

Deep Identity-Aware Transfer of Facial Attributes

Mu Li, Wangmeng Zuo, *Senior Member, IEEE* David Zhang, *Fellow, IEEE* and Jane You, *Member, IEEE*

Abstract—This paper presents a Deep convolutional network model for Identity-Aware Transfer (DIAT) of facial attributes. Given the source input image and the reference attribute, DIAT aims to generate a facial image that owns the reference attribute as well as keeps the same or similar identity to the input image. In general, our model consists of a mask network and an attribute transform network which work in synergy to generate photo-realistic facial image with the reference attribute. Considering that the reference attribute may be only related to some parts of the image, the mask network is introduced to avoid the incorrect editing on attribute irrelevant region. Then the estimated mask is adopted to combine the input and transformed image for producing the transfer result. For joint training of transform network and mask network, we incorporate the adversarial attribute loss, identity aware adaptive perceptual loss and VGG-FACE based identity loss. Furthermore, a denoising network is presented to serve for perceptual regularization to suppress the artifacts in transfer result, while an attribute ratio regularization is introduced to constrain the size of attribute relevant region. Our DIAT can provide a unified solution for several representative facial attribute transfer tasks, *e.g.*, expression transfer, accessory removal, age progression and gender transfer, and can be extended for other face enhancement tasks such as face hallucination. The experimental results validate the effectiveness of the proposed method. Even for the identity-related attribute (*e.g.*, gender), our DIAT can obtain visually impressive results by changing the attribute while retaining most identity-aware features.

Index Terms—Facial attribute transfer, generative adversarial nets, convolutional networks, perceptual loss.

I. INTRODUCTION

Face attributes, *e.g.*, gender and expression, can not only provide a natural description of facial images [1], but also offer a unified viewpoint for understanding many facial animation and manipulation tasks. For example, the goal of facial avatar [2] and reenactment [3] is to transfer the facial expression attributes of a source actor to a target actor. In most applications such as expression transfer, accessory removal and age progression, the animation only modifies the related attribute without changing the identity. But for some other

tasks, the change of some attributes, *e.g.*, gender and ethnicity, will inevitably alter the identity of the source image.

In recent years, a variety of methods have been developed for specific facial attribute transfer tasks, and have achieved impressive results. For expression transfer, approaches have been suggested to create 3D or image-based avatars from handheld video [2], while face trackers and expression modeling have been investigated for offline and online facial reenactment [3], [4]. For age progression, explicit and implicit synthesis methods have been proposed for different image models [5], [6]. Hair style generation and replacement have also been studied in literatures [7], [8].

Convolutional neural network (CNN)-based models have also been investigated for human face generation with attributes. Kulkarni *et al.* [9] propose deep convolution inverse graphic network (DG-IGN). This method requires a large number of faces of a single person for training, and can only generate faces with different pose and light. Gauthier [10] develops a conditional generative adversarial network (cGAN) to generate facial image from a noise distribution and conditional attributes. Yan *et al.* [11] suggest an attribute-conditioned deep variational auto-encoder which extracts the latent variables from a reference image and combines them with attributes to produce the generated image with a generative model. Oord *et al.* [12] propose a conditional image generation model based on PixelCNN decoder for image generation conditioned on an arbitrary feature vector. However, the identity of the generated face is not emphasized in [10], [11], [12], making them not directly applicable to attribute transfer.

Motivated by the strong capability of CNN in modeling complex transformation [13] and capturing perceptual similarity [14], several approaches have also been suggested for facial attribute transfer. Li *et al.* [15] suggest a CNN-based attribute transfer model from the optimization perspective, but both run time and transfer quality is far from satisfying. Considering that it is impracticable to collect labeled data for supervised learning, the generative adversarial net (GAN) framework [16] usually is adopted for handling this task [17], [18], [19], [20]. However, visible artifacts and over-smoothing usually are inevitable in the transfer result for these methods.

In this paper, we present a novel Deep CNN model for Identity-Aware Transfer (DIAT) of facial attributes which can provide a unified solution to several facial animation and manipulation tasks, *e.g.*, expression transfer, accessory removal, age progression, and gender transfer. For each reference attribute label, we train a CNN model for the transfer of the input image to the desired attribute. Note that the reference attribute may be only related to some parts of the image. To avoid

This project is partially supported by the HK RGC/GRF grant (under no. PolyU 5313/12E and PolyU 152212/14E) and the National Natural Science Foundation of China (NSFC) under Grant No. 61671182.

Mu Li and Jane You are with the Department of Computing, Hong Kong Polytechnic University, Hong Kong, e-mail: (csmuli@comp.polyu.edu.hk, csyjia@comp.polyu.edu.hk).

Wangmeng Zuo is with the School of Computer Science and Technology, Harbin Institute of Technology, Harbin, China, e-mail: (cswmzuo@gmail.com).

David Zhang is with the School of Science and Engineering, The Chinese University of Hong Kong (Shenzhen), Shenzhen, China, e-mail: (csdzhang@comp.polyu.edu.hk).

Manuscript received XXX; revised XXX

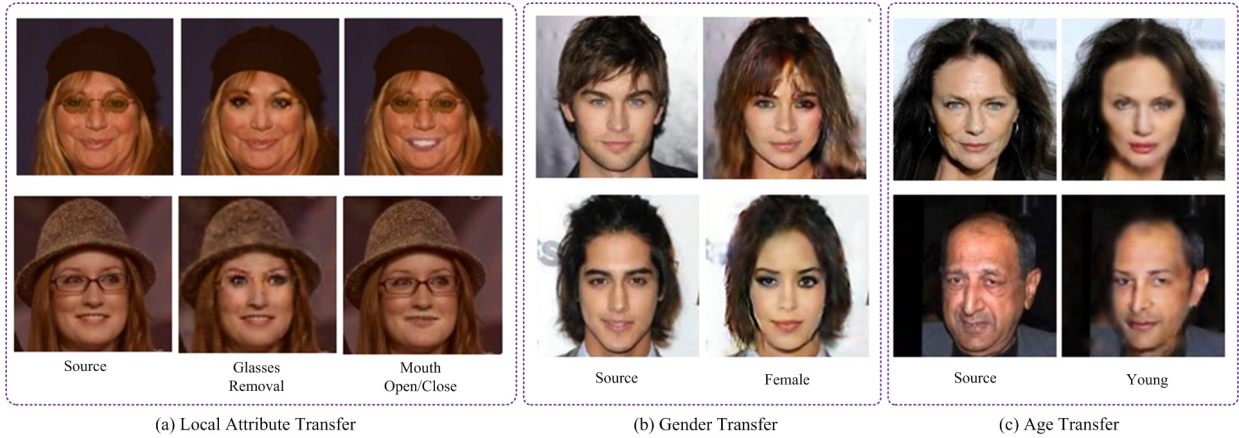


Fig. 1. Illustration of the results by our DIAT on several facial attribute transfer tasks, including *glasses removal*, *mouth open/close*, *gender transfer* and *age transfer*.

the the incorrect editing on attribute irrelevant region, our model consists of a mask network and an attribute transform network. The attribute transform network is presented to edit the input image for generating the desired attribute. While the mask network is adopted to estimate a mask of the attribute relevant region for guiding the combination of the input image and transformed image. Then attribute transform network and mask network work collaboratively to generate final photo-realistic transfer result.

For attribute transfer, the ground truth transfer results generally are very difficult or even impossible to obtain. Therefore, we follow the GAN framework to train the model. As for training data, we only consider the binary attribute labels presented in the large-scale CelebFaces Attributes (CelebA) dataset [21]. To capture the convolutional feature distribution of each attribute, we construct an attribute guided set using all the images with the desired attribute in CelebA. Then, the input set is defined as a set of input images without the reference attribute. Due to the infeasibility of ground truth transfer results, two alternative losses, *i.e.*, adversarial attribute loss and identity-aware perceptual loss, are incorporated for unsupervised training of our DIAT model. Furthermore, two regularizers, *i.e.*, perceptual regularization and attribute ratio regularization, are also introduced in the learning of DIAT.

In terms of attribute transfer, the generated image should have the desired attribute label. Following the GAN framework, we define the adversarial attribute loss on the attribute discriminator to require the generated image to have the desired attribute. As for identity-aware transfer, our DIAT requires that the generated image should keep the same or similar identity to the input image. To this end, the identity-aware perceptual loss is introduced on the convolutional feature map of a CNN model to model the content similarity between the reference face and the generated face. Instead of adopting any pre-trained CNNs, we suggest to define the perceptual loss on the attribute discriminator, which can be adaptively trained along with the learning procedure, and is named as adaptive perceptual loss. Compared with conventional perceptual loss, ours is more effective in computation, tailored to our specific

attribute transfer task, and can serve as a kind of hidden-layer supervision [22] or regularization to ease the training of the DIAT model. To further encourage the identity keeping property, we add an identity loss by minimizing the distance between feature representations of the generated face and the reference. Pre-trained VGG-Face is used to extract the identity related features for face verification.

The model objective of DIAT also consider two regularizers, *i.e.*, perceptual regularization and attribute ratio regularization. To suppress the artifacts, we propose a denoising network to serve for a perceptual regularization on the generated image. To guide the learning of mask network, an attribute ratio regularization is introduced to constrain the size of attribute relevant region. Finally, our DIAT model can be learned from training data by incorporating adversarial attribute loss, adaptive perceptual loss with perceptual regularization and attribute ratio regularization.

Extensive experiments are conducted on CelebA and real images from the website *iStock*¹. As illustrated in Fig. 1, our DIAT performs favorably in attribute transfer with minor or no modification on the identity of the input faces. Even for some identity-related attributes (*e.g.*, gender), our DIAT can obtain visually impressive transfer result while retaining most identity-relevant features. Computational efficiency is also a prominent merit of our method. In the testing stage, our DIAT can process more than one hundred of images within one second. Furthermore, our model can be extended to face hallucination, and is effective in generating photo-realistic high resolution images. A preliminary report of this work is given in 2016 [23]. To sum up, our contribution is three-fold:

- 1) A novel DIAT model is developed for facial attribute transfer. For better preserving of attribute irrelevant feature, our model comprises a mask network and an attribute transform network, which collaborate to generate the transfer result and can be jointly learned from training data.
- 2) Adversarial attribute loss, adaptive perceptual loss, identity loss, perceptual regularization, and attribute ratio

¹<https://www.istockphoto.com/hk>

regularization are incorporated for training our DIAT model. The adversarial attribute loss is adopted to make the transfer result exhibit the desired attribute, and the adaptive perceptual loss is defined on the discriminator for identity-aware transfer while improving training efficiency. Moreover, perceptual regularization and attribute ratio regularization are further introduced for suppressing the artifacts and constraining the mask network.

