

# Ravi Raj Kumar

📞 +1-216-456-7060 ✉ Email 🔗 LinkedIn 🐙 GitHub 🌐 Portfolio 📍 New York

## EDUCATION AND HONORS

Case Western Reserve University

Cleveland, Ohio, Expected May 2025

Master of Science in Computer Science

- Coursework: Machine Learning, Computer Vision, Robotics, High Performant Systems for AI, Probabilistic Graphical Models, Analysis of Algorithms, Computer Networks.
- Observer: Statistical Natural Language Processing, Deep Gen Models, Quantum Computing, Reinforcement Learning.

## TECHNICAL SKILLS

- **Languages & Frameworks:** Python, SQL, Java, C++, JavaScript, TensorFlow, PyTorch, HuggingFace Transformers, FastAPI, Django, Streamlit, Gradio, LangChain, LangGraph, FAISS, Pinecone, VectorDBs
- **Machine Learning & RL:** Classification, Regression, Clustering, SVM, Random Forest, CNN, RNN, Transformers, Policy-Gradient methods (REINFORCE, PPO, GRPO, A2C), Actor-Critic (DDPG, SAC), Multi-Agent RL, Fine-Tuning (LoRA, QLoRA, PEFT)
- **Generative AI:** Diffusion Models, LLMs, multimodal Transformers, GANs, VAEs, Image Generation, Inpainting, Upscaling
- **MLOps & Infrastructure:** MLflow, Docker, Kubernetes, Helm, AWS (EC2, S3, ECR), GCP, Azure, CI/CD (GitHub Actions, GitLab, Jenkins), Monitoring (Prometheus, Grafana), ONNX Runtime, Quantization, Scalable Model Serving

## PROFESSIONAL EXPERIENCE

Tata Consultancy Services

Hyderabad, India [October 2019 - November 2023]

**Machine Learning Engineer (NLP Model Development & MLOps Integration - Banking and Finance domain)**

- Built robust and scalable end-to-end ML pipelines for several finance and Banking Member Complaint Systems on cloud as-well-as on-prem with components like data ingestion, data validation, feature engineering, model training, prediction, and monitoring.
- Implemented data ingestion and data validation components in the pipelines for large-scale data sources like Hadoop, Snowflake, and MongoDB and validated the output artifacts for robustness.
- Leveraged advanced NLP tokenizers, such as **BytePair Encoding (BPE)** and **SentencePiece** for tokenization, trained and finetuned several transformer-based models like **BERT**, **RoBERTa** on tasks such as Complaint Categorization and Prioritization, Named Entity Recognition (NER) for automated information extraction, and Complaint Severity Prediction.
- Built **CI/CD** pipelines with **Jenkins** and integrated **MLflow** for experiment tracking, versioning, and model registry; deployed models using **Docker** and **Kubernetes** with Helm for streamlined, version-controlled production deployment.
- Implemented automated model validation and **drift detection** to ensure reliability, and leveraged **Prometheus** and **Grafana** with custom metrics and alerts for proactive performance and infrastructure monitoring, and efficient issue detection in production.

**Data Engineer (Bank Member Complaints)**

- Catalogued and registered data assets to enforce security policies and enable seamless migration to the DL2 cloud zone. Designed both General-Purpose and Subject-Area marts, and implemented a Data Quality Engine to validate extracts before loading.
- Leveraged the Hera framework, UNIX shell scripting, and advanced SQL in Hive, Snowflake, and DBT to streamline data ingestion and transformation workflows.
- Developed scripts and jobs to apply complex business logic and route outputs to UNIX filesystems, on-premises databases (Netezza, Db2, Hive), and AWS Snowflake—ensuring reliable, end-to-end data delivery.

## RESEARCH EXPERIENCE

**GRPO-Based Efficient Retrieval Optimization in RAG Using Reinforcement Learning**

Supervising Professor: : Dr. Soumya Ray

- Developed a modular **reinforcement learning** framework for **RAG** that trains the search agent via **Group Relative Policy Optimization (GRPO)**, using a **novel reward signal** based on the gain in answer quality from retrieved evidence, avoiding brittle end-to-end fine-tuning and achieves superior performance with lower data requirements across general and medical QA tasks
- This zero-shot approach **shows good performance on other QA domains**, despite training only on general QA, suggests that RL search skills **generalize more reliably** than generation-tuned approaches
- Training only the **retriever agent** significantly outperforms end-to-end optimized RAG systems like **Search-R1**, achieving up to **+8 point gain in QA accuracy** across general and medical domains using a **7B model** with **70× less training data**, by applying **LoRA-based fine-tuning with GRPO**.

