

# ReshapeGAN: Object Reshaping by Providing A Single Reference Image

Ziqiang Zheng\*<sup>1</sup>, Yang Wu<sup>†2</sup>, Zhibin Yu<sup>1</sup>, Yang Yang<sup>3</sup>, Haiyong Zheng<sup>1</sup>, and Takeo Kanade<sup>4</sup>

<sup>1</sup>Ocean University of China.

<sup>2</sup>Nara Institute of Science and Technology, Japan.

<sup>3</sup>University of Electronic Science and Technology of China.

<sup>4</sup>Carnegie Mellon University, United States.

## Abstract

The aim of this work is learning to reshape the object in an input image to an arbitrary new shape, by just simply providing a single reference image with an object instance in the desired shape. We propose a new Generative Adversarial Network (GAN) architecture for such an object reshaping problem, named ReshapeGAN. The network can be tailored for handling all kinds of problem settings, including both within-domain (or single-dataset) reshaping and cross-domain (typically across multiple datasets) reshaping, with paired or unpaired training data. The appearance of the input object is preserved in all cases, and thus it is still identifiable after reshaping, which has never been achieved as far as we are aware. We present the tailored models of the proposed ReshapeGAN for all the problem settings, and have them tested on 8 kinds of reshaping tasks with 13 different datasets, demonstrating the ability of ReshapeGAN on generating convincing and superior results for object reshaping. To the best of our knowledge, we are the first to be able to make one GAN framework work on all such object reshaping tasks, especially the cross-domain tasks on handling multiple diverse datasets. We present here both ablation studies on our proposed ReshapeGAN models and comparisons with the state-of-the-art models when they are made comparable, using all kinds of applicable metrics

\*equal contribution

†equal contribution

that we are aware of. Code will be available at <https://github.com/zhengziqiang/ReshapeGAN>.

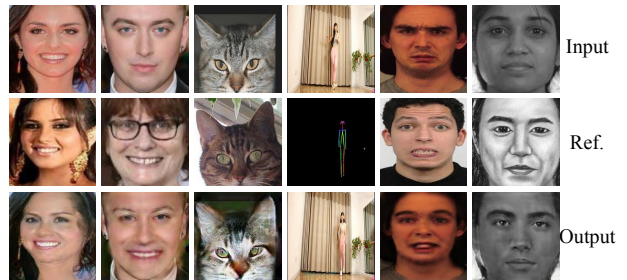


Figure 1: Object reshaping with ReshapeGAN, guided by another object from a single reference image (Ref.). The appearance is preserved (thus being still identifiable) while the reshaping is controllable by just choosing a reference image with the desired shape. The reference can be from an arbitrary domain (same as or different from that of the input).

## 1 Introduction

Many of us may have admired some others' faces or bodies, though it may be hard for us to even dream about changing our own faces/bodies to those desired ones, not to say taking actions for actual reshaping. But what if the reshaping can be seen instantly, as shown in Fig. 1, by just

providing a picture of ours and another one of some other person who has the desired shape? In greater details, if the virtual reshaping can preserve our identity and appearance, while at the same time transferring to the same nice expression or pose as that of the admired person, won't you want to try and check this visible dream?

It is actually even beyond what dreaming can do, because as humans we can hardly imagine the details of the changes caused by such a reshaping if we have never seen similar results before. A skilled artist may be able to do such translation and imagination, but it is definitely not an easy task. Nevertheless, we have the ability to feel how good such reshaping results look like, even without a ground truth. We can be somehow sensitive about how well the identity and appearance are preserved, how well the changed shape matches that of the reference image, and how realistic the generated image is. Clearly, it is fun to try such a task.

In this paper, we propose a framework called **ReshapeGAN** that can automatically learn to do such a reshaping, of which some exemplar results are shown in Fig. 1. ReshapeGAN not only works on faces/bodies, but also applies to other objects like cats.

ReshapeGAN is a type of Generative Adversarial Network (GAN), which was first introduced by Goodfellow et. al. in 2014 [19]. GAN has been developed and proved to be very effective on image generation tasks for many applications [44, 79, 58, 74, 80, 25, 71, 52, 10, 23, 35, 39, 29]. Generally, conditionally generating images could fall into two categories: *supervised image generation* and *unsupervised image-to-image translation*. The former can usually generate higher quality results [44, 68, 1], which are more realistic [24] or with higher resolution [64], but it has the limitation of relying on paired training data (input and ground truth pairs) which may be hard or impossible to collect in many applications. The latter doesn't require such paired training data and thus being more widely applicable. Though unsupervised image-to-image translation is more challenging, its merits have attracted researchers to make a lot of progresses including CycleGAN [80], DualGAN [71], DiscoGAN [25], MUNIT [23], DRIT [35], StarGAN [10], etc [22]. Specially, the cross-domain or cross-dataset unsupervised image-to-image translation is most challenging, as each domain or dataset has its own style and attributes. Image generation have to manipulate and control such styles and attributes

