

BiasConnect: Investigating Bias Interactions in Text-to-Image Models

Pushkar Shukla¹ Aditya Chinchure² Emily Diana³ Alexander Tolbert⁴
 Kartik Hosanagar⁵ Vineeth N. Balasubramanian⁶ Leonid Sigal² Matthew A. Turk¹
¹Toyota Technological Institute at Chicago ²University of British Columbia
³Carnegie Mellon University, Tepper School of Business ⁴Emory University
⁵University of Pennsylvania, The Wharton School ⁶Indian 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 if they are independent, but in reality, they may be deeply interrelated. Addressing bias along one dimension, such as ethnicity or age, can inadvertently influence another dimension, 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. In this paper, we aim to address these questions by introducing BiasConnect, a novel tool designed to analyze and quantify bias interactions in TTI models. Our approach leverages a counterfactual-based framework to generate pairwise causal graphs that reveal the underlying structure of bias interactions for the given text prompt. Additionally, our method provides empirical estimates that indicate how other bias dimensions shift toward or away from an ideal distribution when a given bias is modified. Our estimates have a strong correlation (+0.69) with the interdependency observations post bias mitigation. We demonstrate the utility of BiasConnect for selecting optimal bias mitigation axes, comparing different TTI models on the dependencies they learn, and understanding the amplification of intersectional societal biases in TTI models.

1. Introduction

Text-to-Image (TTI) models such as DALL-E [46], Imagen [51], and Stable Diffusion [48] 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 [10, 11, 60]. These biases manifest across multiple social and 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

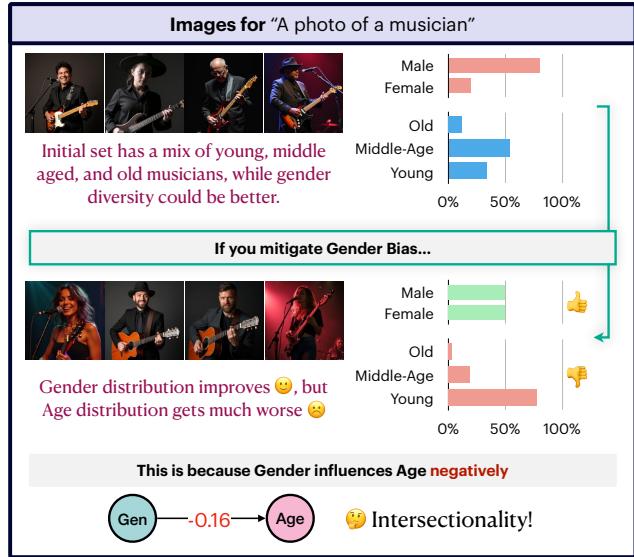


Figure 1. An example output of BiasConnect, revealing the negative impact of bias mitigation along one dimension on another dimension. Here, increasing the gender diversity (GEN) skews age distribution (AGE) for images of musicians generated by Stable Diffusion 1.4 [48].

societal norms [4, 6]. While significant efforts have been made to evaluate and mitigate societal biases in TTI models [5, 10, 11, 18, 23, 62], 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, like age. This relationship between two bias dimensions highlights the intersectional nature of these biases.

The concept of *intersectionality*, first introduced by

Crenshaw [12], 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 [64] (see Supp. A.5). We found that while the targeted biases were reduced in most cases, biases along other axes were negatively affected in over 29.4% 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 consider whether mitigating one bias affects other biases equally or unequally, and if these effects are positive or negative. Answering these questions can improve our understanding of TTI models, enhance their interpretability, and develop better bias mitigation strategies.

To understand how biases in TTI models influence each other, we propose BiasConnect, a first-of-its-kind analysis tool that evaluates societal biases in TTI models while accounting for intersectional relationships, rather than treating them in isolation. Our approach enables us to *identify intersectional relationships and study the positive and negative impacts* that biases have on each other, preventing unintended consequences that may arise from mitigating biases along a single axis. Our tool analyzes intersectionality at the prompt level, but also enables comparative studies across models by aggregating over a set of prompts, thus providing a means to uncover how variations in architecture, datasets, and training objectives contribute to bias entanglement.

Given an input prompt to a TTI model, BiasConnect uses counterfactuals to quantify the impact of bias diversification (intervention) along one bias axis on any other bias axis. We refer to this as a pairwise causal relationship between the axis of intervention and the axis on which the effect is observed, and we visualize these relationships in the form of a *pairwise causal graph*. Additionally, we provide empirical estimates (called Intersectional Sensitivity) that measure how bias mitigation along one dimension influences biases in other dimensions. These empirical estimates serve as weights on the pairwise causal graph. To validate our approach, we show how our estimates correlate with true bias mitigation, and analyze robustness in which different components are systematically modified. Our overall contributions are as follows:

- We propose BiasConnect, a novel intersectional bias analysis tool for TTI models. Through a causal approach, our tool captures interactions between bias axes in a pairwise causal graph and provides empirical estimates of how bias mitigation along one axis affects other axes, through the Intersectional Sensitivity score.
- We show that our empirical estimates strongly correlate (+0.696) with the intersectionality observed post-

mitigation, and through extensive qualitative results, validate the analyses provided by our tool.

- Finally, we demonstrate the usefulness of the tool in conducting audits on multiple open-source TTI models, identifying optimal mitigation strategies to account for intersectional biases, and showing how bias interactions in real-world or a training dataset may change in TTI model generated images.

2. Related Work

2.1. Intersectionality and Bias in AI

Intersectionality, introduced by Crenshaw [12], 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 [13]. In the context of AI, most existing work [15, 26, 33, 34] aligns more closely with the additive model, focusing on quantifying and mitigating biases in intersectional subgroups. This perspective has influenced fairness metrics [16, 19, 22] designed to assess subgroup-level performance, extending across various domains, including natural language processing (NLP) [24, 37, 38, 57] and recent large language models [3, 14, 35, 41], multimodal research [29, 30], and computer vision [56, 63]. 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-driven causal analysis in TTI models. Beyond subgroup analysis, we intervene on a single bias axis to assess its ripple effects on others, revealing independences and interactions.

2.2. Bias in Text-to-Image Models

Extensive research has been conducted on evaluating and mitigating social biases in both image-only models [8, 27, 32, 40, 42, 53, 58, 59, 63] and text-only models [2, 7, 21, 31, 54]. More recently, efforts have expanded to multimodal models and datasets, addressing biases in various language-vision tasks. These investigations have explored biases in embeddings [25], text-to-image (TTI) generation [5, 11, 18, 23, 52, 62, 64], image retrieval [61], image captioning [27, 65], and visual question-answering models [1, 28, 44].

Despite these advances, research on intersectional biases in TTI models remains limited. Existing evaluation frameworks such as T2IAT [62], DALL-Eval [11], and other studies [5, 18, 20, 23] primarily assess biases along predefined axes, such as gender [5, 11, 18, 23, 62], skin tone [5, 11, 18, 23, 62], culture [18, 62], and geographical location [18]. 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 [10], visualizes such interactions, but our approach goes further by systematically quantifying bias interactions and empirically estimating their impact rather than merely identifying correlations.

