

# Chahat Raj

## PhD in Computer Science | George Mason University

🌐 [chahatraj.github.io](https://chahatraj.github.io) @ [craj@gmu.edu](mailto:craj@gmu.edu) 🌐 [github.com/chahatraj](https://github.com/chahatraj) 🎓 Google Scholar

### Research Interests

I work on Responsible AI, Ethics, and Fairness, specifically in the evaluation and mitigation of socio-cultural biases within multilingual and multimodal NLP applications with a focus on LLMs. I have recently published at EMNLP, AIES, and ECIR.

### Education

<b>Present</b> <b>Aug 2022</b>	<b>George Mason University</b> Ph.D. in Computer Science Advisors: <a href="#">Ziwei Zhu</a> , <a href="#">Antonios Anastasopoulos</a>	<b>Fairfax, USA</b>
<b>Aug 2021</b> <b>Aug 2019</b>	<b>Delhi Technological University</b> Masters in Information Systems (Research Track) Advisor: <a href="#">Priyanka Meel</a>	<b>Delhi, India</b>
<b>May 2019</b> <b>Aug 2015</b>	<b>Indira Gandhi Delhi Technical University for Women</b> Bachelors in Computer Science & Engineering	<b>Delhi, India</b>

### Experience

<b>Sep 2024</b> <b>May 2024</b>	<b>University of Washington   The Information School</b> Researcher / Advisor: <a href="#">Aylin Caliskan</a> Project: Social Biases in Visual Question Answering in Generative AI.	<b>Seattle, WA</b>
<b>Aug 2024</b> <b>May 2024</b>	<b>George Mason University   School of Computer Science</b> Graduate Research Assistant / Advisors: <a href="#">Ziwei Zhu</a> , <a href="#">Antonios Anastasopoulos</a> Project: Exploring Stereotypical Associations in Large Vision-Language Models through generative tasks - image generation and question answering.	<b>Fairfax, VA</b>
<b>Aug 2023</b> <b>May 2023</b>	<b>George Mason University   School of Computer Science</b> Graduate Research Assistant / Advisor: <a href="#">Ziwei Zhu</a> Project: Fairness and Interpretability of Transformer Models.	<b>Fairfax, VA</b>
<b>Jul 2022</b> <b>Jun 2021</b>	<b>University of Technology Sydney   School of Computer Science</b> Visiting Scholar / Advisor: <a href="#">Mukesh Prasad</a> Projects: Disaster Fundraising on Social Media Analysis with the Australian Red Cross, AI integration in aboriginal and indigenous communities' lifestyle, Web information pollution detection (cyberbullying, hate, and offensive speech), Machine learning-based decision-making in clinical vertigo diagnosis.	<b>Remote</b>
<b>Aug 2021</b> <b>May 2021</b>	<b>Indian Institute of Management Raipur</b> Researcher / Advisor: <a href="#">Manojit Chattopadhyay</a> Projects: Cryptocurrencies' price prediction, Analysing causal relationships between stock market behavior and societal events.	<b>Remote</b>
<b>Aug 2022</b> <b>Aug 2019</b>	<b>Delhi Technological University</b> Graduate Student Researcher / Advisor: <a href="#">Priyanka Meel</a> Projects: Developing neural network frameworks to detect antisocial web content (fake news, infodemic, misinformation, disinformation, rumor, death hoax, and clickbait), Infodemic impact analysis.	<b>Delhi, India</b>

### Selected Publications

- [C.5] **BiasDora: Exploring Hidden Biased Associations in Vision-Language Models** [PDF | Code]  
[Chahat Raj](#), Anjishnu Mukherjee, Aylin Caliskan, Antonios Anastasopoulos, Ziwei Zhu  
Conference on Empirical Methods in Natural Language Processing [EMNLP'24]
- [C.4] **Breaking Bias, Building Bridges: Evaluation and Mitigation of Social Biases in LLMs via Contact Hypothesis** [PDF | Code]  
[Chahat Raj](#), Anjishnu Mukherjee, Aylin Caliskan, Antonios Anastasopoulos, Ziwei Zhu  
AAAI/ACM Conference on Artificial Intelligence, Ethics, and Society [AIES'24]

