# BIAS ANALYSIS IN UNCONDITIONAL IMAGE GENERATIVE MODELS

Anonymous authors

Paper under double-blind review

#### ABSTRACT

The widespread usage of generative AI models raises concerns regarding fairness and potential discriminatory outcomes. In this work, we define the bias of an attribute (e.g., gender or race) as the difference between the probability of its presence in the observed distribution and its expected proportion in an ideal reference distribution. Despite efforts to study social biases in these models, the origin of biases in generation remains unclear. Many components in generative AI models may contribute to biases. This study focuses on the inductive bias of unconditional generative models, one of the core components, in image generation tasks. We propose a standardized bias evaluation framework to study bias shift between training and generated data distributions. We train unconditional image generative models on the training set and generate images unconditionally. To obtain attribute labels for generated images, we train a classifier using ground truth labels. We compare the bias of given attributes between generation and data distribution using classifier-predicted labels. This absolute difference is named bias shift. Our experiments reveal that biases are indeed shifted in image generative models. Different attributes exhibit varying bias shifts' sensitivity towards distribution shifts. We propose a taxonomy categorizing attributes as *subjective* (high sensitivity) or *non-subjective* (low sensitivity), based on whether the classifier's decision boundary falls within a high-density region. We demonstrate an inconsistency between conventional image generation metrics and observed bias shifts. Finally, we compare diffusion models of different sizes with Generative Adversarial Networks (GANs), highlighting the superiority of diffusion models in terms of reduced bias shifts.

032 033

003

010 011

012

013

014

015

016

017

018

019

021

025

026

028

029

031

034 035

### 1 INTRODUCTION

037

Generative AI models have achieved realistic generation qualities for various modalities including text (Touvron et al., 2023; OpenAI, 2023), image (Ramesh et al., 2022; Rombach et al., 2022; Esser et al., 2024), audio (Kreuk et al., 2023), and video (Ho et al., 2022; Singer et al., 2023). They are consequently employed for commercial uses and are available to every internet user across the world. The widespread use of these high-performing models, along with the potential social biases embedded in their generation, increase the risk of discriminatory outcomes. Taking image generation as an example, Growcoot (2023) and Tiku et al. (2023) report racial and gender biases in popular text-to-image (T2I) systems including DALL-E (Ramesh et al., 2021), Stable Diffusion (Rombach et al., 2022) and Midjourney (https://www.midjourney.com).

We define the bias of an attribute (e.g., gender or race) as the difference between the probability of its presence in the observed distribution and its expected proportion in an ideal reference distribution. The ideal reference distribution may be based on social norms or population statistics, etc. A widely
studied problem is gender or racial bias with respect to occupations (Cho et al., 2023; Bianchi et al., 2023; Luccioni et al., 2023; Friedrich et al., 2024). Depending on the context, previous works use equality or U.S. labor statistics as the ideal reference distribution. In these cases, the typical analysis protocol consists of generating facial images based on text prompts containing occupation information, using pre-trained models to assign gender or racial labels for the generated images, followed by measuring and assessing the degree of biases in the images given a certain occupation.



Figure 1: Illustrations depicting bias shift. The plots represent the distributions of samples with respect to the likelihood of an attribute (solid for training data, dashed for generation). The decision boundary (brown) binarizes the likelihood into positive and negative classes. In each subfigure, the generation distribution is translated from the training. Bias shift is the difference between red and blue areas. When the boundary falls in a low-density region (Figs. 1c and 1d), the bias shifts tend to be small, and vice versa (Figs. 1a and 1b). Detailed discussion is in Section 4.4 with distributions obtained from real datasets.

Other studies have compared social biases between generated images and training datasets of generative AI models, with mixed findings. Friedrich et al. (2024) report that images generated by Stable Diffusion (Rombach et al., 2022) show cases of bias and even bias amplification compared to the training data (LAION-5B) (Schuhmann et al., 2022). On the other hand, Seshadri et al. (2023) conduct similar experiments and discover that bias shift can be mainly attributed to discrepancies between training captions and model prompts.

Although analyzing biases empirically in publicly available generative AI models is of practical sig-077 nificance, identifying the origin of these biases remains a challenge. Modern generative AI systems are complex and generative biases can stem from various sources, such as biased datasets (Schuh-079 mann et al., 2022; Karkkainen & Joo, 2021), the conditioning process (including textual prompts, and guidance (Dhariwal & Nichol, 2021; Ho & Salimans, 2022)), pre-trained modules (including 081 CLIP (Radford et al., 2021) and VAE (Kingma & Welling, 2014)), and inductive bias of the gener-082 ative models (e.g., diffusion process (Ho et al., 2020), generative adversarial training (Goodfellow 083 et al., 2014)). While biases in pre-trained models (Bommasani et al., 2021; Alabdulmohsin et al., 084 2024) and datasets (Schuhmann et al., 2022) have been widely studied, the impact of inductive biases 085 in generative models remains underexplored. Thus, in our experiments, we focus on *unconditional* pixel-level image generative models without any guidance during training or inference.

087 We propose a standardized evaluation framework that employs attribute classifiers to study bias 880 shifts from training to generated data distributions in unconditional image generative models. Train-089 ing the classifiers requires ground-truth labels for the training and validation sets; hence, our frame-090 work is applicable to any supervised learning dataset. We train unconditional image generative mod-091 els using the training set and unconditionally generate images. We then use the trained classifiers to 092 predict attribute labels for each generated image. We compare the bias for each attribute between the training and generated data distributions using classifier-predicted labels. We refer to this absolute difference as the *bias shift*. If bias shift is close to zero, there is no systematic bias exhibited in 094 image generative models. We analyse the bias shifts on two real image datasets, CelebA (Liu et al., 095 2015) and DeepFashion (Liu et al., 2016). 096

Our findings reveal that bias shifts vary in magnitude across different attributes, indicating varying levels of sensitivity to distribution change between generation and training data. We categorize attributes as *subjective* (high sensitivity) and *non-subjective* (low sensitivity) sets, based on the relative sample density at the classifier's decision boundary. If the classifier is confident in its predictions — in other words, the decision boundary lies in a lower-density region (corresponding to *non-subjective* attributes), bias shifts tend to be smaller, and vice versa. Fig. 1 shows translation distribution shift as an example to introduce this idea.

Our bias analysis framework yields the following observations: 1) Biases of attributes shift between training and generation distributions for unconditional image generative models. The magnitude of bias shift is correlated with the *subjectivity* of the attribute. 2) Selecting the checkpoint based on image generation metrics considering quality, diversity, and novelty (FID (Heusel et al., 2017), KID (Binkowski et al., 2018), and FLD (Jiralerspong et al., 2023)) does not guarantee the smallest bias shifts. Bias should be treated as an independent issue when evaluating generations. 3) BigGAN models (Brock et al., 2019) have larger bias shifts compared to diffusion models, despite having similar image generation metrics. 4) Bias shifts in smaller diffusion models tend to increase over extended training steps, even though image generation metrics remain relatively stable.

