

# Gaokai Zhang

+1 217-974-0847 | [gaokaiz@andrew.cmu.edu](mailto:gaokaiz@andrew.cmu.edu) |  
[github.com/yourusername](https://github.com/yourusername) | [linkedin.com/in/yourlinkedin](https://linkedin.com/in/yourlinkedin)

## EDUCATION

---

<b>Carnegie Mellon University</b> <i>M.S. in Intelligent Information Systems (MIIS), Language Technologies Institute (LTI)</i>	Aug 2025 – May 2027 (Expected) <i>Pittsburgh, PA</i>
<b>University of Illinois at Urbana-Champaign</b> <i>B.S. in Computer Engineering, GPA: 3.89/4.0</i>	Aug 2021 – May 2025 <i>Urbana, IL</i>
<b>Zhejiang University</b> <i>B.Eng. in Electrical and Computer Engineering, GPA: 3.95/4.0</i>	Aug 2021 – May 2025 <i>Haining, China</i>

## EXPERIENCE

---

<b>Research Intern – System and Networking Group</b> <i>Microsoft Research Asia (MSRA)</i>	Jul 2024 – Jul 2025 <i>Beijing, China</i>
<ul style="list-style-type: none"><li>– Initiated the development of an efficient Reinforcement Learning recipe for long-context reasoning with LLMs.</li><li>– Contributed to the LongRoPE2 Research project, extending LLM context length to millions of tokens while preserving short-context capabilities; accepted as poster at ICML 2025.</li><li>– Delivered context-extended, downstream-ready LLMs to internal teams, including Microsoft Asia-Pacific R&amp;D.</li><li>– Designed a scalable pipeline for curating large-scale, high-quality supervised fine-tuning (SFT) datasets.</li></ul>	
<b>Research Intern – LLM Systems &amp; Cloud Optimization</b> <i>University of Illinois Urbana-Champaign</i>	Nov 2024 – Present <i>Urbana, IL</i>
<ul style="list-style-type: none"><li>– Designed and evaluated cost-efficient LLM training/inference strategies across heterogeneous accelerators (A100, H100, TPU) using Megatron-LM on CloudLab.</li><li>– Contributed to an automated planner for optimal parallelism and deployment configurations under dynamic SLOs; co-authoring a system paper.</li></ul>	
<b>Research Intern – LLM Safety &amp; Robustness</b> <i>University of Illinois Urbana-Champaign</i>	Mar 2024 – Oct 2024 <i>Urbana, IL</i>
<ul style="list-style-type: none"><li>– Quantified LLM robustness to stochastic attacks (word- and character-level augmentations) using SORRY-Bench, with confidence bounds based on Hoeffding and Clopper-Pearson methods.</li><li>– Co-authored a paper under review at <i>Transactions on Machine Learning Research (TMLR)</i>.</li></ul>	

## TECHNICAL SKILLS

---

**Languages:** Python, C/C++, SQL, x86 Assembly  
**Frameworks & Libraries:** PyTorch, Hugging Face Transformers, Megatron-LM  
**Tools:** Git, Slurm, MySQL, QEMU, CloudLab, Docker  
**Areas:** LLMs, NLP, Distributed Systems, Long-Context Learning, Reinforcement Learning