

SIDDHANT HITESH MANTRI

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

EDUCATION

University of California, San Diego (UCSD) Master of Science, Computer Science and Engineering CGPA 4.0/4	San Diego, CA, USA 2025 – 2027
Mukesh Patel School of Technology Management and Engineering, NMIMS Bachelor's in Technology, Computer Science and Engineering (Cyber Security) CGPA 3.86/4	Mumbai, India 2021 – 2025
Indian Institute of Technology, Madras Bachelor's in Science, Data Science and Applications CGPA 7.39/10, Project CGPA 9.0/10	Chennai, India 2021 – 2025

SKILLS

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

WORK EXPERIENCE

Computation for Indian Language Technology, IIT Bombay <i>Research Intern</i>	Sep. 2024 – Dec. 2025 Mumbai, India
<ul style="list-style-type: none">Engineered a structured-data pipeline converting Hindi WordNet into 1.25M diverse instruction-response pairs, utilizing dynamic chunking and sliding-window algorithms to preserve semantic relationshipsFine-tuned Gemma-3-12B using QLoRA with 4-bit quantization, reducing memory requirements by 75% (48GB to 12GB) while maintaining semantic competence on limited hardwareAchieved superior pedagogical effectiveness (91.0 vs. 79.4 for GPT-4.1) and 86% higher response consistency ($\sigma=1.0$ vs 7.4), validating the system for a specialized Hindi language learning chatbotSynthesized novel findings into a comprehensive research manuscript submitted to ACL 2026, documenting the framework's methodology and evaluation results	
Sudha Gopalakrishnan Brain Centre, IIT Madras <i>Machine Learning Intern</i>	Jan. 2025 – Jun. 2025 Chennai, India
<ul style="list-style-type: none">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 techniquesArchitected end-to-end ML pipeline using PyTorch for production deployment, processing 1,000+ medical images with automated quality validation and real-time inference capabilities	
Volkswagen Group Technology Solutions India <i>Machine Learning Intern</i>	Jun. 2024 – Sep. 2024 Pune, India
<ul style="list-style-type: none">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%	

PROJECTS

AI-Powered Learning Management System (LMS)	Jun. 2024 – Aug. 2024
<ul style="list-style-type: none">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 experienceIntegrated 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	
ShelfSense - Library Management System	Jan. 2024 – Apr. 2024
<ul style="list-style-type: none">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 testingArchitected scalable system with SQLAlchemy ORM, Redis caching for improved performance, and Celery task scheduling for automated reminders and report generation using RabbitMQ message broker	

PUBLICATIONS

ICTCS 2024 (Springer)	2024
<ul style="list-style-type: none">"Architectural Framework for Automated Incident Response: Leveraging LLMs And Classifiers for Rapid Post-Attack Analysis and Reporting"	
ACL 2026 - Under Review	2026
<ul style="list-style-type: none">"From Lexicon to AI: A Structured-Data Pipeline for Specialized Conversational Systems in Low-Resource Languages"	