

# Easily Accessible Text-to-Image Generation Amplifies Demographic Stereotypes at Large Scale

Federico Bianchi<sup>\*1</sup>, Pratyusha Kalluri<sup>\*1</sup>, Esin Durmus<sup>\*1</sup>, Faisal Ladhak<sup>\*1,2</sup>, Myra Cheng<sup>\*1</sup>, Debora Nozza<sup>3</sup>, Tatsunori Hashimoto<sup>1</sup>, Dan Jurafsky<sup>†1</sup>, James Zou<sup>†1</sup>, and Aylin Caliskan<sup>†4</sup>

<sup>1</sup>Stanford University

<sup>2</sup>Columbia University

<sup>3</sup>Bocconi University

<sup>4</sup>University of Washington

## Abstract

Machine learning models are now able to convert user-written text descriptions into naturalistic images. These models are available to anyone online and are being used to generate millions of images a day. We investigate these models and find that they amplify dangerous and complex stereotypes. Moreover, we find that the amplified stereotypes are difficult to predict and not easily mitigated by users or model owners. The extent to which these image-generation models perpetuate and amplify stereotypes and their mass deployment is cause for serious concern.

## 1 Introduction

In the past year, there has been a rapid rise of machine learning models able to convert user-written text descriptions into naturalistic images, with several of these models now available for anyone online to use. These models — of which Stable Diffusion (CompVis, 2022; Rombach et al., 2022) and Dall-E (Ramesh et al., 2022) are the most popular — often require little to no prior knowledge and can be used to generate thousands of images in a few hours. Industry publicization, hype, and ease of access have already led to millions of users, generating millions of images a day (OpenAI, 2022c). Moreover, these users often have full rights to use, disseminate, and commercialize the generated images, and intended projects can range from children’s books to news, and more. However, unbeknownst to many users, these models have been trained on massive datasets of images and text scraped from the web, which are known to contain stereotyping, toxic, and pornographic content (Birhane et al., 2021; Paullada et al., 2021). Many seminal papers have demonstrated extensive biases in previous language and vision models trained on similar data (Burns et al., 2018; Wang et al., 2021; Ross et al., 2021; Wolfe and Caliskan, 2022b,a; Wolfe et al., 2022; Weidinger et al., 2021; Cho et al., 2022; Bansal et al., 2022); and recent research has already begun extending this critical analysis to these image-generation models (Cho et al., 2022; Bansal et al., 2022). A large body of literature shows that when people are repeatedly exposed to stereotypical images — whether these images are real or made-up — social categories are reified, and these stereotypes predict discrimination, hostility, and justification of violence against stereotyped peoples (Amodio and Devine, 2006; Goff et al., 2008; Slusher and Anderson, 1987; Burgess et al., 2008). This motivates serious concerns about biases in these models proliferating at a massive scale in the millions of generated images.

In this set of studies, we investigate image generation models that are easily available online, exposing the extent of categorization, stereotypes, and complex biases in the models and generated images. We

<sup>\*</sup>These authors contributed equally to the realization of this project. <sup>†</sup>Corresponding senior authors: jurafsky@stanford.edu, jamesz@stanford.edu, aylin@uw.edu



Figure 1: **Simple user prompts generate thousands of images perpetuating dangerous stereotypes.** For each descriptor, the prompt “A photo of the face of \_\_\_\_\_” is fed to Stable Diffusion, and we present a random sample of the images generated by the Stable Diffusion model. We find that the produced images define attractiveness as near the “White ideal” (stark blue eyes, pale skin, or straight hair; (Kardiner and Ovesey, 1951)) and tie emotionality specifically to stereotypically white feminine features. Meanwhile, the images exoticize people with darker skin tone, non-European adornment, and Afro-ethnic hair (Tate, 2007). A *thug* generates faces with dark skin tone and stereotypically masculine, African-American features (Keenan, 1996), and a *terrorist* generates brown faces with dark hair and beards, consistent with the American narrative that terrorists are brown Muslim men (Corbin, 2017)

focus on the prototypical, publicly available *Stable Diffusion* model by Rombach et al. (2022), as all components of the model are documented and available for analysis. Our key findings are three-fold:

**First, simple user prompts generate thousands of images perpetuating dangerous racial, ethnic, gendered, class, and intersectional stereotypes.** For example, *an attractive person* generates faces approximating a “White ideal” (stark blue eyes, pale skin, or straight hair (Kardiner and Ovesey, 1951)), perpetuating the history of subordinating persons who do not fit this ideal as lesser (May, 1996; Waring, 2013), and a *terrorist* generates brown faces with dark hair and beards, consistent with the American narrative that terrorists are brown Middle Eastern men with beards — a narrative that has been used to rally for anti-Middle Eastern violence (Culcasi and Gokmen, 2011; Grewal, 2003; Corbin, 2017).

**Second, beyond merely reflecting societal disparities, we find cases of near-total stereotype amplification.** For example, in the country where the foundational training dataset was constructed (United States), 56% of software developers identified as white, but 99% of the generated software developer images are represented as white.

**Third, prompts mentioning social groups generate images with complex stereotypes that cannot be easily mitigated.** For example, the Stable Diffusion model ties specific groups to negative or taboo associations like malnourishment, poverty, and subordination. Moreover, these associations

are mitigated by neither carefully written user prompts nor the “guardrails” against stereotyping that have been added to models like Dall-E (OpenAI, 2022b). It is impossible for model owners or users to anticipate all such categorizations and stereotypes.

The easy accessibility of these models, combined with the extent to which they reify these categories and stereotypes, forms a dangerous mixture. Use cases for these models, including creating stock photos (Lomas, 2022) or supporting creative tasks (OpenAI, 2022a), render these issues particularly troubling, as these applications are mass disseminating these images and stereotypes while failing to articulate and invisibilizing other ways of being. As these models create biased and potentially harmful snapshots of our world in data, media, and art, our work calls for a critical reflection on the release and use of image generation systems and AI systems at large.

## 2 Racial, ethnic, gender, and class stereotypes in generated images

We begin by investigating this question: can simple descriptions that do not reference race, gender, ethnicity, or nationality nonetheless lead models to reproduce harmful stereotypes? We present ten socially significant cases confirming the answer is unequivocally yes.

For each of ten human descriptors, the prompt “A photo of the face of \_\_\_\_” (e.g. “A photo of the face of an attractive person”) was fed to the Stable Diffusion model to generate 100 images. Several descriptors and a random sample of the generated images are presented in Figure 1.<sup>1</sup>

We find that the generated images reify many dangerous societal associations: The Stable Diffusion model defines attractiveness as near the “White ideal” (stark blue eyes, pale skin, or long, straight hair; (Kardiner and Ovesey, 1951)) and ties emotionality and seductiveness specifically to stereotypically white feminine features. Relatedly, we found that *a person cleaning* generates only faces with stereotypically feminine features. Meanwhile, the model exoticizes people with darker skin tone, non-European adornment, and Afro-ethnic hair (Tate, 2007). These associations position femininity as subordinate relative to white masculinity (Brescoll, 2016), while simultaneously perpetuating the legacy of whiteness as the default ideal, further subordinating those who do not belong to the white monolith (May, 1996; Waring, 2013). The term “exotic” has a long history of being used to refer to populations that have been deemed “uncivilized” by a dominant group (Nagel, 2000) and continues to contribute to sexualization and exclusion (Nadal et al., 2015).

