

# Aryasomayajula Ram Bharadwaj

ram.bharadwaj.arya@gmail.com | +91-9108832338 | Bengaluru, India | github.com/rokosbasilisk

## PROFESSIONAL SUMMARY

Senior ML Engineer with 6+ years developing production AI systems and leading engineering teams. Background in AI safety research with published work on model interpretability and LLM inference optimization. Experienced architecting scalable solutions for conversational AI and agentic workflows. Strong expertise in MLOps, distributed systems, and technical leadership.

## EDUCATION

**Bachelor of Technology - Electronics and Communications Engineering**

2015 - 2019

GMR Institute of Technology, Andhra Pradesh

## SKILLS

**Programming:** Python, Scala, Java, R, SQL  
**ML/AI:** PyTorch, JAX, Transformers, LangChain, AutoGen  
**DevOps:** Docker, Kubernetes, CI/CD, GCP  
**Systems:** Redis, Kafka, Microservices

## PROFESSIONAL EXPERIENCE

**Associate Technical Architect - MLOps**

Nov 2024 - Present

Quantiphi Analytics, Bengaluru

- Designed and implemented an AI agent system for automated replies for a telecom client's sales chatbot
- Experimented with LangGraph agent architectures to improve controllability and modularity in conversational flows
- Revamped and automated the scheduled RAG data ingestion pipeline and improved retrieval accuracy by 20%
- Improved retriever latency from 6 seconds to 1.4 seconds through extensive code refactoring

**AI Resident - Lossfunk AI Residency**

April 2025 - May 2025

Lossfunk AI Residency (Remote)

- Selected among 100 applicants for elite 6-week intensive AI residency program with 10 researchers
- Developed STU-PID (Steering Token Usage with PID Control), a novel activation steering technique to reduce redundant reasoning tokens in LLMs
- Conducted experiments on reasoning models like DeepSeek-R1-Distill-Qwen-1.5B with dynamic activation interventions, achieving 32% token reduction and improved reasoning accuracy on GSM8K benchmark
- Published research findings: Steering Token Usage with PID Control

**Senior Developer - Innovation & Development Labs**

June 2019 - Nov 2024

Musigma Business Solutions, Bengaluru

- Led development of multiple high-impact projects in LLM operationalization, automated trading, and MLOps
- Specialized in backend development, DevOps, and ML engineering across various domains
- Received multiple recognitions including Impact Awards and Star Performer of the Team

# KEY TECHNICAL PROJECTS

---

<b>LLM Agent Platform</b> <i>Team Lead, ML Engineer, Backend Developer</i> <ul style="list-style-type: none"><li>Designed and implemented semi-autonomous data analysis platform using Microsoft AutoGen framework</li><li>Developed automated prompt optimization strategies and integrated RAG support</li><li>Built evaluation framework using ensemble of locally hosted LLMs</li><li>Successfully deployed for two major clients, assisted client teams in creating use cases</li></ul>	Dec 2023 - Nov 2024
<b>High-Velocity Trading Platform</b> <i>Team Lead, Backend Developer, DevOps</i> <ul style="list-style-type: none"><li>Refactored backend code and migrated deployment from bare-metal to Kubernetes</li><li>Implemented automated hyperparameter search for ARIMA models</li><li>Rewrote legacy trade-signal generation from R to Scala using Akka framework</li><li>Enabled near real-time metric calculation and portfolio visualization</li></ul>	2021 - Dec 2023
<b>ML Model Operationalization Platform</b> <i>Backend Developer</i> <ul style="list-style-type: none"><li>Developed automatic retraining pipelines for image classification models</li><li>Designed Java microservices for creating and serving Jupyter notebooks</li><li>Engineered ML model deployment service with canary and blue-green deployment strategies</li></ul>	2019 - 2021

# RESEARCH & PUBLICATIONS

---

<b>Steering Token Usage with PID Control</b> <i>Lead Author</i> <p>Novel technique reducing computational overhead in LLMs through activation steering with 32% token reduction on GSM8K.</p>	June 2025
<b>Understanding Hidden Computations in Transformer Language Models</b> <i>Lead Author</i> <p>Investigated internal mechanisms of chain-of-thought reasoning and developed interpretability methods for LLM reasoning.</p>	August 2024

# SELECTED PERSONAL PROJECTS

---

<b>Wiserank.io - AI-Powered Research Discovery</b> <i>Creator &amp; Full-Stack Developer</i> <ul style="list-style-type: none"><li>Built research paper search engine with relevance and 'creativity' ranking algorithms</li><li>Implemented automatic citation generation feature for uploaded manuscripts.</li><li>Deployed scalable architecture serving researchers with enhanced paper discovery and citation workflows</li></ul>	June 2025
---	-----------

# AWARDS & RECOGNITIONS

---

<b>AI Alignment Awards - Winner</b> <i>AI Safety Research Competition</i> <p>Selected among 118 global entries for winning research proposal on "goal misgeneralization" in AI systems.</p>	July 2023
<b>Honorable Mention - Eliciting Latent Knowledge</b> <i>Alignment Research Center</i> <p>Recognized for innovative approach to open research problem in AI interpretability and knowledge extraction.</p>	March 2022
<b>Bronze Medal - Build-on-Redis Hackathon</b> <i>Redis Labs</i> <p>Developed text-to-code search tool using CodeBERT embeddings and Redis Stack for private repository indexing.</p>	February 2021
<b>Excellence Awards (8x)</b> <i>Musigma Business Solutions</i> <p>6 SPOT (Star Performer) awards and 2 Impact Awards for technical leadership, innovation, and delivery excellence.</p>	2019 - 2023