- 3) Experimental results validate the effectiveness and efficiency of our method for identity-aware attribute transfer. Our DIAT can be used for the transfer of either local (*e.g.*, mouth), global (*e.g.*, age progression) or identity-related (*e.g.*, gender) attributes, and can be extended to face hallucination.

The remainder of the paper is organized as follows. Section II gives a brief survey on relevant work. Section III describes the model and learning of our DIAT method. Section IV reports the experimental results on facial attribute transfer and face hallucination. Finally, Section V ends this work with several concluding remarks.

II. RELATED WORK

Deep convolutional neural networks (CNNs) not only have achieved unprecedented success in versatile high level vision problems [24], [25], [26], [27], [28], but also exhibited their remarkable power in understanding, generating, and recovering images [16], [29], [30], [31], [32], [33]. In this section, we focus on the task of facial attribute transfer, and briefly survey the CNN models for image generation and face generation.

A. CNN for image generation

Generative image modeling is a critical issue for image generation and many low level vision problems. Conventional sparse [34], low rank [35], FRAME [36] and non-local similarity [37], [38] based models usually are limited in capturing highly complex and long-range dependence between pixels. For better image modeling, a number of CNN-based methods have been proposed, including convolutional auto-encoder [9], PixelCNN and PixelRNN [39], and they have been applied to image completion and generation.

Several CNN architectures have been developed for image generation. Fully convolutional networks can be trained in the supervised learning manner to generate an image from an input image [40], [41]. The generative CNN model [42] stacks four convolution layers upon five fully connected layers to generate images from object description. Kulkarni *et al.* suggest the Deep Convolution Inverse Graphics Network (DC-IGN), which follows the variational autoencoder architecture [43] to transform the input image into different pose and lighting condition. However, both generative CNN [42] and DC-IGN [43] require many labeled images in training.

To visualize and understand CNN features, several methods have been proposed to reconstruct images by inverting deep representation [30] or maximizing class score [44]. Subsequently, Gatys *et al.* [14] suggest to combine content and style losses defined on deep representation on the off-the-shelf CNNs for artistic style transfer. To improve the efficiency,

alternative approaches have been proposed by substituting the iterative optimization procedure with pre-trained feed-forward CNN [13], [45]. And perceptual loss has also been adopted for style transfer and other generation tasks [13]. Motivated by these works, both identity-aware adaptive perceptual loss and perceptual regularization are exploited in our DIAT model to meet the requirement of facial attribute transfer.

Another representative approach is generative adversarial network (GAN), where a discriminator and a generator are alternately trained as an adversarial game [16]. The generator aims to generate images to match the data distribution, while the discriminator attempts to distinguish between the generated images and the training data. Laplacian Pyramid of GANs is further suggested to generate high quality image in a coarse-to-fine manner [46]. Radford *et al.* [47] extend GAN with the fully deep convolutional networks (*i.e.*, DCGAN) for image generation. To learn disentangled representations, information-theoretic extension of GAN is proposed by maximizing the mutual information between a subset of noise variables and the generated results [48]. In [49], [50], WGAN and WGAN-GP minimize the Wasserstein-1 distance between the generated distribution and the real distribution to improve the stability of learning generator. In this work, we adopt the WGAN framework to learn our DIAT model, and further suggest adaptive perceptual loss for identity-aware transfer and perceptual regularization to suppress visual artifacts.

B. CNN for face generation

Facial attribute transfer has received considerable recent attention. Larsen *et al.* [18] present to combine variational autoencode with GAN (VAE/GAN) for image generation. By modeling the attribute vector as the difference between the mean latent representations of the images with and without the reference attribute, VAE/GAN can provide a flexible solution to arbitrary facial attribute transfer, but is limited in transfer performance. Li *et al.* [15] suggest an attribute driven and identity-preserving face generation model by solving an optimization problems with perceptual loss, which is computationally expensive and cannot obtain high quality results. Perarnau *et al.* [17] adopt an encoder-decoder architecture, where attribute transfer can be conducted by editing the latent representation. Shen *et al.* [20] learn the residual image in the GAN framework, and adopt dual learning to learn two reverse attribute transfer models simultaneously. Zhou *et al.* [19] propose a model to learn object transfiguration from two sets of unpaired images that have the opposite attribute. However, most existing methods cannot achieve high quality transfer results, and visible artifacts and over-smoothing usually are inevitable. In comparison, our DIAT model can achieve much better transfer results than the competing methods [15], [17], [18], [51].

Besides, CNNs have also been developed for other face generation tasks. For painting style transfer of head portrait, Selim *et al.* [52] modify the perceptual loss to balance the contribution of the input photograph and the aligned exemplar painting. Gucluturk *et al.* train a feed-forward CNN with perceptual loss for sketch inversion. Yeh *et al.* [51] apply DCGAN to semantic face inpainting in an optimization manner.

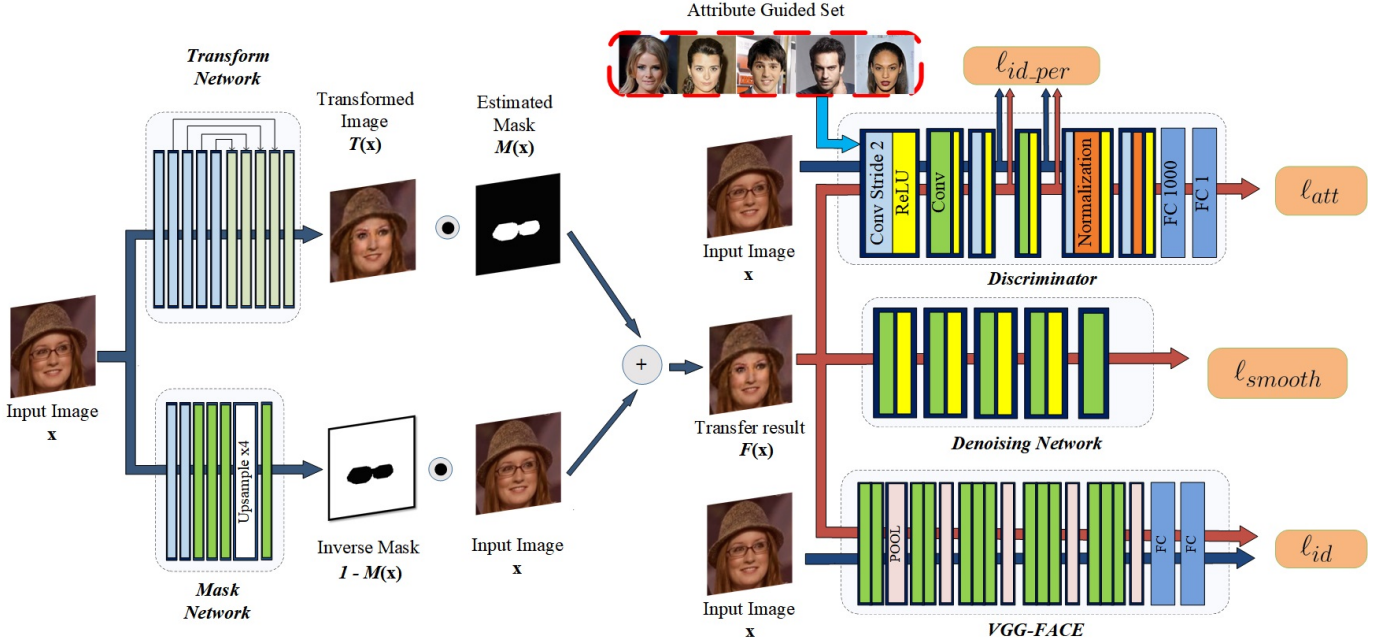


Fig. 2. Schematic illustration of our DIAT model. Here we use *glasses removal* as an example. The whole attribute transfer network $F(\mathbf{x})$ includes two sub-networks, *i.e.*, a mask network to find the attribute relevant region $M(\mathbf{x})$ and an attribute transform network to produce the transformed image $T(\mathbf{x})$. Then $M(\mathbf{x})$ and $T(\mathbf{x})$ collaborate to generate the transfer result $F(\mathbf{x}) = M(\mathbf{x}) \circ T(\mathbf{x}) + (1 - M(\mathbf{x})) \circ \mathbf{x}$. In order to learn $F(\mathbf{x})$ from training data, we incorporate adversarial attribute loss ℓ_{att} , identity loss ℓ_{id} , adaptive perceptual loss ℓ_{id_per} with perceptual regularization ℓ_{smooth} . Besides, an attribute ratio regularization is also adopted to constrain the estimated mask $M(\mathbf{x})$.

III. DEEP CNNs FOR IDENTITY-AWARE ATTRIBUTE TRANSFER

In this section, we present our DIAT model for identity-aware transfer of facial attribute. As illustrated in Fig. 2, our model involves a mask network and an attribute transfer network which collaborate to produce the transfer result. To train our model, we incorporate the adversarial attribute loss, adaptive perceptual loss, perceptual regularization and attribute ratio regularization.

A. Network architecture

Most facial attributes, *e.g.*, expression and accessory, are local-based and only related to part of facial image. Even for global attributes such as age and gender, some parts, *e.g.*, the background, should also keep the same with the source image. For the sake of preserving attribute irrelevant feature, it is natural to only perform attribute transfer in image region related to specific attribute. However, it is not a trivial issue to find the attribute relevant region. One possible solution is to manually specify the relevant region for each attribute given a new transfer task, but it undoubtedly restricts the universality and adaptivity of the solution.

In this work, we aim to provide a unified solution to attribute transfer, which indicates that we only require to prepare training data and retrain the model when a new transfer task comes. To this end, our whole attribute transfer network is comprised of two sub-networks, *i.e.*, mask network and attribute transfer network. Both mask network and attribute transfer network take the source image \mathbf{x} as input. The mask

network is utilized to predict a mask $M(\mathbf{x})$ to indicate the attribute relevant region, while the attribute transfer network is used to produce the transformed image $T(\mathbf{x})$. Given $M(\mathbf{x})$ and $T(\mathbf{x})$, the final transfer result can be obtained by,

$$F(\mathbf{x}) = M(\mathbf{x}) \circ T(\mathbf{x}) + (1 - M(\mathbf{x})) \circ \mathbf{x}, \quad (1)$$

where \circ denotes the element-wise product operator. We also note that both attribute transfer network and mask network can be learned from training data in an end-to-end manner. In the following, we describe the architecture of attribute transfer network and mask network, respectively.