112 113

114

## 2 RELATED WORKS

115 **Bias in Image Generation** Previous studies focus on social biases in image generation, often 116 concluding that these models are unfair (Friedrich et al., 2024; Cho et al., 2023) or fail to reflect 117 real-world biases as observed in U.S. labor statistics (Luccioni et al., 2023; Bianchi et al., 2023). A 118 commonly examined bias is gender or race to occupation. However, different studies select various 119 public models and develop their own evaluation benchmarks. For example, Cho et al. (2023) investi-120 gate minDALL-E (Kim et al., 2021), Karlo (Lee et al., 2022), and Stable Diffusion (Rombach et al., 2022) v1.4, while Luccioni et al. (2023) use Stable Diffusion v.1.4, v.2, and Dall-E 2 (Ramesh et al., 121 2022). However, there are different components in T2I models that may contribute to generative 122 biases. We focus on the unconditional image generative model that directly outputs the generations, 123 without text conditioning or guidance. 124

125 **Bias Shift between Train and Generation** Few studies attempt to compare bias between the 126 generation and training distributions. These efforts often rely on publicly available Stable Diffusion 127 models, comparing generated outputs with the LAION-5B training set (Schuhmann et al., 2022), a 128 large-scale dataset lacking explicit attribute labels. Given a text prompt, Friedrich et al. (2024) select 129 a subset of LAION-5B based on pre-trained image-prompt similarity, then compare the bias between 130 this subset and the images generated using the same prompt. In contrast, Seshadri et al. (2023) select 131 subsets based on keywords in image captions, which may overlook relevant images. To avoid this 132 large-scale dataset search and subset comparison, we train generative models using datasets with labeled attributes, ensuring reliable bias estimation across both the training and generation. 133

Bias-related Attribute Label Prediction To calculate bias in generation, the generated images 135 need to be assigned attribute labels, which is non-trivial in the case of unconditional generation. 136 Some studies (Bianchi et al., 2023) infer the labels in the representation space of self-supervised 137 learning models, for example, CLIP (Radford et al., 2021). Some methods use pre-trained vision 138 language models and conduct zero-shot text generation. Cho et al. (2023) use BLIP-2 (Li et al., 139 2023) and get the label through visual question answering (VQA). Luccioni et al. (2023) use BLIP 140 with VQA task and ViT (Dosovitskiy et al., 2021) with image captioning task. However, pre-trained 141 models introduce their own biases (Bommasani et al., 2021; Alabdulmohsin et al., 2024), rendering 142 the predicted labels unreliable for accurate bias evaluation. Some approaches (Friedrich et al., 2024) 143 train an attribute classifier on other available supervised learning datasets. In our case, we train the 144 classifier on the same dataset used for bias analysis, resulting in more accurate predictions.

145 146 147

159

160

134

## **3** BIAS EVALUATION METHOD

# 148 3.1 BIAS DEFINITION

In this work, bias for an attribute is defined as the difference between the probability of its presencein the observed distribution and its expected proportion in an ideal reference distribution.

Considering a set of binary attributes<sup>1</sup> C for which we want to study bias, each image in the dataset is annotated for every attribute. Given an attribute  $C \in C$ , we consider the positive case C = 1 in the following. We can set an ideal probability  $P^{\text{ideal}}(C = 1)$  for attribute C as the reference probability, depending on the context. We denote the probability of this attribute in the data distribution as  $P^{\text{data}}(C = 1)$ . We can use either  $P^{\text{train}}(C = 1)$  or  $P^{\text{val}}(C = 1)$  as an estimation for  $P^{\text{data}}(C = 1)$ and compare with the reference probability to determine degree of bias. For example, we define the bias of the data distribution relative to  $P^{\text{ideal}}(C = 1)$  as

$$B^{\text{data}}(C=1) = P^{\text{data}}(C=1) - P^{\text{ideal}}(C=1).$$
(1)

<sup>&</sup>lt;sup>1</sup>The use of binary attributes can be extended to K-way attributes by binarizing the K-way attributes as K 1-vs-all binary attributes.



Figure 2: **Bias evaluation framework.** Unconditional generative models are trained on the training set. The pre-trained classifier is fine-tuned on the training set and validated on the validation set using ground truth labels and is then used to classify training, validation, and generation sets. The bias evaluation metrics are calculated based on the classifier-predicted labels.

To get the bias on the generation set, we need to calculate the proportion for this attribute in the generation set  $P^{\text{gen}}(C = 1)$ . We can then measure the bias in the generation

$$B^{\text{gen}}(C=1) = P^{\text{gen}}(C=1) - P^{\text{ideal}}(C=1).$$
<sup>(2)</sup>

We also define the conditional bias. Given a binary anchor attribute  $A \in C$ , the bias of attribute Cconditioned on A = 1 is the conditional probability P(C = 1 | A = 1). Similarly, the conditional bias in the data distribution is  $P^{\text{data}}(C = 1 | A = 1) - P^{\text{ideal}}(C = 1 | A = 1)$ . The conditional bias in the generation distribution is  $P^{\text{gen}}(C = 1 | A = 1) - P^{\text{ideal}}(C = 1 | A = 1)$ .

#### 3.2 BIAS EVALUATION FRAMEWORK

173

174

175

176 177 178

179

180 181

187

188

203

Fig. 2 illustrates our proposed bias evaluation framework. We train image generative models for unconditional image generation using only images from the training set, without feeding ground truth labels into the models. We generate 10,000 images for each checkpoint during training. To calculate the proportion for each attribute in the generation distribution, we require attribute labels for the generated images. We apply a trained classifier, developed using the training and validation sets with ground truth labels, to the generated images to obtain classifier-predicted attribute labels.

The trained classifier inevitably introduces errors, meaning the predicted labels may not match the ground truth labels for all images. To ensure consistent bias estimation across different sets, we use the trained classifier to predict attribute labels for training and validation sets. In addition, we use  $P^{val}(C = 1)$  to estimate the probability of attribute C in the data distribution, as the classifier may overfit to the training set. By adopting these techniques, we aim to minimize the potential bias introduced by the classifier in our bias evaluation framework for generative models.

Given a binary attribute  $C \in C$ , we can therefore define **bias shift** between generation and training data as

$$B_{\text{shift}}(C=1) = |B^{\text{gen}}(C=1) - B^{\text{data}}(C=1)| = |P_{\text{cls}}^{\text{gen}}(C=1) - P_{\text{cls}}^{\text{val}}(C=1)|.$$
(3)

The subscript cls stands for using classifier-predicted labels. In bias shift, the expected probability for positive attribute C = 1 in an ideal reference distribution  $P^{\text{ideal}}(C = 1)$  is canceled out. Bias shift remains the same regardless which ideal bias reference we select. If bias shift is close to 0, then the generation distribution and the training distribution exhibit the same level of bias for the given attribute.

Bias shift evaluates changes in bias between data and generation distribution for each attribute
considered in the study. To provide an overall understanding of the magnitude of bias shift across all attributes, we propose to use the average of bias shift across attributes. Average bias shift (ABS)
evaluates the overall bias shift magnitude across all attributes considered between the training and the generated data distributions. This value represents the absolute difference between probabilities and is expressed as a percentage. We define this metric as

$$ABS = \mathbb{E}_{C \in \mathcal{C}} B_{\text{shift}}(C = 1).$$
(4)



Figure 3: Evaluation metrics for image generation throughout training. In 3a and 3b, FID, KID, and FLD values converge to small values showing the good quality of generated images and good coverage of modes of the training distribution. In 3c, the positive or slightly negative generalization gaps indicate that the trained models do not have severe memorization issues.

