

Improving the Fairness of Deep Generative Models without Retraining

Shuhan Tan¹ Yujun Shen² Bolei Zhou²

¹Sun Yat-Sen University ²The Chinese University of Hong Kong

tanshh@mail2.sysu.edu.cn, {sy116, bzhou}@ie.cuhk.edu.hk

Abstract

Generative Adversarial Networks (GANs) advance face synthesis through learning the underlying distribution of observed data. Despite the high-quality generated faces, some minority groups can be rarely generated from the trained models due to a biased image generation process. To study the issue, we first conduct an empirical study on a pre-trained face synthesis model. We observe that after training the GAN model not only carries the biases in the training data but also amplifies them to some degree in the image generation process. To further improve the fairness of image generation, we propose an interpretable baseline method to balance the output facial attributes without retraining. The proposed method shifts the interpretable semantic distribution in the latent space for a more balanced image generation while preserving the sample diversity. Besides producing more balanced data regarding a particular attribute (e.g., race, gender, etc.), our method is generalizable to handle more than one attribute at a time and synthesize samples of fine-grained subgroups. We further show the positive applicability of the balanced data sampled from GANs to quantify the biases in other face recognition systems, like commercial face attribute classifiers and face super-resolution algorithms.¹

1. Introduction

Artificial Intelligence (AI) is being applied to a wide range of applications in our daily life, such as employee hiring, loan granting, and criminal searching [17]. The decisions made by AI algorithms, which sometimes matter to affect the life of people, are required to be unbiased and trustworthy. Unfortunately, current AI systems are known to have a discriminatory nature due to imbalanced training data [34, 5, 9, 6] and various algorithmic factors [3, 33]. Such biases within the AI models may lead to unfair treatment of people, especially those from minority groups. As a result, fairness, together with other ethical issues,

¹Project page is at <https://genforce.github.io/fairgen>.

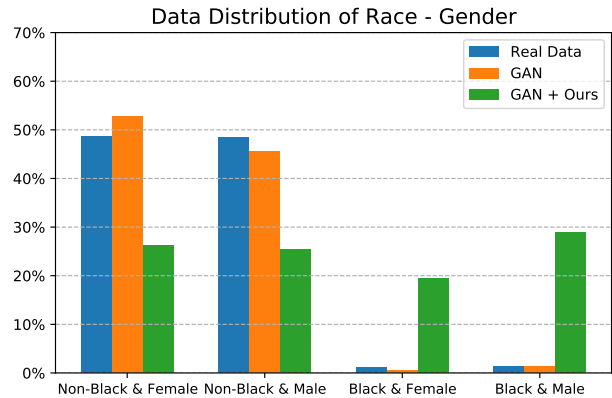


Figure 1. **Top:** The joint distribution of two face attributes, *i.e.*, race and gender, on the training data, the native synthesis from GANs, and the far more balanced synthesis after improving the model fairness with our method. **Bottom:** Visual samples for each subgroup with our proposed sampling scheme. The *fixed* StyleGAN2 model pre-trained on FF-HQ dataset [19] is used.

becomes a crucial AI research topic.

Many studies have been made to analyze and improve the fairness in classification models [14, 29, 28, 10]. Their primary goal is to learn a system that can give equally accurate prediction, with respect to a specific task, on every single group of people. Along with the recent development

of Generative Adversarial Networks (GANs) [12, 4, 18, 19], deep models have shown strong capability in many generative tasks, like style transfer [21], semantic manipulation [8], image super-resolution [35], *etc.* Therefore, it becomes critical to also examine the fairness of these generative models [20]. Concretely, we investigate whether an algorithm can give similar performance on different groups of people from the perspective of whether each group can be generated by the models with equal probability.

In fact, many recent generative models have been found to generate imbalanced images [43, 30, 16]. To improve the fairness of image generation, existing approaches typically fall into two folds, *i.e.*, pre-processing and in-processing. The pre-processing methods try to eliminate the bias from the training data perspective (*e.g.*, collecting new data or selecting a subset of the data collection to make it balanced) [25], while the in-processing methods target at improving the data coverage learned by the model via introducing new training objectives [37, 31, 7, 40]. Both kinds of approaches, however, require training new models, making them hard to apply to the numerous systems that have already been released to the public. On the other hand, as evidenced by some recent works, pre-trained GANs can provide rich information to facilitate various downstream tasks [32, 26, 39, 13, 27, 38]. Accordingly, it would be of great use if we can improve the fairness of existing models without touching the model weights.

In this work, we first conduct an empirical study on the fairness of a state-of-the-art pre-trained face synthesis GAN model. Fig. 1 gives an example on the model bias related to race and gender, where the bias in the training data (blue bars) is carried and amplified by the GAN model (orange bars). Based on this observation, we propose an interpretable baseline method to alleviate the biases in pre-trained GANs, which does not require any retraining or access to the original training data. Instead, it examines the sampling process of GANs and shifts the interpretable semantic distribution in the latent space for each subgroup of interest. As shown in Fig. 1 (green bars) and other experiments, our baseline method can very well improve the fairness of the models.

We summarize our contributions as follows.

- We conduct an empirical study to reveal relationship of the bias in the training data and a well-learned GAN model. We observe that bias in the training data is carried and further amplified by the GAN model.
- We propose an interpretable baseline, which is able to improve the fairness of GANs not only regarding a single attribute but also across multiple attributes.
- We demonstrate the positive applicability of the GAN model calibrated by our method for quantifying the biases in other AI systems, *i.e.*, attribute classifiers and super-resolution algorithms.

2. Related Work

AI Fairness. AI fairness has attracted wide attention in recent years [3, 5, 20]. Most existing works focus on studying the fairness of discriminative models (*e.g.*, face attribute classifiers), where fairness is achieved when the model makes predictions in a non-discriminatory way. For example, an age predictor is expected to achieve similar performance on different genders. Three types of approaches have been proposed: the pre-processing methods try to collect balanced training data [42, 24, 23], the in-processing methods introduce constraints or regularizers into the training process [36, 41, 2], and the post-processing methods attempt to modify the posteriors of pre-trained models [11, 14]. Compared to discriminative models, which are primarily designed to make inference on existing data, generative models are able to create new data, enabling a lot more application scenarios [8, 35, 13, 27]. But little efforts have been made to explore the fairness of generative models and existing methods typically focus on the pre-processing stage and the in-processing stage. This work fills this hole by improving the fairness of well-learned GAN models without retraining.

Fairness of Generative Models. Studying the fairness of generative models is as important as studying that of discriminative models. It has been shown that the data produced by a fair generative model can benefit various downstream classifiers [37, 31]. However, Zhao *et al.* [43] find that GAN will usually carry and amplify the bias existing in the training data. Jain *et al.* [16] further analyze this problem from the perspective of mode collapse. Some attempts have been made to reduce the model bias and learn a more balanced image generation [7, 40, 25]. Choi *et al.* [7] propose to learn a weighting function to reweigh the importance of each instance during the training of GANs. Yu *et al.* [40] mitigate bias by increasing the data coverage learned by GANs. McDuff *et al.* [25] carefully balance the training set such that all groups of interests possess similar number of samples. Our **improvements** are: 1) We propose a fair sampling pipeline by shifting the latent distribution of a well-learned GAN model instead of training new models or tuning the model weights, which can be often costly. In this way, our approach can be applied to numerous existing models with minor effort. 2) We do not rely on the full access of the original training dataset, which is often not available for end-users. 3) Our algorithm is flexible such that we can easily include a new subgroup of interests and still maintain the fairness.

3. Toward Fair Image Generation with GANs

3.1. Problem Setting

Background. We assume there exists an unknown data distribution $P_{\text{data}}(\mathcal{X})$ over observed d -dimensional image

Table 1. Fairness analysis on different datasets with single attribute

Attributes	Eyeglasses	Age	Smiling	Gender	Black	Asian	White
FFHQ Dataset	0.191	0.034	0.002	9.40×10^{-7}	0.576	0.279	0.042
GAN	0.246	0.059	0.040	0.002	0.603	0.319	0.057
StyleFlow [1]	0.101	0.001	0.003	0.004	-	-	-
Ours	5.65×10^{-4}	9.68×10^{-6}	0	1.28×10^{-6}	3.45×10^{-4}	5.13×10^{-6}	2.00×10^{-8}

Table 2. Fairness analysis on different datasets with two attributes.