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_y). BiasConnect, works by systematically altering input prompts and analyzing the resulting distributions of generated images. 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 VQA model (Sec. 3.2). Next, in order to identify whether the intersectional effects of intervening on one bias on another axis is significant, we propose a causal discovery approach, where we employ conditional independence testing (Sec. 3.3) in a pairwise manner between the two bias axes. Finally, to quantify the intersectional effects, and to identify whether these effects are positive or negative, we compute the causal treatment effect, defined as Intersectional Sensitivity (Sec. 3.4).

3.1. Counterfactual Prompts & Image Generation

Given an input prompt P and bias axes $B = [B_1, B_2, \dots, B_n]$, we generate counterfactual prompts $\{CF_i^1, \dots, CF_i^j\}$ using templates from Supp. A.1. 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 [10] and OpenBias [17], where a VQA-based method was used to extract concepts from the generated images. Similar to previous work [10], we use MiniGPT-v2 [9] in a question-answer format to extract attributes from generated images.

For the societal biases we analyze, we have a list of pre-defined questions (Supp. 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). 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. Pairwise Causal Discovery

Given an initial set of bias axes B , we define an intersectional relationship between a pair of biases (B_x, B_y) as $B_x \rightarrow B_y$, indicating that a counterfactual intervention on B_x to mitigate its bias also affects B_y . As a first step, we intervene across all $n \times n$ bias relationships. Using attributes extracted by the VQA, we can count the attributes for a bias axis B_y over any set of images. We construct a contingency table where rows represent the intervened bias axis B_x (e.g., gender with male and female counterfactuals, in Example 2 of Fig. 2), and columns capture the distribution on the target axis B_y (e.g., age with old, middle-aged, and young categories). The values in the contingency tables are the counts of attributes of B_y over the counterfactual image sets of B_x .

Next, we refine these relationships by extracting only statistically significant ones. This ensures that only strong dependencies between different bias pairs are retained. We apply conditional independence testing using the Chi-square (χ^2) test, pruning bias pairs with respect to B_x if their p-value exceeds a predefined threshold (p-value > 0.0001). Bias pairs with a p-value below this threshold are considered strongly dependent, indicating that intervening on B_x results in a significant change in the other bias axis. This process is applied iteratively for all bias axes. This step is referred to as **Pairwise Causal Discovery**, and it returns a set of bias pair relationships where mitigating along one bias axis has led to a strong change in another bias dimension

3.4. Causal Treatment Effect Estimation

While Pairwise Causal Discovery can identify interventions along bias pairs that cause significant changes, it alone is not sufficient to determine whether the impact of interventions on B_x affects B_y in a positive or negative direction with respect to an ideal distribution. This limitation arises because there is no direct comparison to the initial distribution of B_y in the original set of images generated from prompt P , as we had only considered images from the counterfactuals $\{CF_x^1, \dots, CF_x^j\}$ for bias B_x , for pairwise causal discovery. To address this, we propose a metric that quantifies the impact of bias mitigation on dependent biases with respect to an ideal distribution.

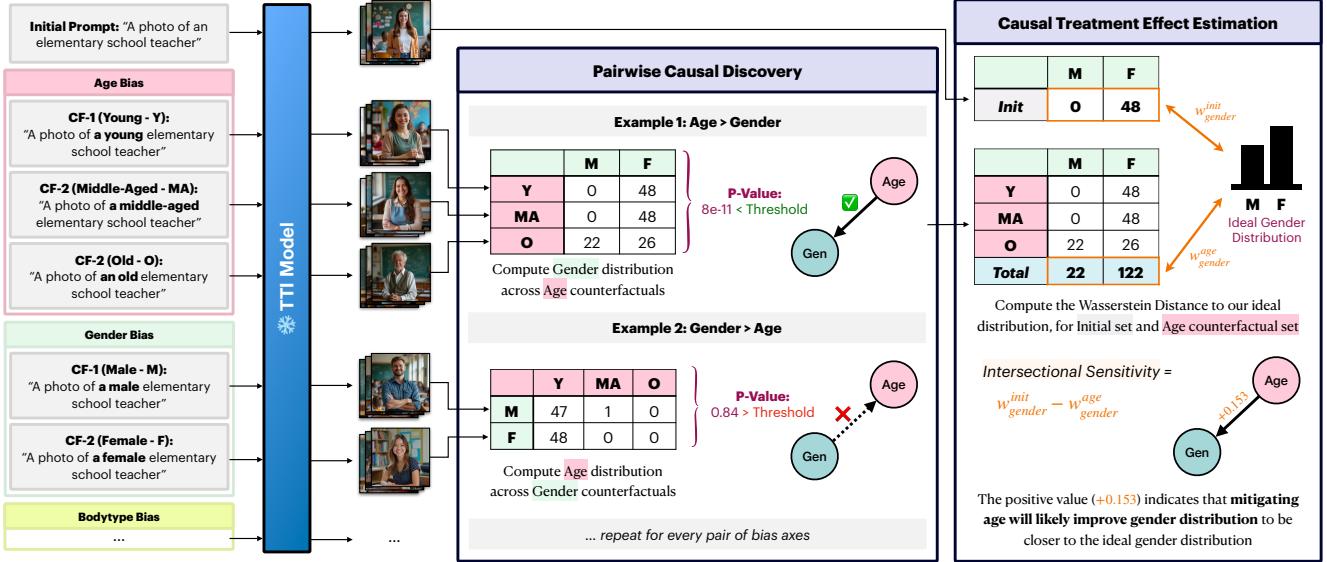


Figure 2. **An overview of BiasConnect.** We use a counterfactual-based approach to measure pairwise causality between bias axes. For dependent axes, we measure the causal effect, estimating how bias mitigation on one axis impacts another.

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 $D_{B_y}^{\text{init}}$. We then compute the Wasserstein distance between this empirical distribution and the ideal distribution:

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

where $W_1(\cdot, \cdot)$ represents the Wasserstein-1 distance. The Wasserstein-1 distance (also known as the Earth Mover's Distance) between two probability distributions D_1 and D_2 is defined as:

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

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.

Intervening on B_x . Next, 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 [64]. 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.

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

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

$$IS_{xy} = w_{B_y}^{\text{init}} - w_{B_y}^{B_x} \quad (4)$$

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. We use Intersectional Sensitivity (IS_{xy}) as a measure of intersectionality for $B_x \rightarrow B_y$.