4 **EXPERIMENTS** 

EXPERIMENTAL SETUP 4.1

**Datasets** We apply our proposed bias evaluation framework to two real datasets – CelebA (Liu et al., 2015) and DeepFashion (Liu et al., 2016). CelebA (Liu et al., 2015) is a large-scale dataset with 200,000 celebrity facial images, each labeled with 40 binary attributes. It covers a wide range of facial features, from details (e.g., earrings, pointy nose) to outlines (e.g., hair color, gender, age). DeepFashion (Liu et al., 2016) is a clothes dataset with over 800,000 diverse fashion images. We use a subset with 26 fine-grained attribute annotations to train the classifier. We then study the bias shift over these fine-grained attributes. For both datasets, we follow the training/validation/test set split from the official release. More details about these datasets are in Appendix A.

**Backbone models in the framework** We follow the setup from Dhariwal & Nichol (2021) to train 245 unconditional ablated diffusion models  $(ADMs)^2$ . We train models of varying sizes by adjusting the 246 number of channels in the U-Net (Ronneberger et al., 2015) bottleneck layer (32 for tiny, 64 for 247 small, and 256 for large), with proportional changes in each layer. In the following sections, we 248 report the results of the large diffusion model if the model is not otherwise specified. We gener-249 ate 10,000 images per checkpoint using 100 inference steps across training. We use a ResNext50 250  $(32x4d)^3$  based image classifier (Xie et al., 2017). We add a linear layer on top as the classification 251 head and fine-tune the last 6 layers of the ResNext50 model. For comparison with a GAN model, 252 we train a BigGAN (Brock et al., 2019) model<sup>4</sup> using the recommended settings. Implementation details are in Appendix B. 253

254

259

261

216

217

218

219

220

222

224 225

226

227

228

229 230 231

232 233

234 235

236

237

238

239

240

241

242

243 244

255 Evaluation metrics for Image Generation We use some common metrics, e.g., FID (Fréchet 256 Inception Distance) (Heusel et al., 2017) and KID (Kernel Inception Distance) (Binkowski et al., 257 2018), to evaluate the generated images. We use FLD (Feature Likelihood Divergence) and general-258 ization gap (Jiralerspong et al., 2023) as two additional metrics to gauge the memorization level of the generative models. FLD provides a comprehensive evaluation considering not only quality and diversity, but also novelty (i.e., difference from the training samples) of generated samples. Positive 260 generalization gap shows no overfitting to the training set. We adopt the implementation<sup>5</sup> of Jiralerspong et al. (2023) and follow their suggestion of using DINOv2 (Oquab et al., 2024) as the feature 262 extractor to calculate FID, KID, and FLD. We also use a conventional FID implementation<sup>6</sup> to give 263 a comparable value of how well the trained models are.

- 264 265 266
- <sup>2</sup>https://github.com/openai/guided-diffusion

<sup>3</sup>The pre-trained model is from torchvision. 267

- <sup>4</sup>https://github.com/ajbrock/BigGAN-PyTorch 268
- <sup>5</sup>https://github.com/marcojira/FLD

<sup>&</sup>lt;sup>6</sup>https://github.com/mseitzer/pytorch-fid



Figure 4: Average bias shift (ABS) for CelebA and DeepFashion. For both datasets, shown in Figs. 4a and 4b, ABS over *subjective* attributes show a much larger bias shift than *non-subjective* ones. Fig. 4c presents that the error coming from sampling of 10K images is small enough, showing that the sampling randomness is not the only cause of bias shifts in generations.

Table 1: Attribute categorization of subjective and non-subjective for each dataset.

Dataset	subjective attributes	non-subjective attributes		
CelebA	Rosy_Cheeks, Big_Nose, No_Beard, Narrow_Eyes, Arched_Eyebrows, High_Cheekbones, Bushy_Eyebrows, Black_Hair, Receding_Hairline, Brown_Hair, Straight_Hair, Bags_Under_Eyes, Pointy_Nose, Big_Lips, Mouth_Slightly_Open, Heavy_Makeup, Attractive, Smiling, Wearing_Lipstick, Wavy_Hair, Young, Oval_Face,	5-o-Clock_Shadow, Bangs, Eyeglasses, Bald, Double.Chin, Wearing.Hat, Male, Blond.Hair, Gray.Hair,Mustache, Chubby, Pale.Skin, Sideburns,Goatee,		
DeepFashion	Floral, Graphic, Embroidered, Solid, Long_sleeve, Short_sleeve, Sleeveless, Knit, Chiffon, Cotton, Maxi_length, Mini_length, No_dress, Crew_neckline,V_neckline, No_neckline, Loose, Tight, Conventional	Striped, Pleated, Leather, Faux, Square.neckline, Lattice, Denim,		

#### 4.2 BACKBONE MODELS PERFORMANCE

**Diffusion Models** Figure 3 shows the image generation evaluation metrics for CelebA and Deep-Fashion datasets. In Figs. 3a and 3b, FID and KID converge to small values showing the good quality of generated images and good coverage of modes of the data distribution. FLD agrees with conventional metrics, showing no severe memorization issues in the generation. In Fig. 3c, the positive or slight negative values of generalization gap indicate that no overfitting is detected in the trained models. More discussions are in Appendix B.1.

303

280

281

282

283 284

287

289 290 291

293

295

296

Classifier For CelebA and DeepFashion datasets, the classification accuracy on the validation set for most attributes is over 80%. Overall, the average accuracy across attributes is 91.7% for CelebA and 90.5% for DeepFashion. Table 4 and Table 5 in Appendix B.2 show in detail the classifier performance for each attribute.

308 309

310

#### 4.3 AVERAGE BIAS SHIFT EVALUATION

Fig. 4 presents the average bias shift (ABS) throughout training. The overall ABS is still perceivable when image generation metrics are small, indicating non-negligible bias shifts from the training to generation distributions. Looking closer into bias shift for each attribute (Figs. 8 and 9 in Section 4.6), we can categorize all attributes into two categories: *subjective* and *non-subjective*.

Taking CelebA as an example, intuitively, *non-subjective* attributes are those where classification judgements are consistent across populations, e.g., eyeglasses, wearing\_hat, bangs, goatee, etc., while *subjective* ones are those where classification judgements differ significantly from one person to another, e.g., heavy\_makeup, arched\_eyebrows, attractive, oval\_face, etc. We present the categorization of attributes in Table 1. In the following section 4.4, we will talk about the criteria for the attributes categorization.

Average bias shift (ABS) for *non-subjective* attributes (purple dashed lines in Fig. 4) converges to
 small values for both datasets, reaching 0.71% for CelebA and 0.98% for DeepFashion. However,
 *subjective* attributes exhibit significantly larger ABS, achieving minima of 3.25% for CelebA and
 4.73% for DeepFashion.



