

# Sai Kogilathota

+1(934)799-0867 | [saiakhilkogilathota@gmail.com](mailto:saiakhilkogilathota@gmail.com) | Stony Brook, New York, United States | [LinkedIn](#)

## SUMMARY

Software Engineer with 3+ years of experience integrating ML solutions into production systems across fintech and telecom. Skilled in NLP, predictive modeling, and API development using Python, Flask, and cloud platforms like AWS and Azure. Proven track record of building scalable pipelines and delivering data-driven features that improve decision-making and user engagement.

## WORK EXPERIENCE

<b>ML Researcher, Stony Brook University, Stony Brook, United States</b>	<b>Aug 2024 — Present</b>
<ul style="list-style-type: none"><li><b>Multimodal Research:</b> Introducing a lightweight probing framework that predicts hallucinations in <b>VLM</b> image captioning by analyzing hidden-states pre-generation; achieved strong results on <b>4 VLMs</b>, culminating in a <b>first-authored EMNLP 2025 submission</b>.</li><li><b>Clinical RAG Assistant:</b> Built and deployed a <b>clinical multimodal RAG system</b> using <b>LangChain, Pinecone, and AWS (EC2, S3)</b> to extract patient data and augment diagnostic processes, enabling faster and more accurate clinical decision support.</li><li><b>Self-Finetuning LLMs:</b> Developing a novel framework for self-finetuning LLMs using Agent-to-Agent (A2A) protocols orchestrated with <b>LangGraph</b> in controlled environments, enabling adaptive model improvement without manual intervention.</li></ul>	
<b>Software Engineer - AI/ML, Lendingkart Finance Limited, Bengaluru, India</b>	<b>Jan 2023 — Aug 2023</b>
<ul style="list-style-type: none"><li><b>LLM Fine-tuning with BERT:</b> Fine-tuned a BERT-Base model on internal search data using <b>Hugging Face &amp; PyTorch</b>, improving document retrieval accuracy by <b>70%</b>. Deployed the model on AWS SageMaker for scalable, real-time search capabilities.</li><li><b>NLP Analytics for Sales:</b> Increased sales conversion by <b>15%</b> by building and deploying an end-to-end NLP pipeline using <b>Google STT API</b> and sentiment models to analyze sales conversations.</li><li><b>KYC Automation:</b> Developed and deployed an <b>OCR-based and face-matching system</b> to accelerate onboarding, reducing onboarding time by <b>80%</b>; leveraged <b>PyTorch, YOLOv8</b> for real-time face detection, with end-to-end hosting and scaling on <b>AWS</b>.</li><li><b>Loan Eligibility Automation:</b> Automated loan eligibility workflows by embedding predictive analytics using <b>Scikit-learn and XGBoost</b> models, reducing manual processing time by <b>20hrs/week</b>, deploying this on <b>AWS EC2</b>.</li><li><b>Engagement Optimization:</b> Built &amp; deployed a personalized discount engine for <b>10,000+ users</b>, leveraging <b>Kafka</b> and <b>AWS SNS</b>, increasing loan completion rates by <b>35%</b>; conducted <b>A/B testing</b> on user behavior to further optimize engagement and conversions.</li></ul>	

<b>Software Engineer - AI/ML, Comviva Technologies, Bengaluru, India</b>	<b>Aug 2020 — Dec 2022</b>
<ul style="list-style-type: none"><li><b>ML Tool:</b> Developed and deployed a <b>Django based no-code ML tool</b> on <b>AWS EC2</b>, enabling business teams to build models there by improving operational efficiency; implemented <b>CI/CD pipelines</b> for seamless updates and faster feature releases.</li><li><b>Recommendation System:</b> Implemented <b>collaborative filtering</b> algorithms to personalize mobile campaign offers, increasing ARPU (Average Revenue Per User) by <b>10%</b>.</li><li><b>ML Model Deployment:</b> Trained and tuned <b>XGBoost, Ensemble models</b> and deployed them on <b>AWS SageMaker</b> for churn and multi sim user prediction, enabling proactive retention strategies that <b>reduced churn by 18%</b> and increased revenue by \$150K+/ quarter.</li><li><b>Real-Time Data Processing Pipelines:</b> Built and deployed real-time <b>ETL workflows</b> on <b>AWS</b> in <b>Python</b> using <b>Pandas</b> for ingesting telecom usage data into ML pipelines, <b>reducing processing lag by 40%</b> and supporting predictive analytics.</li><li><b>A/B Testing Framework:</b> Collaborated with marketing and product teams to develop an experimentation framework to <b>A/B test ML-driven campaign</b> features, improving <b>conversion rates by 17%</b>.</li><li><b>Sentiment Analysis:</b> Applied sentiment analysis on multilingual customer feedback using <b>BERT Models</b> and <b>Google Translate API</b>, surfacing key pain points and enabling UX/UI optimizations across 3 product lines.</li></ul>	

## PROJECTS

### Revize | React.js, Django, SQL, AWS | EMNLP 2025 Submission, [Link](#)

- Developed a smart web application for **spaced repetition learning**, enabling users to add topics, auto-schedule revision dates, track daily tasks, and review progress with flashcard-based content.
- Integrated **OpenAI SDK** for automated summarization into flashcards, and deployed the application on **AWS EC2** with a **Postgres** database for scalable, reliable performance.

### Stock Recommendation System | NLP, LSTM, Pytorch, Galformer, [Link](#)

- Developed a novel **PyTorch-based architecture** leveraging **Galformer**, a specialized transformer for stock systems and NLP, to analyze stock sentiment and generate buy/sell recommendations based on purchase price and quantity.
- Deployed the model with a **Streamlit** front end, achieving **15% better prediction accuracy** compared to a baseline LSTM model.

### TrialMatch | JavaScript, Python, React, MySQL, KMeans, [Link](#)

- A clinical trial recruitment platform with a swipe-based matching interface and integrated **recommendation system**, improving match accuracy for **rare disease trials**, **reducing participant search time by ~40%**, and **accelerating trial enrollment by up to 30%**.

## EDUCATION

### Stony Brook University, Stony Brook, United States - Master's, Data Science

Aug 2023 — May 2025

### Ramaiah Institute of Technology, Bengaluru, India - Bachelor's, Electronics and Communication

Aug 2016 — Jul 2020

## SKILLS

<b>Programming Languages</b>	:	Python, Java, R, MySQL
<b>ML Frameworks</b>	:	Scikit-learn, XGBoost, TensorFlow (basic), Rasa, spaCy, NLTK
<b>Model Development</b>	:	Supervised & Unsupervised Learning, Cross-Validation, Hyperparameter Tuning, A/B Testing
<b>Data Processing &amp; Pipelines</b>	:	ETL Pipelines, Data Wrangling, Feature Engineering, Pandas, NumPy
<b>APIs &amp; Deployment</b>	:	RESTful API Development, Flask, FastAPI, Docker, CI/CD
<b>Cloud Platforms</b>	:	AWS (SageMaker, EC2, S3), Azure
<b>Version Control &amp; Dev Tools</b>	:	Git, JIRA, VS Code, Postman