3.5. Visualization

Following the process above, we have a set of pairwise causal relationships for all significant intersectional bias pairs $B_x \rightarrow B_y$. Furthermore, each pair $B_x \rightarrow B_y$ has an Intersectional Sensitivity score to quantify the intersectional effects. There are many ways to represent these pairwise relationships, including building an $n \times n$ matrix, or

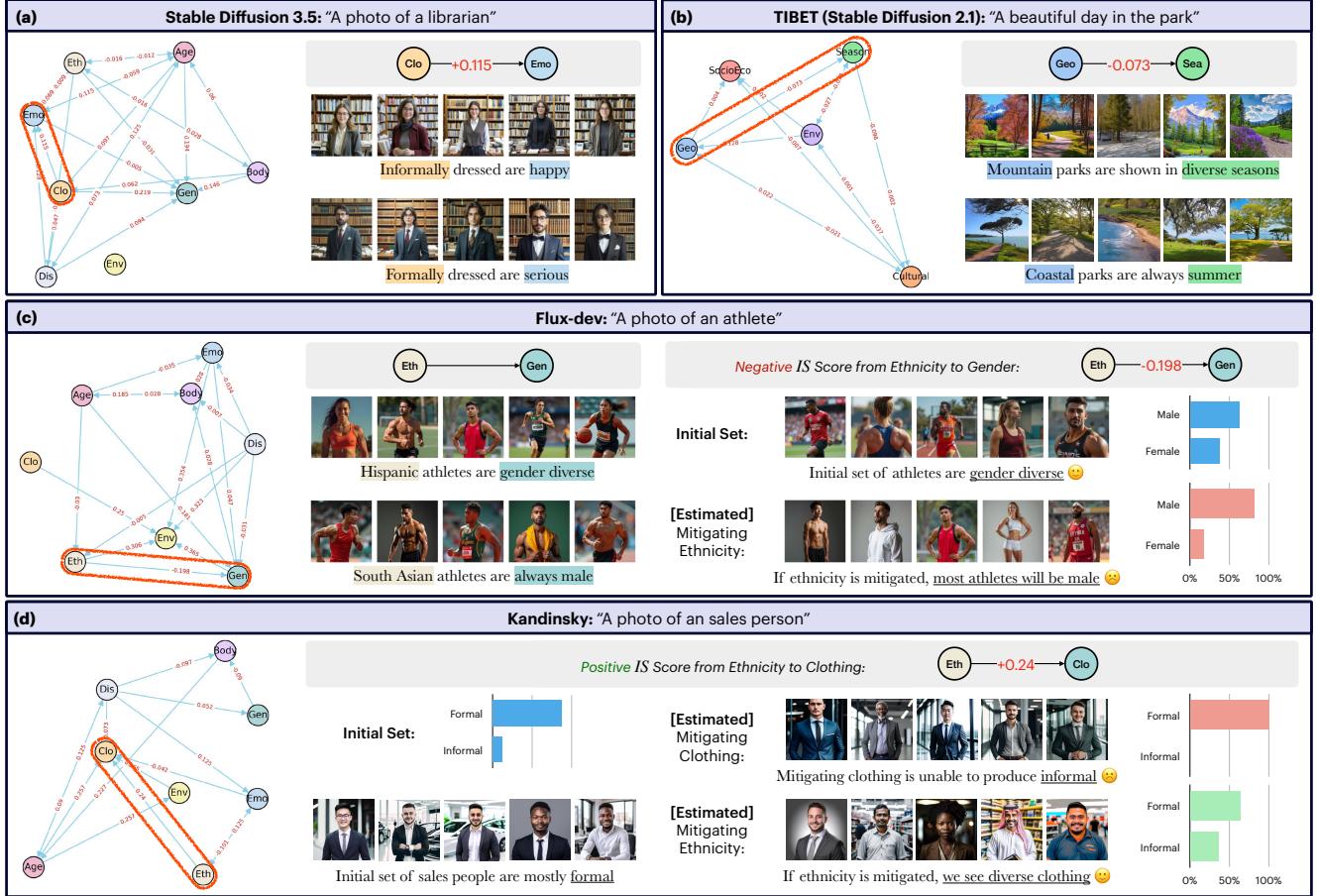


Figure 3. The figure illustrates bias interpretations from Bias Connects, combining all pairwise graphs into one. (a) Shows how mitigating clothing bias also mitigates emotion bias. (b) Explores interactions between non-traditional bias axes in the TIBET dataset. (c) Reveals that generating ethnically diverse athletes reduces gender diversity. (d) Demonstrates that diversifying salesperson clothing is best achieved by increasing ethnic diversity rather than directly specifying clothing variation.

a graph with n nodes and directed edges that represent the relationships between these nodes.

A user of BiasConnect may want to understand all important intersectional effects together. To that end, we adopt a graph representation for our output. This graph is referred to as a **Pairwise Causal Graph** in the rest of the paper. Figures 3 and 5 show examples of such graphs. To interpret this graph, first pick a focal node where the intervention takes place. All outgoing edges from this node indicate intersectional relationships that are statistically significant. The weights of the edges show the Intersectional Sensitivity and can be interpreted as the impact of intervention on the bias axis for the focal node.

4. Causal Interpretations

Pairwise Causal Discovery. Our approach is causal as it involves explicit interventions to measure the effect of one variable on another, aligning with Pearl’s Ladder of Causal-

ity [45]. Rather than analyzing existing images, we actively modify bias attributes (e.g., gender, race, age) in input prompts. However, our pairwise causal discovery pipeline does not capture indirect causal effects between bias axes.

Causal Treatment Effect. The Intersectional Sensitivity metric is a *causal treatment effect* metric because it quantifies how mitigating one bias (B_x) causally influences another bias (B_y) through an intervention-based approach. By actively modifying B_x (e.g., ensuring equal representation across its attributes) and measuring changes in the Wasserstein distance of B_y from an ideal distribution, we estimate the causal impact of debiasing. This aligns with *counterfactual causal inference* [43], where we compare the observed outcome (B_y distribution) with its initial state had no intervention occurred. The method follows *Rubin’s causal model* [49, 50], treating bias mitigation as a *treatment-control experiment*, and can be represented in a *Directed Graph* as $B_x \rightarrow B_y$, making it distinct from mere corre-

lation analysis. The metric IS_{xy} in Equation 4 captures the magnitude of causal influence, providing insights into whether mitigating one bias improves or worsens another.

5. Experiments

In this section, we begin by explaining the two datasets—the occupation prompts and the TIBET dataset—that we use to test BiasConnect (Sec. 5.1). Following that, show the usefulness of BiasConnect by analyzing prompts to study prompt-level bias intersectionality (Sec. 5.2), and validate our Intersectional Sensitivity with the help of a downstream bias mitigation tool, ITI-GEN (Sec. 5.3). Finally, we analyze the robustness of BiasConnect on the number of images generated per prompt, and errors in VQA (Sec. 5.4).

5.1. Models and Datasets

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 Text-to-Image models: Stable Diffusion 1.4, Stable Diffusion 3.5, Flux [36], Playground v2.5 [39] and Kandinsky 2.2 [47, 55]. Further details about the prompts, bias axes, and counterfactuals are provided in the Supp. A.1.

TIBET dataset. The TIBET dataset includes 100 creative prompts with LLM-generated bias axes and counterfactuals [10]. Its diversity of prompts and bias axes, unrestricted to a fixed set, enhances its utility. Additionally, it provides 48 Stable Diffusion 2.1-generated images per initial and counterfactual prompt (See Supp. A.6 for more details).