Figure 5: CelebA classifier's pre-sigmoid logits distributions of selected *subjective* and *non-subjective* attributes. The decision boundary for *subjective* attributes (Fig. 5a, 5b, and 5c) always falls in a high-density region, while for *non-subjective* attributes (Fig. 5d, 5e, and 5f) it falls in a low-density region.



Figure 6: **DeepFashion classifier's pre-sigmoid logits distributions of selected** *subjective* and *non-subjective* attribute. The decision boundary for *subjective* attributes (Fig. 6a, 6b) always falls in a high-density region, while for *non-subjective* attribute (Fig. 6c) it falls in a low-density region.

Bias shifts do not consistently follow the image generation metrics, as illustrated by the comparison between Figs. 3 and 4. This misalignment highlights that models with superior image generation metrics are not necessarily less biased. Bias should be treated as an independent issue, distinct from quality and diversity. While diversity metrics typically assess the coverage of modes in the generated distribution, bias evaluation should focus on the relative proportions of these modes. For CelebA dataset, the bias evaluation metrics plateau between steps 110K and 210K, while the image generation metrics continue to improve. Similarly, for DeepFashion dataset, the image generation metrics continue improving during the whole training, while BSRs for both subjective and non-subjective attributes are stable with slight increases after about 200K steps. 

To demonstrate that sampling 10,000 image generation is sufficient for a reliable statistical estimation, we present ABS between sampled subsets of the validation set and the full validation set on CelebA dataset in Fig. 4c. Additionally, we plot ABS between the generation set at the final checkpoint and the full validation set. The generation set has a much larger ABS compared to the sampled validation set with 10,000 images, emphasizing that the bias shifts observed in Figs. 4a and 4b exceed the variance introduced by the sampling process. This also suggests that using 10,000 images is sufficient to estimate bias shifts with minor errors. 378 FLD of different models for CelebA ABS (overall) ABS (subjective) ABS (non-subjective) large diffusion small diffusion large diffusion small diffusion FLD DINO small large diffusion 379 FLD DINO large 20.00 20.00 g 20.00 small diffusion FLD DINO BigGAN 380 01 40 BigGAN BigGAN BigGAN 10.00 g FLD 10.00 10.00 381 20 bias bias bias 382 0.00 0.0 0.0 n 400 600 200 400 steps (K) 600 600 200 40 steps (K) 400 383 200 400 steps (K) steps (K) 384 (c) ABS subjective (d) ABS non-subjective (a) FLD (b) ABS overall 385

Figure 7: **FLD and ABS of different generative models on CelebA.** The small diffusion model has slightly worse image generation quality but much larger ABS for both *subjective* and *non-subjective* attributes compared to the large diffusion model. BigGAN has a similar FLD as the large diffusion model but has larger bias shifts.

386

387

388

#### 4.4 BIAS SHIFTS' SENSITIVITY RELATES TO DECISION BOUNDARY

In this section, we analyze the classifier to explain why some attributes experience greater bias shifts
 than others, leading to the attribute taxonomy presented in Table 1.

Figs. 5 and 6 show the trained classifier's pre-sigmoid logits distribution for some attributes of CelebA and DeepFashion respectively. The distributions for all attributes are in Appendix B.2. These plots provide visualizations of how the data points are distributed in a projected unidimensional space. To estimate the empirical distributions, we use all the training images, 10,000 images sampled from the validation set, and all the 10,000 images in the generation set.

The main difference between *small bias shift* and *large bias shift* attributes is the density at the decision boundary. The distribution shifts for different attributes can manifest in various ways, but the decision boundaries for *large bias shift* attributes consistently fall in higher density regions compared to those for *small bias shift* ones. We thus use the density where the decision boundary falls in the validation distribution to categorize the attributes. Those with density more than 0.01 are categorized as *subjective*, and vice versa.

407 Bias shifts of subjective attributes are more sensitive to distribution shifts compared to nonsubjective attributes. The distributions for non-subjective attributes still change between training 408 and generation sets, but their effects on bias shifts are small. Since the decision boundary falls in a 409 low-density region, it is more difficult to transport the density mass from one side of the boundary to 410 the other. We find empirically that the distribution shifts between the training and generation distri-411 butions generally have low earth mover's distance (EM distance) (Rubner et al., 1998). Significant 412 reweighting of well-separated modes would constitute a significant EM distance between training 413 and generated distributions. For example, the distribution of male (Fig. 5e) shifts from training 414 to generation, but the shifts are within each side of the decision boundary. This clear classification 415 margin leads to small ABS for *non-subjective* attributes. 416

417 418 4.5 BIAS SHIFT IN DIFFERENT GENERATIVE MODELS

419 In this section, we compare the bias shifts for different sizes of diffusion models by changing the 420 number of channels in the bottleneck layer of U-Net (32 for tiny, 64 for small, and 256 for large), 421 and BigGAN model. Fig. 7 shows image generation metrics and bias evaluation metrics for different 422 generative models. The tiny diffusion model cannot generate realistic images (check Appendix D for 423 sampled images), making it unsuitable for our bias analysis framework. FLD for the small diffusion model is worse than the large diffusion model, while BigGAN achieves a similar FLD as the large 424 diffusion model. However, ABS shows clear differences among generative models (See Figs. 7b, 7c 425 and 7d). 426

427 Diffusion models have matched or even surpassed GAN models regarding image synthesis perfor 428 mance (Dhariwal & Nichol, 2021). We evaluate whether diffusion models also perform better than
 429 BigGAN regarding bias shifts. We observe that BigGAN exhibits a considerably larger ABS com 430 pared to the large diffusion model, despite having only slightly worse image generation performance
 431 according to FLD. This finding may be because the common understanding that GAN models suffer
 from mode collapse issues (Arjovsky et al., 2017).



Figure 8: **Probabilities of selected** *subjective* and *non-subjective* attributes for 3 different random seeds during training. The probabilities of *subjective* attributes (Fig. 8a, 8b, and 8c) present a gap between generation and validation data while *non-subjective* ones (Fig. 8d, 8e, and 8f) do not.



Figure 9: **Probabilities of selected** *subjective* and *non-subjective* attribute during training. The probabilities of *subjective* attributes (Fig. 9a, 9b) present a gap between generation and validation data while *non-subjective* one (Fig. 9c) does not.

We train different sizes of diffusion models to study the influence of the model size on bias shifts.
The *small* diffusion model exhibits larger bias shifts compared to the *large* diffusion model. For the *small* diffusion model, we notice that the average bias shift increases after 300K steps although the
FLD value does not change significantly. While BigGAN has better FLD, the bias shifts are similar
to those of the *small* diffusion model at the end of the training. We also observe more fluctuations
in bias evaluation metrics for the small diffusion model.

We present image samples generated from different models in Appendix D. The images generated
by BigGAN and the small diffusion model are "more washed out" than those produced by the large
diffusion model, showing fewer variations and less details.

478 4.6 Additional Results

Per-attribute Bias Shift Figs. 8 and 9 show probabilities for selected attributes in the generated data during training. Plots for other attributes are in Appendix C. Probabilities of *subjective* attributes generally exhibit values distinct from the classifier-predicted validation probabilities, resulting in bias shifts in Fig. 4.

