

# GAURAV TADKAPALLY

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

## EDUCATION

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

## EXPERIENCE

<b>Pitney Bowes</b>	Connecticut, United States
<b>Data Science Intern</b>	June 2024-August 2024
<ul style="list-style-type: none"><li>Designed an agentic feedback loop with fine-tuned CodeLlama and GPT-4o, using JaCoCo code coverage tool to iteratively optimize test suites and increase code coverage by 15% (<a href="#">Demo</a>)</li><li>Implemented optimized LLM decoding strategies (Speculative Decoding), accelerating inference by 3x, and Abstract Syntax Tree (AST)-based retrieval for precise code context</li><li>Implemented Direct Preference Optimization (DPO) and 4-bit QLoRA quantization, improving model's code generation accuracy</li></ul>	

<b>MUKHAM</b>	Andhra Pradesh, India
<b>Machine Learning Engineer Intern</b>	October 2022-May 2023
<ul style="list-style-type: none"><li>Optimized facial recognition model for edge deployment (mobile application), leveraging <b>knowledge distillation, Post-training Quantization (8-bit quantization) and 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>	

<b>MUKHAM Pvt Ltd</b>	Andhra Pradesh, India
<b>Research Assistant</b>	October 2022 - May 2023
<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 Continual Object Detection model (Incremental Learning) for visually impaired, leading 78% navigational accuracy</li></ul>	

## SKILLS AND CERTIFICATIONS

**Languages:** Python, Java, R, JavaScript  
**ML Stack:** PyTorch, Tensorflow, HuggingFace, LangChain, Keras, OpenCV, Scikit-learn, Pandas, NumPy, Matplotlib  
**Tools & Technologies:** AWS (Cloud Practitioner), Azure (AI Fundamentals), MySQL, MongoDB, Selenium

## ACADEMIC PROJECTS

<b>AK15: Agentic Kubernetes Middleware (Github)</b>
<ul style="list-style-type: none"><li>Devised an LLM-powered middleware that automates Kubernetes cluster read queries, achieving a 93% reduction in contextual token usage through intelligent function calling and agentic 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 (<a href="#">Github</a>)</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 web sources and YouTube transcripts</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

<ul style="list-style-type: none"><li>Sethuraman, S. C., Reddy Tadkapally, G. et al. <b>Simplymime: A dynamic gesture recognition and authentication system for smart remote control</b>. IEEE Sensors Journal (2024). <a href="https://doi.org/10.1109/JSEN.2024.3487070">https://doi.org/10.1109/JSEN.2024.3487070</a></li><li>Sethuraman, Sibi C., Gaurav Reddy Tadkapally, et al. <b>iDrone: IoT-Enabled Unmanned Aerial Vehicles for Detecting Wildfires Using Convolutional Neural Networks</b>. Springer Nature Computer Science (2022). <a href="https://doi.org/10.1007/s42979-022-01160-7">https://doi.org/10.1007/s42979-022-01160-7</a></li></ul>
---