

EDUCATION

The Chinese University of Hong Kong

Ph.D. in Computer Science and Engineering

- Advisor: Prof. Hong Xu
- Research area: Machine Learning System

Hong Kong SAR, China

2022 - 2026 (*expected*)

Northwestern Polytechnical University

B.E. in Computer Science and Technology

- GPA: 93.37/100, Rank: 1/247.

Xi'an, China

2018 - 2022

RESEARCH INTEREST

I am broadly interested in System Design for Machine Learning (Sys4ML), including the following topics:

1. **Distributed Training:** Developing and optimizing strategies for efficient, scalable training of large-scale models.
2. **Efficient Serving Systems:** Designing novel architectures and algorithms for high-performance inference and serving of large models and applications, such as LLMs and diffusion models.

PUBLICATIONS

1. **Xin Tan, Yuetao Chen, Yimin Jiang, Xing Chen, Kun Yan, Nan Duan, Yibo Zhu, Daxin Jiang, Hong Xu, DSV: Exploiting Dynamic Sparsity to Accelerate Large-Scale Video DiT Training.** *ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, 2026.
2. **Xin Tan, Yimin Jiang, Yitao Yang, Hong Xu, Towards End-to-End Optimization of LLM-based Applications with Ayo.** *ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, 2025.
3. **Xin Tan, Jiamin Li, Yitao Yang, Jingzong Li, Hong Xu, Arlo: Serving Transformer-based Language Models with Dynamic Input Lengths.** *ACM International Conference on Parallel Processing (ICPP)*, 2024.

INTERNSHIPS

System Group, StepFun | Beijing, China

2024.08 - present

- Identified attention bottlenecks in video DiTs with long video inputs and analyzed dynamic sparse attention patterns. Co-designed algorithms and systems to scale video-DiT sparse training, implementing the solution to achieve a 3.02x throughput improvement on 128 GPUs.
- Exploring next-generation data center designs for LLM post-training infrastructure (Ongoing).

Network Research Group, Microsoft Research Asia | Remote

2021.10 - 2022.03

- Developed a real-time, non-intrusive monitoring system to collect critical AI infrastructure metrics (GPU utilization, network/NVLink bandwidth), and designed a scalable analytics tool to evaluate ML workloads across large-scale clusters.
- Analyzed six months of datacenter workload data to characterize resource usage and network patterns of various AI tasks, providing actionable recommendations to optimize cloud infrastructure and software stack.

AWARDS AND HONORS

- **Student Travel Grant, ASPLOS 2025** 2025.4
- **Full Postgraduate Scholarship, The Chinese University of Hong Kong** 2022-2026
- **Outstanding Graduate, Northwestern Polytechnical University** 2022
- **National Scholarship, Ministry of Education (China)** 2020
- **National Scholarship, Ministry of Education (China)** 2019
- **Champion, International Underwater Robot Competition** 2020

SKILLS	Languages: Chinese, English Programming: Python, Pytorch, Megatron, Ray, Triton, CUDA, C++
ACADEMIC SERVICES	Reviewers: <i>IEEE Transactions on Network Science and Engineering</i> , Shadow Program Committee: <i>ACM EuroSys 2026</i> , Artifact Evaluation Committee: <i>USENIX OSDI/ATC 2025</i> , <i>ACM CoNEXT 2025</i> , <i>ACM EuroSys 2025 Spring/Fall</i> , <i>USENIX OSDI/ATC 2024</i> ,
TEACHING	Teaching Assistant: <i>CSCI 3150, Introduction to Operating Systems, CUHK. 2023 Spring</i> , <i>CSCI 1120, Introduction to Computing Using C++, CUHK. 2022 Fall</i> .