Subjective attributes exhibit more fluctuations throughout training compared to non-subjective ones.
 While the probabilities for many attributes converge before 300K steps, young (Fig. 8a) still has fluctuations. A similar pattern is also witnessed in DeepFashion, where solid (Fig. 9b), as a



(a) conditioned on Male (b) conditioned on Female (c) conditioned on Young (d) conditioned on Old

Figure 10: **ABS for conditional settings on CelebA.** Bias shifts conditioned on subjective attributes may exhibit different patterns as shown in Fig. 10d.

499

493 494

*subjective* attribute, also exhibits perceivable fluctuations. This suggests that extra caution is needed when handling certain *subjective* attributes using generative models.

We conduct several runs of training using different random seeds on CelebA dataset. There is
randomness across different random seeds as the curves for each random seed vary. However, the
probabilities of each attribute from distinct random seeds generally converge to the same value.
Therefore, we report results for only one seed in other experiments.

504

Bias Shift Evaluation Conditioned on Anchor Attributes Fig. 10 illustrates the conditional setting of bias shift evaluation. We focus on two demographic attributes, gender and age. According to our categorization proxy shown in Table 1, gender is *non-subjective*, while age is *subjective* in CelebA. This categorization may seem counterintuitive at first glance.

We acknowledge that it is not appropriate to naively binarize gender and age. However, due to the constraints of the era when the dataset was created, our analysis is restricted to binary gender and age attributes. By conducting an empirical analysis based on these binary attributes, we aim to highlight the importance of recognizing the fluidity of gender and the variability of age. It is important to note that the *subjective* and *non-subjective* categorization applies specifically to the image-label joint distribution presented in the CelebA dataset and is not universally applicable.

The bias change trends for probabilities conditioned on *non-subjective* attributes exhibit similarities to those of unconditioned probabilities (See Figs. 10a and 10b). However, we observe that the average bias shift for *non-subjective* attributes become larger when conditioning on Old, which is categorized as a *subjective* attribute in CelebA in our study. A possible explanation for this discrepancy is that the classifier-predicted labels of *subjective* attributes are not always accurate. Therefore, when conditioning on *subjective* attributes, classification errors propagate into the bias analysis pipeline, resulting in a distinct pattern of bias shifts.

522

#### 5 CONCLUSION

523 524

This study focuses on bias shifts with regard to inductive biases of unconditional image gener-525 ative models. We propose a standardized bias analysis framework applicable to any supervised 526 learning dataset. Our experimental results show that different attributes have varying bias shifts 527 in response to distribution changes. Attributes for which the classifier's decision boundary falls in 528 a low-density area tend to have small bias shifts. We thus categorize all attributes into *subjective* 529 and non-subjective sets. Our analysis results in the following observations: 1) Biases shift between 530 training and generation distributions for unconditional image generative models. 2) Selecting the 531 checkpoint with the best image generation metrics does not guarantee the smallest bias shifts. 3) 532 BigGAN models and small diffusion models have larger bias shifts compared to large diffusion 533 models, despite having similar image generation metric values.

We hope that our analysis for unconditional generative models can serve as a base framework allowing researchers to add other sources of bias such as conditioning (with ground-truth labels, text, etc.), guidance, pretrained modules, etc. in a gradual manner, and study their effects on bias shift in a systematic way.

539

## 540 REFERENCES

547

565

566 567

568

569

570

583

584

585

- Ibrahim Alabdulmohsin, Xiao Wang, Andreas Peter Steiner, Priya Goyal, Alexander D'Amour, and
   Xiaohua Zhai. CLIP the bias: How useful is balancing data in multimodal learning? In *The Twelfth International Conference on Learning Representations*. OpenReview.net, 2024.
- Martín Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein GAN. *CoRR*, abs/1701.07875, 2017.
- 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. ACM, 2023.
- Mikolaj Binkowski, Danica J. Sutherland, Michael Arbel, and Arthur Gretton. Demystifying MMD
   gans. In *6th International Conference on Learning Representations*. OpenReview.net, 2018.
- 554 Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ B. Altman, Simran Arora, Sydney von Arx, 555 Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson, 556 Shyamal Buch, Dallas Card, Rodrigo Castellon, Niladri S. Chatterji, Annie S. Chen, Kathleen Creel, Jared Quincy Davis, Dorottya Demszky, Chris Donahue, Moussa Doumbouya, Esin Dur-558 mus, Stefano Ermon, John Etchemendy, Kawin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale, Lauren Gillespie, Karan Goel, Noah D. Goodman, Shelby Grossman, Neel Guha, Tatsunori 559 Hashimoto, Peter Henderson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang, Thomas Icard, Saahil Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keel-561 ing, Fereshte Khani, Omar Khattab, Pang Wei Koh, Mark S. Krass, Ranjay Krishna, Rohith Ku-562 ditipudi, and et al. On the opportunities and risks of foundation models. CoRR, abs/2108.07258, 563 2021.
  - Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis. In *7th International Conference on Learning Representations*, 2019.
  - Jaemin Cho, Abhay Zala, and Mohit Bansal. DALL-EVAL: probing the reasoning skills and social biases of text-to-image generation models. In *IEEE/CVF International Conference on Computer Vision*, pp. 3020–3031. IEEE, 2023.
- Prafulla Dhariwal and Alexander Quinn Nichol. Diffusion models beat gans on image synthesis. In
   Advances in Neural Information Processing Systems 34: Annual Conference on Neural Informa tion Processing Systems 2021, NeurIPS 2021, pp. 8780–8794, 2021.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
  Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at
  scale. In *9th International Conference on Learning Representations*, 2021.
- Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, Dustin Podell, Tim Dockhorn, Zion English, and Robin Rombach. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first International Conference on Machine Learning*. OpenReview.net, 2024.
  - Felix Friedrich, Manuel Brack, Lukas Struppek, Dominik Hintersdorf, Patrick Schramowski, Sasha Luccioni, and Kristian Kersting. Auditing and instructing text-to-image generation models on fairness. AI and Ethics, pp. 1–21, 2024.
- Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,
   Aaron C. Courville, and Yoshua Bengio. Generative adversarial networks. *CoRR*, abs/1406.2661,
   2014.

589	Matt	Growcoot.	Which	ai	image	generator	is	the	most	bi-
590	ased?,	, 2023.	UR	L	https:	://petapi:	xel.	com/20	)23/11,	/03/
591	whic	h-ai-image-ger	nerator-i	s-th	e-most-	-biased/.				
592										

593 Melissa Hall, Laurens van der Maaten, Laura Gustafson, and Aaron Adcock. A systematic study of bias amplification. *CoRR*, abs/2201.11706, 2022.