Attributes	Age-Gender	Age-Eyeglasses	Gender-Eyeglasses	Black-Gender	Asian-Gender	Black-Age	Asian-Age
FFHQ Dataset	0.0930	0.2994	0.2228	0.5762	0.2790	0.6097	0.3125
GAN	0.1205	0.4079	0.2808	0.6075	0.3305	0.6622	0.3919
StyleFlow [1]	0.1811	0.1755	0.0980	-	-	-	-
Ours	0.0201	0.0079	0.0013	0.0102	0.0007	0.0033	0.0018

data $\mathcal{X} \subseteq \mathbb{R}^d$. The goal for GAN is to learn distribution $P_\theta(\mathcal{X})$ such that $P_\theta(\mathcal{X}) = P_{\text{data}}(\mathcal{X})$. Given a well-trained GAN model, its generator can be seen as a mapping function $g_\theta : \mathcal{Z} \rightarrow \mathcal{X}$. Here, $\mathcal{Z} \subseteq \mathbb{R}^n$ denotes the n -dimensional latent space, which is usually assumed as Gaussian distribution $\mathcal{N}(\mathbf{0}, \mathbf{I}_d)$. With this model trained on the dataset D_{train} sampled from $P_{\text{data}}(\mathcal{X})$, we are able to obtain a generated distribution $P_\theta(\mathcal{X}) \approx P_{\text{data}}(\mathcal{X})$. By sampling latent codes \mathbf{z} , we can generate a realistic dataset $D_\theta = \{g_\theta(\mathbf{z}_i)\}_{i=1}^N$ with data size N .

Besides, we assume that each of the image sample x in \mathcal{X} contains specific semantic attributes, like age and gender of the face in the image. In the context of fairness, we focus on m target binary attributes \mathcal{A}_t , which we aim to achieve fairness over their distributions. Also, we have m' binary context attributes \mathcal{A}_c that do not require fairness. Suppose we have a attribute classifier as the semantic scoring function $f_S : \mathcal{X} \rightarrow \mathcal{S} \subseteq \mathbb{R}^{m+m'}$, where \mathcal{S} is the score space for all the attributes. It can be mapped to binary labels with unit step function $H(\cdot)$, such that $H(x) = 1$ if $x \geq 0$, otherwise $H(x) = 0$. This allows us to quantify the bias in terms of specific attributes by mapping the latent code \mathbf{z} to its attribute label \mathbf{a} with $\mathbf{a} = H(f_S(g_\theta(\mathbf{z})))$.

3.2. Measuring Data Bias and Generation Bias

In this section, we conduct an empirical study on the relationship between the bias in the training data and in a well-trained GAN model. To this end, we also develop metrics to measure the data generation bias.

Generation Bias. Real-world datasets often follow a long-tail distribution, thus carry various kinds of biases. The GAN model trained on real data D_{train} will therefore produce a dataset D_θ that highly likely carries the biases in the real data. Meanwhile, training GAN is prone to the mode collapse issue [33] which potentially brings in more biases. If the GAN models is biased, then in D_θ the marginal distribution of target attributes $P_{D_\theta}(\mathcal{A}_t)$ will be

highly imbalanced. We measure the imbalance with

$$f_u(D) = KL(P_{D_\theta}(\mathcal{A}_t) \| \mathcal{U}(\mathcal{A}_t)), \quad (1)$$

where KL is the Kullback-Leibler divergence and \mathcal{U} denotes the uniform distribution. A higher $f_u(D)$ indicates a imbalanced distribution in D .

To empirically show the bias introduced by the training process of GAN model, we computed the imbalance measurement f_u for a state-of-the-art face synthesis GAN model (StyleGAN2 [19]), which is well trained on the FFHQ dataset [18]. We first show the bias results w.r.t a single facial attribute (*e.g.* Gender) in Table 1. We observe that: 1) there is a strong correlation between the imbalance degree in the training data and in the GAN’s output, indicating that bias is carried from the training data to the GAN model; 2) the imbalance in the GAN’s output is consistently more significant than in the training data, which shows that bias is amplified during GAN’s training process. Then, we show bias results w.r.t the combination of two and three facial attributes in Table 2 and Table 3 respectively. We find that when we consider multiple facial attributes simultaneously, the bias in the GAN output becomes much more severe. In the next subsection, we will develop an effective baseline method to diminish such biases.

Measuring the Fairness of Image Generation. To reduce the bias in D_θ , one could either acquire a balanced training set [25] or retrain the GAN model with new objectives [7, 40]. However, both approaches are costly and require access to the original training dataset and hyper-parameters. Instead, we propose a new sampling strategy that could generate a fair dataset D_{fair} from the same GAN model.

On one hand, our main objective is to make the marginal distribution of the *target* attributes in the fair dataset $P_{D_{\text{fair}}}(\mathcal{A}_t)$ as close to the uniform distribution as possible, *i.e.*, to reduce $f_u(D_{\text{fair}})$ in Eq. (1). On the other hand, to ensure we do not introduce new bias w.r.t other attributes, we need to preserve the conditional distribution of the

context attributes \mathcal{A}_c in D_θ . Thus we formally define the fairness discrepancy f as

$$f(D_{\text{fair}}) = f_u(D_{\text{fair}}) + \beta KL(P_{D_{\text{fair}}}(\mathcal{A}_c | \mathcal{A}_t) \| P_{D_\theta}(\mathcal{A}_c | \mathcal{A}_t)). \quad (2)$$

where the target fair dataset D_{fair} is sampled from a generative model g_θ w.r.t attributes \mathcal{A}_t and \mathcal{A}_c , and β is used to balance the importance between the main objective of target fairness and the preservation constraint on context attributes. The lower $f(D_{\text{fair}})$ is, the more fair the generated dataset is with respect to \mathcal{A}_t .

3.3. Improving the Fairness of Image Generation via Shifting Latent Distribution

Based on the previous empirical analysis, we can see that pretrained GAN model carries and amplifies the bias in the imbalanced training data. It is impossible to fix all the biases in the models, and we are fully aware that any new debiasing methods might introduce some sort of new bias implicitly. Thus the goal of this work is not to fully solve the bias issue in GANs, instead, we aim at providing a new baseline to improve the fairness of GAN’s image generation to some degree through shifting the latent variable distribution. This baseline method, termed as Latent Distribution Shifting (**LDS**), utilizes the interpretable semantic dimensions identified in the latent space, then shifts the latent distribution for a more fair output.

The objective of our method is to sample a set of latent codes $\mathbf{Z}_{\text{fair}} = \{\mathbf{z}_i\}_{i=1}^N$, such that the GAN generated dataset $D_{\text{fair}} = \{g_\theta(\mathbf{z}_i)\}_{i=1}^N$ achieves low $f(D_{\text{fair}})$. To construct the fair latent code set \mathbf{Z}_{fair} that makes $P_{D_{\text{fair}}}(\mathcal{A}_t)$ close to $\mathcal{U}(\mathcal{A}_t)$, we propose to sample a set of latent codes *conditioned* on each of the possible value $\mathbf{a} \in \mathcal{A}_t$.

Specifically, for \mathcal{A}_t with m binary attributes, we have $K = 2^m$ possible attribute values. For the i th possible attribute value \mathbf{a} , our goal is to sample a set of latent codes \mathbf{Z}_i such that $H(f_S(g_\theta(\mathbf{z}))) = \mathbf{a}, \forall \mathbf{z} \in \mathbf{Z}_i$ and $|\mathbf{Z}_i| = N/K$. Then we can simply compose a fair latent code set with $\mathbf{Z}_{\text{fair}} = \{\mathbf{Z}_i\}_{i=1}^K$ such that $P_{D_{\text{fair}}}(\mathcal{A}_t) = \mathcal{U}(\mathcal{A}_t)$. This process is illustrated in Fig 2. The most critical component of this approach is attribute-conditioned latent code sampling. We split this process into three steps:

1. Create an intermediate code set \mathbf{Z}_{edit} by manipulating random latent codes towards specified condition \mathbf{a} .
2. Filter and fit the distribution of \mathbf{Z}_{edit} with the attribute scoring function and a Gaussian Mixture Model $q_\Phi(\mathbf{z})$.
3. Construct \mathbf{Z}_i by sampling $\mathbf{z} \sim q_\Phi(\mathbf{z})$.

We introduce each step in more details as follows.

3.3.1 Shifting Semantic Distribution in Latent Space

Manipulating Attributes in the Latent Space. InterFaceGAN is based on the assumption that for any binary

semantic attribute $\mathcal{A}_{t,i}$, there exists a hyperplane in the latent space \mathcal{Z} as the separation boundary. Latent codes on the same side of the boundary have the same attribute value, which turns into the opposite when the latent codes cross the boundary.

Assume the boundary for semantic $\mathcal{A}_{t,i}$ has the unit normal vector \mathbf{n}_i , we define its signed distance with a latent code \mathbf{z} as $d(\mathbf{z}, \mathbf{n}_i) = \mathbf{n}_i^T \mathbf{z}$. Given f_i as the scoring function for attribute $\mathcal{A}_{t,i}$, when \mathbf{z} move towards or away from \mathbf{n}_i , both $d(\mathbf{z}, \mathbf{n}_i)$ and $f_i(g_\theta(\mathbf{z}))$ would vary accordingly. Furthermore, when \mathbf{z} move across the boundary, both $d(\mathbf{z}, \mathbf{n}_i)$ and $f_i(g_\theta(\mathbf{z}))$ would change their numerical signs. Therefore, the authors assume a linear relationship such that

$$f_i(g_\theta(\mathbf{z})) = \lambda d(\mathbf{z}, \mathbf{n}_i), \quad (3)$$

where $\lambda > 0$ measures the ratio of changing speeds of semantic score and distance. According to this relationship, to manipulate the attribute score, we can easily vary the original code with $\mathbf{z}_{\text{edit}} = \mathbf{z} + \alpha \mathbf{n}_i$, as $f_i(g_\theta(\mathbf{z}_{\text{edit}})) = f_i(g_\theta(\mathbf{z})) + \lambda \alpha$. Please refer to Sec. 4.1 for the detail about how we obtain \mathbf{n}_i in practice.