Attribute transfer network. We adopt the Unet [53] for attribute transfer due to its good tradeoff between efficiency and reconstruction ability. In general, the Unet architecture involves an encoder subnetwork and a decoder subnetwork, then skip connection and pooling operation are further introduced to exploit multi-scale information. As for attribute transfer, we design a 10-layer Unet, which includes 5 convolution layers for encoding and another 5 convolution layers for decoding. In the encoder, we use convolution with stride 2 for downsampling. In the decoder, a depth to width (DTOW) layer [54] is deployed for upsampling, and the element-wise summation operation is adopted to fuse the feature maps from the encoder and decoder subnetworks. The detailed parameters of the attribute transfer network are summarized in Table I.

Mask network. As for mask network, we first adopt a 5-layer fully convolutional network to generate a 32×32 binary mask for indicating the attribute relevant region. A batch normalization layer is added after each convolution layer. Then, $4 \times$ upsampling is deployed by simply replicating each

TABLE I
ARCHITECTURE OF THE ATTRIBUTE TRANSFORM NETWORK.

Layer	Activation size
Input	$3 \times 128 \times 128$
conv1, $4 \times 4 \times 64$, pad 1, stride 2	$64 \times 64 \times 64$
conv2, $4 \times 4 \times 128$, pad 1, stride 2	$128 \times 32 \times 32$
conv3, $4 \times 4 \times 256$, pad 1, stride 2	$256 \times 16 \times 16$
conv4, $4 \times 4 \times 512$, pad 1, stride 2	$512 \times 8 \times 8$
conv5, $4 \times 4 \times 512$, pad 1, stride 2	$512 \times 4 \times 4$
DTOW, stride 2	$128 \times 8 \times 8$
conv6, $3 \times 3 \times 512$, pad 1, stride 1	$512 \times 8 \times 8$
Element-wise add, conv6 and conv4	$512 \times 8 \times 8$
DTOW, stride 2	$128 \times 16 \times 16$
conv7, $3 \times 3 \times 256$, pad 1, stride 1	$256 \times 16 \times 16$
Element-wise add, conv7 and conv3	$256 \times 16 \times 16$
DTOW, stride 2	$64 \times 32 \times 32$
conv8, $3 \times 3 \times 128$, pad 1, stride 1	$128 \times 32 \times 32$
Element-wise add, conv8 and conv2	$128 \times 32 \times 32$
DTOW, stride 2	$32 \times 64 \times 64$
conv9, $3 \times 3 \times 64$, pad 1, stride 1	$64 \times 64 \times 64$
Element-wise add, conv9 and conv1	$64 \times 64 \times 64$
DTOW, stride 2	$16 \times 128 \times 128$
conv10, $5 \times 5 \times 3$, pad 2, stride 1	$3 \times 128 \times 128$

element in the 32×32 binary mask 4×4 times. In order to make the generated image smooth, we further utilize a 5×5 Gaussian filter with the standard deviation of 1.6 to produce the final mask $M(\mathbf{x})$. To sum up, the details of the mask network are given in Table II.

TABLE II
ARCHITECTURE OF THE MASK NETWORK.

Layer	Activation size
Input	$3 \times 128 \times 128$
conv, $4 \times 4 \times 32$, pad 1, stride 2	$32 \times 64 \times 64$
conv, $4 \times 4 \times 64$, pad 1, stride 2	$64 \times 32 \times 32$
conv, $3 \times 3 \times 64$, pad 1, stride 1	$64 \times 32 \times 32$
conv, $3 \times 3 \times 64$, pad 1, stride 1	$64 \times 32 \times 32$
conv, $3 \times 3 \times 1$, pad 1, stride 1	$1 \times 32 \times 32$
binarization	
$4 \times$ upsampling	$1 \times 128 \times 128$
conv, $5 \times 5 \times 1$, pad 2, stride 1	$1 \times 128 \times 128$

In our mask network, ReLU is adopted for nonlinearity for the first 4 convolution layers. As for the fifth convolution layer, we adopt the sigmoid nonlinearity, and the binarization operation is then used to obtain the binary mask,

$$B(e_{ijk}) = \begin{cases} 1, & \text{if } e_{ijk} > 0.5 \\ 0, & \text{if } e_{ijk} \leq 0.5 \end{cases} \quad (2)$$

where e_{ijk} denotes an element of the feature map. However, the gradient of the binarizer $B(e_{ijk})$ is zero almost everywhere except that it is infinite when $e_{ijk} = 0.5$, making any layer before the binarizer never be updated during training. As a remedy, we follow the straight-through estimator on gradient [55], and introduce a piecewise linear proxy function $\tilde{B}(e_{ijk})$ to approximate $B(e_{ijk})$,

$$\tilde{B}(e_{ijk}) = \begin{cases} 1, & \text{if } e_{ijk} > 1 \\ e_{ijk}, & \text{if } 1 \leq e_{ijk} \leq 0. \\ 0, & \text{if } e_{ijk} < 0 \end{cases} \quad (3)$$

During training, $B(e_{ijk})$ is still used in forward-propagation

calculation, while $\tilde{B}(e_{ijk})$ is used in back-propagation, with its gradient computed by,

$$\tilde{B}'(e_{ijk}) = \begin{cases} 1, & \text{if } 1 \leq e_{ijk} \leq 0. \\ 0, & \text{otherwise} \end{cases}. \quad (4)$$

B. Model objective

By enforcing proper constraints on transfer result $F(\mathbf{x})$, both the attribute transform network and mask network can be learned from training data. However, the ground truth of attribute transfer usually is unavailable and not unique. For example, it is generally impossible to obtain the ground truth of gender transfer in reality. Instead, the training data used in this work includes a guided set \mathcal{X}_G of images with the desired reference attribute and a source set \mathcal{X}_S of input images not with the reference attribute. And we do not require the images from guided set to have the same identity with those from source set.

For the sake of identity-aware attribute transfer, we define two alternative losses: (i) adversarial attribute loss to make the transfer result exhibit the desired attribute, and (ii) adaptive perceptual loss and identity loss to make the generated image keep the same or similar identity to input image. To suppress the visual artifacts of the transfer result, we further include (iii) a perceptual regularization defined on a denoising network. Finally, (iv) an attribute ratio regularization is deployed on $\sum_{i,j,k} B(e_{ijk})$ to guide the learning of mask network. In the following, we provide more details on these losses and regularization terms.

Adversarial attribute loss. Adversarial strategy is a common strategy that is widely used in security problems [56], [57]. For computer vision, an adversarial learning framework called generative adversarial networks (GANs) also shows powerful ability on generating images that fulfil the distribution of a set of images without any ground truth targets. Considering the infeasibility of obtaining the ground truth transfer result, we define the adversarial attribute loss based on the guided set \mathcal{X}_G and the source set \mathcal{X}_S . If the guided set \mathcal{X}_G is of large scale, it can provide a natural representation of the attribute distribution. Therefore, the goal of adversarial attribute loss is to make that the distribution of generated images matches the real attribute distribution. To this end, we adopt the generative adversarial network framework, where the generator is the attribute transfer network $F(\mathbf{x})$, and the discriminator D is used to define the adversarial attribute loss. The details of the discriminator are provided in Table III, which contains 6 convolution layers followed by another two fully-connected layers.

Denote by \mathbf{x} an input image from \mathcal{X}_S , and \mathbf{a} an image from \mathcal{X}_G . Let $p_{source}(\mathbf{x})$ be the distribution of the input images, $p_{att}(\mathbf{a})$ be the distribution of the images with the reference attribute. The discriminator is defined as $D(\mathbf{a})$ to output the probability that the image \mathbf{a} comes from the set \mathcal{X}_G . To train the generator and the discriminator, we take use of the following adversarial attribute loss,

$$\min_F \max_D \mathbb{E}_{\mathbf{a} \sim p_{att}(\mathbf{a})} \log D(\mathbf{a}) + \mathbb{E}_{\mathbf{x} \sim p_{source}(\mathbf{x})} \log[1 - D(F(\mathbf{x}))]. \quad (5)$$

TABLE III
NETWORK ARCHITECTURE OF THE ATTRIBUTE DISCRIMINATOR.

Layer	Activation size
Input	$3 \times 128 \times 128$
conv, $8 \times 8 \times 32$, pad 3, stride 2	$32 \times 64 \times 64$
conv, $3 \times 3 \times 32$, pad 1, stride 1	$32 \times 64 \times 64$
conv, $4 \times 4 \times 64$, pad 1, stride 2	$64 \times 32 \times 32$
conv, $3 \times 3 \times 64$, pad 1, stride 1	$64 \times 32 \times 32$
conv, $4 \times 4 \times 128$, pad 1, stride 2	$128 \times 16 \times 16$
conv, $4 \times 4 \times 128$, pad 1, stride 2	$128 \times 8 \times 8$
Fully connected layer with 1000 hidden units	1000
Fully connected layer with 1 hidden units	1

In order to improve the training stability, the improved Wasserstein GAN is adopted by defining the loss as,

$$\min_F \max_D \mathbb{E}_{\mathbf{a} \sim p_{att}(\mathbf{a})} [D(\mathbf{a})] - \mathbb{E}_{\mathbf{x} \sim p_{source}(\mathbf{x})} [D(F(\mathbf{x}))]. \quad (6)$$

For simplicity, we respectively define the adversarial attribute losses for the generator and discriminator as follows,

$$\min_F \ell_{att,F} = \{-\mathbb{E}_{\mathbf{x} \sim p_{source}(\mathbf{x})} [D(F(\mathbf{x}))]\}, \quad (7)$$

$$\min_D \ell_{att,D} = \{\mathbb{E}_{\mathbf{x} \sim p_{source}(\mathbf{x})} [D(F(\mathbf{x}))] - \mathbb{E}_{\mathbf{a} \sim p_{att}(\mathbf{a})} [D(\mathbf{a})]\}. \quad (8)$$

Adaptive perceptual loss. The adaptive perceptual loss is introduced to guarantee that the transfer result keeps the same or similar identity with the input image. Due to identity is a high level semantic concept, it is not proper to define identity-aware loss by forcing two images to be exactly the same in pixel domain. Instead, we define the squared-error loss on the feature representations of the discriminator, resulting in our adaptive perceptual loss.