594 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. 595 Gans trained by a two time-scale update rule converge to a local nash equilibrium. In Advances in 596 Neural Information Processing Systems 30: Annual Conference on Neural Information Process-597 ing Systems, pp. 6626-6637, 2017. 598 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. CoRR, abs/2207.12598, 2022. 600 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In Advances 601 in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems, 2020. 602 603 Jonathan Ho, Tim Salimans, Alexey A. Gritsenko, William Chan, Mohammad Norouzi, and David J. 604 Fleet. Video diffusion models. In Advances in Neural Information Processing Systems 35: Annual 605 Conference on Neural Information Processing Systems, 2022. 606 Marco Jiralerspong, Avishek Joey Bose, Ian Gemp, Chongli Qin, Yoram Bachrach, and Gauthier 607 Gidel. Feature likelihood score: Evaluating the generalization of generative models using sam-608 ples. In Advances in Neural Information Processing Systems 36: Annual Conference on Neural 609 Information Processing Systems, 2023. 610 611 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 612 on Applications of Computer Vision, pp. 1548–1558, 2021. 613 614 Saehoon Kim, Sanghun Cho, Chiheon Kim, Doyup Lee, and Woonhyuk Baek. mindall-e on con-615 ceptual captions. https://github.com/kakaobrain/minDALL-E, 2021. 616 Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In 2nd International 617 Conference on Learning Representations, 2014. 618 619 Felix Kreuk, Gabriel Synnaeve, Adam Polyak, Uriel Singer, Alexandre Défossez, Jade Copet, Devi 620 Parikh, Yaniv Taigman, and Yossi Adi. Audiogen: Textually guided audio generation. In The 621 Eleventh International Conference on Learning Representations. OpenReview.net, 2023. 622 Donghoon Lee, Jiseob Kim, Jisu Choi, Jongmin Kim, Minwoo Byeon, Woonhyuk Baek, and 623 Saehoon Kim. Karlo-v1.0.alpha on coyo-100m and cc15m. https://github.com/ 624 kakaobrain/karlo, 2022. 625 Junnan Li, Dongxu Li, Silvio Savarese, and Steven C. H. Hoi. BLIP-2: bootstrapping language-626 image pre-training with frozen image encoders and large language models. In International Con-627 ference on Machine Learning, volume 202, pp. 19730–19742. PMLR, 2023. 628 629 Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. 630 In Proceedings of International Conference on Computer Vision (ICCV), December 2015. 631 Ziwei Liu, Ping Luo, Shi Qiu, Xiaogang Wang, and Xiaoou Tang. Deepfashion: Powering robust 632 clothes recognition and retrieval with rich annotations. In Proceedings of IEEE Conference on 633 Computer Vision and Pattern Recognition (CVPR), June 2016. 634 Sasha Luccioni, Christopher Akiki, Margaret Mitchell, and Yacine Jernite. Stable bias: Evaluat-635 ing societal representations in diffusion models. In Advances in Neural Information Processing 636 Systems 36: Annual Conference on Neural Information Processing Systems, 2023. 637 638 Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic mod-639 els. In Proceedings of the 38th International Conference on Machine Learning,, volume 139 of 640 Proceedings of Machine Learning Research, pp. 8162–8171. PMLR, 2021. 641 OpenAI. GPT-4 technical report. CoRR, abs/2303.08774, 2023. 642 643 Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, 644 Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mido Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael 645 Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Jégou, Julien Mairal, Patrick Labatut, Ar-646 mand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision. 647

Trans. Mach. Learn. Res., 2024.

048	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agar-
649	wal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya
650	Sutskever. Learning transferable visual models from natural language supervision. In Proceed-
651	ings of the 38th International Conference on Machine Learning, volume 139, pp. 8748–8763.
652	PMLR, 2021.

- Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139, pp. 8821–8831. PMLR, 2021.
- Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text conditional image generation with CLIP latents. *CoRR*, abs/2204.06125, 2022.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10674–10685, 2022.
- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedi cal image segmentation. In Nassir Navab, Joachim Hornegger, William M. Wells III, and Alejan dro F. Frangi (eds.), *Medical Image Computing and Computer-Assisted Intervention MICCAI*,
   volume 9351 of *Lecture Notes in Computer Science*, pp. 234–241. Springer, 2015.
- Yossi Rubner, Carlo Tomasi, and Leonidas J. Guibas. A metric for distributions with applications to
   image databases. In *Proceedings of the Sixth International Conference on Computer Vision*, pp.
   59–66. IEEE Computer Society, 1998.
- Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. LAION-5B: an open large-scale dataset for training next generation image-text models. In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems*, 2022.
- 677 Preethi Seshadri, Sameer Singh, and Yanai Elazar. The bias amplification paradox in text-to-image
   678 generation. *CoRR*, abs/2308.00755, 2023.
- <sup>679</sup> Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, Devi Parikh, Sonal Gupta, and Yaniv Taigman. Make-a-video: Text-to-video generation without text-video data. In *The Eleventh International Conference on Learning Representations*. OpenReview.net, 2023.
- Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In 9th International Conference on Learning Representations. OpenReview.net, 2021.
- 686 Tiku, Kevin Schaul, and Szu Yu Chen. These fake images re-Nitasha 687 veal how ai amplifies our worst stereotypes, 2023. URL https: 688 //www.washingtonpost.com/technology/interactive/2023/ ai-generated-images-bias-racism-sexism-stereotypes/. 689
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
  Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, ArJoulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation
  language models. *CoRR*, abs/2302.13971, 2023.
- Saining Xie, Ross B. Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual
   transformations for deep neural networks. In 2017 IEEE Conference on Computer Vision and
   Pattern Recognition,, pp. 5987–5995, 2017.
- Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. Men also like
   shopping: Reducing gender bias amplification using corpus-level constraints. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pp. 2979–2989,
   2017.

#### 702 DATASETS А

703 704

CelebA (Liu et al., 2015) is a large-scale face attributes dataset with 200,000 celebrity images, each 705 with 40 attribute annotations. The dataset includes 10,000 celebrities with 20 images for each. 706 These attribute annotations cover a wide variety of facial characteristics, ranging from details (e.g., earrings, pointy noise, etc.) to outlines (e.g., hair color, gender, age, etc.). We list all 40 attributes 708 in Table 2. Before feeding the training images to the model, we centre crop the images and re-709 size them to 128x128 pixels. Because of the crop, some attributes, e.g., Wearing\_Necklace, 710 Wearing\_Necktie, are not visually grounded in the post-process images. Blurry is also an attribute that we do not include since we want the image generation quality to be good. We excluded 711 these attributes in Table 1. We follow the Training/Validation/Test set split in the official release. 712 Training set includes the images of the first eight thousand identities (with 160 thousand images). 713 Validation set contains the images of another one thousand identities (with twenty thousand images). 714 The remaining one thousand identities (with twenty thousand images) go for Test set. In our bias 715 analysis framework, we only use the Training set and the Validadtion set. 716

DeepFashion (Liu et al., 2016) is a clothes dataset with over 800,000 diverse fashion images, in-717 cluding tops and bottoms. No footwears is in this dataset. Each image is associated with 1000 718 coarse attribute annotations about texture, fabric, shape, part, and style of the clothes. These at-719 tribute annotations are scrapped directly from meta-data of the images. They are thus very noisy 720 and not reliable. Most of the attributes have less than 1% positive samples, making the classification 721 problem very imbalanced. This dataset also provides a fine-grained annotation subset, where each 722 image is associated with 26 find-grained attribute annotations. These attributes are presented in Ta-723 ble 2. We train a classifier on this subset and apply this trained classifier to the whole dataset and 724 get classifier-predicted labels for each image. We follow the Training/Validation/Test set split in the 725 official release. Unlike CelebA dataset, the split of DeepFashion dataset is random. 726

