

BIGbench: A Unified Benchmark for Evaluating Multi-dimensional Social Biases in Text-to-Image Models

Anonymous ACL submission

Abstract

Text-to-Image (T2I) generative models are becoming increasingly crucial due to their ability to generate high-quality images, but also raise concerns about social biases, particularly in human image generation. Sociological research has established systematic classifications of bias. Yet, existing studies on bias in T2I models largely conflate different types of bias, impeding methodological progress. In this paper, we introduce **BIGbench**, a unified benchmark for **Biases of Image Generation**, featuring a carefully designed dataset. Unlike existing benchmarks, BIGbench classifies and evaluates biases across four dimensions to enable a more granular evaluation and deeper analysis. Furthermore, BIGbench applies advanced multi-modal large language models to achieve fully automated and highly accurate evaluations. We apply BIGbench to evaluate eight representative T2I models and three debiasing methods. Our human evaluation results by trained evaluators from different races underscore BIGbench’s effectiveness in aligning images and identifying various biases. Moreover, our study also reveals new research directions about biases with insightful analysis of our results. Our work is [openly accessible](#).

1 Introduction

As a key technology in AI-generated content (AIGC), Text-to-Image (T2I) models attract considerable attention (Esser et al., 2024; Song et al., 2024; Li et al., 2024b). However, these models often amplify societal biases, reinforcing harmful stereotypes and perpetuating discrimination (Guilbeault et al., 2024). Studies reveal that T2I models frequently depict high-status professions as white, middle-aged men, even with neutral prompts, reflecting inherent gender and racial biases (Cho et al., 2023). Additionally, they associate professions with specific genders, further entrenching societal stereotypes (Milne, 2023; Bianchi et al.,

2023). Mitigating these biases is essential to prevent AI from exacerbating social inequalities.

While prior efforts to evaluate and decrease biases in T2I models (Luo et al., 2024; Shrestha et al., 2024), existing benchmarks remain several limitations. **(I)** They offer limited prompt diversity and coverage, largely focusing only on occupational biases. For instance, DALL-EVAL (Cho et al., 2023) evaluated biases towards occupations and genders with 252 prompts. **(II)** These benchmarks compare a narrow set of models and do not assess debiasing methods, leaving their broader applicability unknown. **(III)** They focus on specific bias types rather than providing comprehensive evaluations. For example, HRS-Bench (Bakr et al., 2023) considers only the situation where models fail to generate images with specific protected attributes (e.g., age and race), while DALL-EVAL measured gender and skin color diversity without considering protected attributes into prompts. **(IV)** Current benchmarks directly use the general machine learning (ML) bias definition, lacking a tailored classification system for T2I models.

To address the issues, we introduce a unified and adjustable bias benchmark named **Biases of Image Generation Benchmark**, abbreviated as **BIGbench**. The comparative overview of BIGbench against existing benchmarks is shown in Table 1. We establish a comprehensive definition system and classify biases across four dimensions (see Section 2). We construct the dataset with 47,040 prompts, covering occupations, characteristics and social relations. BIGbench employs fully automated evaluations based on the alignment by a fine-tuned multi-modal LLM, featuring adjustable evaluation metrics. The evaluation results cover implicit generative bias, explicit generative bias, ignorance, and discrimination. These characteristics make BIGbench suitable for automated bias evaluation tasks for any T2I model. We evaluate 8 mainstream T2I models and 3 debiased methods with BIGbench. Based on the

| Benchmark | Model | Prompt | Metric | Multi-level |
|-----------|-------|--------|--------|-------------|
| DALL-Eval | 4 | 252 | 6 | no |
| HRS-Bench | 5 | 3000 | 3 | no |
| ENTIGEN | 3 | 246 | 4 | yes |
| TIBET | 2 | 100 | 7 | no |
| BIGbench | 11 | 47040 | 18 | yes |

Table 1: Summary of existing benchmarks as four characteristics are considered for each benchmark.

results, we discuss the performance of the models in different biases and explore the effects of distillation (Meng et al., 2023) and irrelevant attributes on biases. To ensure the reliability of the results, we conduct human evaluations on 1,000 images for alignment, achieving significant consistency. Our contributions are summarized as follows:

- We establish a four-dimensional bias definition system for T2I models based on sociological and machine ethics research, which categorizes biases by acquired, protected attributes, manifestation, and visibility, enabling more precise understanding and mitigation.
- We present BIGbench, a unified benchmark for evaluating T2I model biases. It features a 47,040-prompt dataset and an automated, high-accuracy evaluation pipeline using a fine-tuned multi-modal LLM, providing a versatile and efficient research tool.
- We evaluate 8 mainstream T2I models and 3 debiasing methods, offering the first comparative analysis of debiasing techniques and exploring the impacts of distillation and irrelevant attributes. Our human-validated findings provide guidance for developing fairer AIGC systems.

2 Definition System

To overcome the limitations of existing benchmarks, we propose a new definition and classification system based on sociological and machine-ethical studies on bias (Landy et al., 2018; Varona and Suárez, 2022; Chouldechova, 2017). We consider our definition system from four dimensions: acquired attributes, protected attributes, manifestation of bias, and visibility of bias. Any kind of bias can be represented using these four dimensions.

Acquired Attribute. An acquired attribute is a trait that individuals acquire through their experiences, actions, or choices. It can be changed over time through personal effort, experience, or other

activities. They are used as a reasonable basis for decision-making, but also possible to be related to bias. *Typical protected attributes include occupation, social relation, and personal wealth.*

Protected Attribute. A protected attribute is a shared identity of one social group, which is legally or ethically protected from being used as grounds for decision-making to prevent bias. It is difficult to change as it is usually related to physiological traits. *Typical protected attributes include race, sex, age, and disability status.*

Manifestation of Bias. This definitional dimension (Devine, 1989) stems from a highly influential psychological study, which deconstructed prejudice (analogous to bias in ML) into a combination of automatic and controlled processes, referring to the unconscious neglect or deliberate overemphasis of associations between certain groups and specific concepts. Through analyzing the distributional characteristics of T2I model outputs across protected attributes, we define these two processes as *ignorance* and *discrimination* to reflect their specific manifestations in generative models.

Ignorance refers to the phenomenon where T2I models consistently generate images depicting a specific demographic group, regardless of prompts suggesting positive terms or negative terms. This bias perpetuates a limited, homogenized view of diverse characteristics and roles, reinforcing a narrowed societal perception.

Discrimination refers to the phenomenon where T2I models disproportionately associate positive and high-status terms with images of certain demographic groups while aligning negative and low-status terms with other groups. This bias reinforces stereotypes about certain social groups.

From an ML perspective, Ignorance and Discrimination serve as indicators that reflect the imbalanced distribution of model training data (He and Garcia, 2009; Mehrabi et al., 2021). Ignorance arises when certain groups are severely underrepresented in terms of sample quantity and diversity within the training data. During the training process, models tend to learn the dominant feature distributions in the dataset, resulting in insufficient feature learning for these underrepresented groups. The phenomenon of Discrimination occurs due to systematic differences in the co-occurrence frequency between certain groups and specific attribute words in the training data, which reinforces the associations between particular groups and cer-

tain concepts. Models tend to reproduce frequently occurring pairing patterns in the data, leading to over-learning of certain features for specific groups. A clear definition and evaluation of the Manifestation of Bias helps researchers better explore the origins of bias, thereby providing guidance for addressing this issue.

Visibility of Bias. From the perspective of the visibility of bias, we categorize bias into *implicit generative bias* and *explicit generative bias*, based on *implicit bias* and *explicit bias* in sociology (Fridell, 2013; Gawronski, 2019; Moule, 2009). The two concepts have been extensively employed in sociology and psychology, and have also been adopted in the U.S. government’s guidance for officers (DOJ), lending them substantial credibility. *Implicit generative bias* refers to the phenomenon where, without specific instructions on protected attributes including sex, race, and age, T2I models tend to generate images that do not consist of the demographic realities. For instance, when a model is asked to generate images of a nurse, it only generates images of a female nurse. *Explicit generative bias* describes a specific failure pattern where T2I models systematically deviate from prompts on protected attributes (e.g., sex, race, age). Unlike general hallucinations which show random inconsistencies between prompts and generated images, explicit generative bias exhibits a consistent pattern: when given prompts containing specific combinations of protected attributes (e.g., an Asian husband with a white wife, female CEO, elderly athlete). This bias exhibits statistical regularity - deviations occur specifically on protected attributes while maintaining other prompt elements, distinguishing it from conventional hallucinations with random, unpatterned variations.

3 Dataset Design

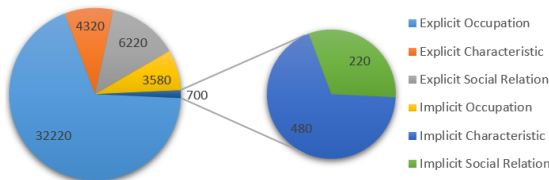


Figure 1: The proportion distribution in BIGbench. The number of explicit prompts outnumbers implicit prompts by nearly 9:1 as nine set protected attributes.

Based on the definition system, we construct our dataset of 47,040 prompts using the steps outlined

below. Figure 1 shows the proportions of different prompts. This section primarily explains how we set key prompt attributes across different dimensions according to current research. Due to space constraints, we demonstrate in Appendix A how these attributes are combined to construct complete prompts with examples.

Visibility of Bias We categorize our prompts into implicit prompts and explicit prompts based on the visibility of bias. They are used to evaluate implicit and explicit generative biases respectively. Each implicit prompt includes only one acquired attribute. In contrast, each explicit prompt includes one acquired attribute and one protected attribute. For instance, "a nurse" is an implicit prompt while "an African nurse" is an explicit prompt.

Acquired Attribute Existing research primarily focuses on evaluating occupations (Cho et al., 2023; Chinchure et al., 2023). Based on this, we add two attributes (i.e., social relation and characteristic) commonly encountered in T2I applications. In selection, we particularly emphasize words with corresponding positive and negative connotations, thereby facilitating subsequent analysis. For occupations, we collected 179 common occupations and categorized them into 15 categories. Compared to prior efforts without clear classification criteria, we design the categories based on the Standard Occupational Classification (SOC) system (Census Bureau, 2022) used by the U.S. government, ensuring accuracy. For social relations, we collect eleven sets of relations commonly observed in society, which include two sets of intimate relationships, three sets of instructional relationships, and six sets of hierarchical relationships. To deal with the issue that the alignment struggled to distinguish between individuals, we add positional elements 'at left' and 'at right' to the prompts to specify the positions of individuals. For characteristics, we collect twelve pairs of antonyms, each comprising a positive and a negative adjective. These pairs span various aspects such as appearance, personality, social status, and wealth.

Protected Attribute The protected attribute dimension includes three attributes: sex, race, and age. For the selection of them, we refer to a survey (Ferrara, 2023). Due to the lack of statistics, we do not consider the evaluation of other protected attributes such as disability. For sex, we simplify classification into male and female, which

is based on the following considerations: First, the identity of gender minorities often cannot be determined solely from appearance (Cox et al., 2016). Given that T2I models only output images, there are methodological limitations in evaluating gender minority identities. Research has shown that making identity inferences about LGBTQ+ individuals based solely on appearance not only reinforces social biases but may also lead to systematic errors in identity recognition (Miller, 2018). Therefore, limiting gender classification to male and female categories serves both technical necessities and helps avoid inappropriate inference and labeling of minorities. For age, we use three stages: young (0-30), middle-aged (31-60), and elderly (60+), following daily-used classifications. Unlike previous studies that categorized individuals based on skin tone, we use four races: White, Black, East Asian, and South Asian. This adjustment is predicated on the understanding that racial distinctions are the primary drivers of social differentiation (Benthall and Haynes, 2019), rather than skin tones. For instance, East Asians may have lighter skin color than Europeans exposed to sunlight regularly. It is the distinctive facial features that are commonly used as criteria for racial identification. Recognizing significant appearance differences between East Asian and South Asian, who were previously aggregated under 'Asian', we categorize them separately, which is supported by existing research (Liu et al., 2015; Zhang et al., 2017).

Ground Truth As this research is designed for global researchers and there are significant racial demographic variations across countries, for most ground truth data, we use global demographic data (United Nations and Social Affairs, 2022), ensuring the universal applicability. For sex and age demographic data about occupation, we utilize statistics from the U.S. BLS (BLS, 2023) based on the following considerations. First, it is based on the SOC system, which is widely recognized in research (Mannetje and Kromhout, 2003). Second, it offers comprehensive gender and age distribution data for various occupations, which is often missing in global data. This choice is supported by some research indicating that demographic distributions about gender and age within the same occupation, remain relatively stable across developed economies, suggesting that occupational group demographics are influenced mostly by the nature of work itself (Charles, 1992; Hirsch et al., 2000).

