



Towards Fairness in AI: Examining and Addressing Bias in Large Language Models

*Samarth Parmanand*¹

¹Computer Science Graduate, Bhilai Institute of Technology Raipur, Raipur, Chhattisgarh, India

ABSTRACT:

Large language models are encountering a significant issue - Bias. That problem impacts various areas, like understanding language naturally and making decisions. Recent studies reveal that biases tied to gender, race, religion, and other social factors are widespread in these models. The size of the model alone does not decide the bias; other factors like training data and model perplexity also play a role. Some methods such as BiQ (Bias Intelligence Quotient) and UniBias have been made to gauge these biases. Instead of seeking new solutions, this review zeroes in on identifying and evaluating existing biases. By observing how gender bias is revealed and exploring social biases, it is evident that these models mirror society's preconceptions. They often exhibit more subtle biases than we previously recognized. It is vital to keep verifying and assessing these biases in large language models to guarantee that AI systems are impartial and equitable.

Keywords: Large Language Models, Fairness, Bias, Explainable AI, Model Perplexity, Social Factors.

1. Introduction

The quick advancements & easy access to Large Language Models (LLMs), along with the rise of Cheap-to-Build Very Large Language Models (CtB-LLMs), has changed natural language processing (NLP). Now, these models are key in a lot of areas. They offer incredible abilities in creating human-like text and making language tasks easier. But even though they are widely used and adaptable, LLMs, including top notch models like GPT-4, LLAMA2, and FALCON, have big problems too—especially their built-in biases. A big problem with LLMs is that they can show biases linked to gender, race, religion, and jobs. This can keep harmful stereotypes and wrong info going. These biases usually come from the training data and model setup. They show historical & societal prejudices. Studies show that it's not just the number of parameters that cause bias; stuff like perplexity matters too.

People are trying to fix these biases with different methods & frameworks. For instance, the Comprehensive Bias Neutralization Framework (CBNF) and the Bias Intelligence Quotient (BiQ) metric have been made to check & lower racial biases in LLMs. Also, techniques like Low-Rank Adaptation (LoRA) have worked well in reducing biases without hurting model performance on other tasks. Fairness in LLMs is more than just fixing bias—it's about fair treatment for all groups. This means looking at taxonomic surveys & fairness-aware rules for classifications, studying in-context learning setups, and using fairness rules when picking models. Also, methods like Uncertainty Quantification (UQ) and Explainable AI (XAI) have been suggested to spot new biases during inference. These help in making AI systems clearer & more trustworthy. Still, finding and fixing hidden biases is a tough job. New frameworks like SOFA for checking social biases in models and UniBias for spotting & removing biased parts of LLMs have made good progress. These new ways aim to find more subtle biases and offer strong solutions to make LLMs fairer. In conclusion, despite the remarkable powers demonstrated by LLMs, ethical and practical issues are raised by their innate biases. To solve these problems, more investigation into sophisticated debiasing methods and fairness frameworks is essential. We can only guarantee that LLMs are used efficiently and ethically across a range of industries by making such efforts.

2. Technical Background

In the realm of artificial intelligence, Large Language Models (LLMs) such as GPT-4, LLAMA, and ChatGPT have brought about a significant transformation. These models, constructed on sophisticated transformer architectures, possess the capability to execute a wide array of tasks ranging from text generation to language translation and sentiment analysis. Their proficiency in generating text that mimics human writing has rendered them indispensable across countless applications. However, the development and utilization of these models entail substantial expenses. Training LLMs mandates humongous volumes of data and considerable computational power, thereby often restricting accessibility to major organizations with ample

resources. This exclusivity raises concerns regarding the democratization of AI and the potential reinforcement of prevailing inequalities in technology access.

2.1 Emergence of Cheap-to-Build Very Large Language Models

The exorbitant costs and resource requirements associated with conventional LLMs have spurred the inception of Cheap-to-Build Very Large Language Models (CtB-LLMs). These models strive to enhance accessibility to advanced AI by curtailing computational and financial impediments. CtB-LLMs like LLaMA and Open Pre-trained Transformers (OPT) achieve exceptional performance with fewer parameters and reduced resource prerequisites. By optimizing model architectures and training methodologies, CtB-LLMs offer a cost-efficient alternative without compromising performance significantly. Consequently, they facilitate wider adoption and innovation within the AI domain.

2.2 Bias in Language Models

Despite their remarkable capabilities, both LLMs and CtB-LLMs may inherit biases inherent in their training data. These biases can manifest in various forms including gender, race, religion, and profession, leading to skewed outputs that mirror societal prejudices potentially reinforcing them further. The existence of bias in language models poses substantial ethical and pragmatic dilemmas. Biased outputs can perpetuate detrimental stereotypes, discriminate against marginalized groups, and undermine the credibility and impartiality of AI systems. Therefore, it is imperative to address bias in language models to uphold equity in AI applications.

