

# RAGHAV KACHROO

(858)-241-1760 | [rkachroo@ucsd.edu](mailto:rkachroo@ucsd.edu) | [linkedin.com/raghavkachroo](https://linkedin.com/raghavkachroo) | [github.com](https://github.com/rkachroo)

## EDUCATION

<b>University of California, San Diego</b>	Sep 2024 – Mar 2026
Master of Science in Data Science (Artificial Intelligence & Machine Learning)	GPA: 3.63/4.0
• Relevant Coursework: Distributed ML Systems, LLM & AI Agents, Data Structures & Algorithms	
<b>Indraprastha Institute of Information Technology, Delhi</b>	Sep 2022 – Sep 2023
Post Graduate Diploma in Data Science & Artificial Intelligence	GPA: 4.0/4.0
<b>University of Delhi</b>	July 2018 – June 2021
Bachelor of Management Studies	GPA: 3.57/4.0

## EXPERIENCE

<b>Amazon</b>	June 2025 – Sep 2025
Software Development Engineer Intern	Bellevue, WA
• Engineered an incident resolution toolkit to cut <b>12+ engineering hours</b> per week by orchestrating AWS Step Functions and DynamoDB to automatically enrich incident tickets and execute routine SOPs.	
• Architected a search microservice to process <b>5GB+</b> of distributed logs in <b>under 45 seconds</b> by utilizing AWS Lambda to concurrently filter streams based on time and keyword inputs.	
• Integrated infrastructure log retrieval into chat workflows via a <b>Model Context Protocol (MCP)</b> server, enabling LLM-powered analysis of production logs through natural language.	
• Established reusable SDK and MCP components to standardize LLM agent development and evaluation workflows across 7 teams.	
<b>Aark Global</b>	Apr 2023 – Sep 2024
Software Developer, AI/ML	Delhi, India
• Orchestrated a high-throughput pipeline processing <b>3,000+ documents daily</b> by utilizing Azure Queue Storage to distribute workloads across a pool of Virtual Machines.	
• Designed a hybrid data layer to deliver <b>sub-100ms P95 latency</b> across production requests by bridging MongoDB for rapid development with Cosmos DB for high performance serving.	
• Architected a full text search capability by designing a Tesseract OCR workflow to populate an Elasticsearch inverted index, delivering query results in <b>under 180ms</b> .	
• Integrated a conversational AI interface over indexed data to enable LLM-driven document analysis and question answering.	

<b>Concentrix</b>	Jun 2022 – Mar 2023
Data Engineer	Gurugram, India
• Built a Sentence-BERT and scikit-learn classification pipeline to automate sentiment analysis, reducing turnaround time from <b>3 days to 6 hours</b> .	
• Scaled data ingestion to increase <b>brand coverage by 60%</b> by replacing legacy scraping with parallelized Airflow and Kafka pipelines capable of monitoring 8 major social and e-commerce platforms.	
• Implemented proactive error detection and logging for ingestion workflows, reducing <b>resolution time by 67%</b> .	

## RESEARCH & PROJECTS

<b>LLM Finetuning for Patient Routing + Clinical Decision Support</b>   <a href="https://ieeexplore.ieee.org/document/11252667">https://ieeexplore.ieee.org/document/11252667</a>
• Developed a patient triage service using BioBERT on AWS Bedrock + FastAPI, achieving <b>91.6%</b> accuracy across <b>1,000+</b> assessments with EHR-integrated REST APIs.
• Optimized pipelines to cut inference time from <b>16s to 5s</b> by using LangChain routing logic and response caching.

<b>Autonomous GPU Rental Agent</b>   <a href="#">GitHub</a>
• Built an autonomous LLM-driven system to plan and execute multi-step GPU rental workflows including vendor selection, budgeting, job execution, monitoring, retries, and payment handling.
• Architected a MongoDB-backed state machine enabling multi-hour workflows with structured planning, crash-safe recovery, and automatic vendor switching across serverless restarts.

## TECHNICAL SKILLS

**AI & ML:** PyTorch, TensorFlow, Transformers, LangChain, Hugging Face, Scikit-learn, LlamaIndex, LangGraph

**LLM Systems:** Agent SDKs, Model Context Protocol (MCP), RAG, Vector Databases

**Programming:** Python, Java, JavaScript, REST APIs, System Design

**Data Engineering:** SQL, MongoDB, Neo4j, ETL Pipelines, Apache Kafka, Data Modeling, Monitoring

**Cloud & Infra:** AWS (Lambda, Bedrock, S3, EC2, CloudWatch), Azure, GCP, Docker, CI/CD, Git, Distributed Systems