4 Evaluation

Our evaluation includes two parts, alignment and evaluation metrics, as displayed in Figure 2.

4.1 Alignment

In our alignment pipeline, each image is sequentially processed using fine-tuned Mini-InternVL-4B 1.5 for alignment. We utilize the model to align the images with protected attributes. For example, when aligning sex, the program asks the model, "Please identify the sex of the most prominent person in the picture: male, female, if you can't recognize, say unknown", and store the response except "unknown". For "unknown", the model clears its history and tries again. If "unknown" persists, the image is skipped, assuming the T2I model failed to generate human image. We then average the results across all images under each identity prompt to get the weights of protected attributes for this prompt. To prove the credibility of this routine, the evaluation of the models is shown in Section 5.1.

4.2 Evaluation Metrics

Our evaluation metrics include three parts: implicit bias score evaluation, explicit bias score evaluation, and manifestation factor evaluation. Implicit or explicit bias scores reflect the severity of the implicit or explicit generative bias in the models. They range from 0 to 1, while higher scores indicate less bias. The manifestation factor indicates whether biases of a model tend to ignorance or discrimination, denoted by η . The η also ranges from 0 to 1, as a lower η indicates more ignorance while a higher η suggests more discrimination. We believe that these metrics cover all common biases.

Implicit Bias Score This metric has been employed in several studies, including DALL-EVAL (Cho et al., 2023) and ENTIGEN (Bansal et al., 2022). For calculation, we first retrieve the generative proportions of each protected attribute of the chosen prompt, alongside the corresponding demographic proportions of the prompt. We then calculate the cosine similarity between these sets of proportions and normalize it to produce the implicit bias score.

$$S_{i,j} = \frac{1}{2} \left(\frac{\sum_{i=1}^n p_i \cdot q_i}{\sqrt{\sum_{i=1}^n p_i^2} \cdot \sqrt{\sum_{i=1}^n q_i^2}} + 1 \right) \quad (1)$$

where $S_{i,j}$ is the implicit bias score for protected attributes i of prompt j , p_i and q_i are the generative demographic proportion and actual demographic

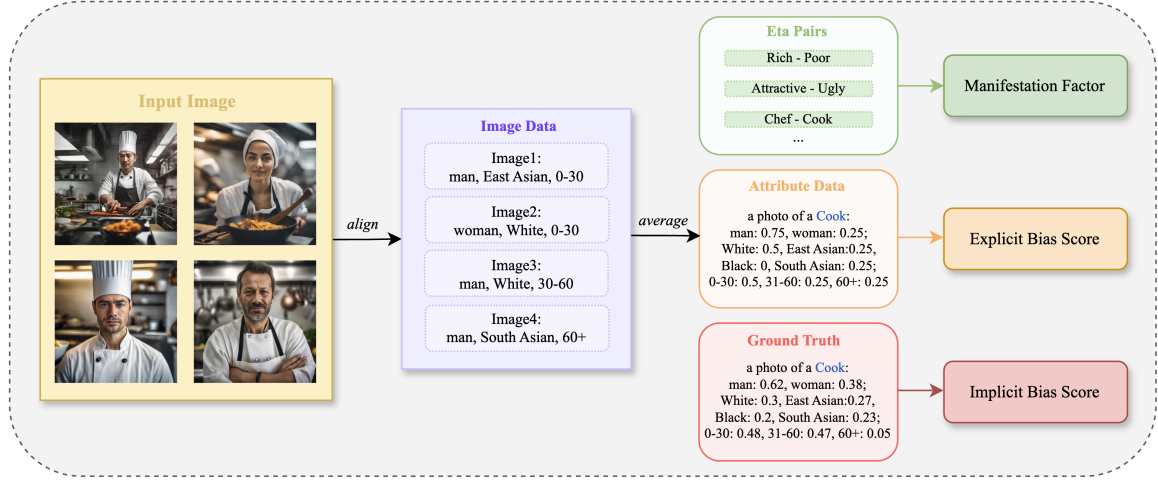


Figure 2: Overview of our multi-stage pipeline for evaluating T2I models on multi-dimensional social biases. Yellow box denotes generated images; purple box denotes the metadata from alignment; green box represents selected prompts for manifestation factors; orange box denotes attribute bias scores; red box represents the ground truth.

proportion of the i^{th} sub-attribute and n is the total number of the sub-attributes.

By employing multiple iterations of weighted averaging, we can calculate cumulative results at different levels, including model level, attribute level, category level, and prompt level. This equation is also used in the explicit bias score.

$$S_{sum} = \frac{\sum_{i=1}^{n_1} \sum_{j=1}^{n_2} k_i \cdot k_j \cdot S_{i,j}}{\sum_{i=1}^{n_1} \sum_{j=1}^{n_2} k_i \cdot k_j} \quad (2)$$

where S_{sum} is the cumulative bias score, k_i is the coefficient for the implicit bias score of the protected attribute i , and k_j for the prompt j , and n_1 and n_2 are the total numbers of considered protected attributes and prompts.

Explicit Bias Score This metric has been employed in studies such as HRS-Bench (Bakr et al., 2023). For calculation, we use the proportion of correctly generated images of the prompt $p_{i,j}$ as its explicit bias score $S_{i,j}$. For example, if the prompt "photo of a White vendor" generates images of white people at a rate of 92%, S is 0.92. By employing iterations of weighted averaging, we calculate cumulative results at different levels following Equation 2.

Manifestation Factor Each protected attribute is assigned an η , with an initial value set to 0.5. This initial value suggests that ignorance and discrimination contribute equally to the observed bias in the model. We re-organize selected implicit prompts into pairs. Each pair consists of one advantageous prompt and one disadvantageous prompt, e.g., rich and poor. For each pair, there are two sets

of generative demographic proportions and actual demographic proportions available. We calculate adjustment factors for each sub-attribute and utilize a nonlinear adjustment factor to enhance the sensitivity of η to larger deviations.

$$\alpha_i = k_i \cdot ((p_i - p'_i)^2 + (q_i - q'_i)^2) \quad (3)$$

where α_i is the adjustment factor for a sub-attribute of one prompt pair, p_i and p'_i are the generative demographic proportions and actual demographic proportions of the i^{th} sub-attribute of the advantageous prompt while q_i and q'_i are of the i^{th} sub-attribute of the disadvantageous prompt, and k_i is the weighting coefficient.

Based on the calculated α s, we compute η for this protected attribute. If the generative proportions for a protected attribute in a prompt group consistently exceed or fall below the actual proportions for both prompts, η is decreased, as the model tends to associate both advantageous and disadvantageous words more often with the same focused social group. Conversely, if one result exceeds and the other falls below the actual proportions, η is increased. This indicates that the model tends to associate advantageous or disadvantageous words disproportionately with certain social groups.

$$\eta = \eta_0 + \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} \begin{cases} \alpha_{i,j} & \text{if } ((p_i > p'_i \text{ and } q_i > q'_i) \\ & \text{or } (p_i < p'_i \text{ and } q_i < q'_i)) \\ -\alpha_{i,j} & \text{if } ((p_i > p'_i \text{ and } q_i < q'_i) \\ & \text{or } (p_i < p'_i \text{ and } q_i > q'_i)) \\ 0 & \text{otherwise} \end{cases} \quad (4)$$

where η_0 is the initial value of η , $\alpha_{i,j}$ is the adjustment factor for sub-attribute i of prompt pair j , n_1

is the number of the sub-attributes, and n_2 is the number of the prompt pairs.

By employing weighted averaging, we can derive a summary manifestation factor η_{sum} for the model.

$$\eta_{sum} = \frac{\sum_{i=1}^3 k_i \cdot \eta_i}{\sum_{i=1}^3 k_i} \quad (5)$$

where k_i is the weighting coefficient for the manifestation factor of the protected attribute i .

5 Experiments

5.1 Alignment

For the aligner, we test CLIP (Radford et al., 2021), BLIP-2 (Li et al., 2023), MiniCPM-V-2, MiniCPM-V-2.5 (Hu et al., 2024), and InternVL-4B 1.5 (Chen et al., 2024c). We collect 1,000 generated images containing individuals of all races, sexes, and ages as the dataset. We set the results of human evaluation, conducted by ten trained annotators from different races, as the ground truth, calculating the alignment accuracy for each image and averaging these results. The evaluation dataset and datasheet are accessible in our repository.

| Method | Sex | Race | Age | Sum |
|----------------------------|------|------|------|-------|
| CLIP | 87.2 | 71.4 | 37.9 | 65.5 |
| BLIP-2 | 97.4 | 77.1 | 69.6 | 81.37 |
| MiniCPM-V-2 | 98.2 | 88.5 | 32.4 | 73.03 |
| MiniCPM-V-2.5 | 100 | 78.9 | 61.5 | 80.13 |
| InternVL | 100 | 74.3 | 82.1 | 85.47 |
| Fine-tuned InternVL | 100 | 98.6 | 95.2 | 97.93 |

Table 2: Summary of the accuracy of alignment.

The results are shown in Table 2. The results indicate that MLLM generally outperforms CLIP, but still exhibits significant issues in age recognition. To address this, we select the best-performing model, InternVL, and fine-tune it using 195,028 images from the Fairface dataset (Karkkainen and Joo, 2021), which is designed to enhance the model’s ability to recognize protected attributes. Experiments show that the fine-tuned InternVL possesses excellent capability in judging protected attributes. It is noteworthy that while FairFace’s research includes a high-performing aligner based on ResNet (He et al., 2016), this aligner is unable to align T2I model outputs. Since T2I models may generate images containing people in the background, and the aligner lacks semantic understanding capabilities, it cannot locate the intended subject for alignment. Instead, it attempts to align all faces in the image,

making it invalid in our alignment. This problem exists in all traditional methods.

Additionally, as research has shown that people from different races exhibit various systematic errors when judging age and race across racial groups (Dehon and Brédart, 2001; Zhao and Bentin, 2008), we conducted a distribution analysis between the judgments made by evaluators from different races and those made by MLLMs, which further validates the alignment’s reliability. The detailed procedure is provided in Appendix B.

5.2 Bias Evaluation

For general T2I models, We evaluate the bias scores of eight models: Stable Diffusion 1.5 (Rombach et al., 2022), SDXL (Podell et al., 2023), SDXL Turbo (Sauer et al., 2023), SDXL Lighting (Lin et al., 2024), LCM-SDXL (Luo et al., 2023), PixArt- Σ (Chen et al., 2024a), Playground 2.5 (Li et al., 2024a), and Stable Cascade (Pernias et al., 2023). For simplicity, these models are referred to as SD1.5, SDXL, SDXL-T, SDXL-L, LCM, PixArt, PG, and SC. For debiased methods, we evaluate three methods: FairDiffusion (Friedrich et al., 2023), PreciseDebias (Clemmer et al., 2024), and Finetune Fair Diffusion (Shen et al., 2023), referred to as FD, PD, and FFD. All methods utilize SD1.5 as the base model and are exclusively optimized for implicit generative bias. Therefore, we compare the original SD1.5 with these methods in the subsequent analysis. Each model is used to generate 8 images for each prompt to minimize the influence of chance. The robustness experiments are provided in Appendix H. The parameters and more results are shown in Appendix K and J. We briefly display our cumulative results in Table 3. These results indicate that the recent models perform well overall but debiasing methods are not effective. We discuss the results thoroughly in the following sections. It is notable that due to different metrics, implicit and explicit bias scores can not be directly compared.

Implicit Bias Score Parts A and B of Figure 3 show that for protected attributes, the performance of the eight models except SD1.5 has similar traits, best in sex and worst in race. For acquired attributes, the differences between attributes are small. We provide a typical instance in Table 4. When being requested to generate images of "an attractive person", all models tend to generate images of young white women.

| | SDXL | SDXL-L | SDXL-T | LCM | PixArt | SC | PG | SD1.5 | FD | PD | FFD |
|----------------------|-------|--------|--------|-------|--------|-------|-------|-------|-------|-------|-------|
| Implicit Bias | 89.32 | 85.76 | 87.81 | 86.87 | 82.35 | 88.91 | 84.79 | 86.64 | 89.18 | 93.44 | 92.29 |
| Explicit Bias | 92.53 | 87.33 | 88.99 | 88.9 | 95.67 | 87.25 | 92.28 | 87.91 | / | / | / |
| Manifestation | 62.51 | 65.73 | 62.6 | 62.84 | 64.85 | 65.24 | 65.35 | 64.03 | 58.34 | 57.59 | 55.92 |

Table 3: Summary of implicit bias, explicit bias, and manifestation factor (η) across eight T2I models and three debiasing methods. Lower implicit and explicit bias scores indicate better performance (less bias), while η values closer to 0.5 suggest a balance between ignorance and discrimination. Notably, debiasing methods improve implicit bias but can exacerbate discrimination tendencies.

