

EXPERIENCE

Kognitos Inc.	San Jose, CA
Machine Learning Engineer Intern	Jun. 2025 - Sep. 2025
<ul style="list-style-type: none"><li>Researched and prototyped an <b>agentic AI pipeline</b> for <b>SOP generation</b> from desktop recordings, integrating <b>key-frame extraction</b> (CV-based), <b>OCR</b>, and <b>multi-agent reasoning</b>; benchmarked <b>GPT-4o</b>, <b>GPT-4.1</b>, <b>GPT-5</b>, and <b>Gemini-2.5</b>.</li><li>Built and deployed a lightweight <b>multimodal web app</b> (video, images, documents) for automated <b>SOP and flow-diagram generation</b>, <b>containerized with Docker</b> and deployed on <b>AWS EC2</b>; enabled internal testing across diverse inputs.</li><li>Developed and productionized a <b>PII masking solution</b> using <b>GLiNER</b> model (supports 40+ entity types), evolving from a <b>REST API</b> on <b>EC2</b> (PoC) to a scalable <b>AWS Lambda</b> service with <b>Docker packaging</b> and <b>role-based access control</b>.</li></ul>	
Samsung R&D Institute India - Bangalore (SRI-B)	Bengaluru, India
Lead Engineer, Machine Learning	Mar. 2023 - Aug. 2024
<ul style="list-style-type: none"><li>Led the development of an <b>automated Python-based Android profiling tool</b> to benchmark rendering uniformly across devices, identifying bottlenecks (scroll janks) and enhancing rendering performance through improvements in the Android framework source.</li><li>Led the development of four <b>edge-based personalization LLM solutions</b>, deriving actionable insights to enhance user experience</li><li>Fine-tuned the <b>FLAN-T5 LLM model</b> for our specific use-case, optimizing performance &amp; achieving a <b>~92%</b> evaluation score.</li></ul>	
Senior Software Engineer, Machine Learning	Mar. 2021 - Feb. 2023
<ul style="list-style-type: none"><li>Designed an <b>edge ML solution</b> to analyze phone usage data and detect boredom with <b>~80% accuracy</b>, enhancing user experience. Built an Android app for <b>real-time inference with under 50ms latency</b>, showcasing the model's effectiveness.</li><li>Pioneered a <b>Federated Learning (FL)</b>-based solution to predict gender and demographic age, <b>enhancing privacy</b> for <b>~10k users</b>. Explored innovative distributed learning techniques and tested diverse FL algorithms across 20+ input and model configurations, laying the groundwork for future privacy-preserving AI advancements.</li></ul>	
Software Engineer, Machine Learning	Jun. 2019 - Feb. 2021
<ul style="list-style-type: none"><li>Engineered a <b>privacy-preserving edge ML solution</b> for <b>Next App Recommendation</b>, published in <b>IEEE ICSC 2022</b>, by designing a memory-efficient model (<b>99% size reduction</b>) to minimize FL bandwidth costs. Trained and deployed the model in Java using <b>DL4J</b>, integrating it on Android edge devices with a User Trial (UT) app for <b>training and inference across 500+ devices</b>.</li></ul>	

TECHNICAL SKILLS

Programming Languages: Python, C/C++, Java, Kotlin, SQL, Golang, Shell Script (Bash)
Frameworks, Libraries & Misc.: TensorFlow, PyTorch, PydanticAI, Scikit-learn, FastAPI, GitHub, DeepLearning4Java (DL4J)

EDUCATION

University of California, San Diego (UCSD)	La Jolla, CA
Master of Science in Computer Science; GPA: 3.95 / 4.0	Sep. 2024 - Present
Indian Institute of Technology, Kanpur (IIT Kanpur)	Kanpur, India
Bachelor of Technology in Computer Science and Engineering; GPA: 9.0 / 10.0	Jul. 2015 - Jun. 2019

PATENTS AND PUBLICATIONS

Methods and Electronic Devices for Behavior Detection using Federated Learning	2023
US 18/191403	
<ul style="list-style-type: none"><li>Novel methodology to identify new behavior trends across users and provide behavioral recommendations using Federated Learning</li></ul>	
System and Method for Distributed Learning of Universal Vector Representations on Edge Devices	2023
US 17/946349	
<ul style="list-style-type: none"><li>Novel framework for learning user behavior embeddings directly on edge devices using Federated Learning</li></ul>	
Memory Efficient Federated Recommendation Model	2022
IEEE 16th International Conference on Semantic Computing (ICSC)	
<ul style="list-style-type: none"><li>Proposed novel model framework for Federated Learning-based recommendation systems which uses a fixed-size encoding scheme for items and users, thus not affected by the number of items/users, hence memory efficient and well suited for large-scale applications.</li></ul>	

RELEVANT PROJECTS

Multi-task Learning with ToolkenGPT Framework   Course Project, UCSD	Sep. 2024 - Dec. 2024
<ul style="list-style-type: none"><li>Re-implemented ToolkenGPT, optimizing efficiency for smaller Llama models and testing multi-task performance.</li></ul>	

AWARDS AND ACHIEVEMENTS

Samsung Excellence Award   SRI-B	2024
<ul style="list-style-type: none"><li>Recognized as <b>Star of the Quarter</b> for excellent contributions in projects</li></ul>	
Key Talent Recognition Program Award   SRI-B	2023
<ul style="list-style-type: none"><li><b>Globally recognized</b> for exemplary teamwork and significant contributions to projects and organizational goals</li></ul>	