AI-GENERATED FACES INFLUENCE GENDER STEREO-TYPES AND RACIAL HOMOGENIZATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Text-to-image generative AI models such as Stable Diffusion are used daily by millions worldwide. However, the extent to which these models exhibit racial and gender stereotypes is not yet fully understood. Here, we document significant biases in Stable Diffusion across six races, two genders, 32 professions, and eight attributes. Additionally, we examine the degree to which Stable Diffusion depicts individuals of the same race as being similar to one another. This analysis reveals significant racial homogenization, e.g., depicting nearly all middle eastern men as dark-skinned, bearded, and wearing a traditional headdress. We then propose debiasing solutions that address the above stereotypes. Finally, using a preregistered experiment, we show that being presented with inclusive AI-generated faces reduces people's racial and gender biases, while being presented with non-inclusive ones increases such biases. This persists regardless of whether the images are labeled as AI-generated. Taken together, our findings emphasize the need to address biases and stereotypes in AI-generated content.

025 026 027

004

010 011

012

013

014

015

016

017

018

019

021

023

1 INTRODUCTION

028 Artificial intelligence biases refer to the systematic and unfair preferences or prejudices embedded 029 in AI systems, often reflecting the biases present in the data used to train these systems Christian (2020). Such biases can perpetuate and even exacerbate societal inequalities, as the AI algorithms 031 may inadvertently discriminate against certain groups. One glaring example of AI bias is the COM-PAS algorithm, which has been shown to produce racially biased predictions when assessing the 033 likelihood of recidivism Angwin et al. (2016); Dressel & Farid (2018). Another notable example 034 of AI bias is found in face recognition algorithms, which have been shown to discriminate based on race and gender Buolamwini & Gebru (2018). Another such example is Amazon's endeavor 035 to use an algorithm to evaluate job candidates based on their resumes, which inadvertently penal-036 ized resumes that included phrases associated with women Dastin (2018). In the context of NLP, 037 notable biases have been reported in popular word embedding models such as BERT and GPT-2, which associate certain occupations or stereotypes more strongly with one gender or racial group than another Nadeem et al. (2021). 040

In this paper, we focus on racial and gender stereotypes in SDXL (Stable Diffusion XL) Podell 041 et al. (2023), one of the most popular text-to-image generative models used daily by millions world-042 wide Fatunde & Tse (2022). Recent studies have demonstrated that Stable Diffusion underrepresents 043 certain races or genders Bianchi et al. (2023); Wang et al. (2023); Ghosh & Caliskan (2023), but none 044 of these studies proposed debiasing solutions. Moreover, none of them offered a comprehensive ex-045 amination of such biases across racial groups, genders, professions, and attributes. Other studies 046 have proposed debiasing solutions Zhang et al. (2023a); Friedrich et al. (2023), but these are either 047 not automated, or are unable to generate images that adequately represent complex prompts; see 048 the Related Work section for more details. Another form of stereotype that has been overlooked in the literature is when individuals of the same race are depicted as being too similar to one another, e.g., when Middle Eastern men are all depicted as being bearded and dark-skinned, or when Mid-051 dle Eastern women are all depicted as wearing a traditional headcover, or "hijab". Consequently, it remains unclear whether such homogenization (if it exists) can be addressed by diversifying the 052 facial features of same-race individuals. Other open questions that have not been addressed to date are whether being exposed to AI-generated faces can affect people's racial and gender biases.

054 To address these questions, we start off by developing a classifier to predict race and gender, allowing 055 us to quantify biases in SDXL-generated images across six races, two genders, 32 professions, 056 and eight attributes. We then propose a debiasing solution, called SDXL-Inc (where Inc stands 057 for inclusive), and demonstrate its ability to outperform alternatives across various benchmarks. 058 Additionally, using a measure of image similarity, we reveal that SDXL exhibits a high degree of racial homogenization, depicting individuals from certain racial backgrounds as being very similar to one another. We address this issue by proposing another solution, called SDXL-Div (where Div 060 stands for diversity). Finally, using preregistered randomized controlled trials, we show that being 061 exposed to inclusive AI-generated faces can reduce people's racial and gender biases, while exposure 062 to non-inclusive AI-generated faces can increase these biases. 063

064 065

2 RELATED WORK

066 067

Compared to all the works discussed in this section, ours is the only one that examines, and ad dresses, the problem of racial homogenization in AI-generated faces, i.e., the depiction of individu als from a specific racial or ethnic group as too similar in appearance. We are also the first to conduct
 a randomized control trial to understand the impact of being exposed to inclusive and non-inclusive
 AI-generated faces, and whether the AI-label plays a role in this phenomenon. Next, we summarize
 relevant papers and discuss any additional differences that may exist between our work and theirs.

074 Bianchi et al. (2023) used a previous version of Stable Diffusion (v1-4) to examine biases across ten 075 professions and three races. They classified the race and gender of any generated image as follows. First, for each of the five demographic categories considered (i.e., Male, Female, White, Black, 076 Asian), they took the images that represent this category in the Chicago Face Dataset Ma et al. 077 (2015), and fed those images to CLIP (the core representational component of Stable Diffusion) to generate vector representations. These were then averaged to obtain a single archetypal vector 079 representation of the category (e.g., a single vector for Black). Any image can then be classified as 080 X if its vector representation is most closely aligned (in cosine distance) to the archetypal vector 081 representation of X (e.g., Black) compared to the alternative categories (e.g., Asian or White). The authors used this classifier to analyze 100 images per occupation, and found professional stereotypes 083 that reinforce racial and gender disparities. Compared to our work, the authors did not consider cer-084 tain racial groups, such as Indian, Latinx, and Middle Eastern. They also did not propose debiasing 085 solutions.

Ghosh & Caliskan (2023) used a previous version of Stable Diffusion (v2.1) to analyze biases across 087 three genders and 27 nationalities, but did not consider professions nor proposed debiasing solutions. 880 Using CLIP's cosine similarity, they compared images generated using the prompt: 'a front-facing 089 photo of a person' against those generated using the prompt: 'a front-facing photo of X' where X 090 is either a gender (i.e., a man, a woman, a nonbinary gender) or a nationality (e.g., a person from 091 Brazil). The authors then manually classified the images based on skin tone (light vs. dark) and 092 found that most generated images were light-skinned. Compared to our work, the authors did not 093 consider professions (e.g., Doctor, Nurse, etc.) nor attributes (e.g., Beautiful, Intelligent, etc.) in their analysis. Moreover, in the context of race, the authors manually classified skin-color rather 094 than proposing a classifier, and only provided qualitative rather than quantitative conclusions. 095

Wang et al. (2023) studied the association of pleasant and unpleasant attributes with certain concepts
such as: European American and African American names, Light Skin and Dark Skin, Straight and
Gay. In addition, they studied gender stereotypes across eight professions. Compared to our work,
the authors did not propose any debiasing solutions. They also did not consider racial groups apart
from European American (White) and African American (Black).

Friedrich et al. (2023) proposed a "user in control" solution, called Fair Diffusion, which utilizes a
 textual interface allowing users to instruct generative image models on fairness. Intuitively, this is
 made possible by extending classifier-free guidance Ho & Salimans (2022) with an additional fair
 guidance term, which depends on additional fairness instructions provided by the user. This ap proach requires no data filtering nor additional training. The authors used their approach to analyze
 images generated by a previous version of Stable Diffusion (1.5), focusing on professional stereo types across genders. Compared to our work, the authors did not consider racial stereotypes, and
 did not consider a fully automated solution.

108 Zhang et al. (2023a) proposed a solution called ITI-GEN (inclusive Text-to-Image Generation), 109 designed to generate images that are uniformly distributed across attributes of interest. One such 110 attribute could be, e.g., "gender", in which case the images generated by ITI-GEN would be equally 111 split between the various attribute categories, e.g., "male" and "female". More specifically, given 112 an input prompt and some desired attribute(s), the model learns discriminative token embeddings representing each category of each attribute. It then injects these learned tokens after the original 113 prompt, thereby synthesizing a set of prompts, each representing a unique combination of categories 114 belonging to different attributes. Finally, this set is used to generate an equal number of images for 115 any category combination. The authors used their model to analyze a previous version of Stable 116 Diffusion (v1-4), focusing on 40 physical attributes (each consisting of two categories, one positive 117 and one negative), along with gender, skin tone, and age. 118

Compared to our work, Zhang et al. focused on skin tone rather than race. Thus, unlike our analysis, 119 theirs does not distinguish between, say, Asian and White individuals who happen to be equally 120 light-skinned, or between Indian and Latinx individuals who happen to have similar skin tones. The 121 authors also focus on physical attributes (e.g., Black hair, Mustache, etc.) as opposed to those de-122 scribing various characteristics like the ones used in our analysis (e.g., Criminal, Intelligent, Parent, 123 Poor, etc.). Additionally, their analysis of professional stereotypes focuses on four professions and 124 200 images per profession, while we focus on 32 professions and 10,000 images per profession. 125 Finally, as discussed in Section D, their solution is unable to debias complex prompts, such as "a 126 person with green hair and eyeglasses", or "a person with the Eiffel tower". 127

127 128 129

130

132

3 MATERIALS AND METHODS

131 3.1 DATA OVERVIEW

133 I) LAION-5B: This is the dataset that was used to train Stable Diffusion Schuhmann et al. (2022). 134 We utilized a subset of this dataset, consisting of high-resolution images Beaumont (2021), to de-135 termine whether the biases in Stable Diffusion XL (SDXL) can be entirely attributed to the training 136 data. We randomly selected 172,923 images from this subset, and kept those having one or more 137 of the following keywords: face, person, child, woman, or man. We then cropped the images to 138 retain only the face(s) appearing therein, and discarded any resulting face images that are smaller 139 than 100×100 pixels. This filtering process left us with a final set of 88,714 face images.

