-level profiling tool to gain insights of performance bottlenecks. ▪ Developed a systematic visualization solution combining temporal and spatial views for effective hotspot identification in the network.	Apr 2022 – Mar 2023
WORK EXPERIENCE	NVIDIA Compute Architect Intern at HPC Architecture Performance Team Enhanced the internal-used simulation tool of GPGPU architecture prototyping and performance analysis ▪ Optimized the simulation process, significantly reducing the time required from more than a day to just several minutes for workloads involving over 400 traces. ▪ Expanded the simulator's capabilities by adding support for the cutting-edge NVIDIA Hopper architecture, enabling performance analysis of latest GPU cards ▪ Introduced new evaluation models for hardware load balance and Streaming Multiprocessors (SM) utilization, boosting the precision and reliability of simulation.	Shanghai, China Mar 2023 – Jul 2023
	ByteDance (TikTok) Software Engineering Intern at Cloud Computing Infrastructure Team Enhanced KubeBrain, a high-performance distributed data management system for Kubernetes ▪ Developed an backup and recovery module in KubeBrain, bolstering its reliability and resilience. ▪ Integrated a hot cache mechanism to efficiently handle time-consuming queries, reducing the request latency by 12x.	Hangzhou, China Jun 2022 – Sep 2022

OPEN SOURCE CONTRIBUTION	OIWiki Contribute to OIWiki, the online wiki for competitive programming competitions <ul style="list-style-type: none"> Developed a plugin to support the compilation of Markdown tab syntax, enabling the OIWiki community to efficiently process documents featuring content tabs Enhanced the export tool-chain for converting into LaTeX and Typst formats, extending export capabilities of OIWiki materials 	
	OpenEuler Contribute to iSulad, the open source container manager written in C++ <ul style="list-style-type: none"> Extended the communication protocol of iSulad, facilitating seamless interaction between the command-line interface (CLI) and the backend server Implemented progress monitoring for container image pulling, enhancing its ability to track and display real-time progress 	
STUDENT CLUSTER COMPETITIONS	ISC Student Cluster Competiton 2023, virtual track Build up and lead a student group of 6, advised by Prof. Xuanhua Shi and Prof. Yao Wan <ul style="list-style-type: none"> Deployed and profiled 3 scientific computing application (Quantam Espresso, FluTAS, POT3D) using several computing nodes in supercomputer centers. Managed to build CPU-GPU hybrid programs and orchestrated them over 4 computing nodes using MPI and InfiniBand (IB). 	Finalist
	ASC Student Supercomputer Challenge 2022 Being one of a student group of 5, advised by Prof. Xuanhua Shi <ul style="list-style-type: none"> Pre-trained a 4.7B language model Yuan-1.0 on the given 100GB dataset using Megatron, orchestrating 3D parallelism for distributed training over 8 nodes Optimized DeepMD training process, achieved 4x speedup on CPU and 2x speedup on GPU 	Second Prize
	3rd “AMD Cup” International Parallel Computing Competition <ul style="list-style-type: none"> Accelerated a high-dimension point cloud feature extraction algorithm, achieving 1200x speedup Accelerated the SOTA graph spectral sparsification algorithm feGRASS, achieving 40x speedup 	Silver Medal
	9th “Intel Cup” Parallel Application Challenge <ul style="list-style-type: none"> Optimized an application for genome sequence analysis, achieving 13x performance improvement 	Bronze Medal
	Programming Languages: C/C++, Python (pytorch), Golang, JavaScript/TypeScript, Verilog Computer Architecture: Software Simulation, High level synthesis (HLS) High Performance Computing: OpenMP, MPI, SSE/AVX, CUDA, Code profiling Large Language Model: Parameter efficient fine-tune (PEFT), Data/Pipeline/Model Parallelism	
HONORS & AWARDS	2023 ISC Student Cluster Competiton	May 2023
	<ul style="list-style-type: none"> Finalist of the online track 	
	2022 ASC Student Supercomputer Challenge	April 2022
	<ul style="list-style-type: none"> Second Prize in the preliminary contest 	
	The 3rd “AMD Cup” International Parallel Computing Competition Silver Medal, China	Nov 2022
	<ul style="list-style-type: none"> Rank 3rd in the final, competed with about 120 teams in China 	
	The 9th “Intel Cup” Parallel Application Challenge Bronze Medal, China	Oct 2021
	<ul style="list-style-type: none"> Rank 5th in the final, competed with about 120 teams in China 	
	Huawei Intelligent Base Scholarship	May 2022
	<ul style="list-style-type: none"> Top 0.5%, for outstanding innovation and strong practical skills of CS students 	
	Scientific and Technology Innovation Scholarship, HUST	Dec 2021
	<ul style="list-style-type: none"> Top 3.0%, award for scientific research and innovation of undergraduates 	