

Deep Identity-aware Transfer of Facial Attributes

Mu Li¹, Wangmeng Zuo², David Zhang¹

¹The Hong Kong Polytechnic University ²Harbin Institute of Technology

csmuli@comp.polyu.edu.hk cswmzuo@gmail.com csdzhang@comp.polyu.edu.hk

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 (i.e., target image) that not only owns the reference attribute but also keep the same or similar identity to the input image. We develop a two-stage scheme to transfer the input image to each reference attribute label. A feed-forward transform network is first trained by combining perceptual identity-aware loss and GAN-based attribute loss, and a face enhancement network is then introduced to improve the visual quality. We further define perceptual identity loss on the convolutional feature maps of the attribute discriminator, resulting in a DIAT-A model. Our DIAT and DIAT-A models can provide a unified solution for several representative facial attribute transfer tasks such as expression transfer, accessory removal, age progression, and gender transfer. The experimental results validate their effectiveness. Even for some identity-related attribute (e.g., gender), our DIAT-A can obtain visually impressive results by changing the attribute while retaining most identity features of the source image.

1. Introduction

Face attributes, e.g., gender and expression, can not only provide a natural description of facial images [24], but also offer a unified viewpoint for understanding many facial animation and manipulation tasks. For example, the goal of facial avatar [15] and reenactment [35] 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 of the source image. But for some other tasks, the modification of some attributes, e.g., gender and ethnicity, will inevitably cause the change of the identity.

In recent years, a variety of methods have been developed for specific facial attribute transfer tasks, and have achieved impressive results. For expression transfer, ap-

proaches have been suggested to create 3D or image-based avatars from hand-held video [15], while face trackers and expression modeling have been investigated for offline and online facial reenactment [18, 35]. For age progression, explicit and implicit synthesis methods have been proposed for different image models [8, 19]. Hair style generation and replacement have also been studied in literature [14, 17].

CNN-based models have also been investigated for human face generation with attributes. Kulkarni et al. [23] proposed 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] developed a conditional generative adversary network (cGAN) to generate facial image from a noise distribution and conditional attributes. Yan et al. [39] developed an attribute-conditioned deep variational auto-encoder which extracts the latent variables from a reference image and combines them with attributes to generate a target image with a generative model. Oord et al. [38] proposed 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, 39, 38], making them not directly applicable to attribute transfer. Li et al. [26] suggested a CNN-based attribute transfer model from the optimization perspective, but both running time and transfer quality is far from satisfactory.

Motivated by the strong capability of CNN in modeling complex transformation [16] and capturing perceptual similarity [9], we present a two-stage Deep CNN model for Identity-Aware Transfer (DIAT) of facial attributes. Instead of developing individual solution to specific task, this paper aims to suggest an unified solution to several facial animation and manipulation tasks that can be viewed as attribute transfer problems. For training data, we only consider the binary attribute labels presented in the large-scale CelebFaces Attributes (CelebA) dataset [27]. For each reference attribute label, we train a two-stage CNN for the transform and enhancement of the image to the desired attribute. Therefore, our DIAT model can provide a general

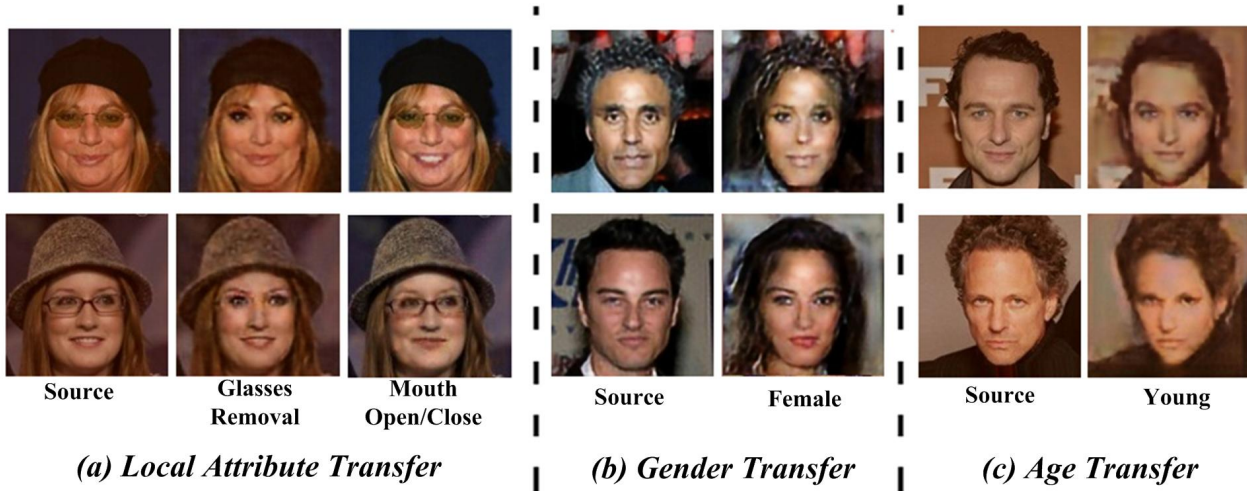


Figure 1. Illustration of our DIAT model for different attribute transfer tasks.

framework for solving several animation and manipulation tasks, such as expression transfer, accessory removal, age progression, and gender transfer.

In the first stage, considering the ground truth transfer results generally are missing, we train the transform network by combining three alternative loss functions, i.e., identity-aware loss, attribute loss, and regularization term. For identity-aware transfer, our DIAT requires that the generated image (i.e., target image) should keep the same or similar identity to the source image. Thus, the identity-aware loss is defined on the convolutional feature maps of a deep CNN model tailored to face recognition, i.e., VGG-Face [30]. For attribute transfer, we require that the target image owns the desired attribute label. Note that the within-class variation is very high and complex for each attribute label. To capture the convolutional feature distribution of each attribute, we construct an attribute guided set for each attribute using all the images with the attribute in CelebA. Then we follow the generative adversarial network (GAN) framework [12] to define the attribute loss as a minimax game played by the transform network and attribute discriminator. The regularization term is also based on perceptual loss to suppress noise and artifacts of the transformed image.

In the second stage, we learn an enhancement network to improve visual quality and remove artifacts of the transformed image. For local attribute transfer, most parts of the input image can be preserved in the transfer result. Thus, we exploit both the high quality input image and the reasonable transformed result to train the enhancement network. For the global attributes, we employ an indirect manner by training a deblurring network for enhancement.

Furthermore, we propose an improved DIAT (i.e., DIAT-

A) model, where the perceptual identity loss in the first stage is defined on the convolutional feature maps of the attribute discriminator. By this way, we can learn an adaptive perceptual identity loss tailored to the specific attribute transfer task. Moreover, adaptive perceptual identity loss can also serve as a kind of hidden-layer supervision or regularization to the discriminator. This is also helpful in improving the convergence performance for GAN training.

Experiments are conducted on CelebA and other real images to evaluate the proposed method. As shown 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 results while retaining most identity-related features of the source image. Computational efficiency is also a prominent advantage of our method. In the test stage, our DIAT can process more than one hundred of images within one second. In summary, the key contributions of this work are:

1. A unified two-stage CNN framework (i.e., transform and enhancement) is suggested which can be used to different attribute transfer tasks. Due to the unavailability of ground truth results in training, we combine perceptual identity loss and GAN-based attribute loss to train the transform network. Then, an enhancement network is deployed to make the transfer results visually more pleasant.
2. An improved DIAT (i.e., DIAT-A) model is proposed by defining perceptual identity loss on the convolutional feature maps of the attribute discriminator. DIAT-A not only allows us learn adaptive perceptual identity loss tailored to specific attribute transfer task,

but also can greatly improve the training efficiency.

