

Vicky Suman

Sr. Data Science Engineer

April 6, 1992

House No.- 4018-3F, Block-H, Ansal Versalia, Sector 67, Gurugram, Haryana 122101

+91 9711404468

jatolia.vicky8@gmail.com

<https://www.linkedin.com/in/vicky-suman-3a22a0123/>

<https://vsuman-ai.github.io/>

<https://medium.com/@vickysuman>

Education

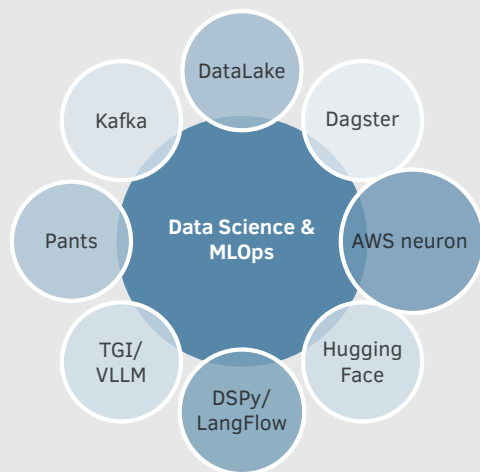
Indian Institute of Technology Bombay
M.Sc. Applied Statistics and Informatics
2015 – 2017

P.G. D.A.V. College, University of Delhi
B.Sc. Statistics (Hons)
2012 – 2015

About Me

AI and Data Engineering leader architecting scalable, production-grade AI platforms. I lead teams of 3 pax, define technical direction and deliver high-impact LLM and ML systems that translate complex data into measurable business value.

Platform & Tools



Working Experience

November, 2022 – present

Sr. Data Science Engineer

Ontic Technologies, Noida

Leading the architecture and delivery of enterprise scale AI platforms across LLM systems, NLP, computer vision, and distributed data infrastructure. Translating complex technical systems into reliable products with measurable business impact.

- ★ **Global Risk Event Intelligence Platform:** Architected and launched a revenue-generating AI-powered Risk Event Intelligence platform (~ 7.5K per client) and processing 10K+ messages per hour.
- Implemented **embedding-based de-duplication and trend clustering with summarization**, reducing alert noise and enabling event-level intelligence aggregation.
- Developed a real-time **embedding platform** (Kafka + Triton Inference Server + gRPC), storing vectors in Qdrant/MongoDB to power semantic search or clustering for trends and de-duplication.
- Engineered **LLM-based contextual location extraction** to identify incident-specific addresses and implemented reverse geocoding for precise geo-mapping.
- Built a **hypothesis-driven LLM severity classification** system with few-shot prompting to prioritize security events across diverse client scenarios.
- ★ **AI Systems & Infrastructure Engineering:** Developed AI systems across LLM applications, NLP pipelines, computer vision services, and inference infrastructure, with emphasis on scalability, reliability, and operational excellence.
- Built an **Agentic AI** framework using DSPy + tool-based reasoning to automate topic creation, web extraction, and topic query generation.
- Architected a **denoising face detection pipeline** (YOLOv11 ONNX + fine-tuned ViT) deployed via Triton with batched inference, significantly reducing false-positive image matches. Eliminated redundant third-party API calls, driving **\$5K/month infrastructure cost savings**.
- Designed a **GenAI metrics** platform powered by **QLoRA** fine-tuned LLMs (served via TGI + DSPy), generating structured in backend API format.
- Built a **camera health monitoring** feature for customer support team to detect offline client cameras and trigger proactive alerts, reducing security blind spots.
- ★ **Data Platform & Architecture:** Build distributed data platforms and processing architectures for large-scale ingestion (millions of records in per hour), transformation, indexing, and retrieval to power enterprise AI and intelligence products.
- Designed and deployed a **production-grade LLM inference** using **(TGI)** on **k8s**, enabling scalable LLM serving across multiple modules.
- Developed a DSPy-based **prompt optimization** framework for LLM classification tasks using **GEPA, SIMBA, and Decision-Tree** optimizers.
- Established a **Monorepo** architecture using **Pants** enabling modular builds, dependency isolation, and faster multi-service development workflows.
- Data lake pipeline** on AWS using **Lambda, Glue (PySpark), and Hudi** on S3, enabling downstream platform use cases in real-time event processing.

