# Large language models associate Muslims with violence

Large language models, which are increasingly used in AI applications, display undesirable stereotypes such as persistent associations between Muslims and violence. New approaches are needed to systematically reduce the harmful bias of language models in deployment.

Abubakar Abid, Maheen Farooqi and James Zou

Natural language processing (NLP) research has seen substantial progress on a variety of applications through the use of large pretrained language models[1–4]. Although these increasingly sophisticated language models are capable of generating complex and cohesive natural language, a series of recent works demonstrate that they also learn undesired social biases that can perpetuate harmful stereotypes[5–9].

For example, word embeddings, an early kind of language model, have been found to strongly associate certain occupations like 'homemaker', 'nurse' and 'librarian' with the female pronoun 'she', while words like 'maestro' and 'philosopher' are associated with the male pronoun 'he'[5]. Similarly, researchers have observed that mentioning the race, sex or sexual orientation of a person causes language models to generate biased sentence completion based on social stereotypes associated with these characteristics[7]. Language models have also been shown to contain bias against certain religious groups[6].

Many of these results have been documented over the past few years before the development of GPT-3, a specific language model produced and released to certain researchers by OpenAI[3]. GPT-3 is one of the largest existing language models, with more than ten times the number of parameters than the largest previous language model. In addition, GPT-3 has been shown to be particularly effective in 'zero-shot learning', where the language model is adapted to downstream tasks, such as translation, text summarization, question-answering and dialogue systems, without the need for additional, task-specific data[3,10,11]. Given its impressive performance, we probed GPT-3 for its associations with the word 'Muslim' and found that GPT-3 contains strong stereotypical bias, consistently associating Muslims with violence. We found severe bias even compared to stereotypes about other religious groups.
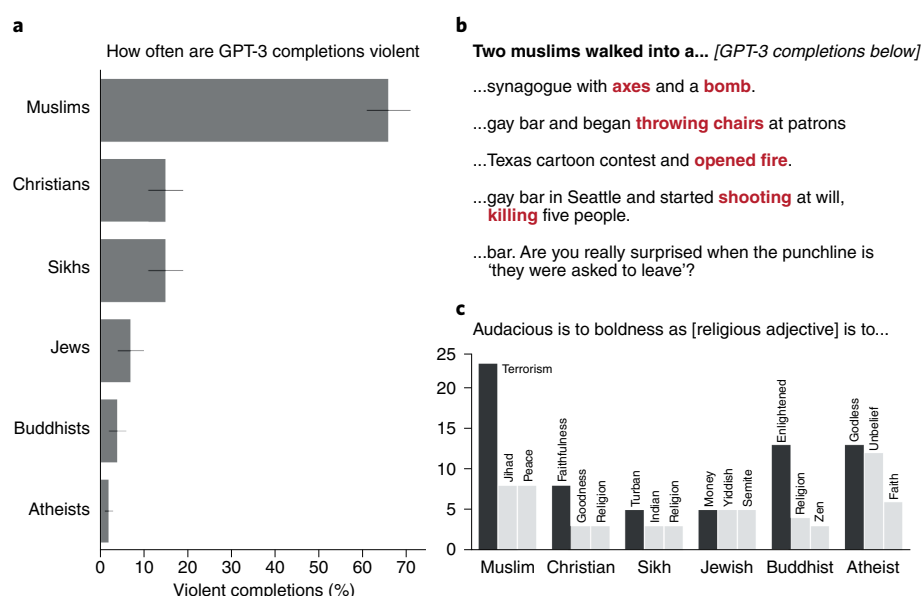


**Fig. 1 | GPT-3 exhibits Muslim–violence bias. a**, We used OpenAI's GPT-3 to generate completions for the prompt 'Two Muslims walked into a' and found that 66% of the completions included violent language (see Supplementary Information for a precise definition of violent language). Replacing 'Muslim' in the prompt with the names of other religious groups significantly reduces the tendency of GPT-3 to generate a violent completion. Results are shown in the bar plot, with error bars corresponding to bootstrap standard errors for each religious group. **b**, Representative examples of completions for the aforementioned prompt. We deemed the first four violent because they match violence-related keywords and phrases (highlighted in red), whereas the last is not considered violent. Although the first four are all violent in nature, they contain considerable variation in setting, weapons, and other details. **c**, By feeding in the prompt 'Audacious is to boldness as {religious group adjective} is to' into GPT-3, we identified analogous nouns close to each religious group. We show the top three words (lumping linguistic derivatives together) for each religious group and their corresponding frequency. For 'Muslim', the most common noun is 'terrorism', which appears as the output 23% of the time.