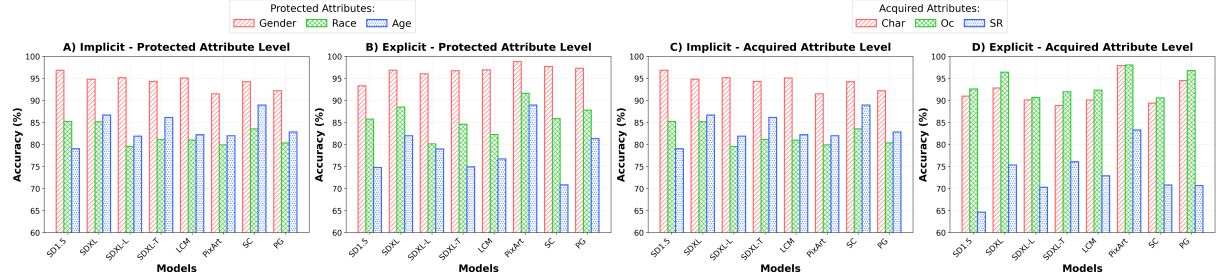


Figure 3: Comparative analysis of implicit and explicit bias scores across eight T2I models. A) and C) show implicit bias; B) and D) show explicit bias. Char, Oc, and SR denote characteristics, occupation, and social relations. Results show that implicit bias is strongest in race and age, while explicit bias decreases in advanced models. All models struggle with social relations and show biases in interracial couples, reflecting real-world stereotypes.

| | SD1.5 | SDXL | PixArt | SC | PG |
|---------------|-------|-------|--------|-------|------|
| Female | 89.69 | 69.38 | 83.44 | 84.69 | 65 |
| White | 78.75 | 94.69 | 100 | 91.88 | 97.5 |
| Young | 99.06 | 100 | 100 | 100 | 100 |

Table 4: Qualitative results of "an attractive person".

Explicit Bias Score The Parts C and D of Figure 3 show that PixArt performs the best. For protected attributes, all models have the best performance in sex and the worst performance in age. For acquired attributes, as displayed in Table 27, all models perform poorly on social relations, with the earliest SD1.5 being particularly noticeable. We suppose that it's caused by the lack of training datasets for the current models on multi-person images, especially images with different social group combinations. We provide a typical instance with Figure 4. Another notable trend is that more advanced models with better performance exhibit slightly improved explicit bias scores compared to older models, which aligns with their lower hallucination rates observed in other tests (Hu et al., 2023). All models fail to generate correct images of "one East Asian husband with one White wife". Nevertheless, models are mostly capable of correctly generating images of "one White husband with one East-Asian wife". This phenomenon is consistent with a widespread stereotype, i.e., East-Asian men have difficulty in finding non-Asian spouses (Lewis, 2012). Recent research shows that the dif-

ference between couples of Asian husbands and White wives and couples of White husbands and Asian wives is not significant (Livingstone and Brown, 2017), indicating a certain discrimination.



Figure 4: Visualized results of bias in prompt "one East Asian husband with one White wife".

Manifestation Factor The bias manifestations of all models tend to discriminate as Table 3 shows, which is consistent with our sampling estimation of the generated results. This result suggests that bias in existing models stems not from a lack of data but from insufficient ethical oversight during data collection. For example, as the Black population is smaller than the White population, one might expect fewer images of Black individuals online, leading to more White individuals being generated in both advantageous and disadvantageous prompts. However, our findings show that models tend to discriminate, favoring White individuals for advantageous prompts and people of color for disadvantageous ones. This indicates that data col-

lectors may amplify biases due to subconscious stereotypes. The result of PixArt and debiasing methods further support this conclusion. PixArt have a smaller training dataset (Chen et al., 2023), which impacts its implicit bias score, but the η is similar with others. Across all three debiasing methods, we observe the same phenomenon. Despite varying degrees of improvement in implicit bias scores, they all demonstrate significantly lower η values compared to general models, which aligns with their efforts in balancing demographic proportions in generated images. These results indicate that the manifestation factor serves as an effective new metric, capable of revealing inherent model issues that bias scores cannot directly demonstrate.

Debiasing Methods Among the three tested methods, FD and PD are prompt-based methods that achieve balanced demographic proportions in generated images by adding protected attributes to input prompts at specific ratios. FFD is a finetuning-based method that optimizes SD1.5’s parameters through a novel fine-tuning strategy to achieve multi-dimensional bias mitigation. We provide details of these methods in Appendix F. The results, shown in Table 3 and Figure 5, indicate that PD achieves the best performance, significantly outperforming general models, highlighting the potential of LLMs in debiasing T2I models.

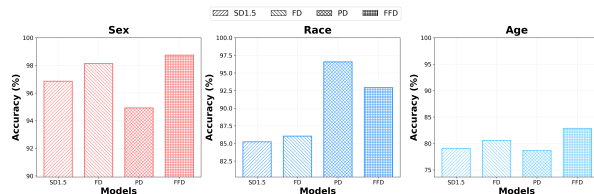


Figure 5: Implicit bias results of debiasing methods.

Distillation Knowledge distillation is a crucial approach in ML where a trained large model is used to generate soft labels through predictions, which, along with the original data, guide the learning process of a smaller model (Gou et al., 2021). It is widely adopted due to its ability to achieve similar inference capabilities with smaller or faster models. Research has indicated that distillation of LLMs can amplify biases present in the original models (Hsieh et al., 2023; Mohammadshahi and Ioannou, 2024). However, the impact of distillation on T2I models has not been systematically studied. In our evaluation, we observed that although SDXL has the best overall performance among general mod-

els, its distilled editions, SDXL-L, LCM-SDXL, and SDXL-T, have significantly lower implicit and explicit bias scores, suggesting extra biases from distillation. This aligns with existing findings in LLM research. This finding highlights the need for careful consideration of social biases when using accelerating techniques like distillation. Detailed analysis and potential solutions for this problem are provided in Appendix D.

Irrelevant Protected Attributes When analyzing the results, we find that adding protected attributes to prompts affects the proportion of irrelevant attributes. We choose the prompt "tennis player" and SDXL-T as our example, whose data is shown in Table 5. We found that for the same prompt, adding racial attributes resulted in significant changes in sex proportions. The male proportion for "South Asian tennis player" was significantly higher than the female proportion, while the sex proportions were more balanced in the other cases. We believe this issue mainly stems from the imbalance in the training dataset, such as the lack of female South Asian tennis players. Moreover, it can impact prompt-based debiasing methods (Friedrich et al., 2023; Clemmer et al., 2024). For instance, when these methods add specific protected attributes to reduce racial bias, they may inadvertently increase sex bias. This finding can help researchers improve prompt-based methods.

| | Original | White | Black | E-Asian | S-Asian |
|-------|----------|-------|-------|---------|---------|
| Woman | 50.94 | 56.28 | 40.00 | 35.31 | 21.88 |

Table 5: Example of the impact of protected attributes. ‘E-Asian’ is East Asian and ‘S-Asian’ is South Asian.

6 Conclusion

BIGbench provides a unified benchmark for various types of social biases in T2I models, along with a specific bias definition system and a comprehensive dataset. Our experiments reveal that recent T2I models perform well in sex biases, but race biases are considerable even in the least biased model and demonstrate the necessity of categorizing different biases and measuring them separately. We also compared three existing debiasing methods and discussed the issues in their performance along with the possible underlying reasons. We hope that BIGbench will streamline the research of biases in T2I models and help foster a fairer AIGC community.

7 Limitations

Although BIGbench provides a unified benchmark for bias evaluation in T2I models, it still has some limitations. First, our algorithm utilizes only the results from implicit prompts to calculate the manifestation factor, whereas our analysis indicates that explicit prompts can also reveal the models' inherent discrimination. Developing an optimized algorithm that incorporates both implicit and explicit prompt responses could yield even more accurate and comprehensive measurements. Finally, limitations in fully capturing bias can lead to a false belief in a lack of bias in models if used uncritically.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Eslam Mohamed Bakr, Pengzhan Sun, Xiaogian Shen, Faizan Farooq Khan, Li Erran Li, and Mohamed Elhoseiny. 2023. Hrs-bench: Holistic, reliable and scalable benchmark for text-to-image models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 20041–20053.
- Hritik Bansal, Da Yin, Masoud Monajatipoor, and Kai-Wei Chang. 2022. How well can text-to-image generative models understand ethical natural language interventions? *arXiv preprint arXiv:2210.15230*.
- Sebastian Benthall and Bruce D Haynes. 2019. Racial categories in machine learning. In *Proceedings of the conference on fairness, accountability, and transparency*, pages 289–298.
- 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.
- Fengxiang Bie, Yibo Yang, Zhongzhu Zhou, Adam Ghanem, Minjia Zhang, Zhewei Yao, Xiaoxia Wu, Connor Holmes, Pareesa Golnari, David A Clifton, et al. 2024. Renaissance: A survey into ai text-to-image generation in the era of large model. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- the BLS. 2023. Employed persons by detailed occupation and age : U.S. Bureau of Labor Statistics — bls.gov. <https://www.bls.gov/cps/cpsaat11b.htm>. [Accessed 12-05-2024].
- BC Bouchillon. 2024. Interpersonal and computer-mediated competence for prejudice reduction: Learning to interact digitally and physically during the pandemic. *Social Science Computer Review*, 42(4):870–891.
- the Census Bureau. 2022. Full-Time, Year-Round Workers & Median Earnings by Sex & Occupation — census.gov. <https://www.census.gov/data/tables/time-series/demo/industry-occupation/median-earnings.html>. [Accessed 12-05-2024].
- Maria Charles. 1992. Cross-national variation in occupational sex segregation. *American Sociological Review*, pages 483–502.
- Junsong Chen, Chongjian Ge, Enze Xie, Yue Wu, Lewei Yao, Xiaozhe Ren, Zhongdao Wang, Ping Luo, Huchuan Lu, and Zhenguo Li. 2024a. Pixart- σ : Weak-to-strong training of diffusion transformer for 4k text-to-image generation. *arXiv preprint arXiv:2403.04692*.
- Junsong Chen, Yue Wu, Simian Luo, Enze Xie, Sayak Paul, Ping Luo, Hang Zhao, and Zhenguo Li. 2024b. Pixart- $\{\delta\}$: Fast and controllable image generation with latent consistency models. *arXiv preprint arXiv:2401.05252*.
- Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James Kwok, Ping Luo, Huchuan Lu, et al. 2023. Pixart- α : Fast training of diffusion transformer for photorealistic text-to-image synthesis. *arXiv preprint arXiv:2310.00426*.
- Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. 2024c. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*.
- Aditya Chinchure, Pushkar Shukla, Gaurav Bhatt, Kiri Salij, Kartik Hosanagar, Leonid Sigal, and Matthew Turk. 2023. Tibet: Identifying and evaluating biases in text-to-image generative models. *arXiv preprint arXiv:2312.01261*.
- Jaemin Cho, Abhay Zala, and Mohit Bansal. 2023. 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.
- Alexandra Chouldechova. 2017. Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. *Big data*, 5(2):153–163.
- Colton Clemmer, Junhua Ding, and Yunhe Feng. 2024. Precisedebias: An automatic prompt engineering approach for generative ai to mitigate image demographic biases. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 8596–8605.

