

# RICHARD SO

📍 Seattle, WA | 📞 347-281-3815 | 📩 richardso2021@gmail.com | 💬 in/richardso21 | 🌐 sorichard.com

## EDUCATION

---

<b>Georgia Institute of Technology</b> <i>M.S. Computer Science, Machine Learning - GPA: 4.0</i>	08/2024 - 05/2025 Atlanta, GA
• Coursework: ML, Deep Learning, Computer Vision, NLP, Databases, Networks, Algorithms Honors	
<b>Georgia Institute of Technology</b> <i>B.S. Computer Science - GPA: 4.0</i>	08/2021 - 05/2024 Atlanta, GA

## WORK EXPERIENCE

---

<b>Amazon</b>   Software Development Engineer	08/2025 - Present
• Maintained <b>data lake infrastructure</b> (AWS Glue, EMR, Lambda, S3) for Amazon Brand Analytics, enabling sellers to monitor <b>sales performance insights across 2B+ products</b> .	
• Orchestrated <b>high-throughput data vending pipelines</b> with Apache Spark to transform ~100 TB/week of raw purchase activity into curated analytical datasets for internal teams and sellers at scale.	
• Developed core components of an <b>explainable ML root cause analysis service</b> , leveraging <b>Shapley value feature attribution</b> to identify drivers of product underperformance relative to benchmarks.	
<b>Data To Insights (D2I) Lab @ Georgia Tech</b>   Research Assistant	01/2025 - 05/2025
• Investigated time-to-first-token (TTFT) reduction in <b>multi-tenant LLMs</b> by overlapping document retrieval and prefill.	
• Led <b>scenario-based benchmarking</b> by simulating web crawler and ANNS retrieval to evaluate latency improvements of our custom token prefilling mechanism under realistic workloads.	
• Analyzed results for and <b>co-authored MLSys research paper</b> demonstrating up to <b>11x faster latencies in TTFT</b> .	
<b>Amazon Web Services</b>   Software Engineering Intern (ML)	05/2024 - 08/2024
• Implemented <b>ML-driven analysis of developer workflows</b> by aggregating behavioral features (keystrokes, UI events, time-on-task), enabling quantitative DX comparison across cloud platforms and <b>influencing AWS roadmap decisions</b> .	
• Re-architected a batch inference pipeline to exploit parallelism, <b>reducing runtime by &gt;85% (hours → minutes)</b> .	
• Automated developer intent and task classification from desktop screenshots using AWS <b>Rekognition, Textract</b> , and <b>multi-modal LLM prompting on Bedrock</b> .	
<b>Tanium</b>   Software Engineering Intern	06/2023 - 08/2023
• Built CRUD logging to Tanium console's PostgreSQL backend and REST endpoints to support a <b>customer-facing audit feature</b> , allowing review of console activity and detection of unauthorized configuration changes.	
• Resolved <b>50+ feature/bug tickets</b> within a 10-week internship maintaining a Knex.js + React TypeScript codebase.	
• Ensured code quality and correctness by applying <b>TDD and validation best practices</b> using Jest, Jasmine, and Joi.	

## PROJECTS

---

<b>Generative Data Augmentation for Image Classification</b> ⚡	04/2024
• Experimented with <b>Stable Diffusion</b> and <b>ControlNet</b> to enhance image classification accuracy when data is scarce.	
• Observed a <b>10% F1 increase</b> for Resnet-50 on a limited dataset when augmented with ControlNet-generated images.	
• Orchestrated <b>large-scale DL experiments</b> , sweeping across 100+ hyperparameter and model configurations.	
<b>LC3Tools</b> ⚡	09/2023 - 05/2024
• <b>Lead maintainer</b> of the educational tooling suite for coding, assembling, and simulating <b>LC-3 assembly programs</b> .	
• Actively collected student & instructor feedback to continuously drive <b>major quality-of-life enhancements</b> .	
• Served <b>1000+ students every semester</b> as one of the core, required tools for Georgia Tech's CS2110 course.	

## SKILLS

---

<b>Programming Languages</b>   Python, TypeScript/JavaScript, Java, Scala, C/C++, Go, Lua
<b>Frameworks &amp; Libraries</b>   NumPy, Pandas, SkLearn, PyTorch, Lightning, Apache Spark, React, Svelte
<b>Databases &amp; Misc.</b>   PostgreSQL, SQLite, DynamoDB, Elasticsearch, Docker, Apache Airflow, AWS CDK