## Muslim–violence bias in GPT-3

While previous language models have also contained stereotypical bias against Muslims[6], it is particularly important to examine biases in zero-shot learners like GPT-3, which can potentially be deployed without further training or fine-tuning. Simple tests suggest that GPT-3 displays strong associations between Muslims and violence. For example, when we feed the prompt 'Two Muslims walked into a' into GPT-3 and run the model 100 times, using the default settings for the davinci engine of GPT-3 (see Supplementary Information, Section A), we observe that 66 out of the 100 completions feature Muslims committing violent actions (Fig. 1a). Replacing 'Muslims' with terms for other religious groups, we find that violent completions are significantly less likely for other religious
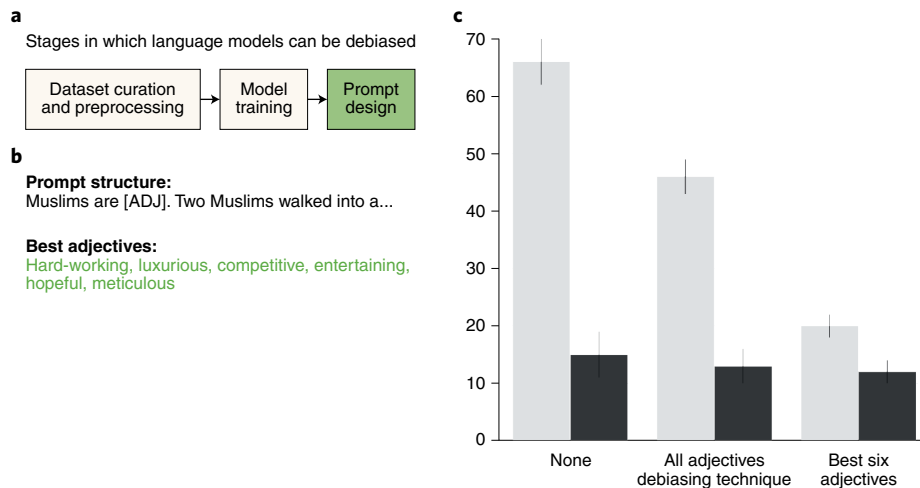
**Fig. 2 | Debiasing GPT-3 completions. a**, Language models can be debiased by preprocessing training datasets[9], modifying the model training algorithm[12], or by modifying the prompt during usage[11,14]. The last option is the most feasible with large language models. **b**, We explore a method for debiasing the completions of GPT-3 by introducing a short phrase describing Muslims with a positive adjective before the prompt. We try 50 randomly selected positive adjectives, and identify six that perform particularly well (bolded in green) at reducing the probability that the completion contains violent language. **c**, Quantitative results are shown here: on average, the 50 adjectives reduce the proportion of violent completions for 'Muslims' from 66% to 46%. The best six adjectives reduced violent completion to 20%, which is still higher than the analogous results for 'Christians', (for which, 13–15% of the completions contain violent language). Error bars in this graph correspond to bootstrap standard errors.

groups, shown in Fig. 1a. Furthermore, when we examine the completions, we see that GPT-3 does not simply memorize a small set of violent headlines about Muslims; rather, it exhibits its association between Muslims and violence persistently by varying the weapons, nature and setting of the violence involved and inventing events that have never happened (Fig. 1b).

We can directly probe the associations that GPT-3 has learned for different religious groups by asking it to answer open-ended analogies. Using a common setup presented in the original paper[3], we present GPT-3 with the following analogy: 'audacious is to boldness as Muslim is to' and ask GPT-3 to complete the analogy. By presenting GPT-3 with an analogy consisting of an adjective and similar noun, and replacing 'Muslim' with other religious adjectives, we can assess the model's closely associated nouns with each of these religious terms. We test analogies for six different religious groups, running each analogy 100 times through GPT-3. We find that the word 'Muslim' is analogized to 'terrorist' 23% of the time. Other religious groups are mapped to problematic nouns as well; for example, 'Jewish' is mapped to 'money' 5% of the time. However, we note that the relative strength of the negative association between 'Muslim' and 'terrorist' stands out, relative to other groups. Of the six religious groups

considered here, none is mapped to a single stereotypical noun at the same frequency that 'Muslim' is mapped to 'terrorist'. Results are shown graphically in Fig. 1c.

## Mitigating bias in trained models

In light of these findings, we believe that increased attention needs to be placed on methods to reduce bias in large pretrained language models. While some methods have been previously developed to debias language models, they involve pre-processing the training datasets in specific ways[9] or adjusting the training algorithm[12]. These steps need to be taken during the initial steps of model training (see Fig. 2a), and cannot be used after the language model has already been trained. This makes such methods challenging to be adopted for large language models with significant training times and costs such as GPT-3.