5.2. Studying prompt-level intersectionality

BiasConnect enables prompt-level analysis of intersectional biases (Fig. 3), helping users identify key bias axes and develop effective mitigation strategies. For instance, in Fig. 3(a), Stable Diffusion 3.5 exhibits a causal link between clothing and emotion bias—informally dressed librarians appear happy, while formally dressed ones seem serious. A strongly positive Intersectional Sensitivity ($IS = 0.115$) indicates that diversifying clothing alone is sufficient to diversify emotion, without explicitly mitigating emotion bias. Conversely, Fig. 3(c) illustrates how ethnicity can negatively impact gender diversity. South Asian athletes, for example, are predominantly depicted as male. The negative Intersectional Sensitivity ($IS = -0.198$) suggests that mitigating ethnicity alone would further skew gender representation toward males. These interpretations of our tool have various applications, including identifying optimal bias mitigation strategies and comparing multiple TTI models, as discussed in Section 6.

Prompt	Edges	Corr.	MaxInf	MaxImp
Pharmacist	12	+0.399	Gender	Age
Scientist	9	+0.600	Clothing	Ethnicity
Doctor	9	+0.638	Age	Disability
Librarian	14	+0.805	Emotion	Age
Nurse	5	<u>+0.997</u>	Age	Disability
Chef	8	+0.757	Bodytype	Ethnicity
Politician	10	+0.782	Emotion	Disability
Overall	-	+0.696	Gender	Age

Table 1. **Correlation Between Estimates and Post-Mitigation Evaluation on ITI-GEN.** The high correlation validates our mitigation estimates. For each prompt, we report one of the most influenced node (MaxInf) and the node with the greatest impact on others (MaxImp).

5.3. Validating Intersectional Sensitivity

Our approach estimates how counterfactual-based mitigation affects bias scores using the Intersectional Sensitivity. To validate this, we debias ITI-GEN, mitigate biases along each dimension, and measure the correlation between pre- and post-mitigation Intersectional Sensitivity values. As shown in Table 1, we achieve an average correlation of +0.696 across occupations, with higher values for specific prompts like Nurse (+0.997). The strong correlation observed between pre- and post-mitigation bias scores suggests that our approach effectively captures the potential impacts of interventions, offering a transparent and data-driven way to evaluate model fairness. More details regarding our experimental setup have been provided in Supp. A.7.

5.4. Robustness of BiasConnect

We analyze the robustness of our method by evaluating the impact of image generation and VQA components on pairwise causal graphs and Intersectional Sensitivity values through experiments on image set size and VQA error rates across occupation prompts. This robustness analysis is useful because it ensures the reliability and stability of our method across varying conditions.

Number of Images. Our method generates 48 images per prompt to study bias distributions reliably. To assess the impact of reducing image count, we analyze changes in the total number of edges in the pairwise causal graph and the percentage change in Intersectional Sensitivity (Fig. 4(a-b)). Removing 8 images (16.6%) results in only 2.4 edge changes and a minor 5.5% shift in Intersectional Sensitivity. Even with 16 images removed (33.3%), only 4.8 edges change, and Intersectional Sensitivity shifts by 8%. This low impact suggests that TTI models consistently generate similar bias distributions (e.g., always depicting nurses as fe-

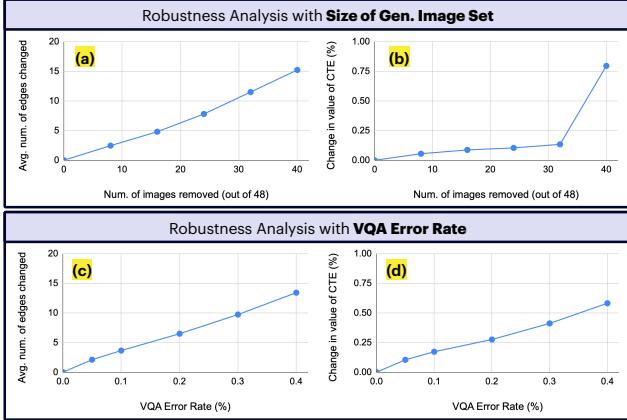


Figure 4. Sensitivity analysis on BiasConnect. We evaluate the robustness of our approach by analyzing the impact of VQA errors and the effect of the number of images on the pairwise causal graph and Intersectional Sensitivity.

males), preserving overall trends despite fewer images. However, excessive pruning significantly affects the analysis—removing 40 images (83%) leads to a sharp 79% change in Intersectional Sensitivity. This demonstrates that our approach is robust to moderate reductions in image count but breaks down when the sample size is too small. While a sufficiently large image set enhances reliability, exceeding 48 images offers only marginal analytical benefits. **VQA Error Rate.** In Fig. 4(c-d), we show the impact of VQA errors on the graph and Intersectional Sensitivity values. We randomly change the VQA answers to a different answer (simulating an incorrect answer) at different thresholds, from 5% to 40% of the time. We observe that with low error rates of 5% and 10%, the impact on number of edges changed is low, with averages of 2.1 and 3.65 edges respectively. However, the small changes in VQA answers does impact Intersectional Sensitivity values, at 10% and 17.3% respectively, as the impact is compounded by the fact that we use both the initial and the counterfactual distributions to obtain this value, 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. Our study shows the graph is more robust if the error rate is below 20%. As VQA models improve, achieving error rates for robust graphs becomes practical.

6. Applications

6.1. Applying BiasConnect to analyze TTI models

To compare bias interactions across models, we aggregate results from all prompts to create a unified representation, enabling a high-level analysis of bias trends (Fig. 5). Details on the aggregation process are in the Supp. A.8.

Identifying high-impact biases. Some biases act as pri-

mary sources, influencing multiple others, while some function as effects, shaped by upstream factors. A node’s impact is measured by its outgoing edges (**MaxImp**, Table 1), while its susceptibility to influence is quantified by incoming edges (**MaxInf**). This helps in model selection based on specific bias priorities.

As an example, let’s analyze how this information can help in selecting appropriate models using the global graphs in Figure 5. If a user prioritizes robustness to age-related bias when selecting a model, Kandinsky 2.2 would be the best choice, as its Age node is the least influenced by other biases in the global analysis. This means that modifying other attributes (e.g., gender or clothing) has minimal unintended effects on age representation, ensuring more stable and independent age depictions across generated images.

Similarly, if the goal is to generate occupation-related images while minimizing unintended bias propagation across other attributes, Playground 2.5 is the optimal choice. In this model, variations in body type have the least impact on other biases, meaning changes in body shape do not disproportionately affect other attributes like gender, ethnicity, or perceived professionalism. This makes Playground 2.5 preferable in scenarios where maintaining fairness across multiple dimensions while altering body type is critical. By analyzing bias influence and susceptibility, users can make informed choices based on fairness priorities, whether aiming for stability in a bias axis or minimizing unintended shifts in related attributes.

6.2. Studying Real-World Biases

BiasConnect can also be used to compare the distribution of images generated by TTI models with real-world data. To demonstrate this, we sampled 48 images of computer programmers, and 48 images of male and female computer programmers each from the internet. We then compared the pairwise causal graph for gender in the real-world distribution to the one generated by Stable Diffusion 3.5. Our analysis (Fig. 6) reveals that in real-world data, gender diversification primarily influences body type in a positive manner. However, in Stable Diffusion 3.5, gender impacts both emotion and body type, with a negative effect on body type, suggesting that increasing gender diversity reduces body type diversity. Studies like these are valuable in identifying discrepancies between generated images and real-world distributions or training datasets, and show how real-wold bias interactions may be amplified in TTI models. Supp A.9 has further details on this process.