Table 2: Labeled attributes in CelebA and DeepFashion datasets. CelebA has 40 attributes and 727 DeepFashion has 26 attributes. 728

730	Dataset	Attributes
731	CelebA	5_o_Clock_Shadow, Arched_Eyebrows, Attractive, Bags_Under_Eyes,
732		Bald, Bangs, Big_Lips, Big_Nose, Black_Hair, Blond_Hair, Blurry,
733		Brown_Hair,Bushy_Eyebrows, Chubby, Double_Chin, Eyeglasses, Goatee,
734		Gray_Hair, Heavy_Makeup High_Cheekbones, Male, Mouth_Slightly_Open,
735		Mustache, Narrow_Eyes, No_Beard, Oval_Face, Pale_Skin, Pointy_Nose,
736		Wavy Hair Wearing Earrings Wearing Hat Wearing Lipstick
737		Wearing_Necklace, Wearing_Necktie, Young
738		floral, graphic, striped, embroidered, pleated, solid, lattice,
739	DeepFashion	long_sleeve, short_sleeve, sleeveless
740		<pre>maxi_length,mini_length,no_dress,</pre>
741		crew_neckline, v_neckline, square_neckline, no_neckline,
742		denim, chilion, cotton, leather, faux, knit,
743		Light, 100se, conventional

744 745

746 747

748

729

#### **TRAINING DETAILS** В

#### **B.1 DIFFUSION MODELS**

749 We follow the training setting of Dhariwal & Nichol (2021) to train the ablated diffusion models 750 (ADMs). Hyperparameters and architecture selections are in Table 3. We train the diffusion using 751 NVIDIA A100 40GB. The batch size per GPU is set to 16, and we use 8 GPUs to train. During 752 training, we save checkpoint for EMA models every 10K steps. We use half precision (FP16) for 753 training and inference. 754

For each saved checkpoint, we employ 100 steps in inference to generate 10K images from the 755 Gaussian noise. We compare the two inference methods used in ADM (Dhariwal & Nichol, 2021),



#### Table 3: Hyperparameters and architecture selection for diffusion models

772 Figure 11: ABS and image generation metrics using different inference methods and inference steps 773 on CelebA dataset. Images generated by DDIM have less bias shifts compared to those by Improved 774 Diffusion Sampler. FID and KID also show the superiority of DDIM sampler.

775 776

756

758

761

763

765

766

767

768

769 770

771

777 one proposed by improved diffusion model (Nichol & Dhariwal, 2021), and DDIM (Song et al., 2021). The results on CelebA dataset are in Fig. 11. Images generated by the improved diffusion 778 sampler exhibit more bias shifts than those from DDIM. Although FLD shows a slight improvement 779 on improved diffusion sampler, DDIM works better in terms of FID and KID using the same steps of inference. Since we want to test less biased generations, we use DDIM with 100 steps during 781 inference in our experiments. 782

783 In Fig. 3c, generalization gaps for CelebA and DeepFashion datasets are different. This is because 784 the split of the dataset is in different ways. In CelebA dataset, the training and validation sets contain 785 the faces of distinct sets of celebrities. In DeepFashion dataset, the training and validation samples are split randomly. The distribution difference between training and validation sets of CelebA is 786 larger than that of DeepFashion. 787

788 789

#### **B.2 RESNET CLASSIFIERS**

790 We employ a pre-trained ResNeXt model as the base model. We add a linear layer to top as the 791 classification layer. We then fine-tune the last 6 layers of the pre-trained model as well as the classi-792 fication layer using CelebA and DeepFashion dataset. We use AdamW optimizer and learning rate 793 at 0.001. We follow a standard training procedure for the classifier training. We train the classifier 794 on the train set (with ground truth labels) and choose the best classifier according to the average performance across all the considered attributes on the valid set (with ground truth labels). We use 796 data augmentations to make the classifier more robust. The data augmentations include random horizontal flip, scaling and resizing, etc. This can help the classifier become more reliable when applied 797 to the generation set. Previous work indicates that classifiers can amplify the discriminative biases 798 in the training set (Zhao et al., 2017; Hall et al., 2022). We use the positive and negative sample 799 ratio to reweigh the cross entropy loss terms. This acts as an upsampling of the minority samples 800 and alleviates the label imbalance issue. We don't see the discriminative biases being amplified for 801 most attributes according to Figs. 15 and 14 comparing the training ground truth probability and 802 the validation classifier-predicted probability. The classifiers' performances for each attribute are 803 listed in Tables 4 and 5. For both dataset, the accuracy for most attributes is over 80%. Figs. 12 and 804 13 show the pre-sigmoid logits distributions for each attribute in CelebA and DeepFashion datasets 805 respectively. 806

- 809





Attr	Accuracy	Precision	Recall	F1	AUPR
Eyeglasses	99.58	97.10	96.82	96.96	94.23
Wearing_Hat	98.98	86.31	93.19	89.62	80.75
Bald	98.92	73.33	74.94	74.13	55.47
Male	98.64	98.47	98.32	98.40	97.53
Gray_Hair	97.74	78.09	74.46	76.23	59.39
Sideburns	97.12	82.88	73.30	77.80	62.59
Goatee	96.61	76.83	77.25	77.04	61.03
Double_Chin	96.51	69.99	50.46	58.64	37.75
Pale_Skin	96.41	60.32	48.83	53.97	31.66
Mustache	95.90	60.78	53.14	56.70	34.66
Blurry	95.86	55.59	62.45	58.82	36.49
Wearing_Necktie	95.66	71.41	67.15	69.21	50.34
No_Beard	95.49	97.87	96.62	97.24	97.34
Chubby	95.35	65.18	51.73	57.68	36.67
Bangs	95.26	82.86	85.39	84.10	72.89
Blond_Hair	95.07	82.75	85.86	84.28	73.23
Rosy_Cheeks	94.64	64.32	48.45	55.27	34.69
Receding_Hairline	94.15	59.84	56.82	58.29	37.11
5-o-Clock_Shadow	93.34	77.82	60.90	68.33	52.00
Mouth_Slightly_Open	92.83	92.97	92.07	92.52	89.42
Wearing_Lipstick	92.08	87.96	95.29	91.48	85.92
Smiling	91.50	90.73	91.80	91.26	87.25
Bushy_Eyebrows	91.42	72.05	65.03	68.36	51.84
Heavy_Makeup	91.19	86.20	92.17	89.08	82.50
Narrow_Eyes	90.97	42.41	56.57	48.48	27.25
Wearing_Earings	90.62	82.10	65.00	72.56	60.04
Black_Hair	89.60	71.52	83.33	76.97	63.07
Wearing_Necklace	86.98	43.51	26.71	33.10	20.46
Young	86.42	90.45	91.47	90.96	89.11
High_Cheekbones	86.09	83.47	86.10	84.76	78.11
Brown_Hair	83.41	66.70	62.42	64.49	50.70
Bags_Under_Eyes	83.33	64.93	42.73	51.54	39.63
Arched_Eyebrows	83.08	72.64	55.40	62.86	51.77
Wavy_Hair	83.06	66.23	79.04	72.07	58.15
Straight_Hair	81.97	56.09	56.70	56.39	40.71
Big_Nose	81.63	69.39	46.81	55.91	45.71
Big_Lips	81.28	37.00	31.57	34.07	22.17
Attractive	80.07	78.42	85.09	81.62	74.48
Pointy_Nose	72.97	52.86	47.24	49.89	40.00
Oval_Face	68.34	44.95	57.86	50.59	37.81

Table 4: Classifier performance on validation set of CelebA.