Denote by D the discriminator, and $D_l(\mathbf{x})$ the feature map of the l -th convolution layer. C_l , H_l and W_l represent the channel number, height, and width of the feature map, respectively. We then define the perceptual loss between \mathbf{x} and $\hat{\mathbf{x}} = F(\mathbf{x})$ on the l -th convolution layer as,

$$\ell_{adaptive}^{D,l}(\hat{\mathbf{x}}, \mathbf{x}) = \frac{1}{2C_l H_l W_l} \|D_l(\hat{\mathbf{x}}) - D_l(\mathbf{x})\|_F^2. \quad (9)$$

And the identity-aware adaptive perceptual loss is further defined as,

$$\ell_{id_per}(\mathbf{x}) = \sum_{l=3}^4 w_l \ell_{adaptive}^{D,l}(T(\mathbf{x}), \mathbf{x}). \quad (10)$$

We note that the discriminator D is learned from training data. Thus, the network parameters of $\ell_{id_per}(\mathbf{x})$ will be changed along with the updating of discriminator, and thus we name $\ell_{id_per}(\mathbf{x})$ as adaptive perceptual loss. In contrast, conventional perceptual loss [13] is defined on the off-the-shelf CNNs (e.g., VGG-Face [25]). Compared with conventional perceptual loss [13], our adaptive perceptual loss generally is more effective in improving the training efficiency and attribute transfer performance:

- 1) The training efficiency of adaptive perceptual loss can be further explained from two aspects. (i) For conventional perceptual loss, the forward and backward calculations are required for both the off-the-shelf CNN and the

discriminator during training. Due to that the adaptive perceptual loss is defined on the discriminator, it is sufficient to only conduct forward and backward calculation on the discriminator, making our DIAT more efficient in training. (ii) For conventional GAN, the generator usually is difficult to be trained. As for our DIAT, it can be trained by both the adversarial attribute loss and adaptive perceptual loss, greatly accelerating the training speed. Actually, the adaptive perceptual loss is defined on the third and fourth convolution layers of the discriminator, which can serve as some kind of hidden-layer supervision and benefit the convergence of network training [22].

- 2) For conventional perceptual loss, the off-the-shelf CNNs generally are pre-trained using other training data and are not tailored to attribute transfer. One plausible choice is the VGG-Face [25], which, however, is trained for face recognition and may not be suitable for identity-aware attribute transfer. In comparison, our adaptive perceptual loss is defined on the discriminator which is trained for modeling $p_{att}(\mathbf{a})$. Such loss can thus provide natural balance between identity similarity and attribute transfer and benefit transfer performance. For example, in terms of gender transfer, the introduction of adaptive perceptual loss will allow the adaptive adjustment on the length of hair.

Identity Loss. The proposed adaptive perceptual loss does help keep the content similarity between the generated face and the reference. However, it cannot guarantee the identity by itself. To further enhance the identity keeping property, we add constrains on the feature representation extracted for face recognition [58] or verification [59]. In face verification task, two faces are from the same person when the distance between two features are smaller than certain threshold. Here, we adopt VGG-Face and model the distance between the features of the generated face and the reference as the identity loss,

$$\ell_{id}(\mathbf{x}) = \|VGG(x) - VGG(T(\mathbf{x}))\|_2^2. \quad (11)$$

Perceptual regularization. Despite the use of adversarial attribute loss and adaptive perceptual loss, visual artifacts are still inevitable in the transfer result. Image regularization is thus required to encourage the spatial smoothness while preserving small scale details of the generated face $F(\mathbf{x})$. One choice is the Total Variation (TV) regularizer which has been adopted in CNN feature visualization [30] and artistic style transfer [13], [14]. However, the TV regularizer is limited in recovering small-scale texture details and suppressing complex artifacts. Moreover, it is a generic model that does not consider the characteristics of facial images.

In this work, we take the facial characteristics into account and train a denoising network for perceptual regularization. To train the denoising network, we generate the noisy image by adding Gaussian noise with the standard deviation of 15 to the clean facial image from CelebA. Inspired by residual learning [31], we train the denoising network through learning the residual between the noise image and the clean image. Taking the noise image \mathbf{y} as input, the denoising network

utilizes a fully convolutional network of 6 layers to predict the residual $DN(\mathbf{y})$. The denoising result can then be obtained by $\mathbf{y} + DN(\mathbf{y})$. The architecture of $DN(\mathbf{y})$ is listed in Table IV.

TABLE IV
NETWORK ARCHITECTURE OF THE DENOISING NETWORK.

Layer	Activation size
Input	$3 \times 128 \times 128$
conv, $3 \times 3 \times 64$, pad 1, stride 1	$64 \times 128 \times 128$
conv, $3 \times 3 \times 64$, pad 1, stride 1	$64 \times 128 \times 128$
conv, $3 \times 3 \times 64$, pad 1, stride 1	$64 \times 128 \times 128$
conv, $3 \times 3 \times 64$, pad 1, stride 1	$64 \times 128 \times 128$
conv, $3 \times 3 \times 64$, pad 1, stride 1	$64 \times 128 \times 128$
conv, $3 \times 3 \times 3$, pad 1, stride 1	$3 \times 128 \times 128$

Denote by $\mathcal{T} = \{(\mathbf{y}_i, \mathbf{x}_i)\}_{i=1}^n$ a training set, where \mathbf{y}_i denotes the i -th noisy image and \mathbf{x}_i the corresponding clean image. The objective for learning $DN(\mathbf{y})$ is given as,

$$\min \|DN(\mathbf{y}) + \mathbf{y} - \mathbf{x}\|_F^2. \quad (12)$$

Given the denoising network and the transfer result $F(\mathbf{x})$, we define the perceptual regularization as,

$$\ell_{smooth}(F(\mathbf{x})) = \max\{0, \|DN(F(\mathbf{x}))\|_F^2 - t\}, \quad (13)$$

where $\|\cdot\|_F$ denotes the Frobenius norm. Note that $DN(F(\mathbf{x}))$ predicts the residual between the latent clean image and $F(\mathbf{x})$. Minimizing $\|DN(F(\mathbf{x}))\|_F^2$ makes $F(\mathbf{x})$ be close to the clean image, and can be used to suppress the noise and artifacts in $F(\mathbf{x})$. Furthermore, the threshold t is introduced for better preserving of small scale details, and we empirically set t be a value in the range of [16, 30]. Note that the regularizer in Eqn. (13) is defined on the denoising network $DN(F(\mathbf{x}))$, and thus is named as perceptual regularization.

Attribute ratio regularization. The size of attribute relevant region varies for different attributes. For example, the region related to *mouth open/close* mainly includes the mouth and should be small. For *glasses removal*, the attribute relevant region includes the two eyes and is relatively large. As for *gender transfer*, all the face region and the hair should be attribute relevant. Therefore, we introduce an attribute ratio regularization term to constrain the size of attribute relevant region. Specifically, such regularization is defined on the binary mask in Eqn. (2). Denote by N the image size, and p the expected ratio of the region for a specific attribute. The attribute ratio regularization is then defined as,

$$\ell_{mask}(M(\mathbf{x})) = \left(\sum M(\mathbf{x}) - pN\right)^2 \quad (14)$$

In our experiments, we set smaller p value for local attribute and larger p value for global attribute.

Objective function. We define the objective function for learning the transfer model F and the discriminator D by combining the adversarial attribute loss, adaptive perceptual loss, perceptual regularization, and attribute ratio regularization. The transfer model $F(\mathbf{x}) = M(\mathbf{x}) \circ T(\mathbf{x}) + (1 - M(\mathbf{x})) \circ \mathbf{x}$ is learned by minimizing the following objective,

$$\min_{M, T} \ell_F(\mathbf{x}) + \mu \ell_{smooth}(F(\mathbf{x})), \quad (15)$$

where $\ell_F(\mathbf{x}) = \ell_{att, F} + \lambda(\ell_{id}(\mathbf{x}) + \ell_{id_per}(\mathbf{x})) +$

$\gamma \ell_{mask}(M(\mathbf{x}))$. λ , γ , and μ are the tradeoff parameters for the adaptive perceptual loss, attribute ratio regularization, and perceptual regularization, respectively. However, it is difficult to set the tradeoff parameter μ . Instead, we empirically find that the transfer model can be stably learned by alternately minimizing $\ell_F(\mathbf{x})$ and $\ell_{smooth}(F(\mathbf{x}))$ during training. Finally, the discriminator D is learned by minimizing the following objective,

$$\min_D \ell_{att, D}. \quad (16)$$

C. Learning algorithm

Generally, both the generator and the discriminator are difficult to converge in GAN. Therefore, we adopt a two-stage strategy for learning the transfer model $F(\mathbf{x})$ and the discriminator D : (i) we first combine the source set \mathcal{X}_S and the guided set \mathcal{X}_G to pre-train for initialization, and (ii) alternate between updating F and D . The procedure for training the transfer model $F(\mathbf{x})$ is summarized in Algorithm 1.

Initialization. For the initialization of the F , we only consider the attribute transform network T , and leave the mask network M be learned in the second stage. Note that the transform network T has the architecture of auto-encoder. Thus, it can be pre-trained by minimizing the following reconstruction objective on \mathcal{X}_S and \mathcal{X}_G ,

$$\ell_{rec} = \sum_{\mathbf{x} \in \mathcal{X}_S \cup \mathcal{X}_G} \|\mathbf{x} - T(\mathbf{x})\|_F^2. \quad (17)$$

As for the initialization of the discriminator, we use the images in \mathcal{X}_S as negative samples and the images in \mathcal{X}_G as positive samples. Then the discriminator can be pre-trained by minimizing the following objective,

$$\ell_{dis} = \sum_{\mathbf{x}_i \in \mathcal{X}_S \cup \mathcal{X}_G} \|y_i - D(\mathbf{x}_i)\|^2, \quad (18)$$

where y_i is 1 for positive image \mathbf{x}_i and -1 for negative image. By this way, the initialization can provide a good start point and benefit the convergence and stability of DIAT training.