6.3. Uncovering Optimal Bias Mitigation Strategies

BiasConnect quantifies the impact of one bias on another, helping identify effective bias mitigation strategies for a given prompt. We illustrate this with three examples:

Clothing and Emotion Bias (Fig. 3(a)) – Stable Diffu-

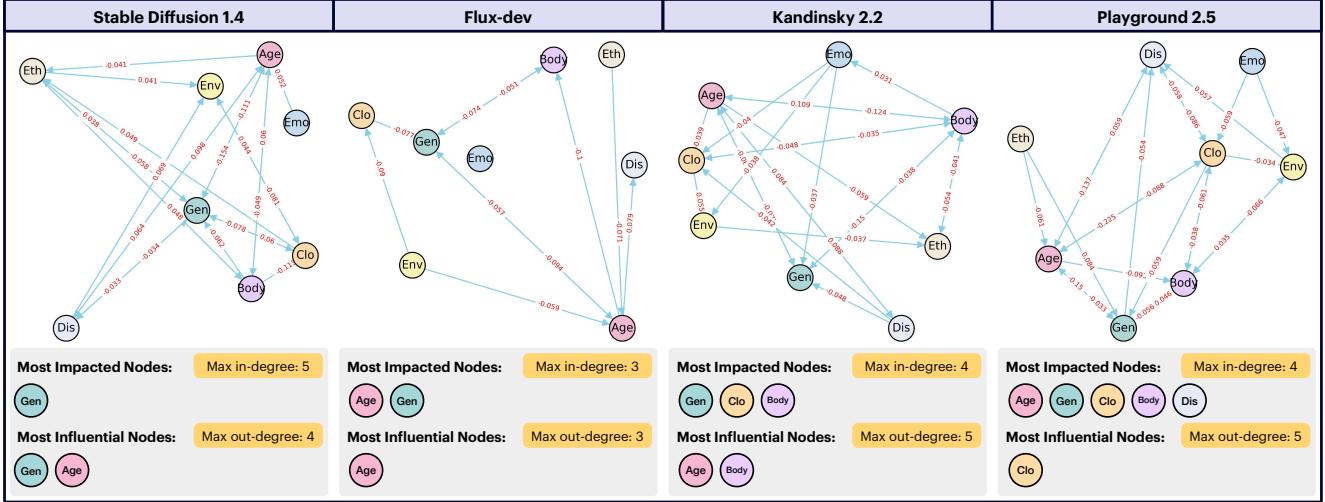


Figure 5. We compare aggregated causal graphs for four models: Stable Diffusion 1.4, Flux-dev, Kandinsky 2.2, and Playground 2.5. These graphs combine pairwise causal relationships across all bias axes, accumulated from occupation prompts in our dataset.

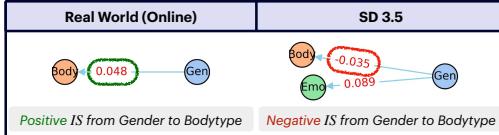


Figure 6. Comparison of real-world and Stable Diffusion 3.5 pairwise causal relationships for gender in computer programmer images. In the real world, gender diversification increases body type diversity, whereas in Stable Diffusion 1.4, it has a negative impact.

sion 3.5 exhibits a causal link where informally dressed librarians appear happy, while formally dressed ones seem serious. A positive Intersectional Sensitivity ($IS = 0.115$) suggests that diversifying clothing alone is sufficient to diversify emotion without explicitly addressing emotion bias.

Ethnicity and Gender Bias (Fig. 3(c)) – South Asian athletes are predominantly depicted as male. A negative Intersectional Sensitivity ($IS = -0.198$) indicates that mitigating ethnicity alone would further skew gender representation toward males. A better approach would involve fine-tuning on a dataset with more female South Asian athletes to improve disentanglement between ethnicity and gender.

Ethnicity and Clothing Diversity (Fig. 3(d)) – For salespersons, the best way to diversify clothing styles is not by directly mitigating clothing bias but by increasing ethnic diversity in generated images. This reveals a hidden mitigation strategy where altering one axis (ethnicity) impacts another (clothing) more than direct intervention.

These examples highlight a key utility of BiasConnect: enabling users to adopt complementary bias mitigation strategies based on their specific needs. In some cases, mitigating one bias naturally diversifies another, reducing

the need for direct intervention. In other cases, addressing one axis may worsen another, requiring a more targeted approach. Finally, certain biases may be best mitigated indirectly by adjusting a different, more influential axis.

7. Conclusion

Our study proposes a tool to investigate intersectional biases in TTI models. While prior research has explored bias detection and mitigation in generative models, to the best of our knowledge, no previous work has focused on understanding how biases influence one another. We believe our work makes a significant contribution by enabling a more nuanced analysis of bias interactions. Beyond academic research, BiasConnect has practical applications, including comparing biases dependencies learned across different models, establishing empirical guarantees for mitigation, and determining optimal mitigation approaches that account for intersectionality. We hope that this tool will facilitate more informed decision-making for AI practitioners, policymakers, and developers, ultimately leading to more equitable and transparent generative models.

While BiasConnect provides a valuable framework, it represents only an initial step toward a more comprehensive causal approach to understanding intersectionality. Our current setup does not allow us to reason about indirect causal effects, or develop an optimal bias mitigation strategy that utilizes our tool to mitigate multiple biases simultaneously. Addressing these challenges presents an important avenue for future research.

Ethical Considerations. We acknowledge that the presence of biases in generative AI models can lead to real-world harms, reinforcing stereotypes and disproportionately

affecting marginalized groups. Our tool is intended to provide researchers and practitioners with a means to better understand and mitigate these biases, rather than to justify or amplify them. Additionally, we recognize that bias analysis can be sensitive to the choice of datasets, evaluation methods, and experimental assumptions, and we encourage future work to refine and expand upon our approach.

