

# Gender-Adversarial Networks for Face Privacy Preserving

Deyan Tang, *Member, IEEE*, Siwang Zhou<sup>✉</sup>, *Member, IEEE*, Hongbo Jiang<sup>✉</sup>, *Senior Member, IEEE*,  
Haowen Chen<sup>✉</sup>, and Yonghe Liu<sup>✉</sup>

**Abstract**—Privacy concerns over face recognition systems have attracted extensive attention in various fields. For gender privacy-preserving work, there are two key challenges: 1) *privacy*, i.e., confusing gender classifiers and 2) *utility*, i.e., maintaining its face verification performance. To address both issues, this article develops a novel gender-adversarial network, referred to as Gender-AN, to impart gender privacy to face images. Gender-AN employs an attribute-independent encoder-decoder GAN-based network to perturb the input face image, training with the assistance of the proper facial attributes. The perturbed image is then able to obfuscate gender classifiers while maintaining identity discriminability. To optimize the generator, a multitask-based loss function is utilized, which includes attribute manipulation loss, face matcher loss, adversarial loss, and reconstruction loss functions. This optimization facilitates our model to achieve the generalization, verification preserve, and natural appearance, simultaneously. Extensive experiments confirm the effectiveness of the proposed model in enhancing gender privacy and preserving face verification utility.

**Index Terms**—Auxiliary attribute, encoder-decoder gender-adversarial network (GAN), face privacy preserving, GAN, gender privacy.

## I. INTRODUCTION

PRIVACY concerns over face recognition systems have been voiced for various areas, such as social media [1], smart phones [2], smart camera networks [3], mobile applications [4], IoT [5], and related application domains [6], [7]. Therefore, face privacy protection has become increasingly important and urgent, and lots of research on it [5], [8]–[12] have emerged in the last decades. On the one hand, some research focus on face deidentification technologies [11], [13], [14], in which the facial image is modified

to conceal the identity information. On the other hand, some studies [15]–[17] aim to obfuscate or maintain a set of attributes that can be inferred from face images, such as obfuscating gender attribute classifiers and preserving face verification ability [18].

An effective privacy preserving algorithm should take into account both *privacy* and *utility* [19], where privacy means minimizing information leakage concerning sensitive data, and utility aims to preserve useful data as much as possible. Extending to gender privacy preserving techniques for face data, privacy involves two aspects: 1) the generalizability on unseen gender classification algorithms and 2) the ability to confuse gender classifiers, not just simply flipping. Simultaneously, utility involves three factors: 1) maintaining the performance of face verification; 2) avoiding affecting the recognition of other attributes as much as possible, and 3) retaining the visibility and quality of the perturbed facial image.

The initial research mainly focused on reducing the recognition rate of specified gender classifiers [15], [20], which led to a failure of generalization on arbitrary classifiers. The generalization is further exploited in [8] and [21], yet the effect is still unsatisfactory. Moreover, the original intention of these privacy preserving algorithms was to reverse gender. That is to say, various schemes were introduced to anonymize gender information in face images, so that their gender classification accuracy is greatly reduced, or even flipped [16], [18]. Therefore, a better method is needed to solve the two issues: 1) generalization and 2) confusion abilities, related to privacy in gender privacy preserving.

Meanwhile, researchers are also committed to solving the utility problem while dealing with the privacy issue. Face morphing technologies were provided in [22] and [23] to suppress gender attribute successfully. However, they diminish the visual quality of perturbed images and face matching utility. The algorithms proposed in [24]–[26] tried to remove sensitive information, such as gender, with little loss in face verification. Boutet *et al.* [17] attempted to sanitize gender attribute and maintain activity recognition. Nevertheless, instead of processing images, these models are mainly designed to deal with face representations or motion sensor data, which would somewhat limit their application scenarios and weaken the utility. The PrivacyNet proposed in [18] provided a generative adversarial network (GAN)-based semiadversarial network (SAN) to confound multiple sensitive attributes in facial images while facilitating verification. However, the performance of other attribute classifiers for adversarial images would be affected.

Manuscript received 23 January 2022; accepted 26 February 2022. Date of publication 3 March 2022; date of current version 7 September 2022. This work was supported in part by the National Natural Science Foundation of China under Grant 62172153, and in part by the Changsha Municipal Natural Science Foundation under Grant kq2014057. (*Corresponding author: Siwang Zhou.*)

Deyan Tang, Hongbo Jiang, and Haowen Chen are with the College of Computer Science and Electrical Engineering, Hunan University, Changsha 410082, China (e-mail: deyantang@hnu.edu.cn; hongbojiang@hnu.edu.cn; hwchen@hnu.edu.cn).

Siwang Zhou is with the College of Computer Science and Electrical Engineering, Hunan University, Changsha 410082, China, and also with the Key Laboratory of Intelligent Computing & Information Processing of Ministry of Education, Xiangtan University, Xiangtan 411105, China (e-mail: swzhou@hnu.edu.cn).

Yonghe Liu is with the Department of Computer Science and Engineering, the University of Texas at Arlington, Arlington, TX 76019 USA (e-mail: yonghe@cse.uta.edu).

Digital Object Identifier 10.1109/IJOT.2022.3155878

Also, a human can easily distinguish the original image from the perturbed one generated by PrivacyNet. Thus, all of these would weaken its last two factors of utility, i.e., influence on other attributes and the appearance quality.

Recently, by adapting the specific facial attributes according to practical requirements, Li *et al.* [3] delivered to conceal the visual appearance while maintaining the identity discriminability. That is, with the assistance of proper facial attributes, the appearance of a face would be anonymized, while retaining the recognition ability. Inspired by this research, we associate that gender anonymity can also be achieved with appropriate facial attributes assistance, without affecting the face matching utility. Consequently, we design a new GAN-based gender-adversarial network, named Gender-AN, for face gender privacy preserving. First, to better confuse gender classifiers, some auxiliary facial attributes are utilized to train the model. Second, to achieve the generalization, verification preserve, and natural appearance, we optimize the face generator with a multitask-based loss function, including attribute manipulation loss, face matcher loss, adversarial loss, and reconstruction loss functions. Third, in order to bring less influence to other attribute classifiers, the generator is designed to be attribute independent. On account of the above three efforts, our proposed Gender-AN model can address the two issues involved in privacy and three issues in utility, simultaneously. Incidentally, benefiting from the multitask optimization work, our model can also protect multiple sensitive attributes privacy to a certain extent, such as gender and age.

The remainder of the article is organized as follows. Section II introduces the related work briefly. Section III details the proposed scheme. Section IV presents the experimental results and discussion. Section V gives the conclusion.

