# Affan Arif Khamse

*Backend Software Engineer | LLM Systems | Scalable Architectures | NYU CS Grad 2025*

New York City, NY · (516) 853-4972 · khamseaffan@gmail.com · github.com/khamseaffan · linkedin.com/in/affan-khamse

## Education

**Master of Science in Computer Science (GPA: 3.8/4)** — New York, USA
New York University — *May 2025*
*Courses: Software Engineering, Cloud Computing, Machine Learning, DSA, Databases*

**Bachelor of Engineering in Computer Engineering (GPA: 9.00/10)** — Mumbai, India
University of Mumbai — *June 2023*

## Technical Skills

**Programming and Scripting Languages:** C++, Java, Python, Node.js, JavaScript
**Frameworks and Libraries:** Flask, Spring Boot, FastAPI (RESTful), Django, ReactJS
**Cloud and DevOps:** Amazon Web Services / AWS, Azure, Docker, Kubernetes, GitHub Actions, Travis CI
**Databases & Query Languages:** SQL(MySQL, PostgreSQL), NoSQL(MongoDB, Firebase)
**Testing & Debugging:** JUnit, PyTest, Django Test, Chrome DevTools
**Tools and Methodologies:** Git, GitHub, Figma, UML, Agile (Scrum), Object-Oriented Design
**Specialized Expertise:** API Design, Microservices Architecture, LLM Applications, Scalable Systems

## Work Experience

**Software Teaching Assistant – Object-Oriented Programming** — New York, USA
Courant Institute of Mathematical Sciences @ New York University — Jan 2024 – May 2025
- Mentored 50+ students in **C++** and **Java**, emphasizing Object-Oriented design , algorithms, and complexity analysis
- Conducted code reviews and provided constructive feedback on programming assignments, enhancing code quality, maintainability, and adherence to best practices

**Software Engineer** — New York, USA
InquisAI (NYU)   |   inquis-ai.com 🔗 — Jun 2024 – March 2025
- Spearheaded development of an AI assistant builder that leveraged **LangChain** and **OpenAI** Embeddings for document-based question answering with vector store integration
- **Redesigned backend architecture**, migrating from Flask to FastAPI, improving request throughput and cutting API latency by 30% with asynchronous processing
- Architected and deployed high-performance, scalable **RESTful APIs**, optimizing data operations to enhance efficiency and scalability, demonstrating the ability to support **1K+ concurrent users**
- Led Agile sprint planning in a 3-person team, using Azure DevOps to drive backend scalability and infrastructure decisions, reducing delivery time by 25%

## Projects

**Home Store – E-Commerce Platform (Backend)** — Jan 2025 - Present
*Spring Boot, Docker, PostgreSQL, React Router, Firebase, Microservices, SwaggerUI* — GitHub
- Engineered a scalable microservices using **Spring Boot** and **Eureka service discovery** for modular service deployment & streamline inter-service communication
- Containerized all microservices using **Docker** and orchestrated multi-service local development environments with Docker Compose for production-like testing
- Designed **RESTful APIs** for seamless **React Router** frontend integration and documented endpoints with **SwaggerUI**

**FlashBids - Full Stack** — Sep 2024 - Dec 2024
*AWS EC2, DynamoDB, S3, Redis, WebSockets, CloudWatch, Flask, Python, REST APIs, DevOps* — Demo | GitHub
- Constructed a real-time bidding platform using AWS EC2 with Auto Scaling, DynamoDB, and Redis, handling over **10K+ concurrent user sessions**
- Developed non-blocking APIs for auction creation and bidding, reducing response **latency by 30%**
- Built a low-latency WebSocket communication layer to broadcast live bids, **reducing client-server latency by 40%** and improving system reliability
- Orchestrated an **event-driven** bidding pipeline using Amazon SQS/SES (bid placement → live updates → auction close → email notifications), **achieving 99% bid placement success**

**VibeCheck** (Full-Stack Project) — Sep 2023 - Dec 2023
*Django, Python, Bootstrap, AWS Elastic Beanstalk, PostgreSQL, Redis, TravisCI* — Demo | GitHub
- Created a social platform that paired users based on real-time Spotify listening data and music preference similarity
- Established a reactive messaging framework using Redis Pub/Sub, reducing end-to-end **chat latency by 30–40%**
- Deployed via AWS Elastic Beanstalk with **CI/CD (Travis CI)**, a**chieving 87% test coverage** ensuring reliability