Network training. After the initialization of T and D , network training is further performed by updating the whole F (including both T and M) and D alternately. Moreover, F is updated by first $tstep$ iterations for minimizing ℓ_F and then $nstep$ iterations for minimizing ℓ_{smooth} . We apply the RMSProp solver [60] to train the transfer network F and the discriminator D with a learning rate of 5×10^{-5} .

D. Extension to face hallucination

Besides facial attribute transfer, our DIAT can also be extended to other face editing tasks. Here we use the $8 \times$ face hallucination as an example. Face hallucination is undoubtedly a global transfer task, and thus we remove the mask network M as well as the attribute ratio regularization ℓ_{mask} , making $F(\mathbf{x}) = T(\mathbf{x})$. Moreover, the input image in face hallucination is of low resolution (LR) while the output image is of high resolution (HR). To be consistent with attribute transfer, we super-resolve LR image to the size of HR image with the bicubic interpolator, which is taken as input to $F(\mathbf{x})$.

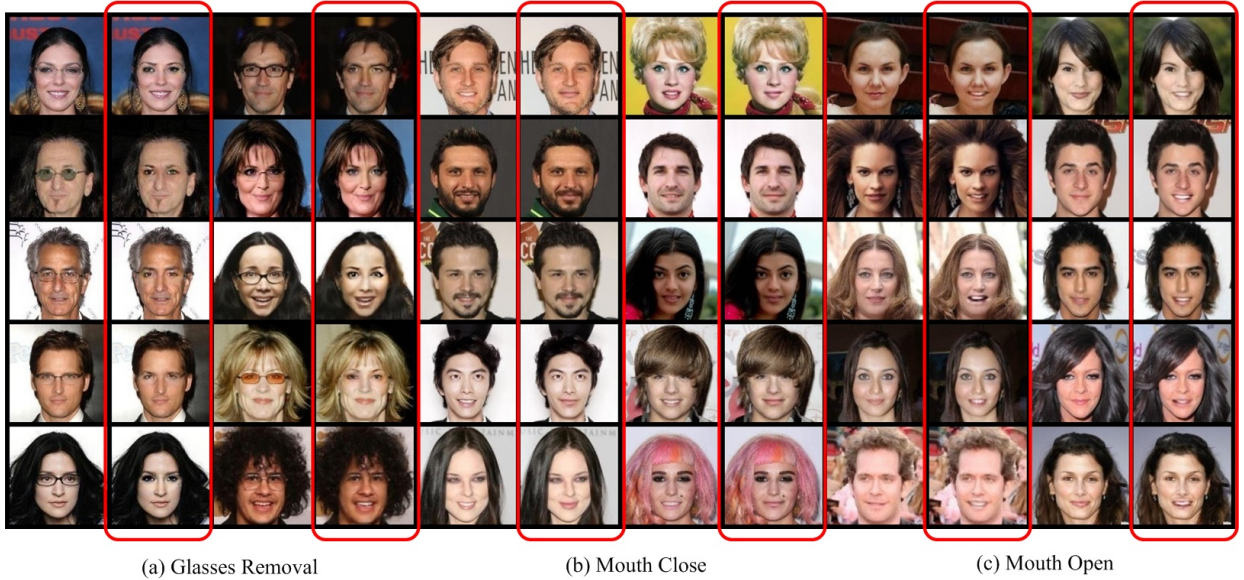


Fig. 3. The results of local attribute transfer. For each task, the left and right columns are the input facial images and the transfer results, respectively.

Algorithm 1 Learning the attribute transfer network

Input: Source set \mathcal{X}_S , guided set \mathcal{X}_G , $dstep = 12$, $tstep = 12$, $nstep = 6$. mini-batch size is 50.

Output: The attribute transfer network F

- 1: Pre-train the transform model T by minimizing the objective in Eqn. (17).
 - 2: Pre-train the discriminator D by minimizing the objective in Eqn. (18).
 - 3: **while** not converged **do**
 - 4: Select a mini-batch \mathbf{X} from \mathcal{X}_S to generate the transfer results, which is further combined with another mini-batch \mathbf{A} from \mathcal{X}_G to form the set \mathbf{X}' for training the discriminator. Here we set $|\mathcal{X}'| = 100$.
 - 5: **for** $i = 1$ to $dstep$ **do**
 - 6: Use the RMSProp solver to update the discriminator D with Eqn. (16) using the mini-batch \mathbf{X}' .
 - 7: **end for**
 - 8: Clip the parameters of the discriminator.
 $D \leftarrow clip(D, -c, c)$
 - 9: **for** $i = 1$ to $tstep$ **do**
 - 10: Use the RMSProp solver to update F by minimizing ℓ_F in Eqn. (15).
 - 11: **end for**
 - 12: **for** $i = 1$ to $nstep$ **do**
 - 13: Use the RMSProp solver to update F by minimizing ℓ_{smooth} in Eqn. (15).
 - 14: **end for**
 - 15: **end while**
 - 16: Return the transfer network F
-

Furthermore, the ground truth HR images can be available to guide the network training for face hallucination. Denote by \mathbf{x} the super-resolved image by bicubic interpolator, and \mathbf{y} the ground truth HR image. Then, the pixel-wise reconstruction loss is defined as,

$$\ell_{rec}(\mathbf{x}) = \|F(\mathbf{x}) - \mathbf{y}\|_F^2. \quad (19)$$

We further modify the definition of ℓ_F by removing the attribute ratio regularization and adding reconstruction loss,

$$\ell_F(\mathbf{x}) = \ell_{att,F} + \lambda \ell_{id}(\mathbf{x}) + \beta \ell_{rec}(\mathbf{x}). \quad (20)$$

where β is the tradeoff parameter for pixel-wise reconstruction loss, and we set $\beta = 0.01$ in our experiment. Given the training data, the models can then be learned by updating F and D alternately.

IV. EXPERIMENTAL RESULTS

In this section, we first describe the experimental settings, including the training and testing data, competing methods, model and learning parameters. Experiments are then performed for local and global attribute transfer. Quantitative metrics and the results on real images are also reported. Moreover, we analyze the effect of adaptive perceptual loss and perceptual regularization. Finally, the results are reported to further assess the performance of our DIAT for face hallucination. The source code will be given after the publication of this work.

A. Experimental settings

Our DIAT models are trained using a subset of the aligned CelebA dataset [21] by removing the images with poor quality. The size of the aligned images is 178×218 . Due to the limitation of the GPU memory, we sample the central part of each image and resize it to 128×128 . For each attribute transfer task, we use all the images with the reference attribute

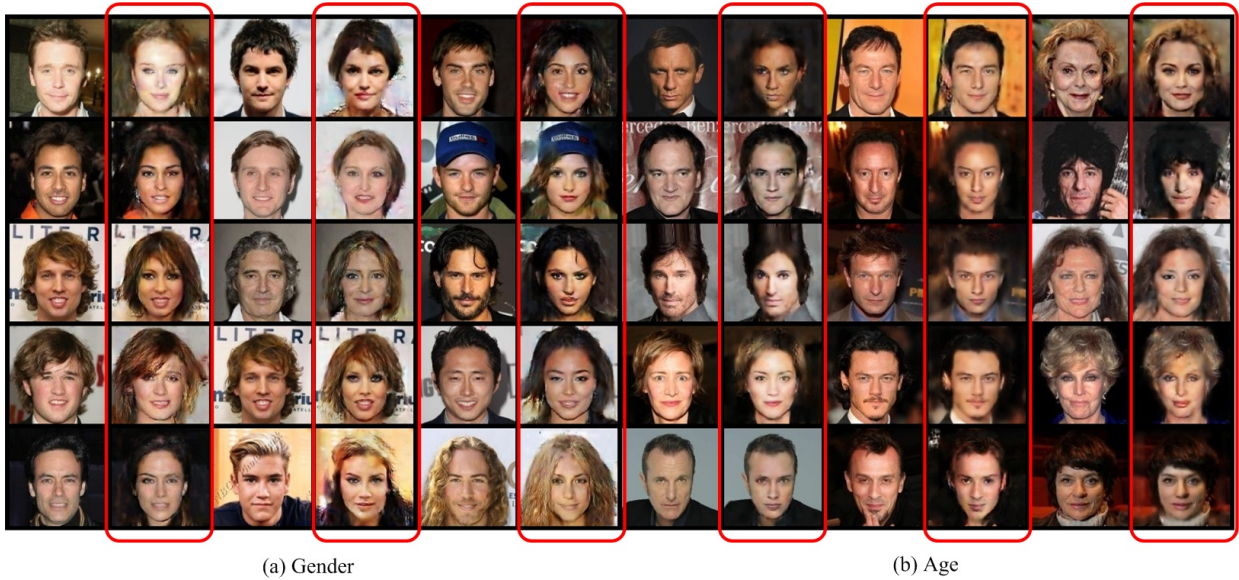


Fig. 4. The results of global attribute transfer. For each task, the left and right columns are the input facial images and the transfer results, respectively.

from training set to form the guided set \mathcal{X}_G , and randomly select 10,000 training images not with the reference attribute as the source set \mathcal{X}_S . After training, 2,000 images apart from the images for training are adopted to assess the attribute transfer performance. And we also test the models on other real images from the website *iStock*.

Only a few methods have been proposed for facial attribute transfer. In our experiments, we compare our DIAT with the convolutional attribute-driven and identity-preserving model (CNIA) [15], IcGAN [17] and VAE/GAN [18] due to that their codes are available. As for IcGAN and VAE/GAN, the original image size is not the same with our DIAT, so we resize the result to the same size with DIAT for comparison. For the task of *glasses removal*, we can first manually detect the region of glasses, and then use some face inpainting methods (e.g., semantic inpainting [51]) to recover the missing pixels. Thus we also compare our DIAT with semantic inpainting [51] for *glasses removal*.

All the experiments are conducted on a computer with the GTX TitanX GPU of 12GB memory. We set the parameters $\lambda = 1$ and $\gamma = 0.01$ for DIAT. For the threshold t in the perceptual regularization ℓ_{smooth} , we set it to be a value in the range of [16, 30]. As for p in the attribute ration regularization, we set it to be (i) $p = 0.16$ for small local attributes (e.g., mouth), (ii) $p = 0.32$ for large local attributes (e.g., eyes), and (iii) $p = 0.62$ for global attributes (e.g., gender and age).

