

# Willis Guo

willisg@cs.cmu.edu · (878) 834-9154 · willisguo14.github.io · linkedin.com/in/willisguo · github.com/willisguo14

## Education

|  |                             |
|--|-----------------------------|
| MS Machine Learning<br>Carnegie Mellon University  | Dec. 2025<br>Pittsburgh, PA |
| BASc Machine Intelligence<br>University of Toronto | Apr. 2024<br>Toronto, CAN   |

## Experience

|   |   |
|---|---|
| Research Engineer Intern, Distributed RL<br>Scale AI  | May 2025 – Aug. 2025<br>San Francisco, CA |
| • Developed Rollout Parallel Attention: a <b>novel distributed training parallelism</b> technique for accelerating <b>online reinforcement learning (RL)</b> for LLMs. Increases GRPO training throughput by 10x.   |   |
| • Built from scratch a new <b>RL post-training library</b> with <b>GPU co-location</b> , <b>FSDP</b> training, vLLM inference, and supports <b>training multimodal models</b> . Planning to open source.  |   |
| • Researched <b>test-time scaling laws</b> for online RL. Training 7B to 32B LLMs using GRPO on math and video reasoning.   |   |
| • Implemented a fused <b>Triton kernel</b> for calculating post-training losses, decreasing memory usage 3x and enabling <b>long-context training</b> up to 96k tokens per GPU.   |   |
| Research Intern, Multimodal<br>Carnegie Mellon University, Ruslan Salakhutdinov   | Sep. 2024 – May 2025<br>Pittsburgh, PA    |
| • Created synthetic data with interleaved <b>multimodal chain-of-thought reasoning</b> traces with <b>visual tool-use</b> for supervised fine-tuning (SFT) <b>vision-language models (VLMs)</b> for video reasoning.  |   |
| • Developed an inference algorithm for video understanding with VLMs by leveraging <b>video diffusion models</b> as a <b>world model</b> .  |   |
| Software Engineer Intern<br>Amazon Web Services   | June 2024 – Aug. 2024<br>Vancouver, CAN   |
| • Built <b>ML infrastructure</b> , including data pipelines and observability tools for analyzing petabytes of AWS resource traffic data. Performed <b>hyperparameter tuning</b> and <b>feature engineering</b> , improving anomaly detection recall by 2%. |   |
| Research Intern, LLM Reasoning<br>University of Toronto, Scott Sanner   | Sep. 2023 – Apr. 2024<br>Toronto, CAN     |
| • Developed a neuro-symbolic, <b>inference-time search</b> algorithm for <b>logical reasoning with LLMs</b> that improves LLM commonsense reasoning accuracy by 13%.  |   |
| • Created an <b>LLM agent for knowledge graph question answering (KGQA)</b> with planning and active <b>retrieval augmentation (RAG)</b> , reducing hallucinations by 79% and outperforming existing KGQA methods by 8%.                                    |   |

## Publications

Many Rollouts Are All You Need: Scaling GRPO to Many Rollouts With Rollout Parallel Attention

Willis Guo, Qin Lyu. *In Preparation*

Active Perception for Efficient Inference-Time Long-Form Video Understanding in Vision-Language Models

Martin Ma, Willis Guo, Aditya Agrawal, Ankit Gupta, Paul Liang, Russ Salakhutdinov, Louis-Philippe Morency. *ICCV 2025 Workshop*

CoLoTa: A Dataset for Entity-based Commonsense Reasoning over Long-Tail Knowledge

Armin Toroghi, Willis Guo, Scott Sanner. *SIGIR 2025*

Verifiable, Debuggable, and Repairable Commonsense Logical Reasoning via LLM-based Theory Resolution

Armin Toroghi, Willis Guo, Ali Pesaranghader, Scott Sanner. *EMNLP 2024*

Right for Right Reasons: Large Language Models for Verifiable Commonsense Knowledge Graph Question Answering

Armin Toroghi, Willis Guo, Mohammad Mahdi Abdollah Pour, Scott Sanner. *EMNLP 2024*

## Projects

### Mini-LLMSys

- Implemented key LLM training techniques on top of MiniTorch: **fused CUDA kernels, data parallelism** and **pipeline parallelism**.

## Skills

Languages Python, Java, Scala, C, C++, SQL, JavaScript, TypeScript, MATLAB

Machine Learning PyTorch, CUDA, Triton, vLLM, verl, Ray, Hugging Face, Apache Spark

Other AWS, PostgreSQL, Docker, React