Table 5: Classifier performance on validation set of DeepFashion.

Attr	Acc	Precision	Recall	F1	AUPR
lattice	99.48	100.00	50.00	66.67	50.52
square_neckline	98.97	0.00	0.00	0.00	1.03
faux	98.45	50.00	33.33	40.00	17.70
leather	97.94	0.00	0.00	0.00	1.03
pleated	97.42	40.00	50.00	44.45	21.03
maxi_length	96.91	96.00	82.76	88.89	82.03
denim	96.91	87.50	58.33	70.00	53.62
striped	96.39	55.56	62.50	58.82	36.27
loose	94.33	60.00	25.00	35.29	19.64
knit	92.27	52.63	62.50	57.14	35.99
mini_length	91.24	75.61	81.58	78.48	65.29
graphic	90.72	69.70	74.19	71.88	55.83
embroidered	90.72	36.36	26.67	30.77	15.37
long_sleeve	90.72	82.54	88.14	85.25	76.36
short_sleeve	90.21	66.67	73.33	69.84	53.01
no_dress	90.21	90.91	94.49	92.66	89.51
solid	88.14	88.89	88.00	88.44	88.41
floral	87.63	61.90	76.47	68.42	51.46
tight	87.63	61.29	61.29	61.29	43.75
chiffon	87.11	57.69	51.72	54.55	37.06
v_neckline	86.60	70.83	47.22	56.67	43.24
sleeveless	86.08	86.79	87.62	87.20	82.75
conventional	80.93	86.54	89.40	87.95	85.62
no_neckline	75.26	71.26	72.94	72.09	63.84
cotton	75.26	81.34	82.58	81.95	79.03
crew_neckline	71.65	59.30	71.83	64.97	52.91

#### 1026 **BIAS SHIFT ANALYSIS PER ATTRIBUTE** С 1027

Figs. 15 and 14 show the bias probability for each attribute in CelebA and DeepFashion datasets 1029 respectively. 1030 1031 floral graphic striped embroidered pleated train cls prob eval cls prob generation pr train cls prob eval cls prob generation pr train cls prob eval cls prob generation pr train cls prob eval cls prob train cls prob eval cls prob 0.2 0.8 0.8 1032 0 0.8 Atilita Probability Probability Ati 0.6 ₹ 0.6 £ 0.6 ₹ 0.6 1033 iedor.4 IEQ 0.4 equ. 1034 0.2 0.3 0.2 0. 0.2 1035 0.0 0.0 0.0 0.0 0.0 20 Steps (K) 20 Steps (K) 20 Steps (K) 30 40 20 Steps (K) 20 Steps (K) 1036 (a) Floral (e) Pleated (b) Graphic (c) Striped (d) Embroidered 1037 solic lattice long\_sleeve short\_sleeve sleeveless 1038 train cls prob eval cls prob generation prob 0.1 0.8 0.8 0.8 0.8 1039 ₹0. ₹ 0.0 ₹0.6 ₹0.e ≩ 0.6 다.4 eq 0.4 Proba 1040 Page 0.4 Q 0.4 train cls prob eval cls prob 0. 0.2 0.: 0. 0. 1041 aene 0.0 0.0 0.0 0.0 0.0 20 Steps (K) 30 40 30 40 20 Steps (K) 20 Steps (K) 20 Steps (K) 10 20 Steps (K) 30 1042 1043 (i) Short Sleeve (j) Sleeveless (f) Solid (g) Lattice (h) Long Sleeve 1044 maxi\_length mini\_length no\_dress rew\_neckline v\_neckline 1 1 train cls prob eval cls prob generation prob train cls prob eval cls prob generation pro train cls prob eval cls prob train cls prob eval cls prob 1045 0.8 0.8 0.8 0. 0.8 Probability Probability ₹0.6 ₹0.e Probability ₹0.e 1046 Probal Proba ٩<u>0</u>.4 train cls prob 1047 0. 0: 0. eval cls prob generation prob 0 3 0 1048 0.0 0.0 0.0 0.0 0.0 20 Steps (K) 20 Steps (K) 40 20 Steps (K) 30 40 20 Steps (K) 20 Steps (K) 1049 (k) Maxi Length (l) Mini Length (m) No Dress (n) Crew Neckline (o) V Neckline 1050 quare\_neckline no\_neckline denim chiffon cotton 1051 1. train cls prob eval cls prob generation prob train cls prob eval cls prob train cls prob eval cls prob generation pro train cls prob eval cls prob generation prol 0. 0. 0. 0. 0. 1052 generation prol ≩ ≣ 0.6 ₹ 0.6 ₹ 0.6 ₹ 0.0 ₹0.e opab o 1053 o o. 0.4 ĝ 0.4 prob. train cls prob eval cls prob generation prot 1054 0.3 0.2 0.2 0. 0.2 0.0 0.0 0. о. 20 Steps (K) 20 Steps (K) 20 Steps (K) 20 Steps (K) 1055 20 Steps (K) 1056 (p) Square Neckline (r) Denim (s) Chiffon (t) Cotton (q) No Neckline 1057 leather faux knit tight loose 1. 1. train cls prob eval cls prob generation prob train cls prob eval cls prob generation pro train cls prob eval cls prob generation prot train cls prob eval cls prob generation prol train cls prob eval cls prob generation pro 1058 0.1 0.8 0.8 0. 0.8 <u>کال 0</u> ₹0.0 ٥.e ₹ 0.6 <u>کا آر</u> 1059 q 0.4 độ 0.4 g 0.4 20.4 20.4 1060 0.3 0.2 0.3 0. 0. 1061 0.0 L 0. 0.0 0. 0.0 20 Steps (K) 20 Steps (K) 20 Steps (K) 40 20 Steps (K) 4( 20 Steps (K) 1062 (w) Knit (u) Leather (v) Faux (x) Tight (y) Loose 1063 conventiona 1064 0.8 1065 <u>م</u> 1066 ĝ 0.4 train cls prob 1067 0. eval cls prob generation prob 0.0 1068 20 Steps (K) 30 1069 (z) Conventional 1070

1071 Figure 14: Probabilities of attributes for DeepFashion dataset during training. Please note that it 1072 might seem like some of the subplots are missing the probability lines; they are actually very close 1073 to the x-axis, especially for Square Neckline and Faux.

1074

1028

- 1075
- 1076

1077



# 1134 D SAMPLES OF GENERATED IMAGES

For different models and different dataset, we sample 80 images from the generation set and present them in Figs. 16, 17, 18, 19 and 20.

Figure 16: Image samples from large diffusion model generations on CelebA dataset.



Figure 17: Image samples from the small diffusion model trained on CelebA dataset.



Figure 18: Image samples from the BigGAN model trained on CelebA dataset.





Figure 20: Image samples from the large diffusion model trained on DeepFashion dataset.