

# SIDDHANT HITESH MANTRI

Email: siddhantmantri328@gmail.com — LinkedIn: siddhant-mantri — GitHub: siddhant-192 — Phone: (+1) 619-986-7810 — Website: www.siddhantmantri.com

## EDUCATION

<b>University of California, San Diego (UCSD)</b> <b>Master of Science, Computer Science and Engineering</b>	San Diego, CA, USA 2025 – 2027
<b>Mukesh Patel School of Technology Management and Engineering, NMIMS</b> <b>Bachelor's in Technology, Computer Science and Engineering (Cyber Security)</b> CGPA 3.86/4	Mumbai, India 2021 – 2025
<b>Indian Institute of Technology, Madras</b> <b>Bachelor's in Science, Data Science and Applications</b> CGPA 7.39/10, Project CGPA 9.0/10	Chennai, India 2021 – 2025

## SKILLS

<b>Programming Languages:</b> Python, Java, C++, JavaScript, SQL, R
<b>ML/AI Frameworks:</b> PyTorch, TensorFlow/Keras, scikit-learn, Hugging Face Transformers, LangChain
<b>Tools &amp; Technologies:</b> Git, Docker, CUDA, Linux, Pandas, NumPy, NLTK, Jenkins, MLflow, Airflow, Selenium, pytest, FastAPI, Node.js
<b>Relevant Coursework:</b> Data Structures and Algorithms, Software Engineering, Database Systems, Computer Networks, Operating Systems, Reinforcement Learning, Deep Learning, Probability and Statistics
<b>Specializations:</b> Large Language Models, Computer Vision, Natural Language Processing, Machine Learning

## WORK EXPERIENCE

<b>Computation for Indian Language Technology, IIT Bombay</b> <i>Research Intern</i>	Sep. 2024 – Present Mumbai, India
<ul style="list-style-type: none"><li>Engineered structured-data pipeline converting Hindi WordNet into 1.25M diverse instruction-response pairs for specialized conversational AI systems, achieving 91.0 pedagogical effectiveness score vs. 79.4-83.6 for general-purpose models</li><li>Fine-tuned Gemma3-12B language model using LoRA with 4-bit quantization, achieving 58% improvement in response consistency and sub-second inference times while reducing memory requirements by 75%</li><li>Built context-aware keyword detection microservice with FastAPI deployment, improving contextual understanding accuracy by 25% for low-resource language processing</li></ul>	
<b>Sudha Gopalakrishnan Brain Centre, IIT Madras</b> <i>Machine Learning Intern</i>	Jan. 2025 – Jun. 2025 Chennai, India
<ul style="list-style-type: none"><li>Trained and enhanced YOLOv9e-seg model on 25-class brain segmentation dataset, achieving 0.97 Dice score and 30% inference speed improvement through novel post-processing techniques</li><li>Architected end-to-end ML pipeline using PyTorch for production deployment, processing 1,000+ medical images with automated quality validation and real-time inference capabilities</li></ul>	
<b>Volkswagen Group Technology Solutions India</b> <i>Machine Learning Intern</i>	Jun. 2024 – Sep. 2024 Pune, India
<ul style="list-style-type: none"><li>Improved RAG-based conversational AI system (Project AISHA 2.0) by optimizing vector databases and caching mechanisms, improving query response time by 40% and reducing memory usage by 35%</li></ul>	

## PROJECTS

<b>AI-Powered Learning Management System (LMS)</b>	Jun. 2024 – Aug. 2024
<ul style="list-style-type: none"><li>Led the development team implementing AI-powered features, including course summaries, peer-driven insights, and coding assistance, leveraging advanced language models such as LLaMa3-70B to enhance the learning experience</li><li>Integrated the vLLM inference engine to improve the efficiency and performance of LLM queries, reducing latency and scaling AI-generated responses across the system. Implemented FastAPI for serving multiple AI-driven endpoints with Git version control and Scrum development methodology for agile project management, utilizing pytest for unit testing and debugging of application components</li></ul>	
<b>ShelfSense - Library Management System</b>	Jan. 2024 – Apr. 2024
<ul style="list-style-type: none"><li>Created full-stack web application using Vue.js frontend and Flask backend, implementing comprehensive library operations including book management, user authentication, and issue tracking with role-based access control and RESTful API design, utilizing Selenium for automated testing</li><li>Architected scalable system with SQLAlchemy ORM, Redis caching for improved performance, and Celery task scheduling for automated reminders and report generation using RabbitMQ message broker</li></ul>	

## PUBLICATIONS

<b>ICTCS 2024 (Springer)</b>	2024
<ul style="list-style-type: none"><li>"Architectural Framework for Automated Incident Response: Leveraging LLMs And Classifiers for Rapid Post-Attack Analysis and Reporting"</li></ul>	
<b>IJCNLP-AAACL 2025 - Under Review</b>	2025
<ul style="list-style-type: none"><li>"From Lexicon to AI: A Structured-Data Pipeline for Specialized Conversational Systems in Low-Resource Languages"</li></ul>	