for a desired output [44, 79, 80, 25, 71, 52, 10, 23, 35]. The proposed ReshapeGAN learns to preserve the appearance of the input object instance and get the shape information from the reference image for the reshaping tasks. Its framework can be tailored for both supervised image generation and unsupervised image-to-image translation, within a domain or across domains/datasets, as shown in Fig 2. As far as we are aware, this is the first time that object reshaping is made possible for all these settings, especially for the cross-domain/cross-dataset setting.

Usually, cross-domain image-to-image translation requires a large amount of training data from each individual for ensuring a reasonably good performance, due to the possibly large style and attribute differences between domains/datasets. Nevertheless, ReshapeGAN is made effective for working with relatively small amount of data from each individual, so that being as generally applicable as possible.

To demonstrate the effectiveness and superiority of ReshapeGAN, we conduct extensive experiments for each of the three main settings, resulting in totally 8 different tasks on 13 datasets. Since there is no unique metric that is widely recognized as the standard for generated image's quality assessment, we adopt all the applicable metrics that we are aware of and use them for performance evaluation and comparison. We compare ReshapeGAN with state-of-the-art methods on all the tasks which they can be applied to or made applicable for, and do ablation studies to verify the effectiveness of each component of ReshapeGAN.

In brief, the main contributions of this work can be summarized as follows.

- It presents as far as we know the first general reference-guided object reshaping framework, which preserves the object's identity and works on most diverse input-reference pairs.
- It introduces tailored models of the proposed framework, which are applicable to both supervised and unsupervised reshaping, for both within-domain and cross-domain or cross-dataset scenarios.
- Extensive experiments have been done on all the settings with many tasks and datasets, to show the effectiveness of the proposal with ablation studies and

