

# Rahul Yadav

+1(716)238-0835 | [rahulyad@buffalo.edu](mailto:rahulyad@buffalo.edu) | [linkedin.com/in/rahul-yadav-704](https://www.linkedin.com/in/rahul-yadav-704) | [Github](#) | [Website](#)

## EDUCATION

<b>University at Buffalo, Buffalo</b> <i>Master of Science in Computer Science and Engineering</i> Coursework: NLP, Computational Linguistics, Pattern Recognition, Machine Learning, Algorithms Design	Aug. 2022 – Dec. 2023 New York, US GPA: 3.82 / 4
<b>Indian Institute of Technology (Banaras Hindu University), Varanasi</b> <i>Bachelor of Technology in Electronics Engineering</i>	Jul. 2014 – May 2018 Uttar Pradesh, India

## EXPERIENCE

<b>Hilabs</b> <i>Senior Data Scientist   Applied Data Science Team</i>	Feb. 2024 – Present Bethesda, MD
---	-------------------------------------

### Legal Document Processing with GenAI

- Led the development of an advanced RAG-based Generative AI service aimed at processing 1000s of large legal documents efficiently for entity retrieval and verification of party agreements
- Leveraged both open-source LLMs (llama, mistral-ai) and enterprise-grade LLMs (GPT-4, Claude) to construct the solution while maintaining a paramount focus on data privacy
- Orchestrated the deployment of a GPU pipeline on Amazon EMR, ensuring inference times of less than 5 seconds using Docker/Terraform. Conducted rigorous load testing with 20 concurrent users to evaluate performance accurately
- Successfully completed a Proof of Concept (POC) which demonstrated the solution's effectiveness and potential for widespread adoption. Currently, integrating the solution for 3 clients, showcasing its impact and scalability

<b>Apexanalytix</b> <i>Data Scientist   ML Engineer Intern   Analytics Team</i>	May. 2023 – Feb. 2024 Remote, New York
--	---

### Generative Knowledge-Specific Chatbot

- Developed advanced Retrieval Augmented Generation(RAG) Chatbot with LLM for intelligent knowledge access across teams
- Engineered custom document chunking and indexing pipeline using LangChain and leveraged MiniLM embeddings
- Designed efficient RAG pipeline with ChromaDB, integrating MMR scoring and Azure OpenAI LLM to deliver precise response
- Achieved 89% approval in human evaluations and integrated technology into 12 internal and 18 external client applications

<b>Oracle</b> <i>Senior Application Engineer(ML)   Oracle Financial Machine Learning Team</i>	Sept. 2020 – Aug. 2022 Bengaluru, India
--	--

### Preemptive Anomaly Prediction in Corporate Billing

- Implemented in-memory multivariate anomaly prediction system for corporate billing, addressing monthly billing challenges
- Leveraged Oracle in-database ML for Semi-Supervised classification with both local and global model explainability
- Optimized service with indexing and parallelism, processing 1.2M bills and 5M segments in 20 mins with 92% precision
- Integrated system with services of Oracle Revenue Management and Billing (ORMB), granted innovation patent at USPTO

<b>Wipro Limited</b> <i>Project Engineer(AI)   AI Research Team</i>	June 2018 – Sept 2020 Bengaluru, India
--	---

### Chatbot Services for Employee Helpline Portal's Ticketing System

- Developed Employee Helpline Portal Chatbot with effective retrieval of historical ticket resolutions to enhance user support
- Improved user experience with precise responses using query intent classification and BERT-powered semantic search
- Achieved 79% accuracy score, reduced user wait times from 18 to 4 minutes, and decreased the workload for support agents

## PATENTS and PUBLICATIONS

<b>Paper Published:</b> Virtual Conversation with Real-Time Prediction of Body Moments/Gestures <i>Gopichand Agnihotram, Rajesh Kumar, Pandurang Naik, <b>Rahul Yadav</b></i>	ICMLIP 2019 <a href="#">springerLink</a>
<b>Patent Granted:</b> Method And System For Multimodal Analysis Based Emotion Recognition <i>Inventors: <b>Rahul Yadav</b>, Gopichand Agnihotram</i>	16/795840 [PDF] Applicant: <b>Wipro Limited</b>

## SKILLS

<b>Languages</b> : Python, Java, C++, C, PLSQL, SQL, MongoDB
<b>Technologies</b> : Retrieval Augmented Generation(RAG), Large Language Models(LLMs), Natural Language Processing(NLP), Machine Learning, Deep Learning, Web Services, Data Structures, Algorithms, Prompt Engineering, Indexing, Quantization
<b>Cloud</b> : Amazon Web Services (AWS), Azure, Google Cloud Platform (GCP), Oracle, Atlas, OpenAI
<b>Frameworks and Libraries</b> : Pytorch, Tensorflow, MLFlow, Kafka, SentenceTransformers(embeddings, re-rankers), Databricks, Snowflake, PySpark, Langchain, LlamaIndex, VectorDBs(ChromaDB, Faiss, elasticsearch), MLFlow, Git, Jenkins, Big Data, Hadoop, Flask, Docker, Kubernetes, Pandas, Streamlit, Flask-RESTful, FastAPI, YOLO, Django, XGBoost, GAN, ActiveMQ, Springboot