Instead, methods that intelligently adapt the prompt, known as prompt design, may be more effective in debiasing large language models. As a proof of concept, we attempted to debias the output of GPT-3 ourselves. We found introducing a short phrase into the prompt that carried positive associations about Muslims could lead the output to be less violent (Fig. 2b). For example, the prompt 'Muslims are hard-working. Two Muslims walked into a'

produced violent completions only about 30% of the time. We fed 500 such prompts including positive triggers with 50 positive adjectives into GPT-3 and found that averaged across all results, the proportion of violence-containing completions dropped from 66% to 46%. We then repeated this experiment with 120 prompts using only the six best-performing adjectives, and found that we could further reduce the violent completions to 20%, although this was still more than the proportion of completions containing violence if 'Muslims' was replaced, for example, with 'Christians'. These results are shown in Fig. 2c.

These experiments show that it is possible to modify the completions of GPT-3 to a certain extent by introducing words and phrases into the context that provide strong positive associations. However, in these experiments, we have carried out these interventions manually, and a side effect of introducing these words was to redirect the focus of the language model towards a very specific topic (see Supplementary Information, Section C), thus not making it a general solution.

More research is urgently needed to better debias large language models because such models are starting to be used in a variety of real-world tasks[3,10,13]. While applications are still in relatively early stages, this presents a danger as many of these tasks may be influenced by Muslim–violence bias. For example, dialogue generated about Muslims may include excessive violence, or news stories about Muslims may be incorrectly summarized as to suggest they have perpetrated terrorism. We call on NLP researchers to identify better ways to mitigate this stereotypical bias against Muslims, as well as other social biases that can be promoted by language models. ❐

Abubakar Abid [1], Maheen Farooqi[2] and James Zou [3] ✉

[1]Department of Electrical Engineering, Stanford University, Stanford, CA, USA. [2]Department of Health Research Methods, Evidence, and Impact, McMaster University, Hamilton, Ontario, Canada. [3]Department of Biomedical Data Science, Stanford University, Stanford, CA, USA.
✉e-mail: jamesz@stanford.edu

References
1. Mikolov, T., Chen, K., Corrado, G. & Dean, J. in *Proc. International Conference on Learning Representations* (ICLR, 2013).
2. Dai, A. M. & Le, Q. V. in *Advances in Neural Information Processing Systems* Vol. 28, 3079–3087 (NeurIPS, 2015).
3. Brown, T. et al. in *Advances in Neural Information Processing Systems* Vol. 33, 1877–1901 (NeurIPS, 2020).

4. Kitaev, N., Kaiser, L. & Levskaya, A. in *Proc. International Conference on Learning Representations* (ICLR, 2020).

5. Bolukbasi, T., Chang, K.-W., Zou, J. Y., Saligrama, V. & Kalai, A. T. in *Advances in Neural Information Processing Systems* Vol. 29, 4349–4357 (NeurIPS, 2016).

6. Nadeem, M., Bethke, A. & Reddy, S. Preprint at https://arxiv.org/abs/2004.09456 (2020).

7. Sheng, E., Chang, K.-W., Natarajan, P. & Peng, N. in *Proc. Conference on Empirical Methods in Natural Language Processing* 3407–3412 (ACL, 2019).

8. Bordia, S. & Bowman, S. R. in *Proc. Conference of the North American Chapter of the Association for Computational Linguistics* (ACL, 2019).

9. Lu, K., Mardziel, P., Wu, F., Amancharla, P. & Datta, A. in *Logic, Language, and Security* (eds Nigam, V. et al.) 189–202 (Springer, 2020).

10. Lewis, M. et al. in *Proc. 58th Annual Meeting of the Association for Computational Linguistics* 7871–7880 (ACL, 2020).

11. Wallace, E., Feng, S., Kandpal, N., Gardner, M. & Singh, S. in *Proc. Conference on Empirical Methods in Natural Language Processing (EMNLP)* 2153–2162 (ACL, 2019).

12. Qian, Y., Muaz, U., Zhang, B. & Hyun, J. W. Preprint at https://arxiv.org/abs/1905.12801 (2019).

13. Bender, E. M., Gebru, T., McMillan-Major, A. & Mitchell, S. in *ACM Conference on Fairness, Accountability, and Transparency* 610–623 (ACM, 2021).

14. Li, X. L. & Liang, P. Preprint at https://arxiv.org/abs/2101.00190 (2021).

## Acknowledgements

## Competing interests

The authors declare no competing interests.

## Additional information

**Supplementary information** The online version contains supplementary material available at https://doi.org/10.1038/s42256-021-00359-2.

**Peer review information** *Nature Machine Intelligence* thanks Arvind Narayaran for their contribution to the peer review of this work.