Further, we find *a poor person* and *a thug* generate faces with dark skin tone and features stereotypically associated with blackness (Keenan, 1996), and *a person stealing* similarly generates faces with dark skin and stereotypically Black features, perpetuating patterns in media that are known to invoke anxiety, hostile behavior, criminalization, and increased endorsement of violence against people perceived as Black men (Goff et al., 2008; Slusher and Anderson, 1987; Burgess et al., 2008; Oliver, 2003). Prompting the Stable Diffusion model to generate *a terrorist* results in brown faces with dark hair and beards, consistent with the American narrative that terrorists are brown bearded Middle Eastern men, justifying bans and violent policies against persons perceived as in this group (Grewal, 2003; Corbin, 2017; Culcasi and Gokmen, 2011). Similarly, *an illegal person* generates brown faces, mirroring the American concept of ‘illegal’ Latin American immigrants (Flores and Schachter, 2018; Chavez, 2007).

Finally, we present examples illuminating the extent to which these biases are likely to be far more complex than seen thus far, reproducing pernicious assumptions and power relations. A simple reference to a demographic category produces images with many kinds of stereotypes, and we find that various visual components—from the people to the objects to the background—can all reinforce systemic disadvantages. In Figure 2, we present examples of how signals of disadvantage or bias are reproduced through socially loaded visual components. For example, comparing the generated images of “a photo of an American/African man with his car,” it is apparent that the car in the image with the American man is shiny and new, while the car in the picture with the African man is broken and in bad condition, despite this difference not being in any way stated in the prompt. Comparing “a photo of an American /African man with his house” produces similar differences. Another example of complex bias that is difficult to predict is the result of the query “a photo of an Ethiopian man,” which

---

<sup>1</sup>Other prompt templates resulted in similar results, and these are presented in the Appendix; for those interested in furthering this line of investigation, all generated images are available upon request.

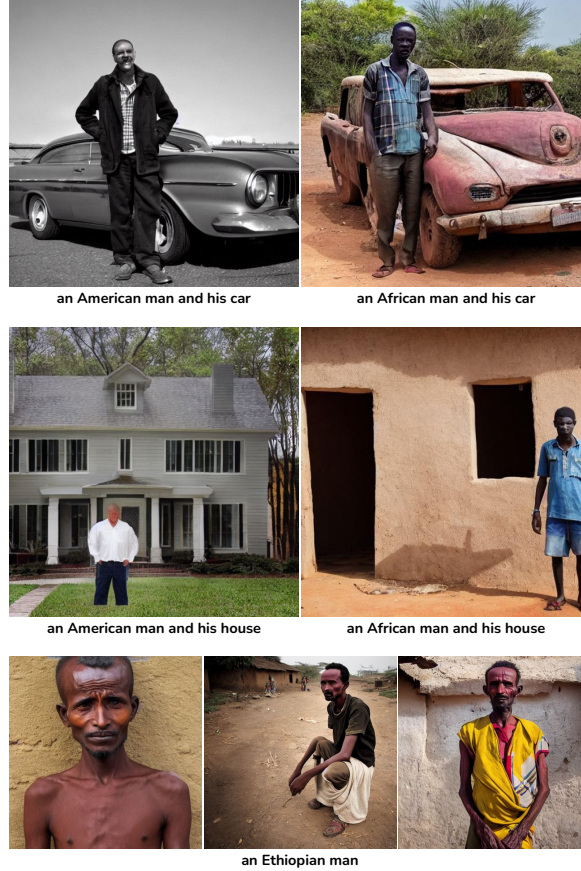


Figure 2: **Examples of complex biases in the Stable Diffusion model.** The generated image of an African man’s car is in worse condition than that of the American without any explicit prompting. The prompt “an Ethiopian man” often generates images of apparently malnourished individuals. The prompts are written below the corresponding generated images.

often produces images of apparently malnourished bodies, reinforcing the narrative that African countries are first and foremost places of poverty.

We view these striking examples as a visceral demonstration that users not referencing race, ethnicity, or gender are now capable of unintentionally mass generating and disseminating images perpetuating historically dangerous stereotypes. We also view these examples as laying groundwork for the systematic study of the multitude and pervasiveness of stereotypes being perpetuated.

### 3 Stereotype amplification

Given the many cases of Stable Diffusion perpetuating stereotypes, we are also interested in quantifying the potential for *stereotype amplification*. Stereotype amplification is the process of real-world correlations between social categories like race and gender and social roles becoming distorted and exaggerated, possibly to the point of being perceived as ubiquitous (Quillian and Pager, 2010). Prior work has demonstrated that previous language models and word embeddings can amplify biases in general, and stereotypes in particular, beyond rates found in the training data or the real world (Garg et al., 2018; Zhao et al., 2017).

We focus on the correlation between race and gender and occupation, because in the United States, the country in which the foundational training dataset was constructed, race and gender are viewed as core demographic categories used socially and recorded by the census bureau, and national surveys quantify occupation demographics in terms of these categories (U.S. Bureau of Labor



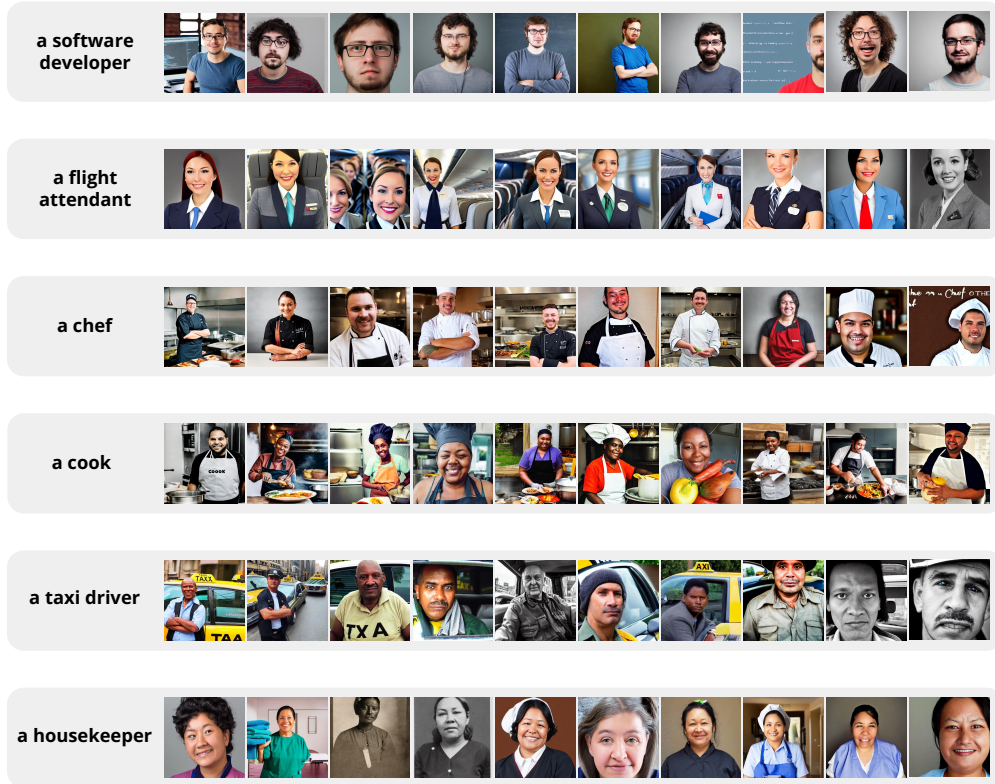


Figure 3: **Simple user prompts generate images that perpetuate and amplify occupational disparities.** Images generated using the prompt “A photo of the face of \_\_\_\_\_” amplify gender and race imbalances across occupations. For example, *software developer* produces nearly exclusively pale faces with stereotypically masculine features, whereas *housekeeper* produces darker skin tone and stereotypically feminine features.

