

Ajay Kumar Reddy Inavolu

+1 (607)313-9336 | inavolu.a@northeastern.edu | [LinkedIn](#) | [GitHub](#)

PROFESSIONAL SUMMARY

Dedicated Engineer with expertise in building scalable, high-performance microservices and real-time data processing solutions. Proficient in designing and deploying distributed cloud and AI platforms, leveraging Java, Spring Boot, Docker, AWS.

WORK EXPERIENCE

Software Engineer, ONE COMMUNITY GLOBAL | San Gabriel, CA

August 2024 – November 2024

Spring Boot, Next.JS, GraphQL, ELK, Kafka, Docker, AWS

- Engineered a **GraphQL** layer over Spring Boot microservices for an aggregator platform, optimizing data retrieval by consolidating client requests into single queries, reducing average network traffic by **30%**, and enhancing system performance and throughput
- Migrated a monolithic employee management system to containerized microservices using **Spring Boot** with **OAuth 2.0** integration for secure authorization, increasing throughput by **30%** and achieving **95%** test coverage via **JUnit**, **PowerMock**, and **Mockito**
- Established a low-latency concurrency mechanism in **Java 17** for the Employee Management App for real-time task updates, achieving sub-second performance for 200+ concurrent tasks using **Multithreading** and **Kafka**-based asynchronous processing
- Crafted **Python** automation scripts for automating the creation of employee timesheets and payroll statements with efficient batch job processing, reducing processing time by **15%** and optimizing financial management workflows

Software Engineer, OMDENA | Palo Alto, CA

July 2022 – October 2022

Spring Boot, Next.JS, Redis, Docker, AWS

- Engineered **RESTful**, containerized microservices leveraging **Spring Boot** for a credit scoring platform hosted on **Elastic Kubernetes Service**, enabling user engagement with score attributes, boosting user base by **10%** and optimizing API response times by **20%**
- Implemented a real-time log analytics platform with **Elasticsearch**, **Logstash**, **Kibana**, and **Apache Kafka** for active log aggregation across distributed microservices, improving L3 issue resolution by **30%**
- Leveraged **Redis** & **Express Rate Limiter** to cache DB query requests & session details with effective handling of **OAuth** tokens to grant endpoints access to external applications, reducing the DB read load by **60%**
- Reduced deployment times by **30%** by integrating Jenkins Core Agent into the CI/CD pipeline with SonarQube, optimizing code quality and deployment efficiency on AWS EC2 instances

Deep Learning Intern, INDIAN INSTITUTE OF TECHNOLOGY (IIT) | Indore, India

June 2021 – August 2021

TensorFlow, Keras, Airflow, PySpark

- Achieved **96% accuracy** in detecting human activities by developing a robust deep learning model using **TensorFlow** and **Keras**, improving model precision by **15%** through iterative threshold optimization and hyperparameter tuning
- Designed and implemented **data pipelines** with **TensorFlow**, improving training and inference efficiency by **10%** and reducing model training time by **20%**
- Orchestrated a fully automated model retraining workflow with **Airflow**, reducing deployment time by **5%** and automating **90%** of model retraining, enhancing system scalability and adaptability to new sensor data
- Processed over **100GB** of sensor data using **PySpark** and **Spark/Scala**, improving data ingestion speed by **25%**, while leveraging the **ELK stack** for real-time monitoring, cutting error detection time by **30%**

EDUCATION

Northeastern University | Khoury College of Computer Science | Boston, MA

Expected: May 2025

Master of Science in Computer Science

GPA: 3.9/4.0

Courses: Design Patterns, Algorithms, Database Management Systems, Cloud Computing, Web Development, Mobile Application Dev

ACADEMIC PROJECTS

Retrieval Augmented Generation API Analyzer | (Python, Langchain, Huggingface, Chroma, AWS) October 2024 – December 2024

- Developed RAG pipeline integrating Llama 3.1, LangChain, AWS Titan embedding model for enhanced API insights, reducing search time by **50%**, and seamlessly deployed on AWS Bedrock designed to support **10,000+** concurrent requests

Distributed Graph Analysis | (PageRank, Scala, Java, Spark, MapReduce, AWS)

January 2024 – March 2024

- Implemented **PageRank algorithm** using **MapReduce** and **Spark** to efficiently analyze a **1-million-node** dense web graph on an **EMR cluster**, with seamless I/O handling via **S3 buckets**
- Leveraged data lineage tracking for enhanced debugging and performance optimization, reducing execution time by **20%**

TECHNICAL SKILLS

Programming/Web Technologies: Java, Python, JavaScript, React, HTML, CSS, SASS, PHP

Framework/API: Spring MVC, Spring Boot, Node.js, Next.JS, TensorFlow, Microservices, Hibernate, JUnit, Mockito, JWT

Database/Cloud: MongoDB, NoSQL, PostgreSQL, RDBMS, SQL, MySQL, Oracle, CloudWatch, MS SQL, AWS EC2

Version Control/Tools: Git, Jenkins, Postman, Maven, Selenium, TestNG, JIRA, IntelliJ, Docker, PySpark, ELK, Kubernetes, Redis