

FairFaceGAN: Fairness-aware Facial Image-to-Image Translation

Sunhee Hwang
sunny16@yonsei.ac.kr

Sungho Park
skrtjdgh18@yonsei.ac.kr

Dohyung Kim
dohkim02@yonsei.ac.kr

Mirae Do
mwwdo109@yonsei.ac.kr

Hyeran Byun*
hrbyun@yonsei.ac.kr

Department of Computer Science
Yonsei University
Seoul, Republic of Korea

Abstract

In this paper, we introduce FairFaceGAN, a fairness-aware facial Image-to-Image translation model, mitigating the problem of unwanted translation in protected attributes (e.g., gender, age, race) during facial attributes editing. Unlike existing models, FairFaceGAN learns fair representations with two separate latents - one related to the target attributes to translate, and the other unrelated to them. This strategy enables FairFaceGAN to separate the information about protected attributes and that of target attributes. It also prevents unwanted translation in protected attributes while target attributes editing. To evaluate the degree of fairness, we perform two types of experiments on CelebA dataset. First, we compare the fairness-aware classification performances when augmenting data by existing image translation methods and FairFaceGAN respectively. Moreover, we propose a new fairness metric, namely *Fréchet Protected Attribute Distance* (FPAD), which measures how well protected attributes are preserved. Experimental results demonstrate that FairFaceGAN shows consistent improvements in terms of fairness over the existing image translation models. Further, we also evaluate image translation performances, where FairFaceGAN shows competitive results, compared to those of existing methods.

1 Introduction

Artificial Intelligence (AI) systems have achieved remarkable success in a broad range of research fields such as computer vision, natural language processing, and audio analysis. However, outputs of the AI systems could be biased since they heavily rely on human-collected datasets which may contain ethically sensitive stereotypes [1]. Research and articles indicated that several AI systems yielded unfair results with respect to protected attributes such as gender, age, or race [2, 3, 4, 5, 6, 7, 8]. This is a critical problem to computer vision

* Corresponding Author

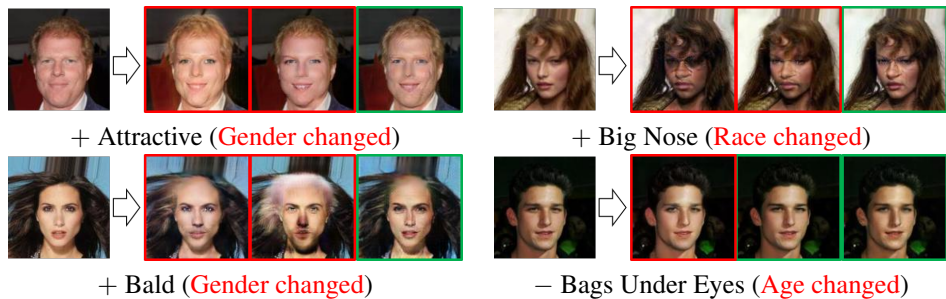


Figure 1: Image translation results on CelebA dataset [20]. For each example, we present four facial images, which are an input image and the results of StarGAN, FixedPointGAN, and FairFaceGAN (ours), respectively (from left to right). + and - denote adding and removing the attribute of the input image, respectively. While Red boxes indicate the occurrence of unwanted translation of protected attributes, Green boxes denote the preservation of protected attributes. Best viewed in color.

systems, which have already been deployed in diverse real world applications without adjusting demographic disparities. For example, PULSE algorithm [24], taking low-resolution faces into high-resolution images, tends to produce racially biased results, *i.e.*, white skin, blue eyes, and brown hair, regardless of input images [30]. Accordingly, in order to resolve the societal bias problem, researchers have directed their attention on developing fair computer vision models [0, 24, 18, 24, 31, 32, 33].

In this paper, we aim to improve fairness in Image-to-Image translation of facial attributes, whose goal is to edit attributes of input images. Even though recent methods based on Generative Adversarial Networks (GANs) [10] succeeded in synthesizing realistic facial images while translating attributes fairly, they might contain unintended discriminative factors. In Figure 1, we present several examples of discriminatory translation results. While translating of target attributes, existing facial attribute editing models [0, 28] unintendedly modify protected attributes (*i.e.*, gender, age, race) as well.

