

Vamsi Krishna Koppala

+1 (940) 999-8706

vamsikvk1234@gmail.com

[Vamsi-Krishna](#)

[Github](#)

[Portfolio](#)

PROFESSIONAL SUMMARY

Software Engineer with 3 years of experience in **AI/ML** and cloud-based application development across **GCP**, **Azure**, and **AWS**. Skilled in Python, REST APIs, LLM integration, and building scalable systems using **Docker** and **CI/CD** pipelines. Strong background in machine learning, semantic search, and intelligent automation workflows.

EDUCATION

Texas Tech University, Lubbock
Master of Sciences in Computer Science
GPA: 3.9/4.0

Jan 2024 – May 2025

EXPERIENCE

Software Engineer, ePATHUSA

Dec 2025 - Present

- Designed and deployed scalable AI-driven applications on **Google Cloud Platform (GCP)** using Cloud Run, Compute Engine, Cloud Functions, and Cloud Storage with secure IAM-based access control.
- Built end-to-end machine learning pipelines including data preprocessing, feature engineering, model training, evaluation, and cloud-based deployment using Python and **TensorFlow/scikit-learn**.
- Developed **LLM**-powered applications integrating semantic search, document processing, and **API**-based inference with containerized deployment using Docker.
- Implemented event-driven architectures using **Pub/Sub** and Cloud Functions to enable real-time data processing and automated microservices communication.
- Designed intelligent automation workflows using **n8n** (self-hosted with Docker) integrating **Webhooks**, **PostgreSQL**, **REST APIs**, and **LLM** services for structured data extraction and orchestration.
- Established **CI/CD** pipelines using Cloud Build and **Docker** to automate deployments, improve **scalability**, and reduce operational overhead.

Research Assistant, Texas Tech University

April 2024 - May 2025

Under the guidance of Professor Dr. Akbar Siami Namin

- Conducted applied research in **Sensitive Data Detection**, integrating **Speech-to-Text (STT)** systems with **Large Language Models (LLMs)** to identify sensitive information in unstructured text and audio.
- Engineered a full-stack **Python** web application combining **Flask**, **HTML5**, **CSS3**, and **Bootstrap** for real-time processing of audio and text data with sensitive information detection.
- Developed a dual-model sensitive data detection framework using **traditional pattern matching** (regex, Word2Vec similarity) and **LLM-based contextual analysis** powered by **Meta-Llama-3-8B-Instruct.Q8_0**.
- Implemented **speech recognition** pipelines utilizing **OpenAI Whisper** with **GPU acceleration** via **PyTorch**, enabling real-time transcription from microphone inputs and audio file uploads.
- Built a **Qdrant vector database** for efficient semantic similarity search of sensitive terms based on **Word2Vec** embeddings, enhancing approximate matching and classification performance.
- Designed and deployed an **adaptive feedback mechanism** that dynamically refines detection models based on user corrections, ensuring continuous learning and accuracy improvement. Integrated text **pre-processing techniques** such as **tokenization**, **stemming**, **lemmatization**, stop-word removal, and TF-IDF vectorization to enhance feature extraction.
- Integrated **Transformer-based NER** models like **BERT** fine-tuned and **DistilBERT model** for entity recognition and **Zero-shot classification** models (**Facebook/BART-large-MNLI**) to augment semantic understanding and context analysis.
- Conducted comparative analysis between traditional methods and LLM-based methods, improving detection accuracy and **reducing false positives** through **confidence scoring** and **overlap analysis**.

Software Engineer, DXC Technology

Nov 2022 - Dec 2023

- The team **executes VM (virtual machine)** and **physical server OS patching operations** for security purposes while maintaining **stability and industry compliance**.
- The software developer has **expertise in Azure VM OS patching** and **troubleshooting** while also **deploying automated deployment** of patches to enhance system performance.
- A total of **95%** of system complex issues on **CentOS** and **RHEL** systems were successfully resolved by my **troubleshooting efforts**. The system **performance** and **reliability** gain improvement through **automated patch deployment** methods.
- The professional maintains a specialization in **infrastructure management** where they excel at **complex infrastructure** design work and maintenance tasks in **Linux**, **VMware**, and **AIX** environments.

