

# YAVAR KHAN

AI/ML ENGINEER

- San Jose, CA • (669)-339-8555 • [yavarkhan1997@gmail.com](mailto:yavarkhan1997@gmail.com)
- [linkedin.com/in/yavar-khan29](https://www.linkedin.com/in/yavar-khan29) • [github.com/yavar29](https://github.com/yavar29) • [yavarkhan.me](https://yavarkhan.me)

## PROFESSIONAL SUMMARY

AI/ML Engineer with 2.8 years at Accenture and a Master's in CS (AI/ML) from SUNY Buffalo. Specialized in production multi-agent systems, RAG pipelines, and LLM-powered applications deployed on AWS and GCP. Proven ability to ship and scale mission-critical systems across global teams.

## EDUCATION

### UNIVERSITY AT BUFFALO, THE STATE UNIVERSITY OF NEW YORK (SUNY)

AUG 2024 - DEC 2025

Master of Science - Computer Science (AI/ML Track) | GPA: 3.87

### AMITY UNIVERSITY, INDIA

JULY 2016 - MAY 2020

Bachelor of Technology - Information Technology

## PROJECTS

### WEATHERWISE.AI

[LIVE DEMO](#)

- Architected a multi-agent conversational weather intelligence system using **LangGraph** with fan-out/fan-in orchestration, dependency-aware parallel execution, and specialized AI agents for data retrieval, computation, research, and visualization.
- Built a Corrective **RAG pipeline** with Claude Vision chart grounding, embedding 429+ vectors from NOAA, NASA, and NCA5 research into Pinecone for scientifically grounded responses.
- Engineered a **sandboxed computational engine** with AST validation for real-time atmospheric analysis (potential temperature, LCL, vapor pressure deficit) applied to live MCP weather data.
- Deployed via Terraform to **GCP Cloud Run** with **LangSmith tracing** for end-to-end agent observability, Anthropic **prompt caching**, real-time **WebSocket streaming**, and dynamic ECharts visualizations.

### DEEP RESEARCH PRO - MULTI-AGENT ASYNC RESEARCH SYSTEM

[GITHUB](#)

- Designed an asynchronous **multi-agent research engine** using the **OpenAI Agents SDK**, featuring an autonomous query planner and multi-step verification for cited synthesis.
- Optimized system throughput by **50x** using **bounded-parallelism** and concurrent search-summarization threads with system-level guardrails.
- Engineered a **dual-layer caching** strategy (LRU + SQLite, 24h TTL) with context-aware bypassing to reduce API overhead and latency.
- Ensured production reliability through **Pydantic schema validation**, automated **retry/fallback logic**, and real-time **observability logging**.

### TRANSFORMER-BASED SENTIMENT ANALYSIS

[GITHUB](#)

- Developed a **Transformer NLP architecture** from scratch in **PyTorch** with **multi-head self-attention** and custom positional encodings.
- Achieved **91.8% accuracy (0.92 F1-score)** on a **560K review** dataset using Dropout and Early Stopping for model regularization.
- Built a scalable **NLP preprocessing pipeline** using **spaCy** and **Gensim** for tokenization, embedding generation, and large-text handling.

## PROFESSIONAL EXPERIENCE

### ACCENTURE, PUNE, INDIA

DEC 2020 - AUG 2023

Software Engineering Analyst (Promoted from Associate, 2022)

- Deployed **50+ production Java APIs** on Apigee and AWS, managing high-concurrency traffic for enterprise-scale consumer applications.
- **Reduced API errors from 25% to under 1%** within 2 months by debugging AWS Lambda/ElasticSearch logs and retiring faulty APIs.
- Built **ElasticSearch/CloudWatch dashboards** with anomaly detection alerting, reducing mean-time-to-detect from hours to minutes.
- **Collaborated** with cross-functional **global teams** on API performance tuning and secure governance (IAM, quotas, rate limits).
- **Led a 6-member team**, driving Agile/Scrum, CI/CD, and TDD practices to achieve on-time delivery for 12+ projects.

## SKILLS

**GenAI & Agents:** Multi-Agent Systems, LangGraph, LangChain, RAG, LLMs (Claude, GPT, Gemini), Anthropic SDK, Claude Code, OpenAI SDK, MCP.

**ML & Frameworks:** PyTorch, TensorFlow, Scikit-learn, Transformers, NLP, Hugging Face, Pandas, NumPy.

**Data & Retrieval:** Vector DBs (Pinecone, FAISS, ChromaDB), ElasticSearch, Data/ETL Pipelines, Postgres, MongoDB, Redis.

**Programming & Backend:** Python, Java, FastAPI, AsyncIO, Pydantic, SQL, REST APIs, Git, Postman.

**Cloud & LLMOps:** AWS (Lambda, S3, IAM, CloudWatch), GCP (Vertex AI, Cloud Run, Cloud Build), Docker, Terraform, LangSmith.

## CERTIFICATIONS

- [AWS Certified Cloud Practitioner](#) | AWS | 2023
- [Machine Learning Specialization](#) | Stanford/ DeepLearning.AI | 2023