3. Experimental results validate the feasibility and efficiency of our method for 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 within less than 0.01s per image.

2. Related work

Deep convolutional networks (CNNs) have shown great success in versatile high level vision problems [22, 30, 11, 2], and also exhibited their remarkable power in understanding and generating images at a fine level of granularity [42]. In this section, we focus on the later, and briefly survey the related CNN models for image generation and face generation.

2.1. CNN for image generation

Generative image modeling is a critical issue for image generation and many low level vision problems. Conventional sparse [1], low rank [13], FRAME [43] and non-local similarity [3] 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 [23], PixelCNN and PixelRNN [37], and they have been applied to image completion and generation. In [37], PixelRNN is evaluated for generating small size images (e.g., 32×32 and 64×64).

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 [6, 28]. The generative CNN model [7] stacks four convolution layers upon five fully connected layers to generate images from object description. Kulkarni et al. suggested the Deep Convolution Inverse Graphics Network (DC-IGN), which follows the variational autoencoder architecture [21] to transform the input image into different pose and lighting condition. However, both generative CNN [7] and DC-IGN [21] require many labeled images to train the CNNs.

To visualize and understand CNN features, several methods have been proposed to reconstruct images that invert deep representation [29] and maximize class score [34]. Subsequently, Gatys et al. [9] introduce the combination of content and style losses defined on deep representation on the off-the-shelf CNNs for artistic style transfer. Li et al. further extend this framework for identity-preserving attribute transfer [26]. However, these methods generally require high computational cost and cannot obtain high quality results for attribute transfer. To improve the efficiency, alternative approaches have been proposed by substituting

the iterative optimization procedure with pre-trained feed-forward CNN [16, 36]. And a number of perceptual losses have also been proposed for style transfer and other generation tasks [16]. Motivated by these works, both perceptual identity loss and perceptual regularization have been exploited in our DIAT model by considering the requirement of attribute transfer.

Another representative approach to train generative CNN is the GAN framework, where a discriminator and a generator are alternatively trained as an adversarial game [12]. 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 [5]. Radford et al. [31] extended GAN with the fully deep convolutional networks (i.e., DCGAN) and applied it to learn a hierarchical representations from object parts. 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 [4]. GAN-based loss has also been adopted in several image generation tasks, such as text-to-image synthesis [32]. Compared with the above models, we combine both perceptual loss and GAN-based loss for identity-aware attribute transfer. We also introduce an enhancement network for improving visual quality of images obtained by GAN-based image generation. Furthermore, we propose to define adaptive perceptual identity loss on the discriminator to improve the efficiency of training GAN.

2.2. CNN for face generation

Besides face synthesis with attributes, several CNN models have also been developed for other face generation tasks. For painting style transfer of head portrait, Selim et al. [33] modified the content loss to balance the contribution of the input photograph and the aligned example painting. Gucluturk et al. trained a feed-forward CNN with perceptual loss for sketch inversion. The most related work is [40], where DCGAN is applied to semantic face inpainting. Different from [40], our DIAT adopts a two-stage CNN solution, and we propose a DIAT-A model by incorporating adaptive perceptual identity loss with GAN-based attribute loss, making the visual quality by our model much better than that by the related methods [26, 40].

3. Deep CNNs for Identity-aware Attribute Transfer

In this section, a two-stage scheme is developed to tackle the identity-aware attribute transfer task. Fig. 3 illustrates the training and testing procedure of our DIAT method. In the first training stage, GAN-based loss and perceptual

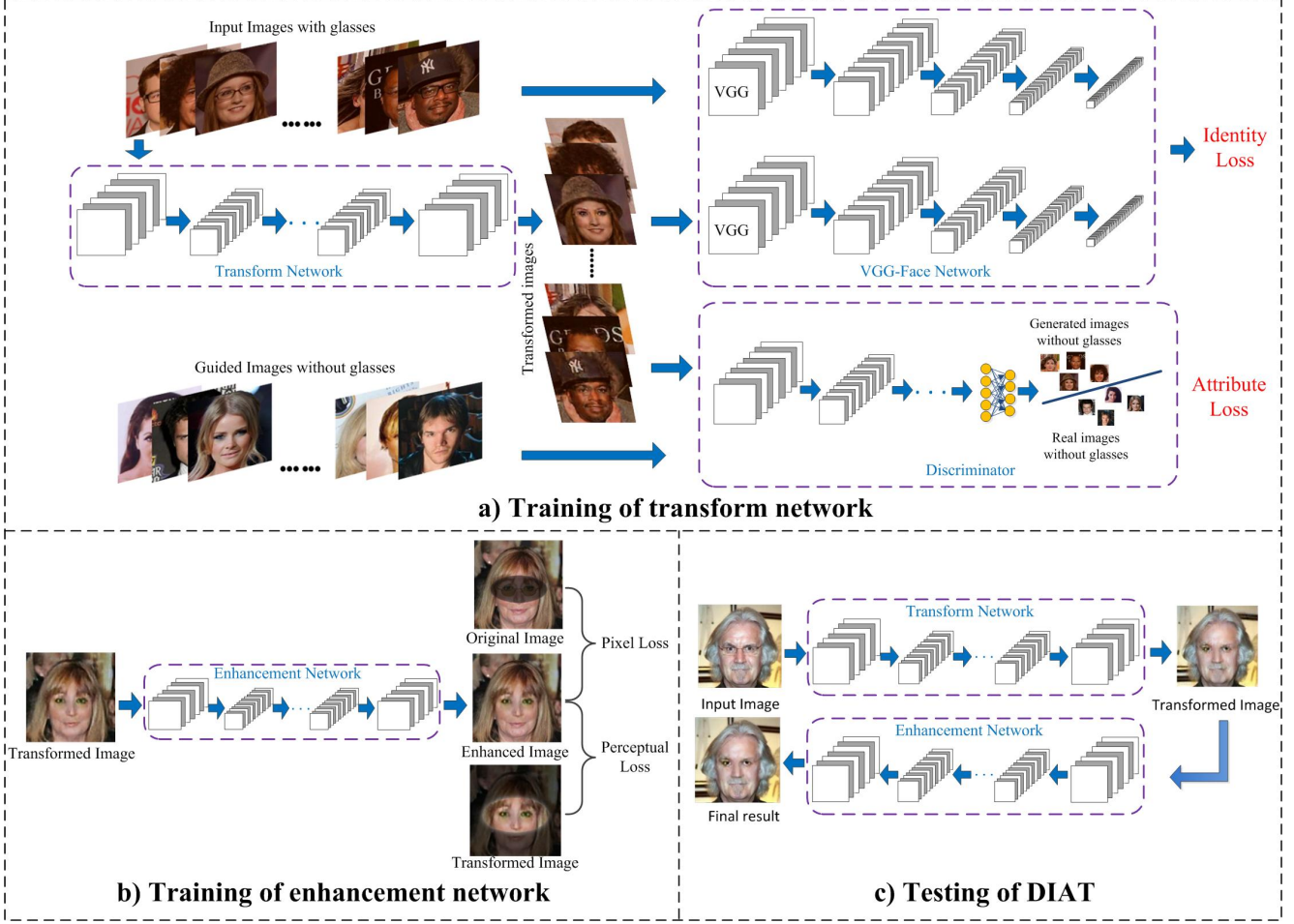


Figure 2. Illustration of our model (Using glasses removal as an example).

identity loss are combined to learn the face transform network for attribute transfer. Note that both identity and attribute are high-level semantic concepts that cannot be defined in the pixel domain, limiting the image quality by the transform network. Thus, in the second training stage, a face enhancement network is trained to further improve visual and color quality of the transformed image. In the test stage, we stack the enhancement network upon the transform network for attribute transfer.