Our goal is to make $f_i(g_\theta(\mathbf{z}_{\text{edit},i})) = \lambda \alpha$, where $\lambda \alpha > 0$ is a predefined scoring threshold such that we can classify the synthesised image with high confidence. Towards this goal, we first move \mathbf{z} onto the decision boundary and then move it away from the boundary with magnitude α . Formally, we have

$$\mathbf{z}_{\text{edit},i} = \mathbf{z} - d(\mathbf{z}, \mathbf{n}_i) \mathbf{n}_i + \alpha \mathbf{n}_i. \quad (4)$$

According to Eq. (3), we can prove that $f_i(g_\theta(\mathbf{z}_{\text{edit},i})) = \lambda \alpha$. The same conclusion holds for $\alpha_i = 0$, where $\alpha < 0$. This enables us to edit \mathbf{a}_i for any $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$.

The next step is to set attribute value \mathbf{a} for all m target attributes \mathcal{A}_m . To this end, we need the normal vectors of any pair of the semantic boundary to satisfy $\mathbf{n}_i^T \mathbf{n}_j \approx 0$, which we find to be the case for the GAN model and attributes we consider. This is because these attributes are already well disentangled in the latent space of the pretrained model. In this way, we can set

$$\mathbf{z}_{\text{edit}} = \mathbf{z} - \sum_{i=1}^m [d(\mathbf{z}, \mathbf{n}_i) \mathbf{n}_i - \alpha_i \mathbf{n}_i], \quad (5)$$

where $\alpha_i > 0$ if $\mathbf{a}_i = 1$ otherwise $\alpha_i < 0$. It is also easy to prove that $f_i(g_\theta(\mathbf{z}_{\text{edit}})) \approx \lambda \alpha_i$ for arbitrary attribute i . In this way, if Eq. (3) is satisfied, we are able to obtain

$$H(f_S(g_\theta(\mathbf{z}_{\text{edit}}))) = \mathbf{a}. \quad (6)$$

For a specified attribute subgroup \mathbf{a} , we can construct an intermediate code set \mathbf{Z}_{edit} with N_{edit} samples. This is done by simply sampling N_{edit} latent codes \mathbf{z} from $\mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ and mapping them into \mathbf{z}_{edit} with Eq. (5).

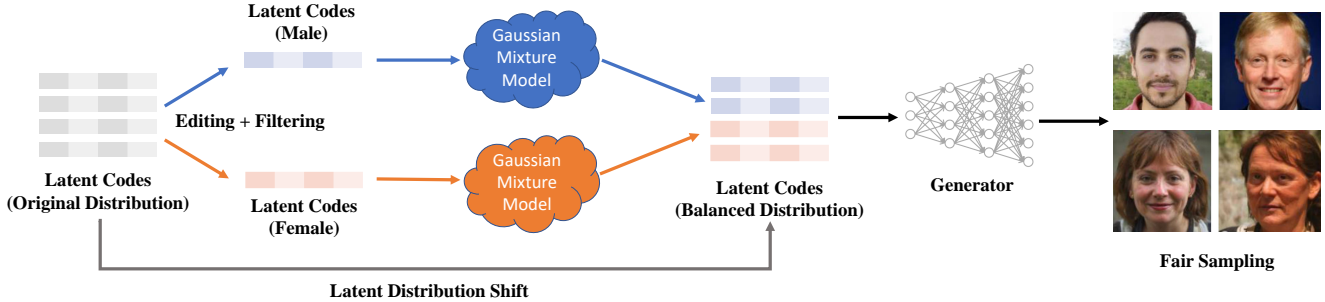


Figure 2. Overview of the proposed baseline to shift the latent distribution of a well-trained GAN model. Starting from the original distribution, we manage to collect a set of latent codes, which are able to synthesize balanced images across all subgroups of interests. The detailed procedure of the latent code collection can be found in Sec. 3.3. Then, for each subset of latent codes, a Gaussian Mixture Model (GMM) is used to fit the sub-distribution. Finally, we integrate these GMMs together as a balanced distribution, which naturally supports conditional sampling for any particular subgroup and hence improves the model fairness.

3.3.2 Conditional Latent Space Modeling

To support re-sampling latent code from a specific subgroup, we utilize a Gaussian Mixture Model (GMM) to fit its distribution in the latent space with the set of latent code \mathbf{Z}_{edit} created with Eq. (5).

However, as the linear relationship in Eq. (3) may not be the case for some latent codes and attributes in practice, not all samples from \mathbf{Z}_{edit} satisfy Eq. (6). To obtain more accurate distribution modeling, we further use f_S to filter out less confident samples for creating a new set: $\mathbf{Z}'_{\text{edit}} = \{\mathbf{z} \in \mathbf{Z}_{\text{edit}} | H(f_S(g_\theta(\mathbf{z}))) = \mathbf{a}\}$.

We then train a GMM model on $\mathbf{Z}'_{\text{edit}}$ with expectation-maximization (EM) algorithm to obtain a probabilistic model of latent codes $q_\Phi(\mathbf{z})$ conditioned on the specified subgroup \mathbf{a} . This model enables us to sample an arbitrary number of high-quality images from a certain subgroup.

To construct the fair latent code set \mathbf{Z}_{fair} , we firstly prepare the GMM model for each of the possible value in \mathcal{A}_t . We then compose \mathbf{Z}_{fair} with the same-size latent code set sampled from each of the GMM models.

Table 3. Fairness analysis on different datasets with three attributes.

Attributes	Age-Glasses-Gender	Black-Age-Gender	Asian-Age-Gender
FFHQ	0.3690	0.6694	0.3721
GAN	0.4771	0.7260	0.4575
StyleFlow [1]	0.3294	-	-
Ours	0.0170	0.0159	0.0077

4. Experiments

We evaluate LDS with a well-trained face synthesis GAN model, StyleGAN2 [19]. Specifically, in Sec. 4.1 we investigate how effective LDS can improve the fairness of an existing GAN model w.r.t different attribute settings.

Then, in Sec. 4.2, to show the potential impact of LDS, we utilize our generated fair dataset to reveal and quantify the biases in two commercial facial classification APIs as well as a state-of-the-art super-resolution model. Finally, in Sec. B, we conduct an ablation study of the design choices used in our model.

Implementation Details. For the GAN model g_θ we study, we use the official StyleGAN2 model² trained on the FFHQ dataset [18]. This model is able to generate realistic human face images with 1024×1024 resolution. For common attributes (*age, gender and eyeglasses*) we use an off-the-shelf classifier trained on Celeb-A dataset [22] as the scoring function f_S ; for the racial attributes, we use a classifier trained on LFW dataset [15]. Indeed these classifiers themselves might be biased or give some wrong predictions; for generality, we assume they are reasonably good to use.

The style-based generator in StyleGAN2 learns to first map the latent code from \mathcal{Z} space to another high-dimensional space \mathcal{W} . As shown in [32], latent codes in \mathcal{W} space have much stronger disentanglement property than in \mathcal{Z} , and latent code manipulation quality is also better in \mathcal{W} space. Therefore, in our framework, we directly operate and sample codes from the \mathcal{W} space, which is simply mapped from \mathcal{Z} space with StyleGAN2 mapping module.

To obtain the semantic boundaries \mathbf{n} used in Eq. (4), we first synthesize a dataset $\mathcal{D}_{\text{direct}}$ with $50K$ images by directly sampling from the original latent space. Then, we use the scoring function f_S to obtain the attribute scores for all the sensitive attributes. For each attribute, we sort the corresponding scores and choose the ones with top 2% highest scores as positive examples and top 2% lowest scores as negative examples. This is to select the most representative examples as the scoring function may not be absolutely accurate. Finally, we use these examples to train a linear SVM to obtain the decision boundary, the normal

²<https://github.com/NVlabs/stylegan2>.