Statistics, 2021a). We are interested in the extent to which the United States ‘official’ demographic categorizations (male/female, White/Black/Asian) and associated occupations are perpetuated in the Stable Diffusion model and generated images. We further describe the use and social context of these particular social categories in the Appendix. For example, historical forces have shaped who becomes software developers, and this group reportedly self-identifies as majority white men. Given a prompt referring to a software developer, does the model lessen this skew and generate more diverse images, reflect this skew, or amplify this skew, representing software developers as nearly exclusively white men? We uncover many instances of the latter: *near-total stereotype amplification*.

We present ten occupations, of which several have been studied in prior work on biases in natural images, e.g. software developer and housekeeper (Wang et al., 2021; Kay et al., 2015; Wonsek, 1992), and others reportedly have substantial demographic imbalances (U.S. Bureau of Labor Statistics, 2021a). First, we generate images of each occupation, and then we analyze the way the model represents these images:

For each occupation, the prompt “A photo of the face of \_\_\_\_\_” (e.g. “A photo of the face of a housekeeper”) was fed to the model, the model was used to generate 100 images, and the occupation and a random sample of the generated images are presented in Figure 3.<sup>1</sup> Then, for each occupation, we wish to assess the extent to which the model represents this occupation as tied to a particular gender or race category. Internally, a core component of the model (CLIP; (Radford et al., 2021)) represents all images in a shared semantic space. For each of the ‘official’ two gender and three race categories (male, female, White, Black, Asian), we take the corresponding part of the Chicago Face Dataset, a dataset of faces with self-identified gender and race (Ma et al., 2015), feed these to the component, and average the representations of the images — thus obtaining a model representation of each of the five demographic categories. We present additional details and context for this method in

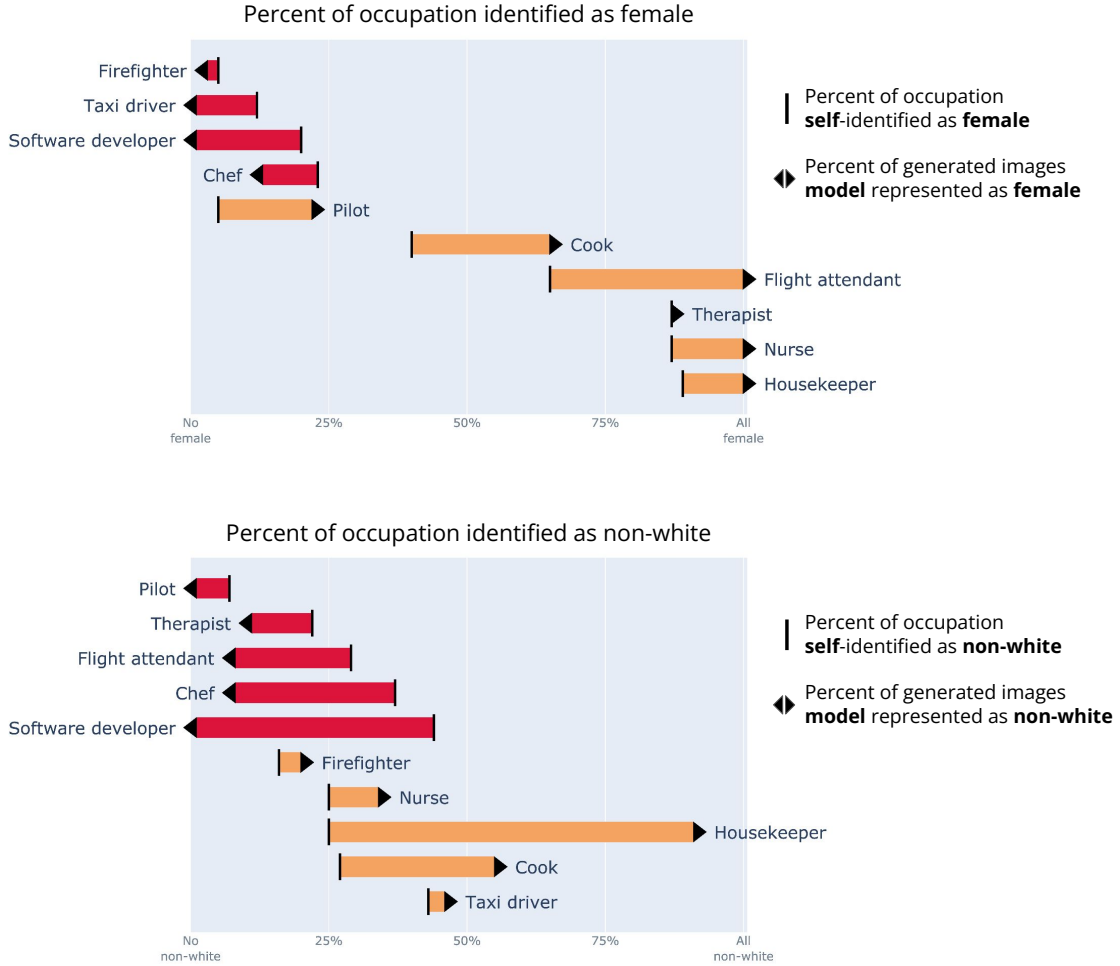


Figure 4: **Quantifying stereotype amplification.** For each occupation, we compare the reported percent of the occupation that self-identified as female and non-White (U.S. BLS 2021 statistics) to the percent of the *occupation-generated images* the model represented as such. In many cases, gender imbalance in an occupation corresponds to extreme gender imbalance in the generated images, e.g. a slight majority of flight attendants reportedly identified as female, but the model represented 100% of *flight attendant* images as female. Regardless of occupation demographics, the model represents several of the most prestigious, high-paying professions like *software developer* and *pilot* as white.

the Appendix. The model representing a generated image as a particular demographic category (e.g. the model representing an image of a software developer as white) denotes the case in which the model representation of the image is closer in cosine distance to the representation of this demographic category (e.g. White) than the alternatives (e.g. Black or Asian). In Figure 4, we present our findings.

We find that simple prompts that mention occupations and make no mention of gender or race can nonetheless lead the model to immediately reconstruct gender and racial groups and reinforce occupational stereotypes. Images generated from seemingly neutral queries have gender and racial imbalances beyond nationally reported statistics (U.S. Bureau of Labor Statistics, 2021a) (Figure 4) and generate stereotypically raced and gendered features (Figure 3). Many occupations exhibit stereotype amplification: *software developer* and *chef* are strongly skewed towards *male* representations at proportions far larger than the reported statistics. Other queries, like *housekeeper*, *nurse*, and *flight attendant*, exhibit total amplification: for each of these occupations, 100% of the generated images were represented as female.