To address this problem, we propose a fairness-aware Image-to-Image translation model, namely FairFaceGAN, which maps input images into target domains while preserving protected attributes. In specific, we introduce a new fair representation learning method that learns two separate latent spaces with different objectives: (i) one is for mapping target attributes adequately; (ii) the other is for preserving information of protected attributes. By employing two decoupled latent spaces, FairFaceGAN successfully prevents unwanted translation during editing target attributes, as shown in the last column of each example of Figure 1. We note that our method can be easily extended to the case of multiple protected attributes as it separates target attributed-related information from the rest. Moreover, another merit of FairFaceGAN is that it does not require protected attribute annotations. Instead, we exploit knowledge related to protected attributes from a pre-trained classification model. We believe that this will largely benefit the application of our method especially in the circumstance where protected attribute labels are not acquirable.

To compare FairFaceGAN with existing image translation models in terms of fairness, we design and perform two kinds of experiments. Specifically, for the first one, we measure how the fairness-aware classification performances are improved when the biased training dataset is augmented by previous translation models and ours respectively. For this, we use standard fairness metrics, *i.e.*, Equality of Opportunity [12] and Equalized Odds [35]. For the

second one, we propose a new fairness metric, *Fréchet Protected Attribute Distance* (FPAD), inspired by Fréchet Inception Distance (FID) [13], to evaluate the protected attribute preservation ability of image translation models. On the both types of experiments, FairFaceGAN shows consistently fairer results over the existing image translation methods. Also, we provide comparisons on the standard image translation metrics, *i.e.*, FID and Kernel Inception Distance (KID), where FairFaceGAN achieves comparable results to the other models.

Our main contributions can be summarized as follows:

- We introduce FairFaceGAN that maps input images into target domain in a fair way with respect to multiple protected attributes.
- To reduce the correlation between protected and target attributes in the mapping, we propose to learn two separate representations with different objectives: target attributes mapping and protected attribute preservation.
- To achieve fairness, we present a knowledge transfer technique for fair translation on the target dataset. It enables our model to mitigate bias related to multiple protected attributes even for the case where annotations for protected attributes are unavailable.
- Through the extensive experiments on CelebA, we demonstrate that FairFaceGAN produces the fairest results in terms of Equality of Opportunity, Equalized Odds, and the proposed FPAD over existing Image-to-Image translation models.

2 Related Work

2.1 Fairness in Computer Vision

In recent years, fairness in computer vision has become a popular research topic. Among various types of fair methods, we briefly introduce two approaches to mitigate bias problems in visual recognition tasks: (1) Reorganizing a biased dataset to the fair dataset (Pre-processing), and (2) Reducing bias through model architecture or algorithm (In-processing).

Pre-processing. Sattigeri *et al.* [27] proposed a fair data generating method based on GANs. They are trained on a biased dataset and generate new data which are fair in terms of the protected attributes. The generated data is utilized to train a fairness-aware face attribute classification model. Quadrianto *et al.* [24] introduced a data-to-data translation method that transforms an original biased dataset into a new fair dataset. In this paper, we also address fairness in the image classification task by generating fair dataset using our FairFaceGAN.

In-processing. Zheng *et al.* [40] proposed a disentangling method that splits feature representation into the two subspaces, one relevant to target labels and the irrelevant one. Similarly, FFVAE [5] aim to represent protected attribute related information and the rest. Park *et al.* [23] proposed a fair disentangling method for representing target, protected attribute, and mutual information of both. Unlike above, Wang *et al.* [31] proposed an adversarial approach to reduce gender bias in a visual recognition model. While, most existing methods take into account a binary protected attribute despite the diversity of demographic groups. In contrast, we introduce a fair method that eliminates multiple protected attributes related biases in computer vision models.

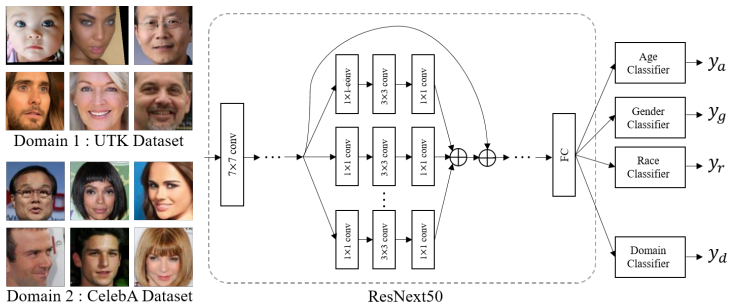


Figure 2: The proposed Protected Attribute Classifier (PAC).