## II. RELATED WORK

The latest advances in machine learning allow for the automatic extraction of sensitive face information, such as gender, race, and age, called soft biometric, from nonchemical data, such as audio, text, and images [27]. Some authoritative organizations even adopted bills to integrate privacy protection methods into technology [25], [27]. Soft biometric privacy is imperative, but it is still a relatively new field of research [18]. In this section, we review the representative soft biometric attribute (such as gender) privacy enhancement work from privacy and utility perspectives.

As mentioned in Section I, both generalization and gender confusion should be taken into consideration for *privacy*. Recently, with the introduction of the notion of adversarial examples [28], researchers began to avail various adversarial machine learning strategies (i.e., adversarial examples, adversarial perturbations, adversarial noise, etc.) to achieve gender privacy protection. Chhabra *et al.* [20] and Mirjalili *et al.* [8], [15], [21] tried to introduce adversarial perturbations into face images to confuse gender classifiers. Chhabra *et al.* [20] introduced an adversarial framework inspired from the Carlini–Wagner L2 attack [29], which adds adversarial perturbations to the input face images, aiming to suppress some predefined attributes classifiers and retain

others. Later, Mirjalili *et al.* [15] developed a SAN to generate perturbed face images, aiming to confound gender information while preserving identity information. In their follow up work [8], the same authors introduced an ensemble of SANs to train the model, and in [21], they tried to combine a variety of SAN transformations to improve the generalization ability of SAN models. Although the models discussed above successfully invalidate the target attribute classifiers, they fail to generalize to arbitrary classifiers. Thus the generalization of *privacy* is diminished.

The PrivacyNet presented in [18] makes up for the limitation of generalizability faultiness. PrivacyNet utilizes a GAN-based SAN to selectively perturb gender, age, and race attributes in face images, while maintaining face matching performance. Nevertheless, the gender classifiers flip after disturbing gender information, weakening the gender confusion performance in privacy. Moreover, the modified images from the PrivacyNet model are not so realistic looking that they would be distinguished from nonmodified ones by a human observer. It would limit the utility discussed next.

For utility, there are three issues worthy of attention: 1) verification preserve; 2) attribute independent; and 3) natural appearance. For some algorithms processing face images directly, Othman and Ross [22] proposed a method to modify the gender attributes of face images while retaining their face matching ability. Wang *et al.* [23] also applied face morphing technologies to synthesize new face images, where the gender information is obfuscated. The main limitation of both methods is major transformation and loss in visual appearance, which reduces the utility. Besides, for some other algorithms processing face representations, Morales *et al.* [25] provided the dedicated deep models, named SensitiveNets, to remove gender and ethnicity information from face representations. Afterward, any gender and ethnicity classifiers will not work anymore on these agnostic representations, while the face matching performance is not affected. Terh rst *et al.* [24] applied an incremental variable elimination (IVE) technology to remove the most important components gradually for the prediction of the target attribute (i.e., age or gender) from a given face representation. The authors demonstrated that a lot of information on sensitive attributes can be discarded, while still maintaining representations recognition ability. Bortolato *et al.* [26] designed a PFRNet to learn a disentangled feature representation, in which the attribute-related information is parted from identity information. The attribute-related part is discarded to remove the sensitive information, while the identity-related part is retained for verification purposes. However, these schemes lack human interpretability since only storing face representations, which limits the utility in many applications.

Therefore, a more suitable generator is designed in the proposed Gender-AN model, to solve both issues in privacy and three problems in utility. Furthermore, Gender-AN refers to the idea of [3]. Li *et al.* [3] offered a face anonymization framework to disguise visual appearance and retain the recognition utility, simultaneously. The model can adaptively discover the specified facial attributes according to practical demands, and then generate the privacy-preserving face conditioned on these attributes with a conditional GAN. It means

that, the appearance of a face will be anonymized with proper facial attributes assistance, while its identity discriminability is preserved. Although, the main purpose of this research is to conceal visual appearance while maintaining the identity information, aided by some appropriate facial attributes, we are still able to apply this idea to gender privacy protection. That is, with the assistance of some specific face attributes, to hide gender information and retain the face matching utility.

### III. PROPOSED METHOD

#### A. Problem Formulation

Given a source face image  $x^a$  with  $n$  binary attributes. Define  $f_G(x^a)$  as the gender classifier, where “Male” is denoted by “1” and “Female” is denoted by “0.” Let  $f_i(x^a)$  represent the remaining  $n - 1$  attribute classifiers, here  $i = 1, \dots, n - 1$ . Define  $f_M(x_1^a, x_2^a)$  as a face matcher to determine whether  $x_1^a$  and  $x_2^a$  come from the same object. The ultimate goal of our paper is to find a model  $\phi$ , so that the generated disturbed image  $x^b = \phi(x^a)$  owns the following characteristics.

- 1) For the gender attribute, any unseen gender classifier  $f_G(x^b)$  can be substantially confused by the modified image  $x^b$ .
- 2) The recognition rate of the gender classifier for  $x^b$  reaches a random probability value, i.e., 0.5.
- 3) For the face matching utility, the face recognition ability must be preserved for the perturbed image  $x^b$ . It means that  $f_M(x_1^b, x_2^b) \approx f_M(x_1^a, x_2^a)$ .
- 4) For the remaining  $n - 1$  attributes, the performance of these attribute classifiers  $f_i(x^b)$  has not to be noticeably impacted either positively or negatively. That is,  $f_i(x^b) \approx f_i(x^a)$ .
- 5) From a human perspective, the modified image  $x^b$  must approximate the source image  $x^a$ . Meanwhile, the modified image should keep its realism.

The above issues can be summarized into two points: 1) and 2) aim to solve the two problems of privacy; and 3)–5) aspire to address the three issues of utility. It will also be reflected in the model design introduced as follows.

#### B. Proposed Gender-Adversarial Network

To resolve the problems of privacy and utility simultaneously, and achieve the objectives described in Section III-A, we devise a GAN-based gender-adversarial network (Gender-AN) for face privacy preserving. First, to address the generalization issue in privacy, the generator is optimized with classification loss (attribute loss). Second, to achieve the goal of random classification probability in privacy, three specific attributes, which are Attractiveness, Heavy\_Makeup, and Age, are utilized to assist in training the model. Third, a face matcher is added to the model, such that the modified image can preserve its face matching ability. It resolves the first issue of utility. Fourth, the second issue in utility is settled by a selective transmission unit (STU) in the generator of the proposed model. With STU, when perturbing gender, the performance of other attributes may be affected as little as possible. Fifth, we employ an encoder-decoder GAN-based network to modify face images directly, addressing the visibility in utility.