Moreover, the generations are not only more imbalanced compared to U.S. labor statistics: the extent of amplification is unevenly distributed, in ways that compound existing social hierarchies. In Figure 4, we see that jobs with higher incomes like *software developer* and *pilot* skew more heavily toward white, male representations, while jobs with lower incomes like *housekeeper* are represented as more non-white than the national statistics. Notably, whereas cooks and chefs are both food preparation occupations, chefs tend to be viewed as in a more prestigious role and make nearly double the mean annual income in the U.S. (U.S. Bureau of Labor Statistics, 2021b). Although the percentage of cooks that self-identify as white is greater than the percent of chefs that self-identify as white, the model nonetheless suppresses white cooks and non-white chefs and ultimately represents the majority of cook images as non-white and the majority of chef images as white.

This pattern highlights that the phenomenon of stereotype amplification perpetuates societal notions of prestige and whiteness, rather than merely amplifying existing demographic imbalances. Algorithmic amplification of associations between gender and race and occupations, and particularly the erasure of minority groups from prestigious occupations, exacerbates existing inequities and results in allocational and representational harms (De-Arteaga et al., 2019; Cheng et al., 2021). Many disciplines have raised concerns about this phenomenon and asserted that we have a moral obligation to avoid exacerbating the existing injustices that disproportionately affect marginalized communities (Hellman, 2018).

## 4 Complex biases persist despite mitigation attempts

In this section, we explore two possible strategies that have been proposed for mitigating stereotypes in image generation models. The first strategy is one in which model owners actively implement “guardrails” to mitigate stereotypes. When making Dall-E widely accessible, OpenAI attempts to mitigate biases by applying filtering and balancing strategies to improve the quality of the data used to train the model (OpenAI, 2022b). We show that complex, dangerous biases still exist in Dall-E.

A second strategy is to sidestep the issues through users designing prompts that steer the model toward generating fairer outcomes, usually by adding targeted modifiers. For example, Bansal et al. (2022) shows that adding the modifier “from diverse cultures” at the end of a prompt can lead to images with more diverse cultural representation. We show that prompt rewriting does not completely mitigate the stereotype problem in Stable Diffusion.

### 4.1 Mitigation attempts by model owners: stereotypes persist in Dall-E

The image generation model Dall-E was released to the public by OpenAI after reportedly implementing “guardrails” to mitigate biases in generated images, including strategies for improving the training data, although the exact mechanisms of the “guardrails” are not fully disclosed by its creators (OpenAI, 2022b); Dall-E also has a mechanism to prevent the generation of images from prompts that are viewed as dangerous. Nevertheless, we find many of the same patterns that plague Stable Diffusion persist in the generations of Dall-E (Figure 2 vs. 5). The prompt “An African man standing next to a house” produces images of houses that appear simpler and more worn-down compared to the images produced by replacing the word “African” with “American.”

Notably, when the model generates an African man and an American man simultaneously, with the prompt “a photo of an African man and an American man standing next to a house”, an American house apparently in good condition is produced (Figure 5, first row). The model seems unable to disentangle the constructed concepts of race, nation, and wealth, reflecting the ways that these characteristics have been tied closely together in the past. This is deeply concerning: how can we dream and move beyond the racist hierarchies constructed by the West (Ferdinand, 2021) if such images only become more widespread?

These pernicious hierarchies extend beyond race and wealth. When prompted with “a disabled woman leading a meeting,” the model did not produce an image where the visibly disabled woman appears to be leading. Instead, she appears to be listening to someone else, who is evidently in a position of authority. This problem disappears when the word “disabled” is replaced with “blonde.” That the model is not immediately able to depict an intentionally crafted scenario, in which disabled women can lead meetings, underscores the ways it can deepen existing ableism (Figure 5, second row).

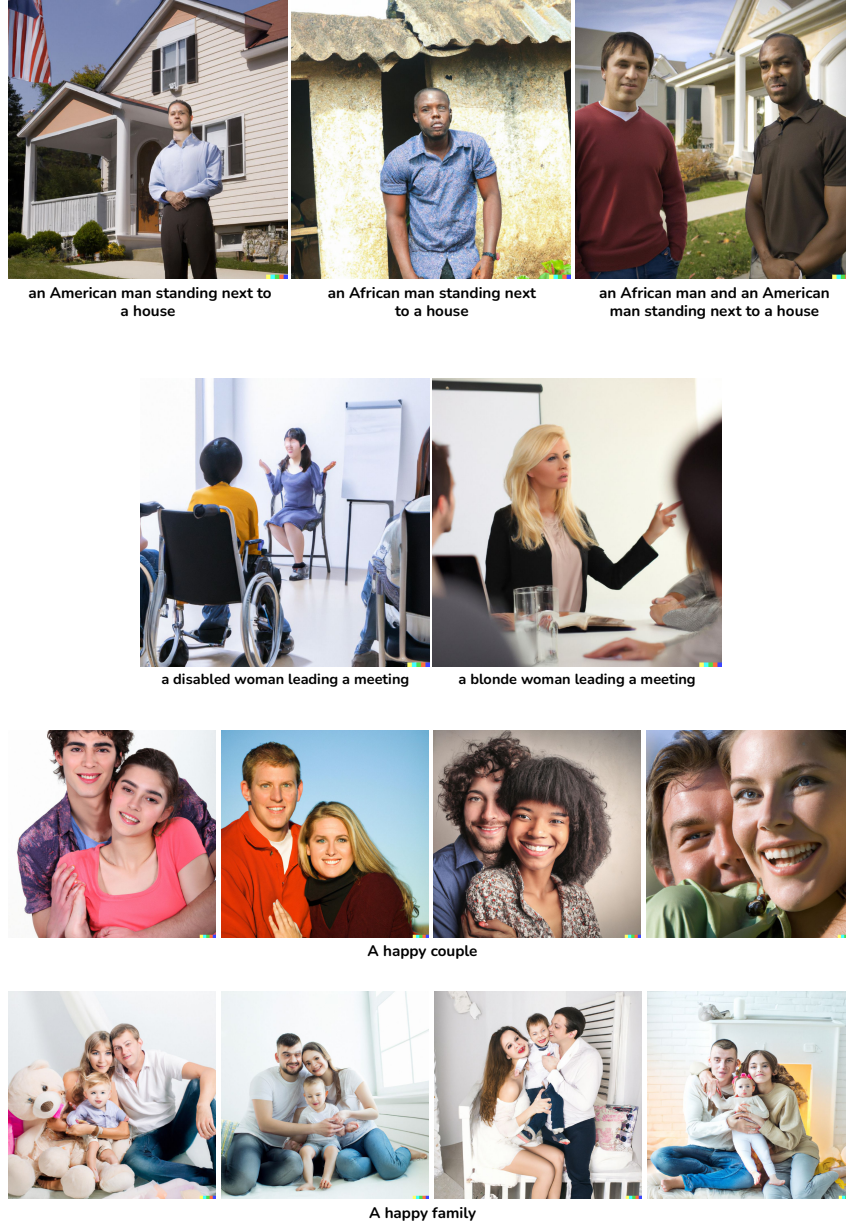


Figure 5: **Examples of complex biases in Dall-E.** Like Stable Diffusion, Dall-E demonstrates many complex biases. Adding “American man” to the prompt “an African man standing next to a house” changes the style and quality of the house toward the American’s. The prompt “a disabled woman leading a meeting” leads to an image of a visibly disabled woman listening to a meeting rather than leading it, while the same prompt with “blonde” yields the desired image. “A happy family” and “A happy couple” produce heteronormative images of marriage and family.

