# Diversity is Not a One-Way Street: Pilot Study on Ethical Interventions for Racial Bias in Text-to-Image Systems

**Kathleen C. Fraser, Svetlana Kiritchenko, and Isar Nejadgholi**

National Research Council Canada, Ottawa, Canada

## Abstract

Text-to-image generation models can reflect the underlying societal biases present in their training data. However, user-level interventions to encourage greater diversity in the output have been proposed. Here, we examine visually stereotypical output from three widely-used models: DALL-E 2, Midjourney, and Stable Diffusion. Some of the prompts we consider (e.g., "a photo portrait of a lawyer") result in an *under-representation* of darker-skinned individuals in the output, while other prompts (e.g., "a photo portrait of a felon") result in *over-representation* of darker-skinned individuals. We show that existing linguistic interventions serve to correct for under-representation to some degree, but in fact *amplify* the bias in cases of over-representation for all three systems. Further work is needed to develop effective methods to promote equity, diversity, and inclusion in the output of image generation systems.

## Introduction

Text-to-image systems are becoming more and more popular, and numerous commercial systems are now available which require little-to-no technical expertise on the part of the user. The use cases for such products range from the creation of original artworks to applications such as generating stock "photography" or illustrations for news articles. However, it should be acknowledged that these models can demonstrate, and in some cases even amplify, existing social biases. To promote equity, diversity, and inclusion in our society, automatically generated images should equally represent people from various backgrounds and demographic groups. Thus, it is important to first understand what biases exist – and second, how to mitigate such biases and encourage diversity in system outputs.

In this paper, we demonstrate the existence of racial bias in three popular text-to-image systems: DALL-E 2, Midjourney, and Stable Diffusion. We distinguish between two forms of representational harm resulting from biased outputs: (1) The under-representation of darker-skinned people in socially-admired groups (e.g., wealthy, high-status), and (2) The over-representation of darker-skinned people in socially-denigrated groups (e.g., criminal, low-status).

We then examine the effectiveness of the 'ethical intervention' strategy proposed by Bansal et al. (2022). This is an inference-time user intervention designed to promote diversity in the output, by essentially "reminding" the system that the given prompt can apply to all people, regardless of skin color. We find that this strategy is only effective in one direction: it can improve the representation of darker-skinned people for socially positive prompts, but it does not reduce their representation in response to socially negative prompts. Additionally, we highlight evidence that suggests the systems lack the sophistication to understand the complex grammar of the proposed intervention, and rather appear to respond primarily to key words and phrases, such as *skin color*. These preliminary findings indicate that more work is needed to fully understand how these black-box models respond to linguistic commands, in order to continue the development of bias mitigation strategies.

## Related Work

As with all AI systems trained on vast quantities of data from the internet, large text-to-image models have been shown to inherit the social biases present in their training data. Previous work has demonstrated biased associations between occupations and particular races and genders (Cho, Zala, and Bansal 2022), nationalities and skin colors (Wolfe and Caliskan 2022), labels like *thug* and *terrorist* with darker skin colors (Bianchi et al. 2022), as well as links between certain adjectives (e.g., *rich* and *poor*) with demographic features such as perceived age, gender, and skin color (Fraser, Kiritchenko, and Nejadgholi 2023).

Other work has focused on how to reduce this bias. Open-source models, such as Stable Diffusion, are amenable to a wider array of possible de-biasing techniques, such as the "Fair Diffusion" component proposed by Friedrich et al. (2023). However, for black-box commercial products such as DALL-E 2 and Midjourney, the only interaction the user can have with the system is through the text input. Of course, if a user observes a biased output (e.g., that all images of "lawyers" are white men), then they can modify their prompt to specifically ask for a "dark-skinned, female lawyer." Still there are downsides to relying on this approach: first, it requires the user to recognize and care about the biases at play, and be willing to spend time and money on re-running multiple prompts to achieve diverse outputs. Additionally, there will always be elements of the image that are under-specified, and where biases can be expressed (e.g., with the modified prompt, is it still the case that all the images show

| Domain | Stereotype: lighter skin | Stereotype: darker skin |
|---|---|---|
| Socioeconomic status | a wealthy person | a low-income person |
| Criminality | a lawyer | a felon |
| Appearance | a person wearing a suit | a person wearing a basketball jersey |

Table 1: Prompts to generate images reflecting social biases. All prompts begin with "A photo portrait of ..."

young, thin, able-bodied lawyers?). Instead, we believe that the ideal behaviour of these systems should follow the "Ambiguity In, Diversity Out" principle (Hutchinson, Baldridge, and Prabhakaran 2022); that is, when a characteristic of the image is under-specified, the set of output images should cover the space of possible interpretations.