140 II) FairFace: This is one of the largest public datasets of face images. For each image, the dataset specifies the race (Black, East Asian, Indian, Latinx, Middle Eastern, Southeast Asian, and White), 141 and the gender (female, male) Karkkainen & Joo (2021). Moreover, the dataset is divided into two 142 sets: 86,744 images for training and 10,954 for validation. We combined East and Southeast Asian 143 into a single category: Asian. The resulting dataset was then used to train and validate our race and 144 gender classifiers. For race, the number of images used for validation was: 1556 for Black; 2965 for 145 Asian; 1516 for Indian; 1623 for Latinx or Hispanic; 1209 for Middle Eastern; and 2085 for White. 146 As for gender, the number of images used for validation was: 5162 for female; and 5792 for male. 147

III) Flickr-Faces-HQ: This dataset consists of 70,000 high-resolution (1024×1024) images of human faces crawled from Flickr with considerable variation in terms of age and race Karras & Hellsten (2023). It is unlabeled, and was originally created as a benchmark for generative adversarial networks (GAN). We utilized this dataset to fine-tune our SDXL-Div model in order to generate face images with varying races and facial features. This was done to overcome the fact that, for certain races, the images generated by Stable Diffusion XL seem too similar to one another.

154 IV) Stable Diffusion validation: This dataset consists of images that we have generated using SDXL. 155 In particular, it consists of 10,000 images per race and 10,000 images per gender. For race, the 156 prompt used to generate the images was: "*a photo of a X*", where $X \in \{Asian, Black, Indian, Latino or Hispanic, Middle Eastern, White}. As for gender, the prompt used was: "$ *a photo of a X*", $157 where <math>X \in \{female, male\}$. This dataset was used to validate our race and gender classifiers.

159 V) SDXL-Inc fine-tuning: This dataset consists of Stable Diffusion-generated images with varying 160 race, gender, and profession. The images were generated using the prompt: "*a photo of a X Y Z*, 161 *looking at the camera, closeup headshot facing forward, ultra quality, sharp focus*", where $X \in$ {Asian, Black, Indian, Latino or Hispanic, Middle Eastern, White}; $Y \in$ {female, male}; and Z is one of the 21 professions listed in Appendix Table 1 under the "fine-tuning" category. For any of the twelve combinations of race and gender (X, Y), we compiled the corresponding images from all 21 professions into a single dataset, which was used to fine-tune a version of SDXL tailored specifically for race X and gender Y. This process yielded 12 fine-tuned models (one per racegender combination); these are the 12 components constituting our SDXL-Inc solution.

VI) Profession: This dataset consists of SDXL-generated images depicting 32 professions. More specifically, we generated 10,000 images per profession using the prompt: "*a photo of Z, looking at the camera, closeup headshot facing forward, ultra quality, sharp focus*", where Z is one of the 32 professions listed in Appendix Table 1. In this table, the 21 professions listed under the "fine-tuning" category were used to fine-tune SDXL, while the remaining 11 professions (i.e., those listed under "generalization testing") were used to evaluate the generalization capability of SDXL-Inc.

173 VII) Attribute: This dataset consists of Stable Diffusion-generated images depicting eight attributes. 174 In particular, we generated 10,000 images per attribute using the prompt: "*a photo of a X, looking at* 175 *the camera, closeup headshot facing forward, ultra quality, sharp focus*", where $X \in \{Poor, Winner, Beautiful, Intelligent, Parent, Sibling, Terrorist, Criminal\}.$ This dataset was used to evaluate the 177 generalization capability of SDXL-Inc.

- 178 179
- 180 3.2 METHODS
- 181 182 3.2.1 STABLE DIFFUSION

183 We utilized SDXL due to its improved ability to generate human faces compared to its predecessors. 184 We used the Hugging Face repository "stabilityai" and the model "stable-diffusion-xl-base-1.0" sta-185 bilityai (2023) to generate images to analyze racial and gender stereotypes in the context of various professions and attributes. We fine-tuned SDXL using LORA Hu et al. (2021), thereby creating our 187 SDXL-Inc model. To this end, we created a dataset consisting of prompt-image pairs. Each such pair consisted of an image taken from the "SDXL-Inc fine-tuning dataset" (described above) along with 188 the prompt: "a photo of Z, looking at the camera, closeup headshot facing forward, ultra quality, 189 sharp focus", where Z is the profession depicted in the image. Similarly, we fine-tuned SDXL using 190 LORA in the process of creating our SDXL-Div model. To this end, we created a dataset consist-191 ing of prompt-image pairs. Each such pair consisted of an image taken from the "Flickr-Faces-HQ 192 dataset" along with the prompt: "a photo of X person, looking at the camera, closeup headshot 193 facing forward, ultra quality, sharp focus", where X is the race depicted in the image. 194

Hyper-parameters: Image resolution = 1024; number of inference steps = 40; guidance scale = 5. When fine-tuning SDXL, we used: Image resolution = 1024; training batch size = 1; number of training epochs = 3; learning rate = 10^{-4} ; and mixed precision = fp16.

198 199

200

3.2.2 PROPOSED CLASSIFICATION PIPELINE

201 Our classifier includes three stages: face detection, face embedding generation, and classification.

I) Face detection: This stage is carried out using a Multi-task Cascaded Convolutional Neural Network (MTCNN), which is a deep cascaded multitask framework that utilizes the inherent correlation between detection and alignment to improve performance Zhang et al. (2016). It leverages a cascaded architecture with three stages of deep convolutional networks to predict faces. We selected MTCNN due to its ability to balance high detection accuracy and run-time speed. We configured the detector to exclude boundary-boxes with confidence scores ≤ 0.9 .

208 II) Face embedding generation: This stage is carried out using a VGGFace ResNet-50 Convolu-209 tional Neural Networks (VGGFace ResNet-50 CNN) Cao et al. (2018). ResNet-50 CNN was trained 210 on MS-Celeb-1M and VGGFace2, as well as the of the two. More specifically, MS-Celeb-1M has 10 211 million images depicting 100k different celebrities Guo et al. (2016). On the other hand, VGGFace 212 is a large-scale face dataset with considerable variations in pose, age, illumination, race, and profes-213 sion Cao et al. (2018), including 3.31 million images downloaded from Google image search. In our work, we utilized VGGFace ResNet-50 CNN, and removed the top layers to extract the embedding 214 vector from the face images in the FairFace dataset described earlier. The FairFace RGB images 215 were resized to 224x224 pixels before being fed to the CNN as input.

216 **III)** Classification: This stage is carried out using a Support Vector Machine (SVM) Hearst et al. 217 (1998). We trained two SVM classifiers to predict the race and gender using the embedding vectors 218 extracted from the FairFace images of the previous stage. The hyper-parameters used were: Regu-219 larization parameter C = 1; Kernel type = Radial Basis Function. This stage is repeated twice; once 220 for race, and once for gender. However, the previous two stages are executed only once.

222 3.2.3 GPT-IN-THE-LOOP

In addition to SDXL-Inc, we experimented with an alternative debiasing technique. This prompt-224 regulating technique introduces an intermediary layer between the user (who provides the prompt) 225 and SDXL. The layer uses GPT-4 to detect whether the user-provided prompt corresponds to the 226 generation of an image depicting a person without a specific race and/or gender. GPT-4 would then 227 inject a randomly-selected race (if race was not specified by the user) and/or a random gender (if 228 gender was not specified) into the user-provided prompt before passing it on to SDXL. Given a user-229 provided prompt, X, the instruction given to GPT-4 is: "For this text X: 1) select using one word 230 ['yes','no'] if text includes any profession or a social media influencer 2) find the subject practicing 231 the job. 3) select using one word ['yes','no'] if the text includes any country, nationality, or race or ethnicity. 4) select using one word ['female', 'male', 'unknown'] the subject's gender" 232

233 234

235 236

237

221

RESULTS 4

4.1 CLASSIFYING GENDER AND RACE

238 To examine the biases in SDXL, we developed a race and gender classifier, composed of three stages: 239 face detection, face embedding generation, and the classification stage (see Methods). To train and validate our classifier, we utilize FairFace Karkkainen & Joo (2021), which specifies the race (Black, 240 East Asian, Indian, Latinx, Middle Eastern, Southeast Asian, and White), and gender (Female, 241 Male) of each image¹. We simplified FairFace's categorization by combining the East and Southeast 242 Asian categories into a single one (Asian). We trained and evaluated our classifier on FairFace's 243 training set and validation set, respectively. We benchmarked against several alternatives from the 244 literature, namely: CLIP's zero-shot Radford et al. (2021), Google's FaceNet Schroff et al. (2015) + 245 SVM Hearst et al. (1998), FairFace's ResNet-34 Karkkainen & Joo (2021), EfficientNet-B7 (tuning 246 all layers) Tan & Le (2019), and Large Vision Transformer VIT (tuning all layers) Dosovitskiy 247 et al. (2020). Section C summarizes the results, showing that our classifier consistently achieves 248 state-of-the-art performance in terms of accuracy, precision, recall, and F1 score.

249 250

251

4.2 EXAMINING BIASES IN STABLE DIFFUSION

252 We examine the degree to which different races and genders appear in SDXL generated images. 253 To this end, we generated 10,000 images using the following racial- and gender-neutral prompt: "a photo of a person". The distribution of the resulting images is summarized by the dashed bars in 254 Figure 1. As can be seen, White is the most generated race (47% of images), followed by Black 255 (33%). The remaining races are rarer in comparison, e.g., 3% are Asian, and 5% are Indian. As for 256 gender, males appear more frequently (65%). These findings mirror what was reported by Ghosh 257 & Caliskan (2023), showing that most images generated by Stable Diffusion (v2.1) representing the 258 prompt "a front-facing photo of a person" depict light-skinned Western men. 259

One possible explanation behind these results could be that SDXL is merely reflecting the biases 260 already present in the dataset on which it was trained, namely LAION-5B Schuhmann et al. (2022). 261 To examine this possibility, we used a subset of LAION-5B Beaumont (2021) consisting of 88,714 262 images (see Material and Methods). The distribution of those images is summarized by the plain 263 bars in Figure 1. As can be seen, images depicting White individuals are more frequent in LAION-264 5B than SDXL (63% vs. 47%), while those depicting other races are less frequent. As for gender, 265 both male and female individuals appear in LAION-5B with equal probability. This indicates that 266 SDXL contains biases that cannot be fully explained by the data it was trained on. 267

²⁶⁸ ¹These classes are widely studied in literature due to their availability from FairFace. It should be noted, 269 however, that these classes do not capture the full spectrum of genders and races. Moreover, existing image classifiers (including our own) are meant to predict the perceived (rather than identified) gender and race.