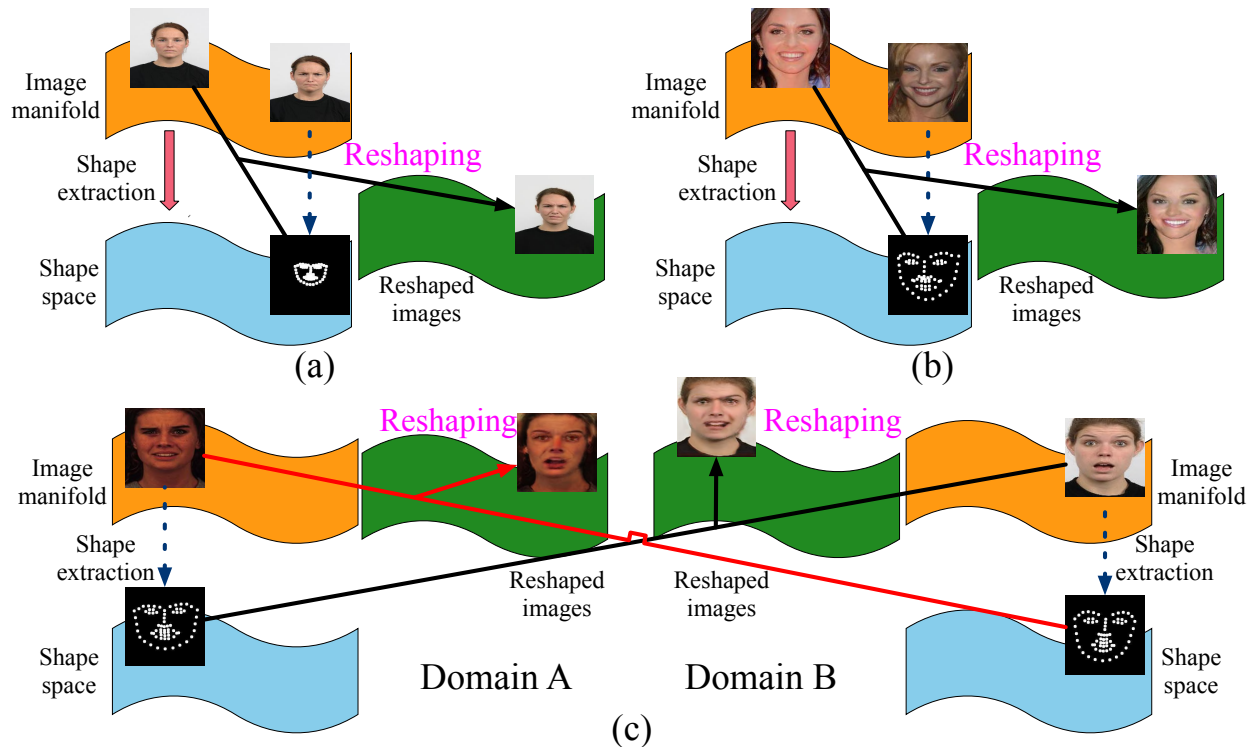


Figure 2: The typical settings that ReshapeGAN handles: (a) reshaping by within-domain shape guidance with paired data, (b) reshaping by within-domain shape guidance with unpaired data, and (c) reshaping by cross-domain shape guidance.

its significant superiority in comparison with state-of-the-art methods when they are or are made to be comparable.

## 2 Related work

### 2.1 Generative Adversarial Networks (GANs)

Thanks to the development of GAN [19], the recent extensions of GAN have achieved a great progress in various vision tasks such as image editing and in-painting [17, 51, 69, 3], super resolution [34, 9, 59], image restoration and enhancement [40, 72], object detection [36, 47] and other applications [12, 79, 74, 43, 55, 46]. Regarding with adversarial training, how to avoid model

collapse and reduce the training instability is an open discussion [39, 11]. To achieve this goal, some studies focus on designing novel network architectures and training mechanisms [8, 48] while more studies tried to solve this problem by updating adversarial loss functions such as Boundary-Seeking Generative Adversarial Networks [13], Wasserstein Generative Adversarial Networks (WGAN) [2], Loss-Sensitive Generative Adversarial Networks (LS-GAN) [53], Least Square Generative Adversarial Networks (LSGAN) [43], etc [29]. In our tasks, we also include some efficient generative adversarial losses to improve the training process and yield more plausible results. We take advantage of WGAN [2] and its extension WGAN-GP [20] to make the adversarial training process more stable. Furthermore, we use Perturbed loss from DRAGAN [28] to improve the synthetic quality.

## 2.2 GAN based image-to-image translation

Benefited from the success of conditional GANs (cGANs) [44], many researches related with GANs focused on image-to-image translation tasks. According to the data type, these tasks can be roughly divided into paired [24, 64, 81] and unpaired [80, 71, 60, 37, 25, 54] image-to-image translation. Paired image-to-image translation approaches usually require the paired images for training. They always need one or more pixel-wise restriction such as L1 norm to implement supervised learning with data pairs [24]. However, paired data are expensive and sometimes impossible to be collected. To overcome this shortage, CycleGAN, DualGAN and DiscoGAN applied a cycle consistency to translate images using unpaired images [80, 25, 71]. Besides, Choi et al. proposed StarGAN [10] that used a single model and latent code to handle one-to-many image-to-image translation tasks. Similar to ACGAN [48], which trained the generator with the discriminator that combines with an auxiliary classifier, StarGAN inherited the structure and the cycle consistency loss from CycleGAN and used an additional classifier to control the synthesized attribute. Besides, combining variational autoencoders (VAE) [27] with GANs is another important branch to translate images between different domains [32] with unpaired data. Fader Networks [30] used the attribute-invariant representations, encoded by the input image, and the latent code for image reconstruction. Zhu et al. proposed BicycleGAN [81], which combines VAE-GAN [32] objects and latent regressor objects [14, 16] for a bijective consistency to obtain more realistic and diverse samples using paired data. These studies have advanced the development of the one-to-one image-to-image translation. However, no existing work can achieve identity-preserved object reshaping for both the case of training with paired data and that with only unpaired data across domains/datasets.

## 2.3 Conditional image editing

Conditional generative models are widely used for image synthesis under one or more given conditions. Many studies were developed based on two pioneer works: conditional variational Autoencoder (CVAE) [68] and conditional Generative Adversarial Network (cGAN) [44]. Lassner et al. proposed a conditional architecture hybridizing VAE with adversarial training to generate full-body people in clothing [33].

Reed et al. [55] presented a generative model to generate bird and flower images conditioned on text descriptions by adding textual information to both generator and discriminator. They further discussed the feasibility to control the features, structure and the locations of the generated images with different conditional text-to-image models [57, 56]. StackGAN [74] and its extension StackGAN++ [73] can also use text descriptions to generate high-quality photo-realistic images. As mentioned before, StarGAN [10] used one-hot encoded biases as conditions to translate images from one domain to another with unpaired data.