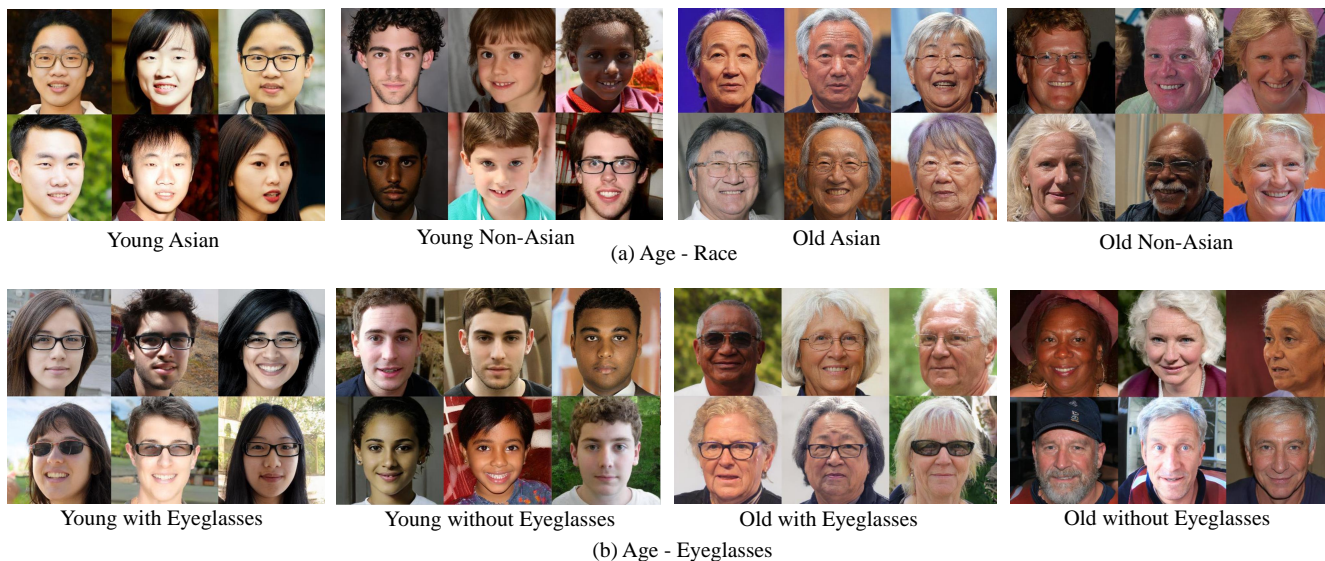


Figure 3. Qualitative results for fair image generation in GANs with two-attribute values.

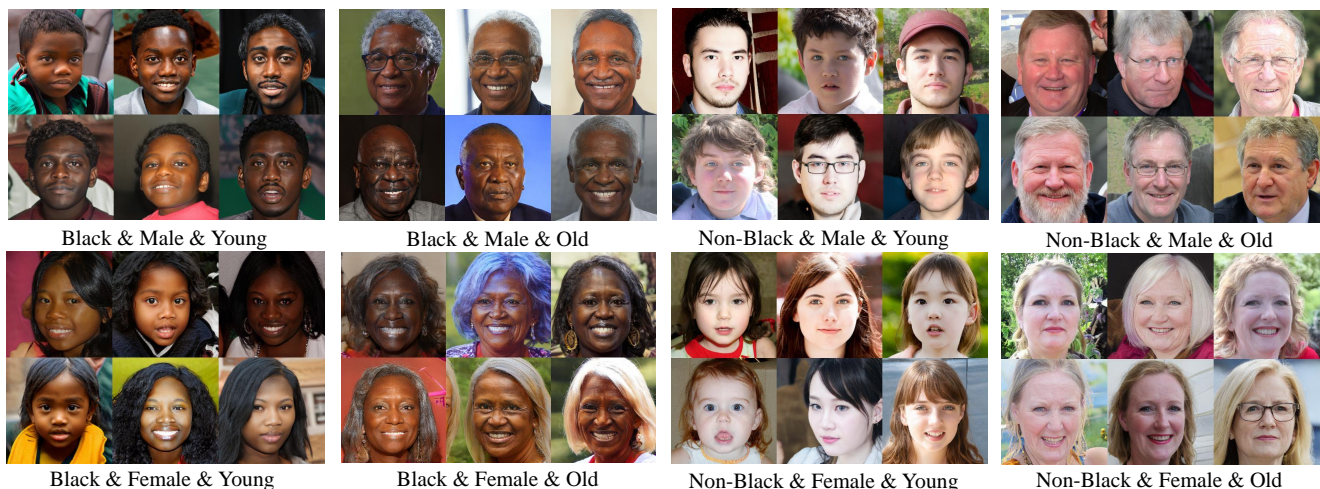


Figure 4. Qualitative results of fair image generation with respect to three different facial attributes, *i.e.*, race, age, and gender.

direction of which results in \mathbf{n} . The SVM is trained to take the latent codes in \mathcal{W} space as input, and output binary labels obtained with f_S .

For InterFaceGAN, we set the magnitude of code editing factor $|\alpha| = 3.0$. For the conditional latent space modeling of each attribute subgroup, we manipulate and generate \mathcal{Z}_{edit} with $N_{edit} = 2.5K$ samples; then we use a GMM model with $k = 10$ components to fit the distribution. We provide an empirical analysis of these hyper-parameters in the **Supplementary Material**.

4.1. Fair Image Generation

We firstly show the existing bias in the GAN model w.r.t different attribute settings. Then, we present both the quantitative and qualitative results that show our LDS method can significantly improve the fairness of image

generation while preserving the image quality.

Experiment Setting. For evaluation, We use 3 common attributes (*age*, *gender* and *eyeglasses*) and the race attributes (*Black*, *Asian* and *White*). To form different sampling tasks, we combine n of the attributes as a pair of target attributes \mathcal{A}_t to form subgroups for each task while leaving the other attributes in the context set \mathcal{A}_c . For each of the compared sampling methods, we sample 10K images in total and evaluate on this generated dataset. We then compute the generative sampling fairness discrepancy f with Eq. (2) for each of the datasets, where we set $\beta = 0.1$ in f . We conduct experiments with $n = 1, 2, 3$, respectively.

Quantitative Evaluation. We compare LDS with two baselines. 1) **GAN**: we directly sample latent codes from the original distribution of the GAN. 2) **StyleFlow** [1], which trains conditional continuous normalizing flows

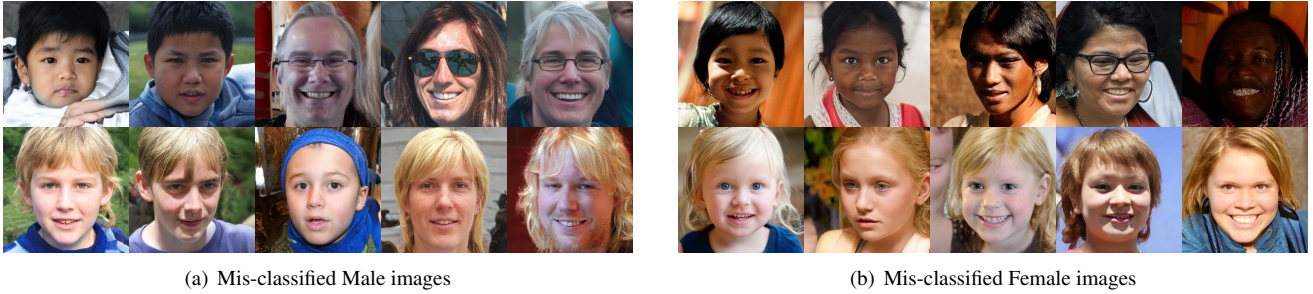


Figure 5. Mis-classified images by APIs.

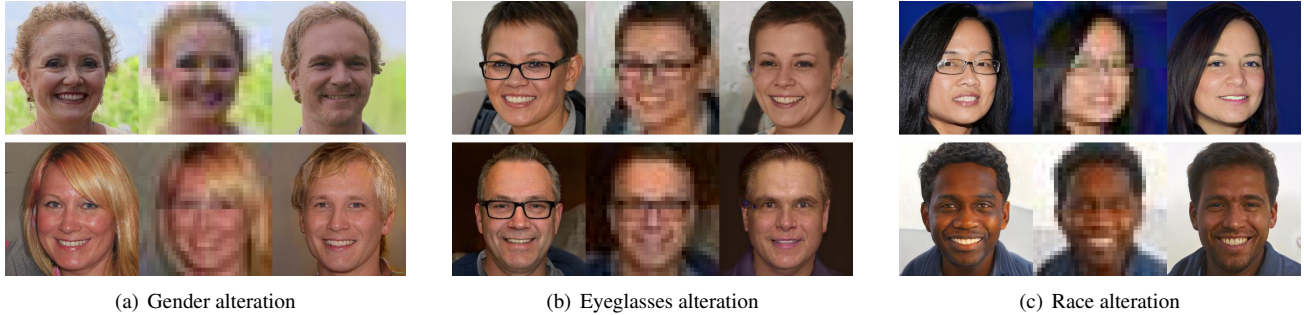


Figure 6. Examples of attribute alteration by the super-resolution model. From left to right we showcase 1) the original image generated by our method; 2) the LR image down-sampled from the original image; 3) HR image output by PULSE given the LR image.

