

# TANG KEN YI

Tel: +65 86618162 | Email: [tangkenyi2001@gmail.com](mailto:tangkenyi2001@gmail.com) | LinkedIn: <https://www.linkedin.com/in/tang-ken-yi-05629421a/> | Personal Portfolio: <https://tangkenyi2001.github.io/>

## EDUCATION

---

<b>NANYANG TECHNOLOGICAL UNIVERSITY</b> <b>Bachelor of Engineering (Computer Engineering)</b> Expected Graduation Date: December 2026	Aug 2022 – Present
<ul style="list-style-type: none"><li>Second Class (Upper) Honours/ Honours (Distinction), Current CGPA: 4.43/5.00</li><li>Dean's List AY 24/25</li><li>Relevant Coursework: Data Structures and Algorithms, Operating Systems, Computer Networks, Object Oriented Programming, Machine Learning, Neural Networks, Natural Language Processing</li></ul>	

## EXPERIENCE

---

<b>Apple</b> <b>Software Engineer Intern</b>	Feb 2026 – August 2026
<ul style="list-style-type: none"><li>Working on the Apple Online Store team to design and engineer AI-powered data pipelines processing 80+ GB of production traffic daily.</li></ul>	
<b>PayPal</b> <b>Software Engineer Intern</b>	July 2025 – December 2025

- Engineered an agentic framework that ensures deterministic code quality by orchestrating Claude Code within autonomous TDD feedback loops, driving the model to iteratively refine outputs until validation passes (20–30% faster development).
- Developed, containerized, and deployed remote MCP (Model Context Protocol) servers on AWS using Docker and Kubernetes.
- Led exploration and deployed the team's first OAuth 2.0 authentication system with SSO for MCP servers.
- Presented solutions to 30+ stakeholders across multiple time zones and authored internal technical blogs to increasing adoption and usage of newly developed tools across engineering teams.

<b>SAP</b> <b>Software Engineer Intern</b>	Jan 2025 – July 2025
<ul style="list-style-type: none"><li>Collaborated with cross-functional teams to design and release a full-stack, end-to-end internal ChatGPT platform, focusing on backend services.</li><li>Engineered key backend features in Node.js, including file upload integration with Amazon S3, persistent chat history storage, and document-based conversational querying.</li><li>Led and coordinated the development of a PowerPoint generation feature leveraging LLMs to dynamically create presentation content; integrated into the internal productivity platform.</li><li>Integrated robust observability by logging metrics and application data to OpenSearch, enabling stakeholders to monitor and evaluate usage and performance trends.</li><li>Deployed and optimized Large Language Models (LLMs) for production utilizing Docker and vLLM, reducing inference costs and improving response latency through efficient model serving.</li></ul>	

## PROJECTS

---

<b>Reducing Cold-Start Latency in LLM Inference (Final-Year Project)</b>	May 2025 – Present
<ul style="list-style-type: none"><li>Forked and extended the vLLM inference engine to prototype a disaggregated inference architecture separating executors from persistent GPU workers to mitigate cold-start latency.</li><li>Implemented a Worker Controller to manage GPU workers, model lifecycle, and resource allocation under dynamic workloads.</li><li>Reduced cold-start engine overhead by 50% per model (from 2.6s to 1.3s) by avoiding repeated CUDA initialization.</li><li>Achieved 10% faster end-to-end cold start over baseline vLLM, with 6.5s savings across 3 models, validated via controlled benchmarks.</li></ul>	

## SKILLS & CERTIFICATIONS

---

Languages: Proficient in English and Chinese, conversant in German

Technical Skills: Python, C, Java, SQL, React, JavaScript, TypeScript, Docker, Git, Kubernetes

Certificates: Amazon Solution Architect (Associate), Generative AI with Large Language Models by DeepLearning.AI

Hackathons: EasyA x BCG Hackathon (4<sup>th</sup> Place), NTU Techfest Hackathon (Participant)