TABLE I  
NETWORK ARCHITECTURE OF GENERATOR

| Encoder( $G_{enc}$ )          | Decoder( $G_{dec}$ )      |
|-------------------------------|---------------------------|
| Conv(16,4,2), BN, Leaky ReLU  | DeConv(256,4,2), BN, ReLU |
| Conv(32,4,2), BN, Leaky ReLU  | DeConv(128,4,2), BN, ReLU |
| Conv(64,4,2), BN, Leaky ReLU  | DeConv(64,4,2), BN, ReLU  |
| Conv(128,4,2), BN, Leaky ReLU | DeConv(32,4,2), BN, ReLU  |
| Conv(256,4,2), BN, Leaky ReLU | DeConv(3,4,2), BN, Tanh   |

The architecture of Gender-AN is comprised of three sub-networks: 1) a generator ( $G$ ) to modify the source image; 2) a discriminator ( $D$ ), to distinguish an image is original or synthetic and predict facial attributes simultaneously; and 3) a face matcher (FM) to recognize face images. In addition to the input image, the attribute labels are used in the generator and discriminator as condition variables, which are expanded into a matrix of the same width and height as the input image ( $224 \times 224$ ). Together, these subnetworks form an encoder-decoder GAN [30], and the overview of the proposed Gender-AN model is illustrated in Fig. 1.

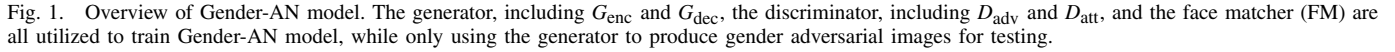
#### C. Neural Network Architecture

This section will introduce the details of the Gender-AN architecture.

1) *Generator*: The generator  $G$ , as shown in Fig. 1, is composed of an encoder  $G_{enc}$  for abstract latent representation and a decoder  $G_{dec}$  for generating a target image. As illustrated in Table I, there are five convolution layers with kernel size 4 and stride 2 in the encoder  $G_{enc}$ , and five transposed convolution layers in the decoder  $G_{dec}$ . Besides, the STU is applied right after each of the first four layers of the encoder to selectively transform the encoder features to make them compatible and complementary with the decoder features, as described in [30].

Given a source face image  $x^a$  with binary attributes  $\mathbf{att}_s = \mathbf{a}$ . The goal of  $G$  is to modify the input image  $x^a$  so that its attributes vector changes to  $\mathbf{att}_t = \mathbf{b}$ . Define the vector  $\mathbf{att}_{diff} = \mathbf{att}_t - \mathbf{att}_s$  as the difference vector between target and source attribute vectors, guiding  $G_{enc}$  which attributes need to be transformed, toward what direction an attribute should be changed. To achieve this goal,  $G_{enc}$  encodes this input image  $x^a$  into a latent representation  $z$ . Then, guided by  $\mathbf{att}_{diff}$ , STUs are deployed to transform encoder features as  $f_i = \{f_i^1, \dots, f_i^4\}$ , that concatenates with the results of transposed convolution layers for  $G_{dec}$  to recover targeted image  $x^b$ . Thus, the output image can be given by  $x^b = G_{dec}(z, f_i)$ , which can be rewritten as  $x^b = G(x^a, \mathbf{att}_{diff})$ .

2) *Discriminator*: The discriminator  $D$ , as shown in Fig. 1, contains two subnetworks: 1) a real/fake adversarial discriminator  $D_{adv}$  and 2) a facial attribute classifier  $D_{att}$ . The backbone of  $D$  is combined by five convolutional layers with kernel size 4 and stride 2, and the backbone is exploited for the feature extraction of an input image. Here, instance normalization (IN) along with the Leaky ReLU function is applied in each convolutional layer, as illustrated in Table II. Then, the CNN backbone is followed by two independent fully connected layers of  $D_{att}$  and  $D_{adv}$ , severally, such as the 6th and 7th layers shown in Table II. When  $G$  attempts to generate a perturbed



| Layer | Discriminator( $D_{adv}$ )    | Attribute Classifier( $D_{att}$ ) |
|-------|-------------------------------|-----------------------------------|
| 1     | Conv(16,4,2), IN, Leaky ReLU  |                                   |
| 2     | Conv(32,4,2), IN, Leaky ReLU  |                                   |
| 3     | Conv(64,4,2), IN, Leaky ReLU  |                                   |
| 4     | Conv(128,4,2), IN, Leaky ReLU |                                   |
| 5     | Conv(256,4,2), IN, Leaky ReLU |                                   |
| 6     | FC(1024), Leaky ReLU          | FC(1024), Leaky ReLU              |
| 7     | FC(1)                         | FC(4), Sigmoid                    |

3) *Face Matcher*: Finally, the auxiliary face matcher is adjusted according to the publicly available pretrained SE-ResNet-50 model [31]. The model receives input face images with a size of  $224 \times 224 \times 3$  and calculates the corresponding face representations with a size of 2048.

In this section, the adversarial, attribute manipulation, reconstruction, and face matcher losses are detailed. All of them are collaborated to train our model.

3) *Reconstruction Loss*: As described in Section III-A, the synthesized image  $\hat{x}^b$  should approximate the input image  $x^a$

in a visual perspective. Reconstruction cannot be guaranteed if only adversarial and classification losses are applied. Thus,  $G_{\text{dec}}$  also needs to learn to reconstruct an image from the latent representations conditioned on source attribute label  $\mathbf{a}$ , i.e.,  $\mathbf{att}_{\text{diff}} = \mathbf{0}$ . Then, the reconstruction loss can be formulated as

$$L_{G,\text{rec}} = \mathbb{E}_x[\|x^{\mathbf{a}} - G(x^{\mathbf{a}}, \mathbf{0})\|_1]. \quad (5)$$

4) *Face Matcher Loss*: As mentioned in Section III-A, the face recognition ability should be preserved for the perturbed image  $x^{\mathbf{b}}$ . For this purpose, the loss of face matcher is introduced to optimize the performance of face matcher FM on  $x^{\mathbf{b}}$ . The face matcher loss is defined as

$$L_{G,\text{fm}} = \mathbb{E}_{x^{\mathbf{a}}, \mathbf{att}_{\text{diff}}}[\|\text{FM}(x^{\mathbf{a}}) - \text{FM}(G(x^{\mathbf{a}}, \mathbf{att}_{\text{diff}}))\|_2^2] \quad (6)$$

where  $\text{FM}(x)$  denotes the normalized face representation of face image  $x$ , and  $\text{FM}(x)$  can be obtained as described in Section III-C-3.