- [C.3] **SALSA: Saliency-Based Switching Attack for Adversarial Perturbations in Fake News Detection Models** [PDF | Code]  
 Chahat Raj\*, Anjishnu Mukherjee\*, Hemant Purohit, Antonios Anastasopoulos, Ziwei Zhu  
*European Conference on Information Retrieval* [ECIR'24]
- [C.2] **Global Voices, Local Biases: Socio-cultural Prejudices across Languages** [PDF | Code]  
 Anjishnu Mukherjee\*, Chahat Raj\*, Ziwei Zhu, Antonios Anastasopoulos  
*Conference on Empirical Methods in Natural Language Processing* [EMNLP'23]
- [C.1] **True and Fair: Robust and Unbiased Fake News Detection via Interpretable Machine Learning** [PDF | Code]  
 Chahat Raj, Anjishnu Mukherjee, Ziwei Zhu  
*AAAI/ACM Conference on Artificial Intelligence, Ethics, and Society* [AIES'23]

## Skills

Natural Language Processing, Generative AI, Multimodal Large Language Models, Model Training, Evaluation & Validation, Instruction-tuning, Prompt Engineering, LLM-human alignment, Vision, Ethics and Bias Mitigation, Explainability.

## Talks

### “Breaking Bias, Building Bridges”

› Talk at AAAI/ACM AIES October 2024 (San Jose, USA)

### “LLM Ethics - Bias and Mitigation”

› Invited Lecture - CS 678 Advanced NLP @ George Mason University October 2024 (Virginia, USA)

### “A Psychological View to Social Biases in LLMs”

› SouthNLP @ Emory University April 2024 (Georgia, USA)

## Awards

**Travel Grant by AAAI/ACM AIES 2024 & 2023** Funded \$1000 to travel to San Jose, USA and \$1000 to travel to Montreal, Canada to attend the AIES conference.

**CAHMP Fellowship @ George Mason University, 2024** Summer research on social biases in question answering.

**Best Paper Award @ MASC-SLL 2024, Johns Hopkins University** For the paper ‘A Psychological View to Social Bias in LLMs: Evaluation and Mitigation’.

**Summer Graduate Research Award @ George Mason University, 2023** For research on ‘Explainability and Robustness of Debiasing Approaches for Transformers’.

**Research Excellence Award @ Delhi Technological University 2023 (×2) & 2022** Received three Commendable Research Awards with a total cash prize of INR 150k by DTU, Delhi, for the research published in reputed scientific journals.

**Graduate Scholarship @ Delhi Technological University 2019** Received INR 297k by AICTE for qualifying GATE.

## Teaching

**CS 108 Introduction to Computer Programming, GMU** Teaching Assistant Spring'24

› Responsible for teaching labs, one-to-one student tutoring, creating lab modules, grading, and exam invigilation.

**CS 478 Natural Language Processing, GMU** Teaching Assistant Fall'23

› Delivered ad-hoc lectures, student tutoring, teaching labs, and assignment grading.

**COMP 502 Mathematical Foundations of Computing, GMU** Teaching Assistant Fall'23

› Responsible for one-to-one student tutoring sessions, and assignment grading.

**CS 112 Introduction to Python Programming, GMU** Teaching Assistant Fall'22, Spring'23

› Responsible for teaching labs, one-to-one student tutoring, creating lab modules, assignment grading, and exam invigilation.

## Academic Service

**Reviewer** ACM Transactions on Intelligent Systems and Technology 2024, LTEDI @ EACL 2024, TrustNLP 2024, CSCW 2024, NeurIPS 2023, Elsevier & Springer Journal papers

**Sub-Reviewer** EMNLP'24, AIES'24