Comparing with these works which use labels and texts as the conditions for image generation, using multiple images as conditions for image synthesis is more challenging. Ma et al. [42] considered both pose images and person images as conditions to guide the network to generate a person image with a specified pose. Yang et al. [70] further extended this idea to generate videos under the constraint of pose series. Zhao et al. [78] explored generating multi-view cloth images from only a single view input, while Ma et al. [41] and Park et al. [50] used an extra image as exemplar for semantic-preserved unsupervised translation, which have a similar motivation to our task. However, most of image guided image editing approaches rely on paired images to implement a pixel-wise supervised training, or limit to some low-level (color, texture, etc.) translation applications. Unlike those methods, our proposed ReshapeGAN can work on unpaired training data for high-level object reshaping tasks. We achieve this by making use of shape and appearance information in a more efficient and flexible way. Our model can generate a desired image that preserves the appearance of a specific object instance while borrowing the shape information from another image, even across domains.

## 3 Network Architecture for Object Reshaping

Given an input image  $x \in X$  and a reference image  $y \in Y$  for geometric (shape) guidance represented by  $s_y$ , object reshaping targets at learning a generator  $G$  which can generate a new image  $G(x, s_y)$  inheriting  $x$ 's appearance

while at the same time changing to  $y$ 's shape  $s_y$ . Our proposal for object reshaping is called ReshapeGAN, which can be tailored for three typical settings (Fig. 2). We introduce each of the settings and our corresponding tailored model in the following subsections.

### 3.1 Reshaping by within-domain guidance with paired data

When the reference image is in the same domain (which usually means the same dataset or style) as the input image and the reference image is about the same object instance (i.e., having the same identity) as the one in the input image, we get the easiest setting for object reshaping, which has paired training data. In this case, ReshapeGAN can be learned by solving the following problem:

$$G^* = \arg \min_G \max_D \mathcal{L}_{ReshapeGAN}^{paired}(G, D), \quad (1)$$

with

$$\begin{aligned} \mathcal{L}_{ReshapeGAN}^{paired}(G, D) = & \mathcal{L}_{adv}^s(G, D) + \gamma \mathcal{L}_{perturb}(D) \\ & + \delta \mathcal{L}_{pixel}^{paired}(G) + \sigma \mathcal{L}_{percep}^{paired}(G), \end{aligned} \quad (2)$$

where  $D$  is its discriminator;  $\mathcal{L}_{adv}^s(G, D)$  is the adversarial loss with shape guidance;  $\mathcal{L}_{perturb}(D)$  is the perturbed loss for regularizing  $D$ ;  $\mathcal{L}_{pixel}^{paired}(G)$  is the paired pixel-level appearance matching loss;  $\mathcal{L}_{percep}^{paired}(G)$  denotes the perceptual loss;  $\gamma$ ,  $\delta$ , and  $\sigma$  are super-parameters for balancing the corresponding losses. Details of the loss functions are explained as follows and the overall model is illustrated in Fig. 3.

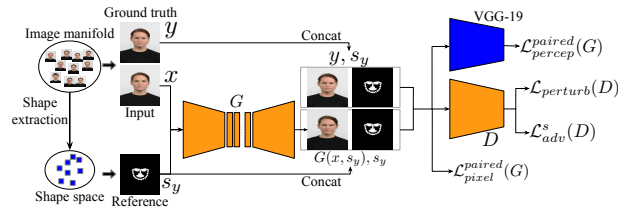


Figure 3: The detailed ReshapeGAN model for reshaping by within-domain guidance with paired data.

**Shape guided adversarial loss  $\mathcal{L}_{adv}^s(G, D)$ .** Based on our shape guided generator  $G$  and following the common definition of adversarial loss, we have  $\mathcal{L}_{adv}^s(G, D)$

defined as

$$\begin{aligned} \mathcal{L}_{adv}^s(G, D) = & \mathbb{E}_y [\log D(y, s_y)] \\ & + \mathbb{E}_{x, y} [\log(1 - D(G(x, s_y), s_y))]. \end{aligned} \quad (3)$$

The adversarial loss is a basic requirement in an adversarial network, which makes sure that the generated image  $G(x, s_y)$  cannot be distinguished from the reference image  $y$  by the discriminator  $D$ , in terms of the shape information, but not the appearance as in existing works.

