

Education

University of Toronto	Toronto, CA
M.Eng., Computer Engineering	Jan. 2024 - Apr. 2025(Expected)
Relevant Courses: Distributed System in Moden Web-Scale Applications, Video Encoding, Blockchain Development, AI Alignment, Natural Language Processing Application.	
Beijing University of Posts and Telecommunications	Beijing, CN
B.Eng., Software Engineering	Sep. 2019 - Jun. 2023

Technical Skills

Python, JavaScript, Java, Machine Learning, Deep Learning, MLOps, LLM Agent, PySpark, Azure, Git, Github, Flask, Django, MySQL, PyTorch, MLflow, Microsoft Fabric, Databricks, Multimedia Encoding, Distributed System.

Work Experience

Microsoft	Beijing, CN
Software Engineer Intern	Articles on <a href="#">Microsoft Learn</a>
Feb. 2023 - Dec. 2023	
<ul style="list-style-type: none"><li>Contributed to <b>FLAML</b>, Microsoft's open-source <b>Automated Machine Learning (AutoML)</b> library, initiated an internal version for the <b>Microsoft Fabric</b> data platform. Submitted <b>100+ PRs</b>, driving a <b>263% increase</b> in user adoption.</li><li>Expanded AutoML model candidate pool by adding <b>15+ Machine Learning models (PyTorch Lightning, PySpark MLlib, Statsmodels, Scikit-learn)</b>.</li><li>Integrated multiple <b>MLOps</b> services (experiment tracking, visualization, artifact logging, deployment) with custom <b>MLflow</b> backend, reducing <b>inference latency</b> by <b>37%</b>.</li><li>Developed the <b>first-ever autonomous feature engineering</b> module using <b>LLM Agent</b> within an AutoML library.</li><li>Benchmarked AutoML solutions against <b>Azure Machine Learning Studio</b> and <b>Azure Databricks</b>; FLAML on Fabric achieved <b>30× faster runtimes</b> and a <b>16%</b> average performance improvement over competing solutions.</li></ul>	
Call Fusion	Toronto, CA
Research Software Engineer	Try our <a href="#">Demo</a>
Sep. 2024 - Present	
<ul style="list-style-type: none"><li>Redesigned the technical architecture for an <b>Audio LLM Agent</b> using OpenAI's Real-time API, reducing end-to-end speech-to-speech <b>reply latency</b> by <b>52%</b>.</li><li>Integrated RNNNoise, a neural network-based noise-reduction algorithm, boosting <b>transcription accuracy</b> by <b>13%</b>.</li><li>Evaluated and migrated from <b>Faiss</b> to <b>ChromaDB</b> for Retrieval-Augmented Generation (RAG), <b>improving retrieval performance</b> by <b>16%</b>, <b>reducing 96% of peak cases</b>.</li></ul>	
National University of Singapore	Singapore, SG
Undergraduate Research Assistant	
Apr. 2022 - Feb. 2023	
<ul style="list-style-type: none"><li>Collaborated on an <b>Event Forecasting with Temporal Knowledge Graph</b> project at the NExT++ Centre, focusing on research design and implementation.</li></ul>	

Side Projects

AutoGen, A programming framework for Agentic AI
<ul style="list-style-type: none"><li>Contributed to <b>high-level system design</b> in an Agentic AI framework with <b>35k+ GitHub stars</b>.</li><li>Early contributor to the <b>RAG</b> module, benchmarking various vector databases and integrating <b>ChromaDB</b>.</li><li>Added support for <b>Anthropic Claude</b> models as an alternative foundation model for LLM Agents.</li></ul>
D2Helper, A Full-Stack Web Application for Destiny 2
<ul style="list-style-type: none"><li>Developed a <b>full-stack</b> web application for the Destiny 2 community with 4,000+ active users.</li><li>Built frontend with <b>JavaScript and CSS</b>; implemented a <b>RESTful</b> backend with <b>Flask</b> and a <b>MySQL</b> cloud database.</li><li>Deployed via <b>Docker</b> on Tencent Cloud's <b>Serverless</b> platform.</li></ul>