# Mitigate One, Skew Another? Tackling Intersectional Biases in Text-to-Image Models

Pushkar Shukla $^{*1}$  Aditya Chinchure $^{*2,3}$  Emily Diana $^4$  Alexander Tolbert $^5$  Kartik Hosanagar $^6$  Vineeth N. Balasubramanian $^7$  Leonid Sigal $^{2,3}$  Matthew A. Turk $^1$ 

<sup>1</sup>Toyota Technological Institute at Chicago
 <sup>2</sup>University of British Columbia
 <sup>3</sup>Vector Institute for AI
 <sup>4</sup>Carnegie Mellon University, Tepper School of Business
 <sup>5</sup>Emory University
 <sup>6</sup>University of Pennsylvania, The Wharton School
 <sup>7</sup>Indian Institute of Technology Hyderabad

{pushkarshukla, mturk}@ttic.edu {aditya10, lsigal}@cs.ubc.ca

#### **Abstract**

The biases exhibited by text-to-image (TTI) models are often treated as independent, though in reality, they may be deeply interrelated. Addressing bias along one dimension—such as ethnicity or age—can inadvertently affect another, like gender, either mitigating or exacerbating existing disparities. Understanding these interdependencies is crucial for designing fairer generative models, yet measuring such effects quantitatively remains a challenge. To address this, we introduce BiasConnect, a novel tool for analyzing and quantifying bias interactions in TTI models. BiasConnect uses counterfactual interventions along different bias axes to reveal the underlying structure of these interactions and estimates the effect of mitigating one bias axis on another. These estimates show strong correlation (+0.65) with observed postmitigation outcomes.

Building on BiasConnect, we propose InterMit, an intersectional bias mitigation algorithm guided by user-defined target distributions and priority weights. InterMit achieves lower bias (0.33 vs. 0.52) with fewer mitigation steps (2.38 vs. 3.15 average steps), and yields superior image quality compared to traditional techniques. Although our implementation is training-free, InterMit is modular and can be integrated with many existing debiasing approaches for TTI models, making it a flexible and extensible solution.

#### 1 Introduction

Text-to-Image (TTI) models such as DALL-E (Ramesh et al., 2021), Imagen (Saharia et al., 2022), and Stable Diffusion (Rombach et al., 2022) have become widely used for generating visual content from textual prompts. Despite their impressive capabilities, these models often inherit and amplify biases present in their training data (Wang et al., 2022b; Chinchure et al., 2024; Cho et al., 2023). These biases manifest across multiple social and

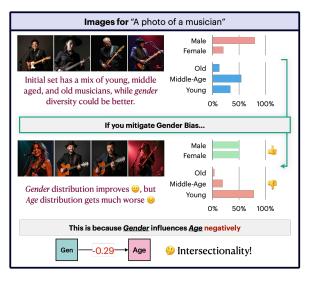


Figure 1: An example for which BiasConnect estimates a negative impact of bias mitigation along one axis on another axis. For this query, increasing the gender diversity (Gen) skews age distribution (Age) for images of musicians generated by Flux-dev.

non-social dimensions – including gender, race, clothing, and age - leading to skewed or inaccurate representations. As a result, TTI models may reinforce harmful stereotypes and societal norms (Bender et al., 2021; Birhane and Prabhu, 2021). While significant efforts have been made to evaluate and mitigate societal biases in TTI models (Wang et al., 2023a; Cho et al., 2023; Ghosh and Caliskan, 2023; Esposito et al., 2023; Bianchi et al., 2023; Chinchure et al., 2024), these approaches often assume that biases along different dimensions (e.g., gender and race) are independent of each other. Consequently, they do not account for relationships between these dimensions. For instance, as illustrated in Figure 1, mitigating gender (male, female) may effectively diversify the gender distribution in a set of generated images, but this mitigation step may negatively impact the diversity of another bias dimension, such as age. This relationship between two bias dimensions highlights the

intersectional nature of these biases.

The concept of intersectionality, first introduced by Crenshaw (Crenshaw, 1989), motivates the need to understand how overlapping social identities such as race, gender, and class contribute to systemic inequalities. In TTI models, these intersections can have a significant impact. As a motivating study, we independently mitigated eight bias dimensions over 26 occupational prompts on Stable Diffusion 1.4, using a popular bias mitigation strategy, ITI-GEN (Zhang et al., 2023) (see A.5). We found that while the targeted biases were reduced in most cases, biases along other axes were negatively affected in over 29% of the cases. This suggests that for an effective bias mitigation strategy, it is crucial to understand which biases are intersectional. Additionally, it is important to strive towards building a more holistic bias mitigation algorithm that can either mitigate multiple biases simultaneously or predict what biases cannot be mitigated together.

To understand how biases in TTI models influence one another, we propose BiasConnect, the first analysis tool that evaluates biases while explicitly modeling their intersectional relationships. Unlike prior methods that treat biases in isolation, BiasConnect identifies how mitigating one bias can positively or negatively affect others. Specifically, BiasConnect uses a novel metric, the Intersectional Sensitivity (IS), to quantify how mitigation along one axis affects others. These IS scores show a strong correlation (+0.65) with observed intersectional outcomes post-mitigation. We validate our approach through robustness studies and qualitative analyses, demonstrating its utility for auditing open-source TTI models.

Furthermore, we extend BiasConnect with a holistic intersectional bias mitigation algorithm, InterMit. While we propose an effective and straightforward implementation in this paper, InterMit is modular and can be integrated with any existing sequential bias mitigation method. Unlike prior approaches that assume fixed ideal distributions and treat all biases equally, InterMit allows users to define arbitrary target distributions, select specific bias axes, and assign custom priority weights to each bias—enabling flexible joint mitigation and informed reasoning about conflicting biases.

In our evaluation, InterMit outperforms existing methods by mitigating biases more effectively, producing higher-quality images, and requiring fewer mitigation steps. Moreover, unlike other methods, it can handle a larger number (> 3) of bias axes and alerts users when mitigation along one axis adversely affects others.

#### 2 Related Work

#### 2.1 Intersectionality and Bias in AI

Intersectionality, introduced by Crenshaw (Crenshaw, 1989), describes how multiple forms of oppression-such as racism, sexism, and classism—intersect to shape unique experiences of discrimination. Two key models define this concept: the additive model, where oppression accumulates across marginalized identities, and the interactive model, where these identities interact synergistically, creating effects beyond simple accumulation (Curry, 2018). In the context of AI, most existing work (Diana and Tolbert, 2023; Kavouras et al., 2023; Kearns et al., 2018) aligns more closely with the additive model, focusing on quantifying and mitigating biases in intersectional subgroups. This perspective has influenced fairness metrics (Diana et al., 2021; Foulds et al., 2020; Ghosh et al., 2021) designed to assess subgroup-level performance, extending across various domains, including natural language processing (NLP) (Lalor et al., 2022; Lassen et al., 2023; Guo and Caliskan, 2021; Tan and Celis, 2019) and recent large language models (Kirk et al., 2021; Ma et al., 2023; Devinney et al., 2024; Bai et al., 2025), multimodal research (Howard et al., 2024; Hoepfinger, 2023), and computer vision (Wang et al., 2020; Steed and Caliskan, 2021). These approaches typically measure disparities across predefined demographic intersections and propose mitigation strategies accordingly. Our work aligns with the interactive model of intersectionality, using counterfactual analysis with TTI models, where we intervene on a single bias axis to assess its ripple effects on others.

#### 2.2 Bias in Text-to-Image Models

Extensive research has been conducted on evaluating and mitigating social biases in both image-only models (Buolamwini and Gebru, 2018; Seyyed-Kalantari et al., 2021; Hendricks et al., 2018; Meister et al., 2023; Wang et al., 2022a; Liu et al., 2019; Joshi et al., 2022; Wang and Russakovsky, 2023) and text-only models (Bolukbasi et al., 2016; Hutchinson et al., 2020; Shah et al., 2020; Garrido-Muñoz et al., 2021; Ahn and Oh, 2021). More recently, efforts have expanded to multimodal models

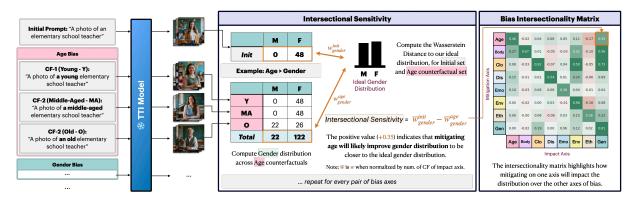


Figure 2: **An overview of BiasConnect**. We use a counterfactual-based approach to measure how interventions along a single bias axes impact other bias axes. Our metric Intersectional Sensitivity estimates how bias mitigation on one axis impacts another. Our results are visualized as a matrix called the Bias Intersectionality Matrix.

and datasets, addressing biases in various languagevision tasks. These investigations have explored biases in embeddings (Hamidieh et al., 2023), text-toimage generation (Cho et al., 2023; Bianchi et al., 2023; Seshadri et al., 2023; Ghosh and Caliskan, 2023; Zhang et al., 2023; Wang et al., 2023a; Esposito et al., 2023), image retrieval (Wang et al., 2022c), image captioning (Hendricks et al., 2018; Zhao et al., 2021), and visual question-answering models (Park et al., 2020; Aggarwal and Bhargava, 2023; Hirota et al., 2022).