References

- [1] Lavisha Aggarwal and Shruti Bhargava. Fairness in ai systems: Mitigating gender bias from language-vision models. *arXiv preprint arXiv:2305.01888*, 2023. [1](#) [2](#)
- [2] Jaimeen Ahn and Alice Oh. Mitigating language-dependent ethnic bias in bert. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 533–549, 2021. [2](#)
- [3] Xuechunzi Bai, Angelina Wang, Ilia Sucholutsky, and Thomas L Griffiths. Explicitly unbiased large language models still form biased associations. *Proceedings of the National Academy of Sciences*, 122(8):e2416228122, 2025. [2](#)
- [4] Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Margaret Shmitchell. 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, 2021. [1](#)
- [5] Federico Bianchi, Pratyusha Kalluri, Esin Durmus, Faisal Ladhak, Myra Cheng, Debora Nozza, Tatsunori Hashimoto, Dan Jurafsky, James Zou, and Aylin Caliskan. 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, 2023. [1](#) [2](#)
- [6] Abeba Birhane and Vinay Uday Prabhu. Multimodal datasets: Misogyny, pornography, and malignant stereotypes. *arXiv preprint arXiv:2110.01963*, 2021. [1](#)
- [7] Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. *Advances in neural information processing systems*, 29, 2016. [2](#)
- [8] Joy Buolamwini and Timnit Gebru. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on fairness, accountability and transparency*, pages 77–91. PMLR, 2018. [2](#)
- [9] Jun Chen, Deyao Zhu, Xiaoqian Shen, Xiang Li, Zechun Liu, Pengchuan Zhang, Raghuraman Krishnamoorthi, Vikas Chandra, Yunyang Xiong, and Mohamed Elhoseiny. Minigpt-v2: large language model as a unified interface for vision-language multi-task learning. *arXiv preprint arXiv:2310.09478*, 2023. [3](#)
- [10] Aditya Chinchure, Pushkar Shukla, Gaurav Bhatt, Kiri Salij, Kartik Hosanagar, Leonid Sigal, and Matthew Turk. Tibet: Identifying and evaluating biases in text-to-image generative models. In *European Conference on Computer Vision*, pages 429–446. Springer, 2024. [1](#) [3](#) [6](#) [2](#)
- [11] Jaemin Cho, Abhay Zala, and Mohit Bansal. Dall-eval: 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, 2023. [1](#) [2](#)
- [12] Kimberle Crenshaw. 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, 1989. [2](#)
- [13] Tommy J Curry. Killing boogymen: Phallicism and the misandric mischaracterizations of black males in theory. *Res Philosophica*, 2018. [2](#)
- [14] Hannah Devinney, Jenny Björklund, and Henrik Björklund. 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, 2024. [2](#)
- [15] Emily Diana and Alexander Williams Tolbert. Correcting underrepresentation and intersectional bias for classification. *arXiv preprint arXiv:2306.11112*, 2023. [2](#)
- [16] Emily Diana, Wesley Gill, Michael Kearns, Krishnaram Kenthapadi, and Aaron Roth. Minimax group fairness: Algorithms and experiments. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, pages 66–76, 2021. [2](#)
- [17] Moreno D’Incà, Elia Peruzzo, Massimiliano Mancini, Dejia Xu, Vudit Goel, Xingqian Xu, Zhangyang Wang, Humphrey Shi, and Nicu Sebe. 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, 2024. [3](#)
- [18] Piero Esposito, Parmida Atighehchian, Anastasis Germanidis, and Deepti Ghadiyaram. Mitigating stereotypical biases in text to image generative systems. *arXiv preprint arXiv:2310.06904*, 2023. [1](#) [2](#)
- [19] James R Foulds, Rashidul Islam, Kamrun Naher Keya, and Shimei Pan. An intersectional definition of fairness. In *2020 IEEE 36th international conference on data engineering (ICDE)*, pages 1918–1921. IEEE, 2020. [2](#)
- [20] Felix Friedrich, Patrick Schramowski, Manuel Brack, Lukas Struppek, Dominik Hintersdorf, Sasha Luccioni, and Kristian Kersting. Fair diffusion: Instructing text-to-image generation models on fairness. *arXiv preprint arXiv:2302.10893*, 2023. [2](#)
- [21] Ismael Garrido-Muñoz, Arturo Montejo-Ráez, Fernando Martínez-Santiago, and L Alfonso Ureña-López. A survey on bias in deep nlp. *Applied Sciences*, 11(7):3184, 2021. [2](#)
- [22] Avijit Ghosh, Lea Genuit, and Mary Reagan. Characterizing intersectional group fairness with worst-case comparisons. In *Artificial Intelligence Diversity, Belonging, Equity, and Inclusion*, pages 22–34. PMLR, 2021. [2](#)
- [23] Sourojit Ghosh and Aylin Caliskan. ‘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, 2023. [1](#) [2](#)

- [24] Wei Guo and Aylin Caliskan. 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, 2021. 2
- [25] Kimia Hamidieh, Haoran Zhang, Thomas Hartvigsen, and Marzyeh Ghassemi. Identifying implicit social biases in vision-language models. 2023. 2
- [26] Ursula Hebert-Johnson, Michael Kim, Omer Reingold, and Guy Rothblum. Multicalibration: Calibration for the (Computationally-identifiable) masses. In *Proceedings of the 35th International Conference on Machine Learning*, pages 1939–1948. PMLR, 2018. 2
- [27] Lisa Anne Hendricks, Kaylee Burns, Kate Saenko, Trevor Darrell, and Anna Rohrbach. Women also snowboard: Overcoming bias in captioning models. In *Proceedings of the European conference on computer vision (ECCV)*, pages 771–787, 2018. 2
- [28] Yusuke Hirota, Yuta Nakashima, and Noa Garcia. Gender and racial bias in visual question answering datasets. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, pages 1280–1292, 2022. 2
- [29] Elizabeth Hoepfinger. Racial and intersectional debiasing of contrastive language image pretraining. Master’s thesis, University of Georgia, 2023. 2
- [30] Phillip Howard, Avinash Madasu, Tiep Le, Gustavo Lujan Moreno, Anahita Bhiwandiwalla, and Vasudev Lal. Social-counterfactuals: Probing and mitigating intersectional social biases in vision-language models with counterfactual examples. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11975–11985, 2024. 2
- [31] Ben Hutchinson, Vinodkumar Prabhakaran, Emily Denton, Kellie Webster, Yu Zhong, and Stephen Denayl. 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, 2020. 2
- [32] Aparna R Joshi, Xavier Suau Cuadros, Nivedha Sivakumar, Luca Zappella, and Nicholas Apostoloff. Fair sa: Sensitivity analysis for fairness in face recognition. In *Algorithmic fairness through the lens of causality and robustness workshop*, pages 40–58. PMLR, 2022. 2
- [33] Loukas Kavouras, Konstantinos Tsopelas, Giorgos Giannopoulos, Dimitris Sacharidis, Eleni Psaroudaki, Nikolaos Theologitis, Dimitrios Rontogiannis, Dimitris Fotakis, and Ioannis Emiris. Fairness aware counterfactuals for subgroups. *Advances in Neural Information Processing Systems*, 36:58246–58276, 2023. 2
- [34] Michael Kearns, Seth Neel, Aaron Roth, and Zhiwei Steven Wu. Preventing fairness gerrymandering: Auditing and learning for subgroup fairness. In *International conference on machine learning*, pages 2564–2572. PMLR, 2018. 2
- [35] Hannah Kirk, Yennie Jun, Haider Iqbal, Elias Benussi, Filippo Volpin, Frederic A. Dreyer, Aleksandar Shtedritski, and Yuki M. Asano. Bias Out-of-the-Box: An Empirical Analysis of Intersectional Occupational Biases in Popular Generative Language Models, 2021. arXiv:2102.04130 [cs]. 2
- [36] Black Forest Labs. Flux. <https://github.com/black-forest-labs/flux>, 2024. 6
- [37] John P Lalor, Yi Yang, Kendall Smith, Nicole Forsgren, and Ahmed Abbasi. 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, 2022. 2
- [38] Ida Marie S Lassen, Mina Almasi, Kenneth Enevoldsen, and Ross Deans Kristensen-McLachlan. 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, 2023. 2
- [39] Daiqing Li, Aleks Kamko, Ehsan Akhgari, Ali Sabet, Lin-miao Xu, and Suhail Doshi. Playground v2.5: Three insights towards enhancing aesthetic quality in text-to-image generation, 2024. 6
- [40] Bingyu Liu, Weihong Deng, Yaoyao Zhong, Mei Wang, Jian Hu, Xunqiang Tao, and Yaohai Huang. Fair loss: Margin-aware reinforcement learning for deep face recognition. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10052–10061, 2019. 2
- [41] Weicheng Ma, Brian Chiang, Tong Wu, Lili Wang, and Soroush Vosoughi. Intersectional stereotypes in large language models: Dataset and analysis. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 8589–8597, 2023. 2
- [42] Nicole Meister, Dora Zhao, Angelina Wang, Vikram V Ramaswamy, Ruth Fong, and Olga Russakovsky. Gender artifacts in visual datasets. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4837–4848, 2023. 2
- [43] Stephen L Morgan and Christopher Winship. *Counterfactuals and causal inference: Methods and principles for social research*. Cambridge University Press, 2014. 5
- [44] Sungho Park, Sunhee Hwang, Jongkwang Hong, and Hyeran Byun. Fair-vqa: Fairness-aware visual question answering through sensitive attribute prediction. *IEEE Access*, 8: 215091–215099, 2020. 2
- [45] Judea Pearl. *Causality*. Cambridge university press, 2009. 5
- [46] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. *arXiv preprint arXiv:2102.12092*, 2021. 1
- [47] Anton Razzhigaev, Arseniy Shakhmatov, Anastasia Maltseva, Vladimir Arkhipkin, Igor Pavlov, Ilya Ryabov, Angelina Kuts, Alexander Panchenko, Andrey Kuznetsov, and Denis Dimitrov. 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, 2023. 6
- [48] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 1

