

**Gao Zhiyuan**

+86 13262224615  
alapha23@gmail.com

## Work Experiences

### Amazon Web Services

*Associate Solutions Architect*

*Jun. 2024 - now*

#### 1. Global Network Infrastructure Optimization Project

- Monitored and analyzed network latency between the USA and Singapore, identifying bottlenecks within individual CloudFront components causing TCP reconnections.
- Developed and implemented a PoP warm-up strategy and an NGINX-based solution to ensure persistent TCP connections.
- Reduced cross-region latency inconsistency from 60% to under 10%.

#### 2. RAG-based Chatbot for AWS Knowledge

- A chatbot using Retrieval Augmented Generation, to answer technical questions specifically focused on AWS services
- Evaluate the responses in terms of hallucination, self-knowledge dependence, Context Utilization and etc.
- The tech stack includes AWS Lambda, RDS, Opensearch, Bedrock/Claude, Sagemaker, Rerank.

#### 3. Multimodal Content Guardrailing

- Employed a model unlearning project to forget certain concepts during inference time
- Moderate Content through AWS Nova, Bedrock Guardrails, and the model unlearning project

#### 4. Improve YOLOv8 Inference Throughput with AWS Inf1 Instance - Using Pytorch to convert the model weight of YOLOv8 for AWS Inf1 Instance. - Gained 1.85x throughput increase whilst lowering the cost down to 25

#### 5. EKS Node Group Migration - For a customer who uses only elastic IPs (EIPs), and wishes to migration from unmanaged node group to managed node group in production, - I created a user data script and used VPC endpoints to make sure the newly started unmanaged node groups can acquire EIPs. Then I guided the customer for pod evictions to complete the migration.

### Trilogy

*Software Engineer*

*Aug. 2022 - Aug. 2023*

#### 1. AWS Cost Optimization project

- Lowered AWS Service Costs from 25 million USD annually down to 5 million USD in a team of 3.
- Designed non-intrusive, sustainable plans to guarantee resource sharing, user resource awareness, automated anomaly detection and permission rules.

#### 2. Social Platform for Education

A social platform that enables 4-12 grade students to memorize curricular knowledge through making posts with the help of AI and memorization techniques.

The tech stack includes Next.js, Prisma/MySQL, GPT-3, Stable-Diffusion, Docker and AWS DevOps.

#### 3. Data Platform based on Azure AD

Construct a graph-based data platform and use it to calculate company employee interaction intensity.

Extract data from MS 365 through Azure Graph API, engineer the data into graph relationships and maintain the in the graph database.

## MISE

*Devops Engineer*

*Jan. 2022 - Sep. 2022*

Addressed latency spikes in an AWS Lambda-based startup application during peak traffic hours.

Diagnosed and resolved performance bottlenecks by analyzing AWS Lambda, API Gateway, and MySQL RDS.

Streamlined deployment processes using Terraform Infrastructure as Code (IaC).

## MetaMUI-SovereignWallet

*Python SDK Engineer*

*Jun. 2021 - Jan. 2022*

Migrate Rust SDK to Python SDK for MetaMUI blockchain built with Substrate framework

## Projects

### Seoul National University

*Research Intern*

City Energy Lab

*Sep. 2023 - Now.*

RAG-based Chatbot for Urban Planning

- Built a RAG chatbot with FAISS, parent-child chunking, keyword-extraction for domain-specific knowledge

- Designed a chatbot with NextJS, Typescript, Prisma, FAISS and OpenAI API

### Seoul National University

2019

Project Name: Accomodate Input Spikes with AWS Lambda Functions

A hybrid system of Dockers and AWS Lambda Functions, where stable workloads are running on docker and input rate spikes are offloaded to AWS Lambda Functions,

Implemented job scheduling/balancing algorithm to coordinate Lambda Functions' life cycle, data relay, and system-level dynamic scaling

### Google

*Apache Foundation*

Summer of Code

*2019*

Engineered Apache Nemo to process single-stage batch data with AWS Lambda Functions

### Chinese Academy of Science

*QEMU Researcher*

PLCT Lab

*2020*

One patch to QEMU upstream to emulate Nuclei RISC-V SoCs with customised interrupt controllers and registers,

Delivered an oral presentation at CRVA 2020

## Education

### National Cheng Kung University

*B.S., Political Science*

Taiwan

*Sep. 2015 - May. 2017*

### Seoul National University

*B.S., Computer Science and Engineering*

Korea

*Sep. 2019 - Aug. 2024*

**Languages:** Chinese, English (TOEFL 100), Japanese (JLPT N1) and Korean

**Skills:** Languages: Python, Java, Typescript, C/C++, NextJS

DevOps: AWS, Redis, K8S, Docker

Others: AI Agent / RAG development, Serverless development, QEMU, React, Selenium