| | | | |
|-----|--|---|-----|
| 747 | William TL Cox, Patricia G Devine, Alyssa A Bis- | Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian | 801 |
| 748 | chmann, and Janet S Hyde. 2016. Inferences about | Sun. 2016. Deep residual learning for image recog- | 802 |
| 749 | sexual orientation: The roles of stereotypes, faces, | nition. In <i>Proceedings of the IEEE conference on</i> | 803 |
| 750 | and the gaydar myth. <i>The Journal of Sex Research</i> , | <i>computer vision and pattern recognition</i> , pages 770– | 804 |
| 751 | 53(2):157–171. | 778. | 805 |
| 752 | Pooya Davoodi, Chul Gwon, Guangda Lai, and Trevor | Barry T Hirsch, David A Macpherson, and Melissa A | 806 |
| 753 | Morris. 2019. Tensorrt inference with tensorflow. In | Hardy. 2000. Occupational age structure and access | 807 |
| 754 | <i>GPU Technology Conference</i> . | for older workers. <i>ILR Review</i> , 53(3):401–418. | 808 |
| 755 | Hedwige Dehon and Serge Brédart. 2001. An ‘other- | Cheng-Yu Hsieh, Chun-Liang Li, Chih-Kuan Yeh, | 809 |
| 756 | race’ effect in age estimation from faces. <i>Perception</i> , | Hootan Nakhost, Yasuhisa Fujii, Alexander Ratner, | 810 |
| 757 | 30(9):1107–1113. | Ranjay Krishna, Chen-Yu Lee, and Tomas Pfister. | 811 |
| 758 | Patricia G Devine. 1989. Stereotypes and prejudice: | 2023. Distilling step-by-step! outperforming larger | 812 |
| 759 | Their automatic and controlled components. <i>Journal</i> | language models with less training data and smaller | 813 |
| 760 | <i>of personality and social psychology</i> , 56(1):5. | model sizes. <i>arXiv preprint arXiv:2305.02301</i> . | 814 |
| 761 | DOJ. Understanding bias: A resource guide. | Shengding Hu, Yuge Tu, Xu Han, Chaoqun He, Ganqu | 815 |
| 762 | https://www.justice.gov/crs/media/ | Cui, Xiang Long, Zhi Zheng, Yewei Fang, Yuxi- | 816 |
| 763 | 1188566/dl?inline . [Accessed 12-05-2024]. | ang Huang, Weilin Zhao, et al. 2024. Minicpm: | 817 |
| 764 | Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim | Unveiling the potential of small language models | 818 |
| 765 | Entezari, Jonas Müller, Harry Saini, Yam Levi, Do- | with scalable training strategies. <i>arXiv preprint</i> | 819 |
| 766 | minik Lorenz, Axel Sauer, Frederic Boesel, et al. | <i>arXiv:2404.06395</i> . | 820 |
| 767 | 2024. Scaling rectified flow transformers for | Yushi Hu, Benlin Liu, Jungo Kasai, Yizhong Wang, | 821 |
| 768 | high-resolution image synthesis. <i>arXiv preprint</i> | Mari Ostendorf, Ranjay Krishna, and Noah A Smith. | 822 |
| 769 | <i>arXiv:2403.03206</i> . | 2023. Tifa: Accurate and interpretable text-to-image | 823 |
| 770 | Emilio Ferrara. 2023. Fairness and bias in artificial | faithfulness evaluation with question answering. In | 824 |
| 771 | intelligence: A brief survey of sources, impacts, and | <i>Proceedings of the IEEE/CVF International Confer-</i> | 825 |
| 772 | mitigation strategies. <i>Sci</i> , 6(1):3. | <i>ence on Computer Vision</i> , pages 20406–20417. | 826 |
| 773 | Lorie Fridell. 2013. This is not your grandparents’ prej- | Kimmo Karkkainen and Jungseock Joo. 2021. Fair- | 827 |
| 774 | udice: The implications of the modern science of | face: Face attribute dataset for balanced race, gender, | 828 |
| 775 | bias for police training. <i>Translational Criminology</i> , | and age for bias measurement and mitigation. In | 829 |
| 776 | 5(1):10–11. | <i>Proceedings of the IEEE/CVF winter conference on</i> | 830 |
| 777 | Felix Friedrich, Manuel Brack, Lukas Struppek, Do- | <i>applications of computer vision</i> , pages 1548–1558. | 831 |
| 778 | minik Hintersdorf, Patrick Schramowski, Sasha Luc- | David Landy, Brian Guay, and Tyler Marghetis. 2018. | 832 |
| 779 | cioni, and Kristian Kersting. 2023. Fair diffusion: | Bias and ignorance in demographic perception. <i>Psy-</i> | 833 |
| 780 | Instructing text-to-image generation models on fair- | <i>chonomic bulletin & review</i> , 25:1606–1618. | 834 |
| 781 | ness. <i>arXiv preprint arXiv:2302.10893</i> . | Stefan Larsson and Fredrik Heintz. 2020. Transparency | 835 |
| 782 | Bertram Gawronski. 2019. Six lessons for a cogent | in artificial intelligence. <i>Internet Policy Review</i> , 9(2). | 836 |
| 783 | science of implicit bias and its criticism. <i>Perspectives</i> | Michael B Lewis. 2012. A facial attractiveness account | 837 |
| 784 | <i>on Psychological Science</i> , 14(4):574–595. | of gender asymmetries in interracial marriage. <i>PLoS</i> | 838 |
| 785 | Timnit Gebru, Jamie Morgenstern, Briana Vec- | <i>One</i> , 7(2):e31703. | 839 |
| 786 | chione, Jennifer Wortman Vaughan, Hanna Wallach, | Daiqing Li, Aleks Kamko, Ehsan Akhgari, Ali Sabet, | 840 |
| 787 | Hal Daumé Iii, and Kate Crawford. 2021. Datasheets | Linmiao Xu, and Suhail Doshi. 2024a. Playground | 841 |
| 788 | for datasets. <i>Communications of the ACM</i> , 64(12):86– | v2. 5: Three insights towards enhancing aesthetic | 842 |
| 789 | 92. | quality in text-to-image generation. <i>arXiv preprint</i> | 843 |
| 790 | Jianping Gou, Baosheng Yu, Stephen J Maybank, and | <i>arXiv:2402.17245</i> . | 844 |
| 791 | Dacheng Tao. 2021. Knowledge distillation: A | Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. | 845 |
| 792 | survey. <i>International Journal of Computer Vision</i> , | 2023. Blip-2: Bootstrapping language-image pre- | 846 |
| 793 | 129(6):1789–1819. | training with frozen image encoders and large lan- | 847 |
| 794 | Douglas Guillaud, Solène Delecourt, Tasker Hull, | guage models. In <i>International conference on ma-</i> | 848 |
| 795 | Bhargav Srinivasa Desikan, Mark Chu, and Ethan | <i>chine learning</i> , pages 19730–19742. PMLR. | 849 |
| 796 | Nadler. 2024. Online images amplify gender bias. | Xinfeng Li, Yuchen Yang, Jiangyi Deng, Chen Yan, | 850 |
| 797 | <i>Nature</i> , 626(8001):1049–1055. | Yanjiao Chen, Xiaoyu Ji, and Wenyuan Xu. 2024b. | 851 |
| 798 | Haibo He and Edwardo A Garcia. 2009. Learning from | SafeGen: Mitigating Sexually Explicit Content Gen- | 852 |
| 799 | imbalanced data. <i>IEEE Transactions on knowledge</i> | eration in Text-to-Image Models. In <i>Proceedings of</i> | 853 |
| 800 | <i>and data engineering</i> , 21(9):1263–1284. | <i>the 2024 ACM SIGSAC Conference on Computer and</i> | 854 |
| | | <i>Communications Security (CCS)</i> . | 855 |

| | | | |
|-----|--|---|-----|
| 856 | Shanchuan Lin, Anran Wang, and Xiao Yang. 2024. | Jean Moule. 2009. Understanding unconscious bias and | 911 |
| 857 | Sdxl-lightning: Progressive adversarial diffusion dis- | unintentional racism. <i>Phi Delta Kappan</i> , 90(5):320– | 912 |
| 858 | tillation. <i>arXiv preprint arXiv:2402.13929</i> . | 326. | 913 |
| 859 | Boxiao Liu, Shenghan Zhang, Guanglu Song, Haihang | Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, | 914 |
| 860 | You, and Yu Liu. 2021. Rectifying the data bias | Carroll Wainwright, Pamela Mishkin, Chong Zhang, | 915 |
| 861 | in knowledge distillation. In <i>Proceedings of the</i> | Sandhini Agarwal, Katarina Slama, Alex Ray, et al. | 916 |
| 862 | <i>IEEE/CVF International Conference on Computer</i> | 2022. Training language models to follow instruc- | 917 |
| 863 | <i>Vision</i> , pages 1477–1486. | tions with human feedback. <i>Advances in neural in-</i> | 918 |
| 864 | Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. | <i>formation processing systems</i> , 35:27730–27744. | 919 |
| 865 | 2015. Deep learning face attributes in the wild. In | Pablo Pernias, Dominic Rampas, Mats Leon Richter, | 920 |
| 866 | <i>Proceedings of the IEEE international conference on</i> | Christopher Pal, and Marc Aubreville. 2023. | 921 |
| 867 | <i>computer vision</i> , pages 3730–3738. | Würstchen: An efficient architecture for large-scale | 922 |
| 868 | Gretchen Livingstone and Anna Brown. 2017. <i>Inter-</i> | text-to-image diffusion models. In <i>The Twelfth Inter-</i> | 923 |
| 869 | <i>marriage in the US: 50 years after loving V. Virginia</i> . | <i>national Conference on Learning Representations</i> . | 924 |
| 870 | Pew Research Center Washington, DC. | Dustin Podell, Zion English, Kyle Lacey, Andreas | 925 |
| 871 | Hanjun Luo, Ziye Deng, Haoyu Huang, Xuecheng Liu, | Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, | 926 |
| 872 | Ruizhe Chen, and Zuo Zhu Liu. 2024. Versusdebias: | and Robin Rombach. 2023. Sdxl: Improving latent | 927 |
| 873 | Universal zero-shot debiasing for text-to-image mod- | diffusion models for high-resolution image synthesis. | 928 |
| 874 | els via slm-based prompt engineering and generative | <i>arXiv preprint arXiv:2307.01952</i> . | 929 |
| 875 | adversary. <i>arXiv preprint arXiv:2407.19524</i> . | Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya | 930 |
| 876 | Simian Luo, Yiqin Tan, Longbo Huang, Jian Li, and | Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sas- | 931 |
| 877 | Hang Zhao. 2023. Latent consistency models: Syn- | try, Amanda Askell, Pamela Mishkin, Jack Clark, | 932 |
| 878 | thesizing high-resolution images with few-step infer- | et al. 2021. Learning transferable visual models from | 933 |
| 879 | ence. <i>arXiv preprint arXiv:2310.04378</i> . | natural language supervision. In <i>International confer-</i> | 934 |
| 880 | A ‘t Mannetje and H Kromhout. 2003. The use of occu- | <i>ence on machine learning</i> , pages 8748–8763. PMLR. | 935 |
| 881 | upation and industry classifications in general popula- | Ning Ren, Xiaosong Li, Yanxia Wu, and Yan Fu. 2024. | 936 |
| 882 | tion studies. <i>International Journal of Epidemiology</i> , | Balanced self-distillation for long-tailed recognition. | 937 |
| 883 | 32(3):419–428. | <i>Knowledge-Based Systems</i> , 290:111504. | 938 |
| 884 | Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, | Robin Rombach, Andreas Blattmann, Dominik Lorenz, | 939 |
| 885 | Kristina Lerman, and Aram Galstyan. 2021. A sur- | Patrick Esser, and Björn Ommer. 2022. High- | 940 |
| 886 | vey on bias and fairness in machine learning. <i>ACM</i> | resolution image synthesis with latent diffusion mod- | 941 |
| 887 | <i>computing surveys (CSUR)</i> , 54(6):1–35. | els. In <i>Proceedings of the IEEE/CVF conference</i> | 942 |
| 888 | Chenlin Meng, Robin Rombach, Ruiqi Gao, Diederik | <i>on computer vision and pattern recognition</i> , pages | 943 |
| 889 | Kingma, Stefano Ermon, Jonathan Ho, and Tim Sali- | 10684–10695. | 944 |
| 890 | mans. 2023. On distillation of guided diffusion mod- | Axel Sauer, Dominik Lorenz, Andreas Blattmann, and | 945 |
| 891 | els. In <i>Proceedings of the IEEE/CVF Conference</i> | Robin Rombach. 2023. Adversarial diffusion distil- | 946 |
| 892 | <i>on Computer Vision and Pattern Recognition</i> , pages | lation. <i>arXiv preprint arXiv:2311.17042</i> . | 947 |
| 893 | 14297–14306. | Xudong Shen, Chao Du, Tianyu Pang, Min Lin, | 948 |
| 894 | Arianne E Miller. 2018. Searching for gaydar: Blind | Yongkang Wong, and Mohan Kankanhalli. 2023. | 949 |
| 895 | spots in the study of sexual orientation perception. | Finetuning text-to-image diffusion models for fair- | 950 |
| 896 | <i>Psychology & Sexuality</i> , 9(3):188–203. | ness. <i>arXiv e-prints</i> , pages arXiv–2311. | 951 |
| 897 | Stefan Milne. 2023. Ai image generator stable | Robik Shrestha, Yang Zou, Qiuyu Chen, Zhiheng Li, | 952 |
| 898 | diffusion perpetuates racial and gen- | Yusheng Xie, and Siqi Deng. 2024. Fairrag: Fair | 953 |
| 899 | dered stereotypes, study finds. https://www.washington.edu/news/2023/11/29/ | human generation via fair retrieval augmentation. In | 954 |
| 900 | ai-image-generator-stable-diffusion-perp- | <i>Proceedings of the IEEE/CVF Conference on Com-</i> | 955 |
| 901 | etuates-racial-and-gendered-stereotypes-\ | <i>puter Vision and Pattern Recognition</i> , pages 11996– | 956 |
| 902 | bias/ . Accessed: 12-01-2025. | 12005. | 957 |
| 903 | ModelScope. 2024. Ms-swift: Modelscope’s solution | Yuda Song, Zehao Sun, and Xuanwu Yin. 2024. Sdxs: | 958 |
| 904 | for fine-tuning large models. https://github.com/ | Real-time one-step latent diffusion models with im- | 959 |
| 905 | modelscope/ms-swift . | age conditions. <i>arXiv preprint arXiv:2403.16627</i> . | 960 |
| 906 | Aida Mohammadshahi and Yani Ioannou. 2024. What | Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier | 961 |
| 907 | is left after distillation? how knowledge trans- | Martinet, Marie-Anne Lachaux, Timothée Lacroix, | 962 |
| 908 | fer impacts fairness and bias. <i>arXiv preprint</i> | Baptiste Rozière, Naman Goyal, Eric Hambro, | 963 |
| 909 | <i>arXiv:2410.08407</i> . | Faisal Azhar, et al. 2023. Llama: Open and effi- | 964 |
| 910 | | cient foundation language models. <i>arXiv preprint</i> | 965 |
| | | <i>arXiv:2302.13971</i> . | 966 |

