

Nirajan Paudel

 [paudelnirajan](#) |  [paudelnirajan.github.io](#) |  nirajan.paudel@colorado.edu |  +1 (720) 595-8207

EDUCATION

University of Colorado Boulder

Master of Science in Computer Science

Relevant Coursework: Data Center Scale Computing, Object-Oriented Analysis & Design

Boulder, CO

May 2027

Institute of Engineering, Tribhuvan University

Bachelor of Engineering in Electronics, Communication & Information Engineering

Pokhara, Nepal

2020–2024

WORK EXPERIENCE

Machine Learning Intern

PalmMind Technology

March 2025 – May 2025

- Architected production-grade agentic RAG chatbots for EV dealers and insurance clients using OpenAI GPT-4, migrating workflows from LangChain to LangGraph to enable advanced tool calling and conversational AI capabilities with orchestrator-worker patterns and parallel processing.
- Implemented LLM routing strategy with GPT-4 for complex queries and GPT-3.5-turbo for simple classification, reducing API costs for 2 production chatbots while maintaining response quality through prompt engineering and context optimization.
- Optimized retrieval pipeline by integrating web scraping, Tesseract OCR, and Redis caching for 500+ customer inquiries, reducing redundant data fetches and improving context accuracy for insurance policy and vehicle specification queries.
- Applied best practices from Anthropic's "Building Effective Agents" and LangChain documentation to design robust agentic systems with improved reliability, error handling, and conversational memory.

Teaching Assistant

Pashchimanchal Campus, Tribhuvan University

Sep 2024 – Feb 2025

- Instructed 90+ students in C programming fundamentals and Information Systems, designing hands-on lab exercises covering cloud computing, distributed systems, neural networks, and data mining with comprehensive code reviews.
- Mentored 15+ semester-end projects, resulting in 85% of students demonstrating proficiency in algorithm design and implementation.

PROJECTS

Music Separation as a Service (MSaaS)

University Course Project

Oct 2025 – Nov 2025

- Architected scalable microservices on Google Kubernetes Engine (GKE) for AI-driven music source separation, processing MP3 uploads into four instrumental stems (vocals, drums, bass, other) via asynchronous pipeline with horizontal pod autoscaling.
- Migrated storage from self-hosted MinIO to Google Cloud Storage (GCS), leveraging 99.99% durability guarantee and eliminating single points of failure while significantly reducing operational overhead.
- Engineered RESTful Flask API with Redis-based job queue decoupling from Demucs ML workers, enabling non-blocking API responses for inference tasks averaging 2-5 minutes and supporting concurrent processing of multiple songs.
- Implemented Kubernetes resource management (6Gi RAM limits) and dynamic JavaScript frontend with real-time polling, preventing OOM pod evictions and delivering seamless UX for long-running background tasks.

Zenco - AI-Powered Code Analysis Tool

Open-Source Python Package (Published on PyPI)

Nov 2025 – Present

- Developed production-ready CLI tool supporting 5 languages (Python, JavaScript, Java, Go, C++) with Tree-sitter AST parsing and multi-provider LLM integration (Groq/OpenAI/Anthropic/Gemini), automating code documentation and analysis workflows.
- Engineered execution priority optimization algorithm using dead code detection to skip unnecessary LLM API calls, reducing token consumption and operational costs while maintaining analysis quality.
- Refactored monolithic codebase using Strategy and Factory design patterns, reducing core module size by 70% (2,057 to 603 lines) and establishing modular architecture for extensible language support and maintainability.
- Published to PyPI with automated GitHub Actions CI/CD pipeline for pytest-based multi-platform testing (Linux, macOS, Windows), implementing comprehensive error handling and mock testing capabilities.

AI-Powered Question Assistant (IOE-GPT)

Personal Project

June 2025 – July 2025

- Built semantic search system for IOE programming questions using LangChain, Groq LLaMA-3, Milvus vector database, and sentence-transformer embeddings (all-MiniLM-L6-v2), deployed as FastAPI microservice with Redis state management and Docker containerization.

TECHNICAL SKILLS

Languages

Python, C/C++, Java, SQL, Bash

ML/AI Frameworks

PyTorch, TensorFlow, Scikit-learn, LangChain, LangGraph, HuggingFace Transformers, OpenAI/Anthropic APIs

ML Techniques

RAG, Transformers, CNNs, RNNs, NLP, Computer Vision, Vector Embeddings (Sentence-Transformers)

Cloud & DevOps

GCP (GKE, GCS), Kubernetes, Docker, Redis, GitHub Actions, CI/CD Pipelines, Microservices

Tools & Databases

FastAPI, Flask, Git, Milvus, PostgreSQL, Tree-sitter, pytest, Jupyter, Pandas, NumPy