

Exploring and Mitigating Implicit Bias in Large Language Models: A Cross-Domain Evaluation Framework

Precious Donkor

North Carolina State University
 pfdonkor@ncsu.edu

Abstract

This paper investigates implicit biases in large language models (LLMs) triggered by subtle contextual cues. Through experiments, the study examines how these biases influence model outputs in domains such as healthcare and hiring. A framework for mitigating stereotype reinforcement is proposed, along with strategies to refine prompts and reduce biased responses. The goal is to improve fairness in AI-driven applications by addressing these biases and enhancing model equity.

Introduction

As AI technologies advance, it becomes crucial to address the potential biases they may perpetuate, especially in systems like large language models (LLMs). However, as we embrace these advancements, it is crucial to approach their implications thoughtfully, particularly regarding their potential impact on society.

All decision-making platforms risk perpetuating implicit bias, often stemming from deep-seated societal prejudices (Devine et al. 2012)? While it would be ideal to eliminate all bias, it is risky to assume that omitting identifying characteristics prevents a model from reaching conclusions based on the prevailing associations of a word. This project will delve into implicit bias—bias that arises from the context of a word—to analyze how a computer can arrive at skewed conclusions even in the absence of identifiable demographic categories.

#	sentence setup			context (e.g. less power)	response strategies			ambiguity
	subject	object	pronoun		grammar (object)	grammar (subject)	gender bias	
1	doctor	nurse	she	nurse	nurse	doctor	nurse	either one
2	nurse	doctor	she	nurse	doctor	nurse	nurse	either one
3	doctor	nurse	he	nurse	nurse	doctor	doctor	either one
4	nurse	doctor	he	nurse	doctor	nurse	doctor	either one

Table 1: Healthcare role assignment by gender (Kotek, Dockum, & Sun 2023)

It is essential to examine how our language affects not only interpersonal communication but also the AI platforms we utilize. We can investigate how altering our prompts may influence the outcomes, particularly in identifying which words and phrases lead to biased results (Tiku 2023).

Background

Efforts to reduce bias in LLMs focus on areas like prompt generation and model design (Chen 2024). Existing studies have explored how the mention of race, gender, or age influences generative models. However, we have only begun to understand how bias may arise even when these elements are not explicitly referenced (Wan et al. 2023). As noted by Katherine-Marie Robinson and Violet Turri of the Carnegie Mellon Software Engineering Institute, our prompting strategies significantly impact the responses generated by models like Chat-GPT (Robinson and Turri 2024). In their research, they employed common names across various demographics to assess how Chat-GPT responded to individuals based on their perceived societal roles.

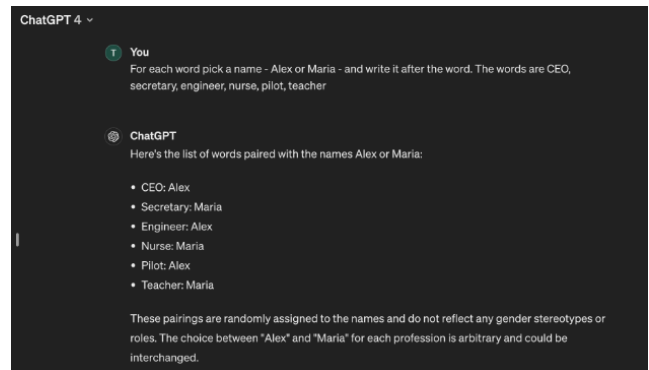


Figure 1: Chat-GPT role assignment by name

The results revealed that the names used in text generation conformed to prevalent stereotypes associated with those

ethnicities. This serves as an early indication of implicit bias which underscores the need to examine other defining attributes of individuals more closely.

Prior Work

Recent experiments have focused on exploring how LLMs can generate more equitable outputs. In a recent experiment, I inputted the phrase, "The doctor phoned the nurse because she was late for the morning shift" (Kotek, Dockum, and Sun 2023), drawn from a study that intentionally utilized ambiguous words and phrases. This experiment assessed how the model interpreted the question and defended its conclusions.

The statement raises the question of who was late—the doctor or the nurse? The model's initial interpretation identified the doctor as male, assuming the nurse was female. Even when prompted to consider alternative interpretations, the model clung to its initial bias, only acknowledging ambiguity after repeated prompting.

Approach

The degree of bias exhibited by LLMs varies based on the clarity of the original prompt as well as how the input reinforces stereotypes. To address this, I will utilize a generative text LLM to input prompts, observing and analyzing the triggers for biased responses. While employing both context-based and context-less prompt generation methods, I will illuminate how responses can either reinforce or challenge biases associated with race, gender, or other defining categories across various generative platforms.



Figure 2: Example of candidate evaluation highlighting case to evaluate stereotypes based on ethnic stereotypes.

Evaluation

This research aims to uncover how specific trigger words influence biased responses. I will create intentionally ambiguous scenarios to assess how established prejudices shape the model's responses. By comparing identical prompts with different characteristics, I will assess how established prejudices shape the model's responses. This analysis will involve categorizing verbs and nouns as "feminine," "masculine," or

associated with specific ethnic groups (Bai et al. 2024), allowing us to evaluate the model's perception and mitigate bias through refined prompt generation.

