

# VISHNU SAI VARDHAN REDDY BASI REDDY GARI

vishnubr@umd.edu | +1 2272751335 | College Park, MD | [LinkedIn](#) | [Github](#) | [Website](#) | [Medium](#)

## EDUCATION

<b>University of Maryland, United States</b> <i>Master's of Science in Applied Machine Learning</i> Courses : NLP, Statistics & Probability	<i>Aug 2025 - Jun 2027</i>
<b>Vellore Institute of Technology, India</b> <i>Bachelor's of Science in Computer Science Engineering with specialization in Data Science</i> Courses : Data Structures & Algorithms, Deep Learning, Python, Database Management, Cloud Computing, Artificial Intelligence	<b>CGPA : 8.36/10</b> <i>Jul 2019 - Jun 2023</i>

## EXPERIENCE

<b>Machine Learning Engineer</b>   <b>Maruti Suzuki India Limited, Bengaluru, India</b>	<i>Jun 2023 - Jul 2025</i>
<b>Agentic Data Analyst</b>	
• Automated Data analysis by Architecting and Building a multi-agentic system using Agency Swarm and <b>Azure Open AI API</b> .	
• Engineered memory using Memgraph for short-term context graphs and ChromaDB+Neo4j for long-term <b>RAG</b> retrieval.	
• Enhanced <b>contextual reasoning</b> with LangChain pipelines, achieving <b>99%</b> retrieval accuracy and <b>30%</b> lower token usage.	
• Built a Streamlit interface to enhance user experience, reducing analysis time from <b>10 days to 10 minutes</b> and human effort by 90%.	
<b>Telematics Chatbot - Natural Language to SQL</b>	
• Designed a GenAI chatbot to democratize data access by translating natural language into optimized SQL for self-service analytics.	
• <b>Fine-tuned Flan-T5</b> using LoRA and PEFT on domain SQL/NL datasets, achieving <b>99%</b> accurate query-to-database mapping.	
• Optimized inference using ONNX Runtime and mixed-precision quantization, reducing GPU memory by <b>25%</b> and latency by 40%.	
• Drove business impact by boosting user engagement by <b>80%</b> and improving cross-team collaboration by 58%.	
<b>Machine Learning Intern</b>   <b>Maruti Suzuki India Limited, Bengaluru, India</b>	<i>Jan 2023 - Jun 2023</i>
<b>Accident Detection using LSTM AutoEncoders</b>	
• Developed an LSTM autoencoder to predict accidents by modeling sequential time-series data and detecting anomalies.	
• Engineered features using rolling statistics and window-based methods to capture temporal patterns and improve prediction accuracy.	
• Validated insights through EDA and hypothesis testing, achieving an <b>F1-score of 82%</b> for reliable deployment.	
• Created a benchmark dataset from processed accident data to enable hotspot dashboards and safety design improvements.	

## PROJECTS

<b>AI Event Reminder</b>   <i>Tesseract, MLX, Swift, Core ML</i>	<a href="#">Github</a>
• Launched an iOS app that extracts text from images via <b>OCR</b> and converts it into structured calendar events.	
• Fine-tuned SmoLLM-360M using <b>MLX</b> + Unified Memory for accurate event entity extraction.	
• Optimized on-device inference for <b>&lt;2s latency</b> , ensuring smooth, real-time iOS performance.	
• Boosted user engagement with <b>50%</b> adoption in the first month, simplifying event management for students.	
<b>Visual Search Engine</b>   <i>Pytorch, FAISS, Cuda, Streamlit, Hugging face, transformers</i>	<a href="#">Github</a>
• Developed a CLIP-powered multimodal search system to recommend visually similar apparel from uploaded images.	
• Engineered object-level detection extensions to identify <b>multiple entities (person,shirt)</b> for context-aware visual similarity retrieval.	
• Optimized embedding pipelines using FAISS, PyTorch, and <b>Transformers</b> , achieving <200 ms latency across 100K+ fashion images.	
• Deployed a scalable, containerized pipeline on AWS SageMaker with Streamlit UI, boosting search speed by 40%.	

## SKILLS

<b>Programming</b> :	Python, SQL, C++, C, CUDA, R
<b>Frameworks</b> :	Pytorch, MLX, Tensorflow, Langchain, Tesseract, Core ML, Agency Swarm, ADK, Crew AI, Unislot AI
<b>Generative AI</b> :	Agentic AI, RAG, LLM Finetuning, Model Profiling, Ollama, OpenAI, Gemini, Distributed Training,ONNX
<b>Databases</b> :	Vector DB - Chroma DB, Pinecone   RDBMS - MYSQL, DB2
<b>Cloud &amp; Web</b> :	Aws, Azure, GCP, Docker, Kubernetes, Git, CI/CD, ML flow, LLMops, Stream-lit, FastAPI

## CERTIFICATES

AWS - Certified Machine Learning Engineer Associate

Azure - Developing AI Applications on Azure