3.1. Face transform network

We adopt a 16-layer fully convolutional network for face transform, which takes use of 3 convolutional layers followed by 5 residual blocks and then another 3 convolutional layers. Each residual block has two convolutional layers. The details of the face transform network are summarized in Table 1. To train the transform network, the groundtruth image to the input image generally is unknown and not unique. In this work, the alternative losses are adopted for model training, which only require a face dataset with the anno-

tated attributes. In this subsection, we define the perceptual identity loss, GAN-based attribute loss and perceptual regularization, and then introduce the learning algorithm for training the face transform network.

| Layer | Activation size |
|---|----------------------------|
| Input | $3 \times 128 \times 128$ |
| $9 \times 9 \times 32$ conv, pad 4, stride 1 | $32 \times 128 \times 128$ |
| $3 \times 3 \times 64$ conv, pad 1, stride 2 | $64 \times 64 \times 64$ |
| $3 \times 3 \times 128$ conv, pad 1, stride 2 | $128 \times 32 \times 32$ |
| Residualblock, 128 filters | $128 \times 32 \times 32$ |
| Residualblock, 128 filters | $128 \times 32 \times 32$ |
| Residualblock, 128 filters | $128 \times 32 \times 32$ |
| Residualblock, 128 filters | $128 \times 32 \times 32$ |
| $3 \times 3 \times 64$ deconv, pad 1, stride 2 | $64 \times 64 \times 64$ |
| $3 \times 3 \times 32$ deconv, pad 1, stride 2 | $32 \times 127 \times 127$ |
| $10 \times 10 \times 3$ deconv, pad 4, stride 1 | $3 \times 128 \times 128$ |

Table 1. Architecture of the face transform network.

3.1.1 Model objective

The goal of this work is to train a feed-forward residual network to transform source image to the desired attribute while keeping the same or similar identity to the source image. Due to the groundtruth output usually is unavailable, our training data includes a guided set of images that are with the desired reference attribute and a set of input images that are not with the reference attribute. Therefore, we define the identity loss and attribute loss by referring to the input image and guided set, respectively.

Identity loss. The identity loss is introduced to guarantee that the generated image has 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 loss by forcing two images to be exactly the same in pixel domain. In our model, we define the squared-error loss on CNN feature representations, resulting in our perceptual identity loss.

For modeling the perceptual identity loss, we take use of the VGG-Face network [30], which is trained on a very large scale face dataset and has been shown to yield excellent performance on face recognition. There are about 40 layers in VGG-Face. We only exploit the first 5 convolutional layers due to their better performance on image reconstruction. Denote by ϕ the VGG-Face network, $\phi_l(\mathbf{x})$ the feature map of the l th convolutional layer with respect to the input image \mathbf{x} . C_l , H_l and W_l represent the channels, height and width of the feature map, respectively. The perceptual loss between two images \mathbf{x} and $\hat{\mathbf{x}}$ on the l th convolutional layer is defined as,

$$\ell_{content}^{\phi,l}(\hat{\mathbf{x}}, \mathbf{x}) = \frac{1}{2C_l H_l W_l} \|\phi_l(\hat{\mathbf{x}}) - \phi_l(\mathbf{x})\|_F^2. \quad (1)$$

Denote by $T(\mathbf{x})$ the output of transform network. The perceptual identity loss is defined as a weighted perceptual loss based on the 4th and 5th convolutional layers of VGG-Face,

$$\ell_{id}(\mathbf{x}) = \sum_{l=4}^5 w_l \ell_{content}^{\phi,l}(T(\mathbf{x}), \mathbf{x}), \quad (2)$$

where w_l is the weight for the l th convolutional layer.

Attribute loss. The attribute loss is introduced to make the generated image have the desired attribute. One natural choice is to adopt the similar way with [7] to define the loss based on the ground truth image with the reference attribute while preserving the identity. Unfortunately, the ground truth image usually is unavailable and not unique for attribute transfer. As mentioned above, what we have is a set of images with the reference attributes, i.e., attribute guided set. The images in the guided set do not require to have the same identity with the input images. For the attribute transfer task, we adopt the CelebA dataset to form the guided set.

The attribute guided set \mathcal{A} provides a natural representation of the attribute distribution. With the attribute loss, our goal is to make that the distribution of transformed images matches the real attribute distribution. To this end, we adopt the generative adversarial network (GAN) framework to train the face transform network by a minimax game. In general, GAN includes a generator and a discriminator. For attribute transfer, the generator is the face transform network, which is used to transform the input image to an image with the desired attribute. The discriminator is used to distinguish the generated images from the real images in the guided set, which contains 6 convolutional layers followed by another two fully connected neural layers. The details of the discriminator is shown in Table 2.

| Layer | Activation size |
|---|---------------------------|
| Input | $3 \times 128 \times 128$ |
| $8 \times 8 \times 32$ conv, pad 3, stride 2 | $32 \times 64 \times 64$ |
| $3 \times 3 \times 32$ conv, pad 1, stride 1 | $32 \times 64 \times 64$ |
| $4 \times 4 \times 64$ conv, pad 1, stride 2 | $64 \times 32 \times 32$ |
| $3 \times 3 \times 64$ conv, pad 1, stride 1 | $64 \times 32 \times 32$ |
| $4 \times 4 \times 128$ conv, pad 1, stride 2 | $128 \times 16 \times 16$ |
| $4 \times 4 \times 128$ conv, 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 |

Table 2. Network structure used for the discriminator.

Let $p_{data}(\mathbf{x})$ be the distribution of the input images, $p_{att}(\mathbf{a})$ be the distribution of the images in the guided set. Denoted by $T(\mathbf{x})$ the transform network is used to transform the input images to the given attribute. The discriminator is defined as $D(\mathbf{a})$ to output the possibility the image \mathbf{a} comes from the set \mathcal{A} rather than the transformed ones. To train the generator and the discriminator, we take use of the following attribute loss:

$$\min_T \max_D \mathbb{E}_{\mathbf{a} \sim p_{att}(\mathbf{a})} [\log D(\mathbf{a})] + \mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})} [1 - \log D(T(\mathbf{x}))] \quad (3)$$

Perceptual regularization. The regularization term is introduced to encourage the spatial smoothness while preserving small scale details of the generated face $T(\mathbf{x})$. One choice is the Total Variation (TV) regulariser $\ell_{TV} T(x)$ which has been adopted in CNN feature visualization [29] and artistic style transfer [9, 16]. Despite its effectiveness in reducing noise while preserving strong edges, the TV regularizer is limited in recovering small-scale texture details. 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 simulate the possible artifacts of the transform network and identity loss to form a perceptual regularizer. To simulate the possible artifacts, we train a reconstruction network $g(\mathbf{x})$ which has the same architecture as the transform network. Denoted by ϕ the VGG-Face network [30], $\phi_l(\mathbf{x})$ the feature map of the l th convolutional layer with respect to the input image \mathbf{x} . By referring to the perceptual

identity loss, $g(\mathbf{x})$ is learned by minimizing the combined perceptual loss defined on the 4th and 5th layers,

$$\min_g \sum_{l=4}^5 w_l \ell_{content}^{\phi,l}(g(\mathbf{x}), \mathbf{x}), \quad (4)$$

where w_l is the weight for the l th convolutional layer.