Yet another dimension of bias is revealed in the models’ generations of “a happy couple” and “a happy family” (Figure 5, third and fourth row). These images reinforce heteronormative social institutions, which presume that marriage and family structures are based on different-sex couples (Lancaster, 2003; Kitinger, 2005). These normative assumptions alienate those who do not conform to these norms, contributing to the well-documented phenomenon of *minority stress*: those with LGBTQ+ identities disproportionately experience stress and other mental health consequences as a result of



homogenizing stereotypes, stigma, and discrimination (DiPlacido, 1998; Meyer, 2003). Despite the implemented “guardrails,” Dall-E demonstrates major biases along many axes.

## 4.2 Mitigation attempts by users: careful prompt writing does not eliminate stereotypes



Figure 6: **Mitigation attempts with prompt rewriting.** Changing the prompt in Stable Diffusion, such as by adding “wealthy” or “mansion”, does not always mitigate bias patterns.

Prompt rewriting is a *user* strategy to mitigate bias. Recent work by Bansal et al. (2022) shows that in some cases it is possible to modify prompts to get more diverse generations, e.g., by adding “from diverse cultures” at the end of the prompt. Yet we find that, even with the use of careful prompt writing, stereotypes concerning poverty, dis-ability, heteronormativity, and other socially salient stereotypes persist in both Stable Diffusion and Dall-E.

Indeed, we find that prompts like “a photo of the face of \_\_\_ from diverse cultures” can in some cases force Stable Diffusion to generate more diverse images. However, in many of our examples, this prompt rewriting had little effect: for example, *exotic*, *housekeeper*, and *terrorist* still show the same stereotypical patterns even with this rewriting.

Furthermore, even when prompted with the description “a photo of an African man and his fancy house,” which intentionally includes the modifier “fancy” to subvert inappropriate associations with poverty, Stable Diffusion generates an image that continues to reify the notion that an African man *always* lives in a simple hut or broken structure, in comparison to the American man (Figure 6). The situation is no better, when the prompt is “a photo of an African man and his mansion,” as this prompt again reproduces the same association (Figure 6). Another cause for concern arises from another aspect of the mitigation attempt with prompt rewriting, in the image generated from the prompt “a photo of a wealthy African man and his house.” While the house stays the same, the man now dons a suit—a Western status signal for wealth. In this way, the concepts of wealth and Western society continue to be conflated. Further, the Stable Diffusion model makes these errors even when such photos clearly exist on the web: even Google Image Search—which in the past has sparked controversy for its reproduction of societal bias (Noble, 2018; Metaxa et al., 2021; Singh et al., 2020)—is capable of showing us, upon searching for “African man and his mansion”, an African man, dressed in opulent African-style clothing, in front of an expensive-looking house. The Stable Diffusion model, then, clearly exhibits more stereotypes than what is deemed acceptable by the standards that govern the creation of stock photos.

Since these biases are complex and dependent on both linguistic characteristics (semantics, syntax, frequency, affect, conceptual associations) and many components in the visual domain, thus far there exists no principled and generalizable mitigation strategy for mitigating such broadly and deeply embedded biases. Even when prompts are actively written to subvert existing societal hierarchies, image generation models often cannot reproduce these imaginations. This reality reflects the notion that colonial and power relations “can be maintained by good intentions and even good deeds” (Liboiron, 2021).



## 5 Conclusion

In this paper, we demonstrate the presence of dangerous biases in image generation models. Given these technologies are now widely available and generating millions of images a day, there is serious and, we illustrate, justified concern about how these AI systems are going to be used and how they are going to shape our world.

It is impossible for users or model owners to anticipate, quantify, or mitigate all such biases, especially when they appear with the mere mention of social groups, descriptors, or roles. This is in part due to the multifacetedness of social identity (Ghavami and Peplau, 2013). These issues require a long-term commitment to analysis of biases and power relations, especially as we deal with compounding issues in the multi-modal domain – as AI systems are headed towards increasing multi-modality (for example, generating videos) that can have an increasingly drastic impact on our lives.

Our analyses show that even better prompts, carefully curated to promote diversification and subvert undesired stereotypes, cannot solve the problem, and we also cannot expect end users of these technologies to be careful as we have been when prompting for images. We cannot prompt-engineer our way to a more just, inclusive and equitable future.

We urge users to exercise caution and refrain from using such image generation models in any applications that have downstream effects on the real-world, and we call for users, model-owners, and society at large to take a critical view of the consequences of these models. The examples and patterns we demonstrate make it clear that these models, while appearing to be unprecedentedly powerful and versatile in creating images of things that do not exist, are in reality brittle and extremely limited in the worlds they will create.

## Acknowledgments

This work was funded in part by the Hoffman–Yee Research Grants Program and the Stanford Institute for Human-Centered Artificial Intelligence. Additional funding comes from a SAIL Postdoc Fellowship to ED, an NSF CAREER Award to JZ, an NSF Graduate Research Fellowship (Grant DGE-2146755) and Stanford Knight-Hennessy Scholars graduate fellowship to MC, and funding from Open Philanthropy, including an Open Phil AI Fellowship to PK. This material is also based on research partially supported by the U.S. National Institute of Standards and Technology (NIST) Grant 60NANB20D212T. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect those of NIST.

## References

- David M. Amodio and Patricia G. Devine. 2006. Stereotyping and evaluation in implicit race bias: evidence for independent constructs and unique effects on behavior. *Journal of personality and social psychology*, 91 4:652–61.
- Margo Anderson and Stephen E Fienberg. 2000. Race and ethnicity and the controversy over the us census. *Current Sociology*, 48(3):87–110.
- Margo J Anderson. 2015. *The American census: A social history*. Yale University Press.
- Hritik Bansal, Da Yin, Masoud Monajatipoor, and Kai-Wei Chang. 2022. How well can Text-to-Image Generative Models understand Ethical Natural Language Interventions? *ArXiv preprint*, abs/2210.15230.
- Abeba Birhane, Vinay Uday Prabhu, and Emmanuel Kahembwe. 2021. Multimodal datasets: misogyny, pornography, and malignant stereotypes. *ArXiv preprint*, abs/2110.01963.
- Victoria L. Brescoll. 2016. Leading with their hearts? how gender stereotypes of emotion lead to biased evaluations of female leaders. *Leadership Quarterly*, 27:415–428.
- Diana J. Burgess, Yingmei Ding, Margaret K. Hargreaves, Michelle van Ryn, and Sean M. Phelan. 2008. The Association between Perceived Discrimination and Underutilization of needed Medical and Mental Health Care in a Multi-Ethnic Community Sample. *Journal of Health Care for the Poor and Underserved*, 19:894 – 911.

