

# Kapil Wanaskar

📍 Cupertino, CA, USA    ✉️ kapilw25@gmail.com    ☎️ +16698379485    🔗 LinkedIn    🏠 Google Scholar

## 📁 Work Experience

**Amazon Web Services (AWS), LLM Optimization / Applied ML** May 2024 – present | Cupertino, CA, USA

- Enabled MoE (DBRX, Mistral) with stability-aware expert routing; improved LLM training efficiency by 18%
- Built eval pipelines for MoE variants (dropless/dropping); reducing convergence errors by 30%
- Optimized PyTorch-based NxD parallelism for Mixtral-style LLMs, enabling 12% faster training
- Developed activation/gradient drift tools (CPU vs. TRN); flagged 95% silent regressions

**Intuitive Surgical, Software (Machine Learning) Engineer** May 2023 – May 2024 | Sunnyvale, CA, USA

- Developed FastAPI-based ML inference pipeline with VectorDB to detect 98% precision "Unknown" attacks
- Supervised fine-tuned LLMs (Claude, Llama) via PEFT (LoRA) on few-shot human-labeled feedback
- Post-trained embedding encoders and re-indexed FAISS for similarity updates, improving security by 13%
- Fine-tuned 130,000+ unsupervised models on 150+ GB data using SageMaker + MLflow; achieved 92% precision and 99.9% accurate training inputs

**Vectorr.in, Software Engineer** Mar 2018 – Jul 2022 | Mumbai, India

- Clustered 10k+ (daily) visits stored in Snowflake database, surging customer satisfaction from 3.1 to 4.8.
- Integrated Apache Kafka and Superset to segment real-time audience data for digital marketing while training Unsupervised models on AWS EC2, amplifying ROI by 23%.
- Deployed Docker via CI/CD for automating deployment, achieving a 43% reduction in data overhead.

## 📄 Research Publications

**Multimodal Benchmarking and Recommendation of Text-to-Image Generation Models** 🔗

- Evaluated 12+ (text-to-image) models (Stable Diffusion, CogView, FLUX, etc.) with ground truth from DeepFashion Multimodal dataset for alignment
- Designed Weighted Score metric combining CLIP-Score, LPIPS, FID, MRR& Recall@3 via min-max normalization
- Integrated metadata features and CLIP embeddings to align generated with ground truth image and prompt context
- Metadata-augmented models (Flux, InContext LoRA) showed ~19% higher Weighted Score & ~15 point FID reduction

**Evaluating Nano Vision-Language Models' Robustness Against Adversarial Attacks** 🔗

- Designed a modular framework to benchmark 3B/4-bit VLMs (e.g., Qwen2.5-VL) under transfer and black-box attacks (PGD, CW, Pixel, SimBA), measuring performance degradation across 15+ adversarial methods
- Built quantized inference pipeline and evaluation suite (VQA accuracy, SSIM-aware constraints); found CW- $L^\infty$  and Pixel attacks degrade accuracy by up to 52.94%, while GeoDA surprisingly improved performance (+5.88%)

**Prompt Recommendations for AI art, IEEE AIKE, California, USA** 🔗

- Extracted features of 5000 images via text embeddings and ensemble models
- Proposed Graph-based evaluation of 3 recommendation Algorithms and Community Detection Algorithms, via analyzing absence of ratings or preference scores

**Surveillance Drone Cloud and Intelligence Service, IEEE (MobileCloud), Greece** 🔗

- Proposed a surveillance drone cloud for efficient utilization of cloud computing and real-time data sharing

**Detection of Cyber Security Threats using IOT Deep Learning** 🔗

- Suggested TensorFlow deep neural system to classify stolen programming with source code literary theft

**Real World Use of Deep Learning Models for Cyber Security in IoT Network** 🔗

- Presented Deep reinforcement learning models for cyber security in IoT (Internet of Things) networks

**Analyzing Effect of Workpiece Stiffness Variation on the Stability in Flank Milling of an Impeller Blade** 🔗

- Scrutinized FFT (fast-fourier-transform) plots, chatter boundary plots, and stability region diagrams

## 🎓 Education

**MS in Artificial Intelligence, Computer Engineering,** CA, USA  
San José State University

**Master of Computer Integrated Manufacturing and Bachelor of Engineering,** Mumbai, India  
Indian Institute of Technology (IIT) Bombay