We note that the reconstruction network $g(\mathbf{x})$ has the same architecture with the transform network. Therefore, the difference $\mathbf{n} = g(\mathbf{x}) - \mathbf{x}$ can also be viewed as an synthesis of the possible artifacts caused by the transform network and perceptual identity loss. Then we train a denoising network $f(g(\mathbf{x}))$ to separate \mathbf{n} from $g(\mathbf{x})$. The denoising network has a simple structure of only two convolutional layers with 3×3 kernel, and it is trained by,

$$\min_f \|f(g(\mathbf{x})) - \mathbf{x}\|_F^2 + \|f(\mathbf{x}) - \mathbf{x}\|_F^2. \quad (5)$$

The term $\|f(\mathbf{x}) - \mathbf{x}\|_F^2$ is introduced to prevent the denoising network from over-smoothing a clear image.

With the denoising network, we design a perceptual regularization term which requires the transformed images to be close to the result of denoising network.

$$\ell_{smooth}(T(\mathbf{x})) = \|f(T(\mathbf{x})) - T(\mathbf{x})\|_F^2. \quad (6)$$

Objective function. Taking the perceptual identity loss, GAN-based attribute loss, and perceptual regularizer into account, the final objective function for attribute transform network is defined as,

$$\min_T \max_D \mathbb{E}_{\mathbf{a} \sim p_{att}(\mathbf{a})} [\log D(\mathbf{a})] + \mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})} [1 - \log D(T(\mathbf{x})) + \lambda \ell_{id}(\mathbf{x}) + \gamma \ell_{smooth}(\mathbf{x})] \quad (7)$$

where γ and λ are the tradeoff parameters for the smooth regularizer and identity loss, respectively.

3.1.2 Learning algorithm

The procedure for training the transform network is summarized in Algorithm 1. The guided set \mathcal{A} with reference attribute is first extracted from CelebA. After pre-training, we alternate between updating the transform network and updating the discriminator.

Initialization. We pre-train the transform network using the whole set of the CelebA dataset \mathcal{D} with about 200,000 face images [27]. Note that the transform network has the structure of auto-encoder. Thus, it can be pre-trained by minimizing the following reconstruction loss,

$$\ell_{rec} = \sum_{\mathbf{x} \in \mathcal{D}} \|\mathbf{x} - T(\mathbf{x})\|_F^2. \quad (8)$$

We also use CelebA to pre-train the discriminator. For given attribute, we use the images with the reference attribute as positive samples, and those with other attributes

Algorithm 1 Training of the transform network

Input: CelebA face image set \mathcal{D} , Attribute label set \mathbf{Y} , Given attribute att .

Output: The transform network T

- 1: Extract all the images with the given attribute att from \mathcal{D} to obtain the guided set \mathcal{A} , and extract a set of images with the other attributes to obtain the input set \mathcal{X} .
 - 2: Pre-train the transform network with the loss in Eqn. (8).
 - 3: Pre-train the discriminator with the loss in Eqn. (9).
 - 4: **while** not converged **do**
 - 5: Generate face images with the transform network from the input image set \mathcal{X} . Combine the transformed images and the guided set \mathcal{A} to get the set \mathcal{X}' for training discriminator.
 - 6: **for** $i = 1$ to $dstep$ **do**
 - 7: Use the ADAM solver to update the discriminator D with Eqn. (7) on the training set \mathcal{X}' .
 - 8: **end for**
 - 9: **for** $i = 1$ to $tstep$ **do**
 - 10: Use the Adam solver to update the transform network T with Eqn. (7).
 - 11: **end for**
 - 12: **end while**
 - 13: Return the transform network T
-

as the negative samples to train the discriminator. Then the discriminator is pre-trained by minimizing the following loss,

$$\ell_{dis} = \sum_{\mathbf{x}_i \in \mathcal{D}} \|y_i - D(\mathbf{x}_i)\|^2 \quad (9)$$

where y_i represents the attribute label of the image \mathbf{x}_i . Note that the transform network is pre-trained to approximate the input (i.e., image without reference attribute), and the discriminator is pre-trained to distinguish the images in the guided set \mathcal{A} from the other images. Thus, our initialization scheme can guarantee the synchronization of transform network and discriminator, and provide a good start point for our DIAT.

Network training. With the pre-trained models, the optimization of the transform network and the discriminator is alternated to minimize the objective in Eqn. (7). Here we apply the ADAM solver [20] to train the transform network and the discriminator with a learning rate of 0.0001.

3.2. Face enhancement networks

Both identity and attribute losses are defined on high-level feature representation, and the training of GAN is difficult to converge. All these factors may result in poor visual quality of the generated images. To remedy this, we further train a face enhancement network $E(\mathbf{x})$ to make the transform result visually more pleasant.

For local attribute transfer, only part of the facial image is required to be modified, while the other parts should keep unchanged with the input face image. Thus we can utilize the high quality input image to train the enhancement network. For each local attribute, we introduce a mask image to denote the unchanged region. Specifically, a facial image is first aligned based the 68 facial landmarks obtained using the landmark detection method in [41]. Then, each local attribute is associated with certain landmarks and a convex hull is generated based on these landmarks. By expanding the convex hull with a certain margin, we define the attribute mask \mathbf{m} by setting the value of 1 to the pixels in the convex hull and 0 to the others.

Based on the attribute mask \mathbf{m} , the loss function of face enhancement network consists of two terms: (i) for the unchanged region, we require the enhanced image should be close to the input image in pixel domain; (ii) for the attribute related region, we require the enhanced image should be close to the generated image $T(\mathbf{x})$ in CNN feature domain of VGG-Face. Then the local enhancement network can be trained by minimizing the following loss,

$$\begin{aligned} \ell_{E_{local}} = & (1 - \mathbf{m}) \circ \|(E(T(\mathbf{x})) - \mathbf{x})\|_F^2 \\ & + \sum_{i=0}^2 \beta_i \ell_{content}^{\phi,i}(\mathbf{m} \circ T(\mathbf{x}), \mathbf{m} \circ E(T(\mathbf{x}))) \end{aligned} \quad (10)$$

For global attribute, it is difficult to define the unchanged region. Here we adopt an indirect way to train the global enhancement network. Despite the complexity of the noise and artifacts caused by the transform network, one can suppress them effectively by blurring the generated images with Gaussian kernel. To enhance the blurry images, we then train a deblurring network on a set of high quality facial images. Thus, we train the enhancement network for global attribute transfer by minimizing the following loss,

$$\ell_{E_{global}} = \|E(B(\mathbf{x})) - \mathbf{x}\|_F^2 \quad (11)$$

where $B(\mathbf{x})$ denotes the blurring result of \mathbf{x} . After training, we can use $E(B(T(\mathbf{x})))$ for the enhancement of the transformed images.

Different network architectures are adopted for local attribute enhancement and global attribute enhancement. For local attribute, the input to the network includes both the source input image and the transformed image, and take use of a fully convolutional network (FCN) of 4 layers. For global attribute, the input only includes the transformed image, and we adopt a more complex network, i.e., the same architecture of the transform network. The ADAM solver [20] is adopted to train the enhancement networks, and we set the learning rate as 0.0001.