Despite these advances, research on intersectional biases in TTI models remains limited. Existing evaluation frameworks such as T2IAT (Wang et al., 2023a), DALL-Eval (Cho et al., 2023), and other studies (Ghosh and Caliskan, 2023; Bianchi et al., 2023; Friedrich et al., 2023) primarily assess biases along predefined axes, such as gender (Wang et al., 2023a; Cho et al., 2023; Esposito et al., 2023; Bianchi et al., 2023), skin tone (Wang et al., 2023a; Cho et al., 2023; Ghosh and Caliskan, 2023; Esposito et al., 2023; Bianchi et al., 2023), culture (Esposito et al., 2023; Wang et al., 2023a), and geographical location (Esposito et al., 2023). While these works offer key insights into single-axis bias detection and mitigation, they lack a systematic examination of how biases on one axis influence another—a core aspect of intersectionality. The closest research, TIBET (Chinchure et al., 2024), visualizes such interactions, but our approach goes further by systematically quantifying bias interactions, and using these interactions for mitigation.

#### 3 Approach

The objective of BiasConnect is to identify and quantify the intersectional effects of intervening

on one bias axis  $(B_x)$  to mitigate that bias, on any other bias axis  $(B_u)$ . BiasConnect works by systematically altering input prompts and analyzing the resulting distributions of generated images (see Fig. 2). To achieve this, we leverage counterfactual prompts by modifying specific attributes (e.g., male and female) along a bias axis (e.g., gender) and examine how these interventions impact other bias dimensions (e.g., age and ethnicity). If modifying one bias axis through counterfactual intervention causes significant shifts in the distribution of attributes along another bias axis, it indicates an intersectional dependency between these axes. We first construct prompt counterfactuals and generate images using a TTI model (Sec. 3.1). Subsequently, to identify bias-related attributes in the generated images, we use a Visual Question Answering (VQA) model (Sec. 3.2). Finally, to quantify the intersectional effects, and to identify whether these effects are positive or negative, we compute the Intersectional Sensitivity (Sec. 3.3).

# 3.1 Counterfactual Prompts & Image Generation

Given an input prompt P and bias axes  $B = \{B_1, B_2, \ldots, B_n\}$ , we generate counterfactual prompts  $\{CF_i^1, \ldots, CF_i^j\}$  for each bias  $B_i \in B$ . These counterfactual prompts may be templated (Appendix Table 3) or LLM-generated. The original prompt P and its counterfactuals are then used to generate images with the TTI model to measure intersectional effects.

#### 3.2 VQA-based Attribute Extraction

To facilitate the process of extracting bias-related attributes from the generated images, we use VQA. This is inspired by previous approaches on bias evaluation, like TIBET (Chinchure et al., 2024) and OpenBias (D'Incà et al., 2024), where a VQA-based method was used to extract concepts from generated images. Following TIBET, we use MiniGPT-v2 (Chen et al., 2023) in a question-answer format to extract attributes from generated images. We select MiniGPT-v2 over other VQA models because it is capable of answering bias-related questions, which other safety-tuned VQA models refuse, and is shown to be reliable at extracting attributes from sets of images in alignment with humans (Chinchure et al., 2024).

For the societal biases we analyze, we have a list of predefined questions (Appendix A.3) corresponding to each bias axis in B, and each question has a choice of attributes to choose from. For example, for the gender bias axis, we ask the question "[vqa] What is the gender (male, female) of the person?". Note that every question is multiple choice (in this example, male and female are the two attributes for gender). For datasets where counterfactuals are dynamically generated (e.g. TI-BET dataset), an LLM-generated set of questions is used instead. The questions asked for all images of prompt P and its counterfactuals  $CF_i^j$  remain the same. With the completion of this process, we have attributes for all images, where each image has one attribute for each bias axis in B.

## 3.3 Computing Intersectional Sensitivity

Our objective is to understand how the impact of interventions on  $B_x$  affects  $B_y$  in a positive or negative direction concerning an ideal distribution. To address this, we propose a metric that quantifies the impact of bias mitigation on dependent biases with respect to an ideal distribution.

**Defining an Ideal Distribution.** We first define a desired (ideal) distribution  $D^*$ , which represents the unbiased state we want bias axes to achieve. This can be a real-world distribution of a particular bias axis, a uniform distribution (which we use in our experiments), or anything that suits the demographic of a given sub-population.

**Measuring Initial Bias Deviation.** Given the images of initial prompt P, we compute the empirical distribution of attributes associated with bias axis  $B_y$ , denoted as  $D_{B_y}^{\rm init}$ . We then compute the Wasserstein distance between this empirical distribution and the ideal distribution:

$$w_{B_y}^{\text{init}} = W_1(D_{B_y}^{\text{init}}, D^*) \tag{1}$$

where  $W_1(\cdot,\cdot)$  represents the Wasserstein-1 dis-

tance. The Wasserstein-1 distance (also known as the Earth Mover's Distance) between two probability distributions  $D_1$  and  $D_2$  is defined as:

$$W_1(D_1, D_2) = \inf_{\gamma \in \Pi(D_1, D_2)} \mathbb{E}_{(x,y) \sim \gamma}[|x - y|]$$
 (2)

where  $\Pi(D_1, D_2)$  is the set of all joint distributions  $\gamma(x, y)$  whose marginals are  $D_1$  and  $D_2$ , and |x-y| represents the transportation cost between points in the two distributions.

We use  $\overline{w}_{B_y}^{init}$  to measure the amount of bias in the image set, where  $\overline{w}_{B_y}$  is computed by normalizing  $w_{B_y}$  based on the number of counterfactuals in  $B_y$ .  $\overline{w}_B^{init} \in [0,1]$  where 1 indicates that the distribution is completely biased and 0 indicates no bias

Intervening on  $B_x$ . Next, say we intervene on  $B_x$  to simulate the mitigation of bias  $B_x$ . This intervention ensures that all counterfactuals of  $B_x$  are equally represented in the generated images. For example, if  $B_x$  is gender bias, we enforce equal proportions of male and female individuals in the dataset. This intervention is in line with most bias mitigation methods proposed for TTI models, like ITI-GEN (Zhang et al., 2023). Using our counterfactuals along  $B_x$ , we sum the distributions on  $B_y$  across all counterfactuals of  $B_x$ . This sum across the counterfactuals of  $B_x$  yields a new empirical distribution of  $B_y$ , denoted  $D_{B_y}^{B_x}$ , simulating the effect of mitigating  $B_x$  (See Fig 2). We compute its Wasserstein distance from the ideal distribution.

$$w_{B_y}^{B_x} = W_1(D_{B_y}^{B_x}, D^*) (3)$$

Computing Intersectional Sensitivity. To quantify the effect of mitigating  $B_x$  on  $B_y$ , we define the metric, Intersectional Sensitivity, as:

$$IS_{xy} = \overline{w}_{B_y}^{\text{init}} - \overline{w}_{B_y}^{B_x} \tag{4}$$

as Wasserstein distance is sensitive to the number of counterfactuals, and  $IS_{xy} \in [-1,1]$ . A positive value  $(IS_{xy}>0)$  indicates that mitigating  $B_x$  improves  $B_y$ , bringing it closer to the ideal distribution, while a negative value  $(IS_{xy}<0)$  suggests it worsens  $B_y$ , moving it further from the ideal. If  $IS_{xy}=0$ , mitigating  $B_x$  has no effect on  $B_y$ . This approach enables us to assess whether addressing one bias (e.g., gender) improves or worsens another (e.g., ethnicity) in generative models, providing a systematic way to evaluate trade-offs and unintended consequences in bias mitigation strategies.

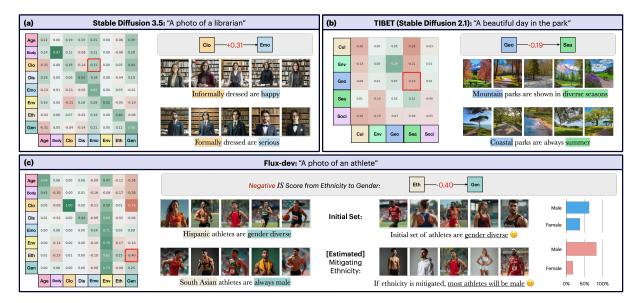


Figure 3: Analyzing bias intersectionality matrices from BiasConnect. (a) Shows how mitigating clothing bias also mitigates emotion bias. (b) Explores interactions between non-traditional bias axes in the TIBET dataset. (c) Reveals that generating ethnically diverse athletes reduces gender diversity. BiasConnect can allow the user the user to understand whether interventions along one dimension impact other dimensions positively or negatively.

#### 3.4 Visualization

To visualize IS scores comprehensively, we use a Bias Intersectionality Matrix S, where each entry  $IS_{ij}$  quantifies the effect of intervening on row  $B_i$  on column  $B_j$  for mitigation. This matrix captures directional dependencies and enables a structured analysis of intersectional bias effects.

# 4 Intersectional Bias Mitigation using BiasConnect

This section introduces an iterative strategy for mitigating intersectional biases, designed to be modular and compatible with existing sequential bias mitigation algorithms. We propose a subjective, user-guided mitigation framework built on top of BiasConnect, called InterMit. The framework allows users to select a subset of bias dimensions from a predefined set and assign mitigation priorities to each. Additionally, users can specify a desired target distribution that the model should conform to, providing both flexibility and control over the mitigation process. Our proposed framework (in Algorithm 1) leverages a matrix S at each step to iteratively reduce bias across multiple axes.

Given a TTI model M, a subset of selected axes by the user  $B^* \subseteq B$ , let  $\mathbf{p} \in \mathbb{R}^{|B^*|}$  be a user-defined *priority vector* that encodes the relative importance of mitigating bias along each axis, where  $|\mathbf{p}|_1 = 1$ . The ideal desired distribution  $D^*$ , for each bias axes in  $B^*$  is also specified by the user.

Given the aforementioned information, we first calculate a bias score for initial model  $M^{(0)}$  by taking the dot product of the priority vector  $\mathbf{p}$  and the initial measures of biases  $B^*$  ( $\overline{w}_{B*}^{init}$ ) computed using Eq. 3. This is denoted by  $\tau = \langle \overline{w}_{B*}^{init}, p \rangle$  and measures the overall bias of the model on  $B^*$  at any timestep. We proceed to mitigation if  $\tau$  is greater than a threshold  $\epsilon$ .

