

# JIAJUN (NICK) HUO

jiajunhuo726@gmail.com · 267-777-8518 · github.com/nickhuo · nickhuo.com · IL (UTC-6)

## EDUCATION

<b>University of Illinois Urbana Champaign</b> <i>Master of Science in Information Science (STEM)</i> GPA 3.9/4.0 · A in Applied Machine Learning & Database Systems · U.S. News #1 Information Systems program	Aug 2024 – May 2026 Urbana, IL
<b>Shenzhen University</b> <i>Bachelor of Science in Mathematics</i> GPA 3.79/4.50 (Top 10 %) · 2× merit scholarships · 4 national awards in ML modeling (supply-chain, finance, computer-vision)	Sep 2019 – Jun 2023 Shenzhen

## WORK EXPERIENCE

<b>Donut Labs</b> <i>Software Engineer Intern – 1st AI Agentic Crypto Browser (Sequoia-backed)</i> • Built GPU-free sentiment pipeline for 50K daily tweets, cut inference cost from \$6 to \$0.03/day using custom ANN model • Scaled 18-route FastAPI microservice to 3× QPS with P95 220 ms; reduced 4xx/5xx 40%, frontend adoption to 66% • Integrated 46 crypto actions across 16 providers; improved agent task completion 22% with constrained decoding • Designed Next.js microservice mesh with 5 MCP endpoints, enabling horizontal scaling and isolated failure handling • Implemented Redis model weight caching and KV reuse strategy, reducing P99 latency 26% and GPU costs \$2K/month	Jun 2025 – Aug 2025 Remote · NYC HQ
<b>Sonic SVM</b> <i>Software Engineer, Data Solution – Bitkraft-backed SaaS for game studios</i> • Led the development of data governance system, drove 3 C-level strategic decisions on product roadmap, report to CPO • Scaled data platform 2× by integrating new sources while maintaining <1hr SLA through incremental processing • Automated PostgreSQL→BigQuery ETL workflows, cutting data processing time 43% and enabling real-time analytics • Created Prometheus + Grafana dashboards tracking KPIs that influenced product roadmap and secured 10+ partnerships	Mar 2023 – Jul 2024 Shenzhen
<b>Tencent</b> <i>Technical Product Manager Intern – The largest WeChat ecosystem platform for universities</i> • Drove the full SDLC for campus app, mitigating risks and unblocking dependencies to launch from scratch in 12 weeks • Translated product requirements into technical specifications for a real-time notification system serving 50k+ students • Leveraged user engagement metrics from A/B tests to guide feature priorities, boosting student engagement by 22%	Nov 2020 - May 2021 Shenzhen

## PROJECTS

<b>ReplicaGenie</b> – Real-time, cost-efficient RAG platform for personalized LLMs • Architected 4-service LLM platform; modular microservices allow independent scaling, isolated failures & fast rollouts • Delivered 4-bit LoRA fine-tuning on SageMaker, trimming GPU training cost 77% and guaranteeing reproducibility • Built MongoDB → RabbitMQ CDC stream pushing changes to feature store in <1 s, enabling real-time personalization • Automated scraping for 4 data sources via Selenium on AWS Lambda, parsing and persisting structured data in MongoDB	
<b>BiteMatch</b> – Smart recipe recommender, ranked top-10% in class • Owned the backend development, creating a RESTful API with 12+ endpoints for user auth, CRUD, and video search • Reduced YouTube API calls by caching videos in Redis and using a quota-aware fetch logic to avoid exceeding limits • Worked with 4-member Agile squad, grooming backlog and delivering on 2-week sprints with 100% sprint goal hit rate	GitHub · Demo
<b>Claude Code Extension for Raycast</b> • Built TypeScript & React Raycast extension (5 modules) transforming Claude CLI into a GUI, accelerating dev workflows • Implemented JSONL session search with real-time filters, plus memory limits and timeout guards for reliable performance	GitHub

## SKILLS

<b>Languages</b>	Python, TypeScript/JavaScript, Go, SQL, Java
<b>Technologies</b>	Backend: FastAPI, Django, Node.js   Frontend: React, Next.js   Data: PostgreSQL, MongoDB, Redis
<b>Cloud/DevOps</b>	AWS (EC2, S3, Lambda, SageMaker), GCP, Docker, Kubernetes, Terraform, Kafka, CI/CD, Prometheus