

# Alfred Dwan

New York, NY • +1 5513768692 • [alfred.dwb@gmail.com](mailto:alfred.dwb@gmail.com)

## EDUCATION

<b>New York University</b> Master of Computer Science & Engineering	GPA: 4.0/4.0	Expected May 2026 New York, NY
<b>City University of Hong Kong</b> Bachelor of Computer Science		2020 - 2024 Hong Kong

## EXPERIENCE

<b>PredictX (Startup)</b> <b>Software Engineer Intern - AI track</b>	May 2025 – Sep 2025 Hong Kong
<ul style="list-style-type: none"><li>In collaboration with the DeepSeek AI team, shipped a user-facing intelligent Q&amp;A feature that increased the self-service success rate by 25% by providing instant, accurate answers to complex user questions.</li><li>Engineered a user-facing AI backend service using FastAPI, Docker, and Kubernetes (HPA, readiness/liveness probes), successfully maintaining a p95 latency of &lt;500ms under load, supporting 10K+ monthly requests via async I/O and backpressure controls.</li><li>Architected the underlying agentic RAG workflow for this feature using LangChain and FAISS (HNSW/IVF), which automated query processing and cut internal support's task resolution time by 40%.</li><li>Optimized LLM inference costs by 35% with no performance regression by implementing LoRA fine-tuning, quantization (BitsAndBytes), and packaging models with TensorRT for efficient serving.</li><li>Established the end-to-end MLOps pipeline that accelerated the delivery of new AI features, reducing the model deployment cycle from days to hours using GitHub Actions and integrating a monitoring stack with Prometheus, Grafana, and OpenTelemetry for proactive alerting.</li></ul>	
<b>ByteDance</b> <b>Software Engineer Co-op Program (Industry Collaboration)</b>	Jun 2024 – Aug 2024 Beijing
<ul style="list-style-type: none"><li>Improved the listening experience for users on low-quality networks by reducing audio streaming latency by 22%, using an agent-driven framework in multi-threaded C++ and Swift (AVFoundation) that optimized buffer pre-warming and adaptive chunking.</li><li>Designed a real-time, multi-agent orchestration layer using Kafka, Redis, and WebSockets to handle concurrent live-stream events, improving system stability by 25% during peak traffic.</li><li>Shipped an AI-powered content discovery feature that increased Daily Active Users (DAU) by 12%, using PyTorch and FAISS to reduce user navigation time by 20%.</li><li>Cut model inference latency and memory usage by over 30% by deploying HuggingFace agents with ONNX Runtime on Kubernetes, managed via an Istio service mesh for intelligent traffic shifting and fault tolerance.</li><li>Architected an automated canary and rollback system for the CD pipeline using Prometheus time-series telemetry, cutting incident recovery time from hours to minutes and safeguarding production SLOs.</li></ul>	
<b>Greenhouse Data (Startup)</b> <b>Software Engineer Intern</b>	Feb 2023 – May 2023 Hong Kong
<ul style="list-style-type: none"><li>Built a distributed web crawler in Node.js microservices using JavaScript and TypeScript; achieved 10K+ pages/min with Express.js, PM2, RabbitMQ (&lt;5 ms); deployed on AWS EC2 via Docker Swarm, Jenkins CI/CD, and Prometheus monitoring with Unix/Linux environments</li><li>Designed MongoDB + Mongoose for 1TB+ data with SQL query optimization; validation/indexes ensured 99.99% consistency and 5 ms queries; added full-text search via Elasticsearch and implemented PostgreSQL for relational data management.</li><li>Boosted throughput with Redis Cluster caching and Lua atomic ops, reaching 10K+ req/sec and cutting API latency by 20%; secured access with OAuth 2.0/session auth and HTTP/HTTPS protocols through security software development practices.</li><li>Shipped a real-time data monitoring dashboard for clients using React, allowing users to track job status and performance from the web application.</li></ul>	
<b>Siemens</b> <b>Software Engineer Intern</b>	Sep 2022 - Jan 2023 Hong Kong
<ul style="list-style-type: none"><li>Built a full-stack monitoring dashboard for station control systems using React, FastAPI, and PostgreSQL, enabling real-time visualization of sensor and telemetry data across 20+ IoT nodes.</li><li>Shipped public-facing RESTful APIs that enabled self-service data access for field engineers, standardizing communication between IoT devices and the cloud.</li><li>Engineered the CI/CD pipeline using Jenkins, Docker, and Google Test, which accelerated the feature delivery lifecycle by 40% for critical mobility projects.</li><li>Led the modernization of legacy C++ modules by migrating them to containerized microservices, improving system scalability and reducing maintenance downtime by 5%.</li></ul>	
<b>Dolby</b> <b>Software Engineer Intern</b>	Mar 2022 - Aug 2022 Beijing
<ul style="list-style-type: none"><li>Built responsive React + TypeScript applications with modern state management (React Query/Redux Toolkit), route-level code splitting, and performance budgets (LCP/CLS/TTI); configured Vite/Webpack optimizations and HTTP caching.</li><li>Accelerated UI development by 30% by implementing a Storybook-driven design system and WCAG compliance (keyboard paths, focus management, contrast), reducing UI inconsistency and review churn.</li><li>Improved reliability with resilient API clients (retry with jitter, timeouts, centralized error boundaries) and added client-side observability (web-vitals plus custom events, correlation IDs); integrated Jest unit tests and Cypress/Playwright E2E into Jenkins pipelines with artifact retention.</li></ul>	
<h2>PROJECTS — Cloud-Native GoPaaS Platform Development</h2>	
<ul style="list-style-type: none"><li>Built a cloud-native GoPaaS on Go-micro v3 and Kubernetes with service mesh using hardware optimization techniques; standardized config, discovery, circuit breaking, and rate limiting; operated scalable clusters with Docker containers, NGINX Ingress, and custom controllers; achieved &lt;50 ms inter-service latency and 99.95% uptime.</li><li>Integrated Istio (mTLS, traffic shifting, retries) and Helm deploy/rollback workflows using market data protocols for financial markets applications; release cycles accelerated by 3x.</li><li>Standardized microservices with gRPC/protobuf using binary protocols; added an API Gateway with Hystrix, load balancing, and rate limiting; Prometheus/Grafana observability cut peak-time failures by 65% and sustained +60% traffic growth.</li><li>Automated Helm-based canaries and rollbacks, tripling deployment throughput while reducing change failure rate and aligning practices for multi-region release hygiene</li></ul>	