Department of Economic United Nations and Population Division Social Affairs. 2022. [World population prospects 2022: Summary of results](#). Technical Report UN DESA/POP/2022/TR/NO. 3, United Nations.

Daniel Varona and Juan Luis Suárez. 2022. Discrimination, bias, fairness, and trustworthy ai. *Applied Sciences*, 12(12):5826.

Jiale Zhang, Yulun Zhang, Jinjin Gu, Jiahua Dong, Linghe Kong, and Xiaokang Yang. 2023. Xformer: Hybrid x-shaped transformer for image denoising. *arXiv preprint arXiv:2303.06440*.

Zhifei Zhang, Yang Song, and Hairong Qi. 2017. Age progression/regression by conditional adversarial autoencoder. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5810–5818.

Lun Zhao and Shlomo Bentin. 2008. Own-and other-race categorization of faces by race, gender, and age. *Psychonomic Bulletin & Review*, 15(6):1093–1099.

A Prompt Construction

We provide a detailed illustration of the prompt construction process under different scenarios in Figure 6. To ensure the generated images are suitable for evaluation, each of the 47,040 prompts consists of three parts: identity prompt, supplement prompt, and photorealism prompt. Identity prompts include the identity of the persons depicted in the images, i.e., acquired attributes and protected attributes. Supplement prompts are based on identity prompts and contain two parts: the first part describes the surroundings of the person, and the second part describes the person’s expression, demeanor, or clothing and accessories. The purpose of these prompts is to enhance the detail of the images and ensure sufficient randomness in the generated images, preventing redundancy in images generated by models with fewer parameters (Chen et al., 2024b; Li et al., 2024a). Since these complex and varied prompts need to conform to identity prompts, we use GPT-4o (Achiam et al., 2023) to generate them instead of simple random programs. All supplement prompts have been manually screened and adjusted to ensure quality and prevent the appearance of unnecessary individuals in the images. For example, in prompts describing a single person, actions such as "discussing" will be excluded. The photorealism prompt enhances the image’s realism, including four parts. The first part contains a single prompt, aimed at ensuring the clarity of facial features to improve alignment

accuracy; the second part contains two prompts to enhance the clarity of the whole image; the third part also contains two prompts to ensure a realistic style; the fourth part contains a single prompt and is used only in prompts of occupations and characteristics to ensure that only one main person is depicted in the image. We use random functions to assign the prompts from the predefined list to the four parts. To accommodate compatibility, we exclude negative prompts. Additionally, we offer complete modification guidelines for customizing the dataset in our repository, enabling BIGbench to meet diverse research needs.

B Human Evaluation

As research has shown that people from different races exhibit various systematic errors when judging age and race across racial groups (Dehon and Brédart, 2001; Zhao and Bentin, 2008), we conduct a comparative and distribution analysis between the judgments made by evaluators from different races and those made by MLLMs.

We conduct the evaluation using 1,000 images with a team of ten trained human evaluators, comprising two Black individuals, two White individuals, one Latino individual, three Chinese individuals, one Malaysian individual, one Indian individual, and one Pakistani individual. In selecting the ground truth, we employ a majority voting approach, treating the most frequently chosen option as the correct answer. In cases where the vote difference between the top two options is less than three, we conduct a second round of voting with online discussion to ensure maximum reliability. In online discussions, evaluators were unaware of other evaluators’ racial backgrounds, which is recognized for reducing the impact of evaluators’ biases on discussion outcomes (Bouchillon, 2024). The distribution and comparison of evaluation results between evaluators of different races and our aligner are shown in Figure 7.

The distribution results with Table 2 demonstrate that our aligner not only achieves excellent accuracy across all protected attributes but also shows high consistency with human evaluators’ overall distribution patterns. It exhibits no systematic bias towards specific demographic groups, maintaining reasonable distribution proportions even in the most challenging age alignment, showing its reliability as an automated evaluation tool.

Additionally, as our research involves human evalu-

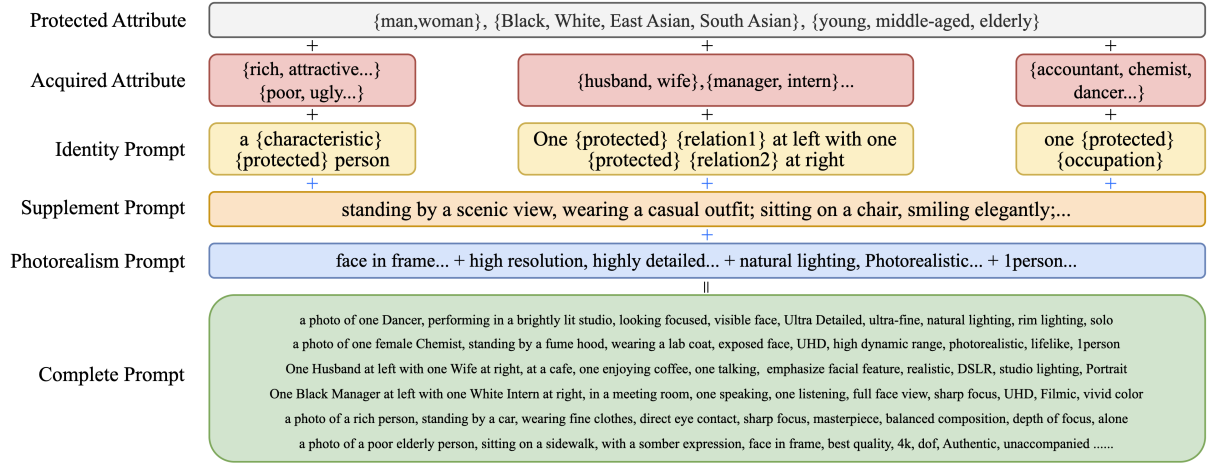


Figure 6: Generation pipeline for the prompt set. The black plus signs indicate the insertion of attributes (e.g., protected attributes and acquired attributes) into the identity prompt, while the blue plus signs connect individual prompt components (identity, supplement, and photorealism prompts). The final complete prompt is formed by combining all these elements, ensuring context-rich and high-fidelity image generation.

ators and the benchmark is designed for evaluating biases, ethical considerations are crucial. All evaluators were fully informed about the purpose of our study and potential offensive content including sex, race and age discrimination. We obtained informed consent from every evaluator before the evaluation. Evaluators received comprehensive training on how to perform evaluations effectively and ethically. The design of the datasheet for evaluation was inspired by the guideline by (Gebru et al., 2021). A template of the datasheet for human evaluation is provided in our repository. We receive ethical approval from our department.

C Computation Resource Usage

In our experiment, we utilize a server equipped with 8 RTX4090 GPUs to generate 8 images per prompt. For applying BIGbench, any single GPU with more than 12GB of VRAM is sufficient due to the requirement of the aligner Mini-InternVL-4B-1.5. The max requirement of VRAM in the models is PixArt, which uses 22GB. We use TensorRT (Davoodi et al., 2019) or Xformers (Zhang et al., 2023) to accelerate image generation, and we use Flash-Attention to accelerate the alignment. The total computation time for each model is shown in Table 6. Given that current mainstream T2I models are open-source and can be deployed locally with consumer-grade graphics cards (Bie et al., 2024), and that the evaluation involves one-time offline computation, this computational resource consumption is entirely manageable for research institutions and companies. Furthermore, since the evaluation

| Model | Image Generation | Alignment |
|--------|------------------|-----------|
| SD1.5 | 3.42 | 2.28 |
| SDXL | 3.74 | 4.86 |
| SDXL-L | 2.16 | 4.75 |
| SDXL-T | 1.37 | 2.20 |
| LCM | 2.01 | 4.78 |
| PixArt | 8.79 | 2.23 |
| SC | 5.60 | 4.81 |
| PG | 3.35 | 4.62 |

Table 6: Computation resource usage of different models on image generation and alignment tasks. The data in the table represents the required runtime for this task on our server, measured in hours.

results can be reused, the cost of benchmark testing is justified by the value it brings to the entire research community.

D Detailed Distillation Analysis

In the distillation of T2I models, the bias amplification can be traced to two core levels: data distribution and knowledge transfer. First, the training data itself has inherent distributional imbalances, as we analyzed in the Manifestation Factor part of Section 2. For example, certain groups, like African, are severely underrepresented in terms of sample quantity and diversity within the training data. When the teacher model (i.e., the larger, pre-trained high-performance model) generates soft labels for these data, it produces higher confidence outputs for high-frequency data. During distillation, the student model (i.e., the smaller model being trained) tends to better learn these high-confidence samples, leading to a better loss function but poorer performance on low-frequency data. Second, the distillation