2.2 Image-to-Image Translation

The main goal of Image-to-Image translation task is to learn how to map images from a source domain into images of a target domain. The methods based on Conditional Generative Adversarial Networks (CGANs) [15, 27] have shown a great success with pixel-wise paired datasets in super-resolution [29], image in-painting [36], image restoration [39], and image segmentation [19]. In addition, Cycle consistency adversarial networks (CycleGANs) [41] are introduced to learn a mapping between unpaired datasets. They train the Image-to-Image translation models in an unsupervised manner. Moreover, Choi *et al.* [9] proposed StarGAN that reduces the computational cost of models based on CycleGAN. The unified and unsupervised Image-to-Image translation model learns a mapping between multiple domains effectively. However, we find out that the learned mapping is biased to protected attributes (See Figure 1). There are some studies [14, 28, 32] that prevent unwanted information translation during mapping. Although Siddiquee *et al.* [28] proposed a FixedPointGAN that generates unchanged images in same-domain translation, it generates biased results in different-domain translation, a still remaining issue. In addition, fair representation methods by semantic constraints [32] and a disentangling method [14] are developed. Inspired by [14, 32], we also aim to train a fairness-aware image translation model by proposing a fair representation learning method.

3 Proposed Method

In this work, we propose two modules: 1) Protected Attribute Classifier (PAC) module, which learns high-level features of multiple protected attributes. 2) FairFaceGAN, which is a fairness-aware facial Image-to-Image translation network to learn a fair mapping of the multiple facial attributes in the multi-domain. The main network for the fairness-aware Image-to-Image translation is FairFaceGAN and PAC module is introduced to train FairFaceGAN without protected attribute annotations. In this section, we explain the modules in sequence.

3.1 Protected Attribute Classifier (PAC)

As illustrated in Figure 2, PAC consists of two branches: one is for predicting protected attributes (gender y_g , age y_a , race y_r) and the other is for predicting the domain labels y_d . The encoder of PAC with a number of convolutional layers is shared by the two branches and

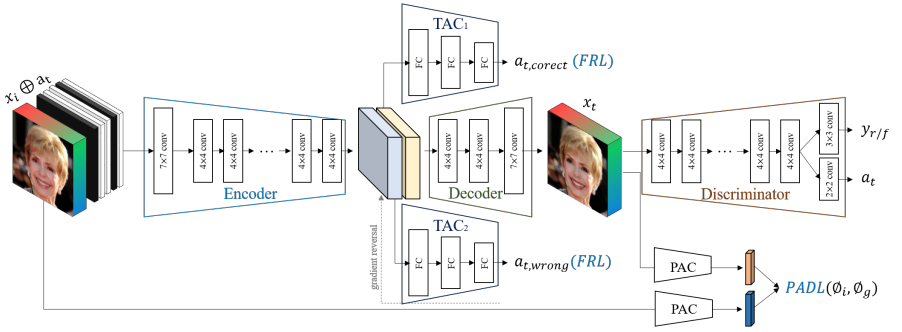


Figure 3: An overview of the proposed FairFaceGAN framework, which consists of Encoder-Decoder Generator, Discriminator, and Target Attribute Classifiers (TACs). Given an image x_i and target attribute a_t , we learn the model fairly work on protected attributes with Fair Representation Loss (FRL) and Protected Attribute Distance Loss (PADL) to generate x_t .

followed by task-specific fully connected layers: f_g (gender classifier), f_a (age classifier), f_r (race classifier), and f_d (domain classifier). We define the objective function for PAC as follows:

$$\mathcal{L}_{PA} = \mathcal{L}_{ce}(y_g|f_g(h)) + \mathcal{L}_{ce}(y_a|f_a(h)) + \mathcal{L}_{ce}(y_r|f_r(h)), \quad (1)$$

where \mathcal{L}_{ce} and h respectively denote a cross-entropy loss and a flattened feature of the last layer from the shared encoder.

In addition, to transfer knowledge related to protected attributes from the learned PAC into the FairFaceGAN, we train a discriminator to fail classification on source domain (UTK dataset [67]) and target domain (CelebA dataset [20]) using a gradient reversal layer like DANN [8] since the representation of PAC and FairFaceGAN are trained on different domains. To do so, we optimize the domain adversarial loss as follows:

$$\mathcal{L}_{PAC} = \mathcal{L}_{PA} - \lambda \mathcal{L}_{ce}(y_d|f_{c_d}(f(x))). \quad (2)$$

Optimization We use Adam optimizer with a learning rate of 0.001 and a batch size of 128. The PAC was optimized before ten epochs on a single 1080Ti GPU.

3.2 FairFaceGAN

FairFaceGAN aims to map input images into target facial attributes using a unified generator. As shown in Figure 3, it contains four components: one encoder-decoder generator, two target attribute classifiers (TACs), and one discriminator.