To choose which bias axis to mitigate on, we extract the submatrix  $\mathbf{S}' \in \mathbb{R}^{n \times |B^*|}$  consisting of the relevant columns from  $\mathbf{S}$  obtained using BiasConnect. For each row  $\mathbf{s}'_i$  of  $\mathbf{S}'$ , we compute a similarity score  $\gamma_i = \langle \mathbf{s}'_i, \mathbf{p} \rangle$ , which quantifies the alignment between the i-th intersectional bias and the desired direction of mitigation. The bias axis  $i^* = \arg\max_i \gamma_i$  with the highest alignment score is selected for targeted mitigation in the current iteration. The model is then updated to reduce bias along the direction corresponding to  $i^*$ , using a mitigation method, giving  $M^{(1)}$ . After mitigation, we generate a new set of images, recompute  $\tau$ , and continue the mitigation process if  $\tau > \epsilon$ .

### 5 Experiments

We evaluate BiasConnect for its ability to study intersectional biases across multiple models and prompts (Section 5.2, 5.3) and its robustness (5.4). Following that, we use InterMit for mitigation, and compare it to an existing strategy (5.5).

Algorithm 1 InterMit: Intersectional Mitigation

**Require:** Relevant bias axes  $B^* \subseteq B$ , priority vector  $\mathbf{p} \in \mathbb{R}^{|B^*|}$  with  $\|\mathbf{p}\|_1 = 1$ , sensitivity matrix  $\mathbf{S} \in \mathbb{R}^{n \times |B^*|}$ , bias threshold  $\epsilon$ , TTI model M

**Ensure:** Final mitigated model  $M^{(t)}$  with  $\tau < \epsilon$ 

1: Initialize model  $M^{(0)}$ , set iteration counter 2: repeat Extract submatrix  $\mathbf{S}' \in \mathbb{R}^{n \times |B^*|}$  from  $\mathbf{S}$ 3: Extract priority vector  $\mathbf{p} \in \mathbb{R}^{|B^*|}$ 4: for i = 1 to n do 5: 6: Compute similarity score  $\gamma_i \leftarrow \langle \mathbf{s}_i', \mathbf{p} \rangle$ 7: Identify target axis:  $i^* \leftarrow \arg \max_i \gamma_i$ 8: Mitigate axis  $i^*$  to update model:  $M^{(t+1)}$ 9: Compute bias score  $\tau^{(t+1)} = \langle \overline{w}_{B*}^{init}, p \rangle$ 10:  $t \leftarrow t + 1$ 11: 12: **until**  $au^{(t)} < \epsilon$ 

# 5.1 Experiment Setup

13: **return**  $M^{(t)}$ 

We conduct experiments on two prompt datasets, across six TTI models:

Occupation Prompts: To facilitate a structured evaluation, we develop a dataset with 26 occupational prompts, along eight distinct bias dimensions: gender, age, ethnicity, environment, disability, emotion, body type, and clothing. We generate 48 images for all initial counterfactual prompts using five TTI models: Stable Diffusion 1.4, Stable Diffusion 3.5, Flux (BlackForestLabs, 2024), Playground v2.5 (Li et al., 2024) and Kandinsky 2.2 (Shakhmatov et al., 2023; Razzhigaev et al., 2023). Further details about the prompts, bias axes, and counterfactuals are provided in the Appendix A.1.

<u>TIBET dataset</u>: The TIBET dataset includes 100 creative prompts with unique LLM-generated bias axes and counterfactuals (Chinchure et al., 2024) for each prompt, helping us test with a diverse array of biases. Additionally, it provides 48 Stable Diffusion 2.1-generated images per initial and counterfactual prompt (see Appendix A.6).

**Mitigation.** InterMit can use any sequential mitigation method, but we consider a simple training-free mitigation method using only prompt modifications (*PM*). At each mitigation step, we modify the initial prompt to introduce counterfactual concepts associated with the mitigated bias axis. Over multiple steps, we create collections of counter-

factual prompts that include all permutations of all mitigated axes (see A.7). We empirically set  $\epsilon=0.35$  for all our experiments. To compare our method to a traditional mitigation approach, we select ITI-GEN (Zhang et al., 2023), as it uses a similar FairToken-based permutation approach.

# 5.2 Studying prompt-level intersectionality

BiasConnect supports prompt-level analysis of intersectional biases (Fig. 3), helping users identify key bias axes and effective mitigation strategies. For example, in Fig. 3(a), Stable Diffusion 3.5 shows a causal link between clothing and emotion—informal attire leads to happier depictions of librarians (IS = 0.31), suggesting clothing changes can diversify emotional portrayal. In contrast, Fig. 3(c) shows ethnicity negatively affecting gender diversity, with South Asian athletes mostly depicted as male (IS = -0.40), indicating that addressing ethnicity alone may worsen gender bias. These insights support model comparison and targeted bias mitigation through InterMit.

#### 5.3 Validating Intersectional Sensitivity

Our approach estimates how counterfactualbased mitigation affects bias scores using the Intersectional Sensitivity. To validate this, we use ITI-GEN and PM to mitigate biases along each dimension, and measure the correlation between pre- and post-mitigation IS values. We achieve an average correlation of | +0.65 | across occupations using ITI-GEN. Certain axes like musician (+0.91), accountant (+0.81) and lawyer (+0.82) have especially high correlations. An average correlation of +0.95 using PM is unsurprising, as it uses similar counterfactual prompts for mitigation. The strong correlation observed between pre- and postmitigation bias scores suggests that our approach effectively estimates the potential impacts of bias mitigation, motivating the need to account for intersectionality in mitigation. More details on our experimental setup are in Appendix A.10.

#### 5.4 Robustness of BiasConnect

We analyze the robustness of our method by evaluating the impact of number of images (Fig. 4(a)) and VQA error rate (Fig. 4(b)) on Intersectional Sensitivity values. Our method uses 48 images per prompt to study bias distributions. Removing 8 images (16.6%) results in 10.5% change, an removing 32 images (66.6%) yields a 31.3%. This sub-linear impact suggests that TTI

Method	quality ↑	real ↑	natural ↑	colorfulness ↑	IsP?↑   MitAmt↓	MitSteps ↓
ITI-GEN (SD1.4) InterMit-PM (SD1.4)	0.73 <b>0.82</b>	0.92 <u>0.98</u>	0.37 <u>0.58</u>	0.45 <u>0.66</u>	92.8%   0.52 <u>99%</u>   <b>0.33</b>	100% <b>75.6%</b> (2.38/3.15)
InterMit-PM (SD3.5)	0.78	0.99	0.92	0.74	<b>100%</b>   0.27*	76%* (2.71/3.57)

Table 1: **Comparing our Mitigation Algorithm to ITI-GEN**. We mitigate a randomly chosen subset of 2-5 biases for prompts in the occupation set, and compute visual quality metrics and mitigation outcomes. We find that our algorithm uses 22% fewer mitigation steps, while still yielding higher mitigation amount and quality. \*Indicates we use a different prompt set and priority on SD3.5, so these should not be compared to SD1.4 results.

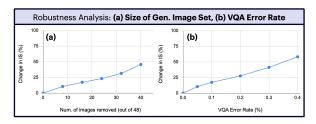


Figure 4: **Sensitivity analysis on BiasConnect**. We evaluate the robustness of our approach by analyzing the impact of VQA errors and the effect of the number of images on Intersectional Sensitivity.

models often generate similar bias distributions (e.g., always depicting nurses as females), preserving overall trends despite fewer images. Therefore, our approach is robust to moderate reductions in image count, but very small sets of images will significantly affect IS values. To test the robustness over VQA errors, we randomly change the VQA answers to a different answer (simulating an incorrect answer), from 5% to 40% of the time. We observe that even with low error rates of 5% and 10%, IS values change by 10% and 17.3% respectively. Here, the impact is compounded twice, because an error can skew the distribution away from one counterfactual towards another, and that a 5% error causes 13,478 answers out of a total of 269,568 answers to be changed, which is substantial. Nonetheless, we note that this impact remains linear. As VQA models improve, achieving low error rates for robustness becomes practical.

#### 5.5 Analyzing InterMit for Mitigation

To evaluate the effectiveness of our bias mitigation approach, we compare it against ITI-GEN (Zhang et al., 2023). ITI-GEN is designed for SD1.4 and is limited in its ability to mitigate more than three axes of bias for any given prompt. We override it in our experiments to facilitate a broader comparison. In contrast, our method combines *PM* with the intersectional mitigation algorithm InterMit.

**Prompts and Metrics.** For SD1.4, we randomly

Method	$\mathbf{MitAmt} \downarrow$	$\mathbf{MitSteps} \downarrow$
Random-PM (SD1.4)	0.340	84.6% (2.75/3.25)
InterMit-PM (SD1.4)	<b>0.339</b>	<b>69.2</b> % (2.25/3.25)

Table 2: **Ablation of mitigation strategy**. We achieve similar mitigation performance with fewer steps.

select subsets of biases to mitigate, assigning equal priority to each bias, in the occupation set. For SD3.5, we use 15 occupation prompts, targeting intersectional biases with some priorities weighted. Details are in Appendix A.8. Table 1 quantifies the mitigation amount (MitAmt: averaging  $\tau^T$  postmitigation across all prompts), and efficiency (MitSteps: ratio of number of biases mitigated to the number of biases in p). Visual quality is evaluated using CLIP-IQA metrics (Wang et al., 2023b) and using the VQA query: "[vqa] Is there a person in the image?" (IsP?).