B. Local attribute transfer

We assess the local attribute transfer models on three tasks, i.e., *mouth open*, *mouth close*, and *eyeglasses removal*. Fig. 3 illustrates the transfer results by our DIAT. It can be seen that our DIAT performs favorably for transferring the input images to the desired attribute with satisfying visual quality. Benefited from the mask network, the results by DIAT can preserve more identity-aware and attribute irrelevant details. Moreover, when

the training data are sufficient, it is feasible to separately train two DIAT models for reverse tasks, e.g., one for *mouth open* and another for *mouth close*.

We further compare our DIAT with three competing methods, i.e., CNIA [15], IcGAN [17] and VAE/GAN [18]. As shown in Fig. 6, the results by our DIAT are visually more pleasing than those by CNIA [15] for all the three local attribute transfer tasks. In terms of run time, CNIA takes about 30 seconds (*s*) to deal with an image, while our DIAT only needs 0.0045 *s*. Fig. 5 further provides the results by IcGAN, VAE/GAN, and our DIAT. In comparison with the competing methods, our DIAT can well address the attribute transfer tasks while recovering more visual details in both attribute relevant and attribute irrelevant regions. Finally, for *glasses removal*, we compare our DIAT with semantic inpainting [51], and the results in Fig. 7 clearly demonstrate the superiority of DIAT.

C. Global attribute transfer

We consider two global attribute transfer tasks, i.e., *gender transfer* and *age transfer*. For *gender transfer*, we only evaluate the model for *male-to-female*. For *age transfer*, we only test the model for *older-to-younger*. Fig. 4 shows the transfer results, and our DIAT is also effective for global attribute transfer. Even *gender transfer* certainly causes the change of the identity, as shown in Fig. 4(a), our DIAT can still retain most identity-aware features, making the transfer result similar to the input image in appearance. Figs. 6 and 5 show the results by our DIAT, CNIA [15], IcGAN and VAE/GAN. Compared with the competing methods, the results by our DIAT well exhibit the desired attribute, and are of high visual quality with photo-realistic details. Finally, we also note that for *gender transfer* our DIAT is able of adjusting the hair length due to the introduction of adaptive perceptual loss.

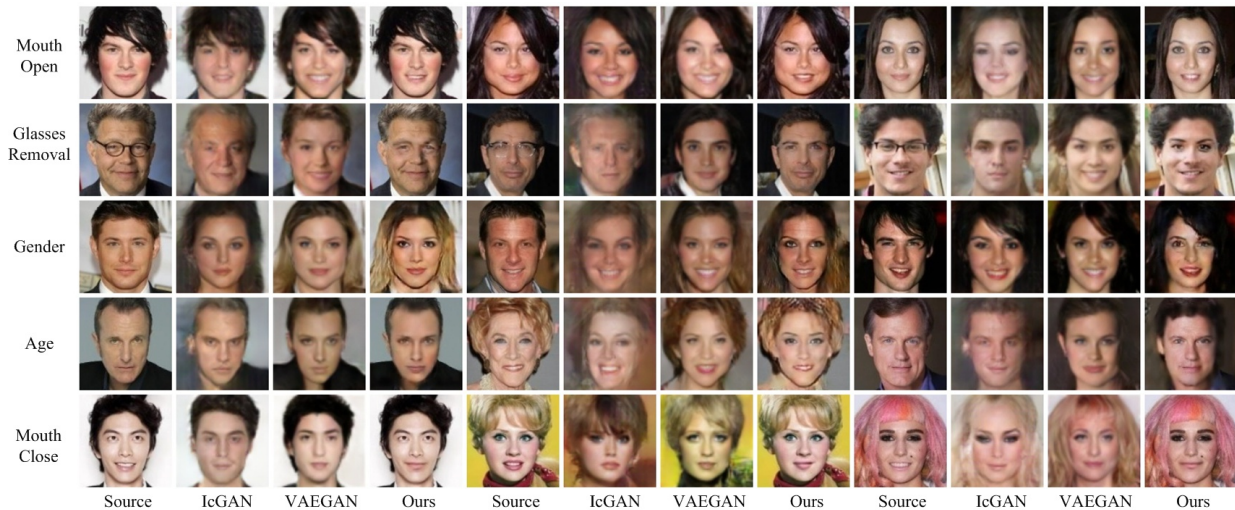


Fig. 5. Comparison of transfer results by our DIAT, IcGAN [17] and VAE/GAN [18].

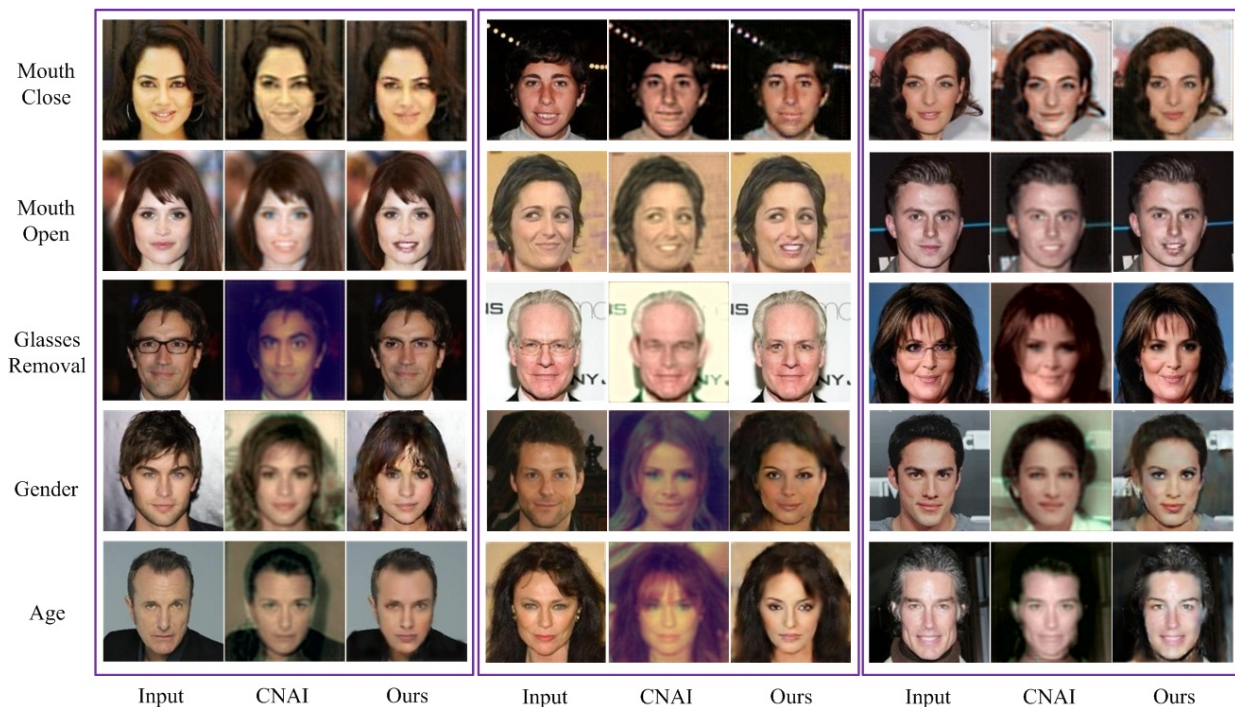


Fig. 6. Comparison of transfer results by our DIAT and CNIA [15].

D. Quantitative evaluation

Given each attribute transfer task, we randomly select 2,000 images without the reference attribute from the testing partition of CelebA to form our testing set. Then, three groups of experiments are conducted to evaluate the transfer performance quantitatively:

- **Attribute classification.** For attribute transfer, it is natural to require the transfer result to exhibit the desired attribute. Thus, we first train a CNN-based attribute classifier (including two convolution layers, three residual blocks and two fully-connected layers) using the training set of CelebA. Given an attribute transfer task, we test

the classification accuracy of the desired attribute for the transfer results of 2,000 testing images. Table V lists the classification accuracy for five attribute transfer tasks, *i.e.*, *mouth open*, *mouth close*, *glasses removal*, *gender transfer*, and *age transfer*. It can be observed that our DIAT achieves satisfying accuracy (*i.e.*, > 0.70) for all the tasks, indicating that the results by our DIAT generally are with the desired attribute.

- **Identity verification.** As for local attribute transfer, we also require the transfer result to preserve the identity of input image. Here we use the open source face

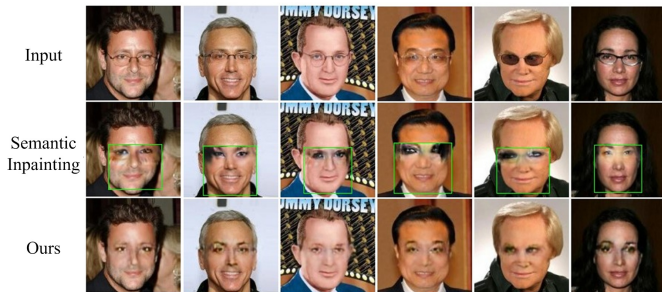


Fig. 7. Comparison of the results by semantic inpainting [51] for *glasses removal*. The green rectangle shows the region of the input subimage for semantic inpainting

recognition platform *Openface*² for matching the input image with the transfer result. By setting the threshold be 0.99, Table VI lists the identity verification accuracy for *mouth open*, *mouth close*, and *glasses removal*. The results demonstrate that our DIAT can well preserve the identity-aware feature for local attribute transfer.

- **Image quality.** Image quality is another crucial metric to assess attribute transfer. However, the ground truth of transfer result is unavailable, making it difficult to perform quantitative evaluation. Here we use a pair of reverse attribute transfer tasks (*i.e.*, *mouth close* and *mouth open*) as an example, and adopt an indirect scheme to compute average PSNR on the 2000 testing images. Specifically, we first perform *mouth open* to the images with *mouth close*, and then perform reverse *mouth close* to the transfer results. Finally, the input images are taken as ground truth and the images after two steps of transfer can be viewed as the generated images. By this way, we obtain the average PSNR of 33.27dB, indicating the effectiveness of our DIAT for local attribute transfer.

TABLE V
ATTRIBUTE CLASSIFICATION ACCURACY FOR TRANSFER RESULTS.

mouth open	mouth close	glasses removal	gender	age
0.821	0.806	0.763	0.684	0.702

TABLE VI
FACE VERIFICATION ACCURACY FOR THE TRANSFER RESULTS.

mouth open	mouth close	glasses removal
0.912	0.903	0.872

E. Results on other real facial images

To assess the generalization ability, we use the DIAT models learned on CelebA to other real facial images from the website *iStock*. Each test image is first aligned with the 5 facial landmarks, and then input to the DIAT models. Taking *mouth open* and *gender transfer* as examples, Fig. 8 gives the transfer results on 15 images for each task, clearly demonstrating the generalization ability of our models to other real facial images.