4. DIAT with Adaptive Perceptual Identity Loss

In Sec. 3.1.1, the perceptual identity loss is defined on the pre-trained VGG-Face. Actually, it may be more effective to define this loss on some CNN trained to attribute transfer. On the other hand, the GAN-based attribute loss is difficult to be optimized, and the incorporation of hidden-layer supervision may be helpful in improving the convergence performance. Thus, we suggest an improved DIAT model, where in the first training stage we incorporate adaptive perceptual identity loss into the discriminator, resulting in a DIAT-A model.

4.1. Adaptive perceptual identity loss

Feature sharing has shown to be effective for multi-task learning. Here we treat identity-preserving and attribute transfer as two related tasks, and define the perceptual identity loss based on the convolutional features of the discriminator. By this way, the network parameters for identity loss will be changed along with the updating of discriminator, and thus we named it as adaptive perceptual identity loss.

Denoted by D the discriminator, $D_l(\mathbf{x})$ is the feature map of the l th convolutional layer. C_l , H_l and W_l represent the channels, height, and width of the feature map, respectively. The semantic perceptual loss between two images \mathbf{x} and $\hat{\mathbf{x}}$ on the l th convolutional layer is then defined as,

$$\ell_{dynamic}^{D,l}(\hat{\mathbf{x}}, \mathbf{x}) = \frac{1}{2C_l H_l W_l} \|D_l(\hat{\mathbf{x}}) - D_l(\mathbf{x})\|_F^2. \quad (12)$$

The adaptive perceptual identity loss is further defined as,

$$\ell_{id}^{dynamic}(\mathbf{x}) = \sum_{l=4}^5 w_l \ell_{dynamic}^{D,l}(T(\mathbf{x}), \mathbf{x}). \quad (13)$$

4.2. DIAT-A

The introduction of adaptive perceptual identity loss is also helpful in generating clear image, and the perceptual regularization term can be removed from the model. Thus, in DIAT-A, the objective function for training transform network is given by,

$$\begin{aligned} \min_T \max_D \quad & \mathbb{E}_{\mathbf{a} \sim p_{att}(\mathbf{a})} [\log D(\mathbf{a})] + \mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})} [1 - \\ & \log D(T(\mathbf{x})) + \lambda \ell_{id}^{dynamic}(\mathbf{x})] \end{aligned} \quad (14)$$

The adaptive perceptual identity loss can also be treated as some kind of hidden-layer supervision to the 4th and 5th layers. As shown in [25], hidden-layer supervision can result in evident performance boost in convergence. For attribute transfer, adaptive perceptual identity loss indeed relieves the difficulty of training GAN, and is also effective in suppressing color distortion.

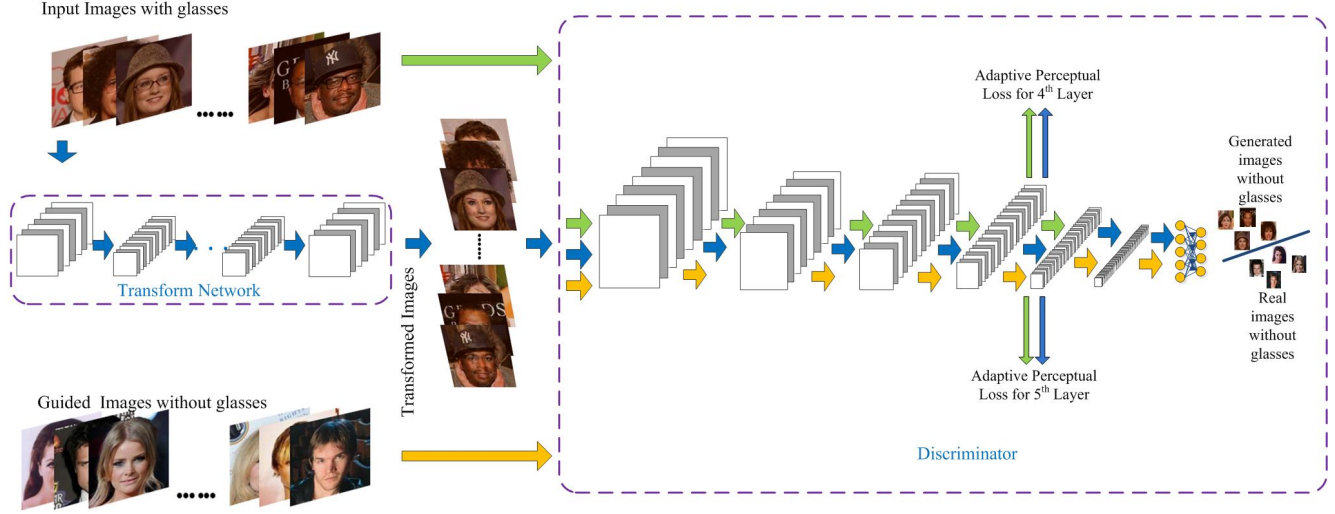


Figure 3. Illustration of transform network training in DIAT-A.

The training algorithm of the DIAT-A is almost the same with the original one. The only difference is that the objective function is replaced with Eqn. (14). Moreover, to prevent the drastic change of adaptive perceptual loss during training, we decrease the learning rate of discriminator and the transform network from 0.0001 to 0.00001. Even with such a small learning rate, DIAT-A converges faster than DIAT.

4.3. Discussion

Besides facial attribute transfer, our DIAT and DIAT-A may also be extended to other interesting applications. For example, in artistic style transfer [9], the ground truth transfer results generally are also unknown, and we expect that the use of GAN loss and adaptive perceptual loss could make the generated image more realistic and visually pleasant. For hand script imitation of Chinese characters, it is also difficult to obtain all the Chinese characters written by someone. Thus, one may train a CNN network to imitate hand script of unseen Chinese characters by design proper GAN and perceptual losses, and use the enhancement network to further improve the visual quality of imitation results.

Even for image deblurring and super-resolution [6], due to their intrinsic ill-posed property, the ground truth image to a given degraded image may be not unique. The introduction of GAN loss can be adopted to overcome the limitation of conventional mean squared error (MSE) based loss, and may result in more visually pleasant results. Moreover, for some low level vision tasks (e.g., rain removal), both the degradation model and ground truth images are unknown in training. One may treat image quality as some kind of attribute, and investigate suitable perceptual content loss and

GAN-based attribute loss to learn CNN for recovering high quality image while preserving the content of input image.

5. Experimental Results

5.1. Experimental setting

Our DIAT models are trained using a subset of the aligned CelebA dataset [27] by removing poor samples. According to the 5 landmarks provided by CelebA [27] for each image, the faces are clustered into 62 clusters using k-means. By observing the representative images of each cluster, we exclude 14 clusters due to that their faces are partially invisible in general. 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 .

We use all the images in the subset to pre-train the transform network and the discriminator. For further training, we use all the images with the reference attribute as the guided attribute set, and randomly select 2,000 images that are not with the reference attribute as the input set. As for the enhancement network, we take use of the same dataset used for training transform network. All the experiments are conducted on a computer with the GTX TitanX GPU of 12GB memory.

