

# Raju Gopi

+91 9353973699    rajugrai29@gmail.com    www.linkedin.com/in/raju-gopi-contact

## SUMMARY

### ML/Gen AI Engineer | 2+ Years of Experience

Results-driven **ML/Gen AI Engineer** with **2+ years of experience** in developing AI-driven solutions for business automation. Skilled in **LangChain, OpenAI, Google Gemini**, and **Hugging Face Transformers**, with expertise in **Retrieval-Augmented Generation (RAG)**, **code modernization**, and **predictive analytics**. Proven ability to integrate scalable AI models, optimize performance, and streamline workflows through intelligent automation. Passionate about delivering high-performance solutions that enhance business efficiency and innovation.

## EDUCATION

RNS Institute of Technology  
Bachelor of Engineering – 8.1 CGPA

08/2018 – 07/2022  
Bengaluru, India

## TECHNICAL SKILLS

**Programming Languages:** Python (Primary), Java, SQL  
**Generative AI & LLMs:** Paid and Open-Source LLMs (OpenAI, Google Gemini Pro, Llama2 and Hugging Face Transformers)  
**AI/ML Technologies:** Fine-tuning with custom data, Vector Embedding, NLP, Prompt Engineering, Apache Spark (PySpark), QLoRA (Quantized Low-Rank Adapters), LoRA (Low-Rank Adaptation), RAG (Retrieval-Augmented Generation)  
**Frameworks & Tools:** LangChain, FastAPI, GitHub, VS code, Jira, DataBricks  
**Deployment & Platforms:** GitHub Action, AWS, Docker, Kubernetes.  
**Visualization & Documentation:** PowerBI, Microsoft Office  
**Data Science & Miscellaneous:** Pandas, NumPy, Scikit-learn, Traditional ML Algorithms, MLflow, Pyspark, pytorch  
**Frontend Development:** HTML, CSS, Bootstrap

## EXPERIENCE

Infosys	08/2022 – Present Bengaluru, India
<ul style="list-style-type: none"><li><b>AI &amp; NLP Model Integration:</b> Developed and integrated AI models like Google GenAI, OpenAI, and Hugging Face Transformers for automating business processes, including natural language understanding, task generation, and data categorization</li><li><b>Deep Learning Framework Implementation:</b> Developed and optimized neural networks using PyTorch for NLP tasks, implementing custom architectures and distributed training workflows that improved model accuracy by 15%</li><li><b>Model Optimization &amp; Fine-tuning:</b> Fine-tuned Machine Learning Models like BART and GPT using domain-specific datasets, applied performance metrics such as sacrebleu for evaluation, and implemented LoRA/QLoRA techniques for efficient model adaptation</li><li><b>LangChain Framework Development:</b> Built advanced <b>LangChain</b> applications for document processing and conversational AI, integrating multiple <b>LLM providers</b> and implementing custom <b>prompt templates</b> for business-specific use cases</li><li><b>Traditional ML Algorithm Implementation &amp; Optimization:</b> Developed and deployed classification and regression models using <b>Scikit-learn</b>, including <b>Random Forest</b>, <b>SVM</b>, and <b>Gradient Boosting</b> algorithms for predictive analytics, implemented comprehensive model evaluation pipelines using <b>GridSearchCV</b> and <b>cross-validation</b> techniques for hyperparameter tuning, achieving 85% accuracy on business forecasting tasks</li><li><b>Containerized ML Deployment &amp; Orchestration:</b> Designed and deployed scalable ML models using Docker containerization and Kubernetes orchestration, implementing auto-scaling policies and CI/CD pipelines that reduced deployment time by 60%</li><li><b>End-to-End MLOps Pipeline Management:</b> Built comprehensive MLOps workflows using MLflow, automating model training, testing, deployment, monitoring, and drift detection with integrated A/B testing frameworks, achieving 95% pipeline success rate</li><li><b>Cloud-Native ML Infrastructure:</b> Architected production-ready ML solutions on AWS using SageMaker, implementing Infrastructure as Code with Terraform for reproducible deployments across multiple environments</li><li><b>Python Backend Development:</b> Developed scalable Python backends using FastAPI for web applications, creating high-performance APIs and robust integrations that improved data accessibility and processing efficiency by 40%</li><li><b>RAG &amp; Vector Database Systems:</b> Created scalable VectorDB integrations for retrieval-augmented generation (RAG) applications, enabling advanced search and knowledge retrieval capabilities</li><li><b>Apache Spark Integration:</b> Designed and implemented large-scale data processing workflows using <b>PySpark</b> and <b>Apache Spark</b>, handling petabyte-scale datasets for ML model training and feature extraction with distributed computing capabilities</li><li><b>Data Pipeline Automation:</b> Designed comprehensive data workflows using SQL for extraction and manipulation across relational databases, with optimized queries supporting ML pipeline integration</li><li><b>Exploratory Data Analysis (EDA):</b> Conducted comprehensive data exploration, handling structured and unstructured datasets to identify trends, correlations, and anomalies, ensuring high-quality input for ML models.</li><li><b>Feature Engineering:</b> Applied feature engineering techniques to optimize data for <b>predictive modeling</b> and performance improvements.</li><li><b>Version Control &amp; Collaboration:</b> Utilized <b>Git</b> and <b>GitHub</b> for version control and collaboration across teams, ensuring smooth development workflows and effective codebase management.</li><li><b>FMECA Analysis:</b> Led <b>Failure Modes and Effects Criticality Analysis (FMECA)</b> for gas turbines using <b>Excel Power Query</b> and <b>Power BI</b>. Automated data extraction with <b>Python</b>, reducing processing time by over 90%. Developed dashboards to highlight critical failure modes.</li></ul>	

## ACCOMPLISHMENTS

- Insta-Award- Client Appreciation:** Designed and developed intuitive and visually compelling Power BI dashboards to present complex, large-scale data in a clear, accessible, and actionable format, enabling stakeholders to easily interpret and make informed decisions.

