

# SAYAK BANERJEE

Email: sayakbanerjee022@gmail.com

Phone Number: (412) 844-1988

[LinkedIn/sayak-banerjee/](#)

[Google-Scholar/sayak-banerjee/](#)

## EDUCATION

### Carnegie Mellon University

Master of Computational Data Science (MCDS) | QPA – 4.00/4.00

Coursework – 11-667: Large Language Models, 11-642: Search Engines, 10-601: Machine Learning, 15-619 Cloud Computing

Teaching Assistant – Applications of NL(X) and LLMs, Fundamentals of Operationalizing AI, Python for Data Science

Pittsburgh, PA

Dec 2025

### Vellore Institute of Technology

B. Tech, Electronics and Communication Engineering | CGPA – 9.41/10

Vellore, India

May 2022

## SKILLS

**Languages & Tools:** Python, C++, PySpark, SQL, Snowflake, Git, Apache Airflow, MLFlow, Kafka, Docker, Kubernetes, AWS

**Libraries and Frameworks:** PyTorch, Huggingface, Tensorflow, NumPy, Pandas, SciPy, Scikit Learn, Matplotlib, SpaCy

**Technologies:** Machine Learning, Natural Language Processing, Large Language Models, Prompt Engineering, Search Engines & Information Retrieval, RAG, Deep Learning Systems, Convolution Neural Networks, Cloud Computing, CI/CD, ETL Pipelines

## EXPERIENCE

### Paylocity

Schaumburg, IL

#### Data Science Intern

May 2025 – Aug 2025

- Collaborating with the Search team to enhance an enterprise chatbot solution utilizing RAG architecture focussing on queries related to the Paylocity portal, IRS-compliance and company handbooks.
- Implemented a dynamic fallback mechanism to switch to Bedrock embeddings during OpenAI latency spikes or outages, after thorough experimentation and benchmarking – **improving the AI assistant's SLA by ~2% for our customers. This work laid the foundation for upcoming load balancing and model routing enhancements.**
- Improved the existing re-ranking pipeline by removing BM25, which resulted in an average **reduction of 5k tokens** across user queries and **reduced latency by ~2.5 seconds.**
- Enhanced the ranking pipeline by finetuning multi-match keyword-based search alongside **HNSW ANN search**, creating a hybrid ranking pipeline, **improving recall by 0.6% and acceptable answer proportion by 1.44%.**
- Developed a workflow to detect and extract tables from table heavy pdfs using **table transformers and GPT-4o mini** and stored them in vector databases in an uniform markdown format for **enhanced retrieval** and query performance.

### Acuity Knowledge Partners

Gurugram, India

#### Associate (Data Engineer)

Dec 2022 – July 2024

#### Senior Analyst

July 2022 – Nov 2022

- Orchestrated and analyzed **alternative financial data** for a US-based Hedge Fund covering various investment portfolios including consumer, healthcare, technology, and e-commerce.
- Spearheaded the integration of Apache Airflow with existing ETL pipelines – creating event-based and time-based triggers, reducing manual efforts to monitor jobs **by more than 60%.**
- Designed ETL pipelines and implemented a dynamic data warehouse utilizing Python, PySpark, Pandas, Airflow, AWS, SQL, and Snowflake to store aggregated KPI data from multiple datasets in a SQL database, empowering PMs **to generate financial models by automating the extraction of over 1 million data points.**
- Implemented a full-fledged analytics pipeline involving product tagging and KPI visualizations using PySpark, Pandas, and Dash, **achieving ~40% reduction in time taken to generate KPI reports.**

## PROJECTS

### Carnegie Mellon University

Pittsburgh, PA

#### Optimizing QA performance with Dense Passage Retrieval (DPR) and RAG

April 2025

- Developed a modular multi-stage response generation pipeline (DPR → Reranking → RAG) integrating dense passage retrieval on FAISS-indexes and Flan-T5 based QA agent, enabling end-to-end neural retrieval and answer generation.
- Integrated Lucene for efficient lexical indexing and FAISS with co-condenser architecture for DPR.
- Conducted extensive experiments combining advanced prompting techniques (CoT, Persona prompting), Learning to Rank features with diverse reranking pipelines (BM25, LTR-based SVMRank/Coordinate Ascent/ListNet, BERT-*n*), **achieving a 25% improvement in exact match (EM) scores** over the BM25 baseline on the SQuAD evaluation dataset.

### Twitter Recommendation Microservice (A Cloud-Native Scalable Solution)

April, 2025

- Engineered a distributed ETL pipeline using PySpark on Databricks to process ~1TB of raw Twitter JSON data; computing scores based on user interactions, keywords, and hashtags for user-recommendation.
- Deployed a fault-tolerant production-grade microservice with MySQL back-end (optimized using denormalization and indexing) on a self-managed Kubernetes cluster (AWS ECS) using docker, helm charts; implemented automated CI/CD pipelines with GitHub actions. The microservice achieved 20K+ RPS over a 3-hour live load test.

### Text2Utility – Talk with your Mobile Apps

Dec, 2024

- Developed an AI agent for task automation from natural language inputs – involving intent classification, structured code generation followed by LLM tool-use to help user's complete tasks such as **booking an Uber, finding restaurant tables.**
- Conducted extensive model evaluation, implementing PEFT (LoRA) along with few-shot learning to benchmark various foundation models from Hugging Face, quantifying task-specific accuracy metrics across different user intents.