For performance evaluation, we consider five variants of our DIAT models, i.e., DIAT-A, DIAT-A without the enhancement network (DIAT-A0), DIAT with only attribute loss (DIAT1), DIAT without perceptual regularization (DIAT2) and enhancement network, DIAT without enhancement network (DIAT3). Generally, DIAT-A outperforms the other variants in our experiments. Unless stated otherwise, we also use DIAT to denote DIAT-A in this sec-

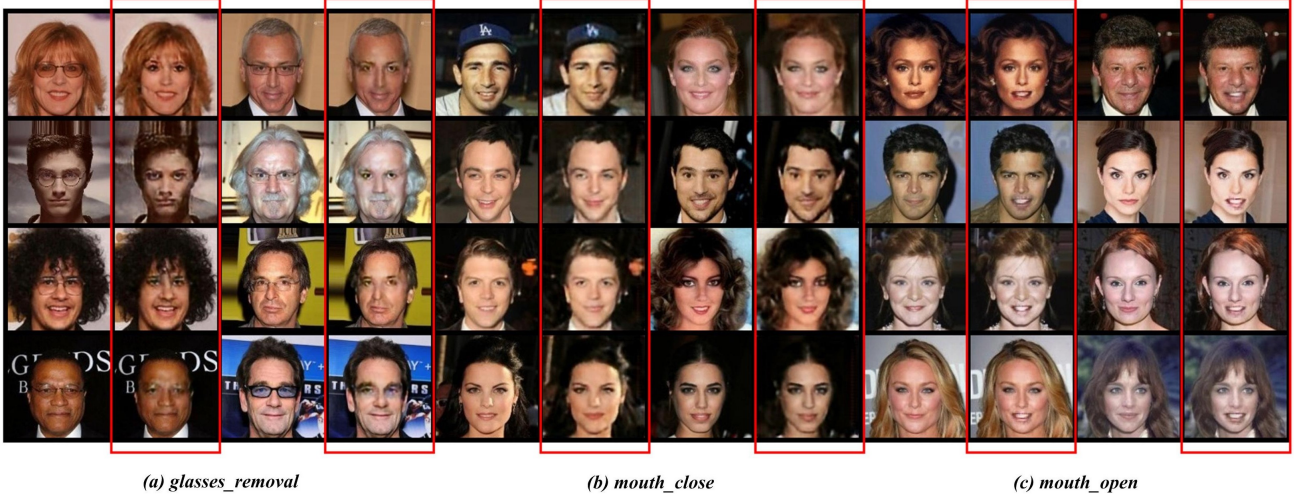


Figure 4. Local attribute transformation. For each attribute, the left column is the input face and the right column is the transfer result.

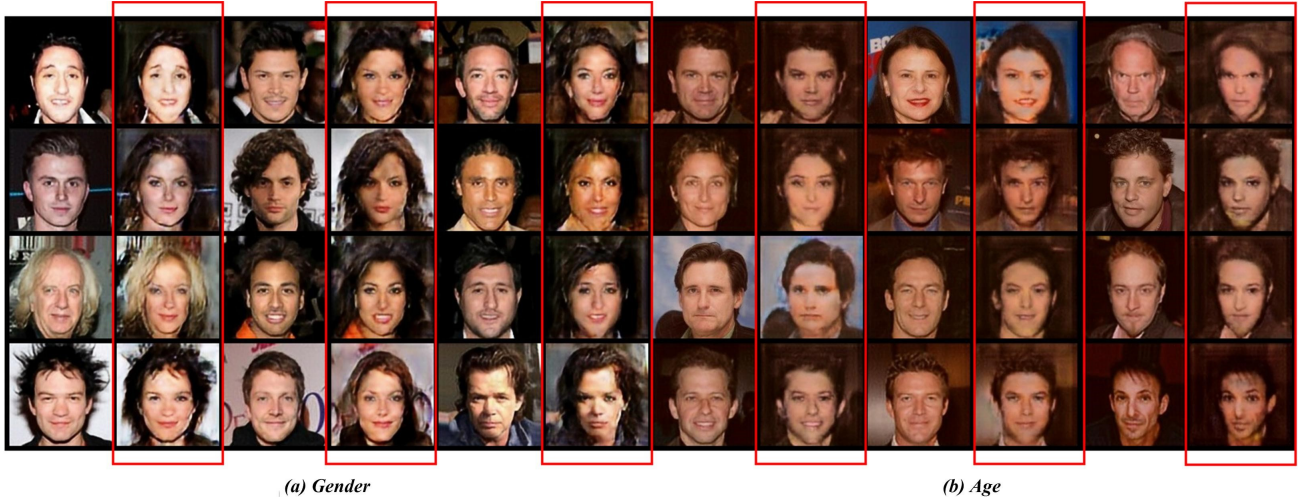


Figure 5. Global attribute transformation. For each attribute, the left column is the input face and the right column is the transfer result.

tion. Actually, almost no methods have been proposed to address different attribute transfer tasks. Thus we only compare our DIAT with the CNIA method [26]. For glasses removal, we can use the mask introduced in Sec. 3.2 to detect and remove the region of glasses, and then use the semantic inpainting method [40] to recover the missing pixels. So we also compare our DIAT with the semantic inpainting method [40] for glasses removal.

In the experiment, we set the parameter $\lambda = 0.1$ and $\gamma = 0.001$ for DIAT3, and set the parameter $\lambda = 0.1$ for DIAT-A. For local enhancement network, we set $\beta_0 = 0.1$, $\beta_1 = 0.5$, and $\beta_2 = 1$. For global enhancement network, we set the standard deviation of Gaussian blur kernel $\sigma = 1.8$.

5.2. Local attribute transfer

Our DIAT is evaluated for three local attribute transfer tasks, i.e., "mouth_open", "mouth_close", and "eyeglasses_remove". Fig. 4 shows the transfer results by DIAT. It can be seen that our DIAT performs favorably for transforming the input images to the reference attribute with satisfying visual quality. As shown in Fig. 4(b) and (c), when the training data are sufficient, we can train a DIAT model for transforming an image with "mouth_open" to "mouth_close", and train another DIAT model for the reverse task.

We further compare our DIAT with CNIA [26]. As shown in Fig. 6, the results by our DIAT are visually more pleasant than those by CNIA [26] for any local attribute transfer tasks. As to running time, our DIAT overwhelm