Given an input image x_i and a target attribute vector a_t , we first depth-wisely concatenate both of them. Then the combined data is fed into the encoder to represent two latent spaces. One is for target attributes and the other is for the rest information. The two features are then concatenated and used as an input of our decoder for generating a fair image x_t with target attributes.

Auxiliary Classifier Generative Adversarial Network Loss. We train FairFaceGAN with an adversarial loss to generate images to be realistic. In addition, we add an auxiliary classification layer on the top of the discriminator to distinguish the target attributes of the input image (a_i) and the generated image (a_t). The adversarial loss with the auxiliary classifier is defined as follows:

$$\min_{\theta_G} \max_{\theta_D} \mathcal{L}_{acgan} = \mathbb{E}_{x_i} [\log D(x_i)] + \mathbb{E}_{x_t} [\log(1 - D(x_t))] - \mathbb{E}_{x_i, a_i} [\log p_{\theta_D}(a_i | x_i)] - \mathbb{E}_{x_t, a_t} [\log p_{\theta_D}(a_t | x_t)]. \quad (3)$$

Reconstruction Loss. For the reconstruction, we use a cycle consistency loss [41] that guarantees the quality of generated images in the unsupervised manner. In addition, inspired by FixedPointGAN [23], we add an identity loss to make the generative model not transfer unnecessary regions in a same-domain translation.

$$\mathcal{L}_{rec} = \mathbb{E}_{x, a} [\|G(\tilde{x}_t, a_i) - x_i\|_1] + \mathbb{E}_{x, a} [\|G(x_i, a_i) - x_i\|_1]. \quad (4)$$

Fair Representation Loss (FRL). During translating target attributes, the high correlation between target attributes and protected attributes causes unwanted protected attribute translation. To prevent it, we separate representation h into target attribute translation (h_{tr}) and protected attribute preservation (h_{tu}) respectively. To this end, we apply a fair representation loss defined as follows:

$$\min_{\theta_{TAC_1}, \theta_{ENC}} \max_{\theta_{TAC_2}} \mathcal{L}_{fp} = \mathbb{E}_{x_i} [-\log p_{\theta_{TAC_1}}(a_t | h_{tr}) + \log p_{\theta_{TAC_2}}(a_t | h_{tu})]. \quad (5)$$

Protected Attribute Distance Loss (PADL). In addition, we propose protected attribute distance loss (PADL) minimizes the protected attribute feature distance between input images (ϕ_i) and generated images (ϕ_g). Since we do not have protected attribute labels in the target dataset, we instead utilize the semantic knowledge of protected attribute from the trained PAC to measure the distance. With Fair Representation Loss (FRL), it explicitly preserves protected attribute information in target attribute translation. The loss is defined as follows:

$$\mathcal{L}_{pad} = \mathbb{E}_x [\|\phi_i - \phi_g\|_1]. \quad (6)$$

Perceptual Loss. On top of that, the perceptual loss [16] is used to improve the quality of outputs. We select the same layers of [16] to measure not only the style loss between input images and reconstructed images but also the content loss between input images and generated images.

Optimization We use WGAN with gradient penalty [14] and Adam for optimizing the parameters of our FairFaceGAN with $\beta_1 = 0.5$ and $\beta_2 = 0.999$. We note that the overall loss function is a weighted sum of all terms. The initial learning rate for both generator and discriminator is set to 0.0001, which is decayed every eight epochs. We obtained the best results before 20 epochs on two 1080-TI GPUs.

Table 1: Protected attribute classification accuracy on UTK dataset [57] (Source Only and DA). DA denotes results of the domain adversarial training.

Attribute [Label]	Source Only	DA	CelebA [20]
Gender [Male, Female]	0.94	0.91	0.92
Race [White, Black, Asian, Indian, Others]	0.87	0.81	N/A
Age [0~9, 10~19, ... , 50+]	0.73	0.65	N/A
Domain Classification [UTK, CelebA]	N/A	0.5	N/A

Table 2: Quantitative comparison on CelebA dataset. f, p, and P indicate the usage of FRL, PADL, and Perceptual Loss. ACC, FID, and KID denote the average of target attribute classification accuracies, Fréchet Inception Distance [13], and Kernel Inception Distance ($\times 100$) [4].