**Results.** On SD1.4, InterMit-PM achieves significantly lower bias (0.33 vs. 0.52 for ITI-GEN), while requiring only 75.6% of the steps (ITI-GEN will always mitigate all biases in the priority vector), and producing higher-quality images (0.82 vs. 0.73). ITI-GEN frequently generates artifacts, with fewer images containing a person (92.8% vs. 99%). For SD3.5, InterMit-PM reduces intersectional biases effectively ( $\tau=0.27$ ) with 76% of the steps. Being training-free, it maintains the original model's image quality, unlike ITI-GEN.

**Ablation.** To ablate on InterMit-PM, we attempt a random baseline, Random-PM where biases are mitigated in a random order, rather than using InterMit's prioritized sequence. Using Stable Diffusion 1.4 on a subset of 12 prompts, we found that InterMit achieved comparable performance with fewer mitigation steps, as shown in Table 2.

# 6 Discussion

**Role of priority vector.** Incorporating user-defined priorities enables flexible and targeted bias mitiga-

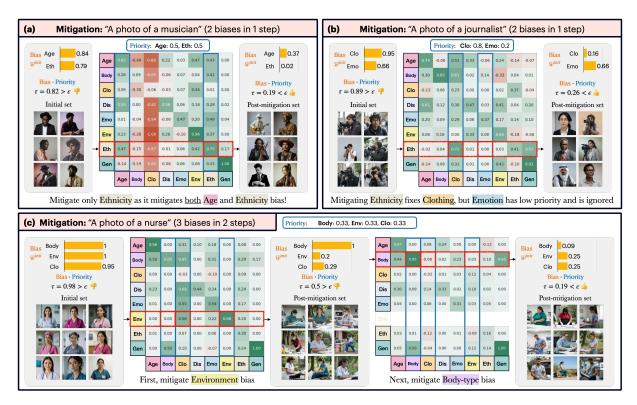


Figure 5: Three examples of mitigation using InterMit. Priority vectors guide the mitigation process. Columns that are a part of sub-matrix S' are in blue. As shown in (a) and (c), the algorithm mitigates multiple biases in fewer steps. (b) shows how user-defined priorities guide the process and when thresholds are met. Mitigating one axis, like ethnicity (Eth), can also affect others like clothing (Clo) and emotion (Emo), revealing bias interdependencies.

tion. For example, in Fig. 5(c), the user assigns equal weights to body type, environment, and clothing, prompting the model to mitigate all three biases equally. In contrast, Fig. 5(b) illustrates a case where the user prioritizes clothing diversity while assigning lower weight to emotion, focusing the mitigation effort accordingly. This flexibility makes our approach adaptable to a wide range of user goals and fairness requirements.

Accounting for different target distributions. Most fairness methods assume a fixed ideal distribution. In contrast, BiasConnect allows users to define a custom target  $D^*$  per bias axis, enabling context-sensitive mitigation. As an experiment, we collect 48 real images of computer programmers, and use this to replace  $D^*$  with  $D^{real}$  for all biases. Now, re-estimating the IS, we observe significant differences (Appendix Fig. A2). Notably, the IS for the effect of mitigating clothing on itself flips from +0.88 to -0.79 in Kandinsky, as the ideal distribution of clothing now reflects the skew towards informal in the real world, rather than a uniform distribution (see Appendix A.11).

Uncovering optimal bias mitigation strategies. InterMit is flexible and supports any set of user-

specified bias axes. As shown in Fig. 5(a) & (c), it often achieves effective mitigation in fewer steps than the user-defined threshold. By leveraging inter-axis relations, it identifies optimal strategies. In Fig. 5(a), when age and ethnicity are equally prioritized, mitigating ethnicity alone can reduce both due to demographic overlap, and a single intervention meets the threshold. In Fig. 5(c) in two mitigation steps, the bias profile progressively aligns with the priority vector (dot product  $\tau$  drops:  $0.98 \rightarrow 0.50 \rightarrow 0.19$ ). Notably, mitigating environment also reduces clothing bias due to strong intersectionality, showing how our method leverages inter-axis relationships for efficient mitigation. Moreover, if InterMit fails to reach the desired bias threshold or if mitigating one axis negatively impacts another, it can alert the user to these trade-offs (Fig. 3), enabling informed decision-making.

**Extension to Other Approaches.** We propose a general framework for mitigating intersectional biases in TTI models. As shown in Alg. 1, our method can be layered on any sequential bias mitigation strategy. At each step, one bias is mitigated, and the intersectionality matrix **S** is recomputed, enabling iterative application.

#### 7 Conclusion

We propose BiasConnect to investigate intersectional biases in TTI models. While prior research has explored bias detection and mitigation, to the best of our knowledge, no previous work has focused on understanding how biases influence one another. Unlike InterMit, prior bias mitigation strategies did not account for intersectional impacts. We believe our work enables a more nuanced analysis of bias interactions and supports informed decision-making for AI users and developers, fostering more equitable and transparent AI.

#### 7.1 Limitations

Studying bias in text-to-image models is difficult because biases shift across real-world contexts and interact in unpredictable ways. Since it is impossible to capture every variation, we restrict our analysis to a tractable subset of counterfactuals rather than testing all attributes across all bias axes. This requires automatic attribute extraction, which we perform using a vision-language model. Yet, using generative models to evaluate TTI systems introduces further challenges, as these models are imperfect and may replicate the very biases they are meant to study. To address this, we rely on a well-tested VQA model (Chen et al., 2023) from TIBET (Chinchure et al., 2024), conduct robustness analyses to measure error rates, and sample 48 images per prompt to stabilize our estimates of Intersectional Sensitivity. These measures improve reliability but cannot fully capture the nuance of real-world biases or eliminate all error sources.

Beyond these technical considerations, our findings also depend on the broader assumptions that guide the study. The choice of dataset, the definition of an "ideal" distribution, and the interpretation of bias by the user all shape the outcomes of the analysis. In our experiments, we assume a uniform distribution for simplicity, but the framework is designed to accommodate any user-specified distribution. While this flexibility improves transparency, it also raises the risk that users could set distributions or priority vectors in ways that reinforce rather than challenge existing biases. Nonetheless, we include these controls intentionally, because making assumptions explicit is better than obscuring them, and because diverse applications demand flexibility. Although the tool does not resolve every limitation, it provides a transparent and scalable framework that uncovers bias in TTI models with

reasonable accuracy and lays the groundwork for more comprehensive approaches.

#### 7.2 Ethics Statement

Our work is motivated by the ethical imperative to make bias in text-to-image (TTI) models visible, auditable, and contestable. We draw on the pluralistic tradition of intersectionality, which spans additive, interactive, and structural perspectives (Collins, 2015; Curry, 2018; Diana and Tolbert, 2023). Any computationally feasible formalization requires simplification, and our use of fixed axes reflects this necessity while acknowledging the complexity of real-world biases. Our framework does not claim to capture the full philosophical depth of intersectionality but instead aims to make it legible and actionable. Following established precedent in fairness research (Buolamwini and Gebru, 2018; Chinchure et al., 2024), we adopt categories pragmatically with the goal of making model behavior auditable, and we design the framework to prioritize transparency rather than treating these categories as fixed or absolute.

More broadly, we view formal modeling as a safeguard: it makes bias dynamics visible, contestable, and accountable. While misuse cannot be prevented outright, open frameworks enable scrutiny and honest auditing. We remain metaphysically agnostic, aiming not to define categories but to reveal how TTI systems encode and reproduce them. Our contribution lies in making intersectional biases numerically measurable and actionable, while acknowledging that our framework is not a definitive account of real-world intersectionality but one of several defensible approaches. By offering a transparent, user-controllable, and computationally grounded method, we hope to advance bias auditing and provide a foundation that future work can extend, challenge, and refine.

# 8 Acknowledgment

This work was funded, in part, by the Vector Institute for AI, Canada CIFAR AI Chair, NSERC CRC, and NSERC DG. Hardware resources used in preparing this research were provided, in part, by the Province of Ontario, the Government of Canada through CIFAR, and companies sponsoring the Vector Institute.

#### References

- Lavisha Aggarwal and Shruti Bhargava. 2023. Fairness in ai systems: Mitigating gender bias from language-vision models. *arXiv preprint arXiv:2305.01888*.
- Jaimeen Ahn and Alice Oh. 2021. Mitigating language-dependent ethnic bias in bert. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 533–549.
- Xuechunzi Bai, Angelina Wang, Ilia Sucholutsky, and Thomas L Griffiths. 2025. Explicitly unbiased large language models still form biased associations. *Proceedings of the National Academy of Sciences*, 122(8):e2416228122.
- Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Margaret Shmitchell. 2021. On the dangers of stochastic parrots: Can language models be too big? *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 610–623.
- Federico Bianchi, Pratyusha Kalluri, Esin Durmus, Faisal Ladhak, Myra Cheng, Debora Nozza, Tatsunori Hashimoto, Dan Jurafsky, James Zou, and Aylin Caliskan. 2023. Easily accessible text-to-image generation amplifies demographic stereotypes at large scale. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, pages 1493–1504.
- Abeba Birhane and Vinay Uday Prabhu. 2021. Multimodal datasets: Misogyny, pornography, and malignant stereotypes. *arXiv preprint arXiv:2110.01963*.
- BlackForestLabs. 2024. Flux. https://github.com/black-forest-labs/flux.
- Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. 2016. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. *Advances in neural information processing systems*, 29.
- Joy Buolamwini and Timnit Gebru. 2018. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on fairness, accountability and transparency*, pages 77–91. PMLR.
- Jun Chen, Deyao Zhu, Xiaoqian Shen, Xiang Li, Zechun Liu, Pengchuan Zhang, Raghuraman Krishnamoorthi, Vikas Chandra, Yunyang Xiong, and Mohamed Elhoseiny. 2023. Minigpt-v2: large language model as a unified interface for vision-language multi-task learning. arXiv preprint arXiv:2310.09478.
- Aditya Chinchure, Pushkar Shukla, Gaurav Bhatt, Kiri Salij, Kartik Hosanagar, Leonid Sigal, and Matthew Turk. 2024. Tibet: Identifying and evaluating biases in text-to-image generative models. In *European Conference on Computer Vision*, pages 429–446. Springer.