**Perturbed loss  $\mathcal{L}_{perturb}(D)$ .** To further improve the robustness of the discriminator, we introduce to our model the perturbed loss originated from Gulrajani et al.'s DRAGAN work [28]. The perturbed loss can be expressed as

$$\begin{aligned} \mathcal{L}_{perturb}(D) = & \mathbb{E}_{\hat{x}, y} [(\|\nabla_{\hat{x}} D(\hat{x}, s_y)\|_2 - 1)^2], \\ \text{with } \hat{x} = & (1 - \alpha)x + \alpha z, \text{ where } \alpha \sim U[0, 1]. \end{aligned} \quad (4)$$

Here  $x$  represents a real image sample and  $z$  means the random noise, which follows the Gaussian distribution;  $\alpha$  is the random hyper-parameter that controls the balance between real and noise and follows the continuous uniform distribution between 0 and 1;  $\nabla_{\hat{x}}$  is the gradient of  $D(\hat{x}, s_y)$ . Please note that  $\hat{x}$  is only used for calculating  $\mathcal{L}_{perturb}(D)$ . We believe that such a perturbed loss makes the convergence more stable and generate images with higher quality. The experiment results in Section 4 also support this point.

**Paired pixel-level appearance matching loss  $\mathcal{L}_{pixel}(G)$ .** It is a widely used loss for ensuring that the generated image matches the ground truth image or has a desire appearance at the pixel-level, and usually the L1 norm is used. When paired training data is available, the loss is simple as:

$$\mathcal{L}_{pixel}^{paired}(G) = \mathbb{E}_{x, y} [\|G(x, s_y) - y\|_1], \quad (5)$$

which is about the generation loss w.r.t. the provided ground truth data  $y$ .

**Paired perceptual loss  $\mathcal{L}_{percep}^{paired}(G)$ .** Unlike the pixel-level loss which matches two images pixel by pixel, the perceptual loss measures the similarity at the feature level. In greater details, we follow Chen et al.'s work [7] and include an extra pre-trained VGG-19 network on ImageNet to compute the perceptual loss [34, 15, 7]. Compared to pixel-level loss, this loss combines the distance at multiple scales of feature representation, which could carry

low-level and high-level information of images. The perceptual loss is described as:

$$\mathcal{L}_{percep}^{paired}(G) = \sum_n \lambda_n \mathbb{E}_{x,y} [\|\Phi_n(y) - \Phi_n(G(x, s_y))\|_1], \quad (6)$$

where  $\Phi_n$  is the feature extractor at the  $n_{th}$  level of the pre-trained VGG-19 network. Following [7], we compute the perceptual loss between outputs at defined  $N = 5$  selective layers. The weights of pre-trained model will not be optimized during the learning process. The hyperparameter  $\lambda_n$  controls the influence of perceptual loss at different scales. The perceptual loss from the higher layer controls the global structure, and the loss from the low layer controls the local details during generation. Thus, the generator should provide better synthesized images to cheat this hybrid discriminator and finally improve the synthesized image quality. Please note, only for the paired training, we compute the perceptual loss between the target images and generated images. The detail information of the proposed method for reshaping by within-domain guidance with paired data could be found in Fig. 3.

### 3.2 Reshaping by within-domain guidance with unpaired data

For the case that the reference image is in the same domain as the input image but they are not about the same instance/identity, the model has to be learned with unpaired training data. In this case, we use a model similar to the one introduced in Section 3.1, but having the pixel-level appearance matching loss and perceptual loss in their unpaired version. In great detail, the overall loss function becomes

$$\begin{aligned} \mathcal{L}_{ReshapeGAN}^{unpaired}(G, D) &= \mathcal{L}_{adv}^s(G, D) + \gamma \mathcal{L}_{perturb}(D) \\ &+ \delta \mathcal{L}_{pixel}^{unpaired}(G) + \sigma \mathcal{L}_{percep}^{unpaired}(G), \end{aligned} \quad (7)$$

where the unpaired version of pixel-level appearance matching loss  $\mathcal{L}_{pixel}^{unpaired}(G)$  and perceptual loss  $\mathcal{L}_{percep}^{unpaired}(G)$  are defined below, and this ReshapeGAN model is illustrated in Fig. 4.