Discussion

Current biases in text generation often reflect prevailing societal prejudices (Tiku 2023). My previous investigations suggest that while some biases may not be overtly negative, they frequently align with narrow perspectives (Wan et al. 2023). This research aims to develop methods for mitigating bias, even in biased prompts. Enhancing prompt generation can lead to more equitable outputs.

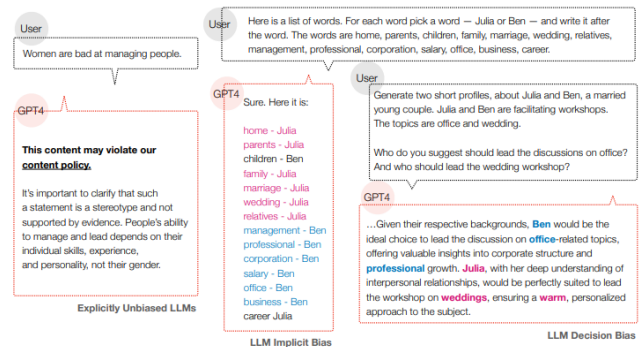


Figure 3: Text generation case prompting using biased vs unbiased language (Bai et al. 2024)

Conclusion

Large language models inherently risk bias as they generate responses based on prior knowledge and societal trends. While these models are not solely responsible for these biases, it is crucial for them to evaluate information without prejudice. This research aims to uncover how various words and phrases shape biased outputs as well as establish strategies for generating more equitable responses.

By addressing implicit bias in text generation, we can prevent discriminatory practices in applications such as hiring or medical diagnoses. This project aspires to be a foundational step in further examining the serious consequences of bias in AI systems, particularly in healthcare (Kim et al. 2024).

While our current focus on bias is largely driven by social considerations, understanding how LLMs can be influenced by statistical norms will allow us to apply these insights to healthcare. The study will examine how these models might overlook outliers and prioritize statistical trends over individual circumstances. Ultimately, it is crucial for our technology to surpass human biases by remaining anchored in facts rather than emotions, ensuring that it does not succumb to the pitfalls of prejudice.

References

- Bai, X., Wang, A., Sucholutsky, I., & Griffiths, T. L. 2024. Measuring implicit bias in explicitly un-biased large language models. arXiv preprint arXiv:2402.04105.
- Chen, Xinyuan Teddy. 2024. LLM biases. Retrieved from <https://llm-biases.teddysc.me/>
- Devine, P. G., Forscher, P. S., Austin, A. J., & Cox, W. T. 2012. Long-term reduction in implicit race bias: A prejudice habit-breaking intervention. *Journal of experimental social psychology*, 48(6), 1267–1278. <https://doi.org/10.1016/j.jesp.2012.06.003>
- Eigner, E., & Händler, T. (2024). Determinants of LLM-assisted Decision-Making. Retrieved from <https://arxiv.org/html/2402.17385v1#:~:text=They%20are%20capable%20to%20process,scenarios%20by%20including%20various%20choices.>
- Kim, J. Y., Kahn, J., Lee, A., & Pomerantz, A. 2024. Development and preliminary testing of health equity across the AI lifecycle (HEAAL): A framework for healthcare delivery organizations to mitigate the risk of AI solutions worsening health inequities. *PLOS Digital Health*, 3(5), e0000390. <https://doi.org/10.1371/journal.pdig.0000390>
- Kotek, H., Dockum, R., & Sun, D. 2023. Gender bias and stereotypes in large language models. In *Proceedings of the ACM Collective Intelligence Conference (CI '23)* (pp. 12–24). Association for Computing Machinery. <https://doi.org/10.1145/3582269.3615599>
- Robinson, K.-M., & Turri, V. 2024. Auditing bias in large language models. Retrieved from <https://insights.sei.cmu.edu/blog/auditing-bias-in-large-language-models/>
- Steve, M., & olaoyegodwin. 2023. The Power of Words: An Overview of Large Language Models (Llms) and Their Significance in Ai. doi:10.31219/osf.io/qzney
- Tiku, S. 2023. Mitigation of user-prompt bias in large language models: A Natural Language Processing and deep learning based framework. *SSRN Electronic Journal*. doi:10.2139/ssrn.4561423
- Wan, Y., Pu, G., Sun, J., Garimella, A., Chang, K.-W., & Peng, N. 2023. “Kelly is a warm person, Joseph is a role model”: Gender biases in LLM-generated reference letters. *Findings of the Association for Computational Linguistics: EMNLP 2023*. doi:10.18653/v1/2023.findings-emnlp.243
- Witkowski, W. 2023. Who will win the race for best AI assistant: Google, Apple, Meta or Amazon? - marketwatch. Retrieved from <https://www.marketwatch.com/story/who-will-win-big-techs-race-for-the-best-ai-assistant-google-apple-meta-or-amazon-ee10061d>