- Jaemin Cho, Abhay Zala, and Mohit Bansal. 2023. Dalleval: Probing the reasoning skills and social biases of text-to-image generation models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3043–3054.
- Patricia Hill Collins. 2015. Intersectionality's definitional dilemmas. *Annual review of sociology*, 41(1):1–20.
- Kimberle Crenshaw. 1989. Demarginalizing the intersection of race and sex: A black feminist critique of antidiscrimination doctrine, feminist theory, and antiracist politics. In *University of Chicago Legal Forum*, pages 139–167.
- Tommy J Curry. 2018. Killing boogeymen: Phallicism and the misandric mischaracterizations of black males in theory. *Res Philosophica*.
- Hannah Devinney, Jenny Björklund, and Henrik Björklund. 2024. We don't talk about that: case studies on intersectional analysis of social bias in large language models. In *Workshop on Gender Bias in Natural Language Processing (GeBNLP), Bangkok, Thailand, 16th August, 2024.*, pages 33–44. Association for Computational Linguistics.
- Emily Diana, Wesley Gill, Michael Kearns, Krishnaram Kenthapadi, and Aaron Roth. 2021. Minimax group fairness: Algorithms and experiments. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, pages 66–76.
- Emily Diana and Alexander Williams Tolbert. 2023. Correcting underrepresentation and intersectional bias for classification. arXiv preprint arXiv:2306.11112.
- Moreno D'Incà, Elia Peruzzo, Massimiliano Mancini, Dejia Xu, Vidit Goel, Xingqian Xu, Zhangyang Wang, Humphrey Shi, and Nicu Sebe. 2024. Openbias: Open-set bias detection in text-to-image generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12225–12235.
- Piero Esposito, Parmida Atighehchian, Anastasis Germanidis, and Deepti Ghadiyaram. 2023. Mitigating stereotypical biases in text to image generative systems. *arXiv preprint arXiv:2310.06904*.
- James R Foulds, Rashidul Islam, Kamrun Naher Keya, and Shimei Pan. 2020. An intersectional definition of fairness. In 2020 IEEE 36th international conference on data engineering (ICDE), pages 1918–1921. IEEE.
- Felix Friedrich, Patrick Schramowski, Manuel Brack, Lukas Struppek, Dominik Hintersdorf, Sasha Luccioni, and Kristian Kersting. 2023. Fair diffusion: Instructing text-to-image generation models on fairness. *arXiv preprint arXiv:2302.10893*.

- Rohit Gandikota, Hadas Orgad, Yonatan Belinkov, Joanna Materzyńska, and David Bau. 2024. Unified concept editing in diffusion models. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 5111–5120.
- Ismael Garrido-Muñoz, Arturo Montejo-Ráez, Fernando Martínez-Santiago, and L Alfonso Ureña-López. 2021. A survey on bias in deep nlp. *Applied Sciences*, 11(7):3184.
- Avijit Ghosh, Lea Genuit, and Mary Reagan. 2021. Characterizing intersectional group fairness with worst-case comparisons. In *Artificial Intelligence Diversity, Belonging, Equity, and Inclusion*, pages 22–34. PMLR.
- Sourojit Ghosh and Aylin Caliskan. 2023. 'person' == light-skinned, western man, and sexualization of women of color: Stereotypes in stable diffusion. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 6971–6985.
- Wei Guo and Aylin Caliskan. 2021. Detecting emergent intersectional biases: Contextualized word embeddings contain a distribution of human-like biases. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, pages 122–133.
- Kimia Hamidieh, Haoran Zhang, Thomas Hartvigsen, and Marzyeh Ghassemi. 2023. Identifying implicit social biases in vision-language models.
- Lisa Anne Hendricks, Kaylee Burns, Kate Saenko, Trevor Darrell, and Anna Rohrbach. 2018. Women also snowboard: Overcoming bias in captioning models. In *Proceedings of the European conference on computer vision (ECCV)*, pages 771–787.
- Yusuke Hirota, Yuta Nakashima, and Noa Garcia. 2022. Gender and racial bias in visual question answering datasets. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, pages 1280–1292.
- Elizabeth Hoepfinger. 2023. Racial and intersectional debiasing of contrastive language image pretraining. Master's thesis, University of Georgia.
- Phillip Howard, Avinash Madasu, Tiep Le, Gustavo Lujan Moreno, Anahita Bhiwandiwalla, and Vasudev Lal. 2024. Socialcounterfactuals: Probing and mitigating intersectional social biases in vision-language models with counterfactual examples. In *Proceed*ings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 11975–11985.
- Ben Hutchinson, Vinodkumar Prabhakaran, Emily Denton, Kellie Webster, Yu Zhong, and Stephen Denuyl. 2020. Social biases in nlp models as barriers for persons with disabilities. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5491–5501.

- Aparna R Joshi, Xavier Suau Cuadros, Nivedha Sivakumar, Luca Zappella, and Nicholas Apostoloff. 2022. Fair sa: Sensitivity analysis for fairness in face recognition. In *Algorithmic fairness through the lens of causality and robustness workshop*, pages 40–58. PMLR.
- Loukas Kavouras, Konstantinos Tsopelas, Giorgos Giannopoulos, Dimitris Sacharidis, Eleni Psaroudaki, Nikolaos Theologitis, Dimitrios Rontogiannis, Dimitris Fotakis, and Ioannis Emiris. 2023. Fairness aware counterfactuals for subgroups. Advances in Neural Information Processing Systems, 36:58246–58276.
- Michael Kearns, Seth Neel, Aaron Roth, and Zhiwei Steven Wu. 2018. Preventing fairness gerrymandering: Auditing and learning for subgroup fairness. In *International conference on machine learning*, pages 2564–2572. PMLR.
- Hannah Kirk, Yennie Jun, Haider Iqbal, Elias Benussi, Filippo Volpin, Frederic A. Dreyer, Aleksandar Shtedritski, and Yuki M. Asano. 2021. Bias Out-of-the-Box: An Empirical Analysis of Intersectional Occupational Biases in Popular Generative Language Models. arXiv preprint. ArXiv:2102.04130 [cs].
- John P Lalor, Yi Yang, Kendall Smith, Nicole Forsgren, and Ahmed Abbasi. 2022. Benchmarking intersectional biases in nlp. In *Proceedings of the 2022 conference of the North American chapter of the association for computational linguistics: Human language technologies*, pages 3598–3609.
- Ida Marie S Lassen, Mina Almasi, Kenneth Enevoldsen, and Ross Deans Kristensen-McLachlan. 2023. Detecting intersectionality in ner models: A data-driven approach. In *Proceedings of the 7th joint SIGHUM workshop on computational linguistics for cultural heritage, social sciences, humanities and literature*, pages 116–127.
- Daiqing Li, Aleks Kamko, Ehsan Akhgari, Ali Sabet, Linmiao Xu, and Suhail Doshi. 2024. Playground v2.5: Three insights towards enhancing aesthetic quality in text-to-image generation. *Preprint*, arXiv:2402.17245.
- Bingyu Liu, Weihong Deng, Yaoyao Zhong, Mei Wang, Jiani Hu, Xunqiang Tao, and Yaohai Huang. 2019. Fair loss: Margin-aware reinforcement learning for deep face recognition. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10052–10061.
- Weicheng Ma, Brian Chiang, Tong Wu, Lili Wang, and Soroush Vosoughi. 2023. Intersectional stereotypes in large language models: Dataset and analysis. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 8589–8597.
- Nicole Meister, Dora Zhao, Angelina Wang, Vikram V Ramaswamy, Ruth Fong, and Olga Russakovsky. 2023. Gender artifacts in visual datasets. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4837–4848.

- Sungho Park, Sunhee Hwang, Jongkwang Hong, and Hyeran Byun. 2020. Fair-vqa: Fairness-aware visual question answering through sensitive attribute prediction. *IEEE Access*, 8:215091–215099.
- Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. 2021. Zero-shot text-to-image generation. *arXiv preprint arXiv:2102.12092*.
- Anton Razzhigaev, Arseniy Shakhmatov, Anastasia Maltseva, Vladimir Arkhipkin, Igor Pavlov, Ilya Ryabov, Angelina Kuts, Alexander Panchenko, Andrey Kuznetsov, and Denis Dimitrov. 2023. Kandinsky: An improved text-to-image synthesis with image prior and latent diffusion. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 286–295.
- 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.
- Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Raphael Gontijo Lopes, and 1 others. 2022. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*.
- Preethi Seshadri, Sameer Singh, and Yanai Elazar. 2023. The bias amplification paradox in text-to-image generation. *arXiv preprint arXiv:2308.00755*.
- Laleh Seyyed-Kalantari, Haoran Zhang, Matthew BA McDermott, Irene Y Chen, and Marzyeh Ghassemi. 2021. Underdiagnosis bias of artificial intelligence algorithms applied to chest radiographs in under-served patient populations. *Nature medicine*, 27(12):2176–2182.
- Deven Santosh Shah, H Andrew Schwartz, and Dirk Hovy. 2020. Predictive biases in natural language processing models: A conceptual framework and overview. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5248–5264.
- Arseniy Shakhmatov, Anton Razzhigaev, Aleksandr Nikolich, Vladimir Arkhipkin, Igor Pavlov, Andrey Kuznetsov, and Denis Dimitrov. 2023. kandinsky 2.2.
- Ryan Steed and Aylin Caliskan. 2021. Image representations learned with unsupervised pretraining contain human-like biases. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 701–713.
- Yi Chern Tan and L Elisa Celis. 2019. Assessing social and intersectional biases in contextualized word