**Unpaired pixel-level appearance matching loss**  $\mathcal{L}_{pixel}^{unpaired}(G)$ . When there is no paired data for training, inspired by CycleGAN [80], we use the cycle consistency

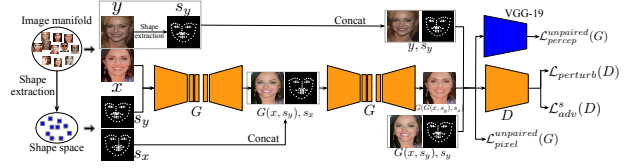


Figure 4: The detailed ReshapeGAN model for reshaping by within-domain guidance with unpaired data.

reconstruction loss with two rounds of reshaping guided by  $s_x$  and  $s_y$  (shapes of  $x$  and  $y$ , respectively) as the unpaired pixel-level appearance matching loss:

$$\mathcal{L}_{pixel}^{unpaired}(G) = \mathbb{E}_{x,y} [\|G(G(x, s_y), s_x) - x\|_1]. \quad (8)$$

**Unpaired perceptual loss**  $\mathcal{L}_{percep}^{unpaired}(G)$ . In the case that there is no paired sample with the given input  $x$ , it is reasonable to assume that the generated image  $G(x, s_y)$  has the same perceptual response as that of  $x$  due to the desired appearance and identity preserving property, those they shall have different shapes ( $s_y$  vs.  $s_x$ ). Unlike the pixel-wise matching loss which requires good alignment and thus minimum shape difference, perceptual loss is about feature-level perception results, which can have some robustness to shape changes. Therefore, we define the unpaired perceptual loss to be the that between the generated image  $G(x, s_y)$  and the input image  $x$ :

$$\mathcal{L}_{percep}^{unpaired}(G) = \sum_n \lambda_n \mathbb{E}_{x,y} [\|\Phi_n(x) - \Phi_n(G(x, s_y))\|_1], \quad (9)$$

where  $\Phi_n$  is the feature extractor at the  $n_{th}$  level of the pre-trained VGG-19 network, the same as the one in Equation 6.

### 3.3 Reshaping by cross-domain guidance

The hardest setting for object reshaping is about cross-domain guidance, i.e., the case when the reference image and input images are from two different domains (e.g. two different styles or datasets). For this setting, we found that it is hard for the unpaired ReshapeGAN model introduced in Section 3.2 to learn stable appearance preserved object reshaping due to the possibly large domain differences. Therefore, we propose a two-stage strategy for dividing

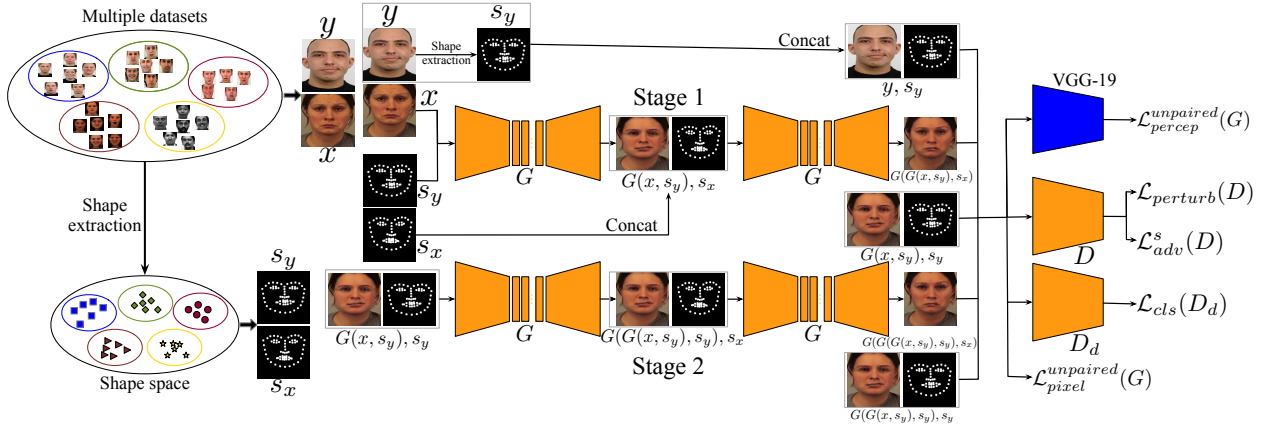


Figure 5: The detailed ReshapeGAN model for reshaping by cross-domain guidance with unpaired data.

the whole task into two sub-tasks: domain-preserved reshaping and refining. The detailed framework is shown in Fig. 5.

For the first stage, we add a *domain classification loss* to the overall loss function of  $\mathcal{L}_{ReshapeGAN}^{unpaired}(G, D)$  to ensure proper domain preservation during the reshaping. In greater details, the overall loss becomes

$$\begin{aligned}
& \mathcal{L}_{ReshapeGAN}^{cross-domain.1}(G, D, D_d) \\
&= \mathcal{L}_{ReshapeGAN}^{unpaired}(G, D) + \lambda \mathcal{L}_{cls}(D_d) \\
&= \mathcal{L}_{adv}^s(G, D) + \gamma \mathcal{L}_{perturb}(D) \\
&+ \delta \mathcal{L}_{pixel}^{unpaired}(G) + \sigma \mathcal{L}_{percep}^{unpaired}(G) + \lambda \mathcal{L}_{cls}(D_d),
\end{aligned} \tag{10}$$

where  $D_d$  is a domain classifier and  $\mathcal{L}_{cls}(D_d)$  is the domain classification loss, with  $\lambda$  as its weight.