- My **expertise** includes detailed understanding of **virtualization** along with **Microsoft Azure** and other services such as **VMware ESXi**, **Hyper-V**, and **IBM AIX virtualization**.
- Experienced in **Azure Virtual Network infrastructure** creation as well as **Network Security Group deployment**, **VM Scaling (VMSS)**, and **Load Balancer configurations** to **optimize** cloud performance.
- Security policies are implemented with three core components: **SSH key management**, **role-based access control (RBAC)**, and **access control** methods.
- Worked with **DevOps tools** such as **Jenkins**, **Docker**, and **Git**, contributing to **CI/CD pipeline development**, containerized application deployment, and **version control** best practices.

PROJECTS

On Demand Professor Q&A Bot

- Deployed and configured the **Qdrant vector database** via **Docker**, establishing a **highly scalable** and **efficient vector storage system** for seamless integration with the **LLM-powered Q&A bot**.
- Integrated **GPT4ALL** as the primary **AI engine**, enabling **localized model training** on knowledge documents and **real-time internet-based query expansion** for comprehensive response generation.
- Designed an optimized document retrieval pipeline using **SentenceTransformers**, ensuring **accurate semantic embedding**, **indexing**, and **page-specific query resolution** to enhance user experience.
- Implemented an **API-driven architecture** to support multi-modal query processing, ensuring efficient retrieval and improved response accuracy for **domain-specific questions**.

Neuro-Symbolic Concept Revision Using Interactive Explanations

- Developed a pipeline leveraging Neuro-Symbolic Explanatory Interactive Learning (NeSy XIL) to improve model interpretability and accuracy by addressing **Clever-Hans behavior**, using **CLEVR-Hans datasets for robust evaluations**.
- Conducted extensive implementation and debugging to reproduce and enhance results from state-of-the-art research, **achieving up to 94.96%** accuracy on complex datasets by integrating **symbolic reasoning with neural network models**.
- Optimized feature selection using **attention-based explainability** methods, improving **model generalization** and reducing **overfitting** in vision-language reasoning tasks.

LLM Privacy Evaluation & Defense Framework (LLM-PBE)

- Implemented four attack vectors are Data Extraction, Jailbreak, Membership Inference, and Prompt Leakage. These are used to rigorously evaluate privacy vulnerabilities across multiple **LLMs (LLaMA2, Mistral, Gemma, Phi, Deepseek-R1)** using **Ollama** for efficient local model inference.
- Developed a **Python-based pipeline** to automate attack execution, logging, and accuracy analysis using **prompt engineering**, **semantic similarity metrics**, and **fuzzy matching** for precision measurement.
- Engineered **scrubbing and defensive prompting** modules to mitigate data leakage, achieving **100% mitigation** on select models and **significant accuracy reductions** on others.
- Optimized execution performance for local LLMs with **GPU offloading**, batch inference strategies, and runtime logging to support reproducibility and scalability on **12GB GPU environments**.

TECHNICAL SKILLS

Languages: Python, HTML, CSS, Javascript, Typescript

Database: MySQL, NoSQL, MongoDB, Qdrant

Platforms: Linux, MacOS, Visual Studio, Eclipse, Windows

Web Development: React, Flask, Bootstrap 5

Cloud Technologies: Azure(Load Balancer, VMSS, Blob Storage, VM's), AWS(Lambda, Redshift, S3, EC2)

Devops & CI/CD: Docker, Jenkins, Kubernetes

NLP Techniques: Named Entity Recognition (NER), Zero-shot Classification, Semantic Similarity Matching

AI/ML Frameworks: Hugging Face Transformers, SentenceTransformers, PyTorch, Scikit-learn

Version Control & Tools: Git, GitHub, GitLab

Visualization Tools: PowerBI, Tableau, Matplotlib, Microsoft Excel

Speech Recognition: OpenAI Whisper, Vosk

Web Technologies: REST APIs, FormData API

PROFESSIONAL CERTIFICATIONS

Az 204 - Microsoft Certified: Azure Developer Associate

Microsoft Certified: Azure Fundamentals (AZ-900)

Oracle Cloud Infrastructure 2022 Certified Foundation Associate