5) *Overall Objective*: Considering all the above loss functions, the objective for discriminator  $D$  can be formulated as

$$L_D = L_{D,\text{adv}} + \lambda_1 L_{D,\text{att}} \quad (7)$$

and that for generator  $G$  is

$$L_G = L_{G,\text{adv}} + \lambda_2 L_{G,\text{att}} + \lambda_3 L_{G,\text{rec}} + \lambda_4 L_{G,\text{fm}} \quad (8)$$

where  $\lambda_1, \lambda_2$ , and  $\lambda_3$  are the hyperparameters that represent the relative weights for the corresponding loss terms and balance the losses.

#### IV. PERFORMANCE EVALUATION

##### A. Details and Database

We evaluate the effectiveness of our proposed algorithm mainly on the image database CelebA, which is widely used by most relevant works [18], [33]. CelebA contains 202 599 images, and each image is annotated with 40 binary attributes (with/without). Approximately, 200k images are used as the training set, the last 1.8k images are used for the testing set, and the remaining images are used for the validation set. In our experiment, each image is resized to  $224 \times 224$  by bicubic interpolation. We select four attributes, Male (Gender), Attractive (Attractiveness), Heavy\_Makeup, and Young (Age), to train our model. Because Heavy\_Makeup has close correlations with gender, meanwhile, Attractiveness and Age are the most related attribute to Heavy\_Makeup.

Due to the memory-intensive training process, the batch size is set to 16. Our proposed Gender-AN model was trained for 1 000 000 iterations. The ADAM [30] optimizer is used to train our model, with  $\beta_1 = 0.5$  and  $\beta_2 = 0.999$ . Set the initial learning rate  $2 \times 10^{-4}$ , and decay it to  $2 \times 10^{-5}$  after 800 000 iterations. The optimal hyperparameter in (7) and (8) is set to  $\lambda_1 = 1$ ,  $\lambda_2 = 10$ , and  $\lambda_3 = \lambda_4 = 100$ . After training the model, as shown in Fig. 1, for the *Test* part, only the generator network  $G$  is utilized to generate adversarial samples that confuse the gender classifiers while retaining the face matcher performance.

Some other designs also employed adversarial networks to impart gender privacy to face images. PrivacyNet [18] is a

TABLE III  
ACCURACY FOR GENDER CLASSIFIERS BEFORE AND AFTER  
PERTURBING GENDER ATTRIBUTE

| Adversarial Model | Gender Classifier | $R_{\text{or}}$ | $R_{\text{adv}}$ | Reduction    |
|-------------------|-------------------|-----------------|------------------|--------------|
| PrivacyNet Model  | SVM               | 0.986           | 0.239            | 1.537        |
|                   | RF                | 0.973           | 0.287            | 1.450        |
|                   | MLP               | 0.982           | 0.258            | 1.502        |
|                   | KNN               | 0.984           | 0.347            | 1.312        |
| Gender-AN Model   | SVM               | 0.986           | <b>0.508</b>     | <b>0.984</b> |
|                   | RF                | 0.973           | <b>0.544</b>     | 0.865        |
|                   | MLP               | 0.982           | <b>0.469</b>     | <b>1.064</b> |
|                   | KNN               | 0.984           | <b>0.584</b>     | 0.826        |

typical design that utilizes GAN-based generator to obfuscate the gender information. Thus, this article will use PrivacyNet as a comparison object.

As illustrated in Section III-A, we assess the Gender-AN model from the following five perspectives: 1) generalize to unknown gender classifiers; 2) confuse gender classifiers; 3) retain face matcher ability; 4) preserve the performance of other facial attributes; and 5) maintain the natural appearance and image quality. Several experiments are conducted to evaluate whether our proposed Gender-AN model can achieve these goals.

##### B. Evaluation on Gender Attribute

To assess the generalizability of our Gender-AN model, we considered four gender classification algorithms in the experiments: 1) support vector machines (SVM); 2) random forests (RF); 3) multilayer perception (MLP); and 4)  $K$ -nearest neighbors (KNN). All these algorithms were trained with 20 000 original images from CelebA database. The publicly available pretrained SE-ResNet-50 [31] was used here to extract embeddings before training or testing the gender classifiers.

To test the confusion ability of our model, we referenced to Reduction from [25], which is defined as

$$\text{Reduction} = (R_{\text{or}} - R_{\text{adv}})/(R_{\text{or}} - R_{\text{random}}) \quad (9)$$

where  $R_{\text{or}}$  and  $R_{\text{adv}}$  represent recognition accuracy of gender attribute of original and adversarial samples, respectively.  $R_{\text{random}}$  is a random probability of gender attribute classification. As we all know, the closer to  $R_{\text{random}}$  the value of  $R_{\text{adv}}$  is, the stronger the privacy performance of the adversarial sample is. Consequently, the closer the Reduction value in (9) is to 1, the better the privacy of the adversarial sample is. Because the value of gender attribute is either male or female, the value of  $R_{\text{random}}$  here is 0.5.

Table III displays the accuracy obtained by each classifier for original and adversarial samples. The adversarial samples are generated with the PrivacyNet model and our proposed Gender-AN model, respectively. When reproducing the experiments of the PrivacyNet, for the sake of fairness, we replace the pretrained VGG-Face CNN auxiliary face matcher in the original paper with the pretrained SE-ResNet-50 face matcher. The results in Table III show that the values of Reduction are close to 1 when Gender-AN is used to generate adversarial samples, which indicates the gender classifiers can basically



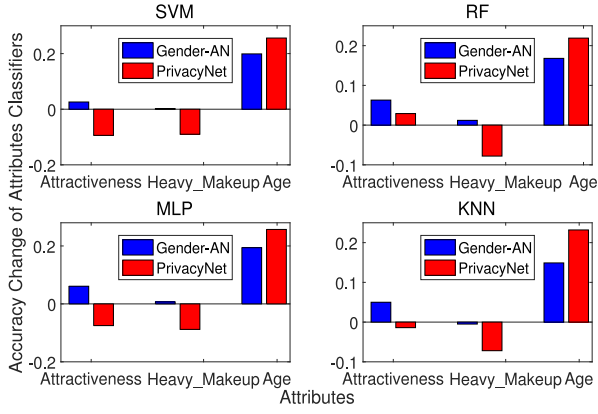


Fig. 2. Accuracy changes of attribute classifiers for gender perturbed images.

reach the level of random classification. Conversely, if the PrivacyNet model is used to generate adversarial samples, the values of Reduction are between 1.3 and 1.5, which is much greater than 1. It indicates that the performance of the gender classifiers has dropped sharply, so does the confusion ability. The experimental results demonstrate that our Gender-AN model achieves the goal of random classification of privacy, while PrivacyNet cannot reach it. Moreover, in our algorithm, the classification accuracy of the four unknown gender classifiers all drops to about 0.5. It shows that our method can generalize in the gender classifiers. That is, Gender-AN achieves both generalization and confusion goals of privacy.

