# John Kitaoka

jckitaok@gmail.com — 507.884.1155 — in/johnkitaoka — johnkitaoka.github.io

AI architect building production systems at scale: fine-tuned LLMs and VLMs, custom evaluations, and agentic architectures. Depth in technical communication, AI/ML systems, and public sector vertical software.

## EXPERIENCE

**Applied AI Architect II** Amazon Web Services (New York, NY) *2024 - Present*

- Design, develop, evaluate, and optimize AI/ML workloads on AWS Cloud for public sector customers, leading 6 end-to-end AI system deployments, from initial design through production, with technical patterns reused and adopted org-wide
- Built end-to-end RAG systems: hybrid retrieval (semantic + keyword) with reranking, agentic query decomposition for multi-hop reasoning and synthesis, and guardrails for PII/policy compliance
- Designed agentic workflows with multi-agent orchestration, custom tooling, data source aggregation, and report generation; for document processing (90% review time reduction) and life science research (60% query writing time reduction)
- Fine-tuned Qwen3-4B with LoRA for long-context QA (64K tokens); built custom industry-specific evaluation benchmark, achieved 70% accuracy improvement over base model at 90% lower cost than Claude Haiku

**Solutions Architect II** Amazon Web Services (New York, NY) *2022 - 2024*

- Design, manage, and optimize cloud software architectures for public sector organizations through direct and scaled engagement
- Built first public RAG workshop (Semantic Search with Amazon OpenSearch) across major cloud providers; presented at Re:Invent 2023/24/25 to 1,000+ practitioners, with reference architectures adopted in production by multiple organizations
- Led technical evaluations against hyperscaler and frontier lab alternatives; designed proof-of-concept implementations that won 8 competitive assessments including 3 platform migrations
- Architected 100TB data platform migrations, replacing self-managed single-node databases with managed stack: Aurora Postgres, zero-ETL replication to Redshift, BI layer—deployed across 4 organizations, 70%+ reduction in data engineering work
- Created enablement curriculum on AI responsibility for 150 solutions architects; content scaled org-wide to 750+ practitioners

## PROJECTS

**Multimodal AI Form Filling** — GitHub: aws-samples/sample-bedrock-form-filling *December 2025*

- Built a serverless, event-driven pipeline for multimodal-to-text extraction and eForm automation using AWS Step Functions and Amazon Nova LLMs, deployed via SAM/CloudFormation, with CloudFront and Cognito

**Open-Source Software** *2023 - Present*

- Led LangChain-AWS features (embeddings support: models, parameter abstractions, performance optimizations, multimodal content), AWS code samples and projects (AWSLabs, aws-samples), contributor to RL frameworks (SkyRL, Gymnasium)

**Building a Production Text2SQL System for Public Health**— Customer Success Story *October 2025*

- Architected NL2SQL system for health reporting: Bedrock query generation, PostgreSQL vector store for embedded schemas, ECS/Fargate serving; delivered POC and technical enablement, migrating from OpenAI; 50% reduction in report development

**Blind Spots in Vision-Language Models** — Amazon internal CV conference (ACVC 24) *August 2024*

- First author on internal paper: introduced spatial awareness evaluation methodology for LLaVA-style VLMs; training modifications improved peripheral region accuracy 30%; framework integrated into production model evaluation pipeline

**Fine-Tuning and Hosting LLMs in Isolated Environments** — AWS blog series *June 2024*

- Led development on a technical blog series on fine-tuning LLMs with QLoRA and deploying an inference server with HuggingFace libraries: PEFT, transformers, and datasets for customers in national security, defense, and federal software

**Generative AI in the Public Sector** — AWS eBook *January 2024*

- Co-authored official AWS advisory guide on generative AI strategy for public sector; 15K+ downloads, used by 3,000+ tech leaders for production AI planning, writing chapters on foundation model selection and public sector reference architecture

## EDUCATION

**University of Wisconsin-Madison** Madison, WI

*BSc, Mathematics and Computer Science; Minor in Business* *GPA: 3.92/4.0*

- Top 5% of class, Distinction in both majors, D.E. Shaw Nexus Fellowship, Men's Water Polo, Mayo Foundation Scholarship

---

**Languages and Tools**: Python, Java, C, C++, SQL

**Frameworks**: pandas, PyTorch, TensorFlow, LangChain, CUDA, Neuron, spark, transformers, accelerate

**Certifications**: AWS Solutions Architect Professional, AWS Machine Learning Specialty, AWS SysOps Administrator