2.3 Mitigation Strategies for Bias Reduction

To combat biases ingrained in language models, researchers have devised several debiasing techniques aimed at curbing the impact of biased data on model outputs while preserving overall performance integrity. One method involves incorporating anti-stereotypical sentences during training to counteract ingrained biases effectively. Another approach known as fine-tuning with Low-Rank Adaptation (LoRA) entails adjustments to specific parameters for bias minimization purposes. These techniques offer promise in reducing bias levels and enhancing fairness within language models. However, continual refinement is necessary to tackle the dynamic nature of bias within AI applications effectively.

2.4 Few-Shot Fairness Paradigm in Large Language Models

Few-shot learning has surfaced as a propitious avenue for augmenting fairness within LLMs by facilitating efficient task execution with minimal input-output examples allowing swift adaptation to new tasks and contexts. Studies indicate that models like GPT-4 can achieve heightened accuracy & impartiality when provided with additional context along with clearly defined fairness criteria. Leveraging this few-shot learning mechanism harnesses the flexible nature of LLMs thereby elevating their adaptability across diverse scenarios resulting in greater equity across various applications.

2.5 Comprehensive Bias Neutralization Framework (CBNF)

The Comprehensive Bias Neutralization Framework, or CBNF, represents a significant step forward in assessing & mitigating bias in Large Language Models (LLMs). This framework introduces a novel metric, the Bias Intelligence Quotient (BiQ), to gauge the extent of bias in model outputs. CBNF builds upon prior methodologies by incorporating diverse data sources, such as Black history & culture, to provide a more inclusive viewpoint. This approach aims to reduce biases by ensuring that training data reflects a wide range of experiences & perspectives. By concentrating on diverse data sources & robust evaluation metrics, CBNF contributes to the advancement of fairer & more accurate AI systems.

2.6 Detection of Unanticipated Bias

Unanticipated biases in large language models (LLMs) present a significant challenge for AI developers & users alike. These biases can stem from subtle, often overlooked elements of training data & model architecture. Advanced techniques like Adversarial Testing, Uncertainty Quantification (UQ), and Explainable AI (XAI) are being explored to detect and uncover these biases. UQ methods assess the certainty of model decisions, helping pinpoint areas where bias may influence outputs. XAI techniques aim to enhance the transparency of LLMs' internal decision-making, enabling researchers to understand & address non-obvious biases. By improving the transparency & interpretability of AI systems, these methods contribute to more reliable, equitable and fair AI applications.

2.7 Social Bias Probing

The Social Bias Probing framework introduces a new benchmark called SOFA (Social Fairness) to evaluate social biases in large language models (LLMs). Unlike traditional methods that use binary association tests with limited datasets, which can oversimplify social identities and stereotypes, SOFA takes a broader approach. It incorporates a wide range of identities and stereotypes, offering a more thorough analysis of bias. This framework assesses biases across various dimensions, providing insights into how different demographic groups are represented by language models. By broadening the evaluation of bias, SOFA enhances our understanding of social biases in AI and aids in creating more equitable language models.

2.8 Mitigation of Gender Bias

Gender bias in large language models (LLMs) is a significant concern, as it can reinforce harmful stereotypes and spread misinformation. Traditional methods for detecting gender bias typically rely on direct gender references or predefined stereotypes, which may not fully uncover the extent of the bias. To tackle this issue, researchers have introduced an indirect probing framework based on conditional generation. This technique uses naturally sourced and model-generated inputs to prompt LLMs to exhibit their biases without explicitly mentioning gender or stereotypes. Additionally, several mitigation strategies, including Hyperparameter Tuning, Instruction Guiding, and Debias Tuning, have been explored to minimize both explicit and implicit gender biases in AI outputs. These approaches aim to create fairer and more balanced representations in language models, supporting the ethical deployment of AI technologies.

3. Understanding Bias

Bias in Large Language Models (LLMs) can be attributed to a variety of variables, including the data used for training, the model's architecture, and human influence during fine-tuning. Training data bias arises when the data reflects historical and societal preconceptions, resulting in models that reinforce these biases. Training data bias occurs when the data reflects historical and societal prejudices, leading to models that perpetuate these biases. For example, if a dataset disproportionately features male doctors and female nurses, the LLM is likely to reinforce such stereotypes. Embedding bias is another critical issue, where word embeddings, designed to capture semantic relationships, inadvertently introduce societal biases. This can result in models associating certain professions or characteristics with specific genders or races, thereby reflecting, and reinforcing societal stereotypes. Additionally, label bias can emerge during instruction tuning, influenced by the subjective judgments of human annotators, further embedding personal and societal biases into the model. These biases have significant implications, leading to representational harms, which reinforce negative stereotypes, and allocational harms, which result in unequal resource distribution or opportunities. Detecting and mitigating these biases is challenging due to their subtle and hidden nature. Advanced debiasing techniques, such as fine-tuning with anti-stereotypical data or using comprehensive frameworks like the Bias Intelligence Quotient (BiQ), are essential. These methods focus on evaluating and reducing biases without compromising the model's performance, aiming to create fairer and more inclusive AI systems. Understanding and addressing these biases is crucial for developing ethical and responsible LLMs that align with societal values of fairness and equity.