### C. Evaluation on Auxiliary Attributes

The performance of auxiliary attributes will be evaluated from the following aspects: 1) recognition accuracy diversification; 2) MSE; and 3) cross-entropy of attributes vectors.

1) *Recognition Accuracy Diversification*: As mentioned above, after perturbing gender attribute, auxiliary attributes should not be disturbed. Fig. 2 gives the performance of each attribute classifier on adversarial images generated with our model and PrivacyNet. The performance is measured with the change of attribute classifier (CAC), where  $CAC = R_{or} - R_{adv}$ ,  $R_{or}$  and  $R_{adv}$  represent classification accuracy of each attribute of original and adversarial samples, respectively. Three auxiliary attributes include Attractiveness, Heavy\_Makeup, and Age. The CACs of our model are shown in blue bars, and those of PrivacyNet are shown in red bars. The scale of the y-axis is the change in the recognition rates of each attribute classifier before and after gender perturbation. i.e., 0.1 denotes 10%. As discussed above, Heavy\_Makeup attribute has closer relationship with gender than Attractiveness and Age. Therefore, as long as the gender attribute changes, the Heavy\_Makeup attribute will change more than Attractiveness and Age too in a high probability. The results in Fig. 2 also confirm it, the change of Heavy\_Makeup reaches near 20%, while the changes of Attractiveness and Age are much less, even near zero. Moreover, the change of Heavy\_Makeup of PrivacyNet is much higher than 20%, it is also much higher than that of our Gender-AN. From Fig. 2, we can see that our method outperforms PrivacyNet in most cases. However, for

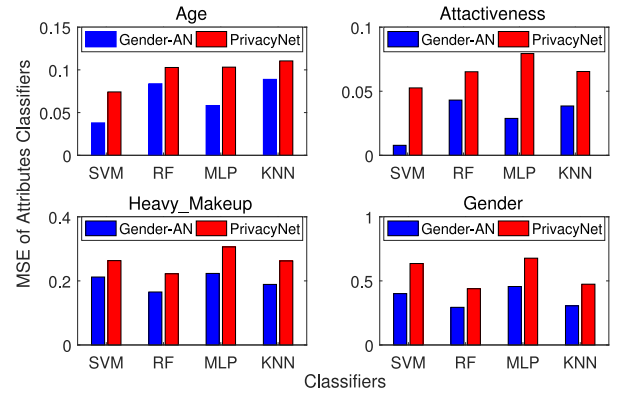


Fig. 3. MSE of attribute classifiers before and after perturbing the gender attribute.

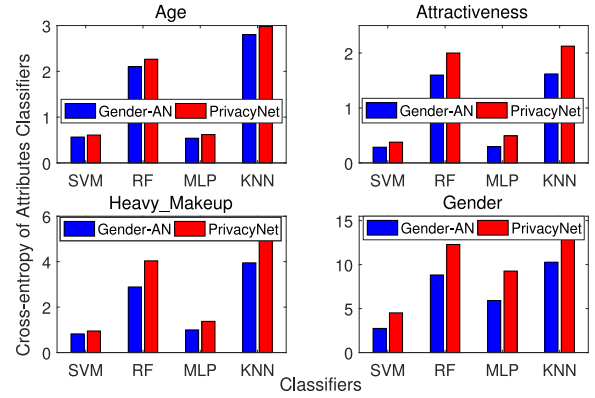


Fig. 4. Cross-entropy of attribute classifiers before and after perturbing the gender attribute.

Attractiveness, the performance of our method is not so good as PrivacyNet in KNN and RF. This is probably because each classifier has different sensitivity to changes in Attractiveness attribute. Fortunately, the diversification is less than 6%, which is within the acceptable range. Furthermore, our method performs better than PrivacyNet on both other auxiliary attributes: Heavy\_Makeup and Age. Therefore, although some categories are worse than PrivacyNet, in general, our model works better than the PrivacyNet in the vast majority of cases. The results also illustrate the success of our model in retaining the performance of auxiliary attributes classifications, which is the second goal of utility.

2) *MSE and Cross-Entropy of Attributes Vectors*: Figs. 3 and 4 display the MSE and Cross-entropy of each auxiliary attribute vector before and after perturbing the gender attribute, respectively. In these figures, the MSE and Cross-entropy results of Gender-AN are shown in blue bars, while those of PrivacyNet are shown in red bars. It can be seen from Figs. 3 and 4 that our Gender-AN obtains lower MSE and cross-entropy values than PrivacyNet. It indicates that our Gender-AN model can provide a better preservation quality of the auxiliary attributes than PrivacyNet.

### D. Evaluation on Face Matching Preservation

This section will assess the verification-preservation utility of the proposed Gender-AN model. To this end, we test

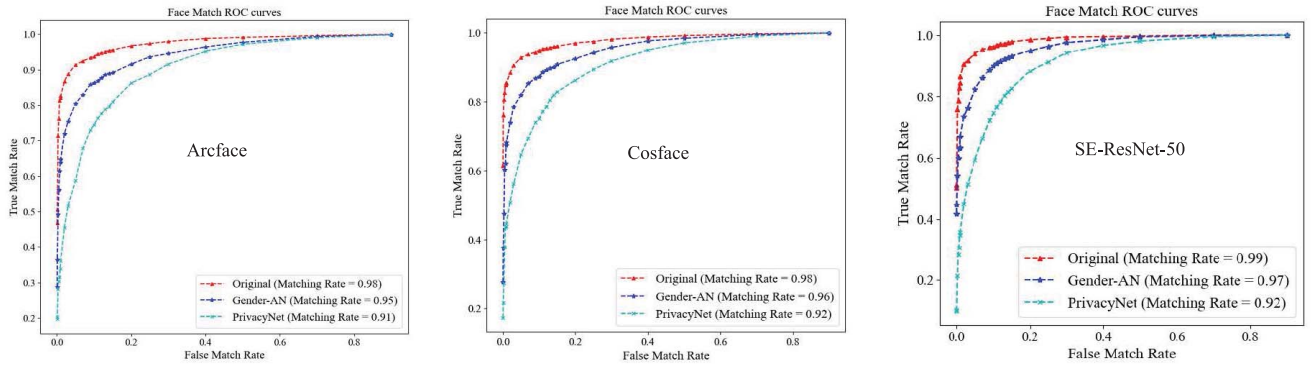


Fig. 5. ROC curves showing the performance of face matchers on original images, for our method and PrivacyNet. The results show that ROC curves of our method have smaller deviation from the ROC curve of original images, suggesting that the performance of face matching is minimally impacted, which is desired.