**Domain classification loss  $\mathcal{L}_{cls}(D_d)$ .** The domain classifier  $D_d$  is expected to be able to help ensuring the domain preserving, leading to the effect of appearance preserving, i.e., maintaining the appearance and identity of the object during the reshaping. More specifically,  $\mathcal{L}_{cls}(D_d)$  is defined as

$$\begin{aligned}
& \mathcal{L}_{cls}(D_d) \\
&= \mathbb{E}_x [-\log D_d(x, d_x)] + \mathbb{E}_y [-\log D_d(y, d_y)] \\
&+ \mathbb{E}_{x,y} [-\log D_d(G(x, s_y), d_x)],
\end{aligned} \tag{11}$$

where  $d_x$  and  $d_y$  denote the domain labels of  $x \in X$  and  $y \in Y$ , respectively. Here we derive this loss from StarGAN [10], and we use the one-hot encoding to capture

the domain distribution. Our strategy is similar to that in the work of MUNIT [23], which tries to learn disentangled representation to get more expressive and representative style information. Following the CRN [7], we use the channel-wise concatenation to integrate the image and domain label (e.g.  $x$  and  $d_x$ ) for the classifier  $D_d$ .

In the second stage, we apply the object reshaping again by making use of the  $\mathcal{L}_{ReshapeGAN}^{unpaired}(G, D)$ , so that a second round of reshaping (should be significantly minor than that in the first stage) can be performed. We take the output of the first stage as the input of this stage. The overall loss for this stage is just

$$\begin{aligned}
& \mathcal{L}_{ReshapeGAN}^{cross-domain.2}(G, D, D_d) = \mathcal{L}_{ReshapeGAN}^{unpaired}(G, D) \\
&= \mathcal{L}_{adv}^s(G, D) + \gamma \mathcal{L}_{perturb}(D) \\
&+ \delta \mathcal{L}_{pixel}^{unpaired}(G) + \sigma \mathcal{L}_{percep}^{unpaired}(G).
\end{aligned} \tag{12}$$

The reason for splitting a difficult cross-domain object reshaping task to two easier sub-tasks is that we can get better performance than that of a one-stage generation (using only the first stage). In the one-stage generation (the first stage of our proposal), two types of information (domain/style information  $d$  and shape/geometric information  $s$ ) are used as constrains, making it hard to ensure quality reshaping. Ablation study will be given on verifying this in the following experiment section.

### 3.4 Implementation details

Our backbone network derives from StarGAN [10], but different from [10], we concatenate geometric guidance and raw input to get a new input for both generator and discriminator. Besides, in order to capture domain label effectively, we devise a new architecture for generator. In our experiments, we find that the model performs well when we choose  $\gamma = \eta = 1$  and  $\delta = 10$ . Our final loss is the sum of these losses as described in Eq. 7, 10 and 12. We apply this final loss to the generator with Adam optimizer of learning rate 0.0002. Our code will be made available at <https://github.com/zhengziqiang/ReshapeGAN>. All detailed implementation could be found in our code.

## 4 Experiments and results

### 4.1 Evaluation metrics

**Learned Perceptual Image Patch Similarity** (LPIPS) is first proposed by [75], which computes the perceptual similarity (actually in terms of distance) between two image patches. Lower LPIPS means the two image patches have higher perceptual similarity. Considering two image domains, we can compute the LPIPS metric (averaged over sampled patch pairs) to evaluate the perceptual similarity between them.

**Fréchet Inception Distance** (FID) computes the similarity between the generated sample distribution and real data distribution. This metric is consistent and robust for evaluating the quality of generated images [39, 4], and it can be calculated by:

$$\text{FID} = \|\mu_x - \mu_g\|_2^2 + \text{Tr} \left( \sum_x + \sum_g - 2(\sum_x \sum_g)^{\frac{1}{2}} \right), \quad (13)$$

where  $(\mu_x, \sum_x)$  and  $(\mu_g, \sum_g)$  are pairs of the mean and covariance of the sample embeddings from the real data distribution and generated data distribution, respectively. Lower FID means smaller distribution difference between the generated and the target images and therefore higher quality of generated images.