Table 4. Gender classification error rate (in percentage) in different attribute subgroups combined with gender.

Sensitive Attributes Subgroups	Gender		Age		Eyeglasses		Race		
	Male	Female	Young	Old	With	Without	Black	Asian	White
Face++ Detect API	0.81	4.19	2.02	1.76	0.58	3.40	3.50	0.64	2.28
Azure Facial Recognition	0.71	2.52	1.08	1.00	0.47	2.36	1.62	0.84	0.44

to support attribute-conditional sampling from the GAN model. Here we use the official-released StyleFlow model³ to uniformly sample latent codes from each attribute subgroup in a way similar to our method. Note that this model currently does not support latent code sampling conditioned on the race attributes.

We show the fairness results with different numbers of attributes in Tab. 1, Tab. 2, and Tab. 3 respectively. The results show that: 1) LDS can significantly improve the fairness of GAN model across various attribute combinations. Furthermore, it is able to easily handle more than one attribute at one time. For example, for *age-eyeglasses* task in Tab. 2, LDS significantly decreases the bias score f of GAN model from 0.4079 to 0.0079, showing the high effectiveness of LDS. 2) LDS consistently outperforms StyleFlow across multiple attribute combinations.

Qualitative Evaluation. Fig. 3 and Fig. 4 plot the generated images from the codes sampled by LDS for subgroups regarding two and three face attributes. It suggests that LDS performs well to generate images for all the subgroups with correct attributes while retaining high image quality and diversity. This quality enables us to use the generated bal-

³<https://github.com/RameenAbdal/StyleFlow>.

anced dataset to examine the fairness of other visual tasks and models. Please refer to the **Supplementary Material** for image examples from more attribute subgroups.

4.2. Quantifying Bias in Existing Models

The fair image generation achieved by LDS can be useful in many applications. Here we apply our method to reveal and quantify the potential biases in the existing face classifiers and a super-resolution model.

Bias in Face Classifiers. We first study the bias in existing face classifiers. To make it more practical, we select two state-of-the-art commercial face attribute classification APIs (Face++ Detect API and Azure Facial Recognition). To analyze the bias problem, we focus on gender classification under different attribute conditions: *age*, *eyeglasses* and *race*.

Specifically, for each attribute subgroup (e.g., Young Male), we use our method to generate a subgroup dataset with 2.5K images. Then, we run the face classifiers on these datasets and compare the error rate of different groups.

We show the results in Tab. 4, which well quantifies the bias problem in the two APIs. In this table, besides the error rate of gender in each subgroup, we also compute

Table 5. Percentage of attribute alternation by the super-resolution model.

Attributes Groups	Gender		Age		Glasses		Race		
	Male	Female	Young	Old	With	Without	Black	Asian	White
Value alternation rate	3.4	7.8	3.3	29.5	89.8	0.3	76.4	34.3	0.5

the average error rates for male and female people over all the subgroups. We first observe that for both APIs, the gender classification accuracy for females is significantly lower than for the males. Also, people with black skin color are more likely to be wrongly classified, while accuracy is more balanced w.r.t age. We show some of the failure cases observed by our model in Fig. 5.

Bias in Super-resolution Model. In this part, we study the bias problem of super-resolution method PULSE [26]. This recent neural network model takes a low-resolution (LS) image as input and outputs a high-resolution (HS) image. It has been found that for certain minority groups, their rare attribute values in the LS input will often be changed to more common values in the HS output [20]. We aim to use the images generated with LDS to examine which attributes will be more likely to be altered by PULSE.

To this end, we first generate images from a certain subgroup with LDS, and then input its down-sampled LS (32×32) version to PULSE to obtain a HS (1024×1024) output. Then, we use the scoring function f_S to obtain the attribute values of the HS images. Finally, we compare this result with the original image’s attribute value, and compute the rate of PULSE alternation for each of the attributes. Here we select four attributes *Gender*, *Eyeglasses*, *Age* and *Race*. We exhibit some examples of the image attribute alternated by PULSE in Fig. 6. We can see the PULSE wrongly alters the gender, glasses and race attributes through the super-resolution process.

We show the rate of attribute alternation of PULSE in Tab. 5. Here each number represents the alternation rate in a subgroup. We have the following fairness analysis: 1) For the race attributes, the alternation rate is much higher for people with Black (76.4%) and Asian (34.3%) race, while few images (0.5%) with White race people are alternated. This indicates PULSE is prone to output people with white race. 2) For the gender attribute, the alternation rate is higher for female people, indicating PULSE is prone to output male faces. 3) For the eyeglasses attribute, we observe a very high alternate rate (89.8%) for input images with eyeglasses, which shows that PULSE often fails to preserve the eyeglasses.

It is worth noting that none of the above bias analysis requires additional human labeling or careful dataset balancing for each of the attributes, making the bias analysis based on LDS much more convenient to run than the previous bias analysis works [5, 25]. With the support to sample images from any attribute subgroup, LDS enables us to do such detailed analysis with low cost on the biases

in existing models w.r.t different subgroups. In addition, we are able to easily extend the analysis to other attributes when provided with scoring function f_S of new attributes. This makes LDS a general and flexible tool for studying the bias of visual models.

4.3. Ablation Study

In this section, we analyze the impact of different design choices in LDS on the fairness score. Specifically, we focus on three variations of LDS by ablating different components: 1) **Ours w/o Edit**: We use f_S to filter latent codes for different subgroups with the latent code *directly* sampled from the original distribution; then fit GMMs with these latent codes; 2) **Ours w/o Filter**: We *skip* the subgroup attribute filter before using GMM to fit the shifted latent codes for each subgroup; 3) **Ours w/o GMM**: We generate the dataset by directly using the latent coded generated in Sec. 3.3.1.

We show the results on three tasks in Figure 6, which proves that all the components in our method are important to fairness performance of the generated dataset. Particularly, we find the latent code shifting (editing) process is very important in all the cases. This is because we are able to obtain more diverse sets of latent codes for the rare subgroups compared with original latent code distribution.

Table 6. Ablation study of our method.

Attributes	age- eyeglasses	gender- eyeglasses	black- gender
Ours w/o Edit	0.0497	0.0139	0.0292
Ours w/o Filter	0.0200	0.0053	0.0404
Ours w/o GMM	0.0122	0.0031	0.0267
Ours	0.0079	0.0013	0.0102

5. Conclusion

In this work, following the empirical study on the bias in the training set and the trained generative model, we develop a baseline method to eliminate the bias within a well-trained GAN model by discreetly altering the sampling strategy yet retaining the model weights. Our method has shown great potential in turning numerous publicly available GAN models unbiased and helping detect the unfairness in other AI systems. We hope our work can raise the awareness of the fairness in generative models, and our method will be a starting point to diminish and eventually eliminate all the potential biases in generative models.