Fig. 6. Seven example face images from the CelebA data set along with their gender disturbed version using the PrivacyNet model and our Gender-AN model.

three famous face matchers: 1) Cosface [34]; 2) Arcface [35]; and 3) SE-ResNet-50, which has trained on the VGGFace2 data set [31]. Each face matcher is evaluated on the original test subdata set of CelebA, and the corresponding gender adversarial images generated with our Gender-AN model and PrivacyNet model, respectively. Fig. 5 exhibits their ROC curves. Compared to PrivacyNet, the ROC curve of our proposed model is much closer to the one obtained from original images. It indicates that our Gender-AN model can better maintain the face recognition ability of perturbed images. That is, it guarantees to preserve the face matching ability in utility.

#### E. Image Quality Assessment

1) *Visual Quality*: Retaining the realism of adversarial images is the last objective of this work. Fig. 6 shows the visual quality comparison results. The images shown in the first line are the test samples. The second line and third line show the corresponding adversarial samples generated by the PrivacyNet model and our Gender-AN model, respectively. As can be seen from Fig. 6, the gender adversarial images generated by our model look more realistic and natural than those generated by the PrivacyNet model. Especially, when transferring female image to male image, the hair of the generated male images becomes shorter, and even the beard appears. These are the characteristics of most real male images. However, the images generated by the PrivacyNet model do not have these features. All these phenomena

TABLE IV  
IMAGE QUALITY COMPARISON

| Metric     | LPIPS        | NIQE          | PI            |
|------------|--------------|---------------|---------------|
| PrivacyNet | 0.487        | 5.6142        | 5.1416        |
| Gender-AN  | <b>0.476</b> | <b>5.2063</b> | <b>4.7154</b> |

indicate that our model can guarantee the visibility and realism of the generated adversarial images, and it works better than the PrivacyNet model.

2) *Quantitative Evaluation*: To quantitatively evaluate the quality of gender perturbing images, three typical perceptual metrics are introduced: 1) learned perceptual image patch similarity (LPIPS) [36]; 2) natural image quality evaluator (NIQE) [37]; and 3) perceptual index (PI) [38]. LPIPS is a popular perception-oriented evaluation metric, measuring the average deep feature distance between generated samples and ground-truth samples. NIQE is a famous nonreference image quality evaluation method, which is used to assess the real image restoration without ground truth, to provide a quantitative comparison. PI is a combination of no-reference image quality of NIQE and Ma [39], formulated as  $PI = 1/2(10 - Ma) + NIQE$ . Lower LPIPS, NIQE, and IP scores indicate higher perceptual quality. The quantitative results are presented in Table IV. The LPIPS, NIQE, and PI scores of Gender-AN are lower than those of PrivacyNet, affirming that the proposed Gender-AN wins on all three metrics.

#### F. Discussion on Picking Auxiliary Attributes and Multiple Sensitive Attributes Protection

1) *Picking Auxiliary Attributes*: As described in Section III-B, to achieve the goal of confusing gender classifications in *privacy*, some proper attributes are utilized to assist in training. Thus, this section will discuss which attributes are chosen to help for training our Gender-AN model. There are 40 attributes for each image. We first selected Heavy\_Makeup as an auxiliary attribute in our model for two reasons: 1) Heavy\_Makeup has a close correlation with gender as mentioned in [16] and 2) it is a global attribute. However, if only Heavy\_Makeup is used to assist training, the

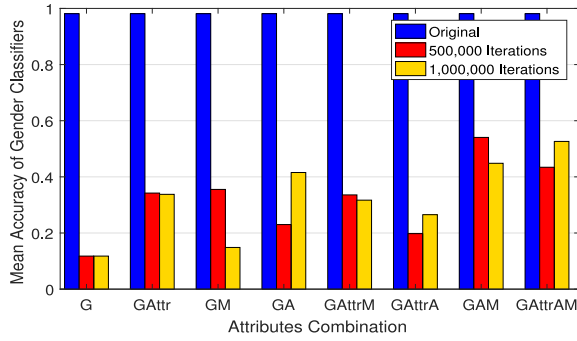


Fig. 7. Performance of gender classifiers on different attributes combinations.

gender will be flipped. This is not conducive to completely confusing gender, i.e., the probability of gender classification is not close to 0.5. Therefore, we chose Attractiveness and Age, which are closely related to Heavy\_Makeup, to train the model together. Then, we test all three attributes whether they can contribute to hiding gender information.

There are eight attribute combinations that can be used to train the Gender-AN model, i.e.: 1) Gender (G); 2) Gender & Attractiveness (GAttr); 3) Gender & Heavy\_Makeup (GM); 4) Gender & Age (GA); 5) Gender & Attractiveness & Heavy\_Makeup (GAttrM); 6) Gender & Attractiveness & Age (GAttrA); 7) Gender & Age & Heavy\_Makeup (GAM); and 8) Gender & Attractiveness & Heavy\_Makeup & Age (GAttrAM). After training, the models are used to generate adversarial samples with gender perturbed on the testing set. Then, these adversarial samples are tested on four trained gender classifiers mentioned in Section IV-B, which are SVM, RF, MLP, and KNN.

Fig. 7 shows the performance of the gender classifiers on diverse attributes combinations, where the accuracy of gender classifier is the average classification accuracy of the four classifiers. In each combination, the blue bar gives the gender classification accuracy of the original images, the orange bar shows that of adversarial samples when the model is trained 500 000 iterations, and the third bar is that of adversarial samples when the model is trained 1 000 000 iterations. From Fig. 7, it can be seen that, in the combination of G, GAttr, GM, GAttrM, and GAM, the accuracy of them drops sharply, even flips. This result is not conducive to the protection of gender privacy. Fortunately, in the last combination GAttrAM, the accuracy will first flip, then rebound back and stay at a stable value around 0.5, which can truly confuse the gender classifiers. Consequently, the last combination is selected to train our Gender-AN model, where Attractiveness, Heavy\_Makeup, and Age attributes are regarded as auxiliary attributes for gender perturbing.

