

# Sidhartha Reddy Potu

+1 (917) 914-1526 | sp7835@nyu.edu | GitHub | LinkedIn | Portfolio

## EDUCATION

<b>New York University, Tandon School of Engineering</b> <i>Master of Science in Computer Science; GPA: 3.83/4</i>	New York, NY Sep 2023 – May 2025
<b>Indian Institute of Technology (IIT), Indore</b> <i>Bachelor of Technology in Electrical Engineering, Minor in Astronomy; GPA: 8.49/10</i>	Indore, IN Jul 2019 – May 2023

## EXPERIENCE

<b>SWE Intern</b> <i>Granica AI</i>	Mountain View, CA Sep 2025 – Present
<ul style="list-style-type: none"><li>Implemented a dual-lane FIFO router for memory-bound workers, doubling throughput (<b>10</b> to <b>20+</b> req/worker).</li><li>Prevented runaway Kubernetes API calls under <b>20k</b> concurrency by caching scalar state in Go.</li><li>Fixed lock contention and races in load balancer/consumer goroutines, improving throughput (<b>28</b> to <b>31–32</b> cores).</li><li>Built scale-validation tests for multi-consumer routing and integration checks for compression pipelines.</li></ul>	
<b>MTS Intern</b> <i>Neopolis AI</i>	Mountain View, CA Aug 2025 – Sep 2025
<ul style="list-style-type: none"><li>Refined LiveKit-based LLM voice AI Agent prompts to improve instruction-following and downstream tool inputs.</li><li>Used Langfuse observability to inspect agent outputs and guide prompt iteration.</li><li>Added session inactivity timeouts to prevent idle conversations and improve reliability.</li></ul>	
<b>Graduate Student Researcher, CILVR, NYU</b> <i>3D Computer Vision under Prof. David Fouhey</i>	New York, NY Jan 2024 – May 2025
<ul style="list-style-type: none"><li>Integrated COLMAP SfM with HaMeR to recover global 3D hand trajectories on Epic-Kitchens and EgoExo4D.</li><li>Optimized hand-motion trajectories using temporal smoothing with velocity and jerk regularization.</li><li>Recovered full-perspective 3D from weak-perspective predictions using SfM intrinsics and frame transforms.</li><li>Developed visualization workflows with Rerun Viewer and scaled experiments on NYU HPC using SLURM.</li></ul>	
<b>Course Assistant, NYU</b> <i>Computer Vision, Deep Learning, LLVM courses</i>	New York, NY Sep 2024 – May 2025
<ul style="list-style-type: none"><li>Led office hours, grading, and project support across CV/ML courses; supported <b>450+</b> students.</li><li>Course staff for DL (Prof. Chinmay Hegde) and LLVM (Prof. Saining Xie); Advanced CV (Prof. David Fouhey).</li></ul>	
<b>B.Tech Thesis, PRIA Lab, IIT Indore</b> <i>Multimodal Crowd Counting using Vision Transformers [Report]</i>	Indore, IN Jun 2022 – Dec 2022
<ul style="list-style-type: none"><li>Trained a ViT-based model in a weakly supervised setting, achieving MAE 15.78 and MSE 30.70.</li><li>Designed a Pyramid Vision Transformer (PVT v2) model, achieving MAE 19.218 and MSE 32.940.</li><li>Ranked <b>top 10 of 250+</b> B.Tech theses; nominated for Best Thesis award at IIT Indore.</li></ul>	

## PROJECTS

Fine-Tuning Medical QA Models — GitHub	Oct 2024 – Dec 2024
<ul style="list-style-type: none"><li>Optimized Meta LLaMA 3.2 3B for medical QA using LoRA, QLoRA, and GaLore.</li><li>Achieved best USMLE accuracies: Step 1 - <b>42%</b> (GaLore), Step 2 - <b>38%</b> (QLoRA), Step 3 - <b>45%</b> (GaLore).</li><li>Improved MedQuAD F1 scores to <b>70%</b> (QLoRA) and 69% (GaLore) with efficient fine-tuning.</li></ul>	
VioletPass: Ticket Booking Platform — GitHub	Oct 2024 – Dec 2024
<ul style="list-style-type: none"><li>Designed a robust, scalable ticket booking platform on AWS to prevent double booking.</li><li>Implemented Redis distributed locks with TTL for seat reservations and PostgreSQL for finalization.</li><li>Integrated QR code-based authentication for secure and efficient event check-ins.</li></ul>	
AnimeVerse: Recommendation System — GitHub	Apr 2024 – May 2024
<ul style="list-style-type: none"><li>Built a personalized anime recommendation system; analyzed broadcast types, genres, and scores.</li><li>Analyzed <b>300k</b> user-anime interactions with PySpark and trained ALS, TF-IDF, and ChromaDB models.</li></ul>	

## TECHNICAL SKILLS

<b>Languages:</b> Python, Go, C/C++, SQL
<b>Systems:</b> Linux, Git, Docker, Kubernetes, AWS, GCP, CI/CD, Weights & Biases (W&B)
<b>ML/Data:</b> PyTorch, Transformers (Hugging Face), TensorFlow, NumPy, Pandas, Scikit-learn, PySpark, OpenCV