²<https://github.com/cmusatyalab/openface>

F. Evaluation on adaptive perceptual loss and perceptual regularization

We also implement a variant of DIAT (*i.e.*, DIAT-1) by replacing the adaptive perceptual loss with the conventional perceptual loss defined on VGG-Face [25]. Taking gender transfer as an example, Fig. 9 compares DIAT with DIAT-1. It can be observed that DIAT converges very fast and can generate satisfying results after 4 epochs of training. In comparison, DIAT-1 requires much more epochs in training, and the gender just begins to be modified after 18 epochs. Moreover, the adoption of adaptive perceptual loss also benefits the transfer performance, and adaptive adjustment on the hair length can be observed on the transfer results by DIAT. Furthermore, Fig. 10 shows the transfer results by DIAT with the perceptual regularization and the TV regularization. It can be clearly seen that the perceptual regularization is more effective on suppressing noise and artifacts while preserving sharp edges and fine details.

TABLE VII
COMPARISON OF PSNR (IN DB) FOR FACE HALLUCINATION.

	Bicubic	Unet [?]	DIAT
PSNR	29.68	30.12	28.85
SSIM	0.606	0.672	0.643

G. Results of the learnt mask

Fig. 11 gives the masks generated by the mask network for different task. For the local attribute transformation tasks such as *glasses removal* and *closing mouth*, the generated masks accurately cover the local facial part which is related to the attribute. For global transformation like gender transformation, the mask covers most of the face and keep the background out.

H. Experiments on face hallucination

Finally, we evaluate the performance of DIAT for $8\times$ face hallucination. Table VII lists the average PSNR and SSIM [61] values on the 2,000 testing images by DIAT, bicubic interpolator, and Unet, while Fig. 12 shows the super-resolved images. Even our DIAT achieves lower PSNR/SSIM than the baseline Unet, it is much better in terms of visual quality, and can generate hallucinated image with rich textures and sharp edges.

V. CONCLUSION

A deep identity-aware transfer (*i.e.*, DIAT) model is presented for facial attribute transfer. Considering that some attributes may be only related with parts of facial image, the whole transfer model consists of two subnetworks, *i.e.*, mask network and attribute transform network, which work collaboratively to produce the transfer result. In order to train the model, we further incorporate adversarial attribute loss, adaptive perceptual loss with perceptual regularization and attribute ratio regularization. Experiments show that our model can obtain satisfying results for both local and global attribute

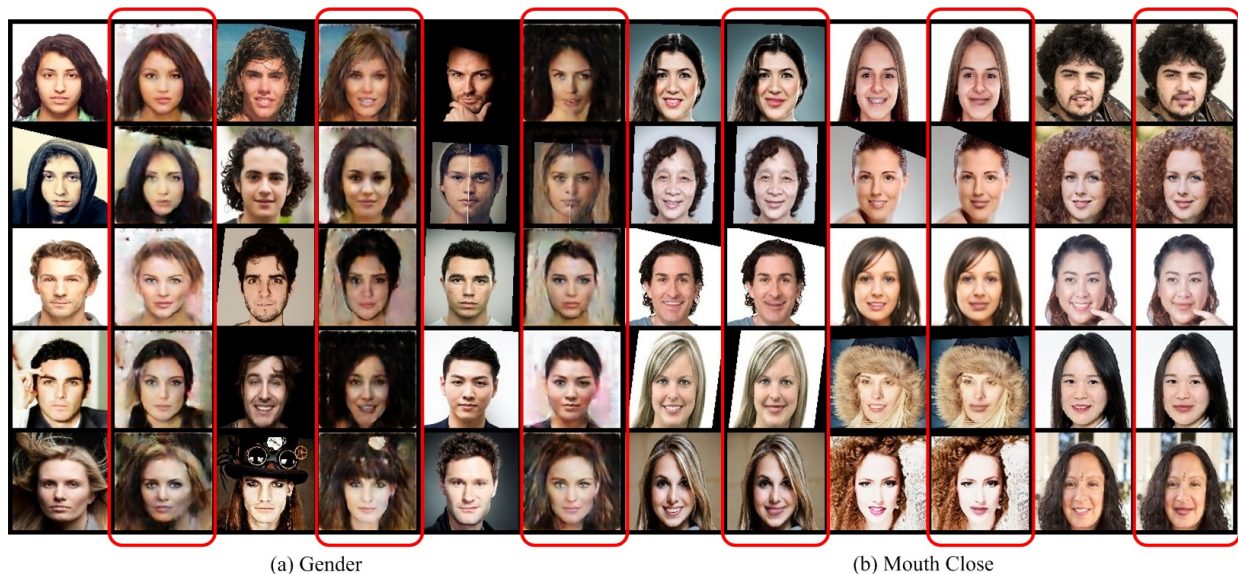


Fig. 8. Local attribute transfer (*mouth open*) and global attribute transfer (*gender transfer*) on images from the website *iStock*. For each task, the left and right columns are the input facial images and the transfer results, respectively.

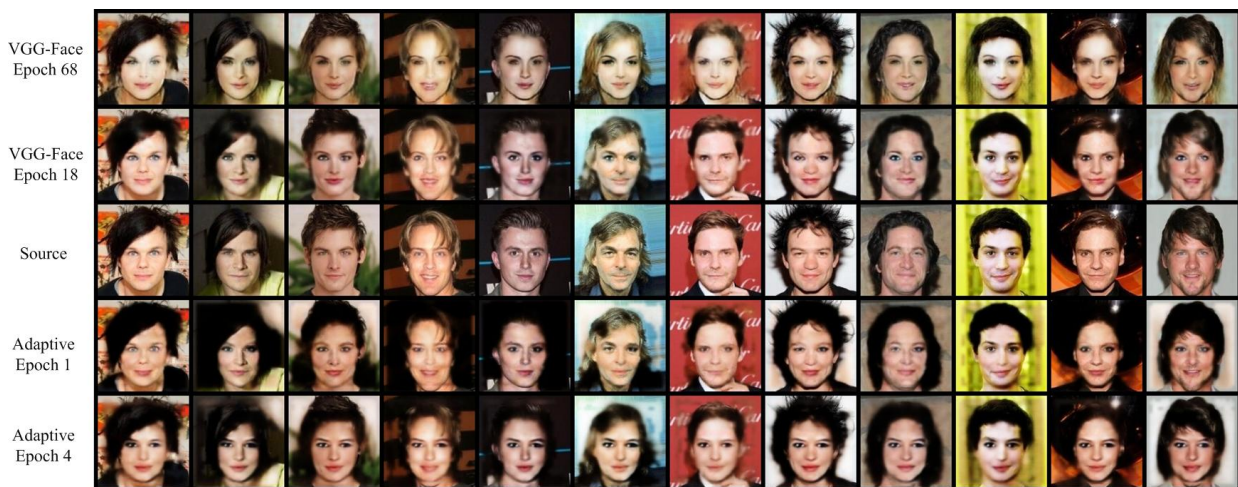


Fig. 9. Comparison between the adaptive perceptual loss and the VGG-Face based perceptual loss.

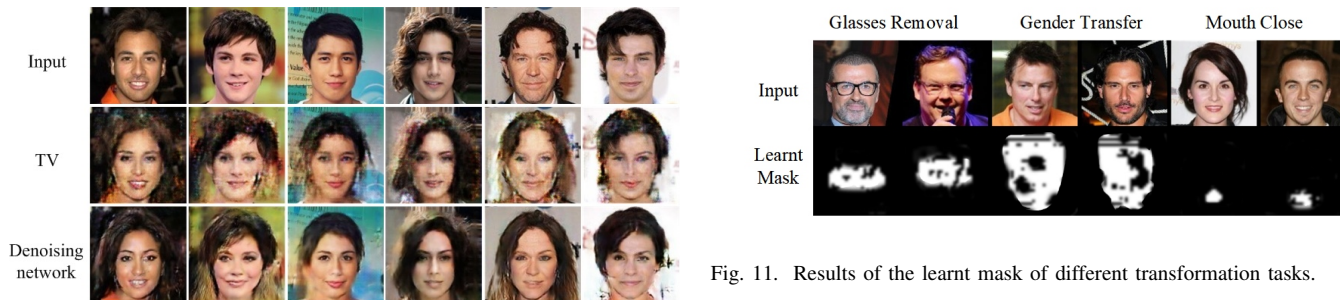


Fig. 10. Comparison between the perceptual regularization and the TV regularization.

transfer. Even for some identity-related attributes (e.g., gender transfer), our DIAT can obtain visually impressive results with



Fig. 11. Results of the learnt mask of different transformation tasks.

minor modification on identity-related features. Our DIAT can also be extended to face hallucination and performs favorably in recovering facial details. In future work, we will further improve the visual quality and diversity of the transfer results, and extend our model to arbitrary attribute transfer.



Fig. 12. Results of face hallucination by different methods.