- representations. Advances in neural information processing systems, 32.
- Angelina Wang, Alexander Liu, Ryan Zhang, Anat Kleiman, Leslie Kim, Dora Zhao, Iroha Shirai, Arvind Narayanan, and Olga Russakovsky. 2022a. Revise: A tool for measuring and mitigating bias in visual datasets. *International Journal of Computer Vision*, 130(7):1790–1810.
- Angelina Wang, Vikram V Ramaswamy, and Olga Russakovsky. 2022b. Towards Intersectionality in Machine Learning: Including More Identities, Handling Underrepresentation, and Performing Evaluation. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '22, pages 336–349, New York, NY, USA. Association for Computing Machinery.
- Angelina Wang and Olga Russakovsky. 2023. Overwriting pretrained bias with finetuning data. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3957–3968.
- Jialu Wang, Xinyue Gabby Liu, Zonglin Di, Yang Liu, and Xin Wang. 2023a. T2iat: Measuring valence and stereotypical biases in text-to-image generation. In *Findings of the Association for Computational Linguistics: ACL 2023*. Association for Computational Linguistics.
- Jialu Wang, Yang Liu, and Xin Wang. 2022c. Assessing multilingual fairness in pre-trained multimodal representations. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2681–2695.
- Jianyi Wang, Kelvin CK Chan, and Chen Change Loy. 2023b. Exploring clip for assessing the look and feel of images. In *AAAI*.
- Zeyu Wang, Klint Qinami, Ioannis Christos Karakozis, Kyle Genova, Prem Nair, Kenji Hata, and Olga Russakovsky. 2020. Towards fairness in visual recognition: Effective strategies for bias mitigation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8919–8928.
- Cheng Zhang, Xuanbai Chen, Siqi Chai, Chen Henry Wu, Dmitry Lagun, Thabo Beeler, and Fernando De la Torre. 2023. Iti-gen: Inclusive text-to-image generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3969–3980.
- Eric Zhao, De-An Huang, Hao Liu, Zhiding Yu, Anqi Liu, Olga Russakovsky, and Anima Anandkumar. 2021. Scaling fair learning to hundreds of intersectional groups.

Axis	Counterfactuals
gender bias	A photo of a male <i>occupation</i> A photo of a female <i>occupation</i>
age bias	A photo of an old occupation A photo of a middle-aged occupation A photo of a young occupation
ethnicity bias	A photo of a white occupation A photo of a african american occupation A photo of an asian occupation A photo of a south asian occupation A photo of a middle eastern occupation A photo of a hispanic occupation
bodytype bias	A photo of a obese occupation A photo of a slim occupation A photo of a muscular occupation
environment bias	A photo of a <i>occupation</i> working indoors A photo of a <i>occupation</i> working outdoors
clothing bias	A photo of a <i>occupation</i> in formal attire A photo of a <i>occupation</i> in informal attire
emotion bias	A photo of a <i>occupation</i> who is happy A photo of a sad <i>occupation</i> who is sad A photo of a <i>occupation</i> who is serious A photo of a <i>occupation</i> who is tired
disability bias	A photo of a <i>occupation</i> who is fit A photo of a blind <i>occupation</i> A photo of a <i>occupation</i> with a hearing aid A photo of a <i>occupation</i> on a wheelchair

Table 3: Occupation Dataset.

# A Appendix

#### A.1 Occupation Dataset

Table 3 shows the occupation prompt templates along each of the bias axes. We consider 26 occupations: computer programmer, elementary school teacher, librarian, announcer, pharmacist, chef, chemist, police, accountant, architect, lawyer, philosopher, scientist, doctor, nurse, engineer, musician, journalist, athlete, social worker, sales person, politician, farmer, mechanic, firefighter, gardener.

#### A.2 Models

We test BiasConnect on several different TTI models. Unless mentioned otherwise, all the settings for image generation are the default settings provided for these models on Huggingface. We use a single Nvidia A40 GPU for all our experiments. We do not measure exact GPU hours as we do not do not train any large scale models, and our usage is dependent on the number of biases we are analyzing or mitigating for each prompt.

For prompts in the Occupation dataset, we use the following models:

• Stable Diffusion 1.4: We pick this model

- because of its compatibility with most bias mitigation tools, like ITI-GEN. We use *CompVis/stable-diffusion-v1-4*.
- Stable Diffusion 3.5: This is the latest model in the Stable Diffusion family. We specifically use the *stabilityai/stable-diffusion-3.5-large-turbo* model.
- Flux-dev: The Flux series of models is another popular set of open-source TTI models. We use the *black-forest-labs/FLUX.1-dev* variant of this family, and conduct inference with guidance scale = 3.5 and num inference steps = 30.
- Playground 2.5: This model is trained to produce aesthetically pleasing images. We use *playgroundai/playground-v2.5-1024px-aesthetic*, with guidance scale = 3 and num inference steps = 50.
- Kandinsky 2.2: We use the *kandinsky-community/kandinsky-2-2-decoder* model, and use the default negative prompt = "low quality, bad quality".

For the TIBET dataset, we use the images already provided in the dataset as is. These images were generated using Stable Diffusion 2.1.

# A.3 VQA

For a given set of images and a set of axes B, our goal is to find distributions  $D_{B_i}$  for all bias axes  $B_i \in B$ . In order to find this distribution, we must begin by first identifying the attributes related to  $B_i$  in every image of the image set. We use VQA for this process.

For every image in the set, we first start by asking the VQA the question Is there a person in the image (yes or no)? for the Occupation prompts dataset. This allows us to filter images where we will be unable to extract bias-related attributes due to low quality generation. For the images that have a person, we have the following set of questions to extract all bias-related attributes:

- gender bias: What is the gender (male, female) of the person?
- age bias: What is the age group (young, middle, old) of the person?
- ethnicity bias: What is the ethnicity (white, black, asian, south asian, middle eastern, hispanic) of the person?

- bodytype bias: What is the body type (fat, slim, muscular) of the person?
- environment bias: What is the environment (indoor, outdoor) of the person?
- clothing bias: What is the attire (formal, informal) of the person?
- emotion bias: What is the emotion (happy, sad, serious, tired) of the person?
- disability bias: Is this person blind (yes or no)?; Is this person wearing a hearing aid (yes or no)?; Is this person on a wheelchair (yes or no)?

Note that all questions are multiple choice. Furthermore, for disability bias, we split the question into three parts, and run each part through the VQA model independently. If none of the parts are answered as 'yes', then the person in the image is 'fit' and does not have one of those disabilities.

In terms of error rate for robustness, we believe that our MCQ-based VQA approach would yield a lower than 18% error rate observed in TIBET (Chinchure et al., 2024), which uses the same VQA model. Empirically speaking, we observe that our VQA performs near-perfectly on axes such as gender, environment and emotion, but may sometimes return incorrect guesses among other axes in more ambiguous scenarios. As VQA models improve, our method can utilize them in a plug-and-play manner.

#### A.4 TIBET Data

TIBET dataset contains 100 prompts, their biases and relevant counterfactuals, and 48 images for each initial and counterfactual prompt. Because of the dynamic nature of these biases (they vary from prompt to prompt), we use the VQA strategy in the TIBET method instead of our templated questions from above to extract concepts.

#### A.5 Bias Mitigation Study

We conduct a study using ITI-GEN to measure how often a bias mitigation might yield negative effects on other bias axes. We define a negative Intersectional Sensitivity score  $(IS_{xy} < 0)$  to suggest that mitigating bias axis  $B_x$  reduces the diversity of attributes of axis  $B_y$ .

In this study, for all 26 occupations and across all bias axes listed in Table 3, we mitigate every bias axis independently. We then compute IS,

where the initial distribution  $D_{B_y}^{B_x}$  in equation 3 is replaced by  $D_{B_y}^{mit(B_x)}$ , which is based on the VQA extracted attributes for bias axis  $B_y$  in the newly generated set of images post-mitigation of axis  $B_x$  with ITI-GEN. This score is defined as:

$$w_{B_y}^{B_x} = W_1(D_{B_y}^{mit(B_x)}, D^*)$$
 (5)

$$IS_{xy}^{mit(x)} = \overline{w}_{B_y}^{init} - \overline{w}_{B_y}^{B_x}$$
 (6)

We compute the percentage of  $IS_{xy}^{mit(x)}$  for all possible pairs of biases,  $B_x$  and  $B_y$ , where mitigation of  $B_x$  led to  $IS_{xy}^{mit(x)} < 0$ . We find that a substantial number of times, 29.4% of all mitigations, led to a negative effect.

#### A.6 Additional prompt-level examples

We show additional examples of prompt-level intersectional analysis in Fig A3 below. For TIBET, Fig A3(c) shows how diversifying on an axis like Geography can help diversify the Ethnicity distribution.

# A.7 Prompt-Modification Based Mitigation (PM)

We use the simplest possible mitigation method: using prompt modification. We choose to use a prompt-modification based method over other methods like ITI-GEN because it gives us the capability to mitigate biases sequentially, store intermediate results, and evaluate the effectiveness of our mitigation algorithm InterMit. Moreover, it is training-free, and can leverage compute optimizations like quantization and Flash Attention to reduce computational costs.

