

# SWAYAM GUPTA

2024 Graduate · Material Science & Engineering  
Indian Institute of Technology, Kanpur

✉ gupta.swayam123@gmail.com · ☎ (+91) 8957370095  
👤 swayamg20 · 💬 guptaswayam20  
🔗 swayamgupta20.medium.com · 🌐 swayamg20.github.io

## EDUCATION

Year	Degree	Institution	Score
2024	B.Tech, MSE	Indian Institute of Technology, Kanpur	6.32/10.0
2019	CBSE – XII	Pt. Deen Dayal Vidyalaya, Kanpur (CBSE)	89%
2017	CBSE – X	Pt. Deen Dayal Vidyalaya, Kanpur (CBSE)	9.4/10.0

## ACHIEVEMENTS

- 2nd place in ixigo hackweek'25, built GenAI travel product **2025**
- Won ixigo Premio Award for top contributor Engineers **2025**
- Silver Medal in PS by ISRO hosted by Inter IIT Tech Meet **2022**

## WORK EXPERIENCE

ixigo	<i>Software Engineering 2 - AI Products</i>	<i>Aug'24 - present</i>
• Architected a real-time ASR-LLM-TTS conversational pipeline, solving Voice AI UX challenges like <b>turn detection</b> and <b>barge-in</b> using open source <b>vad</b> like <b>TEN/ Silero</b>		
• Built an internal WebRTC and event-driven streaming framework, achieving sub-500ms latency for real-time voice AI		
• Did POCs with <b>ElevenLabs</b> , <b>conversational agent pipelines</b> , and <b>LiveKit WebRTC</b> , validating the shift from custom WebRTC streams to a scalable managed real-time media platform		
• Implemented custom FFT/PCM audio analysis adding accurate speech boundary detection, <b>silence trimming</b> , and <b>queuing</b>		
• Designed multi agent orchestration layer supporting dynamic tool registration, workflows, proactive triggers, and LLM observability		
• Built the Conversational AI <b>Client SDK</b> from scratch, a cross platform, event-driven integration layer enabling <b>realtime multimodal omnipresent AI</b> across products		
• Developed a <b>low-latency contextual memory system</b> using Redis and Zep, achieving millisecond-level retrieval for long and short-term context		
• Optimized voice AI infra by reducing overall latency 50% and compute cost 40% through semantic caching and streaming optimizations, TTS caching, LLM function tool call optimizations		
• Built LLM based post-call analytics, including summarization, structured intents, call scoring, sentiment analysis, and NPS		
• Developed ixigo's <b>public facing MCP server</b> end-to-end with proxied OAuth 2.0, persistent memory, and production ready reliability, enabling ticket booking, hotel queries, and complex agentic workflows via AI agents		

Overlayy AI	<i>Apr'24 - Aug'24</i>
<i>Fouding Engineer</i>	
• Built scalable production-ready LLM backend pipeline from scratch using advanced RAG techniques, improving search latency while managing complete startup technology infrastructure	
• Managed comprehensive AWS DevOps infrastructure for multiple microservices, implementing CI/CD automation and integrating AI-enhanced Plausible analytics using Elixir/Phoenix framework	
• Developed secure real-time WebSocket messaging system with robust API authorization, architecting custom analytics solution using ClickHouse data warehouse for optimal performance	

## PREVIOUS INTERNSHIP EXPERIENCE

Fischer Jordan (offered PPO)	<i>Jun'23 - Feb'24</i>
Built React.js client management portal with Redux state management, RESTful APIs, and WebSocket real-time features	

Llama Mindfulness	<i>Apr'23-Jun'23</i>
Developed Node.js backend with Slack Bolt integration, configured Moodle LMS plugins, created JavaScript admin portal	

Vaticinari Solutions	<i>Jul'23-Sep'23</i>
Performed cybersecurity threat analysis using Neo4j graph databases, Bloodhound AD mapping, and CVE data processing	

C3i Hub, IIT Kanpur	<i>Jul'21-Jul'22</i>
Integrated Drupal CMS with PHP Composer, built React.js Threatmap frontend using React.js and Express APIs	

## KEY PROJECTS

Reel2Trip	<i>March'25</i>
🕒 ixigo Hackweek'25 - <i>Live on Production</i>	
• Built <b>FastAPI</b> backend integrating <b>Google Gemini 2.0</b> and OpenAI GPT-4 for video-to-itinerary conversion system	
• Developed Instagram DM method using Facebook <b>Graph API</b> with real-time <b>webhook</b> processing, currently live production	
• Implemented intelligent caching layer reducing AI processing costs by 70% through reel deduplication and optimization	
• Created context aware chat system with <b>ZEP memory</b> management enabling iterative itinerary editing and personalization	
MCP Framework (Auth based)	<i>Jul'25 - present</i>
👤 Self Project	
• Architected a production ready <b>OAuth 2.0</b> framework for <b>MCP servers</b> , reducing implementation complexity by <b>90%</b> for DX	
• Built enterprise authentication system with PKCE flow, AES-256 encryption, automated token refresh, and multi-provider support	
• Published open-source npm package with comprehensive documentation, TypeScript definitions, and developer-friendly APIs following modular architecture	
• Led e2e development of first-to-market OAuth solution for MCP ecosystem, establishing production standards and compatibility	
ISRO Web Based X-RAY Burst Automation	<i>Feb'22 - Apr'22</i>
📅 Inter IIT Tech Meet 10.0	
• Developed <b>web automated system</b> for identifying <b>solar</b> bursts, utilising the <b>ReactJS framework</b> with data segregation features	
• Organised curve data by saving specific details in <b>react hooks</b> in <b>JSON</b> to ensure smooth management and API calls to the server	
• Led end-to-end <b>product</b> documentation and orchestrated the development of a user journey and research, ensuring project success	
• Won a <b>Silver Medal</b> in the <b>mid-prep</b> event, with IIT Kanpur's contingent securing the <b>2nd position</b> in the <b>overall Meet</b>	
Journal Scrapper for Data Mining	<i>Jul'22 - Dec'22</i>
Supervised By- Proff. Shikhar Krishan Jha	
• <b>Automated</b> system for fetching <b>.ris/.bib</b> files, extracting author details programmatically, and optimizing data retrieval processes	
• Proficiently conducted <b>NER (Named Entity Recognition)</b> to extract author entities and <b>DOI URLs</b> from scholarly articles	
• Used Scopus <b>Cited By API</b> to implement citation counts of <b>input journals</b> , enhancing our tool's functionality & data accuracy	
• Seamlessly connected <b>Flask app</b> and frontend was established through <b>HTTP</b> or other <b>REST-based API</b> calls via backend	
SKILLS	

**Languages:** TypeScript, JavaScript, Python, Elixir, C++

**Backend & Realtime Systems:** Node.js, FastAPI, WebRTC, WebSockets, SSE, RPC

**AI / ML:** Agent Orchestration, Multi-Agent Systems, RAG, Semantic Caching, LLM Observability, Memory Layer, Vector DBs

**Databases:** PostgreSQL, MongoDB, SQLite, Neo4j

**Infra:** Docker, CI/CD, AWS (EC2, RDS, S3), nginx

## POSITIONS OF RESPONSIBILITY

Head Web & App, Techkriti'23	<i>Jul'22 - May'23</i>
• Led a 400+ cross-functional team across engineering, finance, design, and operations, ensuring alignment and smooth execution	
• Re-architected server infrastructure and built 10+ unified portals, improving scalability, reliability, and user experience	
• Managed 350K+ monthly traffic and 30K+ daily users, maintaining a 15K+ user database with consistent performance	