

SOM SAGAR

Email: ssagar6@asu.edu | Portfolio: somsagar07.github.io | LinkedIn: linkedin.com/in/somsagar

EDUCATION	Arizona State University <i>Ph.D. in Computer Science</i> • Advisor: Ransalu Senanayake • GPA : 3.9/4.0 • Relevant Coursework : Natural Language Processing, Data Mining, Planning Learning Methods in AI, Statistical Machine Learning, Knowledge Representation, Reinforcement Learning	Tempe, Arizona Aug. 2023 - Present
	Indian Institute of Information Technology (IIIT) <i>B.Tech (Honors) in Computer Science</i> • CGPA : 8.72/10.0 • Relevant Coursework : Machine Learning, Deep Learning, Python, Object Oriented Programming, Linear Algebra, Big Data, Data Structures and Algorithm, Data Warehousing and Data Mining, Applied Predictive Analysis, Probability and Statistics, Calculus I, II	Kottayam, Kerala Aug. 2019 - May 2023
RESEARCH INTEREST	Agentic AI, Reinforcement Learning, Foundation Models, Failure Detection and Mitigation, Generative AI, Explainability	
PREPRINTS & PUBLICATIONS	*denotes equal contribution <ol style="list-style-type: none">Atharva Gundawar*, Som Sagar*, Ransalu Senanayake. PAC Bench: Do Foundation Models Understand Prerequisites for Executing Manipulation Policies? <i>Conference on Neural Information Processing Systems (NeurIPS)</i>, 2025.Som Sagar*, Aditya Taparia*, Harsh Mankodiya, Pranav Bidare, Yifan Zhou, Ransalu Senanayake. Trustworthy Explanations for Robot Behaviors. <i>International Conference on Intelligent Robots and Systems (IROS)</i>, 2025.Aditya Taparia, Som Sagar, Ransalu Senanayake. Explainable Concept Generation through Vision-Language Preference Learning for Understanding Neural Networks' Internal Representations <i>International Conference on Machine Learning (ICML)</i>, 2024.Som Sagar, Jiafei Duan, Sreevisakh Vasudevan, Heni Ben'Amor, Dieter Fox, Ransalu Senanayake. From Mystery to Mastery: Failure Diagnosis for Improving Manipulation Policies. <i>RSS Workshop on Out-of-Distribution Generalization in Robotics</i>, 2025.Som Sagar, Aditya Taparia, Ransalu Senanayake. Failures Are Fated, But Can Be Faded: Characterizing and Mitigating Unwanted Behaviors in Large-Scale Vision and Language Models. <i>International Conference on Machine Learning (ICML)</i>, 2024. (Spotlight)Som Sagar*, Aditya Taparia*, Harsh Mankodiya, Pranav Bidare, Yifan Zhou, Ransalu Senanayake. Trustworthy Conceptual Explanations for Neural Networks in Robot Decision-Making. <i>NeurIPS Workshop on Safe and Trustworthy Agents</i>, 2024.Som Sagar, Aditya Taparia, Ransalu Senanayake. LLM-Assisted Red Teaming of Diffusion Models through "Failures Are Fated, But Can Be Faded" <i>NeurIPS Workshop on Red Teaming GenAI: What Can We Learn from Adversaries?</i>, 2024.Aditya Taparia, Som Sagar, Ransalu Senanayake. Explainable Concept Generation through Vision-Language Preference Learning. <i>NeurIPS Workshop on Interpretable AI: Past, Present and Future</i>, 2024.Joshua Tint, Som Sagar, Aditya Taparia, Caleb Liu, Kelly Raines, Bimsara Pathiraja, Ransalu Senanayake. ExpressivityArena: Can LLMs Express Information Implicitly?. <i>NeurIPS Workshop on Behavioral Machine Learning</i>, 2024.	

10. **Som Sagar**, Swani Sundara Didde, Cinu S Killilor. Embedding Based Analysis for Simplification of Legal Terms. *International Conference on Computing Science, Communication and Security*, 2023.

EXPERIENCE	Research Assistant	Aug. 2023 - Present
	Laboratory for Learning Evaluation of autoNOMous Systems (LENS Lab), ASU <ul style="list-style-type: none"> • Conducting research at the intersection of reinforcement learning, foundation models, and robotics, with a focus on improving model adaptability and robustness in real-world applications. • Working on developing a framework that enhance the interpretability and trustworthiness of AI systems in dynamic environments. • Collaborating with interdisciplinary teams to address key challenges in explainability, preference learning, and failure detection in machine learning models. 	
	Machine Learning Research Intern	May 2025 - Aug. 2025
	LinkedIn Corporation, Agents Platform Team <ul style="list-style-type: none"> • Built intelligent preference learning tool by combining reinforcement learning (RL), preference optimization, and supervised fine-tuning (SFT) to train domain-specialized agents for real-world tasks. • Improved multi-agent communication by designing faster, more effective agent-to-agent interaction protocols and scalable training pipelines. • Enhanced agent alignment through structured reward signals and human feedback, enabling generative agents to collaborate and adapt across diverse environments. 	
AWARDS AND HONORS	• Spotlight (Top 3 %) , International Conference on Machine Learning	2024
	• ASU Graduate & Professional Student Association Travel Award , Arizona State University	2024
	• Interpretable AI NeurIPS Travel Grant , NeurIPS	2024
	• Graduate College Travel Award , Arizona State University	2024
	• SCAI Conference Award , School of Computing and Augmented Intelligence	2024
	• Prime Minister Scholarship , Government of India	2019-23
	• Inter IIT Hackthon Winner , Indian Institute of Information Technology	2022
SERVICE	Reviewer: <i>International Conference on Learning Representations (ICLR)</i> 2026, <i>Conference on Robot Learning (CoRL)</i> 2025, <i>International Conference on Intelligent Robots and Systems (IROS)</i> 2024, <i>International Conference on Learning Representations (ICLR)</i> 2025, <i>Conference on Neural Information Processing Systems (NeurIPS)</i> 2024, <i>International Conference on Intelligent Robots and Systems (IROS)</i> 2025	
TEACHING	• Instructor , FSE 100 : Introduction to Engineering, ASU	Fall 2023, 2024, 2025
	• Teaching Assistant , CSE 598 : Operational Deep Learning, ASU	Spring 2024
	• Teaching Assistant , CSE 100 : Introduction to C++, ASU	Spring 2024
SKILLS	Programming Languages: Python, C, C++, Dart, JavaScript	
	Frameworks: PyTorch, NumPy, Pandas, Captum, Stable Baselines, Diffusers, NLTK, Gymnasium, Gradio, TensorFlow, Sckit-learn, Keras, TRL	
	Simulation and Environment Tools: MuJoCo, Isaac Gym/Sim/Lab, CARLA, OpenAI Gym, RLBench	
	Databases and Cloud Services: MySQL, AWS, Firebase	
	Development Tools: Visual Studio Code, Spyder, Android Studio, Git, Docker.	