

# Xiayang Jin

+1 (734) 274-0486 | Ann Arbor, MI | [ausummer@umich.edu](mailto:ausummer@umich.edu)

## EDUCATION

<b>University of Michigan, Ann Arbor</b> <i>Bachelor's of Science, Computer Science</i>	Aug 2024 — Now <i>Ann Arbor, MI</i>
• Cumulative GPA: 4.0/4.0	
<b>Shanghai Jiaotong University</b> <i>Bachelor's of Science, Mechanical Engineering</i>	Sep 2022 — Now <i>Shanghai, China</i>

## EXPERIENCES

<b>Research Intern</b> OrderLab, University of Michigan	May 2025 — Present <i>Ann Arbor, MI</i>
• Conduct research at the System–ML boundary on reliability and performance of machine learning systems	
• Develop system support for failure modeling, tolerance and recovery strategies of ML workloads that preserve correctness while minimizing overhead	
• Evaluate techniques on LLM/RLHF workloads across multi-GPU setups; analyze trade-offs among time-to-recovery, throughput, and data integrity	
<b>Research Intern</b> Shanghai Jiaotong University	May 2024 — Aug 2024 <i>Shanghai, China</i>
• Research in MINLP optimizations and their applications in Computer Networking	
• Designed an Machine Learning approach to boost traditional MINLP solvers, and applied it to cutting-edge networking algorithms	
<b>Teaching Assistant</b> Shanghai Jiaotong University	May 2024 — Aug 2024 <i>Shanghai, China</i>
• Assisted in presenting a course about game development and software engineering	
• Developed auto grading and feedback programs for student projects	
• Conducted weekly lab and office hours for over 50 students	
<b>Instructional Aide</b> University of Michigan	Expected Jan 2026 — May 2026 <i>Ann Arbor, MI</i>
• Will assist in presenting a course about Introduction to Operating Systems	

## PROJECTS

<b>PhoenixML</b>	May 2025 — Present
• Developed a stage-level automatic recovery system for distributed machine learning training jobs with transient CUDA faults without checkpoint/restore	
• Designed an runtime interposition layer to transparently trace CUDA alloc/launch calls and maintain a VA to tensor registry	
• Implemented CUDA Virtual Memory Management with a proxy context to enable zero-copy remapping across context respawns	
<b>CudaProxy</b>	Aug 2025 — Present
• Developed a CUDA runtime proxy for machine learning inference to use CUDA Graphs and Persistent Kernels for accelerated performance without brittle fusion	
• Implemented Bucketing, Automatic Padding, and Static Memory Pooling to handle dynamic workloads with minimal overhead	
• Implemented automatic routing for kernel launches based on workload characteristics	
<b>OS Kernel Implementation</b>	Jan 2025 — May 2025

- Developed a C++ thread library and virtual memory pager on UNIX-like systems
- Implemented synchronization primitives like mutexes and condition variables using UNIX context management techniques.
- Managing Thread life cycle, Scheduler, Context Switching, CPU Booting, Process creation, Process forking, Page fault handling, and Process destruction
- Implemented optimizations including copy-on-write and page sharing

### Dynamic Typed Compiler

May 2024 — Aug 2024

- Used Rust to develop a compiler for a simple language supporting dynamic typing and heap allocation on x86-64 architecture
- Implemented frontend checking, middleend Single Static Assignment (SSA) forming, and backend code generation with System V ABI
- Implemented optimizations including register allocation and assertion removal

### Hyprtile

May 2025 — Present

- Developed a high-performance tiling window management mechanism for Linux using C++ under Wayland protocol, targeting efficient multitasking across multiple projects
- Features include workspace management, windows overview, dynamic window tiling, and multi-monitor support

### Deep Learning Approach to MINLP Problems in Wireless Networks

May 2024 — Aug 2024

- Designed and implemented an imitation learning approach to accelerate Branch and Bound algorithm for solving Mixed-Integer Nonlinear Programming optimization problems
- Applied the approach to Inteferece Graph Estimation in full-duplex millimeter-wave backhaul networks

## SKILLS

---

- **Programming Languages:** C/C++, Python, CUDA, Rust, HTML, CSS, Java, SQL, Bash, Elm
- **Technologies:** Linux Kernel, CUDA, DeepSpeed, Pytorch, Tensorflow, Git, UNIX, Docker