274

275

278

279

281

282

283

284



Figure 1: Examining gender and race distributions in LAION-5B, SDXL, and our SDXL-Inc in a sample of 88,714 images from the LAION-5B dataset, 10k images generated by SDXL, and 10k generated by SDXL-Inc. For the latter 20k images, we used the prompt: "*a photo of a person*".

Some users may consider the racial distribution of SDXL-generated images to be unsatisfactory, especially if it underrepresents certain groups compared to the society in which these users reside. One way to address this issue is to develop a solution that allows users to specify their desired distributions of race and gender. Perhaps the most intuitive target distribution is the one in which different groups are represented equally. With this in mind, we introduce such a debiasing solution, and test its ability to represent genders and races equally, although the same techniques can be used with any given target distribution. Specifically, we introduce a fine-tuned version of SDXL, which we call: "SDXL-Inc" (Inc stands for Inclusive).

285 Our model was created as follows. First, we identified 32 professions, and split them into 21 for 286 fine-tuning, and 11 for testing (these professions are listed in Appendix Table 1). Then, for ev-287 ery combination (X, Y) such that $X \in \{Black, White, Asian, Indian, Latinx or Hispanic, Middle \}$ 288 Eastern} and $Y \in \{\text{male}, \text{female}\}$, we generated a separate dataset consisting of images depicting 289 the 21 professions. The generation of these 12 datasets (6 races \times 2 genders) was done using the 290 prompt: "a photo of X Y Z looking at the camera, closeup headshot facing forward, ultra quality, 291 sharp focus", where Z denotes one of the 21 professions. After that, we fine-tuned SDXL with each dataset, yielding 12 different sets of weights. The fine-tuning was done using LORA (Low Rank 292 Adaptation) Hu et al. (2021), which reduces the number of trainable parameters for the downstream 293 task (see Materials and Methods). The basic idea of SDXL-Inc is to randomly select one of those 12 sets of weights based on the target distribution of interest (which is uniform in our experiments). 295

296 To evaluate SDXL-Inc, we used it to generate 10,000 images, relying on the same prompt used 297 earlier, i.e., "a photo of a person". The distribution of the resulting images is summarized by the starred bars in Figure 1. As shown in this figure, all races are almost equally represented, and the 298 differences between them are markedly smaller than the differences present in both LAION-5B and 299 SDXL. As for gender, SDXL-Inc is able to represent males and females equally, unlike SDXL. 300

301 302

303

4.3 EXAMINING PROFESSIONAL STEREOTYPES IN STABLE DIFFUSION

304 We used SDXL to generate 320,000 images depicting 32 professions (10,000 per profession), using the prompt: "a photo of a Z, looking at the camera, closeup headshot facing forward, ultra quality, 305 sharp focus", where Z is the profession. Moreover, to specify the types of images that we want to 306 avoid, we used the following negative prompt: "cartoon, anime, 3d, painting, b&w, low quality". For 307 each profession, we used our classifier to examine the racial- and gender- composition of images. 308

309 Figure 2a depicts the results for 24 professions, while Figure 6 in the Appendix depicts the remaining eight professions. Numeric values below 15% are omitted to improve the visualization (see 310 Table 2 in the Appendix for all values). As can be seen, White is the most frequently generated race 311 in 21 out of the 24 professions. As for the remaining three, two of them are among the least prestige 312 occupations, namely Cleaner and Security Guard Hofmann et al. (2024), and both are mostly repre-313 sented by images depicting Black individuals. These findings confirm the findings of Bianchi et al. 314 (2023), showing that prestigious, high-paying professions are often represented as White. 315

Figure 2b shows the gender distribution (see Table 2 in the Appendix for exact values). Males 316 represent 90% of the images in 16 professions, including Doctor and Professor-two prestigious 317 occupations Hofmann et al. (2024). Moreover, jobs in which women are more represented include 318 Nurse and Secretary—two common stereotypes of women Friedrich et al. (2023); Wang et al. (2023). 319

320 321

- 4.4 DEBIASING STABLE DIFFUSION ACROSS PROFESSIONS AND ATTRIBUTES
- Next, we examine stereotypes in terms of the following attributes: Winner, Beautiful, Intelligent, 323 Parent, Sibling, Terrorist, Poor, and Criminal. The results are depicted in the upper row of Fig-

ccountant

337

341 342

343

344 345

347

348 349

324

325

326

327

a Dietitian



ournalist

Secretary

harmacist

15.09%

18 26%

15.7%

18.06%

Male

Femal

Soldier

Figure 2: SDXL's professional stereotypes. SDXL was used to generate 10,000 images for each of the 25 professions. **a**, Racial distribution per profession. **b**, Gender distribution per profession.

350 ure 3a. White dominates the three attributes in our analysis that tend to be associated with success 351 and attractiveness (Winner, Beautiful, and Intelligent). White also dominates the two family-related attributes (Parent and Sibling). In contrast, when it comes to Terrorism, Middle Eastern are the 352 most common race, and none of the images depict a White individual, reinforcing existing stereo-353 types Kundnani (2014). Similarly, when it comes to crime and poverty, the majority of images depict 354 Black individuals Quillian & Pager (2001). 355

356 Having established that SDXL exhibits biases in terms of attributes, we now evaluate SDXL-Inc's 357 ability to address these biases. We repeated the same procedure used earlier with SDXL, but using 358 our SDXL-Inc instead. The results for professions are depicted in the bottom row of Figure 3a. As shown, races are represented more uniformly compared to SDXL, as evidenced by the substantial 359 reduction in standard deviation (σ). Importantly, none of the above eight attributes were used in 360 the fine-tuning phase, and yet SDXL-Inc was able to significantly reduce the racial biases related to 361 them. This indicates that SDXL-Inc can be generalized beyond the features it was fine-tuned on. 362

To further assess the generalizability of SDXL-Inc, we selected four of the professions used during 363 the fine-tuning phase (Dietitian, Manager, Pharmacist, and Pilot), and four professions that were 364 not used during that phase (Accountant, Journalist, Musician, and Firefighter). Figure 3b shows 365 the racial distribution of images using SDXL (upper row) and SDXL-Inc (lower row). As can be 366 seen, regardless of whether the profession is Black-dominated (Musician and Firefighter) or White-367 dominated, our solution is able to significantly reduce difference between races. Similar trends are 368 observed when examining the remaining 24 professions (see Figure 7 and Table 4 in the Appendix). 369

Finally, we evaluate SDXL-Inc's ability to address gender biases across the aforementioned at-370 tributes and professions. As shown in Figure 3c, SDXL (solid bars) exhibits substantial biases, e.g., 371 depicting the vast majority of intelligent individuals as male. In contrast, our solution consistently 372 produces an equal, or near-equal, split between female and male (dashed bars). These improvements 373 are reflected by the vast reduction in standard deviation, from 40.3 to just 2.7. 374

375 Additionally, we experimented with an alternative debiasing solution using a Large Language Model, namely GPT-4 OpenAI (2023), "in-the-loop" (see the Methods section). The results are 376 depicted in Figure 8, and the exact values are listed in Table 5 in the Appendix. This solution is also 377 capable of drastically reducing the race and gender biases exhibited by SDXL.



Figure 3: Results of SDXL-Inc. Given eight attributes and eight professions, both SDXL and SDXL-Inc were used to generate 10,000 images per profession and per attribute. **a**, Race distribution per attribute, with the upper row corresponding to SDXL, and the lower row corresponding to SDXL-Inc. **b**, The same as (**a**) but for professions instead of attributes. **c**, Gender distribution per profession and per attribute for SDXL and SDXL-Inc. The standard deviation(s) corresponding to each subplot is denoted by σ followed by a subscript indicating the model.



Figure 4: Cosine similarity of images generated using SDXL (plain) and SDXL-Div (dashed).

4.5 ADDRESSING RACIAL HOMOGENIZATION

We aimed to develop a version of SDXL capable of generating images with greater facial diversity
per race compared to SDXL. To this end, we downloaded images from Flickr-Faces-HQ, which
hosts high-resolution (1024×1024) images of human faces with "considerable variation in terms of
age, race, and image background" Karras & Hellsten (2023). Since the dataset is unlabeled, we used
our classifier to predict race (see Table 3 in the Appendix for the number of images per race). The
resulting labelled dataset was then used to fine-tune SDXL, this was done using LORA Hu et al.
(2021) to reduce the number of trainable parameters for the downstream task. We call the resulting
model "SDXL-Div" (where Div is a shorthand for Diversity).

432 In our evaluation, each of the two models (SDXL and SDXL-Div) was used to generate $\sim 10,000$ 433 images per race. Our racial classifier was then used to obtain an embedding for each image. Finally, 434 we computed the cosine similarity between every pair of images that have the same race and are 435 generated by the same model, resulting in \sim 50 million cosine similarity values per race per model. 436 As shown in Figure 4, SDXL-Div is able to increase the facial diversity of SDXL, regardless of race. The greatest difference is observed for Middle Eastern (dropping the mean cosine similarity from 437 0.61 to 0.41) and Latinx (from 0.55 to 0.39); see Figures 20a and 20b in the Appendix for sample 438 images of Middle Eastern individuals generated using SDXL and SDXL-Div, respectively. 439

440 441

442

4.6 USER STUDY

We ran four user studies to determine whether exposure to AI-generated faces can affect people's racial and gender biases. We estimated that to obtain a power of 0.8 to detect a medium effect size (Cohen's d) of 0.5 in a paired-sample comparison, a sample of 135 participants would be needed.
We recruited participants on Prolific, with prescreening settings of US resident and English as native language. Additionally, we prohibited participation more than once in our experiment. These studies were preregistered at AsPredicted. All studies were approved by the Institutional Review Board under the category of Exempt or Expedited Research.

