

# SOM SAGAR

ssagar6@asu.edu

<https://somsagar07.github.io>

EDUCATION	<b>Arizona State University</b> <i>Ph.D. in Computer Science</i> • Advisor: Ransalu Senanayake • GPA : 4.0/4.0 • Relevant Coursework : Natural Language Processing, Data Mining, Planning, Learning Methods in AI, Statistical Machine Learning <b>Indian Institute of Information Technology (IIIT)</b> <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	Tempe, Arizona Aug. 2023 - Present  Kottayam, Kerala Aug. 2019 - May 2023
RESEARCH INTEREST	Deep Reinforcement Learning, Foundation Models, Failure Detection and Mitigation, Generative AI, Explainability	
PREPRINTS & PUBLICATIONS	*denotes equal contribution <ol style="list-style-type: none"><li>1. <b>Som Sagar</b>, 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. (<b>Spotlight</b>)</li><li>2. <b>Som Sagar*</b>, 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.</li><li>3. <b>Som Sagar</b>, 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.</li><li>4. Aditya Taparia, <b>Som Sagar</b>, Ransalu Senanayake. Explainable Concept Generation through Vision-Language Preference Learning. <i>NeurIPS Workshop on Interpretable AI: Past, Present and Future</i>, 2024.</li><li>5. Joshua Tint, <b>Som Sagar</b>, 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.</li><li>6. <b>Som Sagar</b>, Swani Sundara Didde, Cinu S Killilior. A Sentiment Word2Vec Approach for Simplification of Legal Terms. <i>International Conference on Computing Science, Communication and Security</i>, 2023.</li></ol>	
EXPERIENCE	<b>CS Researcher</b> Laboratory for Learning Evaluation of autoNomous Systems (LENS) • 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.	Aug. 2023 - Present

AWARDS AND HONORS	• <b>Spotlight (Top 3 %)</b> , International Conference on Machine Learning	2024
	• <b>Graduate College Travel Award</b> , Arizona State Univeristy	2024
	• <b>SCAI Conference Award</b> , School of Computing and Augmented Intelligence	2024
	• <b>Prime Minister Scholarship</b> , Government of India	2019-23
	• <b>Inter IIT Hackthon Winner</b> , Indian Institute of Information Technology	2022
SERVICE	<b>Reviewer for:</b> <i>International Conference on Learning Representations (ICLR) 2025</i> , <i>Conference on Neural Information Processing Systems (NeurIPS) 2024</i> , <i>International Conference on Intelligent Robots and Systems (IROS) 2024</i> ,	
TEACHING	• <b>Instructor</b> , FSE 100 : Introduction to Engineering, ASU	Fall 2023, Fall 2024
	• <b>Teaching Assistant</b> , CSE 598 : Operational Deep Learning, ASU	Spring 2024
	• <b>Teaching Assistant</b> , CSE 100 : Introduction to C++, ASU	Spring 2024
SKILLS	<b>Programming Languages:</b> Python, C, C++, Dart, JavaScript. <b>Frameworks:</b> PyTorch, NumPy, Pandas, Captum, Stable Baselines, Diffusers, NLTK, Gymnasium, Gradio, TensorFlow, Sckit-learn, Keras. <b>Simulation and Environment Tools:</b> MuJoCo, CARLA, OpenAI Gym, RLBench, Coppeliasim <b>Databases and Cloud Services:</b> MySQL, AWS, Firebase <b>Development Tools:</b> Visual Studio Code, Spyder, Andriod Studio, Git, Docker.	