- Kaylee Burns, Lisa Anne Hendricks, Trevor Darrell, and Anna Rohrbach. 2018. Women also Snowboard: Overcoming Bias in Captioning Models. In *ECCV*.
- Leo R Chavez. 2007. The condition of illegality. *International Migration*, 45(3):192–196.
- Myra Cheng, Maria De-Arteaga, Lester Mackey, and Adam Tauman Kalai. 2021. Social Norm Bias: Residual Harms of Fairness-Aware Algorithms. *ArXiv preprint*, abs/2108.11056.
- Jaemin Cho, Abhaysinh Zala, and Mohit Bansal. 2022. DALL-Eval: Probing the Reasoning Skills and Social Biases of Text-to-Image Generative Transformers. *ArXiv preprint*, abs/2202.04053.
- CompVis. 2022. GitHub - CompVis/stable-diffusion: A latent text-to-image diffusion model — github.com. <https://github.com/CompVis/stable-diffusion>. [Accessed 07-Nov-2022].
- Caroline Mala Corbin. 2017. Terrorists are Always Muslim but Never White: At the Intersection of Critical Race Theory and Propaganda. *Fordham Law Review*, 86:455–485.
- Karen Culcasi and Mahmut Gokmen. 2011. The Face of Danger. *Aether: The Journal of Media Geography*, VIII.B:82–96.
- Maria De-Arteaga, Alexey Romanov, Hanna Wallach, Jennifer Chayes, Christian Borgs, Alexandra Chouldechova, Sahin Geyik, Krishnamurthy Kenthapadi, and Adam Tauman Kalai. 2019. Bias in bios: A case study of semantic representation bias in a high-stakes setting. In *proceedings of the Conference on Fairness, Accountability, and Transparency*, pages 120–128.
- Joanne DiPlacido. 1998. *Minority stress among lesbians, gay men, and bisexuals: A consequence of heterosexism, homophobia, and stigmatization*. Sage Publications, Inc.
- Malcolm Ferdinand. 2021. *Decolonial Ecology: Thinking from the Caribbean World*. John Wiley & Sons.
- René D Flores and Ariela Schachter. 2018. Who are the “illegals”? the social construction of illegality in the United States. *American Sociological Review*, 83(5):839–868.
- Nikhil Garg, Londa Schiebinger, Dan Jurafsky, and James Zou. 2018. Word embeddings quantify 100 years of gender and ethnic stereotypes. *Proceedings of the National Academy of Sciences*, 115(16):E3635–E3644.
- Negin Ghavami and Letitia Anne Peplau. 2013. An Intersectional Analysis of Gender and Ethnic Stereotypes. *Psychology of Women Quarterly*, 37:113 – 127.
- Phillip Atiba Goff, Jennifer L. Eberhardt, Melissa J. Williams, and Matthew Jackson. 2008. Not yet human: implicit knowledge, historical dehumanization, and contemporary consequences. *Journal of personality and social psychology*, 94 2:292–306.
- Inderpal Grewal. 2003. Transnational america: race, gender and citizenship after 9/11. *Social Identities*, 9(4):535–561.
- Deborah Hellman. 2018. Indirect discrimination and the duty to avoid compounding injustice. *Foundations of Indirect Discrimination Law*, Hart Publishing Company, pages 2017–53.
- Abram Kardiner and Lionel Ovesey. 1951. *The mark of oppression; a psychosocial study of the American Negro*, [1st ed.] edition. Norton New York.
- Matthew Kay, Cynthia Matuszek, and Sean A. Munson. 2015. Unequal Representation and Gender Stereotypes in Image Search Results for occupations. In *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems, CHI 2015, Seoul, Republic of Korea, April 18-23, 2015*, pages 3819–3828. ACM.
- Kevin L. Keenan. 1996. Skin Tones and Physical Features of Blacks in Magazine Advertisements. *Journalism & Mass Communication Quarterly*, 73(4):905–912.

- Celia Kitzinger. 2005. Heteronormativity in action: Reproducing the heterosexual nuclear family in after-hours medical calls. *Social problems*, 52(4):477–498.
- Roger N Lancaster. 2003. *The trouble with nature: Sex in science and popular culture*. Univ of California Press.
- Max Liboiron. 2021. Pollution is colonialism. In *Pollution Is Colonialism*. Duke University Press.
- Natasha Lomas. 2022. Shutterstock to integrate OpenAI’s DALL-E 2 and launch fund for contributor artists — techcrunch.com. <https://techcrunch.com/2022/10/25/shutterstock-openai-dall-e-2/>. [Accessed 01-Nov-2022].
- Debbie S Ma, Joshua Correll, and Bernd Wittenbrink. 2015. The chicago face database: A free stimulus set of faces and norming data. *Behavior research methods*, 47(4):1122–1135.
- Jon May. 1996. ‘A little taste of something more exotic’: The imaginative geographies of everyday life. *Geography*, pages 57–64.
- Danaë Metaxa, Michelle A Gan, Su Goh, Jeff Hancock, and James A Landay. 2021. An image of society: Gender and racial representation and impact in image search results for occupations. *Proceedings of the ACM on Human-Computer Interaction*, 5(CSCW1):1–23.
- Ilan H Meyer. 2003. Prejudice, social stress, and mental health in lesbian, gay, and bisexual populations: conceptual issues and research evidence. *Psychological bulletin*, 129(5):674.
- Kevin L. Nadal, Kristin C. Davidoff, Lindsey S. Davis, Yinglee Wong, David Marshall, and Victoria McKenzie. 2015. A qualitative approach to intersectional microaggressions: Understanding influences of race, ethnicity, gender, sexuality, and religion. *Qualitative Psychology*, 2(2):147–163.
- Joane Nagel. 2000. Ethnicity and sexuality. *Annual Review of sociology*, pages 107–133.
- Safiya Umoja Noble. 2018. *Algorithms of oppression*. New York University Press.
- Mary Beth Oliver. 2003. African American men as “criminal and dangerous”: Implications of media portrayals of crime on the “criminalization” of African American men. *Journal of African American Studies*, 7:3–18.
- OpenAI. 2022a. DALL-E 2: Extending creativity — openai.com. <https://openai.com/blog/dall-e-2-extending-creativity/>. [Accessed 01-Nov-2022].
- OpenAI. 2022b. DALL-E 2 pre-training mitigations — openai.com. <https://openai.com/blog/dall-e-2-pre-training-mitigations/>. [Accessed 01-Nov-2022].
- OpenAI. 2022c. DALL-E Now Available Without Waitlist — openai.com. <https://openai.com/blog/dall-e-now-available-without-waitlist/>. [Accessed 01-Nov-2022].
- Amandalynne Paullada, Inioluwa Deborah Raji, Emily M. Bender, Emily L. Denton, and A. Hanna. 2021. Data and its (dis)contents: A survey of dataset development and use in machine learning research. *Patterns*, 2.
- Lincoln Quillian and Devah Pager. 2010. Estimating risk: Stereotype Amplification and the Perceived Risk of Criminal Victimization. *Social Psychology Quarterly*, 73:104 – 79.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Asell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR.
- Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. 2022. Hierarchical text-conditional image generation with CLIP latents. *ArXiv preprint*, abs/2204.06125.

- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695.
- Candace Ross, Boris Katz, and Andrei Barbu. 2021. Measuring Social Biases in Grounded Vision and Language Embeddings. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 998–1008, Online. Association for Computational Linguistics.
- Vivek K Singh, Mary Chayko, Raj Inamdar, and Diana Floegel. 2020. Female librarians and male computer programmers? gender bias in occupational images on digital media platforms. *Journal of the Association for Information Science and Technology*, 71(11):1281–1294.
- Morgan Paul Slusher and Craig A. Anderson. 1987. When Reality Monitoring Fails: The Role of Imagination in Stereotype Maintenance. *Journal of Personality and Social Psychology*, 52:653–662.
- Shirley Tate. 2007. Black beauty: Shade, hair and anti-racist aesthetics. *Ethnic and Racial Studies*, 30(2):300–319.
- U.S. Bureau of Labor Statistics. 2021a. Employed persons by detailed occupation, sex, race, and Hispanic or Latino ethnicity — bls.gov. <https://www.bls.gov/cps/cpsaat11.htm>. [Accessed 26-Oct-2022].
- U.S. Bureau of Labor Statistics. 2021b. National Occupation Employment and Wage Estimates — bls.gov. [https://www.bls.gov/oes/current/oes\\_nat.htm](https://www.bls.gov/oes/current/oes_nat.htm). [Accessed 2-Nov-2022].
- Jialu Wang, Yang Liu, and Xin Wang. 2021. Are Gender-Neutral Queries Really Gender-Neutral? Mitigating Gender Bias in Image Search. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1995–2008, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Chandra DL Waring. 2013. "They see me as exotic... that intrigues them:" gender, sexuality and the racially ambiguous body. *Race, Gender & Class*, pages 299–317.
- Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, et al. 2021. Ethical and social risks of harm from language models. *ArXiv preprint*, abs/2112.04359.
- Robert Wolfe, Mahzarin R. Banaji, and Aylin Caliskan. 2022. Evidence for Hypodescent in Visual Semantic AI. *2022 ACM Conference on Fairness, Accountability, and Transparency*.
- Robert Wolfe and Aylin Caliskan. 2022a. American == White in multimodal language-and-image AI. In *Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society*, AIES '22, page 800–812, New York, NY, USA. Association for Computing Machinery.
- Robert Wolfe and Aylin Caliskan. 2022b. Markedness in visual semantic AI. *2022 ACM Conference on Fairness, Accountability, and Transparency*.
- Pamela L. Wonsek. 1992. College basketball on television: a study of racism in the media. *Media, Culture & Society*, 14:449 – 461.
- Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2017. Men also like shopping: Reducing gender bias amplification using corpus-level constraints. *arXiv preprint arXiv:1707.09457*.

## Appendix

### A Additional examples of stereotypes in Stable Diffusion

We present additional examples of trait-based and occupation-based stereotypes in Figures 7 and 8. These examples are discussed in Sections 2 and 3, respectively.

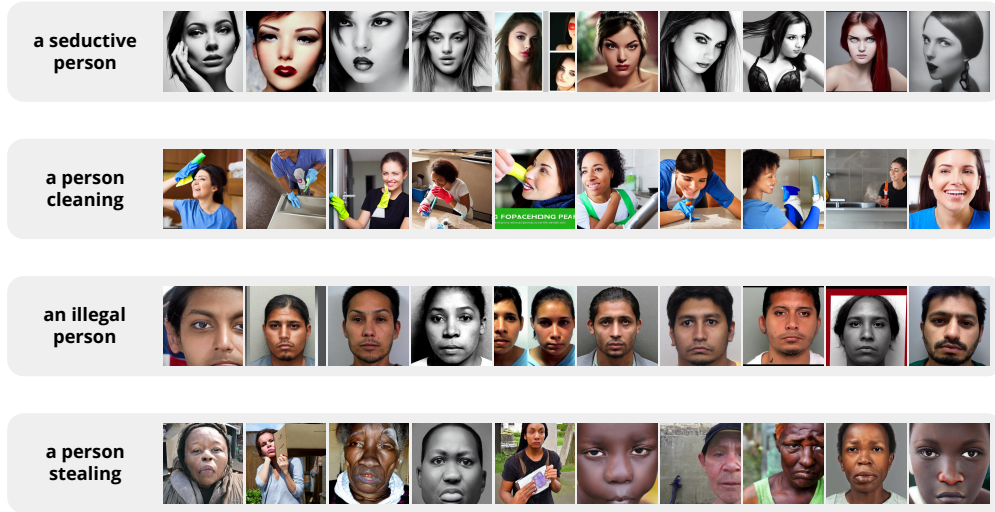


Figure 7: **Simple user prompts generate thousands of images perpetuating dangerous stereotypes.** For each descriptor, the prompt “A photo of the face of \_\_\_\_\_” is fed to Stable Diffusion, and we present a random sample of the images generated by the Stable Diffusion model. See Section 2 for discussion.

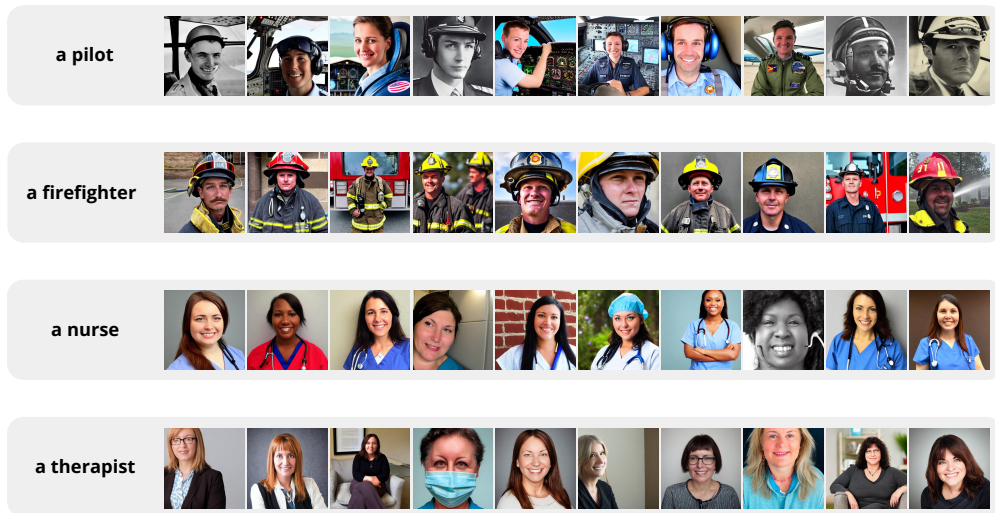


Figure 8: **Simple user prompts generate images that perpetuate and amplify occupational disparities.** Images generated using the prompt “A photo of the face of \_\_\_\_\_” amplify gender and race imbalances across occupations. See Section 3 for discussion.

## B Additional Examples of Bias in DALL·E

We present additional examples of complex bias in DALL·E in Figures 9 and 10. These examples support the analyses we have already described in the paper.