	Star GAN [4]	FixedPoint GAN [23]	Ours (f)	Ours (p)	Ours (f+p)	Ours (f+p+P)
ACC \uparrow	92.07	91.01	90.55	92.11	89.71	90.66
FID \downarrow	10.23	6.91	10.66	6.98	9.98	9.8
KID \downarrow	1.94 \pm 0.29	2.06 \pm 0.41	2.33 \pm 0.28	1.47\pm0.35	2.13 \pm 0.3	1.89 \pm 0.27

4 Experiments

4.1 Dataset

PAC. We train PAC on UTK Face [57] and CelebA [20] datasets. CelebA dataset is utilized only for domain adversarial training and UTK Face dataset is leveraged for protected attribute (gender, race, and age) classification training as well as domain adversarial training. We randomly select 19,708, 2,000, and 2,000 images of UTK dataset for training, validation, and test, respectively, where 200,599 images of CelebA dataset are set to the domain adversarial training. All images are resized to 128×128 . Results with ranges of age and race for the classification are shown in Table 1.

FairFaceGAN. For training FairFaceGAN, we use only CelebA dataset without protected attribute annotation. Instead, by transferring knowledge from pre-trained PAC on UTK dataset, we utilize the protected attribute related semantic information. Training and test datasets are composed of 200,599 and 2,000 respectively. We pre-process all images by randomly cropping (178×178) and resizing into 128×128 . The five target attributes (*attractive, blond hair, bags under eyes, bald, big nose*) are selected manually. While we conduct both qualitative and quantitative evaluation for the *gender* attribute, we only conduct qualitative evaluation for the age and race attributes since their labels are not included in CelebA dataset.

Table 3: User study results.

	StarGAN [4]	FixedPointGAN [23]	Ours
Quality	30.78	20.97	48.25
Fairness	11.31	34.46	54.23

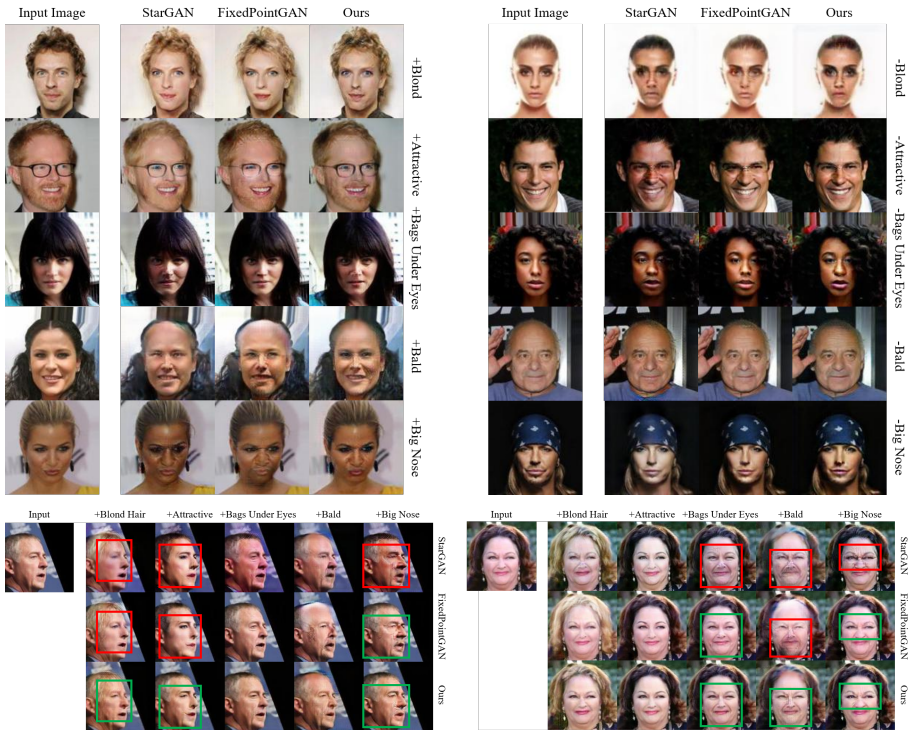


Figure 4: Image-to-Image translation results compare to StarGAN [4] and FixedPointGAN [28]. + and - denote the case of target attribute is added or removed. Red and Green boxes indicate the discriminative outputs and fairly mapped results respectively.

4.2 Evaluation

Qualitative evaluation. As shown in Figure 4, FairFaceGAN generates better quality images compared to StarGAN [4] and FixedPointGAN [28]. The models tend to change the skin color, add mustache on female images, apply makeup on male images, or make them aged, even though those are not the target attributes. Unlike their results, FairFaceGAN prevents the unwanted translation of protected attributes better.

Protected Attribute Classification. Table 1 shows the protected attribute classification accuracy of PAC on UTK and CelebA datasets. We fine-tune the ImageNet [25] pre-trained ResNext50 [52], one of the state-of-the-art image classification networks. The result demonstrates that our PAC encodes representations informative to the protected attributes on both UTK and CelebA datasets.