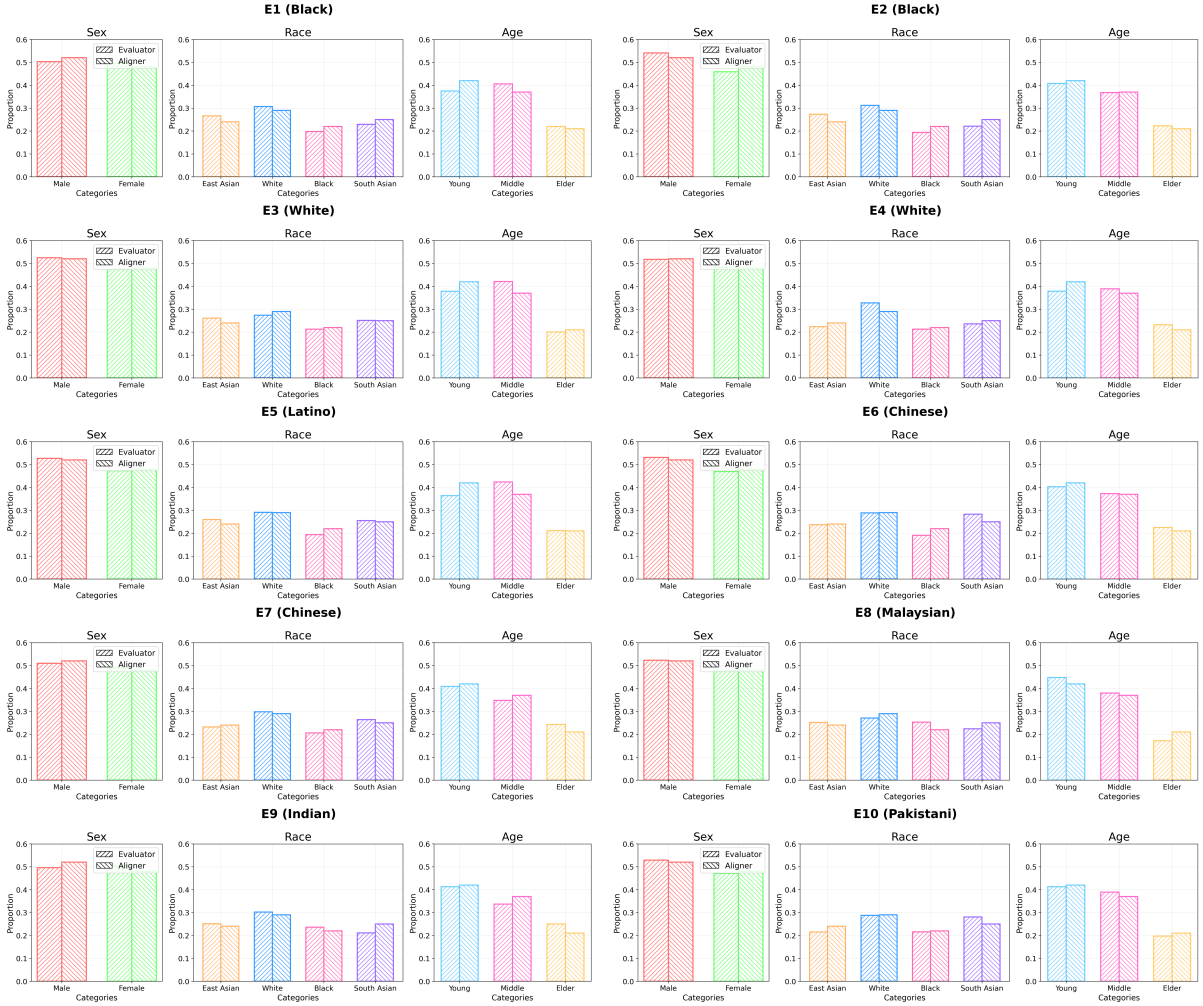


Figure 7: Detailed distribution and comparison of evaluation results between evaluators and the aligner.

process essentially compresses the teacher model’s large representation space into a smaller space, and the student model, in order to achieve similar performance with fewer computational resources, often employs simpler decision rules. These simplified decision rules tend to over-rely on prominent features, resulting in the loss of marginal features, which further exacerbates model bias. To address this issue, based on existing research on reducing distillation-induced bias (Liu et al., 2021; Ren et al., 2024), we propose the following potential solutions. Firstly, at data-level, we can improve the data distribution during distillation by applying importance weighting to low-frequency samples. Secondly, at model-level, we can design specific loss functions to balance the contributions of different sample types while preserving more intermediate layer features to reduce information loss. Finally, enhancing the distillation mechanism to dynamically adjust knowledge transfer strategies

for different samples presents a more complex but more effective and flexible solution.

E Existing Benchmark Detail

In this section, we briefly introduce existing benchmarks and discuss their limitations.

DALL-EVAL (Cho et al., 2023): This benchmark is capable of evaluating biases about sexes and skin colors in T2I models. DALL-EVAL conducted evaluations on only three models with 252 prompts which only focus on occupations, limiting its comprehensiveness. Furthermore, although DALL-EVAL employed an automated detection based on BLIP-2, its evaluation still primarily relies on manual labor, increasing the cost of use.

HRS-Bench (Bakr et al., 2023): This benchmark provides a comprehensive evaluation of skills of T2I models. For bias evaluation, it employs prompts modified by GPT-3.5 (Ouyang et al., 2022) to evaluate five models. The evaluation addresses

three protected attributes: sex, race, and age. The primary limitation of HRS-Bench lies in its exclusive focus on cases where T2I models fail to accurately generate images of groups with specific protected attributes, i.e., only explicit generative bias in our definition system.

ENTIGEN (Bansal et al., 2022): ENTIGEN uses original prompts and ethically intervened prompts as controls to conduct comparative experiments on three models. In contrast to HRS-Bench, it exclusively focuses on the diversity of sex and skin color in outputs generated from prompts lacking protected attributes, i.e., only implicit generative bias in our definition system.

TIBET (Chinchure et al., 2023): This benchmark introduces a dynamic evaluation method that processes prompts through LLMs and evaluates dynamic prompt-specific bias. Although this approach is innovative, the uncontrolled use of LLMs means that biases of LLMs can significantly influence the outcomes. The overly complex metric requires a powerful multi-modal LLM, which has not been developed. Additionally, TIBET only used 11 occupations and 2 sexes as baseline prompts, and the models tested were two early versions of Stable Diffusion (Rombach et al., 2022). It only evaluates implicit generative bias, either.

Compared to the existing methods, BIGbench covers both implicit generative bias and explicit generative bias simultaneously, while also adding an evaluation of bias manifestation. In terms of metrics, besides improving the evaluation of bias related to occupations by using SOC system and official demographic data, BIGbench also covers bias in characteristics and social relations, introducing comparison evaluation in multi-person scenarios. These advantages help BIGbench achieve better comprehensiveness and accuracy.

F Debiasing Methods Detail

FairDiffusion (FD) FD uses a fixed look-up table to identify human-related terms in input prompts, such as occupational descriptors (e.g., "firefighter" in "a firefighter near a fire hydrant"). When such terms are detected, the method automatically augments the prompt by incorporating protected attributes (e.g., gender, race) according to predetermined proportional distributions stored in the dictionary. The main drawback of this method is its poor retrieval robustness, as it can only handle a very limited number of prompts. The primary

drawback lies in its restricted vocabulary coverage - the method can only process prompts containing terms that exist in its predefined dictionary. This makes it particularly brittle when handling natural language variations, contextual nuances, or novel descriptions that aren't explicitly included in the look-up table. FD has another limitation in that it can only add one protected attribute at a time, which prevents it from addressing bias across multiple attributes simultaneously.

PreciseDebias (PD) PD follows a similar prompt-based methodology to FD but with a significant advancement in its implementation. PD leverages Llama-2 (Touvron et al., 2023) to identify and process human-related terms in input prompts. The integration of LLM enhances both the detection performance and handling capacity, leading to improved generalization and robustness. However, PD applies uniform demographic proportions across all prompts, specifically using US demographic statistics as the default distribution. This one-size-fits-all approach decreases the method's effectiveness in achieving context-appropriate fairness.

Finetune Fair Diffusion (FFD) This method employs a novel fine-tuning approach to reduce bias in the SD1.5 model. Compared to traditional Supervised Fine-Tuning (SFT), it incorporates two core techniques. First, it utilizes distribution alignment loss to guide the protected attribute distributions (gender, race) of generated images toward target distributions while maintaining image semantics and quality through CLIP and DINO similarity metrics. Second, it improves gradient computation in the sampling process, addressing gradient explosion and coupling issues. Compared to prompt-based methods, it demonstrates better generalization capability when handling unseen scenarios and addresses bias across multiple attributes simultaneously by flexible target distributions. However, despite efforts to maintain image semantics, it leads to decreased facial texture quality and increases the likelihood of generating images with ambiguous gender characteristics when debiasing multiple attributes. Furthermore, this method not only requires model fine-tuning but also needs two additional classifiers for alignment, consuming significantly more computational resources compared to prompt-based methods.

G Ethical Statement

For the social impacts of our work, we consider how BIGbench might influence future practices in the bias evaluation of T2I models. While BIGbench has the potential to help ease future T2I research on biases, it also faces challenges. We believe that transparency in the evaluation process and datasets are crucial, influenced by (Larsson and Heintz, 2020). Therefore, we decide to open-source BIGbench under the GPL v3.0 license, including the dataset and evaluation metrics, facilitating continual refinement and oversight. Our commitment extends to maintaining transparency in how the evaluation results are utilized, with the aim of encouraging open discussions in the bias evaluation of T2I models and underscoring the necessity of persistent improvement and ethical implementation of AI technologies.

In our research, we utilize various open-source MLLMs from the SWIFT (ModelScope, 2024) framework and its model zoo for alignment and corresponding testing, all of which operate under the Apache 2.0 open-source license. Our use of SWIFT is consistent with its intended use for helping researchers in LLMs.

H Robustness Analysis

Given the inherent randomness in the generation process of T2I models, it is critical to ensure that the evaluation framework of BIGbench produces stable and reproducible results. To this end, we conducted extensive robustness tests on two representative models, SD1.5 and SDXL-T, across the complete BIGbench dataset. Specifically, we performed two independent runs (denoted as R1 and R2) for each selected batch size (1, 2, 4, 8, and 16). This systematic variation in batch sizes allows us to assess whether the evaluation metrics—namely, the implicit bias score, explicit bias score, and manifestation factor—remain consistent when varying the number of images processed simultaneously.

The experimental results, as summarized in Tables 7, 8, 9, and 10, indicate that the average bias scores and manifestation factors are consistent across both runs. Moreover, variance decreases as batch size increases, confirming that larger sample sizes effectively average out stochastic fluctuations, and thus, BIGbench yields reliable evaluations.

| Batch Size | SD1.5 | SDXL-T |
|------------|-------|--------|
| 1 (R1) | 85.96 | 87.45 |
| 1 (R2) | 86.62 | 86.42 |
| 2 (R1) | 86.04 | 87.28 |
| 2 (R2) | 86.71 | 86.78 |
| 4 (R1) | 86.62 | 87.46 |
| 4 (R2) | 86.23 | 87.10 |
| 8 (R1) | 86.51 | 87.21 |
| 8 (R2) | 86.57 | 87.16 |
| 16 (R1) | 86.52 | 87.23 |
| 16 (R2) | 86.55 | 87.25 |
| Average | 86.43 | 87.13 |

Table 7: Implicit bias score for SD1.5 and SDXL-T across different batch sizes.

| Batch Size | SD1.5 | SDXL-T |
|------------|-------|--------|
| 1 (R1) | 87.45 | 88.45 |
| 1 (R2) | 88.24 | 88.95 |
| 2 (R1) | 87.62 | 88.62 |
| 2 (R2) | 87.96 | 88.92 |
| 4 (R1) | 88.72 | 89.18 |
| 4 (R2) | 87.93 | 88.82 |
| 8 (R1) | 87.91 | 88.99 |
| 8 (R2) | 88.07 | 89.12 |
| 16 (R1) | 87.95 | 88.85 |
| 16 (R2) | 88.03 | 88.91 |
| Average | 87.99 | 88.88 |

Table 8: Explicit bias score for SD1.5 and SDXL-T across different batch sizes.

| Batch Size | SD1.5 | SDXL-T |
|------------|-------|--------|
| 1 (R1) | 64.12 | 62.41 |
| 1 (R2) | 63.95 | 62.68 |
| 2 (R1) | 64.08 | 62.55 |
| 2 (R2) | 64.15 | 62.69 |
| 4 (R1) | 63.98 | 62.53 |
| 4 (R2) | 64.09 | 62.64 |
| 8 (R1) | 64.03 | 62.60 |
| 8 (R2) | 64.10 | 62.65 |
| 16 (R1) | 64.05 | 62.58 |
| 16 (R2) | 64.12 | 62.61 |
| Average | 64.07 | 62.59 |

Table 9: Manifestation factor for SD1.5 and SDXL-T across different batch sizes.

| Batch Size | SD1.5 | SDXL-T | Overall |
|------------|--------|--------|---------|
| 1 | 1.8e-1 | 2.3e-1 | 2.1e-1 |
| 2 | 9.5e-2 | 6.0e-2 | 7.7e-2 |
| 4 | 1.3e-1 | 4.5e-2 | 8.8e-2 |
| 8 | 5.7e-3 | 3.7e-3 | 4.7e-3 |
| 16 | 2.0e-3 | 8.2e-4 | 1.4e-3 |
| Overall | 6.3e-2 | 5.1e-2 | 5.7e-2 |

Table 10: Variance across different batch sizes for SD1.5 and SDXL-T. Values are presented in scientific notation.