- [49] Donald B Rubin. Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of educational Psychology*, 66(5):688, 1974. 5
- [50] Donald B Rubin. Bayesian inference for causality: The importance of randomization. In *The Proceedings of the social statistics section of the American Statistical Association*, page 239. American Statistical Association Alexandria, VA, 1975. 5
- [51] 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, et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022. 1
- [52] Preethi Seshadri, Sameer Singh, and Yanai Elazar. The bias amplification paradox in text-to-image generation. *arXiv preprint arXiv:2308.00755*, 2023. 2
- [53] Laleh Seyyed-Kalantari, Haoran Zhang, Matthew BA McDermott, Irene Y Chen, and Marzyeh Ghassemi. Under-diagnosis bias of artificial intelligence algorithms applied to chest radiographs in under-served patient populations. *Nature medicine*, 27(12):2176–2182, 2021. 2
- [54] Deven Santosh Shah, H Andrew Schwartz, and Dirk Hovy. 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, 2020. 2
- [55] Arseniy Shakhmatov, Anton Razzhigaev, Aleksandr Nikolich, Vladimir Arkhipkin, Igor Pavlov, Andrey Kuznetsov, and Denis Dimitrov. kandinsky 2.2, 2023. 6
- [56] Ryan Steed and Aylin Caliskan. 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, 2021. 2
- [57] Yi Chern Tan and L Elisa Celis. Assessing social and intersectional biases in contextualized word representations. *Advances in neural information processing systems*, 32, 2019. 2
- [58] Angelina Wang and Olga Russakovsky. Overwriting pre-trained bias with finetuning data. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3957–3968, 2023. 2
- [59] Angelina Wang, Alexander Liu, Ryan Zhang, Anat Kleiman, Leslie Kim, Dora Zhao, Iroha Shirai, Arvind Narayanan, and Olga Russakovsky. Revise: A tool for measuring and mitigating bias in visual datasets. *International Journal of Computer Vision*, 130(7):1790–1810, 2022. 2
- [60] Angelina Wang, Vikram V Ramaswamy, and Olga Russakovsky. 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*, pages 336–349, New York, NY, USA, 2022. Association for Computing Machinery. 1
- [61] Jialu Wang, Yang Liu, and Xin Wang. Assessing multilingual fairness in pre-trained multimodal representations. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2681–2695, 2022. 2
- [62] Jialu Wang, Xinyue Gabby Liu, Zonglin Di, Yang Liu, and Xin Wang. 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, 2023. 1, 2
- [63] Zeyu Wang, Clint Qinami, Ioannis Christos Karakozis, Kyle Genova, Prem Nair, Kenji Hata, and Olga Russakovsky. 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, 2020. 2
- [64] Cheng Zhang, Xuanbai Chen, Siqi Chai, Chen Henry Wu, Dmitry Lagun, Thabo Beeler, and Fernando De la Torre. Iti-gen: Inclusive text-to-image generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3969–3980, 2023. 2, 4
- [65] Eric Zhao, De-An Huang, Hao Liu, Zhiding Yu, Anqi Liu, Olga Russakovsky, and Anima Anandkumar. Scaling fair learning to hundreds of intersectional groups. 2021. 2

BiasConnect: Investigating Bias Interactions in Text-to-Image Models

Supplementary Material

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 <i>occupation</i> A photo of a middle-aged <i>occupation</i> A photo of a young <i>occupation</i>
ethnicity bias	A photo of a white <i>occupation</i> A photo of a african american <i>occupation</i> A photo of an asian <i>occupation</i> A photo of a south asian <i>occupation</i> A photo of a middle eastern <i>occupation</i> A photo of a hispanic <i>occupation</i>
bodytype bias	A photo of a obese <i>occupation</i> A photo of a slim <i>occupation</i> A photo of a muscular <i>occupation</i>
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 2. Occupation Dataset.

A. Supplementary Material

A.1. Occupation Dataset

Table 2 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 gener-

ation are the default settings provided for these models on Huggingface.

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 [10], 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 [10] instead of our templated questions from above. Moreover, in the causal discovery process, because tibet concepts are more diverse than the fixed attributes we use with occupation prompts, our p-value threshold changes to 0.05.

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 2, we mitigate every bias axis independently. We then compute Intersectional Sensitivity 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:

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

$$IS_{xy}^{mit(x)} = w_{B_y}^{\text{init}} - w_{B_y}^{B_x} \quad (6)$$

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 8 below. Fig 8(b) shows how diversifying on an axis like Geography can help diversify the Ethnicity distribution.

A.7. Validating Mitigation Effect Estimation

Our approach provides empirical estimates of how a counterfactual-based mitigation strategy may influence an intersectional relationship $B_x \rightarrow B_y$ 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. For all 26 occupations, we consider all intersectional relationships $B_x \rightarrow B_y$, 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:

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

$$IS_{xy}^{mit(x)} = w_{B_y}^{\text{init}} - w_{B_y}^{B_x} \quad (8)$$

Note that these equations are the same as the ones we used in Supp. A.5, with the main difference being that, in this case, we only consider the scores for the intersectional relationships $B_x \rightarrow B_y$ found through causal discovery. 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 $B_x \rightarrow B_y$ present for each prompt.