**DEEP GENERATIVE MODELS TO ENHANCE SEMI SUPERVISED LEARNING**

[\[Github Link\]](#)

Supervising Professor: : Dr. Soumya Ray

- Explored the integration of **deep generative modeling** approaches to semi-supervised learning by replicating and analyzing **Kingma et al.'s M2 VAE** framework, focusing on classification under limited-label regimes using MNIST and CIFAR-10.
- Conducted a rigorous mathematical and experimental analysis of the **Evidence Lower Bound (ELBO)** within a variational inference framework, identifying and resolving **three critical issues in the paper**: entropy penalization for sharper decision boundaries; mutual information maximization to strengthen input-label coupling; and smoothed-label integration into the classification loss.

- Achieved a **+4%** test accuracy improvement on MNIST and **+2.5%** on CIFAR-10 over the M2 baseline, along with a **15% reduction in prediction entropy**, demonstrating better confidence and calibration in semi-supervised settings.

## MULTIMODAL TRANSFORMER FOR IMAGE-CONDITIONED TEXT GENERATION

[\[Github Link\]](#)

*Independent Research Project*

- Designed and implemented a multimodal transformer integrating a SigLIP-style Vision Transformer with a Gemma-based causal decoder for image-grounded generation tasks such as captioning and visual question answering.
- Engineered autoregressive decoding with KV caching and Rotary Positional Embeddings, reducing inference latency by **40%** per token on long sequences.
- Built a PaLI-inspired tokenizer pipeline with image-token prefixing and robust image normalization, improving input alignment and reducing token mismatch errors by **23%**.
- Enabled diverse generation using temperature-scaled top-p sampling; achieved a **+3.7 BLEU-4 score improvement** over greedy decoding on benchmark prompts.

## Latent Diffusion Transformer for Text-Conditioned Image Synthesis

[\[GitHub\]](#)

*Independent Research Project*

- Developed a complete multimodal pipeline integrating CLIP-based language encoder, variational autoencoder (VAE), and U-Net-based denoising diffusion model for high-fidelity text-to-image generation.
- Achieved **4× faster inference** over pixel-space models by operating in latent space ( $64 \times 64$  vs.  $512 \times 512$ ), reducing memory and compute by over **85%**.
- Implemented **Classifier-Free Guidance** to enhance prompt alignment, improving semantic accuracy of generated images by **8%** (measured via CLIPScore).
- Enabled multimodal capabilities: **text-to-image**, **image-to-image translation**, and **inpainting**, through unified conditional diffusion framework.

## 3D POINT CLOUD SEGMENTATION USING 2D IMAGE SEGMENTATION

[\[Github Link\]](#)

*Supervising Professor: : Dr. Yu Yin*

- Developed a novel 3D point cloud segmentation framework leveraging state-of-the-art 2D image segmentation models (OneFormer) and a voting-based approach to project 2D semantic and panoptic labels onto 3D point clouds, achieving real-time segmentation with reduced computational overhead.
- Utilized RGB images, depth maps, and LiDAR data captured with the iPhone 13 Pro, integrating segmentation masks generated by OneFormer with a voting mechanism to accurately transfer semantic labels to 3D point clouds.
- Achieved segmentation accuracy of 96.5%, matching PointFormer, while significantly reducing computational overhead and memory usage demonstrating the model's scalability and efficiency.

# PROJECTS

## MULTI-AGENT MEDICAL APPOINTMENT SYSTEM

[\[Github Link\]](#)

- **Built a modular multi-agent AI system using LangChain and LangGraph**, implementing a supervisor-agent architecture to route user queries to specialized agents for doctor information and appointment management.
- **Engineered a full-stack solution with FastAPI for scalable backend services and Streamlit for real-time, interactive user interfaces**, enabling seamless doctor appointment management via natural language queries.
- **Improved query-to-response turnaround time by 60%** through automated slot filtering and decision routing, reducing manual filtering from 30 seconds to under 12 seconds across 4,000+ simulated booking records.

## RAG-POWERED CUSTOMER SUPPORT AGENT

[\[Github Link\]](#)

- Built an end-to-end ETL pipeline with Pandas to parse Flipkart reviews into LangChain documents, then ingested them into AstraDB Vector Store for high-performance, semantic top-k retrieval.
- Integrated Google Gemini-1.5-Pro embeddings and LangChain's ChatPromptTemplate to orchestrate context-aware retrieval and natural-language answer generation for real-time product recommendations.
- Developed a RESTful FastAPI backend with environment-driven configuration with PyYAML, python-dotenv and integrated secure secret management, paired with a modular AJAX chat UI using Jinja2, Bootstrap 4 and jQuery.