Let us assume we want to mitigate environment bias, and then clothing bias for "nurse". We will assume that our ideal distribution is the uniform distribution across all counterfactuals of each of these axes. The prompt modification process (*PM*) works as follows:

• Environment bias has two counterfactuals, 'indoor' and 'outdoor'. If our total set of images is 48, mitigating it would mean generating 50% images indoor, and 50% outdoor. Therefore, during mitigation, our initial prompt A photo of a nurse is replaced by a combination of two initial prompts, A photo of a nurse working indoors, A photo of a nurse working outdoors. This is our mitigated model. At this stage, any counterfactual prompt (for

BiasConnect) will also account for this mitigated prompt set. So the gender counterfactuals, at this step, will be [A photo of a male nurse working indoors, A photo of a male nurse working outdoors], [A photo of a female nurse working indoors, A photo of a female nurse working outdoors].

- Next, we want to mitigate clothing bias. Clothing, again, has two counterfactuals: 'formal' and 'informal'. Our new initial prompt set will be A photo of a nurse working indoors dressed formally, A photo of a nurse working outdoors dressed formally, A photo of a nurse working indoors dressed informally, A photo of a nurse working outdoors dressed informally.
- All future mitigation steps will add to these permutations, and an equal number of images are generated for each prompt in the set, such that the total is 48 (or more) images.

#### **A.8** Mitigation Prompts

Table 4 contains all the occupation prompts, and their mitigation results, for the SD1.4 model using InterMit-PM and ITI-GEN. Table 5 has all prompts and mitigation results for the SD3.5 experiments. This table also shows cases where we did mitigation based on a priority vector.

#### A.9 Alternate Mitigation Method: UCE

Alongside ITI-GEN, we also considered concept editing methods—primarily Unified Concept Editing (UCE) (Gandikota et al., 2023)—as potential baselines for bias mitigation. We implemented UCE using the publicly available GitHub code for Stable Diffusion 1.4 and applied it iteratively to debias images across various axes. However, we encountered major challenges in both effectiveness and interpretability. While UCE initially maintained image plausibility (e.g., the "isP?" VQA metric—Is this a person?—remained above 90% for the first 6 iterations), this rapidly deteriorated. After around 15 iterations, the outputs devolved into noisy and uninterpretable patterns, with "isP?" scores dropping to 0%, despite not yet reaching bias mitigation convergence. This failure was especially pronounced for occupations such as "nurse" and along axes like age and ethnicity. The optimization seemed to push toward abstract visual patterns that might improve numerical bias metrics but yielded informationless and unusable images.

We show examples of UCE-based debiasing in Figure A1. Based on these findings, we determined that UCE is not a viable baseline for effective and interpretable model debiasing.

#### A.10 Validating Mitigation Effect Estimation

Our approach provides empirical estimates of how a counterfactual-based mitigation strategy may influence an intersectional relationship  $B_x \to B_y$  in the form of the Intersectional Sensitivity score. To validate these estimates, we conduct an experiment where we actually perform mitigation on SD 1.4 using ITI-GEN and SD3.5 using PM. For all 26 occupations, we consider all intersectional relationships  $B_x \to B_y$ , and mitigate all  $B_x$  independently. To compute the new Intersectional Sensitivity post mitigation, we replace the initial distribution  $D_{B_y}^{B_x}$  in equation 3 with  $D_{B_y}^{mit(B_x)}$ , which is based on the VQA extracted attributes for bias axis  $B_y$  in the newly generated set of images post-mitigation of axis  $B_x$  with ITI-GEN. This new score can be defined as:

$$w_{B_y}^{B_x} = W_1(D_{B_y}^{mit(B_x)}, D^*)$$
 (7)

$$IS_{xy}^{mit(x)} = \overline{w}_{B_y}^{init} - \overline{w}_{B_y}^{B_x}$$
 (8)

Note that these equations are the same as the ones we used in Section A.5. To quantify the effectiveness of BiasConnect we measure the average correlation between the Intersectional Sensitivity scores before  $IS_{xy}$  and after mitigation  $IS_{xy}^{mit(x)}$  across all intersectional relationships  $B_x \to B_y$  present for each prompt.

The high correlations (0.65 for ITI-GEN, 0.95 for *PM*) suggest that our method effectively estimates the potential impacts of bias interventions without actually doing the mitigation step itself, which can be computationally expensive.

Such empirical guarantees provide users with valuable insights into whether altering bias along a particular dimension will lead to meaningful improvements in fairness across other bias dimensions. By estimating how counterfactual-based interventions influence overall bias scores, our approach helps researchers and practitioners predict the effectiveness of mitigation techniques before full deployment.

#### A.11 Using Real World Biases

Other than understanding biases in TTI models BiasConnect can be used to compare bias de-



Figure A1: We could not use UCE (Gandikota et al., 2024) as one of our baselines, as it is not effective at bias mitigation across some axes, like ethnicity.

pendencies in images generated by Text-to-Image (TTI) models with a reference real-world image distribution. Instead of assuming a uniform distribution as the baseline for bias sensitivity calculations, we consider the empirical distribution of the reference dataset as the initial distribution.

Given a prompt P (e.g., "A computer programmer"), let  $B = [B_1, B_2, ..., B_n]$  represent the set of bias axes (e.g., gender, age, race). For each bias axis  $B_y$ , we define:

- $D_{B_y}^{\text{real}}$ : real-world distribution of  $B_y$  (from a dataset or observed statistics).
- $D_{B_y}^{\text{init}}$ : distribution of  $B_y$  in TTI-generated images. This is the same as in BiasConnect.

The Wasserstein-1 distance between real-world and TTI-generated distributions quantifies how far the TTI bias distribution is from real-world data is:

$$w_{B_y}^{\text{init}} = W_1(D_{B_y}^{\text{init}}, D_{B_y}^{\text{real}}) \tag{9}$$

To measure the impact of intervening on  $B_x$ , we compute the post-intervention Wasserstein distance:

$$w_{B_y}^{B_x} = W_1(D_{B_y}^{B_x}, D_{B_y}^{\text{real}})$$
 (10)

The Intersectional Sensitivity Score  $IS_{xy}$  for the effect of changing  $B_x$  on  $B_y$  measures the difference between  $w_{B_y}^{\rm init}$  and  $w_{B_y}^{B_x}$  similar to the one calculated in Eq 4.

In our experiment, we obtain a real-world distribution of all biases for a computer programmer by sampling 48 images from Google Images. We measure the real world distributions using VQA. Now, by using BiasConnect with the real-world distributions of each bias set to be the ideal distribution  $D^*$ , we recompute IS as described above. Fig. A2 show how significant the effect of this is. Notably, in both SD3.5 and Kandinsky images, we observe that mitigating clothing (to diversify clothing) actually has a negative IS value, as this would move

away from our ideal (real-world) distribution of computer programmers mostly wearing informal clothes.

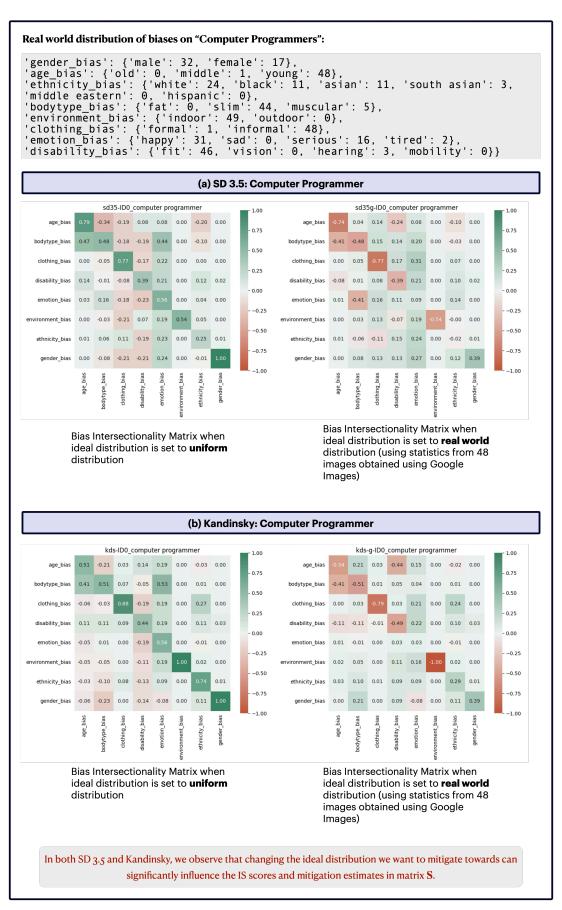


Figure A2: Modifying  $D^*$  (ideal distribution) in BiasConnect can have a significant effect on the Intersectional Sensitivity values.

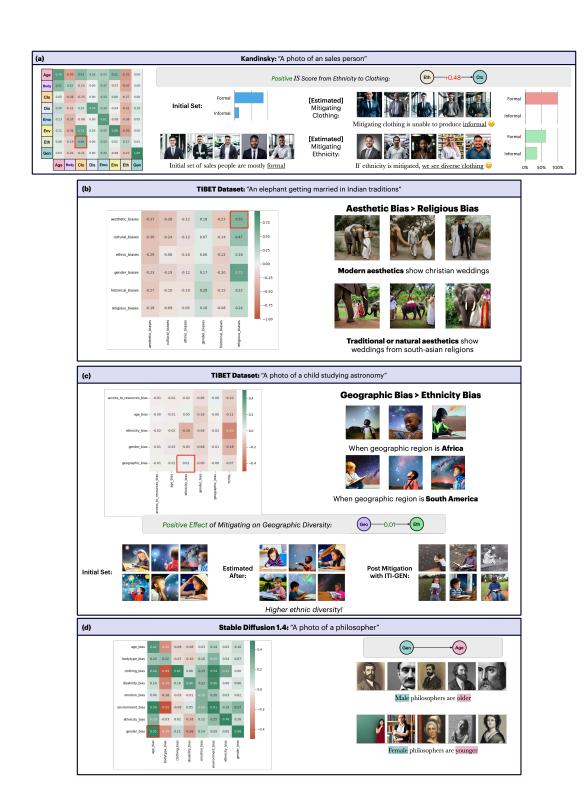


Figure A3: Additional examples on TIBET (b,c) and Occupation prompt (a,d) on prompt-level analysis provided by BiasConnect.