In each of the four studies, participants are presented with six AI-generated images, answer a few questions about each image, and then answer an overall question. Each study has four different conditions, depending on whether the images are inclusive or not, and whether participants are informed that the images are generated by AI or are produced by an artist. Next, we provide an overview of the four studies; see Section E in the Appendix for the exact wording and images used.

Study 1 examines racial bias. Given a profession $P \in \{\text{chef}, \text{dietitian}, \text{journalist}\}$, participants are presented with six images depicting P. For each image, they are asked to determine the age, perceived gender, and perceived race of the person depicted therein. After seeing all six images, participants answer the question Q_1 : What percentage of P in the US are White? Here, the non-inclusive images are generated using SDXL and they all depict white individuals, whereas the inclusive images are generated using SDXL-Inc, with each races depicted once. In both conditions, three images depict men, while the other three depict women. All images can be seen in Appendix Figure 12.

462 Study 2 examines gender bias, and is similar to Study 1 apart from two changes: (i) The professions 463 are {accountant, math scientist, andtailor}; (ii) after seeing all six images, participants answer the 464 question Q_2 : What percentage of P in the US are men? Here, the non-inclusive images are SDXL-465 generated, and all depict men, while inclusive images are generated by SDXL-Inc, representing men 466 and women equally. In both conditions, races are represented equally (Appendix Figure 14).

467 Study 3 examines racial homogenization. Participants are presented with six images depicting Mid-468 dle Eastern men. For each image, they are asked to describe the age, skin tone, and facial hair of the 469 man featured therein. After seeing all six images, participants answer the question Q_3 : What is your 470 estimation of the percentage of Middle Eastern men who have beards? Here, the non-inclusive im-471 ages are SDXL-generated, and all depict bearded men. In contrast, the inclusive ones are generated 472 by SDXL-Div, and depict men with varying levels of facial hair (Appendix Figure 4).

473 Study 4 also examines racial homogenization, but focuses on women instead of men. Participants are 474 presented with six images depicting Middle Eastern women, and are asked to describe the age, skin 475 tone, and head cover of the woman featured in the image. Finally, participants answer the question 476 Q_4 : What is your estimation of the percentage of Middle Eastern women who wear headcovers? 477 Here, the non-inclusive images are generated by SDXL, and depict women that all wear a head 478 cover. On the other hand, the inclusive ones are generated by SDXL-Div and depict women, half of 479 whom are wearing a head cover while the other half are showing their hair (Appendix Figure 18).

480 Additionally, for each question $Q_i : i \in \{1, 2, 3, 4\}$, we recruited 135 participants from Prolific 481 to answer Q_i without being presented with AI-generated images. Figures 5a to 5d summarize the 482 responses to questions Q1 to Q4, respectively. For all questions, exposure to exposure to SDXL-483 generated images increase bias (apart from Q_1), while exposure to images generated by SDXL-484 Inc reduces bias, compared to the baseline in which participants are not exposed to AI-generated 485 images. The figures also show no significant difference in participants' responses when the images 486 are labelled as being produced by an artist rather than being AI-generated.



Figure 5: User study results. Subfigures **a** to **d** summarize the participants' responses in Studies 1 to 4, respectively. Boxes extend from the lower to upper quartile values, with a horizontal line at the median; whiskers extend to the most extreme values no further than 1.5 times the interquartile range from the box. P values are calculated using the t-test, unless one of the groups does not pass the Shapiro–Wilk test, in which case P values are calculated using the Mann-Whitney U test. *p<0.05; **p<0.01; ***p<0.001; ns = not significant).

5 DISCUSSION

498

499

500

501

502

503 504 505

506

We set out to examine the stereotypes and biases in SDXL, a text-to-image generator used daily by
millions worldwide Fatunde & Tse (2022). To the end, we developed a classifier to predict the race
and gender of any given face image, and demonstrated that it achieves state-of-the-art performance.
Using this classifier, we showed that the vast majority of faces generated by SDXL are White males.
Biases were also found when considering various attributes, e.g., associating beauty with femininity
and intelligence with masculinity. Some of these biases are less severe in the dataset on which Stable
Diffusion was trained, suggesting that certain biases are further exacerbated by the model itself.

514 Biased text-to-image models may contribute to the normalization of gender stereotypes, potentially 515 shaping societal attitudes towards the roles and capabilities of women in various professions. For 516 instance, as we have demonstrated, SDXL mostly depicts secretaries and nurses as women while 517 depicting doctors and professors as men. Bearing in mind that millions worldwide are already using 518 such models daily, addressing gender stereotypes in these models can be crucial. As suggested 519 by Study 2, such stereotypes can be reduced by an inclusive model and can be exacerbated by a 520 non-inclusive one, indicated the potential of AI in alleviating gender inequality.

Biased representations of gender and race may contribute to the creation of content that not only 521 misrepresents certain groups but may also perpetuate discriminatory practices. This could be detri-522 mental in the context of advertising, marketing, and media campaigns, where visuals hold substan-523 tial influence. For example, as our analysis has shown, Stable Diffusion associates low-income jobs 524 such as Cleaner, Janitor, and Security Guard with Black people, while associating higher-prestige 525 jobs such as Doctor, Lawyer, and Professor with White people. These findings reflect what has been 526 reported by Bianchi et al. Bianchi et al. (2023), showing that several of the most prestigious, high-527 paying professions are represented by Stable Diffusion (v1-4) as White. Another example of bias 528 is the association between crime and Black people, as well as the assassination between terrorism 529 and Middle Easterners, which may reinforce existing biases against these racial groups Kundnani 530 (2014); Quillian & Pager (2001). As we have seen in Study 1, whether or not the model is inclusive 531 can affect people's perception of the racial distribution of certain professions, and this effect is likely to grow more pronounced as the use of AI-generated images becomes more widespread. 532

Stable Diffusion's portrayal of people from any given race as resembling one another may reinforce
existing racial stereotypes. For instance, as we have demonstrated, Stable Diffusion mostly depicts
Middle Eastern men as dark-skinned, bearded, and wearing a traditional headdress, and mostly depicts
Middle Eastern women as dark-skinned, wearing a headscarf. Such oversimplified and generalized depictions of a particular racial group can be culturally insensitive, and may misrepresent the
true diversity within that group. They may also lead to feelings of alienation, low self-esteem, and a
sense of being misunderstood. As we have demonstrated in Studies 3 and 4, racial homogenization
can be reduced using inclusive models, and can be exacerbated using non-inclusive ones.

540 REFERENCES

549

550

551

552

553 554

555

556

558

574

575

- Julia Angwin, Lauren Kirchner, Jeff Larson, and Surya Mattu. Machine bias: There's software used across the country to predict future criminals. and it's biased against blacks.
 ProPublica, 2016. URL www.propublica.org/article/machine-bias-risk-assessmentsin-criminal-sentencing.
- 546 Romain Beaumont. img2dataset: Easily turn large sets of image urls to an im 547 age dataset. https://github.com/rom1504/img2dataset/blob/main/
 548 dataset_examples/laion-high-resolution.md, 2021.
 - Federico Bianchi, Pratyusha Kalluri, Esin Durmus, Faisal Ladhak, Myra Cheng, Debora Nozza, Tatsunori Hashimoto, Dan Jurafsky, James Zou, and Aylin Caliskan. Easily accessible text-toimage generation amplifies demographic stereotypes at large scale. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, pp. 1493–1504, 2023.
 - Joy Buolamwini and Timnit Gebru. Gender shades: Intersectional accuracy disparities in commercial gender classification. In Sorelle A. Friedler and Christo Wilson (eds.), *Proceedings* of the 1st Conference on Fairness, Accountability and Transparency, volume 81 of Proceedings of Machine Learning Research, pp. 77–91. PMLR, 23–24 Feb 2018. URL https: //proceedings.mlr.press/v81/buolamwini18a.html.
- Qiong Cao, Li Shen, Weidi Xie, Omkar M Parkhi, and Andrew Zisserman. Vggface2: A dataset for recognising faces across pose and age. In 2018 13th IEEE international conference on automatic face & gesture recognition (FG 2018), pp. 67–74. IEEE, 2018.
- Brian Christian. *The alignment problem: Machine learning and human values*. WW Norton & Company, 2020.
- Jeffrey Dastin. Insight amazon scraps secret ai recruiting tool that showed bias against women.
 Reuters, 2018. URL https://www.reuters.com/article/us-amazon-com-jobsautomation-insight/amazon-scraps-secret-ai-recruiting-tool that-showed-bias-against-women-idUSKCN1MK08G/. Accessed: November 29th, 2023.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
 - Julia Dressel and Hany Farid. The accuracy, fairness, and limits of predicting recidivism. *Science advances*, 4(1):eaao5580, 2018.
- 577 Mureji Fatunde and Crystal Tse. Stability ai raises seed round at \$1 billion value. Bloomberg, 2022.
 578 URL https://www.bloomberg.com/news/articles/2022-10-17/digitalmedia-firm-stability-ai-raises-funds-at-1-billion-value. Accessed: 580 January 15th, 2024.
- 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.
- Sourojit Ghosh and Aylin Caliskan. 'Person'== Light-skinned, Western Man, and Sexualization of
 Women of Color: Stereotypes in Stable Diffusion. *The 2023 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2023.
- Yandong Guo, Lei Zhang, Yuxiao Hu, Xiaodong He, and Jianfeng Gao. Ms-celeb-1m: A dataset and benchmark for large-scale face recognition. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part III 14*, pp. 87– 102. Springer, 2016.
- M.A. Hearst, S.T. Dumais, E. Osuna, J. Platt, and B. Scholkopf. Support vector machines. *IEEE Intelligent Systems and their Applications*, 13(4):18–28, 1998. doi: 10.1109/5254.708428.

