

# Gayathri Akkinapalli

(413) 466-1166 | [gayathri.akkinapalli2000@gmail.com](mailto:gayathri.akkinapalli2000@gmail.com) | [linkedin.com/in/gayathri-akkinapalli/](https://linkedin.com/in/gayathri-akkinapalli/) | [github.com/akkina10gaya](https://github.com/akkina10gaya)

## Education

University of Massachusetts Amherst (UMASS)   MS in Computer Science   CGPA: 4.0/4.0	Sep 2023 - May 2025
Indian Institute of Information Technology (IIIT)   B.Tech in Computer Engineering   CGPA: 9.07/10.0	Aug 2017 - May 2021

## Publications & Preprints

- |  |                                  |
|--|----------------------------------|
| 1. LS-GAN: Human Motion Synthesis with Latent-space GANs                       | <a href="#">IEEE WACVW '25</a>   |
| 2. Iterative Critique-Refine Framework for Enhancing LLM Personalization       | <a href="#">arXiv:2510.24469</a> |
| 3. Safe to Serve: Aligning Instruction-Tuned Models for Safety and Helpfulness | <a href="#">arXiv:2412.00074</a> |
| 4. Automated Model Selection for Tabular Data                                  | <a href="#">arXiv:2401.00961</a> |

## Technical Skills

**Programming Languages:** Python, R, C, HTML, SQL, MySQL, NoSQL | Familiar: Java, C++, PHP, CSS, JavaScript, C#

**Tools/Libraries:** LangChain, MCP, Unsloth, TRL, OpenAI, Wandb, Keras, PyTorch, TensorFlow, PySpark, SKLearn, NumPy, vLLM, Copilot

**Software/Frameworks:** Git, Docker, Kubernetes, Kafka, Flask, FastAPI, gRPC, AWS, Splunk, MongoDB, Tableau, CI/CD, Snowflake, FAISS

## Professional Experience

<b>Graduate Student Researcher</b>   Cisco	Jan 2025 – May 2025
<ul style="list-style-type: none"><li>Built <b>PerFine</b>, a training-free critique-refine <b>Agentic RAG</b> framework for enhancing personalization in long-text using <b>LangChain</b></li><li>Evaluated using <b>LLM-as-a-Judge</b> (G-Eval), improving personalization by <b>13%</b> and Meteor by <b>10%</b> over RAG-based baselines</li><li>Enhanced outputs using profile-grounded critic feedback retrieved via <b>Pinecone</b> (Vector DB), <b>FAISS</b>, <b>MCP</b> to refine style and content relevance.</li></ul>	
<b>Research Assistant</b>   UMass IESL Lab	Aug 2024 – Dec 2024
<ul style="list-style-type: none"><li>Built an <b>autoregressive</b> model that performs lookahead by decoding in superposition with just two forward passes using cross attention.</li><li>Applied the approach to <b>machine translation</b>, improving <b>BLEU</b> score and generation quality of the <b>MT5</b> model by approximately <b>15%</b>.</li></ul>	
<b>Machine Learning Engineer</b>   Carelon Global Solutions	Jun 2021 – Jul 2023
<ul style="list-style-type: none"><li>Built <b>Recommendation Systems</b> using <b>NER</b>, SpaCy &amp; ML models like <b>XGBoost</b>, LightGBM in <b>PySpark</b> improving NDCG@5 by <b>75%</b>.</li><li>Deployed models into <b>ENSO ML pipeline</b> using RabbitMQ, <b>Kubernetes</b>, and Kafka, cutting deployment time by <b>~70%</b> via <b>CI/CD</b> pipelines.</li><li>Integrated <b>REST APIs</b> with <b>Flask</b> for end-to-end model automation using <b>Hive</b>, <b>MongoDB</b>, <b>Redis</b>, reduced care plan creation time by <b>~2hrs</b>.</li><li>Created <b>Splunk</b> Dashboard for user feedback KPIs &amp; ran A/B testing on recommendations, driving a <b>60%</b> improvement in model performance.</li><li>Designed Aspect-Based Sentiment Analysis on call transcripts using <b>RoBERTa</b>, BERT, SpaCy, achieved <b>85%</b> accuracy &amp; <b>0.81 F1-score</b>.</li><li>Used AWS <b>SageMaker</b>, Kubeflow, S3, &amp; GlueDB for model development and built <b>Conversational AI</b> bot resolving <b>65%</b> of patient queries.</li><li>Integrated web-scraped healthcare articles into <b>Elasticsearch</b> with <b>ranking</b> optimization, reducing content retrieval time to <b>~120ms</b>.</li><li>Developed an <b>IBM Watson</b> Assistant chatbot that handled <b>100+</b> queries daily, and automated real-time insights through metric generation.</li><li>Built <b>Lambda</b> functions to extract conversational data from IBM Object Storage to <b>S3</b> and moved it into <b>DynamoDB</b> using <b>AWS Glue</b>.</li><li>Deployed <b>ETL pipelines</b> and ML models using Google Cloud <b>Vertex AI</b> with <b>Airflow</b> for orchestration, and <b>Docker</b> containerization.</li></ul>	
<b>AI Engineer Intern</b>   SensorDrops Networks (STEP at IIT Kharagpur)	Aug 2020 – Sep 2020
<ul style="list-style-type: none"><li>Designed real-time Social Distance Monitoring system using <b>YOLOv3</b> with live feed, bounding boxes, and <b>90%</b> detection accuracy.</li><li>Deployed the application on <b>AWS EC2</b> with <b>Docker</b>, enabling real-time analytics with <b>~200ms</b> latency via socket-based data transfer.</li></ul>	
<b>AI Engineer Intern</b>   Centre for Development of Advanced Computing (C-DAC)	May 2020 – Aug 2020
<ul style="list-style-type: none"><li>Developed a prototype of a customized deep CNN model to identify COVID-infected chest X-rays with an <b>accuracy of around 92%</b>.</li><li>Trained the model on <b>High Performance Computing (HPC)</b> for three chest X-ray classes, achieving a validation <b>F1-score of 0.9</b>.</li></ul>	