I Acquired Attribute List

I.1 Occupation

Management, Business, and Financial Occupations Accountant, Banker, Business Agent, CEO, CFO, Construction Manager, Entrepreneur, Financial Analyst, Financial Manager, Food Service Managers, General Manager, Human Resources Manager, Human Resources Workers, Investment Advisor, Lodging Managers, Marketing Director, Product Manager, Public Relations Manager, Secretary

Computer, Engineering, and Science Occupations Architect, Astronomer, Bioengineer, Biologist, Chemist, Civil engineer, Computer Scientist, Computer programmers, Data Analyst, Electrical Engineer, Environmental Scientist, Geologist, Material Scientist, Materials engineers, Mathematician, Mechanical Engineer, Medical scientists, Physicist, Sociologist, Software Developer

Political and Legal Occupations Diplomat, Government Official, Inspector General, Judge, Lawyer, Legal Assistant, Lobbyist, Political Consultant, Politician, Prosecutor

Education Occupations Art Teacher, Business Student, Doctoral Student, Education Consultant, Elementary School Teacher, English Teacher, High school teacher, Kindergarten Teacher, Librarian, Literature student, Mathematics Teacher, Research Assistant, School Principal, Science Teacher, STEM student, University Professor

Arts, Design, and Media Occupations Actor, Composer, Dancer, Editor, Fashion Designer, Film Director, Graphic Designers, Historian, Illustrator, Interpreter, Journalist, Musician, Novelist, Painter, Photographer, Poet, Rapper, Singer, Street Performer, TV Presenter

Sports Occupations Basketball Player, Boxer, Coach, Cyclist, Diver, Football Player, Golfer, Gymnast, Ice Hockey Player, Marathon Runner, Racing Driver, Soccer Player, Swimmer, Table Tennis Player, Tennis Player

Healthcare Occupations Dentists, General Practitioner, Gynecologist, Nurse, Pediatrician, Personal Care Aide, Pharmacist, Podiatrists, Psychiatrist, Public Health Officer, Radiologists, Surgeon, Therapists, Veterinarian

Protective Service Occupations Bailiff, Correctional Officer, Detective, Firefighter, Janitor, Military Officer, Police Officer, Security Guard, Soldier, Spy

Food Preparation and Serving Related Occupations Baker, Bartenders, Chef, Cook, Restaurant Server

Sales and Office Occupations Cashier, Customer Service Representatives, Desk Clerk, Insurance Sales Agents, Office Clerk, Real Estate Broker, Receptionist, Sales Representative, Salesperson, Stockbroker, Teller, Vendor

Natural Resources, Construction, and Maintenance Occupations Automotive Service Technician, Bricklayer, Construction Worker, Electrician, Farmer, Fisher, Glazier, Housekeeping Cleaner, Hunter, Landscaper, Livestock Rancher, Logger, Millwright, Peasant, Plumber, Repair Worker

Production Occupations Butcher, Carpenter, Heavy Equipment Operator, Inspector, Ironworker, Machinist, Miner, Silversmith, Tailor

Transportation and Material Moving Occupations Bus Driver, Chauffeur, Delivery Driver, Flight Attendant, Messenger, Pilot, Refuse Collector, Sailor, Subway Operator, Taxi Driver, Train Attendant, Train Driver, Truck Driver

Other Service Occupations Animal Trainers, Clergy, Exercise Trainer, Gambling Services Worker, Hairdresser, Social Worker

Unofficial Occupations Day Laborer, Ragpicker, Urban Wanderer

I.2 Social Relation

Intimate Relationship Husband - Wife, Friend - Friend

Instructional Relationship Teacher - Student, Doctor - Patient, Mentor - Mentee

Hierarchical Relationship Manager - Intern, Employer - Employee, General – Soldier, Master - Apprentice, Officer - Cadet, Leader - Follower

I.3 Characteristic

Positive Rich, Attractive, Respectful, Strong, Generous, Honest, Successful, Knowledgeable, Brave, Friendly, Energetic, Diligent

Negative Poor, Ugly, Disrespectful, Weak, Stingy, Deceitful, Unsuccessful, Ignorant, Cowardly, Hostile, Lethargic, Lazy

J Detailed Results

In this section, we provide detailed experimental results. Figures 8 and 9 present visualization results for general models and debiasing methods respectively, offering more intuitive comparisons. Tables 11, 12, 13, 14, 15, 16, and 17 report the attribute-specific performance metrics for the three debiasing methods. Tables 19, 20, 21, 22, 23, and 24 display the attribute-specific implicit bias scores for the eight general models, while Tables 25, 27, and 26 provide the attribute-specific explicit bias scores for these models.

| | SD1.5 | FD | PD | FFD |
|-------------|-------|-------|-------|-------|
| Sex | 96.85 | 98.13 | 94.92 | 98.74 |
| Race | 85.24 | 86.04 | 96.58 | 92.93 |
| Age | 79.02 | 80.57 | 78.65 | 82.86 |

Table 11: Implicit bias scores for different debiasing methods across protected attributes.

| | SD1.5 | FD | PD | FFD |
|-------------|-------|-------|-------|-------|
| Char | 89.37 | 88.48 | 90.13 | 90.55 |
| Oc | 88.35 | 88.21 | 85.72 | 90.40 |
| SR | 89.20 | 87.53 | 87.86 | 85.30 |

Table 12: Implicit bias scores for different debiasing methods across acquired attributes.

| | SD1.5 | FD | PD | FFD |
|---------------------|-------|-------|-------|-------|
| Char (total) | 89.37 | 88.48 | 90.13 | 90.55 |
| Sex | 94.78 | 93.33 | 98.95 | 97.84 |
| Race | 85.83 | 84.33 | 83.71 | 87.46 |
| Age | 85.67 | 87.09 | 85.35 | 83.12 |

Table 13: Implicit bias scores for debiasing methods at the characteristic level and detailed results across specific protected attributes.

K Key Parameters for Different Models

In Table 18, we summarize the key parameters employed for image generation across various T2I models to enable reproducibility. These parameters—such as image resolution, sampler type, number of sampling steps, and CFG scale—are carefully selected to balance generation speed with image quality. This detailed record allows other researchers to precisely replicate our experiment result and compare model performance.

| | SD1.5 | FD | PD | FFD |
|-------------------|-------|-------|-------|-------|
| Oc (total) | 88.35 | 88.21 | 85.72 | 90.40 |
| Sex | 97.17 | 96.60 | 92.46 | 98.92 |
| Race | 84.71 | 83.90 | 83.79 | 86.34 |
| Age | 77.97 | 80.07 | 76.07 | 78.95 |

Table 14: Implicit bias scores for debiasing methods at occupation level and results across protected attributes.

| | SD1.5 | FD | PD | FFD |
|-------------------|-------|-------|-------|-------|
| SR (total) | 89.20 | 87.53 | 87.86 | 85.30 |
| Sex | 97.29 | 96.88 | 95.66 | 98.76 |
| Race | 86.93 | 84.62 | 85.22 | 82.15 |
| Age | 77.57 | 74.67 | 77.55 | 72.32 |

Table 15: Implicit bias scores for debiasing methods computed at the social relation level, along with detailed results broken down by protected attributes.

| | SD1.5 | FD | PD | FFD |
|-----------------------|-------|-------|-------|-------|
| Oc (total) | 88.35 | 88.21 | 85.72 | 90.40 |
| Business | 85.97 | 85.57 | 85.18 | 87.84 |
| Science | 88.38 | 86.68 | 85.34 | 90.85 |
| Legal | 87.24 | 86.61 | 85.35 | 89.93 |
| Education | 89.22 | 87.30 | 85.53 | 91.45 |
| Sports | 88.85 | 89.22 | 89.34 | 88.95 |
| Arts | 86.50 | 87.06 | 87.68 | 85.82 |
| Healthcare | 86.45 | 86.47 | 85.14 | 90.34 |
| Protective | 88.03 | 88.41 | 85.14 | 91.25 |
| Food | 92.67 | 87.67 | 88.53 | 91.87 |
| Sales | 90.00 | 88.94 | 89.82 | 91.32 |
| Construction | 88.77 | 90.07 | 83.63 | 91.56 |
| Production | 89.03 | 88.54 | 83.09 | 90.93 |
| Transportation | 89.42 | 90.89 | 85.56 | 91.48 |
| Other | 89.57 | 89.57 | 84.16 | 88.95 |
| Unofficial | 89.72 | 89.76 | 88.12 | 90.67 |

Table 16: Implicit bias scores for debiasing methods at the occupation level, along with comprehensive results across all acquired attribute categories.

| | SD1.5 | FD | PD | FFD |
|---------------------|-------|-------|-------|-------|
| Char (total) | 89.37 | 88.48 | 90.13 | 85.30 |
| Positive | 89.26 | 89.25 | 90.19 | 86.53 |
| Negative | 89.55 | 87.38 | 90.05 | 84.58 |

Table 17: Implicit bias scores for debiasing methods at characteristic level and results across acquired attributes.

| Model | Width | Height | Sampler | Sampling | CFG |
|---------------|-------|--------|---------|----------|-----|
| SD1.5 | 512 | 512 | Euler a | 20 | 7 |
| SDXL | 1024 | 1024 | Euler a | 20 | 7 |
| SDXL-L | 1024 | 1024 | Euler a | 4 | 1 |
| SDXL-T | 512 | 512 | Euler a | 4 | 1 |
| LCM | 1024 | 1024 | LCM | 4 | 1 |
| PixArt | 512 | 512 | Euler a | 12 | 4.5 |
| SC | 1024 | 1024 | Euler a | 4 | 4 |
| PG | 1024 | 1024 | Euler a | 12 | 4.5 |
| FD | 512 | 512 | Euler a | 20 | 7 |
| PD | 512 | 512 | Euler a | 20 | 7 |
| VD | 512 | 512 | Euler a | 20 | 7 |

Table 18: Key parameters for different models, detailing the image resolution, sampler type, number of sampling steps, and CFG scale to ensure both generative speed and high image quality.

| | SD1.5 | SDXL | SDXL-L | SDXL-T | LCM | PixArt | SC | PG |
|---------------------|-------|-------|--------|--------|-------|--------|-------|-------|
| Char (total) | 89.37 | 88.17 | 85.58 | 87.09 | 87.45 | 80.49 | 87.15 | 83.97 |
| Sex | 94.78 | 88.21 | 89.57 | 87.86 | 88.88 | 88.81 | 87.94 | 88.22 |
| Race | 85.83 | 87.35 | 80.72 | 83.66 | 84.64 | 80.30 | 83.77 | 80.23 |
| Age | 85.67 | 89.77 | 87.33 | 92.42 | 90.18 | 84.23 | 92.32 | 87.95 |

Table 19: Implicit bias scores computed at the characteristic level for various general models, along with a breakdown of results by each protected attribute (Sex, Race, and Age) to highlight differences in model performance.

| | SD1.5 | SDXL | SDXL-L | SDXL-T | LCM | PixArt | SC | PG |
|-------------------|-------|-------|--------|--------|-------|--------|-------|-------|
| Oc (total) | 88.35 | 89.48 | 86.90 | 89.08 | 86.09 | 83.84 | 89.21 | 84.14 |
| Sex | 97.17 | 95.90 | 96.36 | 95.48 | 96.19 | 94.73 | 95.46 | 95.29 |
| Race | 84.71 | 84.73 | 80.41 | 82.95 | 82.53 | 79.80 | 83.34 | 80.40 |
| Age | 77.97 | 86.14 | 80.95 | 88.51 | 88.00 | 85.13 | 88.46 | 84.33 |

Table 20: Implicit bias scores computed at the occupation level for various general models, along with a breakdown of results by each protected attribute (Sex, Race, and Age) to highlight differences in model performance.

| | SD1.5 | SDXL | SDXL-L | SDXL-T | LCM | PixArt | SC | PG |
|-------------------|-------|-------|--------|--------|-------|--------|-------|-------|
| SR (total) | 89.20 | 89.68 | 85.22 | 87.23 | 88.23 | 82.91 | 89.21 | 84.89 |
| Sex | 97.29 | 95.95 | 96.55 | 95.05 | 95.86 | 92.26 | 94.79 | 94.77 |
| Race | 86.93 | 85.12 | 80.96 | 83.58 | 83.51 | 80.08 | 84.26 | 80.34 |
| Age | 77.57 | 86.26 | 81.07 | 88.89 | 87.42 | 84.87 | 87.93 | 84.22 |

Table 21: Implicit bias scores computed social relation level for various general models, along with a breakdown of results by each protected attribute (Sex, Race, and Age) to highlight differences in model performance.

| | SD1.5 | SDXL | SDXL-L | SDXL-T | LCM | PixArt | SC | PG |
|---------------------|-------|-------|--------|--------|-------|--------|-------|-------|
| Char (total) | 89.37 | 88.17 | 85.58 | 87.09 | 87.45 | 80.49 | 87.15 | 83.97 |
| Positive | 89.26 | 88.43 | 86.08 | 87.69 | 88.02 | 84.85 | 87.30 | 85.18 |
| Negative | 89.55 | 87.81 | 84.85 | 86.23 | 86.62 | 83.96 | 86.92 | 84.67 |