References

- [1] Rameen Abdal, Peihao Zhu, Niloy Mitra, and Peter Wonka. Styleflow: Attribute-conditioned exploration of stylegan-generated images using conditional continuous normalizing flows. *arXiv e-prints*, pages arXiv–2008, 2020. 3, 5, 6
- [2] Alekh Agarwal, Alina Beygelzimer, Miroslav Dudík, John Langford, and Hanna Wallach. A reductions approach to fair classification. In *Proceedings of Machine Learning Research*, 2018. 2
- [3] Solon Barocas, Moritz Hardt, and Arvind Narayanan. Fairness in machine learning. *NIPS Tutorial*, 2017. 1, 2
- [4] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018. 2
- [5] Joy Buolamwini and Timnit Gebru. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on Fairness, Accountability and Transparency*, 2018. 1, 2, 8
- [6] Irene Chen, Fredrik D Johansson, and David Sontag. Why is my classifier discriminatory? In *Adv. Neural Inform. Process. Syst.*, 2018. 1
- [7] Kristy Choi, Aditya Grover, Trisha Singh, Rui Shu, and Stefano Ermon. Fair generative modeling via weak supervision. In *ICML*, 2019. 2, 3
- [8] Yunjey Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2018. 2
- [9] Sam Corbett-Davies and Sharad Goel. The measure and mismeasure of fairness: A critical review of fair machine learning. *arXiv preprint arXiv:1808.00023*, 2018. 1
- [10] Elliot Creager, David Madras, Jörn-Henrik Jacobsen, Marissa A Weis, Kevin Swersky, Toniann Pitassi, and Richard Zemel. Flexibly fair representation learning by disentanglement. In *Proceedings of Machine Learning Research*, 2019. 1
- [11] Michael Feldman. *Computational fairness: Preventing machine-learned discrimination*. PhD thesis, 2015. 2
- [12] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Adv. Neural Inform. Process. Syst.*, 2014. 2
- [13] Jinjin Gu, Yujun Shen, and Bolei Zhou. Image processing using multi-code gan prior. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020. 2
- [14] Moritz Hardt, Eric Price, and Nati Srebro. Equality of opportunity in supervised learning. In *Adv. Neural Inform. Process. Syst.*, 2016. 1, 2
- [15] Gary B. Huang, Manu Ramesh, Tamara Berg, and Erik Learned-Miller. Labeled faces in the wild: A database for studying face recognition in unconstrained environments. Technical Report 07-49, University of Massachusetts, Amherst, October 2007. 5
- [16] Niharika Jain, Alberto Olmo, Sailik Sengupta, Lydia Manikonda, and Subbarao Kambhampati. Imperfect image augmentation: Implications of gans exacerbating biases on facial data augmentation and snapchat selfie lenses. *arXiv preprint arXiv:2001.09528*, 2020. 2
- [17] Matthew Joseph, Michael Kearns, Jamie H Morgenstern, and Aaron Roth. Fairness in learning: Classic and contextual bandits. In *Adv. Neural Inform. Process. Syst.*, 2016. 1
- [18] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019. 2, 3, 5, 11
- [19] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020. 1, 2, 3, 5, 11
- [20] Andrey Kurenkov. Lessons from the pulse model and discussion. <https://thegradient.pub/pulse-lessons/>, 2020. 2, 8
- [21] Yijun Li, Chen Fang, Jimei Yang, Zhaowen Wang, Xin Lu, and Ming-Hsuan Yang. Universal style transfer via feature transforms. In *Adv. Neural Inform. Process. Syst.*, 2017. 2
- [22] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Int. Conf. Comput. Vis.*, 2015. 5
- [23] Christos Louizos, Kevin Swersky, Yujia Li, Max Welling, and Richard Zemel. The variational fair autoencoder. In *Int. Conf. Learn. Represent.*, 2016. 2
- [24] Kristian Lum and James Johndrow. A statistical framework for fair predictive algorithms. *arXiv preprint arXiv:1610.08077*, 2016. 2
- [25] Daniel McDuff, Shuang Ma, Yale Song, and Ashish Kapoor. Characterizing bias in classifiers using generative models. In *Adv. Neural Inform. Process. Syst.*, 2019. 2, 3, 8
- [26] Sachit Menon, Alexandru Damian, Shijia Hu, Nikhil Ravi, and Cynthia Rudin. Pulse: Self-supervised photo upsampling via latent space exploration of generative models. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020. 2, 8, 22
- [27] Xingang Pan, Xiaohang Zhan, Bo Dai, Dahua Lin, Chen Change Loy, and Ping Luo. Exploiting deep generative prior for versatile image restoration and manipulation. In *Eur. Conf. Comput. Vis.*, 2020. 2
- [28] Hee Jung Ryu, Hartwig Adam, and Margaret Mitchell. Inclusivefacenet: Improving face attribute detection with race and gender diversity. *arXiv preprint arXiv:1712.00193*, 2017. 1
- [29] Hee Jung Ryu, Margaret Mitchell, and Hartwig Adam. Improving smiling detection with race and gender diversity. *arXiv preprint arXiv:1712.00193*, 2017. 1
- [30] Joni Salminen, Soon-gyo Jung, Shammur Chowdhury, and Bernard J Jansen. Analyzing demographic bias in artificially generated facial pictures. In *Extended Abstracts of the 2020 CHI Conference on Human Factors in Computing Systems*, 2020. 2
- [31] Prasanna Sattigeri, Samuel C Hoffman, Vijil Chenthamarakshan, and Kush R Varshney. Fairness gan: Generating datasets with fairness properties using a generative adversarial network. *IBM Journal of Research and Development*, 2019. 2
- [32] Yujun Shen, Ceyuan Yang, Xiaoou Tang, and Bolei Zhou. Interfacegan: Interpreting the disentangled face representation

- learned by gans. *IEEE Trans. Pattern Anal. Mach. Intell.*, 2020. [2](#), [5](#)
- [33] Akash Srivastava, Lazar Valkov, Chris Russell, Michael U Gutmann, and Charles Sutton. Veegan: Reducing mode collapse in gans using implicit variational learning. In *Adv. Neural Inform. Process. Syst.*, 2017. [1](#), [3](#)
 - [34] Antonio Torralba and Alexei A Efros. Unbiased look at dataset bias. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2011. [1](#)
 - [35] Xintao Wang, Ke Yu, Shixiang Wu, Jinjin Gu, Yihao Liu, Chao Dong, Yu Qiao, and Chen Change Loy. Esrgan: Enhanced super-resolution generative adversarial networks. In *Eur. Conf. Comput. Vis. Worksh.*, 2018. [2](#)
 - [36] Blake Woodworth, Suriya Gunasekar, Mesrob I Ohannesian, and Nathan Srebro. Learning non-discriminatory predictors. In *Proceedings of Machine Learning Research*, 2017. [2](#)
 - [37] Depeng Xu, Shuhan Yuan, Lu Zhang, and Xintao Wu. Fairgan: Fairness-aware generative adversarial networks. In *IEEE Int. Conf. Big Data*, 2018. [2](#)
 - [38] Yinghao Xu, Yujun Shen, Jiapeng Zhu, Ceyuan Yang, and Bolei Zhou. Generative hierarchical features from synthesizing images. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021. [2](#)
 - [39] Ceyuan Yang, Yujun Shen, and Bolei Zhou. Semantic hierarchy emerges in deep generative representations for scene synthesis. *Int. J. Comput. Vis.*, 2020. [2](#)
 - [40] Ning Yu, Ke Li, Peng Zhou, Jitendra Malik, Larry Davis, and Mario Fritz. Inclusive gan: Improving data and minority coverage in generative models. In *Eur. Conf. Comput. Vis.*, 2020. [2](#), [3](#)
 - [41] Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez Rogriguez, and Krishna P Gummadi. Fairness constraints: Mechanisms for fair classification. In *Artificial Intelligence and Statistics*, 2017. [2](#)
 - [42] Rich Zemel, Yu Wu, Kevin Swersky, Toni Pitassi, and Cynthia Dwork. Learning fair representations. In *ICML*, 2013. [2](#)
 - [43] Shengjia Zhao, Hongyu Ren, Arianna Yuan, Jiaming Song, Noah Goodman, and Stefano Ermon. Bias and generalization in deep generative models: An empirical study. In *Adv. Neural Inform. Process. Syst.*, 2018. [2](#)

Appendix

A. Data Distribution Comparisons

In this section, we show more comparisons of the data distributions of different datasets w.r.t the target attributes. Specifically, we compare the data distribution of the real training dataset FFHQ [18], the dataset directly sampled with GAN [19], and the dataset sampled with our method.

In Figure 7 and 9 we show results for the other four tasks in addition to the task we show in the paper. We observe that our method is able to consistently remove the bias in the GAN model without retraining across multiple tasks.