In Table 1 in the main paper, we show these findings, and highlight that we achieved an average correlation of +0.696, suggesting that our method effectively estimates the potential impacts of bias interventions without actually doing the mitigation step itself, which often requires fine-tuning some or all parts of the diffusion model.

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.8. Global Aggregations

In order to do a comparative analysis of intersectionality across models over a dataset of prompts, we perform an aggregation step. For the 26 occupation prompts, we first start by using counterfactuals and VQA to identify attributes over all bias axes in B . Now, in the Causal Discovery step, we build contingency tables that aggregate attributes over all CF prompts across all the occupations. For example, when considering the intersectional relationship $Gender \rightarrow Age$, we consider all images for male *occupation* and female *occupation* for all occupations for the rows of the contingency matrix, and count over the *Age* attributes young, middle-aged, old to find the overall global distribution. This gives us the global contingency table for any bias pair. We follow the steps in Sec. 3.3 to obtain this list of bias intersectionality relationships that are significant. Next, in order to compute Intersectional Sensitivity, we use the same contingency table and sum across its columns to get $D_{B_y}^{global(B_x)}$. For the initial distribution, we accumulate attributes across all initial prompt images for all occupations, to give us $D_{B_y}^{global(init)}$. We can now compute Intersectional Sensitivity as:

$$w_{B_y}^{global(init)} = W_1(D_{B_y}^{global(init)}, D^*) \quad (9)$$

$$w_{B_y}^{global(B_x)} = W_1(D_{B_y}^{global(B_x)}, D^*) \quad (10)$$

$$IS_{global(xy)} = w_{B_y}^{global(init)} - w_{B_y}^{global(B_x)} \quad (11)$$

Given the large number of images (as we aggregate over multiple sets), we choose to use a p-value threshold of 0.00005, and we further discard edges in the pairwise causal graph where the $-0.03 > IS_{global(xy)} > 0.03$.

A.9. Studying Real World Biases

BiasConnect can be used to compare bias dependencies 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 $B = [B_1, B_2, \dots, B_n]$ 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).
- $D_{B_y}^{\text{TTI}}$: distribution of B_y in TTI-generated images.

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

$$w_{B_y}^{\text{init}} = W_1(D_{B_y}^{\text{TTI}}, D_{B_y}^{\text{real}}) \quad (12)$$

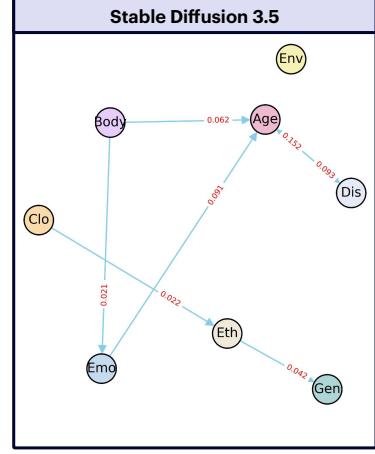


Figure 7. Global graph for Stable Diffusion 3.5. This graph is not included in the main paper due to space constraints.

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

$$w_{B_y}^{B_x} = W_1(D_{B_y}^{B_x}, D_{B_y}^{\text{real}}) \quad (13)$$

The Intersectional Sensitivity Score IS_{xy} for the effect of changing B_x on B_y measures the difference between $w_{B_y}^{\text{init}}$ and $w_{B_y}^{B_x}$ similar to the one calculated in Eq 4. To measure overall intersectional bias amplification, we compute:

$$\mathcal{I} = \sum_{x \neq y} |IS_{xy}| \quad (14)$$

where a high \mathcal{I} indicates strong intersectional bias amplification, while a low \mathcal{I} suggests minimal entanglement.

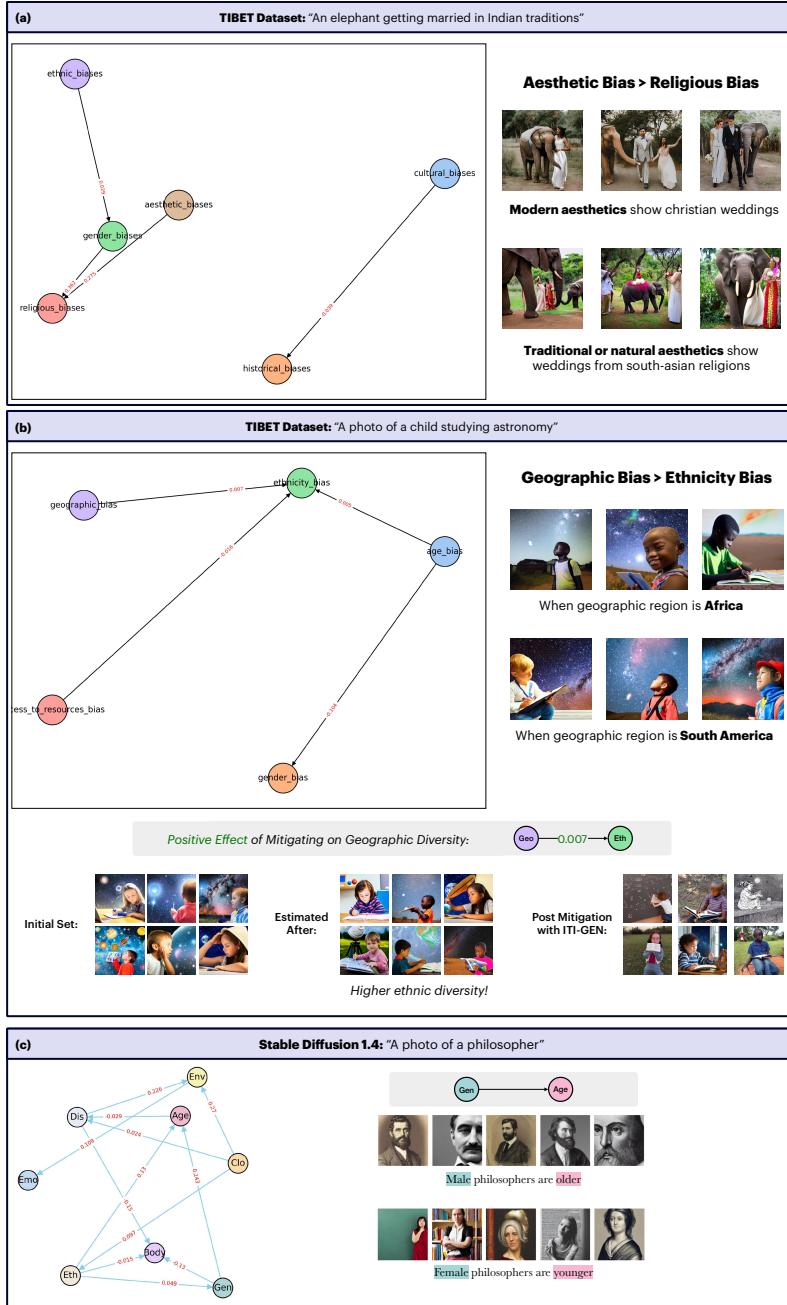


Figure 8. Additional examples on TIBET (a-b) and Occupation prompt (c) on prompt-level analysis provided by BiasConnect.