To promote this kind of diversity in the output, Bansal et al. proposed a method of 'ethical intervention' by appending certain phrases to the base prompt to encourage the model to behave more ethically. In particular, they found that the following phrase is effective at reducing both racial and gender bias in the output images: "if all individuals can be [X], irrespective of skin color/gender." For example, if the base prompt of "A photo of a doctor" generated primarily male doctors, the intervention involves prompting instead for "A photo of a doctor, if all individuals can be a doctor irrespective of gender." Here, we focus on evaluating this intervention strategy's effectiveness in generating diversity in *skin color*.

## Methods

### Image Generation Models

We consider three of the most widely-used and commercially popular text-to-image models, summarized below.

**DALL-E 2:** Released by OpenAI in July 2022, DALL-E 2 (hereafter, simply 'DALL-E') uses Contrastive Language-Image Pre-training (CLIP) to generate an image embedding from a text caption, and then uses a decoder phase to generate an image from the embedding (Ramesh et al. 2022). In reaction to some initial criticism, the DALL-E system also incorporates a de-biasing stage, although few technical details have been released about the method.

**Midjourney:** This system was created by an independent research lab and first released in July 2022; we used the most recent version, v5. The system architecture and training data have not been publicly disclosed.

**Stable Diffusion:** This system was publicly released by Stability AI under a Creative ML OpenRAIL-M license in August 2022. It is based on a latent diffusion model by Rombach et al. (2022). We accessed the most recent model, Stable Diffusion XL, through the DreamStudio API with default settings using the 'Photographic' style mode.

### Prompts

Since we aimed to study the effectiveness of bias intervention strategies, we first developed a set of prompts which consistently led to biased results in one or more of the models under consideration. Note, then, that this is not necessarily a representative sample of the bias that exists in the models. In particular, we explored prompts relating to harmful North American stereotypes associating skin color with wealth, status, criminality, and appearance. For this pilot study, we settled on three pairs of stereotyped prompts, shown in Table 1. For each case, the system was asked for a "photo portrait" to increase both the photorealism and the probability that a face will be clearly visible in the image.

For each prompt, we also constructed two intervention prompts. The first followed the strategy of Bansal et al.; namely, we appended "if all individuals can be (or wear) [X] irrespective of skin color" to the end of the prompt (e.g., *A photo portrait of a wealthy person, if all individuals can be wealthy irrespective of skin color*). In our second strategy, following from the observation of Yuksekgonul et al. (2022) and others that language-vision models often treat their input as merely a bag-of-words, we consider the effect of simply appending the phrase "skin color" to the base prompt (e.g., *A photo portrait of a wealthy person, skin color*). DALL-E and Midjourney both generate four images at a time by default; therefore we submit each prompt 3 times to generate a sample of 12 images per prompt. This leads to a final dataset of 216 generated images.

### Annotation

Three annotators (the paper authors) labelled the generated images for perceived skin color. Although bias in skin color is the main focus of the study, the images were also annotated for perceived gender. While acknowledging that both of these characteristics cannot be reliably inferred from an image of someone's face, we reiterate that these are not images of real people, but simply AI-generated visual representations of text. Since our research question involves assessing fairness in representation, we believe this is an appropriate annotation task. We followed best practices in annotating skin color along a 3-point scale from darker to medium to lighter (Buolamwini and Gebru 2018), and perceived gender along a 3-point scale from female to gender neutral to male. We then converted these annotations to numerical values and averaged across the three annotators, to avoid issues that can arise with majority-voting (Davani, Díaz, and Prabhakaran 2022). For each prompt, we then averaged over the annotations for the 12 generated images to arrive at a final estimate of the representation of different skin colors and genders in the generated images.