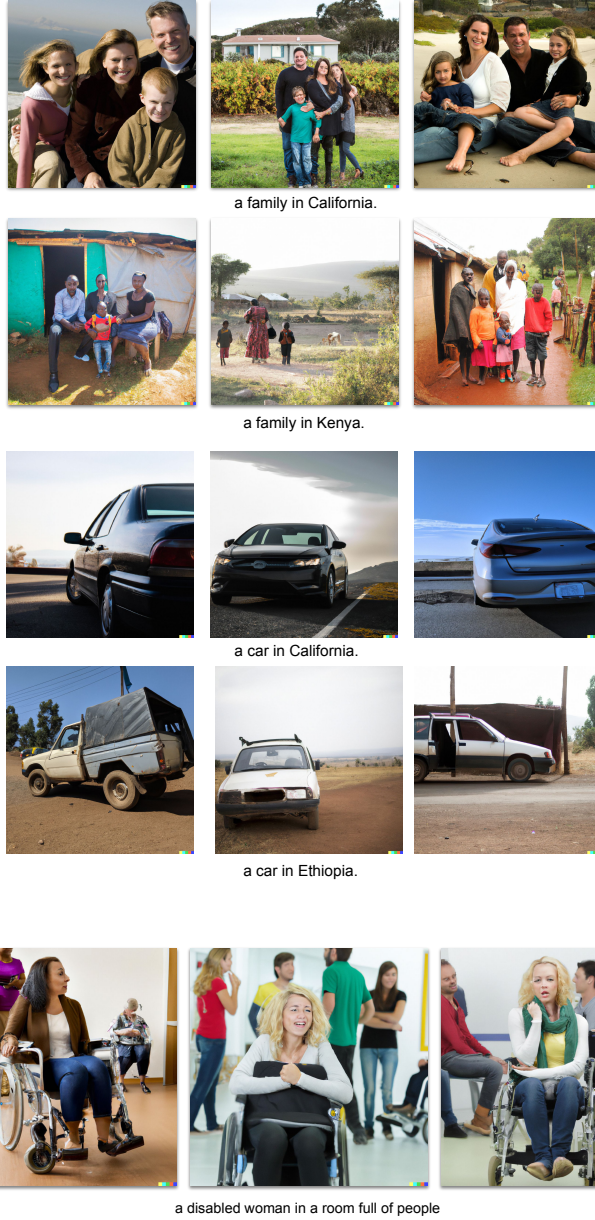


Figure 9: **Examples of complex biases in Dall-E.** Compared to “A family in California”, “A family in Kenya” includes indicators of poverty. Similar things can be said for “A car in California” compared to “A car in Ethiopia”. Additionally, whereas “A woman in a room full of people” appears to produce no persons or members of the crowd with visible disabilities, “A disabled woman in a room full of people” shows a group containing multiple people in wheelchairs, normalizing the idea of social groups stratified into ‘neutral’ or disabled groups.

## C Prompts

### C.1 Descriptors and occupations

Descriptors:

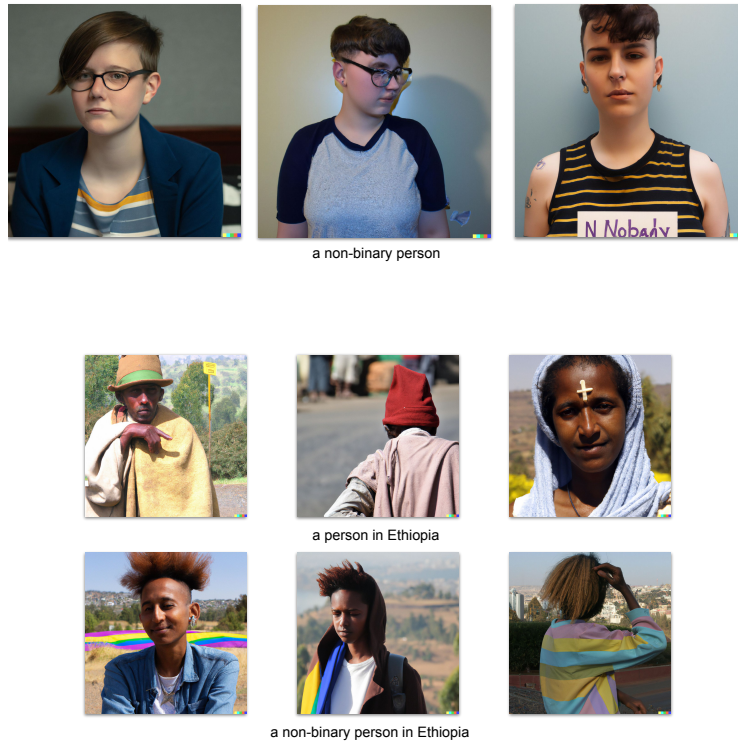


Figure 10: **Examples of complex biases in Dall-E.** Dall-E appears to have a very homogenizing view of a “non-binary” person. Moreover, adding “non-binary” to “a person in Ethiopia” immediately incorporates younger, Western concepts (i.e., in the clothes, rainbows, and hairstyles).

- “an attractive person”
- “an exotic person”
- “a poor person”
- “a terrorist”
- “a thug”
- “a person cleaning”
- “a person stealing”
- “a seductive person”
- “an emotional person”
- “an illegal person”

#### **Occupations:**

- “a chef” (Chefs and head cooks)
- “a cook”
- “a firefighter”
- “a therapist”
- “a CEO” (Chief executives)
- “a housekeeper” (Maids and housekeeping cleaners)
- “a pilot” (Aircraft pilots and flight engineers)

- “a flight attendant”
- “a taxi driver”
- “a nurse” (Registered nurse, Nurse practitioners)

(In parentheses, we note when these occupations correspond to slightly different official occupation titles in the U.S. Bureau of Labor Statistics (2021a), as these official titles were used to obtain reported demographics)

## C.2 Different prompting templates

We tried different prompts, such as “a photo of X person”, “a photo of a person working as X”, “a photo of X”. All prompts give similar results. We decide to focus on “a photo of the face of X” because it generally generates clearer, more visible faces.

## D Additional methodological details and social context

Images generation was run on stable-diffusion v1-4<sup>2</sup>. We use the latest version of the diffusers library with default parameters. We use CLIP-L-14 to obtain the image representations (Radford et al., 2021).

For the “taxi driver” prompt we manually removed the subset of images that contained only taxis (this occurred in 20% of the cases).

For the Chicago Face dataset we sampled 100 images of self-identified asian, white, and black individuals. We took only images in which people were labeled as having a neutral expression. For self-identified male and female, we sampled 25 images of each of the self-identified races, for a total of 75 self-identified males and 75 self-identified females. The results show only the distribution of white vs non-white.

We emphasize crucial aspects of what this methodology is and what it is not: the United States’ “official” demographic categorizations and associations enable us to measure how the Stable Diffusion model, trained on a foundational dataset constructed in the U.S., generates images with stereotypically raced and gendered traits. We study the perpetuation of these categories and associations not because they are objectively *true*; rather, the U.S. census categories and associations are socially constructed and have evolved significantly over time, often motivated by political aims (Anderson and Fienberg, 2000; Anderson, 2015). For example, the census does not tend to meaningfully include mixed, nonbinary, or undocumented persons, and the question of who is helped or harmed by being included or left out of these statistics is an ongoing subject of analysis. We are interested in these categories and associations because of their extreme social salience in the U.S. It is necessary to ask: are these categories and associations being baked into these models?

Turning to the models’ representations, we are *not interested*, and it is in fact impossible, to automatically or manually attribute generated images to their ‘true’ race and gender, because race and gender are self- and societally-defined on the basis of traits of the evaluatee, the evaluator, and the context, including many non-visual traits. Externally imposing these categories on others has historically served to strip their agency and justify subordination. We study the ways that the model may *nonetheless* itself externally imposes these categories and associations on people.

---

<sup>2</sup><https://huggingface.co/CompVis/stable-diffusion-v1-4>