

# Manas Jain

📞: +1 (858)-220-1295 | 📩: manasjain26@gmail.com | 🌐: manasjain26 | 🌐: manasjain26.github.io | 🌐: Manas Jain | 🌐: GScholar

## EDUCATION

### University of California San Diego

San Diego, CA

*Master of Science in Data Science*

Sept. 2025 - June 2026

- Coursework: Probability & Statistics, Statistical Models, Recommender Systems, ML Systems, Safety in GenAI

### Indian Institute of Technology Bombay

Mumbai, India

*Bachelor of Technology in Civil Engineering, Double minor in CS, AI & Data Science (CMInDS)*

July 2017 - June 2021

- GPA: 8.7 (major) — 8.75 (minor)

## TECHNICAL SKILLS

**Programming:** C/C++, Python, Go, Java, React, R, Julia, SQL

**Software/Frameworks:** Pytorch, OpenCV, Jax, Scope, Git, AWS, Langchain, LangGraph, MCP, TensorRT-LLM, vLLM, Tinker

## INDUSTRY EXPERIENCE

### Microsoft R&D

Bangalore, India

*Data & Applied Scientist 2 — Bing Autosuggest (Microsoft AI)*

Mar. 2024 - Sept. 2025

- Designed & deployed **T5 & mT5 based SLM model** - **Multilingual Generative Query Suggestions** for a user prefix in production, achieving **+200k DAU**, improving coverage (**+2 SBS**) and reducing latency (**-6 ms**).
- Developed Generative Suggestions for **Mid-Query Reformulation**, improving query intent handling through **correction aware modeling** leading to **+20k DAU** & winning the **Best Hack Award** among 30+ submissions.
- Fine-tuned **SLMs** (Phi-3-mini-3.8B, Llama-3.2-3B) with **SFT/instruction tuning** using PEFT technique **LoRA** for EYP
- Optimized inference with **TensorRT-LLM**, **pruning**, **vocab trimming** to handle **100K QPS** within **30 ms** SLA.
- Built & managed **DeepSeek-R1-Distill-Qwen-32B** based offline evaluation pipeline with **vLLM**, enabling **large-scale benchmarking**, supporting competitor comparisons and A/B testing.

### Wadhwani Institute of AI

Bangalore, India

*Associate ML Scientist — Agriculture NLP Team*

Dec. 2022 - Mar. 2024

- Led development of **multilingual NLP pipelines & RAG based LLM-powered chatbots** (news monitoring **Krishi 24/7**, farmer support), including **fine-tuning**, **evaluation**, & **deployment** of models for classification & entity extraction.

### HiLabs Inc.

Pune, India

*Research Data Scientist — Biomedical NLP Team*

July 2021 - Dec. 2022

- Built and deployed ML pipelines for **ICD10 Medical Code prediction** and NER from provider contracts using **BioBERT**, **SciSpacy**, **TF-IDF**, **OCR** (OpenCV, Tesseract), improving accuracy of healthcare AI products.

### Daikin Industries

Osaka, Japan

*Research and Development AI Intern — DigiNavi, ICT Group*

May 2020 - June 2020

- Built an **end-to-end NLP pipeline** for video tagging, captioning, and summarization using **TF-IDF**, **BERT**, **LDA** **Topic Modelling**, **S-BERT**, with transcripts generated via **GCP Speech-to-Text API** from recorded videos.

## PUBLICATIONS

### Zephyrus: An Agentic Framework for Weather Science



*S. Varambally, M. Fisher, ..., M. Jain, ..., T. Berg-Kirkpatrick, D. Watson-Parris, Y. Ma, R. Yu*

Proceedings of **ICLR '26**; Main Track

### EnhanceMyPrompt: Rewriting Chat Queries for Effective Response Generation from LLMs



*M. Jain\*, T. Abhishek\*, S. Hardia, S. Suriyanarayanan, S. Anil, R. Gandhi, M. Gupta (\* denotes equal contribution)*

Proceedings of **CIKM '25**; Applied Research Track

### Natural Answer Generation: From Factoid Answer to Full-length Answer using Grammar Correction



*M. Jain, S. Saha, P. Bhattacharyya, G. Chinnadurai, M. Vatsa*

Proceedings of **ICON '24**; Main Track

## KEY RESEARCH EXPERIENCE

### Multimodal LLM Agents for Scientific & Weather Discovery

Ongoing

*Prof. Rose Yu, Spatiotemporal Lab, CSE*

*University of California San Diego*

- Pioneered a **Bayesian Optimization evaluation task** integrated into ZephyrusBench using a **Neural GCM simulator**
- Enhancing the Zephyrus agent's ability to solve multi-step scientific problems via **subgoal decomposition** and RL

### Post Training Alignment Against Prompt Injection

Fall 2025

*Prof. Yu-Xiang Wang, HDSI*

*University of California San Diego*

- Implemented a sequential alignment pipeline (**SFT** followed by **DPO**) to robustify Llama-3-3B and Qwen-3-4B using Tinker
- Reduced Attack Success Rate (ASR) by **70%** (9.0% → 2.7%) on **Qwen-3-4B** while simultaneously improving benign helpfulness scores by **34%** compared to baseline.
- Engineered an automated "**LLM-as-a-Judge**" evaluation harness to benchmark model safety against adversarial datasets (WildJailbreak, AdvBench, JailbreakBench).