

Education

Nanyang Technological University	Aug 2021 – Jun 2025
Bachelor of Engineering (Computer Engineering)	Singapore
<ul style="list-style-type: none">Expected: Honours (Highest Distinction); GPA: 4.60 / 5.0Relevant Modules: Operating Systems (A+), Computer Network (A+), Computer Architecture and Organisation (A+), Advanced Computer Architecture (A)	

Skills Summary

- Languages:** Golang, C++, Python, CUDA, ROCm
- Tools:** Docker, Kubernetes, Knative, Kafka, Clickhouse, Grpc, Ray

Work Experience

TikTok Pte. Ltd.	Singapore
Backend Engineer Intern (Video Infrastructure)	Jan 2024 – May 2024
<ul style="list-style-type: none">Co-Designed and implemented metrics metadata discover and manage system, bridged the gap between development teams and SRE teams concerning the monitoring metricsImplemented persistent global SLI monitor and manage system, monitored and managed the compliance of SLI metrics across all global regions, contributing to improvements in full-link stability	

Open Source Projects

ServerlessLLM	https://github.com/ServerlessLLM/ServerlessLLM
Core Contributor	Jun 2024 – Current
<ul style="list-style-type: none">Supported ROCm for <code>sllm-store</code>, the internal library of ServerlessLLM which provides high-performance model loadingIntegrated vLLM backend, enabling ServerlessLLM project to perform inference through vLLMExplored 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	

Co-Curricular Activities

NTU ICPC Team	https://icpc.global/ICPCID/B15T259WIX3C
Team Member / Captain	Dec 2021 – Mar 2024
<ul style="list-style-type: none">Represented the school in ICPC (International Collegiate Programming Contest) and solved complex algorithm problems	

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
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

Research Experience

LLM Inference in Serverless Systems	Mar 2024 – Current
Supervised by Dmitrii Ustiugov	
<ul style="list-style-type: none">Investigated cluster-level scheduling for large language model inference in serverless systemsExplored optimal scaling policies and mechanisms for serverless LLM environmentsUtilized GPU memory usage for a memory-centric scheduling LLM inference system	
MIG-based GPU Partitioning and Performance Analysis	Jun 2023 – Nov 2023
Supervised by Dmitrii Ustiugov	
<ul style="list-style-type: none">Explored the use of MIG (Multi-Instance GPU) to physically partition a single GPUAnalyzed memory and PCIe bandwidth utilization across multiple MIG instances	