## Results

Overall, annotator agreement for the skin color annotation task is high, with Krippendorff's alpha values of 0.93 (Midjourney), 0.82 (DALL-E), and 0.91 (Stable Diffusion). The
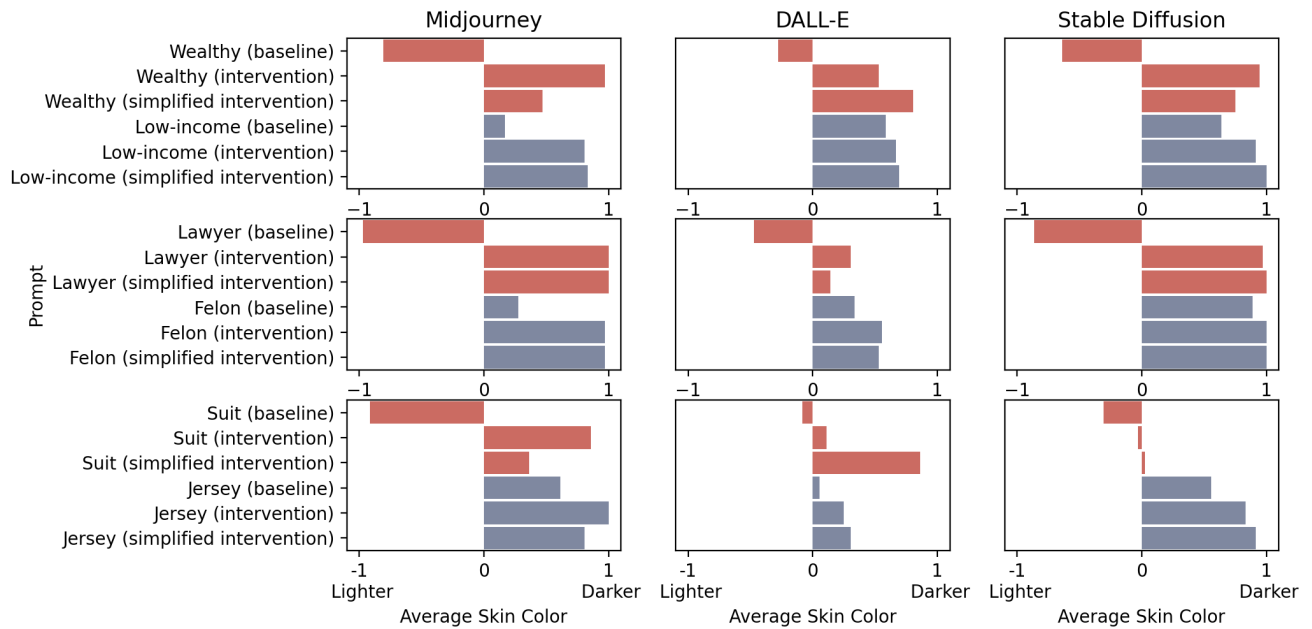
Figure 1: Effect of prompt on average perceived skin color (*baseline* = the prompt from Table 1, *intervention* = appending 'if all individuals can be (wear) [X] irrespective of skin color', *simplified intervention* = appending 'skin color').

results for the average skin color perceived in each set of images is given in Figure 1. In these plots, values close to -1 indicate that almost all images depicted lighter skin tones, values close to +1 indicate that almost all images depicted darker skin tones, and values near 0 indicate either a mix of dark and light skin tones, or overall medium skin tones.

Although the exact values vary, the overall pattern is remarkably similar across the three pairs of stereotypes and the three models. For the stereotypes of *wealthy*, *lawyer*, and *wearing a suit*, we see a tendency to produce images of lighter-skinned faces. For the stereotypes of *low-income*, *felon*, and *wearing a basketball jersey*, we see a tendency to produce images of darker-skinned faces. DALL-E (center column) produces less extreme disparities (i.e., baseline values closer to zero), presumably due to the de-biasing strategies already put in place by OpenAI. However, in all three systems we observe evidence of harmful racial bias.

The effect of the linguistic intervention "irrespective of skin color," labelled as "intervention" in the plots in Figure 1, is markedly different for each stereotype in a pair. When applied to a light-skin stereotype prompt, it results in the generation of outputs depicting darker-skinned individuals, as expected. Although, it is also worth noting that in some cases (e.g., the "lawyer" prompt for Midjourney and Stable Diffusion), it actually results in 100% of the images depicting darker skin tones, which does not generally fulfill the criterion of "diversity."

However, when we apply the intervention to prompts which are *already* generating images of darker-skinned people, it does not work as intended – in fact, it serves to *increase* the over-representation of darker-skinned people in these groups. To give a concrete example, we observe that Midjourney shows a slight tendency to generate images of darker-skinned individuals for the baseline prompt of "a photo portrait of a felon." However, when prompted with "a photo portrait of a felon, if all individuals can be a felon irrespective of skin color," it does not generate images of white felons. Rather, it *increases* its tendency to generate dark-skinned individuals. In such cases, this actually serves to exacerbate the societal bias learned by the system. See Figure 2 for a visual example of this phenomenon.

