

# GAURAV TADKAPALLY

Los Angeles, CA | (213) 913-7899 | [gaurav.tadkapally@usc.edu](mailto:gaurav.tadkapally@usc.edu) | [linkedin.com/in/gauravreddy08](https://www.linkedin.com/in/gauravreddy08) | [gauravreddy08.github.io](https://gauravreddy08.github.io)

## EDUCATION

<b>University of Southern California</b> <b>Master of Science in Computer Science: 3.7/4.0</b> <ul style="list-style-type: none"><li>Served as a Teaching Assistant (TA) for the graduate course Applied Machine Learning for Natural Language Processing (ITP 459)</li></ul>	California, United States June 2023 - December 2024
<b>Vellore Institute of Technology</b> <b>Bachelor of Technology in Computer Science and Engineering: 8.94/10</b>	Andhra Pradesh, India May 2019 - May 2023

## EXPERIENCE

<b>Pitney Bowes</b> <b>Data Science Intern</b> <ul style="list-style-type: none"><li>Designed agentic code assistant for software testing, leveraging <b>Speculative Decoding</b> to accelerate inference speed by 3x and <b>Abstract Syntax Tree (AST) based retrieval</b> (tree-sitters) for document indexing (<b>Demo</b>)</li><li>Leveraged <b>Direct Preference Optimization (DPO)</b> and <b>4-bit QLoRA</b> quantization to finetune codellama &amp; gpt4o, improving model's generative accuracy by 15% (tested via Mutational Testing)</li><li>Implemented retrieval methodologies (BM25-FTS, Contextual Embedding, and Reranking algorithms) to enhance efficiency and accuracy in retrieving relevant codebase context</li><li>Integrated JaCoCo and Mutational Testing (PIT) to automatically evaluate code coverage &amp; test effectiveness of generated unit tests</li></ul>	Connecticut, United States June 2024-August 2024
---	---

<b>MUKHAM</b> <b>Machine Learning Engineer Intern</b> <ul style="list-style-type: none"><li>Optimized facial recognition model for edge deployment (mobile application), leveraging <b>knowledge distillation</b>, <b>Post-training Quantization (8-bit quantization)</b> and <b>Automatic Mixed Precision</b>, decreasing model size by 75%</li><li>Designed a Presentation Attack Detection system (facial spoof detection) utilizing the Lucas Kanade algorithm for motion analysis, achieving a 80% success rate in identifying spoofed faces</li></ul>	Andhra Pradesh, India October 2022-May 2023
--	--

<b>MUKHAM Pvt Ltd</b> <b>Research Assistant</b> <ul style="list-style-type: none"><li>Developed a UAV-based wildfire detection algorithm utilizing the EfficientNetB0 architecture, incorporating <b>Neural Architecture Search (NAS)</b> for model optimization, resulting in a 98% precision rate</li><li>Engineered smart glasses with a Object Detection model (Incremental Learning) for visually impaired, leading 78% accuracy</li></ul>	Andhra Pradesh, India October 2022 - May 2023
--	--

## SKILLS AND CERTIFICATIONS

**Languages:** Python, TypeScript, JavaScript

**ML Stack:** PyTorch, Tensorflow, HuggingFace, LangChain, Keras, OpenCV, Scikit-learn, Pandas, NumPy

**Tools & Technologies:** AWS (Cloud Practitioner), Azure (AI Fundamentals), MySQL, MongoDB, Selenium, Redis

## ACADEMIC PROJECTS

<b>AK15: Agentic Kubernetes Middleware (Github)</b> <ul style="list-style-type: none"><li>Devised an LLM-based middleware that automates Kubernetes cluster read queries, achieving a 93% reduction in contextual token usage through agentic function calling and context retrieval</li><li>Implemented 15 specialized API functions enabling the LLM to perform human-like, context-aware interactions with Kubernetes, optimizing and reducing API costs by leveraging targeted data retrieval strategies</li></ul>
<b>GlancyAI: Consumer Product Research Assistant (Github)</b> <ul style="list-style-type: none"><li>Developed an AI agent using GPT-4 and Agentic Retrieval Augmented Generation (RAG), with a vector database for optimized query retrieval, automating the extraction of data from YouTube</li><li>Integrated summarization module condenses extensive online information into concise insights, streamlining the product recommendation process and significantly reducing user research time</li></ul>
<b>Original Vision Transformer Implementation from Scratch (Github)</b> <ul style="list-style-type: none"><li>Implemented ViT components including MultiHeadAttention, Image Patch Embedding, and MLP layers, achieving a one-to-one parameter match (86 million) with the original proposed model</li></ul>

## PUBLICATIONS

- Sethuraman, S. C., Reddy Tadkapally, G. et al. **Simplymime: A dynamic gesture recognition and authentication system for smart remote control**. IEEE Sensors Journal (2024). <https://doi.org/10.1109/JSEN.2024.3487070>
- Sethuraman, Sibi C., Gaurav Reddy Tadkapally, et al. **iDrone: IoT-Enabled Unmanned Aerial Vehicles for Detecting Wildfires Using Convolutional Neural Networks**. Springer Nature Computer Science (2022). <https://doi.org/10.1007/s42979-022-01160-7>