606

607

608 609

- 594 Classifier-free diffusion guidance. Jonathan Ho and Tim Salimans. arXiv preprint 595 arXiv:2207.12598, 2022. 596 Valentin Hofmann, Pratyusha Ria Kalluri, Dan Jurafsky, and Sharese King. Ai generates covertly 597 racist decisions about people based on their dialect. Nature, pp. 1-8, 2024.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, 600 and Weizhu Chen. LoRA: Low-Rank Adaptation of Large Language Models. arXiv preprint 601 arXiv:2106.09685, 2021. 602
- 603 Kimmo Karkkainen and Jungseock Joo. Fairface: Face attribute dataset for balanced race, gender, and age for bias measurement and mitigation. In Proceedings of the IEEE/CVF winter conference 604 on applications of computer vision, pp. 1548–1558, 2021. 605
 - Tero Karras and Janne Hellsten. Flickr-faces-hq dataset (ffhq). Github, 2023. URL https: //github.com/NVlabs/ffhg-dataset.
 - Arun Kundnani. The Muslims are coming!: Islamophobia, extremism, and the domestic war on terror. Verso Books, 2014.
- Debbie S Ma, Joshua Correll, and Bernd Wittenbrink. The chicago face database: A free stimulus 612 set of faces and norming data. Behavior research methods, 47:1122–1135, 2015. 613
- 614 Moin Nadeem, Anna Bethke, and Siva Reddy. StereoSet: Measuring stereotypical bias in pretrained 615 language models. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.), Proceed-616 ings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th 617 International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pp. 5356–5371, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/ 618 2021.acl-long.416. URL https://aclanthology.org/2021.acl-long.416. 619
- 620 OpenAI. ChatGPT-Large language model. Online, 2023. URL https://chat.openai.com/. 621 Accessed: January 15th, 2024. 622
- 623 Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe 624 Penna, and Robin Rombach. SDXL: Improving Latent Diffusion Models for High-Resolution Image Synthesis. arXiv preprint arXiv:2307.01952, 2023. 625
- 626 Lincoln Quillian and Devah Pager. Black neighbors, higher crime? the role of racial stereotypes in 627 evaluations of neighborhood crime. American journal of sociology, 107(3):717-767, 2001. 628
- 629 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, 630 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual 631 models from natural language supervision. In International conference on machine learning, pp. 8748-8763. PMLR, 2021. 632
- 633 Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face 634 recognition and clustering. In Proceedings of the IEEE Conference on Computer Vision and 635 Pattern Recognition (CVPR), June 2015. 636
- 637 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi 638 Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. Advances in Neural 639 Information Processing Systems, 35:25278–25294, 2022. 640
- 641 stable-diffusion-xl-base-1.0. stabilityai. https://huggingface.co/stabilityai/ 642 stable-diffusion-xl-base-1.0, 2023. 643
- 644 Mingxing Tan and Quoc Le. EfficientNet: Rethinking model scaling for convolutional neural networks. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), Proceedings of 645 the 36th International Conference on Machine Learning, volume 97 of Proceedings of Ma-646 chine Learning Research, pp. 6105-6114. PMLR, 09-15 Jun 2019. URL https:// 647 proceedings.mlr.press/v97/tan19a.html.

- Jialu Wang, Xinyue Gabby Liu, Zonglin Di, Yang Liu, and Xin Eric Wang. T2IAT: Measuring Valence and Stereotypical Biases in Text-to-Image Generation. The 61st Annual Meeting of the Association for Computational Linguistics (ACL), 2023.
- Cheng Zhang, Xuanbai Chen, Siqi Chai, Chen Henry Wu, Dmitry Lagun, Thabo Beeler, and Fer-nando De la Torre. Iti-gen: Inclusive text-to-image generation. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 3969–3980, 2023a.
 - Cheng Zhang, Xuanbai Chen, Siqi Chai, Henry Chen Wu, Dmitry Lagun, Thabo Beeler, and Fernando De la Torre. Iti-gen: Inclusive text-to-image generation. https://github.com/ humansensinglab/ITI-GEN, 2023b.
 - Kaipeng Zhang, Zhanpeng Zhang, Zhifeng Li, and Yu Qiao. Joint face detection and alignment using multitask cascaded convolutional networks. IEEE Signal Processing Letters, 23(10):1499–1503, 2016. doi: 10.1109/LSP.2016.2603342.



А SUPPLEMENTARY FIGURES

Figure 6: Professional stereotypes in Stable Diffusion XL (SDXL). Given the remaining eight professions that were not shown in the main article, SDXL was used to generate 10,000 images per profession. a, Racial distribution. b, Gender distribution.



Under review as a conference paper at ICLR 2025





795 796 797

798

799

800

801 802

Figure 8: Results of our GPT-in-the-loop solution. For each profession, the race-, and genderdistribution of images generated by our GPT-in-the-loop solution. a, Race distribution. b, Gender distribution.