4. Literature Review

Bias in Large Language Models (LLMs) has become increasingly prominent in recent studies, with scholars exploring various forms of bias and methods to counter them. In their comprehensive analysis, Ranaldi et al. [1] investigate biases in Cheap-to-Build Very Large-Language Models (CtB-LLMs) like LLaMA, OPT, and BLOOM. They find that biases related to gender, race, religion, and profession are prevalent and that these biases correlate more with model perplexity when compared with the number of parameters. The authors show that fine-tuning with anti-stereotypical sentences and slurs can significantly reduce bias without degrading the model's performance. For example, debiasing the OPT model using LoRA reduced the bias by up to 4.12 points in the normalized stereotype score, highlighting the critical role of well-curated training data in developing fairer models. Chu, Wang, and Zhang [2] offer a taxonomic survey on fairness in LLMs, breaking it down into metrics for quantifying biases, algorithms for mitigating biases, and resources for evaluating biases. They discuss various metrics for bias quantification, such as embedding-based, probability-based, and generation-based metrics. Their review of fairness-promoting algorithms spans pre-processing, in-training, intra-processing, and post-processing stages. They emphasize the need for a systematic survey to consolidate recent advances and stress the importance of developing a clear framework that aligns fairness notions with corresponding methodologies. The Bias Neutralization Framework (CBNF) by Narayan et al. [3] advances the detection and mitigation of racial biases in LLMs, introducing the Bias Intelligence Quotient (BiQ) metric and demonstrating the effectiveness of Latimer AI in bias detection and mitigation compared to ChatGPT 3.5. The CBNF approach is scalable and adaptable to various AI applications, underscoring the importance of continuous evaluation and adaptation to ensure equitable AI systems. Kruspe [4] addresses unanticipated biases in LLMs, proposing Uncertainty Quantification (UQ) and Explainable AI (XAI) methods for detection. UQ assesses model decision certainty, while XAI aims to make internal decision-making processes transparent. The paper highlights representational harms, which reinforce detrimental stereotypes, and allocational harms, which lead to unequal distribution of resources. Kruspe suggests that while bias mitigation strategies guide models toward reduced bias, they typically do not eliminate it entirely, stressing the importance of developing fairer and more transparent AI systems. Chhikara et al. [5] explore bias in LLMs and their potential for fairness-aware classification. They introduce a framework that incorporates fairness rules into in-context learning, showing that GPT-4 achieves superior results in both accuracy and fairness compared to Llama-70b and Gemini. Their study demonstrates that providing additional context and defining fairness criteria can improve LLM outputs. This approach is particularly beneficial for smaller companies/organizations that lack resources for extensive model fine-tuning. Manerba et al. [6] critique prior bias evaluation methods and propose a new framework for probing social biases in LLMs. They introduce SOFA, a large-scale benchmark designed to assess disparate treatment across a diverse range of identities and stereotypes. Their findings reveal nuanced biases within LLMs, with religious disparities being particularly pronounced. This research emphasizes the need for a comprehensive investigation into biases across multiple dimensions, contributing significantly to the understanding and evaluation of social biases in LLMs. Zhou, H. et al. [7] explore ways to minimize bias in large-scale datasets used for LLM training. They highlight the importance of using diverse and representative data, demonstrating that carefully curated datasets can significantly reduce model output biases. They advocate for continuous dataset evaluation and refinement. Reif et al. [8] examine how reinforcement learning can mitigate bias, showing that these techniques can fine-tune LLMs to reduce biased predictions. They emphasize the potential of adaptive learning methods to create more equitable AI systems.

Dong, X. et al. [9] explore the role of adversarial training in debiasing LLMs, presenting evidence that adversarial training techniques can help in identifying and mitigating biases. Their study highlights the need for continuous evaluation and adaptation of adversarial training methods to ensure their effectiveness. Kotek et al. [10] conduct a cross-linguistic analysis of biases in multilingual LLMs, finding that biases in one language often persist across others. They call for further research into cross-linguistic bias mitigation strategies and stress the importance of inclusive language representation in multilingual models. Prakash et al. [11] discuss the ethical implications of biases in LLMs, stressing the need for ethical frameworks to guide the development and deployment of these models. They propose guidelines for ethical AI development, focusing on transparency, accountability, and the importance of addressing biases to ensure fair and responsible AI use.