Quantitative Comparisons. To compare quantitative results of generated images of ours and existing models, we measure the target attribute classification accuracy, Fréchet Inception Distance (FID) [13], and Kernel Inception Distance (KID) [2]. In this experiment, we also conduct an ablation study of the proposed loss functions as follows: 1) Fair Representation Loss (FRL) only. 2) FRL and Protected Attribute Distance Loss (PADL). 3) FRL, PADL, and VGG Perceptual Loss. Firstly, to evaluate target attribute classification accuracies on the

Table 4: Fréchet Protected Attribute Distance (FPAD) of generated images. BUE denotes Bags Under Eyes. ($- \rightarrow +$) denotes without attribute into with attribute, and vice versa.

Gender (transform)	Attribute	StarGAN [4]	FixedPoint GAN [28]	Ours (f)	Ours (p)	Ours (f+p)	Ours (f+p+P)
Male ($- \rightarrow +$)	BlondHair	56.32	24.55	31.05	32.54	4.86	5.63
	Bald	11.68	11.90	6.67	14.24	5.19	8.30
	BUE	6.38	2.60	2.18	3.41	1.41	3.44
	BigNose	16.20	7.05	4.62	9.99	1.51	4.94
	Attractive	11.32	3.49	4.84	3.39	2.94	3.79
Male ($+ \rightarrow -$)	BlondHair	41.37	21.04	20.11	32.01	9.96	8.91
	Bald	17.79	3.71	3.97	8.51	2.19	9.02
	BUE	21.29	6.87	9.23	13.32	3.13	5.23
	BigNose	2.66	2.02	2.19	3.75	1.11	1.63
	Attractive	7.85	4.43	4.09	13.92	1.35	6.7
Female ($- \rightarrow +$)	BlondHair	135.7	108.71	72.98	104.13	4.75	17.39
	Bald	60.33	131.48	22.18	57.83	21.00	24.79
	BUE	3.25	3.10	1.71	3.08	1.55	4.02
	BigNose	22.42	12.18	4.98	8.97	2.22	3.47
	Attractive	13.85	7.29	6.17	3.05	2.78	5.00
Female ($+ \rightarrow -$)	BlondHair	29.80	94.38	35.49	55.39	5.17	5.94
	BUE	6.06	4.19	9.57	4.42	2.29	3.74
	BigNose	5.77	3.10	4.95	4.5	2.18	3.86
	Attractive	22.79	13.36	17.30	19.2	7.12	11.70
Average		25.94	24.50	13.91	20.82	4.35	7.24

Table 5: Fair Classification Results. TPR, FPR, $Eq.Opp.$, and $Odds$ indicate Classification Accuracy, True Positive Rates, False Positive Rate, Equality of Opportunity [4], and Equalized Odds [35]. O and G indicate the subset of original images in test dataset for the image translation model and the generator. Last three rows present results of data augmentation.

Training Dataset	Male		Female		Fairness Score	
	TPR	FPR	TPR	FPR	$Eq.Opp.$	$Odds$
O	64.10	18.40	86.36	49.00	22.26	26.43
$G_{ours}(O)$	79.49	29.45	90.40	53.00	10.92	17.23
$O+G_{StarGAN}(O)$ [4]	64.10	15.34	91.41	43.00	27.31	27.49
$O+G_{FixedPointGAN}(O)$ [28]	56.41	19.63	87.88	42.00	31.47	26.92
$O+G_{ours}(O)$	74.36	22.70	85.35	45.00	10.99	16.65

generated images, we re-train the ImageNet [25] pre-trained ResNext50 [34] to classify the target attributes on CelebA dataset. As shown in Table 2 (first row), the generated images from ours achieve the best result (92.11%) over other models, where original testset achieves the accuracy of 88.88%. We also measure FID and KID values to evaluate our model with standard metrics. As shown in Table 2 (second and third rows), our model shows the best KID and competitive FID. Meanwhile, our final model shows slightly lower accuracy than others since there is a trade-off between fairness and the image generation ability [4, 26]. Note that our goal focuses on improving fairness of the translation model.

User Study. We also present results of a user study to compare the fairness and visual quality of generated images of ours, StarGAN [9], and FixedPointGAN [28]. We randomly select 24 sets, four images per set (Input, Results of StarGAN, FixedPointGAN, and ours), and request 73 participants to choose the best produced (Quality) and the best protected attribute preserved (Fairness) images. As shown in Table 3, our model achieves the best scores for both image quality and fairness.