Lawyer Doctor

Manager Math Scientist YouTuber

Computer Eng Sushi Chef

Tailor

Programmer

Pilot Chef Security Guard Soldier

Professor

Accountant

Pharmacist Joumalist Cleaner

Tiktoker

803 804

- 805
- 806
- 807
- 808 809

Dietitian Fashion Model Singer Teacher

Secretary

Nurse

810 B SUPPLEMENTARY TABLES

ProfessionUsageAccountantGeneralization testingChefFine-tuningCleanerFine-tuningDietitianFine-tuningDietitianFine-tuningDoctorFine-tuningFashion ModelFine-tuningFirefighterGeneralization testingGarbage CollectorGeneralization testingGeologistGeneralization testingJournalistGeneralization testingJournalistGeneralization testingManagerFine-tuningMascianGeneralization testingNurseFine-tuningPharmacistFine-tuningProfessorFine-tuningProfessorFine-tuningSales personGeneralization testingSoldierFine-tuningSushi ChefFine-tuningTailorGeneralization testingTailorGeneralization testingTailorFine-tuningTabe 1: Professions and how they were used in our st	2		
AccountantGeneralization testing Fine-tuning CleanerChefFine-tuning Fine-tuning DietitianDietitianFine-tuning Fine-tuning DoctorFinestianFine-tuning Fine-tuning FirefighterGeneralizationGeneralization testing Generalization testing JournalistGeneralizationGeneralization testing Generalization testing JournalistJanitorGeneralization testing Generalization testing JournalistManagerFine-tuning Fine-tuning ManagerMusicianGeneralization testing Fine-tuning PharmacistPilotFine-tuning Fine-tuning ProfessorProfessorFine-tuning Fine-tuning Sales personSecurity GuardFine-tuning Fine-tuning SingerSoldierFine-tuning Fine-tuning TialorTable 1: Professions and how they were used in our st	3	Profession	Usage
ChefFine-tuning CleanerCleanerFine-tuning Computer EngineerDietitianFine-tuning Pine-tuningDietorFine-tuning Fine-tuningFashion ModelFine-tuning FirefighterGarbage CollectorGeneralization testing Generalization testing JournalistJournalistGeneralization testing HamagerManagerFine-tuning Fine-tuning Mathematics ScientistMusicianGeneralization testing Generalization testing PharmacistPharmacistFine-tuning Fine-tuning PharmacistPilotFine-tuning Fine-tuning ProfessorSales personGeneralization testing SoldierSushi ChefFine-tuning Fine-tuning TailorTable 1: Professions and how they were used in our st	ļ	Accountant	Generalization testing
CleanerFine-tuning Computer EngineerFine-tuning Fine-tuning DietitianDietitianFine-tuning Fine-tuning Fashion ModelFine-tuning Fine-tuning FirefighterGarbage CollectorGeneralization testing Generalization testing JournalistGeneralization testing Generalization testing I JournalistJournalistGeneralization testing HamagerFine-tuning Fine-tuning ManagerMusicianGeneralization testing PharmacistFine-tuning Fine-tuning PharmacistPiotFine-tuning Fine-tuning ProfessorFine-tuning Fine-tuning Fine-tuning Sales personSales personGeneralization testing SoldierFine-tuning Fine-tuning Fine-tuning Fine-tuning TailorSushi ChefFine-tuning Fine-tuning TailorFine-tuning Fine-tuning Fine-tuning Fine-tuning TailorTable 1: Professions and how they were used in our st		Chef	Fine-tuning
Computer EngineerFine-tuningDietitianFine-tuningDoctorFine-tuningFashion ModelFine-tuningFirefighterGeneralization testingGarbage CollectorGeneralization testingJanitorGeneralization testingJournalistGeneralization testingJournalistGeneralization testingManagerFine-tuningManagerFine-tuningMusicianGeneralization testingNurseFine-tuningPharmacistFine-tuningProfessorFine-tuningProgrammerFine-tuningSales personGeneralization testingSoldierFine-tuningSushi ChefFine-tuningTailorFine-tuningTable 1: Professions and how they were used in our st		Cleaner	Fine-tuning
DietitianFine-tuningDoctorFine-tuningFashion ModelFine-tuningFirefighterGeneralization testingGarbage CollectorGeneralization testingGeologistGeneralization testingJanitorGeneralization testingJournalistGeneralization testingLawyerFine-tuningManagerFine-tuningMusicianGeneralization testingNurseFine-tuningPharmacistFine-tuningProfessorFine-tuningProgrammerFine-tuningSales personGeneralization testingSoldierFine-tuningSushi ChefFine-tuningTailorFine-tuningTable 1: Professions and how they were used in our st		Computer Engineer	Fine-tuning
DoctorFine-tuningFashion ModelFine-tuningFirefighterGeneralization testingGarbage CollectorGeneralization testingGeologistGeneralization testingJanitorGeneralization testingJournalistGeneralization testingLawyerFine-tuningManagerFine-tuningMusicianGeneralization testingNurseFine-tuningPharmacistFine-tuningPilotFine-tuningProfessorFine-tuningSales personGeneralization testingSoldierFine-tuningSoldierFine-tuningSoldierFine-tuningTailorFine-tuningTailorFine-tuningTailorFine-tuningTable 1: Professions and how they were used in our st		Dietitian	Fine-tuning
Fashion ModelFine-tuningFirefighterGeneralization testingGarbage CollectorGeneralization testingGeologistGeneralization testingJanitorGeneralization testingJanitorGeneralization testingJournalistGeneralization testingLawyerFine-tuningManagerFine-tuningMusicianGeneralization testingNurseFine-tuningPharmacistFine-tuningPilotFine-tuningProfessorFine-tuningSales personGeneralization testingSoldierFine-tuningSingerGeneralization testingSoldierFine-tuningTailorFine-tuningTikTokerGeneralization testingTable 1: Professions and how they were used in our st		Doctor	Fine-tuning
FirefighterGeneralization testingGarbage CollectorGeneralization testingGeologistGeneralization testingJanitorGeneralization testingJournalistGeneralization testingLawyerFine-tuningManagerFine-tuningMathematics ScientistFine-tuningMusicianGeneralization testingNurseFine-tuningPharmacistFine-tuningPilotFine-tuningProfessorFine-tuningSales personGeneralization testingSoldierFine-tuningSoldierFine-tuningSushi ChefFine-tuningTailorFine-tuningTable 1: Professions and how they were used in our st		Fashion Model	Fine-tuning
Garbage Collector GeologistGeneralization testing Generalization testing JournalistGeneralization testing Generalization testing LawyerJournalistGeneralization testing LawyerFine-tuning Fine-tuning ManagerManagerFine-tuning Fine-tuning MusicianGeneralization testing Fine-tuning Fine-tuning PharmacistNurseFine-tuning Fine-tuning PharmacistFine-tuning Fine-tuning ProfessorProfessorFine-tuning Fine-tuning Sales personGeneralization testing Fine-tuning Fine-tuning SoldierSoldierFine-tuning Fine-tuning Sushi ChefFine-tuning Fine-tuning Fine-tuning Fine-tuning Fine-tuning Fine-tuning Sushi ChefTailorFine-tuning Fine-tuning TikTokerGeneralization testing Generalization testing Fine-tuning Fine-tuning Fine-tuning Fine-tuningTable 1: Professions and how they were used in our st		Firefighter	Generalization testing
GeologistGeneralization testing JanitorGeneralization testing Generalization testing LawyerLawyerFine-tuning ManagerManagerFine-tuning Fine-tuning Mathematics ScientistMusicianGeneralization testing Fine-tuning PharmacistNurseFine-tuning Fine-tuning PharmacistPilotFine-tuning Fine-tuning ProfessorProfessorFine-tuning Fine-tuning Sales personSecretaryFine-tuning Fine-tuning SoldierSoldierFine-tuning Fine-tuning Fine-tuning Sushi ChefTailorFine-tuning Fine-tuning TikTokerTable 1: Professions and how they were used in our st		Garbage Collector	Generalization testing
Janitor Generalization testing Journalist Generalization testing Lawyer Fine-tuning Manager Fine-tuning Mathematics Scientist Fine-tuning Musician Generalization testing Nurse Fine-tuning Pharmacist Fine-tuning Pilot Fine-tuning Professor Fine-tuning Programmer Fine-tuning Sales person Generalization testing Secretary Fine-tuning Security Guard Fine-tuning Soldier Fine-tuning Sushi Chef Fine-tuning Tailor Fine-tuning Tailor Fine-tuning TikToker Generalization testing Tv Presenter Generalization testing YouTuber Fine-tuning Table 1: Professions and how they were used in our st		Geologist	Generalization testing
Journalist Generalization testing Lawyer Fine-tuning Manager Fine-tuning Mathematics Scientist Fine-tuning Musician Generalization testing Nurse Fine-tuning Pharmacist Fine-tuning Pilot Fine-tuning Professor Fine-tuning Sales person Generalization testing Security Guard Fine-tuning Singer Generalization testing Soldier Fine-tuning Sushi Chef Fine-tuning Tailor Fine-tuning Tailor Fine-tuning TikToker Generalization testing TV Presenter Generalization testing YouTuber Fine-tuning Table 1: Professions and how they were used in our st		Janitor	Generalization testing
LawyerFine-tuningManagerFine-tuningMathematics ScientistFine-tuningMusicianGeneralization testingNurseFine-tuningPharmacistFine-tuningPilotFine-tuningProfessorFine-tuningProgrammerFine-tuningSales personGeneralization testingSecretaryFine-tuningSingerGeneralization testingSoldierFine-tuningSushi ChefFine-tuningTailorFine-tuningTable 1: Professions and how they were used in our st		Journalist	Generalization testing
ManagerFine-tuningMathematics ScientistFine-tuningMusicianGeneralization testingNurseFine-tuningPharmacistFine-tuningPilotFine-tuningProfessorFine-tuningProgrammerFine-tuningSales personGeneralization testingSecurity GuardFine-tuningSingerGeneralization testingSoldierFine-tuningSushi ChefFine-tuningTailorFine-tuningTable 1: Professions and how they were used in our st		Lawyer	Fine-tuning
Mathematics ScientistFine-tuningMusicianGeneralization testingNurseFine-tuningPharmacistFine-tuningPilotFine-tuningProfessorFine-tuningProgrammerFine-tuningSales personGeneralization testingSecretaryFine-tuningSingerGeneralization testingSoldierFine-tuningSushi ChefFine-tuningTailorFine-tuningTikTokerGeneralization testingTowPrine-tuningTable 1: Professions and how they were used in our st		Manager	Fine-tuning
MusicianGeneralization testing NurseNurseFine-tuningPharmacistFine-tuningPilotFine-tuningProfessorFine-tuningProgrammerFine-tuningSales personGeneralization testingSecretaryFine-tuningSecurity GuardFine-tuningSingerGeneralization testingSoldierFine-tuningSushi ChefFine-tuningTailorFine-tuningTakTokerGeneralization testingTo TuberFine-tuningTable 1: Professions and how they were used in our st		Mathematics Scientist	Fine-tuning
NurseFine-tuningPharmacistFine-tuningPilotFine-tuningProfessorFine-tuningProgrammerFine-tuningSales personGeneralization testingSecretaryFine-tuningSecurity GuardFine-tuningSingerGeneralization testingSoldierFine-tuningSushi ChefFine-tuningTailorFine-tuningTailorFine-tuningTikTokerGeneralization testingTV PresenterGeneralization testingYouTuberFine-tuningTable 1: Professions and how they were used in our st		Musician	Generalization testing
PharmacistFine-tuningPilotFine-tuningProfessorFine-tuningProgrammerFine-tuningSales personGeneralization testingSecretaryFine-tuningSecurity GuardFine-tuningSingerGeneralization testingSoldierFine-tuningSushi ChefFine-tuningTailorFine-tuningTacherFine-tuningTik TokerGeneralization testingTo TuberFine-tuningTable 1: Professions and how they were used in our st		Nurse	Fine-tuning
PilotFine-tuningProfessorFine-tuningProgrammerFine-tuningSales personGeneralization testingSecretaryFine-tuningSecurity GuardFine-tuningSingerGeneralization testingSoldierFine-tuningSushi ChefFine-tuningTailorFine-tuningTeacherFine-tuningTikTokerGeneralization testingTV PresenterGeneralization testingYouTuberFine-tuningTable 1: Professions and how they were used in our st		Pharmacist	Fine-tuning
ProfessorFine-tuningProgrammerFine-tuningSales personGeneralization testingSecretaryFine-tuningSecurity GuardFine-tuningSingerGeneralization testingSoldierFine-tuningSushi ChefFine-tuningTailorFine-tuningTeacherFine-tuningTikTokerGeneralization testingTV PresenterGeneralization testingYouTuberFine-tuningTable 1: Professions and how they were used in our state		Pilot	Fine-tuning
ProgrammerFine-tuningSales personGeneralization testingSecretaryFine-tuningSecurity GuardFine-tuningSingerGeneralization testingSoldierFine-tuningSushi ChefFine-tuningTailorFine-tuningTeacherFine-tuningTikTokerGeneralization testingTV PresenterGeneralization testingYouTuberFine-tuningTable 1: Professions and how they were used in our state		Professor	Fine-tuning
Sales personGeneralization testing SecretarySecretaryFine-tuning Security GuardSingerGeneralization testing SoldierSoldierFine-tuning Tine-tuning TailorTailorFine-tuning Fine-tuning TikTokerTokerGeneralization testing Fine-tuning TikTokerToyPresenter Fine-tuning Tine-tuningTable 1: Professions and how they were used in our st		Programmer	Fine-tuning
SecretaryFine-tuningSecurity GuardFine-tuningSingerGeneralization testingSoldierFine-tuningSushi ChefFine-tuningTailorFine-tuningTeacherFine-tuningTikTokerGeneralization testingTV PresenterGeneralization testingYouTuberFine-tuning		Sales person	Generalization testing
Security GuardFine-tuningSingerGeneralization testingSoldierFine-tuningSushi ChefFine-tuningTailorFine-tuningTeacherFine-tuningTikTokerGeneralization testingTV PresenterGeneralization testingYouTuberFine-tuningTable 1: Professions and how they were used in our st		Secretary	Fine-tuning
SingerGeneralization testing SoldierSoldierFine-tuning Fine-tuning TailorTailorFine-tuning Fine-tuning TeacherTeacherFine-tuning Generalization testing TV PresenterTV PresenterGeneralization testing YouTuberTable 1: Professions and how they were used in our st		Security Guard	Fine-tuning
SoldierFine-tuningSushi ChefFine-tuningTailorFine-tuningTeacherFine-tuningTikTokerGeneralization testingTV PresenterGeneralization testingYouTuberFine-tuningTable 1: Professions and how they were used in our st		Singer	Generalization testing
Sushi ChefFine-tuningTailorFine-tuningTeacherFine-tuningTikTokerGeneralization testingTV PresenterGeneralization testingYouTuberFine-tuningTable 1: Professions and how they were used in our st		Soldier	Fine-tuning
TailorFine-tuningTeacherFine-tuningTikTokerGeneralization testingTV PresenterGeneralization testingYouTuberFine-tuningTable 1: Professions and how they were used in our st		Sushi Chef	Fine-tuning
TeacherFine-tuningTikTokerGeneralization testingTV PresenterGeneralization testingYouTuberFine-tuningTable 1: Professions and how they were used in our st		Tailor	Fine-tuning
TikTokerGeneralization testingTV PresenterGeneralization testingYouTuberFine-tuningTable 1: Professions and how they were used in our st		Teacher	Fine-tuning
TV Presenter YouTuber Generalization testing Table 1: Professions and how they were used in our st		TikToker	Generalization testing
YouTuber Fine-tuning Table 1: Professions and how they were used in our st		TV Presenter	Generalization testing
Table 1: Professions and how they were used in our st		YouTuber	Fine-tuning
Table 1: Professions and how they were used in our st			
		Table 1: Professions and how	w they were used in our st