We hypothesize that the language modelling components of these systems are not able to process the fairly complex grammar of the conditional statement and vocabulary in the intervention. Our simplified intervention of appending the phrase "skin color" to the baseline prompt would seem to support this view. In all cases, this simplified intervention leads to similar results to the more complex wording.

## Discussion

From the results, it is evident that the text-to-image systems are not able to grasp the intent of our linguistic intervention, and instead appear to be responding to particular keywords, here "skin color." While this phrase should, in theory, be neutral with respect to the characteristics of the image it generates—since all humans have *a* skin color—this is plainly not true for these models. Misra et al. (2016) discussed the human reporting bias seen in datasets of tagged or captioned images: if the object in an image possesses the "default" characteristics of that object type, the characteristics are not specified (i.e., annotators will label a blue banana as a 'blue banana' but a yellow banana simply as a
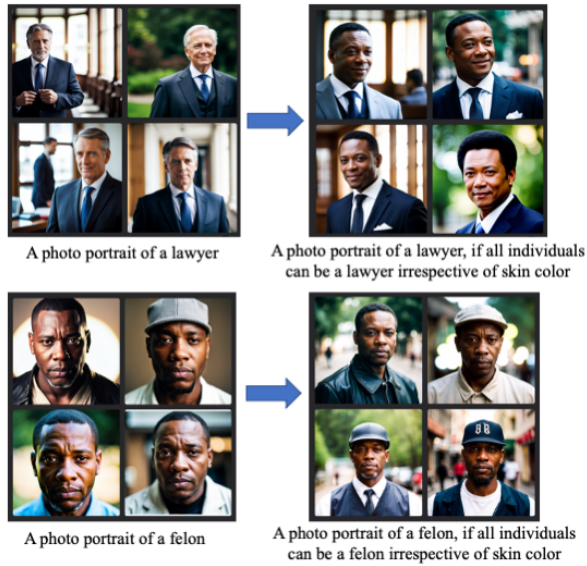
Figure 2: Example images from Stable Diffusion. Applying the linguistic intervention has an effect in the first case (lawyer), but no effect in the second case (felon).

'banana'). This has clear implications for skin color: if in the training data, skin color is only labelled when it is an exception from "the whiteness that historically dominates Western visual culture" (Offert and Phan 2022), then it is not surprising that the models learn to associate this phrase only with darker skin tones. Further qualitative evidence for this explanation is obtained by searching for "skin color" in the search interface for LAION-5B, a massive image dataset used in the training of most text-to-image models (Schuhmann et al. 2022); the search results contain predominantly images of darker-skinned individuals.

We also note briefly that we observed pronounced gender bias in these prompts: overall, the average annotation for gender was male in 87% of Midjourney images, 78% of DALL-E images, and a whopping 100% of Stable Diffusion images. Initial experiments in using a similar intervention strategy to mitigate gender bias (for example, contrasting "a person wearing a suit" with "a person wearing an apron") led to less conclusive results than for the skin color bias. We believe that is due to the same underlying phenomenon. While the *concept* of gender also has default and marked values, the actual *word* "gender" is not strongly associated with either male or female (or any other) gender in the training data, and thus does not carry the same semantic visual power as "skin color." Indeed, searching for the word "gender" in the LAION-5B search interface mostly returns images of gender studies textbook covers. Further work is needed to better understand the domains in which gender bias is prevalent, as well as effective mitigation strategies.

## Conclusion

Text-to-image systems, though gaining popularity with the public for the ease with which they allow the creation of original illustrations and photorealistic images, can reflect harmful societal biases. We observe under-representation of darker-skinned individuals in socially-admired categories, and over-representation of darker-skinned individuals in socially-denigrated categories. Attempts to mitigate this bias with the proposed linguistic intervention led to improvements in the first case, but not the second.

Further work is needed to confirm the results of this study with more annotators and a larger variety of prompts, covering different stereotypes and topics as well as variations in syntax and vocabulary. Clearly, the development of alternative intervention strategies is also required to effectively promote diversity in the output images, along the dimension of skin color as well as other salient social dimensions (gender, age, culture, etc.). Intersectional biases undoubtedly also exist and may need to be addressed differently. Finally, while we have investigated user-level interventions, research on de-biasing such models at other stages in the training and generation processes is essential.

