

Work Experience

AI Engineer

Mad Scientist

Sep. 2025 – Present

- Led development of a scalable **hybrid RAG** system with **adaptive semantic chunking**, parallelised document ingestion, and optimised embedding strategies across diverse content types.
- Built a robust processing architecture leveraging parallel processing, dynamic chunking thresholds, and metadata enrichment for high-precision, efficient retrieval.
- Implemented an **NL2Query with Filter Enrichment** using SLMs, and shipped a baseline system with strong performance; currently enhancing with new chunking methods, BM25 hybrid search, and FastAPI-based production API.

Software Engineer

Data Marshall

Apr. 2025 – Sep. 2025

- Debugged and **optimised** an **end-to-end** medical coding **pipeline**, stabilising parallel scraping by resolving file overwrite and variable collision issues, and **boosted accuracy** from **~50%** to **~99%** while reducing execution time from 90s to 2–4s (**~30x speedup**) through logic refactoring, table schema redesign, and chain-of-thought prompting.
- **Developed** a **custom** in-house **dataset** to **fine-tune LLMs** for automated medical coding, transitioning from regex-based + LLM hybrid systems to a fully LLM-driven pipeline.
- Architected a proof-of-concept **contract price estimation** system using Retrieval-Augmented Generation (**RAG**), combining custom chunking for tabular data, LangChain for text processing, and Chroma as the vector store.
- Explored **semantic web technologies** (RDF, RDFS, SPARQL) to prototype a **domain-specific ontology** for contract understanding and price inference using Neo4j and Cypher.

Intern

Sep. 2024 – Nov 2024

T- Hub

- Designed dynamic **data visualisations** to analyse startup performance metrics and maintained **Zoho CRM** including program trackers and cohort records to streamline operational workflows.
- Contributed to event planning and logistics for major initiatives such as **Mobility Demo Day** and **Innovate UK**, created marketing content (posters, videos), and supported inventory management, lab audits, and mentorship coordination.

Data Analyst

Jan. 2023 – July 2023

Panamax Infotech

- Developed a **robust API** to efficiently **fetch and decode data** from Ethereum nodes, streamlining data retrieval processes for **analysis**.
- Implemented advanced **anomaly detection mechanisms** within the API, ensuring **timely notifications** to the database admin in case of any **irregularities**.
- Leveraged API capabilities to **directly update** the company's database with **accurate and validated** blockchain data, contributing to **data integrity**.

Projects

Caption Kraft

Python, TensorFlow, Docker, Kubernetes, Terraform, GitHub Actions, GCP

- Engineered an advanced **image captioning** model using **CNN** and **LSTM** architecture with **Multi-Head Attention** and **Multi-Head Cross-Attention** layers, seamlessly **integrated** with **MLOps practices** like Terraform, Kubernetes, Docker, and GitHub Actions for efficient deployment and scalability on Google Cloud Platform.
- Orchestrated an interactive **Streamlit UI** for the **dataset creation tool**, leveraging Terraform for infrastructure provisioning, Kubernetes for container orchestration, and Docker for containerisation, ensuring **cost-effective** and **user-friendly image dataset creation**.

Vital Vector *Python, Ollama*

- Built a **custom Retrieval-Augmented Generation** (RAG) pipeline using domain-specific fitness and nutrition textbooks, featuring **semantic chunking, vector search**, and a **local LLM** for private, personalized health guidance.
- Designed an assistant that **tailors exercise** and **nutrition** plans **based on user** profile; integrated with **Gradio UI** and currently implementing **custom tool-calling** to fetch **real-time video content** for **enhanced user interaction**

Jarvis *Python, LangGraph, chainlit, Docker, Qdrant*

- Orchestrated a **multimodal WhatsApp AI assistant**, integrating with **LangGraph** for branch-based workflow, **LLaMA-3** for reasoning, Whisper, ElevenLabs, **Qdrant** and image processing pipelines for vision-to-language and image generation.
- **Containerized** and **deployed** the full pipeline using **Chainlit, Docker**, and **GCP Cloud Run**, supporting real-time, context-aware interactions with persistent memory, auto-scaling capabilities, and a fully orchestrated CI/CD flow.

Image Forensics *Python, TensorFlow*

- Created an advanced system using **CNNs** to **classify authentic images, AI-generated content, and edited visuals** with **96% accuracy**.
- **Improved system accuracy** by incorporating **synthetic images** through **DCGANs**, enhancing the **detection of AI-generated content**, and conducted comprehensive unit testing for **robust performance** across diverse scenarios.

FormFix *Python, OpenCV, Mediapipe*

- Built **real-time** exercise form **app** with OpenCV, Mediapipe, **CNNs** for accurate **landmark tracking**, instant **visual feedback**, and corrective prompts.
- **Enhanced** exercise safety by tailoring **pose estimation** enabling **immediate form correction** for **injury prevention**.

Education

Matrusri Engineering College B.E, Computer Science and Engineering, 7.73 CGPA	Nov. 2020 – Aug. 2024
NPTEL Minor in Data Science	Jan. 2023 – Apr. 2024
Narayana Junior College PCM, 96% TSBIE Boards	June 2017 – Apr. 2019
Little Flower High School SSC, 9.3 CGPA	June 2005 – Apr. 2017

Technical Skills and Competencies

- Languages: Python, C, C++, SQL, R, HTML, CSS
- Frameworks : TensorFlow, PyTorch, LangChain, Ollama, LangGraph
- Developer Tools: GitHub, VS Code, Jupyter Notebook, Docker, Terraform, Tableau, WandB, Llama.cpp, n8n
- Libraries: Pandas, NumPy, Matplotlib, OpenCV, Gradio, Streamlit, huggingface_hub, sklearn, nltk, uv, ruff
- Skills : Foundation Language Model Development, LLM Finetuning, Prompt Engineering

Extracurriculars

Competitions: 5+ hackathons (Offline + Online), 4th Place in the Regional Round of the GFG Solving For India
Recognition: Indian Army SSB Recommended (Medical Out); Microsoft Certified : Azure AI Fundamentals (AI900), Azure Data Fundamentals (DP900); Amazon ML Summer School'23; GCP Cloud Bootcamp'23; Meta Social Media Marketing Professional, Meta Certified Digital Marketing Associate; NPTEL Certified : Social Networks, Probability For Computer Science, Programming, Data Structures and Algorithms using Python; Problem Solving: 6 Stars @ HackerRank.

Volunteering: DeepLearning.AI Mentor and Course Tester, Contributor at Unify.AI, GFG Student Chapter Founder and Lead, Matrusri Developer Space Club Lead, IEEE Chapter Secretary.