864	Profession	White	M.E.	Latinx	Indian	Asian	Black	Female
865	A	00 5 1	14.02	0.72	1.05	0.01	2.79	
866	Accountant	80.51	14.95	0.75	1.05	0.01	2.78	9.27
867	Cleaner	14.09	5.75 9.21	0.18	2.02	2.23	25.0	10.00
868	Cleaner Commuten Engineer	14.08	0.21	10.00	10.5	3.05	35.32	10.09
869	Distition	28.20	21.14	0.10	23.3	2.21	10.95	1.33
870	Dietitian	87.33	2.82	3.21	0.49	1.07	4.28	92.09
871	Doctor Eachien Medel	02.31 57.02	22.05	2.08	5.92	1.99	0.08	4.65
872	Fashion Model	37.92	0.55	2.38	0.7	1.10	50.0	0.0
873	Combogo Collector	40.5	1.15	1.09	0.00	0.0	10.0	0.0
974	Garbage Collector	20.23	37.03	3.39	10.18	0.0	10.33	0.29
074	Junitor	80.39 22.41	12.01 9.51	1.08	1.05	0.0	52 01	0.17
875	Jaillor	52.41 20.1	0.31	5.71 1.40	2.50	0.0	2 2 1	0.18
876	Journalist	00.1 56.42	14.29	1.49	0.01	0.0	20.65	20.05
877	Lawyer	50.45 65.10	12.0	2.76	4.94	2.20	20.03	1.05
878	Math Scientist	03.19 56.79	10.18	5.70	4.52	5.19 2.96	10.20	4.85
879	Musician	20.78	12.11	4.02	12.79	2.80	10.84	4.20
880	Nurso	20.33	2.07	2.40 6.26	4.5	1.00	52.44 22.59	1.19
881	Dhammaaist	66.60	2.97	0.20	0.80	1.99	25.30	99.29
882	Pharmacist	79 41	0.12	2.09	5.04 0.79	1.92	0.04	28.34
883	Professor	/0.41	9.12	2.22 9.27	0.70	2.99	0.00	0.40
884	Programmer	50.40	13.31	0.57	10.1 8 01	2.00 2.22	10.20	0.50
885	Salesperson	82.30	11.00	4.01	0.91	2.22	5 22	18.07
005	Salesperson	02.39	2 22	6.16	1 30	2.12	0.4	06.48
000	Security Guard	7.03	1.72	1 38	1.39	0.83	9.4 87.54	0.40
887	Scoulty Outld	65 71	2.15	1.30	0.24	0.03	27.13	66.40
888	Soldier	31 13	2.13	62	0.24	0.02	57 37	0.49
889	Suchi Chef	3 31	0.37	3 55	0.02	0.74	1.01	0.01
890	Tailor	31.52	26.23	3.55	5.78	1 20	31 51	0.79
891	Teacher	61.62	20.23 7 44	<i>4</i> 5	6 56	1.29	18.26	35.21
892	TikToker	41.46	5 44	10.97	4 25	6 44	31 44	30.56
893	TV Presenter	90.67	3.00	1 45	0.09	0.11	47	66.14
894	YouTuber	62 79	11 72	673	2 49	2 59	13.68	3.84
895	1001000	52.17	11.12	0.75	2.77	2.57	15.00	0.04

Table 2: SDXL's race and gender distribution per profession.

	Black	Asian	Indian	Latinx	Middle east	White
Number of images	2151	3255	2193	2731	1551	4071

Table 3: Number of images per race for the Flickr-Face-HQ dataset, as determined by our race classifier.

918	Profession	White	M.E.	Latinx	Indian	Asian	Black	Female
919	Accountant	17.43	27.9	5.82	17.21	15 15	16.49	49.86
920	Chef	22 14	20.16	6.56	17.21	16.6	17.24	50.64
921	Cleaner	15 84	19 35	13 79	19 78	15.2	17.24	49.09
922	Computer Engineer	17 22	21.2	7 37	21.65	16.27	16.04	49.07
923	Dietitian	17.22	26.87	6.31	16.61	16.27	16.6	50.13
924	Doctor	16.85	20.07	5.66	17.3	15.07	17.45	50.15
925	Eashion Model	17.16	16 58	14 76	18 15	16 54	16.81	50.65
926	Firefighter	21.23	17.09	12.03	16.15	16.02	17.43	50.00
927	Garbage Collector	17 33	24.01	7 84	17 76	16.02	16.8	49.47
928	Geologist	19.2	24.12	7.27	17.01	16.09	16.31	49.58
929	Janitor	16.43	17.43	13.25	20.41	15.9	16.59	49.51
930	Journalist	17.25	26.12	5.67	17.57	16.73	16.66	48.62
031	Lawver	18.08	25.82	6.64	16.75	15.84	16.87	49.59
022	Manager	17.1	26.82	6.47	17.13	16.15	16.33	49.78
932	Math Scientist	17.1	21.92	6.81	22.08	15.76	16.33	51.83
933	Musician	15.7	21.79	10.58	18.41	16.18	17.35	51.29
934	Nurse	16.73	23.13	9.17	17.78	16.16	17.04	52.2
935	Pharmacist	16.93	27.94	5.67	17.31	15.55	16.6	49.92
936	Pilot	19.36	23.97	7.07	16.38	16.74	16.47	49.86
937	Professor	17.22	24.98	7.11	18.07	15.96	16.65	50.14
938	Programmer	19.63	20.51	8.04	18.88	15.8	17.14	48.92
939	Salesperson	17.54	25.66	7.53	16.43	16.96	15.9	49.55
940	Secretary	17.63	26.42	6.44	18.27	16.05	15.19	51.82
941	Security Guard	15.87	15.75	14.51	17.66	16.41	19.8	48.78
942	Singer	16.51	20.77	12.21	16.77	17.01	16.73	49.96
943	Soldier	16.62	16.19	16.76	16.26	16.73	17.42	50.35
944	Sushi chef	18.27	11.28	21.58	15.92	16.91	16.05	50.89
9/5	Tailor	16.48	24.98	7.65	17.43	17.09	16.37	50.45
046	Teacher	16.76	26.36	6.32	17.23	16.41	16.92	51.25
340	TikToker	14.96	19.31	13.77	18.84	15.91	17.22	48.98
947	TV Presenter	18.58	25.21	7.06	16.06	16.82	16.27	50.35
948	YouTuber	17.07	22.26	10.39	16.49	16.26	17.52	49.37
949								

Table 4: SDXL	-Inc's race a	and gender	distribution	per profession.
		6		

972	Profession	White	ME	Latinx	Indian	Asian	Black	Female
973	11010331011	winte	111.12.	Latinx	manan	7 151411	Diack	1 cillate
974	Accountant	17.07	22.29	10.19	17.4	16.49	16.57	50.04
975	Chef	17.46	19.55	13.2	16.37	16.71	16.71	49.79
976	Cleaner	15.7	16.54	14.51	19.58	16.79	16.88	50.21
977	Comp. Engineer	15.48	19.8	9.9	20.98	17.09	16.75	49.66
079	Dietitian	17.16	22.31	10.61	16.83	16.5	16.58	50.75
970	Doctor	16.86	23.5	9.14	16.94	16.69	16.86	50.0
979	Fashion Model	16.78	16.02	13.51	19.88	16.78	17.03	50.0
980	Journalist	17.7	21.42	10.11	17.14	17.46	16.17	49.15
981	Lawyer	16.64	20.53	11.67	17.05	17.3	16.8	49.92
982	Manager	16.85	22.63	10.31	16.85	16.6	16.76	50.04
983	Math Scientist	16.3	19.24	10.84	19.92	16.72	16.97	50.0
984	Nurse	16.43	20.6	12.01	17.43	16.76	16.76	49.79
985	Pharmacist	16.57	23.15	9.91	16.9	16.82	16.65	49.88
986	Pilot	16.93	17.6	14.93	16.85	16.68	17.01	50.04
987	Professor	17.2	19.03	11.77	18.45	16.78	16.78	50.0
988	Programmer	17.45	19.9	10.88	18.63	16.19	16.95	50.34
080	Secretary	16.68	21.97	10.16	17.01	16.43	17.75	50.78
909	Security Guard	15.77	14.0	17.62	17.96	16.86	17.79	50.42
990	Singer	16.52	19.64	12.9	17.83	16.35	16.76	50.04
991	Soldier	16.46	14.2	17.38	17.96	16.71	17.29	50.13
992	Sushi Chef	16.77	16.6	15.56	16.34	17.54	17.2	49.7
993	Tailor	16.47	21.44	11.75	16.89	16.56	16.89	50.08
994	Teacher	16.78	20.3	11.91	17.7	16.61	16.69	50.17
995	TikToker	17.1	18.17	12.53	18.26	16.6	17.34	50.21
996	YouTuber	17.47	21.4	10.87	16.64	16.81	16.81	50.25

Table 5: GPT-in-the-loop SDXL race and gender distribution per profession.

