

Sai Paresh Karyekar

 Sunnyvale, CA  (404) 200-3907  [saikaryekar@gmail.com](mailto:sai.karyekar@gmail.com)  github.com/saikaryekar  linkedin.com/in/sai-karyekar

EDUCATION

Georgia Institute of Technology, Atlanta, GA	May 2025
<i>Master of Science in Electrical and Computer Engineering</i>	<i>GPA: 3.9/4</i>
Veermata Jijabai Technological Institute (VJTI), India	June 2023
<i>Bachelor of Technology in Electronics & Telecommunication Engineering</i>	<i>CGPA: 9.62/10</i>

EXPERIENCE

Amazon, Palo Alto Software Development Engineer	June 2025 - Present
<ul style="list-style-type: none">Scaled Amazon Visual Search infrastructure (25M+ MAU, 1K+ TPS) to 10+ global marketplaces using AWS CloudFormation (IaC) and CloudWatch-based CI/CD pipelines for automated deployment and monitoring.Led cross-region production migration of a core request aggregation service with zero downtime through load testing, auto-scaling validation, and canary rollouts across Python and Node.js (TypeScript) services.Designed the system architecture for integrating Visual and Text Search, enabling scalable APIs and fault-tolerant service decoupling.Enhanced platform reliability and security by deploying a MapToken-based hybrid authentication system and improving observability via Log4j logging, CloudWatch metrics, and alerting.	
NVIDIA, Santa Clara Software Engineering Intern	
June 2024 - Aug 2024	

NVIDIA, Santa Clara Software Engineering Intern	June 2024 - Aug 2024
<ul style="list-style-type: none">Built an internal Generative AI content generation tool using RAG, reducing blog creation time by $2.5\times$.Developed and deployed a vector search pipeline with Milvus for efficient document retrieval.Configured and secured AWS EC2 infrastructure with firewall rules, VPN access, and load-balanced scaling.Containerized and orchestrated deployment with Docker and AWS ECS for scalable production rollout.	
PROJECTS	

ApplyCation: Automated Job Application Platform	 Oct 2024
<ul style="list-style-type: none">Developed a Flask backend with Selenium automation and Claude API for keyword-based resume tailoring.Leveraged GCP Storage for persistent user data and deployed a Streamlit UI for user input and job tracking.	

Enhancing Mathematical Reasoning in Small Language Models (SLMs)	 Sept 2024 - Dec 2024
<ul style="list-style-type: none">Fine-tuned a T5-small transformer using LoRA and Chain-of-Thought prompting, improving model reasoning accuracy by 12%.Exposed inference via Flask APIs for seamless integration with downstream applications.	

SKILLS

Languages:	Python, Java, TypeScript, SQL
Systems Backend & Cloud:	Flask, Node.js, REST/gRPC APIs, GCP, Docker, Distributed Systems
Infra & AI:	OpenSearch, Terraform (IaC), CloudWatch, PostgreSQL, DynamoDB, Redis, PyTorch, TensorFlow, LangChain, MLflow, OpenCV