# Arnav Panigrahi

📞 +1 951-347-2065   ✉ arnav.panigrahi@gmail.com   🏠 Riverside, US
💼 Arnav Panigrahi   🔗 notquitethereyet   🌐 arnavpanigrahi.com

## Professional Summary

Results-driven Machine Learning Engineer with an MSCS and 2+ years of experience in the end-to-end development of AI-powered applications. I have designed and deployed scalable systems using Python, LangChain, and PyTorch, with a strong focus on Generative AI, RAG, and NLP solutions deployed on GCP and AWS

## Experience

**AllCheer (Fateh and Sudha Inc), San Francisco, CA**

*Software Developer (Full-time)*                                                      *Mar 2025 – Present*

- Engineered a fully automated Release of Information (ROI) system using a **LangChain AI agent** to orchestrate document generation, e-signature collection, and secure data sharing, saving **$20K/year** and **reducing processing time by 60%**.
- Deployed a HIPAA compliant production-grade **FastAPI backend** to ingest real-time data from Rethink BH, building out an observability stack on **Google Cloud Run** that **cut costs by 90%** and eliminated 24-48h latency.
- Developed an internal dashboard integrating a GPT-based NLP API for **named-entity recognition (NER)** and text summarization to structure unformatted therapist notes, improving **data accuracy to 92%**.

**AllCheer (Fateh and Sudha Inc), San Francisco, CA**

*Software Developer Intern*                                                         *Jun 2024 – Dec 2024*

- Built a distance matrix generation system to schedule trips for the staff using **Google Route / Distance Matrix API**.
- Prototype the ROI automation using (**Make.com**) and validate workflows with legal/operations teams.

**Pinnacle Consulting LLC (PCON Utilities Pvt Ltd), Bhubaneswar, India**

*Junior Software Developer*                                                          *Jan 2022 – Aug 2023*

- Integrated **ArcGIS** mapping with REST APIs and SQL Server for a $75M utility provider in North Carolina to visualize infrastructure and sync real-time updates from **mobile field apps**, reducing oversight errors by **30%**.
- Delivered **3 ERP systems** using **Angular** and **Node.js** for internal teams to manage **timesheets**, **communication**, and **KPI reporting**; implemented **JWT authentication**, automated invoicing, and messaging, improving operational efficiency by **25%**.

## Technical Skills

- **Programming & Data Science:** Python, NumPy, Pandas, JavaScript, TypeScript
- **AI & Machine Learning:** PyTorch, TensorFlow, LangChain, LangGraph, RAG, LLMs, Prompt Engineering, NLP
- **Data & Backend:** FastAPI, Flask, RESTful APIs, NumPy, Pandas
- **Cloud & Deployment:** AWS, GCP, Docker, Kubernetes
- **Databases:** PostgreSQL, MongoDB, Redis, Pinecone

## Projects

- **bedtime.ai – Multi-Modal AI Storytelling Platform**
  A multi-modal AI pipeline that converts drawings to narrated stories; utilized **EfficientNet** for image feature extraction, a **fine-tuned phi-3 LLM** for story generation, and **coqui-ai TTS** for voice cloning.

- **Agentic Job Tracker & Helper**
  a stateful, multi-agent system using **LangGraph** and **GPT-4o** to automate job search and application tracking; implemented **NLP-driven** resume tailoring and managed state using **FastAPI** backend and **Supabase**.

- **RAG for Document Intelligence**
  Built a RAG pipeline for document Q&A; implemented semantic chunking and embedding generation to index PDFs in **Pinecone**, enabling a **LLM** interface with 35% higher relevance over baseline keyword search.

## Education

**University of California, Riverside**

*Master of Science in Computer Science*                                                    *2023 - 2025*
GPA: 3.61