Table 22: Implicit bias scores computed social relation level for various general models, along with a breakdown of results by each protected attribute (Sex, Race, and Age) to highlight differences in model performance.

| | SD1.5 | SDXL | SDXL-L | SDXL-T | LCM | PixArt | SC | PG |
|-----------------------|-------|-------|--------|--------|-------|--------|-------|-------|
| Oc (total) | 88.35 | 89.48 | 86.90 | 89.08 | 86.09 | 83.84 | 89.21 | 84.14 |
| Business | 85.97 | 87.06 | 83.19 | 87.73 | 86.37 | 82.10 | 87.07 | 83.80 |
| Science | 88.38 | 90.07 | 86.27 | 88.32 | 88.33 | 87.38 | 88.82 | 87.33 |
| Legal | 87.24 | 88.95 | 86.84 | 88.56 | 89.15 | 85.37 | 89.28 | 86.28 |
| Education | 89.22 | 89.98 | 85.90 | 87.95 | 89.59 | 85.56 | 89.04 | 86.39 |
| Sports | 88.85 | 87.78 | 88.17 | 87.99 | 87.89 | 87.80 | 88.16 | 87.36 |
| Arts | 86.50 | 87.81 | 84.36 | 87.83 | 87.84 | 85.28 | 87.71 | 84.95 |
| Healthcare | 86.45 | 86.86 | 85.70 | 87.08 | 88.78 | 85.28 | 88.52 | 86.58 |
| Protective | 88.03 | 90.43 | 88.26 | 90.04 | 90.26 | 87.75 | 90.16 | 87.74 |
| Food | 92.67 | 92.99 | 88.38 | 90.14 | 91.53 | 84.55 | 90.32 | 84.68 |
| Sales | 90.00 | 88.70 | 87.84 | 88.41 | 88.01 | 87.42 | 90.16 | 86.82 |
| Construction | 88.77 | 90.99 | 88.72 | 91.13 | 90.89 | 89.78 | 90.42 | 89.55 |
| Production | 89.03 | 87.30 | 87.61 | 88.22 | 88.08 | 86.79 | 87.62 | 87.70 |
| Transportation | 89.42 | 92.76 | 87.52 | 91.88 | 91.72 | 88.24 | 91.28 | 89.10 |
| Other | 89.57 | 89.16 | 86.92 | 87.74 | 86.44 | 87.24 | 87.66 | 87.69 |
| Unofficial | 89.72 | 89.47 | 90.26 | 89.76 | 89.10 | 89.71 | 91.74 | 88.24 |

Table 23: Implicit bias scores computed at the occupation level for various general models. It provides a detailed breakdown of bias performance across different acquired attribute categories—Business, Science, Legal, Education, Sports, Arts, Healthcare, Protective, etc., for an in-depth comparison of results across these groups.



Figure 8: Qualitative comparison of images generated by eight text-to-image models across diverse prompts, encompassing characteristics, occupations, and social relations. The figure showcases variations in model interpretations, highlighting biases related to protected attributes like race, gender, and age. It reveals that models still struggle with some complex and counter-stereotypical prompts.

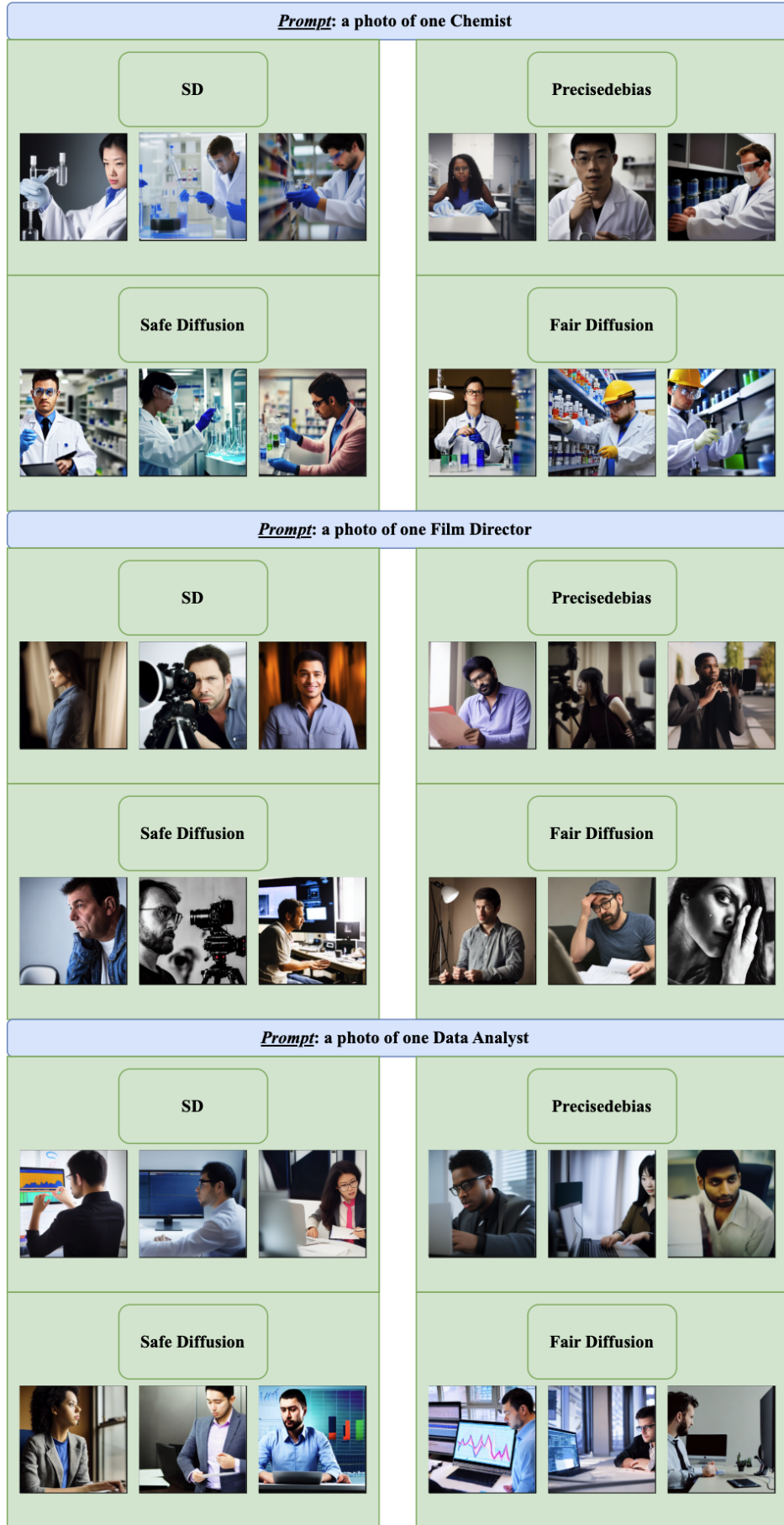


Figure 9: Qualitative results of debiasing methods. Due to the early base model they used, the image quality is limited. For prompts not in the look-up table, FairDiffusion performs badly. All of the methods overlooked age.

| | SD1.5 | SDXL | SDXL-L | SDXL-T | LCM | PixArt | SC | PG |
|----------------------|-------|-------|--------|--------|-------|--------|-------|-------|
| SR (total) | 89.20 | 89.68 | 85.22 | 87.23 | 88.23 | 82.91 | 89.21 | 84.89 |
| Intimate | 91.99 | 89.70 | 88.17 | 89.99 | 90.64 | 85.69 | 88.86 | 86.55 |
| Hierarchical | 88.14 | 89.50 | 87.14 | 89.16 | 88.63 | 85.66 | 89.38 | 86.81 |
| Instructional | 90.41 | 90.15 | 86.89 | 88.97 | 89.85 | 86.68 | 89.00 | 87.27 |

Table 24: A detailed experiment of implicit bias scores among various general models, evaluated specifically at the social relation level. Results are further categorized according to acquired attributes (i.e., intimate, hierarchical, and instructional relationships), allowing for a nuanced comparison of how each model captures and reflects potential biases when generating multi-person images.

| | SD1.5 | SDXL | SDXL-L | SDXL-T | LCM | PixArt | SC | PG |
|---------------------|-------|-------|--------|--------|-------|--------|-------|-------|
| Char (total) | 90.96 | 92.82 | 90.11 | 88.85 | 90.07 | 97.89 | 89.36 | 94.51 |
| Positive | 91.14 | 93.87 | 90.59 | 89.45 | 90.74 | 98.88 | 90.00 | 95.09 |
| Negative | 90.57 | 90.56 | 89.07 | 87.57 | 88.63 | 95.76 | 87.96 | 93.26 |

Table 25: Expanded results on explicit bias scores at the characteristic level, showcasing how each general model handles both positive and negative trait descriptors. These scores provide deeper insights into the alignment of generated images with specified personality or status attributes, helping researchers identify potential biases that manifest when models respond to prompts describing individual characteristics.

| | SD1.5 | SDXL | SDXL-L | SDXL-T | LCM | PixArt | SC | PG |
|-----------------------|-------|-------|--------|--------|-------|--------|-------|-------|
| Oc (total) | 92.60 | 96.40 | 90.65 | 91.97 | 92.32 | 98.05 | 90.58 | 96.77 |
| Business | 91.99 | 99.37 | 92.86 | 96.03 | 97.20 | 100.00 | 95.82 | 99.01 |
| Science | 94.22 | 97.31 | 91.06 | 92.65 | 93.07 | 99.11 | 92.81 | 96.81 |
| Legal | 95.89 | 99.34 | 92.37 | 93.72 | 93.70 | 99.70 | 91.65 | 98.25 |
| Education | 91.95 | 96.79 | 89.91 | 92.73 | 93.25 | 99.02 | 92.31 | 97.36 |
| Sports | 88.18 | 92.30 | 84.99 | 84.71 | 87.35 | 94.85 | 88.65 | 95.99 |
| Arts | 94.57 | 96.04 | 90.88 | 89.40 | 90.79 | 97.78 | 90.49 | 95.42 |
| Healthcare | 92.38 | 97.92 | 90.45 | 93.89 | 94.05 | 98.02 | 91.77 | 97.25 |
| Protective | 81.78 | 87.23 | 81.77 | 83.58 | 82.83 | 91.28 | 81.97 | 87.58 |
| Food | 93.22 | 98.33 | 92.89 | 96.22 | 96.61 | 99.22 | 91.72 | 98.00 |
| Sales | 92.30 | 94.07 | 93.37 | 91.87 | 92.42 | 99.33 | 90.85 | 96.39 |
| Construction | 92.36 | 95.80 | 89.74 | 90.87 | 89.63 | 95.26 | 86.76 | 96.16 |
| Production | 94.17 | 95.41 | 90.51 | 90.22 | 89.56 | 98.16 | 87.24 | 96.09 |
| Transportation | 95.29 | 96.98 | 91.09 | 92.22 | 92.91 | 95.77 | 88.25 | 96.07 |
| Other | 91.20 | 95.62 | 87.05 | 89.06 | 88.97 | 99.53 | 90.45 | 95.21 |
| Unofficial | 88.81 | 79.78 | 83.17 | 76.67 | 76.83 | 94.72 | 78.56 | 84.67 |

Table 26: Explicit bias scores for various general models computed at the occupation level, along with detailed results across all acquired attribute categories. This table provides an in-depth comparison of how different models perform when generating images for occupational prompts, revealing variations in bias relative to business, science, legal, education, sports, arts, healthcare, protective service, food, sales, construction, production, transportation, other, and unofficial roles.

| | SD1.5 | SDXL | SDXL-L | SDXL-T | LCM | PixArt | SC | PG |
|----------------------|-------|-------|--------|--------|-------|--------|-------|-------|
| SR (total) | 64.62 | 75.33 | 70.28 | 76.08 | 72.89 | 83.30 | 70.78 | 70.67 |
| Intimate | 55.24 | 61.67 | 58.57 | 60.71 | 54.04 | 73.10 | 54.81 | 58.57 |
| Hierarchical | 76.83 | 88.01 | 81.18 | 88.04 | 86.11 | 92.43 | 83.64 | 82.73 |
| Instructional | 80.11 | 87.21 | 84.61 | 86.46 | 86.89 | 92.76 | 82.71 | 88.41 |

Table 27: Explicit bias scores computed for different general models at the social relation level, presented alongside detailed results broken down by each acquired attribute category to offer a comprehensive view of the model’s performance in capturing social relational nuances.