

# Yash Thakkar

[thakkaryash21@gmail.com](mailto:thakkaryash21@gmail.com) • [yash-thakkar.com](http://yash-thakkar.com) • [linkedin.com/in/yashthakkar21/](https://linkedin.com/in/yashthakkar21/) • 412-699-0335

## EDUCATION

<b>Carnegie Mellon University</b> , Pittsburgh, PA	Dec 2025
Master of Information Systems Management	<b>GPA - 4.04</b>
<b>Relevant Coursework:</b> Cloud Computing, Distributed Systems, LLM Applications, Machine Learning, MLOps	
<b>Teaching Assistant:</b> Engineering Data-Intensive Scalable Systems, Data Science & Big Data, Agile Methods	
<b>University of Mumbai</b> , Dwarkadas J. Sanghvi College of Engineering	May 2022
Bachelor of Engineering, Mechanical Engineering	<b>GPA - 3.92</b>

## SKILLS

<b>Languages &amp; Libraries:</b> Python, TypeScript, JavaScript, Java, SQL, NumPy, Pandas, Bash/Shell
<b>AI &amp; ML:</b> PyTorch, Hugging Face, LangChain, Prompt Engineering, RAG, Scikit-learn, MLflow
<b>Infra &amp; DB:</b> AWS, GCP, Azure, Docker, Kubernetes, Terraform, Kafka, MySQL, MongoDB, Redis, Neo4j
<b>Web:</b> FastAPI, Flask, React.js, Redux, Node.js, Next.js, GraphQL, REST API, Pytest, Tailwind, WebRTC, HTML/CSS

## EXPERIENCE

<b>Newfront Insurance</b> , San Mateo, CA	May 2025 – Aug 2025
Software Engineering Intern	
○ <b>LLM Document Processing:</b> Built automated pipelines to extract & review benefits plan data using <b>FastAPI, GraphQL, PostgreSQL, &amp; OpenAI APIs</b> , <b>reducing processing time by over 90% &amp; saving \$175k annually</b>	
○ <b>GraphQL Performance Optimization:</b> Improved GraphQL performance by optimizing data retrieval paths & streamlining database queries, <b>reducing server load &amp; improving latency by 450 ms</b>	
○ <b>AI Pipeline Enhancement:</b> Improved text extraction accuracy through enhanced <b>prompt engineering &amp; context-aware document chunking</b> , <b>achieving 7% higher recall &amp; 5% better precision</b>	
○ <b>System Monitoring:</b> Implemented comprehensive monitoring with <b>over 20 Datadog alerts &amp; an analytics dashboard</b> to track system adoption, usage patterns & cost savings, <b>boosting operational reliability</b>	

<b>Advect Systems</b> , Mumbai, India	Jan 2022 – Aug 2024
Founding Software Engineer	
○ <b>Full-Stack Product Development:</b> Led end-to-end development of Nanofactory, a decentralized 3D printer management & automation app using <b>Svelte, Flask, &amp; WebRTC</b> , to enable <b>remote manufacturing workflows</b>	
○ <b>CI/CD &amp; Automation:</b> Built comprehensive CI/CD pipelines with <b>Docker, GitHub Actions</b> , Cloudflare & Python Integrations, <b>reducing time-to-market by 40% &amp; improving team productivity by 34%</b>	
○ <b>Test Coverage:</b> Built end-to-end & unit test suites with <b>Playwright &amp; PyTest</b> , achieving <b>80% code coverage</b> across critical application workflows	

## PROJECTS

<b>Real-Time Air Quality Forecasting System</b>	Oct 2025
○ <b>Kafka Streaming Inference:</b> Developed a real-time inference pipeline using <b>Kafka with over 100 engineered temporal features</b> to predict CO levels, achieving an <b>R<sup>2</sup> accuracy of 0.96</b>	
○ <b>Model Registry &amp; Drift Monitoring:</b> Integrated <b>MLflow</b> for versioned deployments & <b>Evidently</b> for daily/weekly drift detection, <b>improving model governance</b>	
<b>AI-Powered Learning Portal</b>	
○ <b>RAG Pipeline:</b> Engineered a modular retrieval system with <b>query reformulation &amp; ChromaDB</b> vector search, improving multi-turn grounding & quiz generation across 20 curated sources	
○ <b>CoT Evaluation:</b> Implemented a structured 5-step Chain-of-Thought, evaluated via LLM-as-a-Judge, <b>achieving 4.47/5 reasoning score</b>	
<b>Twitter Social Graph Analysis</b>	
○ <b>Graph Processing:</b> Analyzed a Twitter social graph using <b>Spark's iterative execution</b> model and PageRank algorithm over a <b>10.4 GB dataset</b> , identifying influential users with a 15-minute end-to-end job runtime	
○ <b>Performance Debugging:</b> Diagnosed Spark job bottlenecks by inspecting <b>YARN logs</b> and <b>Spark Web UI</b> , applying targeted tuning under constrained cluster resources to <b>improve execution efficiency by 74%</b>	