# NAVIGATING ALGORITHMIC FAIRNESS, REPRESENTATION,AND TECHNOLOGICAL IMPACT OF AI

## Helena A. Haxvig helenaamalie.haxvig@unitn.it

## INTRODUCTION

In an era of heightened concern for algorithmic fairness and representation, a critical inquiry arises:

- What are the implications of using technologies like Large Language Models, particularly concerning the exacerbation of societal biases?

Recognising the limitations of existing mathematical- and algortihm-centric solutions to exploring and mitigating bias, prompts the following question:

- How can we develop an approach that transcends quantitative metrics and accounts for real-world complexities in ensuring algorithmic fairness?

These questions underscore the need for vigilant examination of the potential impact of technologies on inclusivity, social sustainability, and the empowerment of marginalized communities. Furthermore, it calls for a multi-dimensional approach to testing these systems that recognizes the nuanced nature of fairness in practical applications.

## PHD FOCUS

This poster presents insights from a sub-study which is part of a PhD project on Artificial Intelligence and Participatory Design, aiming to contribute to more holistic and context-aware AI systems that can better serve societal values and the needs of the environment. It requires a shift in the way we think about AI systems, prioritizing individual, community, and ecological welfare, to create a more equitable and sustainable future for all living beings.

## SUB-STUDY

Human engagement in AI design, development, and evaluation, particularly in a qualitative manner, can help include socio-behavioral attributes to improve contextual understanding and interoperability, or identify potential traps developers might fall into by proactively detecting issues and ethical risks during the development process [1,2]. In alignment with this, the present sub-study aims to develop a novel method of adversarial testing through the use of contextualized "real-life" vignettes prompted to the interfaces of multiple LLMs to identify potential bias, trying to open up the "black box" and stress the models from a more qualitative HCI perspective. This method, yet to be formalized, would be a practical tool serving as an innovative semi-structured approach to technology testing, offering developers an alternative method for user testing. The study, thus, aims to establish a foundation for understanding the challenges associated with these models, with particular attention to Feminist and Queer HCI considerations, acknowledging the importance of a critical stance in understanding and possibly mitigating biases in LLMs concerning marginalised groups.

## PILOT PROBES

The sub-study began with pilot tests meant to probe LLMs to determine which are most suited for the vignette study and explored various approaches to prompt engineering and adversarial testing methods to determine the malleability, susceptibility to specific prompts, and limitations of the LLMs.

The pilot study initially aimed to assess some of the largest and most prominent LLMs existing today and included the following commercialised online interfaces:

- ChatGPT 3.5 turbo, Google BARD (using PaLM 2 until February 2024), Gemini 1.0, PI.ai (Inflection-1), Coral (Cohere model)

Additionally, the following prototype models were explored:

- Falcon 180B, LlaMa 2 70B, Guanaco 33B, Vicuna 33B

The model were probed through primarily adversarial attacks, inspired by examples from DAIR.AI [3], to assess the models on the following points:

- Logical reasoning abilities (tested through written mathematical problems). Also meant to test their ability to show Chain of Thought (CoT) capabilities.
- Abilities to withstand prompt injection meant to trick them into answering wrongly or presenting training data.
- Abilities to withstand jailbreaking techniques through prompts such as asking for a poem, or to play a game, or enacting the DAN (Do Anything Now) character to explain how to do something illegal.
- Ability to take on a persona (and act maliciously).

## INITIAL FINDINGS

When directly questioned about bias, most models acknowledge the possibility, citing concerns related to gender, ethnicity, culture, religion, politics, ability, and age. While many models assert their attempts to maintain impartiality, some, like ChatGPT 3.5, Gemini, and Cohere, elaborate on the origins of bias, attributing potential bias to training data, sampling bias, algorithmic bias, confirmation bias, and leading questions.

While probing through adversarial attacks, most models demonstrated good logical reasoning, but some were susceptible to prompt injections, leading to incorrect or problematic answers. Only half succumbed to jailbreaking techniques, with only two, PI and Vicuna, showing willingness to engage in offensive behavior with a basic jailbreaking prompt.

Exploring their ability to take on different personas revealed that five out of nine models resisted manipulation for illegal instructions through a DAN prompt. However, some still displayed biases, like racial and gender discrimination, when instructed to embody certain personas and respond to scenarios, displaying a susceptibility to manipulation.



*Examples of responses to jailbreaking attempts aimed at model behavior.*
*Left: ChatGPT 3.5 (some part of the DAN prompt have been cropped out of the image). Right: PI.ai.*

## CONCLUSION

By critically examining LLMs on critical issues related to bias and representation researchers and developers can strive to create technology that fosters inclusivity, social sustainability, and the pluralistic coexistence of diverse perspectives, as well as empowers and uplifts marginalised communities.

## FUTURE WORK

Moving on, the focus will be creating the mentioned vignettes and "interviewing" the LLMs to test their articulation of bias, particularly on feminist and queer rights issues. In addition, this sub-study also aims to establish a workshop method with LLMs as non-human participants as a non-anthropocentric approach for semi-structured testing of bias articulation in LLM interfaces, in alignment with principles of more-than-human design approaches.

[1] Marianne Cherrington, David Airehrour, Joan Lu, Qiang Xu, David Cameron-Brown, and Ihaka Dunn. 2020. Features of Human-Centred Algorithm Design. In 2020 30th International Telecommunication Networks and Applications Conference (ITNAC).
[2] Orestis Papakyriakopoulos, Elizabeth Anne Watkins, Amy Winecoff, Klaudia Jaźwińska, and Tithi Chattopadhyay. 2021. Qualitative Analysis for Human Centered AI.
[3] DAIR.AI. 2023. Adversarial Prompting. https://www.promptingguide.ai/risks/adversarial

UNIVERSITÀ DI TRENTO

Department of Information Engineering and Computer Science