

# Jie Xin

[\[email\]](#) [\[github\]](#) [\[homepage\]](#)

## Experience

---

### NVIDIA

*Senior Compute Architect*

Shanghai, China

2022 – Present

- Researching automated end-to-end model performance optimization on next-gen GPUs (Hopper, Blackwell, Rubin)
- Developed deep learning compilers (CuTile, Triton-to-CuTile) and MLPerf benchmark optimization
- Collaborating with PyG, Megatron-LM, and OpenAI Triton teams on compiler infrastructure

## Education

---

### Huazhong University of Science and Technology

*M.S. in Computer Science*

Wuhan, China

2019 – 2022

### Huazhong University of Science and Technology

*B.S. in Computer Science*

Wuhan, China

2015 – 2019

## Projects

---

### Agent for Compiler

Designing a lightweight harness system specifically for compiler optimization.

2026 – Present

### Triton-to-CuTile

OpenAI incubation project; compiler bridge connecting Triton frontend to CuTile backend.

Enabling Triton users to target NVIDIA's optimized tile programming infrastructure.

2024 – 2025

[\[code\]](#)

### CuTile

CUDA Tile programming standard for NVIDIA GPUs.

Designed and implemented optimization passes for deep learning workloads.

2023 – Present

[\[code\]](#) [\[blog\]](#)

### NvFuser

PyTorch JIT compiler achieving 1.2x–2.5x speedup on GNN models.

Enabled gather/scatter/index\_select graph operations for sparse workloads.

2022 – 2023

[\[code\]](#) [\[blog\]](#)

### OpenFold2 Training

MLPerf HPC v3.1 #1 submission with 7.5x e2e speedup.

Optimized protein structure prediction training on large-scale GPU clusters.

2023 – 2023

[\[code\]](#) [\[paper\]](#) [\[blog\]](#)

### GPT-3 Training

MLPerf Training v4.0; contributed 1.02x training speedup on 175B model.

Optimized distributed training pipeline for large language models.

2023 – 2023

[\[code\]](#)

### PyTorch Geometric TorchScript

Added TorchScript support for PyG community, enabling JIT compilation for GNN models.

Contributed core infrastructure for production deployment of graph neural networks.

2022 – 2022

[\[code\]](#)

### SpMM on GPU – Champion, MIT/IEEE Graph Challenge

High-performance sparse matrix multiplication; up to 652x speedup over cuSPARSE.

Won 1st place in the 2021 MIT/IEEE/Amazon Graph Challenge competition.

2021 – 2021

[\[code\]](#) [\[paper\]](#)

## Current Focus

---

My research interest has fully shifted to AI agents – exploring how LLMs can autonomously optimize compilers, automate performance tuning, and transform traditional software development workflows.