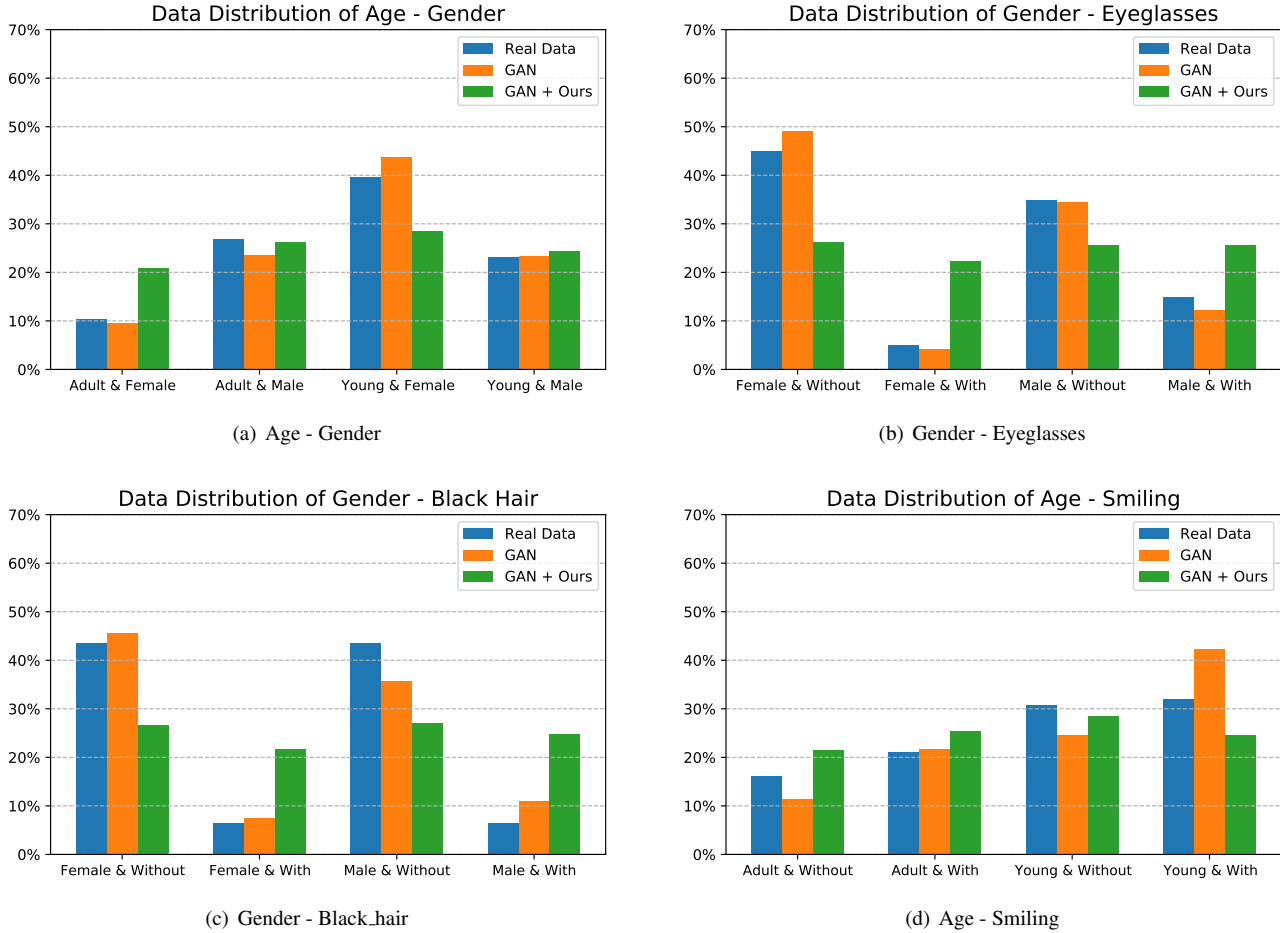


Figure 7. Comparisons of data distributions in different datasets.

B. Hyper-parameter Study

In this section, we analyze the impact of hyper-parameter choices in FairGen on the fairness score. In particular, we analyze 1) the magnitude of InterfaceGAN manipulation $|\alpha|$; 2) the size of edited latent code set N_{edit} ; 3) the number of components used in GMM models k .

We plot the results in Table 8. We observe that we need a large enough $|\alpha|$ to make sure the semantics of the shifted latent codes are correct. Secondly, the size of edited codes N_{edit} has a relatively small impact on the fairness score compared to other parameters. Finally, The more components we have in GMM, the better result we will normally obtain as more components provide a more accurate approximation of code distribution.

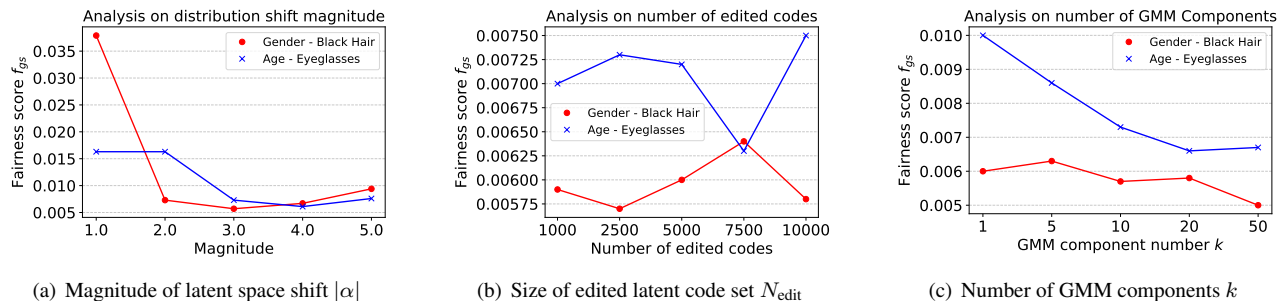
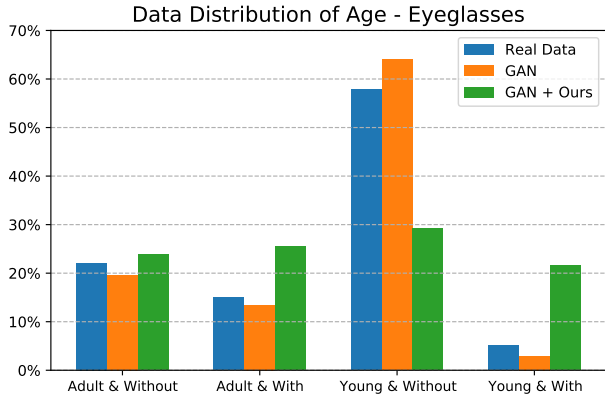
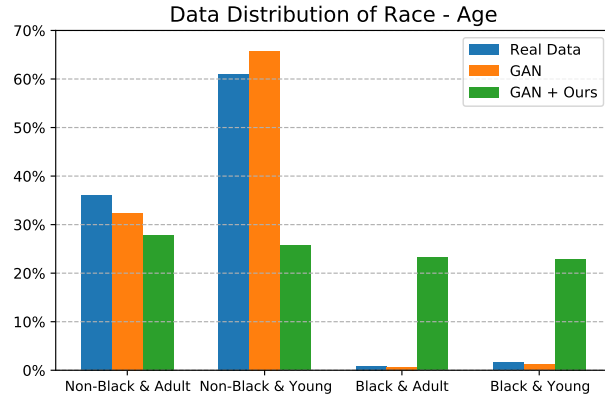


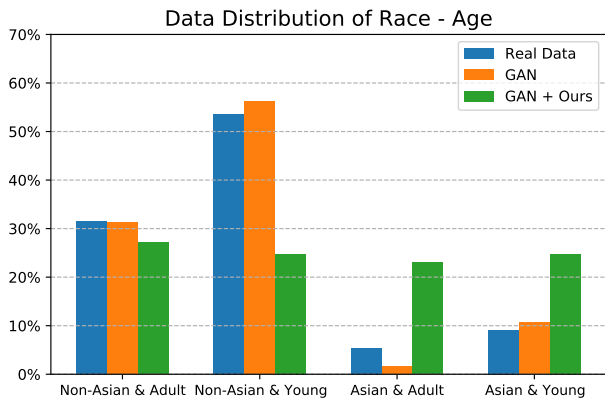
Figure 8. Results of the analysis of each of the hyper-parameters in our framework.



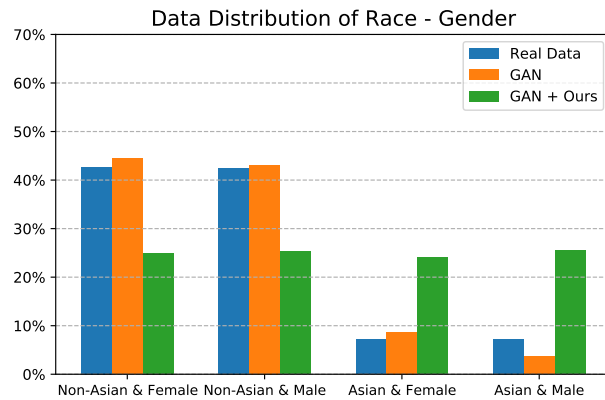
(a) Age - Eyeglasses



(b) Black - Age



(c) Asian - Age



(d) Asian - Gender

Figure 9. Comparisons of data distributions in different datasets (continued).

C. Commercial API details

We detail how we test the commercial APIs in our experiments.

For the evaluation of facial gender classification, we first obtain access to two commercial APIs: MEGVII’s Face++ Detect API (<https://www.faceplusplus.com/face-detection/>) and Microsoft’s Azure Facial Recognition service (<https://azure.microsoft.com/en-us/services/cognitive-services/face/>). Then, we utilize FairGen to generate facial images in different subgroups. We use the gender attribute value of each subgroup as the ground-truth gender label for the images in that subgroup. After that, we use both APIs to detect and analyze the face in each image, and then compare the predicted gender attribute with the ground-truth label to obtain the accuracies in the paper. Note that the APIs might fail to detect the face in a few images (around 0.5%), which are ignored during gender classification accuracy computation.

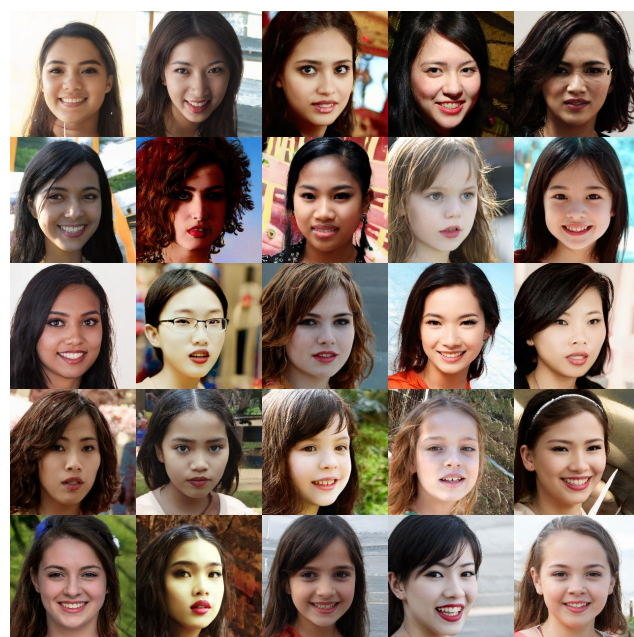
It is worth noting that here we are not claiming the defects or bugs in the commercial products. We just would like to raise the awareness of the potential bias in the existing applications through this small and humble academic work.