2) *Multiple Sensitive Attributes Protection*: Besides gender attribute, more attributes are considered to be the sensitive information for a face image, such as age and race. Therefore, we explored the privacy concern of multiple sensitive attributes, mainly protecting gender and age attributes, simultaneously. To produce an adversarial image jointly concealing the gender and age information, the labels of gender and age are both varied in the generator. To evaluate the privacy performance of these images, whose gender and age

TABLE V  
ACCURACY FOR GENDER AND AGE CLASSIFIERS BEFORE AND AFTER PERTURBING BOTH SENSITIVE ATTRIBUTES

| Sensitive Attribute | Classifier | $R_{or}$ | $R_{adv}$    | Reduction    |
|---------------------|------------|----------|--------------|--------------|
| Gender              | SVM        | 0.986    | 0.308        | 1.395        |
|                     | RF         | 0.973    | 0.348        | 1.321        |
|                     | MLP        | 0.982    | 0.316        | 1.382        |
|                     | KNN        | 0.984    | 0.365        | 1.279        |
| Age                 | SVM        | 0.893    | <b>0.537</b> | <b>0.906</b> |
|                     | RF         | 0.853    | <b>0.507</b> | <b>0.980</b> |
|                     | MLP        | 0.868    | <b>0.463</b> | <b>1.101</b> |
|                     | KNN        | 0.867    | 0.604        | 0.717        |

are both perturbed, the Reduction metric in (9) is referred here. The results are exhibited in Table V, which displays the accuracy obtained by each classifier for gender and age attributes. For gender attribute,  $R_{random}$  is equal to 0.5 as mentioned in Section IV-B. While for the age attribute, its  $R_{random}$  is equal to 0.5 too, since age attribute is either young or old in CelebA database. It can be seen from Table V that after disturbing gender and age attributes jointly, gender classifiers are almost flipped, while age classifiers obtain the accuracy close to 0.5, and the Reductions of age are around 1. These phenomena indicate that our Gender-AN model can protect gender and age attributes simultaneously to a certain extent, but it cannot achieve the effect of completely confusing both attributes' classifiers at the same time, which is also a focus of our future work.

## V. CONCLUSION

We developed a Gender-AN model in this work, which utilizes a deep network model to impart gender privacy to face images. Gender-AN adopts an attribute-independent GANS to transfer the input face images, where the model is trained with the assistance of three proper attributes. Then, from the perspective of privacy, the gender attribute is concealed and has the generalization capacity in unseen gender classifiers. Furthermore, from the perspective of utility, our model can achieve the following three goals: 1) the face matching is facilitated; 2) other face attributes are preserved; and 3) the appearance is nature. The experimental results illustrate the efficacy of our proposed Gender-AN model in perturbing gender attribute, while not adversely impacting the verification, attribute independent, and image quality.

## REFERENCES

- [1] J. Zhang, J. Sang, X. Zhao, X. Huang, Y. Sun, and Y. Hu, "Adversarial privacy-preserving filter," in *Proc. 28th ACM Int. Conf. Multimedia (MM)*, 2020, pp. 1423–1431.
- [2] S. A. Osia *et al.*, "A hybrid deep learning architecture for privacy-preserving mobile analytics," *IEEE Internet Things J.*, vol. 7, no. 5, pp. 4505–4518, May 2020.
- [3] J. Li *et al.*, "Identity-preserving face anonymization via adaptively facial attributes obfuscation," in *Proc. 29th ACM Int. Conf. Multimedia (MM)*, 2021, pp. 3891–3899.
- [4] M. Jaiswal and E. M. Provost, "Privacy enhanced multimodal neural representations for emotion recognition," in *Proc. 32th AAAI Conf. Artif. Intell. (AAAI)*, 2020, pp. 7985–7993.
- [5] Z. Ma, Y. Liu, X. Liu, J. Ma, and K. Ren, "Lightweight privacy-preserving ensemble classification for face recognition," *IEEE Internet Things J.*, vol. 6, no. 3, pp. 5778–5790, Jun. 2019.