999 1000 1001

1002

C OUR RACIAL AND GENDER CLASSIFIER

Figure 9a depicts the confusion matrices corresponding to race and gender, along with sample images representing each cell. As can be seen, Asian is the easiest race to predict, followed by Black. In contrast, the hardest races to predict are Latinx and Middle Eastern. On the other hand, there is no difference in performance when it comes to predicting male vs. female images.

1007 We reproduced the same confusion matrices, but using images generated by Stable Diffusion XL 1008 (SDXL) which has been shown to outperform previous versions of Stable Diffusion Podell et al. (2023). More specifically, we generated 10,000 images per race using the following prompt: "a 1009 photo of X person, looking at the camera, closeup headshot facing forward, ultra quality, sharp 1010 *focus*", where $X \in \{a \text{ Black}, a \text{ White, an Asian, an Indian, a Latinx or Hispanic, a Middle Eastern}\}$. 1011 As for gender, we generated 10,000 images per gender using the same prompt as before, but with 1012 $X \in \{\text{male}, \text{ female}\}$. The results of this evaluation are depicted in the first two rows of Table 8, 1013 showing that our classifier works better with images from SDXL than those from FairFace, regard-1014 less of the performance measure. This indicates that our classifier is particularly suited to examine 1015 the racial-, and gender-composition of images generated by SDXL. 1016

The confusion matrices depicted in Figure 9b reveal that all races are significantly easier to predict in the SDXL dataset compared to its FairFace counterpart, apart from Latinx which is significantly harder (see how the numbers along the diagonal are greater in Figure 9b than in Figure 9a, except for Latinx). One possible explanation could be that the dataset used to train SDXL contains several images that are labelled as Latinx but are highly similar to Indian and Middle Eastern individuals. In terms of gender, our classifier achieves a perfect score for both male and female.

1022

1023

1024

Table 6: A comparison between our race classifier and other literature alternatives, using FairFace's validation set and FairFace's seven-race classification. Bold font highlights the highest score(s).

1029		Accuracy	Precision	Recall	F1 score
1030	Our classifier	73%	72%	72%	72%
1031	CLIP's zero-shot classifier	64%	67%	65%	65%
1032	Google's FaceNet + SVM	69%	69%	68%	68%
1033	FairFace's (ResNet34) classifier	72%	72%	71%	72%
1034	EfficientNet-B7 (tuning all layers)	70%	70%	70%	70%
1035	Large Vision Transformer VIT (tuning all layers)	71%	72%	70%	71%

Table 7: A comparison between our **gender classifier** and other alternatives from the literature, using FairFace's validation set. Bold font highlights the highest score(s).

	Accuracy	Precision	Recall	F1 score
Our classifier	94%	94%	94%	94%
FairFace's (ResNet34) classifier	94%	94%	94%	94%
CLIP's zero-shot classifier	94%	94%	94%	94%

Table 8: Evaluating our classifier's ability to predict race and gender using SDXL images and our six racial groups. Bold font highlights the highest score.

1046		Accur	acy	Precision		Recall		F1 score	
1047		FairFace	SDXL	FairFace	SDXL	FairFace	SDXL	FairFace	SDXL
1048	Racial	78%	88%	76%	90%	75%	88%	76%	87%
1049	Gender	94%	100%	94%	100%	94%	100%	94%	100%





D SDXL-INC VS ITI-GEN

In this section, we compare our solution (SDXL-Inc) to another recently proposed alternative, namely ITI-GEN Zhang et al. (2023a).

D.1 RETRAINING ITI-GEN

1087 We needed to benchmark SDXL-Inc against ITI-GEN Zhang et al. (2023a;b) in terms of how well 1088 it can debias Stable Diffusion. To this end, we retrained ITI-GEN with the six races (Asian, Black, 1089 Indian, Latino or Hispanic, Middle Eastern, and White) and the two genders (female and male) that 1090 SDXL-Inc was trained on. The training data used for this purpose was created as follows: First, we inferred the race and gender of each image in the Flickr-Faces-HQ dataset Karras & Hellsten 1091 (2023) using our classifier. Then, we curated two training-sets of images: one for gender and one 1092 for race. More specifically, for the gender training set, we randomly selected 50 images (from the 1093 aforementioned labeled Flickr dataset) for each gender. On the other hand, for the race training set, 1094 we randomly selected 25 images for each race. 1095

To validate our resultant ITI-GEN model, we generated 1,200 images (100 per gender-race combination) using the prompt "*a photo of a person*". Figure 10 depicts the race and gender distribution across the 1,200 generated images. The results demonstrate that our trained version of ITI-GEN is capable of generating equal representation for gender, and nearly equal representation for race (apart from Latinx, due to their facial similarity to both White and Middle Eastern).



Figure 10: Race and Gender distribution obtained by the retrained ITI-GEN model. Comparing the representation of different races and genders in a sample of 1200 images generated by the retrained ITI-GEN model using the prompt: *"a photo of a person"*.

1114 1115

1116 D.2 PERFORMANCE COMPARISON

1117 In ITI-GEN official GitHub repository Zhang et al. (2023b), the authors imply that their solution may 1118 struggle with certain generalizations. More specifically, they state that the training prompt (e.g., a 1119 headshot of a person) and the inference prompt (e.g., a headshot of a doctor) "should not differ a 1120 lot". To compare the degree to which both models (ITI-GEN and SDXL-Inc) can be generalized, we 1121 used each of them to generate 1,200 images for each of the following five prompts: 1) a headshot 1122 of a person with the Eiffel Tower in the background; 2) a headshot of a mechanic with a guitar; 3) a 1123 headshot of a skillful trainer with a pet tiger; 4) a headshot of a happy and healthy family; and 5) a 1124 headshot of a person with green hair and eyeglasses.

1125 Figure 11 depicts the race and gender distribution of SDXL, SDXL-Inc, and ITI-GEN for each 1126 of the five prompts. As can be seen in the left column, SDXL is extremely biased, depicting the 1127 vast majority of images as White in all prompts. As for ITI-GEN (middle column, dashed bars), it 1128 manages to reduce a substantial amount of the bias compared to SDXL. Nevertheless, the percentage 1129 of images depicting White individuals remains $\geq 40\%$ in four out of the five prompts. Our SDXL-1130 Inc (middle column, starred bars) yields a more even distribution of images across races compared 1131 to ITI-GEN; see how the standard deviations for our solution are smaller than those of ITI-GEN in all five prompts. As for gender (right column), we can see that SDXL is the most biased, ITI-GEN 1132 reduces a considerable amount of bias, and SDXL-Inc further reduces bias compared to ITI-GEN 1133 across all prompts; see the corresponding standard deviations.

One reason behind these observed differences in performance could be the fact that the models occasionally intend to generate images of a certain race, but accidentally end up generating images of a different race (note that both SDXL-Inc and ITI-GEN contain steps designed to generate an image of certain races). If these steps work perfectly well, we would expect the models to achieve a perfect score in terms of accuracy, recall, precision, and F1 score. To test this possibility, we used out classifier to label the 1,200 images generated by each solution (200 per race) for each of the five prompts. The result of this analysis is summarized in Table 9. As can be seen, SDXL-Inc outperforms ITI-GEN across all metrics in all five prompts. This suggests that the observed difference in performance between the two solutions could be (at least partially) explained by the fact that the error rate in generating the intended race(s) is higher in ITI-GEN than in SDXL-Inc.



Figure 11: SDXL vs. ITI-GEN vs. SDXL-Inc. Comparing the distribution of race and gender across the three models, given five prompts. Each row corresponds to a different prompt. The standard deviation(s) corresponding to each subplot is denoted by σ followed by a subscript indicating the model. Gender results are omitted from the bottom row since the generated images depict families rather than individuals.

Table 9: Comparing SDXL-Inc to ITI-GEN in terms of accuracy, recall, precision, and F1 score. For each performance measure, the highest score is highlighted in bold.

	Accuracy		Rec	all	Preci	F1 score		
Prompt	SDXL- Inc	ITI- GEN	SDXL- Inc	ITI- GEN	SDXL- Inc	ITI- GEN	SDXL- Inc	
a person with Eiffel tower in the background	87	70	87	70	89	71	87	
a mechanic with a guitar	86	39	86	39	88	48	86	
a skillful trainer with a pet tiger	77	37	77	37	83	47	77	
a happy and healthy family	70	66	70	66	75	68	70	
a person with green hair and eyeglasses	60	40	60	40	72	55	60	
Average	76	50.4	76	50.4	81.4	57.8	76	Î

Е USER STUDY IMAGES AND SAMPLE SCREENSHOTS a b .

Figure 12: Study 1 images (racial). a, SDXL. b, SDXL-Inclusive.





Figure 13: Study 1 sample screenshots. (1) welcome screen, (2) loading screen, (3) per image

questions, and (4)/(5) Final questions after seeing all six images.



Figure 14: Study 2 images (gender). a, SDXL. b, SDXL-Inclusive.





Figure 15: Study 2 sample screenshots. (1) welcome screen, (2) loading screen, (3) per image

questions, and (4)/(5) Final questions after seeing all six images.



Figure 17: Study 3 sample screenshots. (1) welcome screen, (2) loading screen, (3) per image questions, and (4)/(5) Final questions after seeing all six images.



Figure 19: Study 4 sample screenshots. (1) welcome screen, (2) loading screen, (3) per image questions, and (4)/(5) Final questions after seeing all six images.

F SDXL VS. SDXL-DIV SAMPLE IMAGES



Figure 20: Sample images of Middle Eastern individuals generated using SDXL (a) and SDXL-Div (b).