Occupation	Mitigation Priority Vector	Bias Mitigated	# of Biases	# Steps	MitAmt	ITI-GEN MitAmt
computer programmer	{'age_bias': 0.33, 'emotion_bias': 0.33, 'emotion_bias': 0.33	['emotion_bias', 'ethnicity_bias',	3	3	0.46	0.64
elementary school teacher	(bodytype_bias): 0.33, ethnic-ity bias): 0.33, oethnic-ity bias): 0.33	age_bias ] gender_bias', 'bodytype_bias']	3	2	0.28	0.51
librarian	Adjustation of the control of the co	['disability_bias', 'age_bias', 'ethnic-ity_bias']	4	8	0.32	09.0
announcer	{'age_bias': 0.5, 'environment_bias': 0.5.	['age_bias', 'environment_bias']	2	2	0.04	0.18
pharmacist	('bodytype_bias': 0.25, 'disability-bias': 0.25, 'emotion_bias': 0.25, 'achaicity bias': 0.25,	['age_bias', 'bodytype_bias', 'disabil-ity_bias', 'ethnicity_bias']	4	4	0.40	0.77
chef	{'disability_bias': 0.5, 'emotion_bias': 0.5}	['emotion_bias', 'disability_bias']	2	2	69.0	0.82
chemist	{age_bias': 0.33, 'clothing_bias': 0.33, 'mortion bias': 0.33}	['emotion_bias', 'bodytype_bias']	3	2	0.33	0.40
police	(2.32, cmoton_bas : 0.33) { age_bias': 0.33, 'emotion_bias': 0.33 'athnicity bias': 0.33}	['age_bias', 'ethnicity_bias', 'emo-	3	3	0.39	0.58
accountant	(**ge_bias': 0.25, 'clothing_bias': 0.25, 'environment_bias': 0.25, 'environment_bias': 0.25, 'eth-nicity' bias': 0.25, 'eth-nicity' bias': 0.25	['ethnicity_bias', 'environment_bias']	4	2	0.28	0.28
architect	(bodytype_bias: 0.2, emotion_bias: 0.2, environment_bias: 0.2, ethnicity bias: 0.2, environment_bias: 0.2, ethnicity bias: 0.2, ender bias: 0.2)	['emotion_bias', 'gender_bias', 'body-type_bias']	S	8	0.34	0.51
lawyer	{'emotion_bias': 0.5, 'gender_bias': 0.5}	['emotion_bias']	2	-	0.33	0.33
philosopher	{clothing_bias': 0.5, 'disability_bias': 0.5}	['clothing_bias']	2	1	0.33	0.56
scientist doctor	(age_bias): 0.5, 'emotion_bias): 0.5, 'age_bias': 0.2, 'bodytype_bias': 0.2, 'clothing_bias': 0.2, 'ethnicity_bias': 0.7, 'ender bias': 0.3,	['bodytype_bias', 'disability_bias'] ['gender_bias', 'bodytype_bias']	2 %	7 7	0.38	0.49
nurse	('bodytype_bias': 0.33, 'cloth-ing_bias': 0.33, 'environment_bias': 0.33,	['gender_bias', 'environment_bias']	ю	7	0.31	0.42
engineer musician	(*age_bias': 0.5, 'disability_bias': 0.5, 'age_bias': 0.2, 'disability_bias': 0.2, 'emotion bias': 0.2, 'ethnicity_bias': 0.2, 'scander bias': 0.3, 'scander	['bodytype_bias', 'disability_bias'] ['age_bias', 'gender_bias', 'ethnic-ity_bias', 'disability_bias']	270	0.4	0.53	0.61
journalist	{'emotion_bias': 0.5, 'ethnicity_bias': 0.5}	['emotion_bias', 'ethnicity_bias']	2	2	0.33	0.39
athlete	(**ag_bias': 0.25, 'bodytype_bias': 0.25, 'disability_bias': 0.25, 'gender hias': 0.25	['gender_bias', 'age_bias', 'disabil-ity_bias']	4	3	0.30	0.56
social worker	{bodytype_bias}: 0.33, 'ethnic-ity bias': 0.33 'oender bias': 0.33	['emotion_bias', 'bodytype_bias']	3	2	0.24	0.18
sales person politician	(age_bias): 0.5, 'ethnicity_bias': 0.5} {     (age_bias): 0.3, 'ethnicity_bias': 0.5} {     (age_bias): 0.33, 'bodytype_bias': 0.33, 'ethnicity_bias): 0.33, 'ethnicity_bias': 0.33	['age_bias', 'ethnicity_bias'] ['ethnicity_bias']	2 6	1 2	0.22	0.48
farmer mechanic	{'age_bias': 0.5, 'ethnicity_bias': 0.5} {'bodytype_bias': 0.2, 'clothing_bias': 0.2, 'environ-bias': 0.2, 'environ-bias': 0.2, 'serviron-bias': 0.3, 'ser	['gender_bias', 'ethnicity_bias'] ['bodytype_bias', 'emotion_bias', 'clothing_bias', 'age_bias', 'ethnic-ity, bias']	2.2	2.2	0.27	0.46
firefighter	fage_bias: 0.25, 'clothing_bias': 0.25, 'environment_bias': 0.25, 'gen-der bias': 0.25, 'gen-	['age_bias', 'clothing_bias', 'gen-der_bias']	4	3	0.30	0.67
gardener	{age_bias: 0.33, 'environment_bias': 0.33, 'ethnicity_bias: 0.33}	['environment_bias', 'ethnicity_bias']	8	2	0.33	0.63

Table 4: All 26 occupational prompts and their associated biases that were mitigated (randomly selected), on Stable Diffusion 1.4. We show the priority vector, along with the actual list of biases that were mitigated, and the corresponding number of steps. Finally, we also mention MitAmt using our approach and ITI-GEN. The aggregate results of this table are in the main paper, Table 1.

Occupation	Occupation Mitigation Priority Vector	Bias Mitigated	# of Biases # Steps MitAmt	# Steps	MitAmt
announcer	clothing_bias,0.5;gender_bias,0.5	['gender_bias']	2	-	0.13
politician	clothing_bias,0.5;environment_bias,0.5	['ethnicity_bias', 'environment_bias']	2	7	0.40
musician	age_bias,0.33;ethnicity_bias,0.33	['ethnicity_bias']	2	_	0.20
mechanic	age_bias,0.33;bodytype_bias,0.33;clothing_bias,0.33	['age_bias', 'bodytype_bias', 'clothing_bias', 'disability_bias']	33	4	0.28
nurse	clothing_bias,0.33;environment_bias,0.33;bodytype_bias,0.33	['bodytype_bias', 'environment_bias']	33	7	0.19
gardener	age_bias,0.33;bodytype_bias,0.33;gender_bias,0.33	['gender_bias']	3	1	0.28
sales	clothing_bias,0.33;environment_bias,0.33;gender_bias,0.33	['environment_bias', 'gender_bias']	3	7	0.22
journalist	gender_bias,0.25;clothing_bias,0.25;emotion_bias,0.25;ethnicity_bias,0.25	['ethnicity_bias']	4	1	0.35
engineer	bodytype_bias,0.33;gender_bias,0.33;age_bias,0.33;clothing_bias,0.33	['age_bias', 'clothing_bias', 'gender_bias']	4	ж	0.28
computer	age_bias,0.33;bodytype_bias,0.33;clothing_bias,0.33;disability_bias,0.33	['disability_bias', 'bodytype_bias', 'clothing_bias']	4	ю	0.34
athlete	age_bias,0.33;bodytype_bias,0.33;clothing_bias,0.33;disability_bias,0.33;gender_bias,0.33	['disability_bias', 'bodytype_bias', 'age_bias', 'gender_bias', 'clothing_bias']	5	5	0.32
doctor	age_bias, 0.33; ethnicity_bias, 0.33; clothing_bias, 0.33; emotion_bias, 0.33; gender_bias, 0.33	['ethnicity_bias', 'age_bias', 'disability_bias', 'gender_bias', 'bodytype_bias']	5	5	0.29
teacher	bodytype_bias,0.33;emotion_bias,0.33;gender_bias,0.33;age_bias,0.33;clothing_bias,0.33	['environment_bias', 'clothing_bias', 'bodytype_bias']	5	ю	0.28
chef	bodytype_bias,0.33;emotion_bias,0.33;gender_bias,0.33;age_bias,0.33;clothing_bias,0.33	['environment_bias', 'age_bias', 'bodytype_bias', 'gender_bias', 'emotion_bias']	S	5	0.28
librarian	environment_bias,0.3;clothing_bias,0.7	['environment_bias', 'clothing_bias']	2	2	0.15
annonncer	emotion_bias,0.6;ethnicity_bias,0.4	['ethnicity_bias', 'emotion_bias']	2	2	0.50
journalist	clothing_bias,0.8;emotion_bias,0.2	['ethnicity_bias']	2	1	0.26
accountant	clothing_bias,0.20;environment_bias,0.15;gender_bias: 0.65	['clothing_bias', 'ethnicity_bias', 'emotion_bias']	3	ю	0.25
sales person	age_bias,0.5;gender_bias,0.3;bodytype_bias,0.2	['gender_bias', 'age_bias']	С	7	0.29

Table 5: Selected occupational prompts and their associated biases that were mitigated (manually selected), on Stable Diffusion 3.5. We show the priority vector, along with the actual list of biases that were mitigated, and the corresponding number of steps. Finally, we also mention MitAmt using our algorithm. The aggregate results of this table are in the main paper, Table 1. The bottom section of the table has examples where the priority vector was weighted.