REFERENCES

- [1] N. Kumar, A. C. Berg, P. N. Belhumeur, and S. K. Nayar, "Attribute and simile classifiers for face verification," in *ICCV*. IEEE, 2009, pp. 365–372.
- [2] A. E. Ichim, S. Bouaziz, and M. Pauly, "Dynamic 3d avatar creation from hand-held video input," *ACM Trans. Graph.*, vol. 34, no. 4, pp. 45:1–45:14, Jul. 2015.
- [3] J. Thies, M. Zollhöfer, M. Stamminger, C. Theobalt, and M. Nießner, "Face2face: Real-time face capture and reenactment of rgb videos," in *Proc. Computer Vision and Pattern Recognition (CVPR)*, IEEE, 2016.
- [4] I. Kemelmacher-Shlizerman, A. Sankar, E. Shechtman, and S. M. Seitz, *Being John Malkovich*, 2010, pp. 341–353.
- [5] Y. Fu, G. Guo, and T. S. Huang, "Age synthesis and estimation via faces: A survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 32, no. 11, pp. 1955–1976, Nov 2010.
- [6] I. Kemelmacher-Shlizerman, S. Suwajanakorn, and S. M. Seitz, "Illumination-aware age progression," in *2014 IEEE Conference on Computer Vision and Pattern Recognition*, June 2014, pp. 3334–3341.
- [7] L. Hu, C. Ma, L. Luo, and H. Li, "Single-view hair modeling using a hairstyle database," *ACM Trans. Graph.*, vol. 34, no. 4, pp. 125:1–125:9, Jul. 2015.
- [8] I. Kemelmacher-Shlizerman, "Transfiguring portraits," *ACM Trans. Graph.*, vol. 35, no. 4, pp. 94:1–94:8, Jul. 2016.
- [9] T. D. Kulkarni, W. F. Whitney, P. Kohli, and J. Tenenbaum, "Deep convolutional inverse graphics network," in *NIPS*, 2015, pp. 2530–2538.
- [10] J. Gauthier, "Conditional generative adversarial nets for convolutional face generation," *Class Project for Stanford CS231N: Convolutional Neural Networks for Visual Recognition, Winter semester*, vol. 2014, 2014.
- [11] X. Yan, J. Yang, K. Sohn, and H. Lee, "Attribute2image: Conditional image generation from visual attributes," in *ECCV*, 2016, pp. 776–791.
- [12] A. van den Oord, N. Kalchbrenner, O. Vinyals, L. Espeholt, A. Graves, and K. Kavukcuoglu, "Conditional image generation with pixelcnn decoders,"
- [13] J. Johnson, A. Alahi, and L. Fei-Fei, "Perceptual losses for real-time style transfer and super-resolution," *arXiv preprint arXiv:1603.08155*, 2016.
- [14] L. A. Gatys, A. S. Ecker, and M. Bethge, "A neural algorithm of artistic style," *arXiv preprint arXiv:1508.06576*, 2015.
- [15] M. Li, W. Zuo, and D. Zhang, "Convolutional network for attribute-driven and identity-preserving human face generation," *arXiv preprint arXiv:1608.06434*, 2016.
- [16] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *NIPS*, 2014, pp. 2672–2680.
- [17] G. Perarnau, J. van de Weijer, B. Raducanu, and J. M. Álvarez, "Invertible Conditional GANs for image editing," in *NIPS Workshop on Adversarial Training*, 2016.
- [18] A. B. L. Larsen, S. K. Sønderby, H. Larochelle, and O. Winther, "Autoencoding beyond pixels using a learned similarity metric," in *Proceedings of the 33rd International Conference on International Conference on Machine Learning*, 2016.
- [19] S. Zhou, T. Xiao, Y. Yang, D. Feng, Q. He, and W. He, "Genegan: Learning object transfiguration and attribute subspace from unpaired data," *arXiv preprint arXiv:1705.04932*, 2017.
- [20] W. Shen and R. Liu, "Learning residual images for face attribute manipulation," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2017, pp. 1225–1233.
- [21] Z. Liu, P. Luo, X. Wang, and X. Tang, "Deep learning face attributes in the wild," in *Proceedings of International Conference on Computer Vision (ICCV)*, 2015.
- [22] C.-Y. Lee, S. Xie, P. Gallagher, Z. Zhang, and Z. Tu, "Deeply-supervised nets," in *AISTATS*, vol. 2, no. 3, 2015, p. 6.
- [23] M. Li, W. Zuo, and D. Zhang, "Deep identity-aware transfer of facial attributes," *arXiv preprint arXiv:1610.05586*, 2016.
- [24] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *NIPS*, 2012, pp. 1097–1105.
- [25] O. M. Parkhi, A. Vedaldi, and A. Zisserman, "Deep face recognition," in *British Machine Vision Conference*, 2015.
- [26] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *CVPR*, 2014, pp. 580–587.
- [27] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. L. Zitnick, and D. Parikh, "Vqa: Visual question answering," in *ICCV*, 2015.
- [28] L. Zhang and P. N. Suganthan, "Visual tracking with convolutional random vector functional link network," *IEEE transactions on cybernetics*, vol. 47, no. 10, pp. 3243–3253, 2017.
- [29] J.-Y. Zhu, P. Krahenbuhl, E. Shechtman, and A. A. Efros, "Learning a discriminative model for the perception of realism in composite images," in *ICCV*, 2015, pp. 3943–3951.
- [30] A. Mahendran and A. Vedaldi, "Understanding deep image representations by inverting them," in *CVPR*. IEEE, 2015, pp. 5188–5196.
- [31] K. Zhang, W. Zuo, Y. Chen, D. Meng, and L. Zhang, "Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising," *IEEE Transactions on Image Processing*, 2017.

- [32] K. H. Jin, M. T. McCann, E. Froustey, and M. Unser, "Deep convolutional neural network for inverse problems in imaging," *IEEE Transactions on Image Processing*, vol. 26, no. 9, pp. 4509–4522, Sept 2017.
- [33] B. Du, W. Xiong, J. Wu, L. Zhang, L. Zhang, and D. Tao, "Stacked convolutional denoising auto-encoders for feature representation," *IEEE transactions on cybernetics*, vol. 47, no. 4, pp. 1017–1027, 2017.
- [34] M. Elad and M. Aharon, "Image denoising via sparse and redundant representations over learned dictionaries," *IEEE Transactions on Image processing*, vol. 15, no. 12, pp. 3736–3745, 2006.
- [35] S. Gu, Q. Xie, D. Meng, W. Zuo, X. Feng, and L. Zhang, "Weighted nuclear norm minimization and its applications to low level vision," *International Journal of Computer Vision*, pp. 1–26, 2016.
- [36] S. C. Zhu, Y. Wu, and D. Mumford, "Filters, random fields and maximum entropy (frame): Towards a unified theory for texture modeling," *International Journal of Computer Vision*, vol. 27, no. 2, pp. 107–126, 1998.
- [37] A. Buades, B. Coll, and J.-M. Morel, "A non-local algorithm for image denoising," in *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, vol. 2. IEEE, 2005, pp. 60–65.
- [38] K. Dabov, A. Foi, V. Katkovnik, and K. Egiazarian, "Image denoising by sparse 3-d transform-domain collaborative filtering," *IEEE Transactions on Image Processing*, vol. 16, no. 8, pp. 2080–2095, Aug 2007.
- [39] A. van den Oord, N. Kalchbrenner, and K. Kavukcuoglu, "Pixel recurrent neural networks," *arXiv preprint arXiv:1601.06759*, 2016.
- [40] C. Dong, C. C. Loy, K. He, and X. Tang, "Learning a deep convolutional network for image super-resolution," in *European Conference on Computer Vision*. Springer, 2014, pp. 184–199.
- [41] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3431–3440.
- [42] A. Dosovitskiy, J. Tobias Springenberg, and T. Brox, "Learning to generate chairs with convolutional neural networks," in *CVPR*, 2015, pp. 1538–1546.
- [43] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.
- [44] K. Simonyan, A. Vedaldi, and A. Zisserman, "Deep inside convolutional networks: Visualising image classification models and saliency maps," *arXiv preprint arXiv:1312.6034*, 2013.
- [45] D. Ulyanov, V. Lebedev, A. Vedaldi, and V. Lempitsky, "Texture networks: Feed-forward synthesis of textures and stylized images," *arXiv preprint arXiv:1603.03417*, 2016.
- [46] E. L. Denton, S. Chintala, R. Fergus *et al.*, "Deep generative image models using a laplacian pyramid of adversarial networks," in *Advances in neural information processing systems*, 2015, pp. 1486–1494.
- [47] A. Radford, L. Metz, and S. Chintala, "Unsupervised representation learning with deep convolutional generative adversarial networks," *arXiv preprint arXiv:1511.06434*, 2015.
- [48] X. Chen, Y. Duan, R. Houthoofd, J. Schulman, I. Sutskever, and P. Abbeel, "Infogan: Interpretable representation learning by information maximizing generative adversarial nets," *arXiv preprint arXiv:1606.03657*, 2016.
- [49] M. Arjovsky, S. Chintala, and L. Bottou, "Wasserstein generative adversarial networks," in *International Conference on Machine Learning*, 2017, pp. 214–223.
- [50] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. Courville, "Improved training of wasserstein gans," *arXiv preprint arXiv:1704.00028*, 2017.
- [51] R. Yeh, C. Chen, T. Y. Lim, M. Hasegawa-Johnson, and M. N. Do, "Semantic image inpainting with perceptual and contextual losses," in *CVPR*, 2017.
- [52] A. Selim, M. Elgharib, and L. Doyle, "Painting style transfer for head portraits using convolutional neural networks," *ACM Transactions on Graphics (TOG)*, vol. 35, no. 4, p. 129, 2016.
- [53] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2015, pp. 234–241.
- [54] W. Shi, J. Caballero, F. Huszár, J. Totz, A. P. Aitken, R. Bishop, D. Rueckert, and Z. Wang, "Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 1874–1883.
- [55] M. Courbariaux, I. Hubara, D. Soudry, R. El-Yaniv, and Y. Bengio, "Binarized neural networks: Training deep neural networks with weights and activations constrained to+ 1 or-1," *arXiv preprint arXiv:1602.02830*, 2016.
- [56] A. Garnaev, M. Baykal-Gursoy, and H. V. Poor, "Security games with unknown adversarial strategies," *IEEE transactions on cybernetics*, vol. 46, no. 10, pp. 2291–2299, 2016.
- [57] F. Zhang, P. P. Chan, B. Biggio, D. S. Yeung, and F. Roli, "Adversarial feature selection against evasion attacks," *IEEE transactions on cybernetics*, vol. 46, no. 3, pp. 766–777, 2016.
- [58] J. Huo, Y. Gao, Y. Shi, W. Yang, and H. Yin, "Heterogeneous face recognition by margin-based cross-modality metric learning," *IEEE transactions on cybernetics*, vol. 48, no. 6, pp. 1814–1826, 2018.
- [59] L. Zheng, S. Duffner, K. Idrissi, C. Garcia, and A. Baskurt, "Pairwise identity verification via linear concentrative metric learning," *IEEE transactions on cybernetics*, vol. 48, no. 1, pp. 324–335, 2018.
- [60] G. Hinton, N. Srivastava, and K. Swersky, "Rmsprop: Divide the gradient by a running average of its recent magnitude," *Neural networks for machine learning, Coursera lecture 6e*, 2012.
- [61] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE Transactions on Image Processing*, vol. 13, no. 4, pp. 600–612, April 2004.