Fréchet Protected Attribute Distance (FPAD). To evaluate the fairness of our proposed model, we propose a new metric FPAD inspired by FID [13]. We leverage our PAC model to extract a protected attribute feature and measure feature distance of input images X_i and generated images X_g . We compute $\|M_i - M_g\|^2 + \text{Tr}(C_i + C_g - 2(C_i C_g)^{1/2})$ in given (M_i, C_i) and (M_g, C_g) which are the mean and covariance of protected attribute features from X_i and X_g . As shown in Table 4, our model achieves the lowest FPAD compared to the prior models. In other words, our generative model best preserves the protected attributes during the mapping. Although there is a slight performance drop, we compensate it by applying the perceptual loss that improves visual quality of generated images.

Fair Classification. Furthermore, to evaluate our model using standard fairness metrics, we conduct an attractiveness classification task. We compare the performances when augmenting data by existing image translation models [9, 28] and FairFaceGAN respectively. For the evaluation, we leverage the two fairness metrics: Equality of Opportunity and Equalized Odds ($Eq.Opp. = |TPR_{male} - TPR_{female}|$, $Odds = \frac{1}{2}[|FPR_{male} - FPR_{female}| + |TPR_{male} - TPR_{female}|]$). Details are in our supplementary material. We fine-tune ImageNet pre-trained ResNext50 [52] using the testset of FairFaceGAN, divided into 1,200 (O), 300, and 500 images for training, validation, and test, respectively. As shown in Table 5, we verify whether generated images of FairFaceGAN can be utilized for the classification model to be trained more fairly on gender compare to existing image translation models.

5 Conclusion

In this paper, we introduced a novel fairness-aware facial Image-to-Image translation model to avoid the problem of translating unwanted attributes. Through Fair Representation Loss (FRL) and Protected Attribute Distance Loss (PADL), our model learns fair representations in terms of multiple protected attributes (age, gender, and race). To demonstrate the ability of FairFaceGAN, we conducted an extensive evaluation of image translation and fairness. Overall, our experimental results showed that FairFaceGAN is fairer in terms of Equality of Opportunity, Equalized Odds, and the proposed FPAD over the existing Image-to-Image translation models.

Acknowledgements. This work was supported by Institute for Information & communications Technology Planning & Evaluation (IITP) grant funded by the Korea government (MSIT) (Development of framework for analyzing, detecting, mitigating of bias in AI model and training data) under Grant 2019-0-01396 and (Artificial Intelligence Graduate School Program (YONSEI UNIVERSITY)) under Grant 2020-0-01361.

We thank **Pilhyeon Lee**, **Seogkyu Jeon**, and **Jijoong Kim** for the thorough reviews and the constructive feedback.

References

- [1] Lisa Anne Hendricks, Kaylee Burns, Kate Saenko, Trevor Darrell, and Anna Rohrbach. Women also snowboard: Overcoming bias in captioning models. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 771–787, 2018.
- [2] Mikołaj Bińkowski, Dougal J Sutherland, Michael Arbel, and Arthur Gretton. Demystifying mmd gans. In *International Conference on Learning Representations*, 2018.
- [3] Martim Brandao. Age and gender bias in pedestrian detection algorithms. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2019.
- [4] 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 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8789–8797, 2018.
- [5] Elliot Creager, David Madras, Jorn Jacobsen, Marissa Weis, Kevin Jordan Swersky, Toniann Pitassi, and Richard Zemel. Flexibly fair representation learning by disentanglement. In *Thirty-sixth International Conference on Machine Learn (ICML)*, 2019.
- [6] Terrance de Vries, Ishan Misra, Changhan Wang, and Laurens van der Maaten. Does object recognition work for everyone? In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2019.
- [7] Sanghamitra Dutta, Dennis Wei, Hazar Yueksel, Pin-Yu Chen, Sijia Liu, and Kush R Varshney. Is there a trade-off between fairness and accuracy? a perspective using mismatched hypothesis testing. In *ICML 2020*, July 2020.
- [8] Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In *Proceedings of the 32nd International Conference on International Conference on Machine Learning - Volume 37, ICML'15*, page 1180–1189. JMLR.org, 2015.
- [9] R Stuart Geiger, Kevin Yu, Yanlai Yang, Mindy Dai, Jie Qiu, Rebekah Tang, and Jenny Huang. Garbage in, garbage out? do machine learning application papers in social computing report where human-labeled training data comes from? In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pages 325–336, 2020.
- [10] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.
- [11] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved training of wasserstein gans. In *Advances in neural information processing systems*, pages 5767–5777, 2017.
- [12] Moritz Hardt, Eric Price, and Nati Srebro. Equality of opportunity in supervised learning. In *Advances in neural information processing systems*, pages 3315–3323, 2016.

