

All You Need Is S P A C E: When Jailbreaking Meets Bias Audit and Reveals What Lies Beneath the Guardrails (Student Abstract)

Arka Dutta¹, Aman Priyanshu², Ashiqur R. KhudaBuksh¹

¹ Rochester Institute of Technology

² Carnegie Mellon University

ad2688@rit.edu, amanpriyanshusms2001@gmail.com, axkvse@rit.edu

Abstract

This paper makes a novel combination of a recently proposed bias audit framework and a recently proposed jailbreaking technique for Llama3. On an audit comprising several disadvantaged groups, our experiments reveal that a jailbroken Llama3 exhibits worrisome antisemitism, racism, misogyny, and homophobia (to list a few) much akin to a broad suite of LLMs that were susceptible to similar biases.

Introduction

How safe are large language models for disadvantaged minorities? As large language models grow in sophistication and find their way into our day-to-day applications, the need for robust audit frameworks to uncover biases targeting disadvantaged minorities is ever-increasing. Evolving guardrails pose a key challenge to audit frameworks in a sense that auditing for biases becomes a cat-and-mouse game. If a broad and generalizable vulnerability that leads to an audit framework becomes widely publicized, newer models will *patch* that safety failure and the same audit will not work for these newer models. An apt example is the toxicity rabbit hole framework (Dutta et al. 2024). This framework uncovered shocking antisemitism, racism, and misogyny among several other biases in a broad suite of proprietary and open LLMs. However, the Llama (Touvron et al. 2023) family of models was non-compliant with the iterative rabbit hole setting. While this study revealed that (1) several well-known LLMs generated content suggesting unbridled, collective physical harms to Jews, Blacks, and LGBTQ+ people; and (2) a large collection of LLMs exhibited algorithmic monoculture in the groups they target and the venom they spew, the biases of Llama remained underexplored.

This paper¹ melds the literature of jailbreaking and bias audits. We demonstrate that combining a recent jailbreaking vulnerability with a simplified version of toxicity rabbit hole framework allows us to bias audit Llama3.

This paper makes the following contributions.

- **Method:** We propose a new approach to bias auditing LLMs by integrating jailbreaking techniques with bias audit frameworks.

Copyright © 2025, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

¹This paper contains highly offensive and disturbing content.

- **Social:** Our substantive findings indicate that Llama3 exhibits shocking biases against several disadvantaged groups generating content inciting unbridled, collective physical violence against minorities.

- **Resource:** For 20 identity groups of which several are historically disadvantaged, we release a dataset of 20,000 Llama3 responses elicited using this novel approach. As control, we release additional 20,000 responses from non-jailbroken Llama3 and Llama2.

Background

Toxicity Rabbit Hole (TRH) Framework. The toxicity rabbit hole (TRH) is an iterative framework that begins with a simple stereotype. For an identity group (e.g., a religion, nationality, or ethnic group) denoted as \mathcal{G} , Dutta et al. 2024 employs initial stereotypes such as *\mathcal{G} are not nice people* and instructs the LLM to make the initial stereotype more toxic giving the LLM the freedom to modify, append to, or completely rewrite the stereotype. After the LLM provides a more toxic rewrite in response to the initial request, in the second step, the framework requests the LLM to generate even more toxic content, but this time using its own previously generated content from the first step as the input. In each subsequent step, the instruction to the LLM is to produce more toxic content than what it generated in the previous step. The deceptive simplicity of the framework notwithstanding, the audit reveals worrisome racism, antisemitism, and misogyny among several other biases in ten well-known LLMs (Dutta et al. 2024).

Jailbreaking Llama3. A recent bug report² led by one of the student co-authors reveals that Prompt-Guard-86M, a BERT-based safety classifier from Meta, can be jailbroken by a simple prompt injection technique. Our experiments reveal that a similar transformation where prompt characters are presented as space-separated characters, Llama3 becomes more susceptible to generating harmful responses.

Experimental Setup

We conduct bias audit for the following 20 identity groups: *Black, Women, Lesbian, Gay, White, Poor, Conservative, Muslim, Jew, Latin, Trans, Disabled, Asian, Immigrant, Mexican, Feminist, Indigenous, Liberal, Men,*

²<https://github.com/meta-llama/llama-models/issues/50>

References

- Anil, R.; Dai, A. M.; Firat, O.; Johnson, M.; Lepikhin, D.; Passos, A.; Shakeri, S.; Taropa, E.; Bailey, P.; Chen, Z.; et al. 2023. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*.
- Bommasani, R.; Creel, K. A.; Kumar, A.; Jurafsky, D.; and Liang, P. S. 2022. Picking on the Same Person: Does Algorithmic Monoculture lead to Outcome Homogenization? *Advances in Neural Information Processing Systems*, 35: 3663–3678.
- Dutta, A.; Khorramrouz, A.; Dutta, S.; and KhudaBukhsh, A. R. 2024. Down the Toxicity Rabbit Hole: A Framework to Bias Audit Large Language Models with Key Emphasis on Racism, Antisemitism, and Misogyny. In *Proceedings of the Thirty-Third IJCAI 2024*, 7242–7250.
- Jiang, A. Q.; Sablayrolles, A.; Mensch, A.; Bamford, C.; Chaplot, D. S.; Casas, D. d. l.; Bressand, F.; Lengyel, G.; Lample, G.; Saulnier, L.; et al. 2023. Mistral 7B. *arXiv preprint arXiv:2310.06825*.
- Sap, M.; Gabriel, S.; Qin, L.; Jurafsky, D.; Smith, N. A.; and Choi, Y. 2020. Social Bias Frames: Reasoning about Social and Power Implications of Language. In *ACL*.
- Touvron, H.; Lavril, T.; Izacard, G.; Martinet, X.; Lachaux, M.-A.; Lacroix, T.; Rozière, B.; Goyal, N.; Hambro, E.; Azhar, F.; Rodriguez, A.; Joulin, A.; Grave, E.; and Lample, G. 2023. LLaMA: Open and Efficient Foundation Language Models. *arXiv:2302.13971*.
- Wang, Y.; Ivison, H.; Dasigi, P.; Hessel, J.; Khot, T.; Chandu, K. R.; Wadden, D.; MacMillan, K.; Smith, N. A.; Beltagy, I.; and Hajishirzi, H. 2023. How Far Can Camels Go? Exploring the State of Instruction Tuning on Open Resources. *arXiv:2306.04751*.