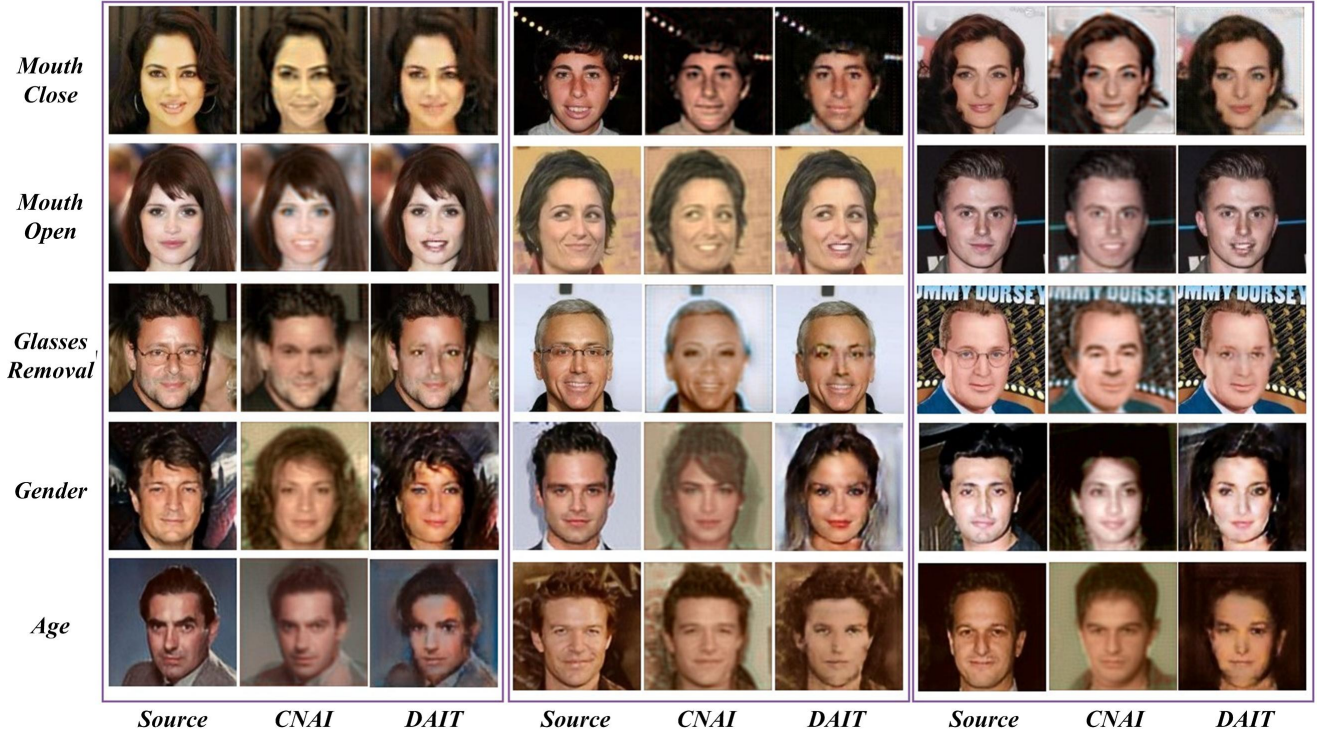


Figure 6. Comparison with CNIA [26].

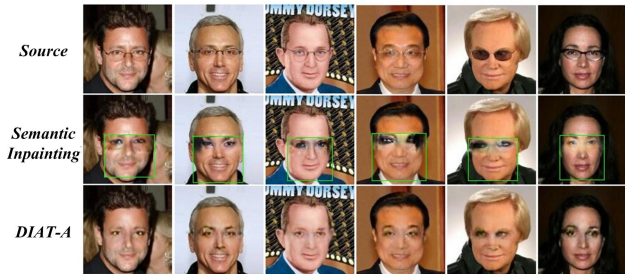


Figure 7. Comparison with semantic inpainting [40] for glasses removal. The green rectangle shows the input (partial face) for semantic inpainting

the counterpart [26]. CNIA takes about 30 seconds to generate an image. While our DIAT only needs 0.0045 second which makes it a real-time transformer.

Moreover, for glasses removal, we also compare our DIAT with semantic inpainting [40]. From Fig. 7, our DIAT can obtain much better transfer results.

5.3. Global attribute transfer

We consider two global attribute transfer tasks, i.e., "gender" and "age". For "gender" we only consider the case "male_to_female". And for "age" we only consider the case "older_to_younger". Fig. 5 shows the transfer re-

sults by DIAT, which validate the effectiveness of DIAT for global attribute transfer. Even gender transfer certainly causes the modification of the identity, as shown in Fig. 5(a), our DIAT-A can still retain most identity features, making the transfer result similar with the input image by identity. Fig. 6 compares our DIAT with CNIA [26], and the results by DIAT are visually more impressive.

5.4. Results on other real facial images

To assess the generalization ability, we train the DIAT models on CelebA and then use them for attribute transfer on other real facial images from LFW and PubFig. Each test image is first aligned with the 5 facial landmarks, and then input to the DIAT models. Here we consider one local (i.e., "mouth_open") and one global (i.e., "gender") attribute transfer tasks. Fig. 8 and Fig. 9 present the transfer results on the two tasks, which clearly validate the generalization ability of our models to real facial images.

5.5. Comparison of DIAT variants

Ablation studies are conducted to assess the contributions of several components of DIAT, including enhancement network, adaptive perceptual identity loss, perceptual identity loss and perceptual regularization. In the first experiment, we evaluate the enhancement network by comparing DIAT-A and DIAT-A0. We use two tasks, i.e., glasses



Figure 8. Local attribute transfer ("mouth_open") on LFW and Pubfig. The left column is the input face and the right column is the transfer result.

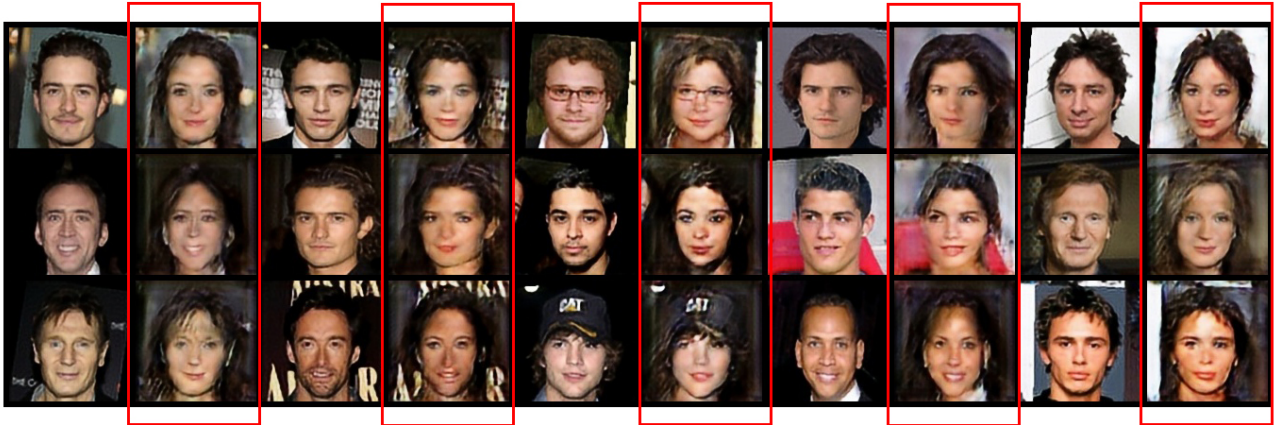


Figure 9. Gender transfer on LFW and Pubfig. The left column is the input face and the right column is the transfer result.

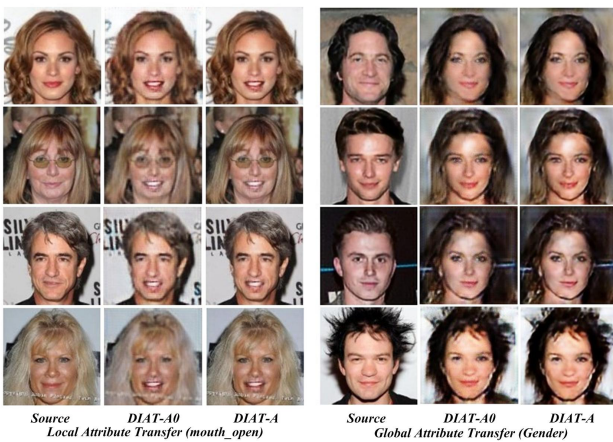


Figure 10. Comparison of DIAT-A and DIAT-A0.

removal and gender transfer. As shown in Fig. 10, the enhancement network is effective in improving the visual

quality while suppressing complex noise introduced by the transform network.

In the second experiment, we compare DIAT-A and DIAT to analyze the role of adaptive perceptual identity loss. Here we use gender transfer as an example. As shown in Fig. 11, the use of adaptive perceptual identity loss (i.e., DIAT-A) can bring two advantages: (i) better preserving the color of input image, and (ii) adaptive adjustment of hair length to match gender. Moreover, adaptive perceptual identity loss is also helpful in improving the convergence. The training time of DIAT-A is only about one third of DIAT.

