

KAMALESH KUMAR

413-472-6246 | kamalvidya.2002@gmail.com | kkk2002 | kkamaleshkumar | kkamaleshkumar.github.io

Education

University of Massachusetts Amherst <i>Master of Science in Computer Science (GPA 4.00/4.00)</i>	Sep. 2024 – May 2026 Amherst, MA
Indian Institute of Technology (IIT) Madras <i>B.Tech in Civil Engineering, Minor in Artificial Intelligence & Machine Learning</i>	Jul. 2020 – May 2024 Chennai, India

Publications

- Breaking Free from Hand-Crafted Rewards: A Genetic Programming Framework for End-Goal-Driven Reinforcement Learning (*under review at 2026 IEEE World Congress on Computational Intelligence (WCCI)*)
- SketchCleanGAN: A generative network for improving 3D CAD model retrieval systems (*Computer & Graphics'24*) [DOI]
- SketchCADGAN: A generative approach for query sketch completion of 3D CAD models. (*Computer & Graphics'23*) [DOI]

Technical Skills

Languages: Python, Rust, C++, MATLAB, L^AT_EX, C, SQL

Libraries: PyTorch, TensorFlow, JAX, vLLM, triton, transformers, VeRL, DeepSpeed, TensorRT-LLM, OpenCV

Technologies/Frameworks: LangChain, ROS, Ray, Linux, Git, Spark, Hadoop, Kubernetes, Docker, Azure

Experience

Machine Learning Intern <i>KLA Corporation</i>	May 2025 – Aug. 2025 Milpitas, California
<ul style="list-style-type: none">• Mitigated catastrophic forgetting by ~99–100%, retaining ≥95–99% of baseline defect-count performance on previously seen wafer processes during sequential fine-tuning, compared to >200–300% degradation under vanilla fine-tuning.• Reduced raw defect count by ~95–98% on previously seen wafers (~85K → ~1.5K–4K) while preserving performance on new wafers using interference-aware replay and gradient-constrained optimization.• Cut fine-tuning time by 50% (7 → 3.5 min) by freezing ~90% of model parameters, identifying variation-sensitive layers via Fisher Information Matrix over 443K parameters, with negligible loss in accuracy across old and new wafers.	
Research Intern (RL) <i>Mitacs Globalink</i>	May 2024 – Aug. 2024 Antigonish, Canada
<ul style="list-style-type: none">• Developed a genetic programming framework that improved agent fitness scores by up to 218% in high-dimensional MuJoCo environments (e.g., Humanoid, HalfCheetah) compared to standard human-engineered reward functions.• Optimized PPO's learning efficiency, enabling agents to reach peak performance in ~200,000 time-steps versus the >800,000 required by the baseline, effectively reducing training time by 80%.• Validated the statistical significance of the results ($p < 10^{-4}$) across the 11 tested environments, and surpassed all competing baselines in 82% of tasks, and achieved 22× gains in task alignment coefficient (TAC) over the PPO baseline.	

Projects

Cost-efficient Agentic LLM Workflows via Reinforcement Learning UMass Amherst	Dec. 2024 – present
<ul style="list-style-type: none">• Ideated a Pareto-frontier-based framework to optimize agentic LLM workflows under accuracy-latency constraints.• Enabling query-adaptive workflow selection by composing sub-agent cost-accuracy trade-offs and cost-aware RL.	
Continual Reinforcement Learning with Average Reward Criterion UMass Amherst	Feb. 2025 – May 2025
<ul style="list-style-type: none">• Investigated non-stationary environments in reset-free, continual RL settings requiring lifelong agent adaptation.• Showed theoretical connections with average-reward POMDPs for modeling partial observability in infinite-horizon tasks.	
Autonomous Object Following Robot using ROS and DeepSORT UMass Amherst	Feb. 2025 – May 2025
<ul style="list-style-type: none">• Built a ROS-based object-following robot using YOLO-v3 and DeepSORT for real-time tracking and re-identification.• Designed a Dockerized ROS Noetic environment on Triton enabling CUDA-accelerated inference, and real-time control.	
Real-Time Fake News Detection in Articles Using Apache Flink UMass Amherst	Sep. 2024 – Dec. 2024
<ul style="list-style-type: none">• Developed a real-time streaming pipeline with Apache Flink and ONNX-optimized DistilBERT for fake news detection.• Optimized system performance for throughput, latency, fault tolerance, and resource efficiency in a scalable deployment.	
Improving Sketch Queries for Robust Retrieval of 3D CAD Models IIT Madras	Aug. 2022 – Dec. 2023
<ul style="list-style-type: none">• Designed a two-stage cascaded GAN architecture to facilitate sketch completion of incomplete query sketches.• Proposed a novel three-branch factorization based on conditional Wasserstein Generative Adversarial Network (GAN) to clean defective sketches and thus improvised a dataset of 58000 CAD sketches.	