

# SHAO Siyang

siyang.shao@outlook.com | +86-15021988618 | [github.com/SiyangShao](https://github.com/SiyangShao)

## Education

<b>Georgia Institute of Technology</b> PhD student in Computer Science • Research Area: Large Language Model Inference, Distributed System	Aug 2025 – Current Atlanta, Georgia, USA
<b>Nanyang Technological University</b> Bachelor of Engineering (Computer Engineering) • Honours (Highest Distinction); GPA: 4.63 / 5.0	Aug 2021 – Jun 2025 Singapore, Singapore

## Research Experience

<b>Cluster Level Scheduling for LLM Inference</b> Supervised by Prof. Dmitrii Ustiugov • Investigated cluster-level scheduling for large language model inference in serverless systems • Explored optimal scaling policies and mechanisms for serverless LLM environments	Mar 2024 – May 2025
---	---------------------

## Open Source Projects

<b>ServerlessLLM</b> Developer • Supported ROCm for <code>sllm-store</code> , the internal library of ServerlessLLM which provides high-performance model loading • Integrated vLLM backend, enabling ServerlessLLM project to perform inference through vLLM • Explored methods to enable vLLM backend to benefit from high-performance model loading via <code>sllm-store</code> • Maintained the controller of the ServerlessLLM project, which manages the lifecycle of the inference backends	<a href="https://github.com/ServerlessLLM/ServerlessLLM">https://github.com/ServerlessLLM/ServerlessLLM</a> Jun 2024 – May 2025
---	--

## Work Experience

<b>Jane Street Asia Limited</b> Software Engineer Intern • Designed a library to support callee version conversion for JSON-RPC, enabling seamless alignment between JSON-RPC and async-RPC. This allows both RPC types to be registered using a unified declaration and automates version conversion, simplifying integration and maintainance. • Designed a database prototype that incrementally synchronizes an internal database (defined by a custom DSL) with a new SQL database. This approach simplifies usage for users and achieves approximately 5x faster query performance.	Hong Kong SAR May 2025 – Jul 2025
<b>TikTok Pte. Ltd.</b> Backend Engineer Intern (Video Infrastructure) • Co-Designed and implemented metrics metadata discover and manage system, bridged the gap between development teams and SRE teams concerning the monitoring metrics • Implemented persistent global SLI monitor and manage system, contributing to improvements in full-link stability	Singapore Jan 2024 – May 2024

## Co-Curricular Activities

<b>NTU ICPC Team</b> Team Member • Represented the school in ICPC (International Collegiate Programming Contest) and solved complex algorithm problems	<a href="https://icpc.global/ICPCID/B15T259WIX3C">https://icpc.global/ICPCID/B15T259WIX3C</a> Dec 2021 – Mar 2025
--	--

## Awards

• 2022 ICPC Asia Manila Regional Ranked 2	Dec 2022
• 2023 ICPC Asia Jakarta Regional Ranked 13	Dec 2023
• 2024 ICPC Asia Pacific Championship Ranked 22	Mar 2024
• 2025 ICPC Asia Jakarta Regional Ranked 11	Dec 2024
• 2025 ICPC Asia Pacific Championship Ranked 24	Mar 2025
• Dean's List in Academic Year 2022-23 (Top 5% of cohort)	Aug 2023
• NTU President Research Scholar in Academic Year 2023-24	Aug 2024