## Academic & Research Projects

<b>Multi Agent Study Assistant</b>   UMass   <a href="#">Github</a>	Aug 2025 - Present
<ul style="list-style-type: none"><li>Built a <b>multi-agent</b> RAG study assistant using LangChain, <b>ChromaDB</b>, LLMs (Flan-T5, MiniLM) that process web content, documents, user queries through specialized NLP agents (retriever, scraper, summarizer, note-maker), automating note creation with <b>85%</b> retrieval accuracy.</li><li>Designed an async agent pipeline with <b>FastAPI</b> backend and <b>Streamlit</b> UI, reducing study note generation time by <b>70%</b>.</li></ul>	
<b>Aligning LLMs towards safety and helpfulness</b>   UMass   <a href="#">Github</a>	Feb 2024 - May 2024
<ul style="list-style-type: none"><li>Aligned <b>LLaMa-2</b> toward safety using <b>LoRA</b>, QLoRA on PKU-SafeRLHF benchmark with SFT, RAFT, <b>RLHF</b>, <b>DPO</b> in Unsloth &amp; TRL.</li><li>Scored <b>93% safe</b> on DPO (40% SFT) with Llama-Guard on I-CoNa. Implemented <b>LLM-as-a-judge</b> to evaluate safety and helpfulness.</li></ul>	
<b>Human Motion Synthesis with Latent-space GANs</b>   UMass   <a href="#">Github</a>	Feb 2024 - May 2024
<ul style="list-style-type: none"><li>Generated text-to-motion sequences in latent space utilizing <b>GANS</b>, VAE, <b>CLIP</b> on HumanML3D with <b>Distributed training</b> in lightning.</li><li>Achieved a <b>FID of 0.48</b> with GAN in the latent space with <b>91% in FLOPs reduction</b> compared to Latent <b>Diffusion Model</b> on HumanML.</li></ul>	
<b>Optimization in Reinforcement Learning</b>   UMass   <a href="#">Github</a>	Sep 2023 - Dec 2023
<ul style="list-style-type: none"><li>Programmed Reinforce with baseline, <b>Actor-Critic</b>, <b>PPO</b> &amp; Semi Gradient n-step SARSA for Acrobat and Cartpole environments.</li><li>Attained stabilized <b>mean rewards of 470 (max = 500), -100 (max = 0)</b> on Cartpole &amp; Acrobat respectively using Reinforce, Actor-Critic.</li></ul>	