D. Conditional Generated Samples

In this section, we show more examples of the attribute subgroup images generated by our method for each task.



(a) Young Male



(b) Young Female



(c) Old Male



(d) Old Female

Figure 10. Qualitative results for fair image generation in GANs with Age and Gender.



(a) Young with Eyeglasses



(b) Young without Eyeglasses



(c) Old with Eyeglasses



(d) Old without Eyeglasses

Figure 11. Qualitative results for fair image generation in GANs with Age and Eyeglasses.



(a) Male with Eyeglasses



(b) Male without Eyeglasses



(c) Female with Eyeglasses



(d) Female without Eyeglasses

Figure 12. Qualitative results for fair image generation in GANs with Gender and Eyeglasses.



(a) Young Black



(b) Old Black

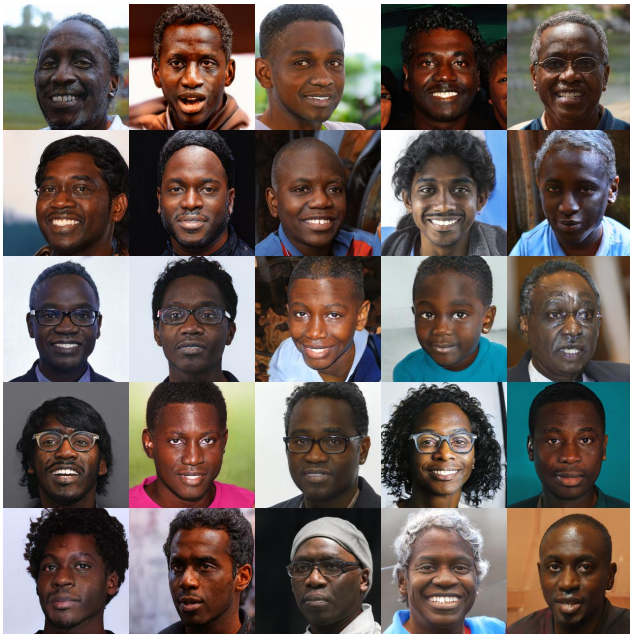


(c) Young Non-Black



(d) Old Non-Black

Figure 13. Qualitative results for fair image generation in GANs with Black and Age.



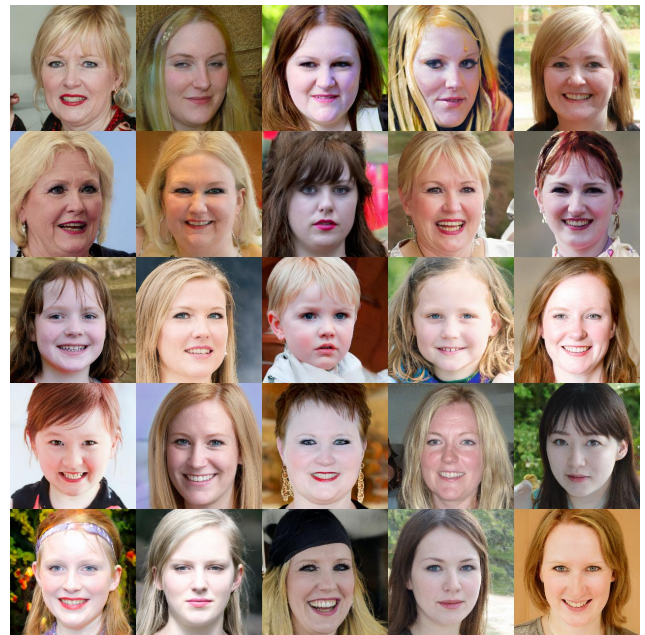
(a) Male Black



(b) Female Black



(c) Male Non-Black



(d) Female Non-Black

Figure 14. Qualitative results for fair image generation in GANs with Black and Gender.



(a) Young Asian



(b) Old Asian



(c) Young Non-Asian

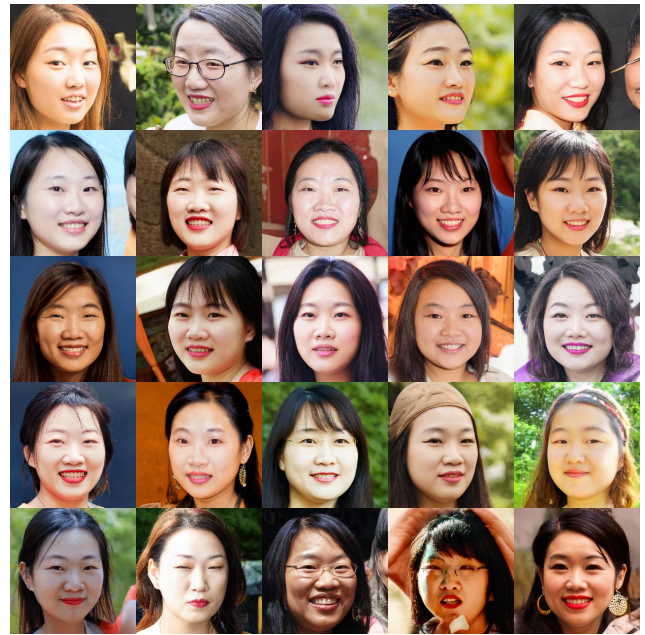


(d) Old Non-Asian

Figure 15. Qualitative results for fair image generation in GANs with Asian and Age.



(a) Male Asian



(b) Female Asian



(c) Male Non-Asian



(d) Female Non-Asian

Figure 16. Qualitative results for fair image generation in GANs with Asian and Gender.



(a) Male with Black Hair



(b) Male without Black Hair



(c) Female with Black Hair



(d) Female without Black Hair

Figure 17. Qualitative results for fair image generation in GANs with Gender and Black Hair.



(a) Young with Smiling



(b) Young without Smiling



(c) Old with Smiling



(d) Old without Smiling

Figure 18. Qualitative results for fair image generation in GANs with Age and Smiling.

E. Bias in Super-Resolution Models

In this section, we show more examples of the failure cases we expose from the super-resolution model PULSE [26].

For each set of images, from left to right we showcase 1) the original image generated by our method; 2) the LR image subsampled from the original image; 3) HR image output by PULSE given the LR image.



Figure 19. Examples of Gender attribute alteration by the super-resolution model.



Figure 20. Examples of Eyeglasses attribute alteration by the super-resolution model.

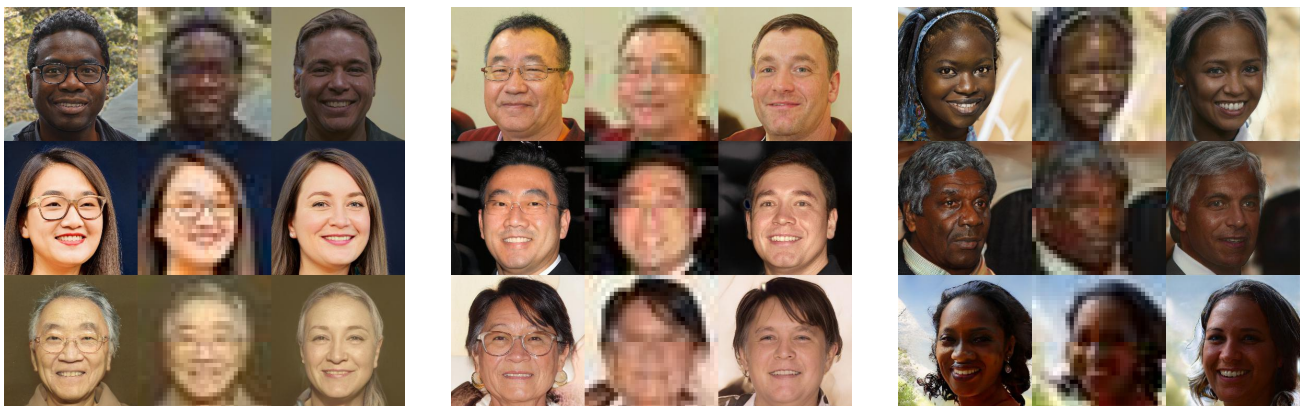


Figure 21. Examples of Race attribute alteration by the super-resolution model.

F. Bias in Gender Classification Models

In this section, we show more examples of the failure cases of gender classification for the two APIs we show in the paper. Specifically, we show the failure cases where our generated male faces are classified as female and vice versa.



Figure 22. Mis-classified Male Images by the Commercial Gender Classification APIs



Figure 23. Mis-classified Female Images by the Commercial Gender Classification APIs