- [13] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in neural information processing systems*, pages 6626–6637, 2017.
- [14] Sunhee Hwang and Hyeran Byun. Unsupervised image-to-image translation via fair representation of gender bias. In *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1953–1957, 2020.
- [15] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1125–1134, 2017.
- [16] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European conference on computer vision*, pages 694–711. Springer, 2016.
- [17] Surya Mattu Julia Angwin, Jeff Larson and Lauren Kirchner. Machine bias. URL <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>.
- [18] Byungju Kim, Hyunwoo Kim, Kyungsu Kim, Sungjin Kim, and Junmo Kim. Learning not to learn: Training deep neural networks with biased data. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9012–9020, 2019.
- [19] Issam H Laradji, David Vazquez, and Mark Schmidt. Where are the masks: Instance segmentation with image-level supervision. In *British Machine Vision Conference (BMVC)*, Cardiff, UK, 2019.
- [20] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.
- [21] Sachit Menon, Alexandru Damian, Shijia Hu, Nikhil Ravi, and Cynthia Rudin. Pulse: Self-supervised photo upsampling via latent space exploration of generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2437–2445, 2020.
- [22] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014.
- [23] Sunggho Park, Dohyung Kim, Sunhee Hwang, and Hyeran Byun. Readme: Representation learning by fairness-aware disentangling method. *arXiv preprint arXiv:2007.03775*, 2020.
- [24] Novi Quadrianto, Viktoriia Sharmanska, and Oliver Thomas. Discovering fair representations in the data domain. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.

- [25] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015.
- [26] Sivan Sabato and Elad Yom-Tov. Bounding the fairness and accuracy of classifiers from population statistics. In *ICML 2020*, July 2020.
- [27] 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*, 63(4/5):3–1, 2019.
- [28] Md Mahfuzur Rahman Siddiquee, Zongwei Zhou, Nima Tajbakhsh, Ruibin Feng, Michael B Gotway, Yoshua Bengio, and Jianming Liang. Learning fixed points in generative adversarial networks: From image-to-image translation to disease detection and localization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 191–200, 2019.
- [29] Vagia Tsiminaki, Wei Dong, Martin R. Oswald, and Marc Pollefeys. Joint multi-view texture super-resolution and intrinsic decomposition. In *British Machine Vision Conference (BMVC)*, Cardiff, UK, 2019.
- [30] James Vincent. What a machine learning tool that turns obama white can (and can't) tell us about ai bias, 2020. URL <https://www.theverge.com/21298762/face-depixelizer-ai-machine-learning-tool-pulse-stylegan-oba>
- [31] Tianlu Wang, Jieyu Zhao, Mark Yatskar, Kai-Wei Chang, and Vicente Ordonez. Balanced datasets are not enough: Estimating and mitigating gender bias in deep image representations. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019.
- [32] Yaxing Wang, Abel Gonzalez-Garcia, Joost van de Weijer, and Luis Herranz. Controlling biases and diversity in diverse image-to-image translation. *arXiv preprint arXiv:1907.09754*, 2019.
- [33] Benjamin Wilson, Judy Hoffman, and Jamie Morgenstern. Predictive inequity in object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2019.
- [34] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1492–1500, 2017.
- [35] Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez Rodriguez, and Krishna P Gumadi. Fairness beyond disparate treatment & disparate impact: Learning classification without disparate mistreatment. In *Proceedings of the 26th international conference on world wide web*, pages 1171–1180, 2017.
- [36] Ruonan Zhang, Yurui Ren, Jingfei Qiu, and Ge Li. Base-detail image inpainting. In *British Machine Vision Conference (BMVC)*, Cardiff, UK, 2019.

- [37] Zhifei Zhang, Yang Song, and Hairong Qi. Age progression/regression by conditional adversarial autoencoder. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5810–5818, 2017.
- [38] 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*, pages 2941–2951, 2017. URL <https://www.aclweb.org/anthology/D17-1319>.
- [39] Yupei Zheng, Xin Yu, Miaomiao Liu, and Shunli Zhang. Residual multiscale based single image deraining. In *British Machine Vision Conference (BMVC)*, Cardiff, UK, 2019.
- [40] Zhilin Zheng and Li Sun. Disentangling latent space for vae by label relevant/irrelevant dimensions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 12192–12201, 2019.
- [41] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232, 2017.