In the third experiment, we test DIAT only with attribute loss (DIAT1). As shown in Figure 12, the attribute can still be transferred, but the transfer result does not keep the same or similar identity with the input image. Thus, the identity loss is indispensable for DIAT.

Finally, to analyze the role of perceptual regularization, experiments are conducted to compare DIAT2 and DIAT3.



Figure 11. Comparison of DIAT-A and DIAT.



Figure 12. Results by DIAT only with attribute loss (DIAT1).



Figure 13. Comparison of DIAT2 and DIAT3.

As shown in Fig. 13, the removal of perceptual regularization clearly has adverse effect, and the transfer results by DIAT2 are much poorer than those by DIAT3.

5.6. Limitation

Although we have evaluated our DIAT model on several attribute transfer tasks. Due to the limitation of CelebA, we do not test our DIAT on all possible tasks. For example, there are only 8500 images with hat in CelebA, which are insufficient for modeling the complex distributions of

facial images with different hats. Thus, we do not conduct the experiments for wearing hats to facial images. Similarly, sunglasses and reading glasses are mixed for the attribute "eyeglasses", making it difficult to train a DIAT for wearing glasses to facial images. Moreover, in CelebA all the attributes are annotated with binary labels. But the facial attributes should be continuous for facial reenactment and age progression. To address these issues, we will either construct some suitable large scale dataset for facial attribute transfer, or develop the new DIAT model which can be trained in more natural way and with fewer samples.

6. Conclusion

We present a unified deep convolutional network model (i.e., DIAT) for identity-aware facial attribute transfer. DIAT adopts a two-stage scheme and consists of a transform network and an enhancement. To train the transform network, we adopt the GAN framework, and define the adaptive perceptual identity loss on the feature maps of the attribute discriminator. Two kinds of enhancement networks are suggested to cope with local and global attribute transfer tasks, respectively. Experiments show that our model can obtain satisfying results for both local and global attribute transfer. Even for some identity-related attribute (e.g., gender transfer), our DIAT can obtain visually impressive results with minor modification on identity-related features. In future work, we will extend our model to hand script imitation and some low level vision applications.

References

- [1] M. Aharon, M. Elad, and A. Bruckstein. k -svd: An algorithm for designing overcomplete dictionaries for sparse representation. *IEEE Transactions on Signal Processing*, 54(11):4311–4322, Nov 2006.
- [2] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. L. Zitnick, and D. Parikh. Vqa: Visual question answering. In *ICCV*, 2015.
- [3] 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)*, volume 2, pages 60–65. IEEE, 2005.
- [4] 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.
- [5] 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*, pages 1486–1494, 2015.
- [6] 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*, pages 184–199. Springer, 2014.

- [7] A. Dosovitskiy, J. Tobias Springenberg, and T. Brox. Learning to generate chairs with convolutional neural networks. In *CVPR*, pages 1538–1546, 2015.
- [8] Y. Fu, G. Guo, and T. S. Huang. Age synthesis and estimation via faces: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(11):1955–1976, Nov 2010.
- [9] L. A. Gatys, A. S. Ecker, and M. Bethge. A neural algorithm of artistic style. *arXiv preprint arXiv:1508.06576*, 2015.
- [10] J. Gauthier. Conditional generative adversarial nets for convolutional face generation. *Class Project for Stanford CS231N: Convolutional Neural Networks for Visual Recognition, Winter semester*, 2014, 2014.
- [11] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *CVPR*, pages 580–587, 2014.
- [12] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *NIPS*, pages 2672–2680, 2014.
- [13] 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*, pages 1–26, 2016.
- [14] L. Hu, C. Ma, L. Luo, and H. Li. Single-view hair modeling using a hairstyle database. *ACM Trans. Graph.*, 34(4):125:1–125:9, July 2015.
- [15] A. E. Ichim, S. Bouaziz, and M. Pauly. Dynamic 3d avatar creation from hand-held video input. *ACM Trans. Graph.*, 34(4):45:1–45:14, July 2015.
- [16] J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. *arXiv preprint arXiv:1603.08155*, 2016.
- [17] I. Kemelmacher-Shlizerman. Transfiguring portraits. *ACM Trans. Graph.*, 35(4):94:1–94:8, July 2016.
- [18] I. Kemelmacher-Shlizerman, A. Sankar, E. Shechtman, and S. M. Seitz. *Being John Malkovich*, pages 341–353. 2010.
- [19] I. Kemelmacher-Shlizerman, S. Suwajanakorn, and S. M. Seitz. Illumination-aware age progression. In *2014 IEEE Conference on Computer Vision and Pattern Recognition*, pages 3334–3341, June 2014.
- [20] D. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [21] D. P. Kingma and M. Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [22] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, pages 1097–1105, 2012.
- [23] T. D. Kulkarni, W. F. Whitney, P. Kohli, and J. Tenenbaum. Deep convolutional inverse graphics network. In *NIPS*, pages 2530–2538, 2015.
- [24] N. Kumar, A. C. Berg, P. N. Belhumeur, and S. K. Nayar. Attribute and simile classifiers for face verification. In *ICCV*, pages 365–372. IEEE, 2009.
- [25] C.-Y. Lee, S. Xie, P. Gallagher, Z. Zhang, and Z. Tu. Deeply-supervised nets. In *AISTATS*, volume 2, page 6, 2015.
- [26] M. Li, W. Zuo, and D. Zhang. Convolutional network for attribute-driven and identity-preserving human face generation. *arXiv preprint arXiv:1608.06434*, 2016.
- [27] 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.
- [28] 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*, pages 3431–3440, 2015.
- [29] A. Mahendran and A. Vedaldi. Understanding deep image representations by inverting them. In *CVPR*, pages 5188–5196. IEEE, 2015.
- [30] O. M. Parkhi, A. Vedaldi, and A. Zisserman. Deep face recognition. In *British Machine Vision Conference*, 2015.
- [31] A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- [32] S. Reed, Z. Akata, X. Yan, L. Logeswaran, B. Schiele, and H. Lee. Generative adversarial text to image synthesis. *arXiv preprint arXiv:1605.05396*, 2016.
- [33] A. Selim, M. Elgharib, and L. Doyle. Painting style transfer for head portraits using convolutional neural networks. *ACM Transactions on Graphics (TOG)*, 35(4):129, 2016.
- [34] K. Simonyan, A. Vedaldi, and A. Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*, 2013.
- [35] 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.
- [36] 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.
- [37] A. van den Oord, N. Kalchbrenner, and K. Kavukcuoglu. Pixel recurrent neural networks. *arXiv preprint arXiv:1601.06759*, 2016.
- [38] A. van den Oord, N. Kalchbrenner, O. Vinyals, L. Espeholt, A. Graves, and K. Kavukcuoglu. Conditional image generation with pixelcnn decoders.
- [39] X. Yan, J. Yang, K. Sohn, and H. Lee. Attribute2image: Conditional image generation from visual attributes. *arXiv preprint arXiv:1512.00570*, 2015.
- [40] R. Yeh, C. Chen, T. Y. Lim, M. Hasegawa-Johnson, and M. N. Do. Semantic image inpainting with perceptual and contextual losses. *arXiv preprint arXiv:1607.07539*, 2016.
- [41] Z. Zhang, P. Luo, C. C. Loy, and X. Tang. Facial landmark detection by deep multi-task learning. In *European Conference on Computer Vision*, pages 94–108. Springer, 2014.
- [42] 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*, pages 3943–3951, 2015.
- [43] 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*, 27(2):107–126, 1998.