- [6] A. Singh, S. Fan, and M. Kankanhalli, "Human attributes prediction under privacy-preserving conditions," in *Proc. 29th ACM Int. Conf. Multimedia (MM)*, 2021, pp. 4698–4706.
- [7] M. Dusmanu, J. L. Schonberger, S. N. Sinha, and M. Pollefeys, "Privacy-preserving image features via adversarial affine subspace embeddings," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2021, pp. 14267–14277.
- [8] V. Mirjalili, S. Raschka, and A. Ross, "Gender privacy: An ensemble of semi adversarial networks for confounding arbitrary gender classifiers," in *Proc. 9th IEEE Int. Conf. Biometr. Theory Appl. Syst. (BTAS)*, 2018, pp. 1–10.
- [9] W. Zhang, S. Zhou, D. Peng, L. Yang, F. Li, and H. Yin, "Understanding and modeling of WiFi signal-based indoor privacy protection," *IEEE Internet Things J.*, vol. 8, no. 3, pp. 2000–2010, Feb. 2021.
- [10] M. Zhang, Y. Chen, and W. Susilo, "PPO-CPQ: A privacy-preserving optimization of clinical pathway query for e-healthcare systems," *IEEE Internet Things J.*, vol. 7, no. 10, pp. 10660–10672, Oct. 2020.
- [11] V. Chandrasekaran, C. Gao, B. Tang, K. Fawaz, S. Jha, and S. Banerjee, "Face-off: Adversarial face obfuscation," *Proc. Priv. Enhancing Technol.*, vol. 2021, no. 2, pp. 369–390, 2021.
- [12] A. Li, J. Guo, H. Yang, F. D. Salim, and Y. Chen, "Deepobfuscator: Obfuscating intermediate representations with privacy-preserving adversarial learning on smartphones," in *Proc. 6th Int. Conf. Internet Things Design Implement. (IoTDI)*, 2021, pp. 28–39.
- [13] T. Li and L. Lin, "AnonymousNet: Natural face de-identification with measurable privacy," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR) Workshops*, 2019, pp. 56–65.
- [14] M. Maximov, I. Elezi, and L. Leal-Taixé, "CIAGAN: Conditional identity anonymization generative adversarial networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2020, pp. 5446–5455.
- [15] V. Mirjalili, S. Raschka, A. Namboodiri, and A. Ross, "Semi-adversarial networks: Convolutional autoencoders for imparting privacy to face images," in *Proc. Int. Conf. Biometr. (ICB)*, 2018, pp. 82–89.
- [16] A. Rozsa, M. Günther, E. M. Rudd, and T. E. Boulton, "Facial attributes: Accuracy and adversarial robustness," *Pattern Recognit. Lett.*, vol. 124, pp. 100–108, Jun. 2019.
- [17] A. Boutet, C. Frindel, S. Gambs, T. Jourdan, and R. C. Nogueira, "DySan: Dynamically sanitizing motion sensor data against sensitive inferences through adversarial networks," in *Proc. ACM Asia Conf. Comput. Commun. Secur. (ASIA CCS)*, 2021, pp. 672–686.
- [18] V. Mirjalili, S. Raschka, and A. Ross, "PrivacyNet: Semi-adversarial networks for multi-attribute face privacy," *IEEE Trans. Image Process.*, vol. 29, pp. 9400–9412, 2020.
- [19] Q. Li, J. S. Gundersen, R. Heusdens, and M. G. Christensen, "Privacy-preserving distributed processing: Metrics, bounds and algorithms," *IEEE Trans. Inf. Forensics Security*, vol. 16, pp. 2090–2103, 2021.
- [20] S. Chhabra, R. Singh, M. Vatsa, and G. Gupta, "Anonymizing k-facial attributes via adversarial perturbations," in *Proc. 27th Int. Joint Conf. Artif. Intell. (IJCAI)*, 2018, pp. 656–662.
- [21] V. Mirjalili, S. Raschka, and A. Ross, "FlowSAN: Privacy-enhancing semi-adversarial networks to confound arbitrary face-based gender classifiers," *IEEE Access*, vol. 7, pp. 99735–99745, 2019.
- [22] A. Othman and A. Ross, "Privacy of facial soft biometrics: Suppressing gender but retaining identity," in *Proc. Eur. Conf. Comput. Vis. (ECCV) Workshop*, 2014, pp. 682–696.
- [23] S. Wang, U. M. Kelly, and R. N. Veldhuis, "Gender obfuscation through face morphing," in *Proc. IEEE Int. Workshop Workshop Forensics (IWWF)*, 2021, pp. 1–6.
- [24] P. Terhörst, N. Damer, F. Kirchbuchner, and A. Kuijper, "Suppressing gender and age in face templates using incremental variable elimination," in *Proc. Int. Conf. Biometr. (ICB)*, 2019, pp. 1–8.
- [25] A. Morales, J. Fierrez, R. Vera-Rodriguez, and R. Tolosana, "SensitiveNets: Learning agnostic representations with application to face images," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 43, no. 6, pp. 2158–2164, 2021.
- [26] B. Bortolato et al., "Learning privacy-enhancing face representations through feature disentanglement," in *Proc. 15th IEEE Int. Conf. Automat. Face Gesture Recognit. (FG)*, 2020, pp. 495–502.
- [27] B. Meden et al., "Privacy-enhancing face biometrics: A comprehensive survey," *IEEE Trans. Inf. Forensics Security*, vol. 16, pp. 4147–4183, 2021.
- [28] S. Baluja and I. Fischer, "Learning to attack: Adversarial transformation networks," in *Proc. 32th AAAI Conf. Artif. Intell. (AAAI)*, 2018, pp. 2687–2695.
- [29] N. Carlini and D. Wagner, "Towards evaluating the robustness of neural networks," in *Proc. 38th IEEE Symp. Secur. Privacy (SP)*, 2017, pp. 39–57.
- [30] M. Liu et al., "STGAN: A unified selective transfer network for arbitrary image attribute editing," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2019, pp. 3673–3682.
- [31] Q. Cao, L. Shen, W. Xie, O. M. Parkhi, and A. Zisserman, "VGGFace2: A dataset for recognising faces across pose and age," in *Proc. 13th IEEE Int. Conf. Automat. Face Gesture Recognit. (FG)*, 2018, pp. 67–74.
- [32] M. Arjovsky, S. Chintala, and L. Bottou, "Wasserstein generative adversarial networks," in *Proc. 34th Int. Conf. Mach. Learn. (ICML)*, 2017, pp. 214–223.
- [33] X. Ding, H. Fang, Z. Zhang, K.-K. R. Choo, and H. Jin, "Privacy-preserving feature extraction via adversarial training," *IEEE Trans. Knowl. Data Eng.*, early access, May 26, 2020, doi: 10.1109/TKDE.2020.2997604.
- [34] H. Wang et al., "CosFace: Large margin cosine loss for deep face recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2018, pp. 5265–5274.
- [35] J. Deng, J. Guo, N. Xue, and S. Zafeiriou, "ArcFace: Additive angular margin loss for deep face recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2019, pp. 4690–4699.
- [36] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2018, pp. 586–595.
- [37] A. Mittal, R. Soundararajan, and A. C. Bovik, "Making a 'completely blind' image quality analyzer," *IEEE Signal Process. Lett.*, vol. 20, no. 3, pp. 209–212, Mar. 2013.
- [38] Y. Blau, R. Mechrez, R. Timofte, T. Michaeli, and L. Zelnik-Manor, "The 2018 PIRM challenge on perceptual image super-resolution," in *Proc. Eur. Conf. Comput. Vis. (ECCV) Workshop*, 2018, pp. 334–355.
- [39] C. Ma, C.-Y. Yang, X. Yang, and M.-H. Yang, "Learning a no-reference quality metric for single-image super-resolution," *Comput. Vis. Image Understand.*, vol. 158, pp. 1–16, May 2017.

**Deyan Tang** (Member, IEEE) received the B.S. degree from Inner Mongolia University, Hohhot, China, in 2010, and the M.S. degree from Hunan University, Changsha, China, in 2013, where she is currently pursuing the Ph.D. degree with the College of Computer Science and Electronic Engineering.

Her research interests include pattern recognition, signal processing, and machine learning.

**Siwang Zhou** (Member, IEEE) received the B.S. degree from Fudan University, Shanghai, China, in 1995, the M.S. degree from Xiangtan University, Xiangtan, China, in 2004, and the Ph.D. degree from Hunan University, Changsha, China, in 2007.

He has been a Professor with the College of Computer Science and Electronic Engineering, Hunan University. His research interests include image compressive sensing, deep learning, and Internet of Things.

**Hongbo Jiang** (Senior Member, IEEE) received the Ph.D. degree from Case Western Reserve University, Cleveland, OH, USA, in 2008.

He is currently a Full Professor with the College of Computer Science and Electronic Engineering, Hunan University, Changsha, China. His research concerns signal processing and computer networking, especially algorithms and protocols for wireless and mobile networks.

**Haowen Chen** received the B.S. and M.S. degrees in computer science from Hunan University, Changsha, China, in 2002 and 2007, respectively, and the Ph.D. degree from the College of Computer Science and Electronic Engineering, Hunan University.

He is currently an Associate Professor with Hunan University. His research interests include security and privacy issues in social network.

**Yonghe Liu** received the B.S. and M.S. degrees from Tsinghua University, Beijing, China, in 1998 and 1999, respectively, and the Ph.D. degree from Rice University, Houston, TX, USA, in 2004.

He is an Associate Professor with the Department of Computer Science and Engineering, University of Texas at Arlington, Arlington, TX, USA. His current research interests include wireless networking, security, system integration, and signal processing.