5. Future Direction

Research into how to cut down bias in Large Language Models (LLMs) has many exciting paths. One key area is making better & more flexible debiasing techniques. The methods we have now can show promise, but they often do not fully remove bias. Sometimes, they even create new kinds of bias. Future efforts should aim to build algorithms that can change with different data inputs & better reflect our evolving societal norms. Another important direction is to broaden the diversity of datasets used for training LLMs. Including a wide variety of languages, cultures, & perspectives helps to lessen the biases that come from groups that are not represented enough. This makes the models much more inclusive. Collaborations that cross international & disciplinary boundaries can help gather these rich datasets. We also must focus on improving transparency and explainability in LLMs. It's vital for models to explain their decisions clearly, which helps users spot and understand any biases at play. This calls for further research in Explainable AI (XAI) and interpretable machine learning models, integrating them with the core of LLMs. There is a real need for standard benchmarks and metrics to assess bias & fairness in LLMs too. Setting these standards allows researchers to keep track of progress consistently & compare how effective different approaches are. Lastly, it is essential to involve ethicists, sociologists, and voices from impacted communities in the research process. Their insights into how biased AI systems affect society can lead to the development of AI technologies that adhere to ethical guidelines and push for social justice.

6. Conclusion

Addressing bias in Large Language Models (LLMs) is a complex but an essential task. These models are used in many areas now. So, we can't overlook their ability to spread & worsen social biases. Our review points out that we've made significant strides in spotting, understanding, & reducing biases in LLMs. Yet, there still are challenges that need attention. The methods we have right now—like fine-tuning with varied datasets & using fairness measures—look promising. They still aren't perfect. Issues come from biased training data, sensitivity to context, & the complicated nature of language itself. This means we need to keep improving these approaches. Also, working together across different fields is key. For example, lessons from ethics, sociology, and linguistics can give us a better grip and a deeper understanding on bias. Plus, involving diverse communities helps ensure that AI systems cater to a wider array of viewpoints and necessities. In the end, pursuing fairness in LLMs isn't just a tech issue; it's a matter of society's wellbeing. As researchers & practitioners, we have a duty that goes beyond building effective models. We must make certain these technologies advance equity & justice. Ongoing research along with strong ethical checks will be crucial for reaching this vital aim and creating more inclusive AI systems for what lies ahead.

References

- [1] Ranaldi, L., Ruzzetti, E. S., Venditti, D., Onorati, D., & Zanzotto, F. M. (2023). A Trip Towards Fairness: Bias and De-Biasing in Large Language Models. arXiv preprint arXiv:2305.13862.
- [2] Chu, Z., Wang, Z., & Zhang, W. (2024). Fairness in Large Language Models: A Taxonomic Survey. Proceedings of the Special Interest Group on Knowledge Discovery and Data Mining (SIGKDD '24).
- [3] Narayan, M., Pasmore, J., Sampaio, E., Raghavan, V., & Waters, G. (2024). Bias Neutralization Framework: Measuring Fairness in Large Language Models with Bias Intelligence Quotient (BiQ). arXiv.
- [4] Kruspe, A. (2024). Towards detecting unanticipated bias in Large Language Models. arXiv preprint arXiv:2404.02650.
- [5] Chhikara, G., Sharma, A., Ghosh, K., & Chakraborty, A. (2024). Few-Shot Fairness: Unveiling LLM's Potential for Fairness-Aware Classification. arXiv.
- [6] Marchiori Manerba, M., Stańczak, K., Guidotti, R., & Augenstein, I. (2024). Social Bias Probing: Fairness Benchmarking for Language Models. arXiv preprint arXiv:2311.09090.
- [7] Zhou, H., Feng, Z., Zhu, Z., Qian, J., & Mao, K. (2024). UniBias: Unveiling and Mitigating LLM Bias through Internal Attention and FFN Manipulation. arXiv.
- [8] Reif, Y., & Schwartz, R. (2024). Beyond Performance: Quantifying and Mitigating Label Bias in LLMs. arXiv preprint arXiv:2405.02743.
- [9] Dong, X., Wang, Y., Yu, P. S., & Caverlee, J. (2024). Disclosure and Mitigation of Gender Bias in LLMs. arXiv preprint arXiv:2402.11190.
- [10] Kotek, H., Dockum, R., & Sun, D. Q. (2023). Gender bias and stereotypes in Large Language Models. Collective Intelligence Conference, ACM
- [11] Prakash, N., & Lee, R. K.-W. (2023). Layered Bias: Interpreting Bias in Pretrained Large Language Models. Proceedings of the 6th BlackboxNLP Workshop: Analyzing and Interpreting Neural Networks for NLP.