

EDUCATION

- **University of Edinburgh** Edinburgh, UK
PhD in Computer Science *August 2023 - Present*
Research Field: AI-Systems, Distributed ML, ML-oriented Architecture, Serverless
- **University of Science and Technology of China** Hefei, China
Bachelor of Computer Science; The school of Gifted Young; *July 2017 - June 2021*
Courses: Operating Systems, Artificial Intelligence, Principles of Compiler, Computer Architecture, High Performance Architecture

PUBLICATIONS

- 1 Yao Fu, Leyang Xue, **Yeqi Huang**, et al. "ServerlessLLM: Locality-Enhanced Serverless Inference for Large Language Models." In *18th USENIX Symposium on Operating Systems Design and Implementation (OSDI 24)*, 2024.
- 2 Yao Fu, Yinsicheng Jiang, **Yeqi Huang**, et al. "MoE-CAP: Cost-Accuracy-Performance Benchmarking for Mixture-of-Experts Systems." *arXiv preprint*, 2024.
- 3 Congjie He, **Yeqi Huang**, Pei Mu, et al. "LLMoC: Large Language Model Inference at Wafer Scale." *Under Review*, Submitted to OSDI 2025.
- 4 Xiao-Long Chen, Lin-Feng Wang, **Yeqi Huang**, et al. "Symplectic structure-preserving particle-in-cell whole-volume simulation of tokamak plasmas." In *SC21: International Conference for High Performance Computing, Networking, Storage and Analysis*, 2021.

PROJECTS

- **AutoReader Project (<https://autoreader.ed-aisys.com/>):**
 - Built a vector database for the latest ArXiv papers for daily semantic search and subscription.(Full-stack AI application)
 - *Implemented high-performance retrieval algorithm on GPU.*
 - Achieving **60× faster indexing than NVIDIA cuVS** library and **13× faster retrieval than Milvus**.
 - Extended system to support bioRxiv and PubMed papers with specialized biology-focused features.
 - Biology patched version is demonstrating in Dartmouth university, will submit to Nature this year.
- **Vector Search on Cerebras (Plan to submit to NIPS 2025):**
 - Implemented KNN/ANN-based vector search algorithm on Cerebras chip.
 - Developing Retrieval based RAG and RetrievalAttention for long context LLM recently.
 - Achieved **800× faster GEVV operations with 40× better energy efficiency**.
- **Cerebras 2D-Mesh AI Research WaferLLM(submitted to OSDI 2025):**
 - Implemented Llama by an assembly-like programming language, get **39× faster with 1.7× energy efficiency**.
 - Developed novel GEMM algorithms surpassing traditional approaches on specialized hardware.
 - Achieved **606× faster GEMV operations with 22× better energy efficiency**.
- **ServerlessLLM(published on OSDI 2024):**
 - Developed a novel serverless system reducing latency by 10 - 200x across various LLM inference workloads.
 - Implemented efficient scheduling strategies for multi-agent systems on distributed infrastructure.
- **BioLLM Research Project:**
 - Collaborated with Dartmouth and Harvard Medical School on biological literature interpretation.
 - Developed advanced chart interpretation capabilities including CoT training, ReFT, and long context with RAG.
 - Achieved 23% improvement over GPT-4o in double-blind evaluations.
- **QED Simulation on GPU:**
 - Implemented high-performance Monte-Carlo simulation achieving 1500x to 5000x speedup.
 - Received recognition from Nobel Laureate Frank Wilczek at APS conference.
 - Demonstrated expertise in CUDA programming and GPU micro-architecture optimization.
- **More Projects (<https://yeqi-huang.com/>):**
 - I have plenty of AI related projects on my personal page.

TALKS

- 1 **Yeqi Huang**. "InfiniTensor: A Tensor-Friendly, Efficient Parallel Programming Library for Accelerator-Centric Clusters." *UKSys 2024*.
- 2 **Yeqi Huang**. "Why we need a new clustering benchmark in AI retrieval?" *International Workshop on Efficient Generative AI*, 2024.

AWARDS

International Awards:

- International Supercomputing Conference Student Cluster Competition Champion - 2021
- Asian Supercomputer Conference Student Cluster Competition First Prize - 2021

National Awards:

- Best Chinese Supercomputing Application of the Year - 2022
- Huawei Pioneer Developer (4 in China) - 2021
- National Compiler Designing Competition Second Prize - 2021

SKILLS

- **AI-Related:**
 - In-depth understanding of LLMs, with experience porting high-performance inference and training frameworks to various hardware platforms, including Apple Silicon, GPUs, and Cerebras.
 - Strong understanding of vector retrieval, having developed and tested graph and vector databases on multiple platforms.
 - In-depth knowledge of Multi-Agent applications and developed efficient development components for such applications.
- **Computer Science Related:**
 - Highly skilled in **C++**, **Python**, **Rust**; familiar with Go, JavaScript, and Latex
 - Proficient with **CUDA**, **Intel ONEAPI**, **OpenMP** and related parallel and distributed programming
 - Well-versed in **LLVM**, frequently participating in LLVM Forum online discussions
 - Extensive experience working with Linux, including usage of eBPF
 - Rich experience in distributed systems and distributed machine learning frameworks
 - Strong engineering development experience, mastering various compiler-related tools and static analysis tools
- **Physics & Math:**
 - Highly skilled in Computational Fluid Dynamics and Molecular Dynamics
 - Strong knowledge in numerical methods and linear algebra
 - Proficient in Quantum Mechanics and Quantum Electrodynamics, with some knowledge of quantum computing

SPECIAL EXPERIENCE

- **Open Source Enthusiast:** Contributing to notable projects like GNOME, LAMMPS, and LLVM, I've used GitHub since 2019 to showcase my development journey and ideas.
- **Running a Science-Themed Cafe:** Leveraging software development income, I established Quantum Coffee near my school—a collaborative space encouraging students to explore and discuss scientific topics across disciplines.
- **UNICEF Charity Projects:** I have donated 10% of all personal income to children's charities since 2019.