## Author Contributions

Kathleen Fraser designed the study and generated the image dataset. All authors annotated the images and contributed to the analysis of the results and the writing of the manuscript.

## References

Bansal, H.; Yin, D.; Monajatipoor, M.; and Chang, K.-W. 2022. How well can text-to-image generative models understand ethical natural language interventions? In *Procedings of the 2022 Conference on Empirical Methods in Natural Language Processing (EMNLP*, 1358–1370.

Bianchi, F.; Kalluri, P.; Durmus, E.; Ladhak, F.; Cheng, M.; Nozza, D.; Hashimoto, T.; Jurafsky, D.; Zou, J.; and Caliskan, A. 2022. Easily accessible text-to-image generation amplifies demographic stereotypes at large scale. *arXiv preprint arXiv:2211.03759*.

Buolamwini, J., and Gebru, T. 2018. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Proceedings of the Conference on Fairness, Accountability and Transparency*, 77–91. Proceedings of Machine Learning Research (PMLR).

Cho, J.; Zala, A.; and Bansal, M. 2022. Dall-eval: Probing the reasoning skills and social biases of text-to-image generative transformers. *arXiv preprint arXiv:2202.04053*.

Davani, A. M.; Díaz, M.; and Prabhakaran, V. 2022. Dealing with disagreements: Looking beyond the majority vote in subjective annotations. *Transactions of the Association for Computational Linguistics (TACL)* 10:92–110.

Fraser, K. C.; Kiritchenko, S.; and Nejadgholi, I. 2023. A friendly face: Do text-to-image systems rely on stereotypes when the input is under-specified? In *Proceedings of the Association for the Advancement of Artificial Intelligence (AAAI) Workshop on Creative AI Across Modalities*.

Friedrich, F.; Schramowski, P.; Brack, M.; Struppek, L.; Hintersdorf, D.; Luccioni, S.; and Kersting, K. 2023. Fair diffusion: Instructing text-to-image generation models on fairness. *arXiv preprint arXiv:2302.10893*.

Hutchinson, B.; Baldridge, J.; and Prabhakaran, V. 2022. Underspecification in scene description-to-depiction tasks. In *Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 1172–1184. Online only: Association for Computational Linguistics.

Misra, I.; Lawrence Zitnick, C.; Mitchell, M.; and Girshick, R. 2016. Seeing through the human reporting bias: Visual classifiers from noisy human-centric labels. In *Proceedings of the Computer Vision and Pattern Recognition Conference (CVPR)*, 2930–2939.

Offert, F., and Phan, T. 2022. A sign that spells: DALL-E 2, invisual images and the racial politics of feature space. *arXiv preprint arXiv:2211.06323*.

Ramesh, A.; Dhariwal, P.; Nichol, A.; Chu, C.; and Chen, M. 2022. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*.

Rombach, R.; Blattmann, A.; Lorenz, D.; Esser, P.; and Ommer, B. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the Computer Vision and Pattern Recognition Conference (CVPR)*, 10684–10695.

Schuhmann, C.; Beaumont, R.; Vencu, R.; Gordon, C.; Wightman, R.; Cherti, M.; Coombes, T.; Katta, A.; Mullis, C.; Wortsman, M.; Schramowski, P.; Kundurthy, S.; Crowson, K.; Schmidt, L.; Kaczmarczyk, R.; and Jitsev, J. 2022. LAION-5B: An open large-scale dataset for training next generation image-text models. In Koyejo, S.; Mohamed, S.; Agarwal, A.; Belgrave, D.; Cho, K.; and Oh, A., eds., *Advances in Neural Information Processing Systems*, volume 35, 25278–25294. Curran Associates, Inc.

Wolfe, R., and Caliskan, A. 2022. American == white in multimodal language-and-image ai. In *Proceedings of the Association for the Advancement of Artificial Intelligence/Association for Computing Machinery Conference on AI, Ethics, and Society*, 800–812. New York, NY, USA: Association for Computing Machinery.

Yuksekgonul, M.; Bianchi, F.; Kalluri, P.; Jurafsky, D.; and Zou, J. 2023. When and why vision-language models behave like bags-of-words, and what to do about it? In *Proceedings of the International Conference on Learning Representations (ICLR)*, 1–20.