**Geometrical Consistency** is required for our object reshaping tasks to evaluate the matching level geometrically. In order to know whether the model can synthesize

images with desired geometric information, we use available state-of-the-art geometry estimation model to extract the geometric information (such as landmarks and poses) from both the synthesized images and the target images, and compare the results from sample pairs in the shape space. For facial expression generation, we use dlib [26] to extract landmarks and measure the similarity between two images by computing SSIM (see Traditional Evaluation Metrics below for details) and LPIPS scores using the landmark information, marked as **SSIM (Landmark)** and **LPIPS (Landmark)** respectively.

**Identification Distance** is important for our object reshaping tasks to evaluate whether the appearance and identity information are preserved while the objects are reshaped. We adopt object instance recognition (i.e., identification) algorithms to evaluate the similarity or distance between generated images and input images. In the case of faces, we use an effective open-source face recognition algorithm<sup>1</sup> to do that. A significantly large distance according to the classifier, i.e., more than the normal cut-off 0.6, indicates that the two faces are from two different identities. That is, if the distance between two faces is lower than 0.6, we can consider that the two faces have the same identity. Moreover, theoretically, for the distance below 0.6, larger distance might indicate better reshaping with preserved identity, while too small distance may mean that the model fails to do reshaping so that the generated image is almost identical to the input image. This can also be proved by computing the identification distances on RaFD [31] dataset, where each identity has 8 different emotional expressions captured from 5 different angles/viewpoints ( $0^\circ$ ,  $\pm 45^\circ$  and  $\pm 90^\circ$ ), and the average distance between the neutral faces and the other emotional expressions from  $0^\circ$  viewpoint of all identities is about 0.3, while, the average distance between the neutral faces from  $0^\circ$  viewpoint and the other emotional expressions from  $\pm 45^\circ$  viewpoint is about 0.5, and the average distance between the neutral faces from  $0^\circ$  viewpoint and the other emotional expressions from  $\pm 90^\circ$  viewpoint is about 0.6.

**User Study (Identity / Shape)** is the human evaluation applicable for all the generated images. We also use this golden standard for our object reshaping tasks. Here we ask 20 volunteers (users) to give a True/False judgement

<sup>1</sup>[https://github.com/ageitgey/face\\_recognition](https://github.com/ageitgey/face_recognition)



on the output images. Specifically, we give the user an input image, a reference image and a synthesized image, and ask them to judge true or false whether the synthesized image keep the identity information (Identity) and whether the synthesized image has the same shape expression with the reference image (Shape). The users are given unlimited time to make the decision. For each comparison, we randomly generate 100 images and each image is judged by at least 2 different users. Only higher identity metric or only higher shape metric can not provide confident performance for our object reshaping tasks, since they always require both identity preservation and reshaping ability. So higher votes to both same identity and consistent shape could indicate better reshaping performance.

**Traditional Evaluation Metrics.** For the cases where paired data exists, we can use some traditional image quality assessment metrics for performance evaluation, including **MSE**, **RMSE**, **PSNR** and **SSIM**. MSE (Mean Square Error) and RMSE (Root Mean Square Error) [67] compute pixel-wise errors between synthesized images and real images, lower MSE and RMSE represent higher image generation quality. PSNR (Peak Signal-to-Noise Ratio) [21] can roughly evaluate the image quality independently, and usually the higher PSNR the better. SSIM (Structural Similarity Index Measure) [21] measures the similarity between two images, and higher SSIM denotes higher structural similarity between generated images and real images.

## 4.2 Reshaping by within-domain guidance with paired data

### 4.2.1 Facial expression generation

Table 1: Quantitative comparison of facial expression generation on KDEF dataset. The symbol  $\uparrow$  ( $\downarrow$ ) indicates that the larger (smaller) the value, the better the performance.

Method	SSIM $\uparrow$	PSNR $\uparrow$	MSE $\downarrow$	RMSE $\downarrow$	LPIPS $\downarrow$	FID $\downarrow$
Pix2pix [24]	0.8673	18.0138	0.0160	0.1133	0.2603	140.7901
PG <sup>2</sup> [42]	0.9118	19.2469	<b>0.0122</b>	<b>0.0946</b>	0.2032	97.2435
ReshapeGAN	<b>0.9227</b>	<b>19.7374</b>	0.0123	0.0949	<b>0.1808</b>	<b>84.1437</b>

Table 2: Quantitative comparison of facial expression generation on RaFD dataset. The symbol  $\uparrow$  ( $\downarrow$ ) indicates that the larger (smaller) the value, the better the performance.

Method	SSIM $\uparrow$	PSNR $\uparrow$	MSE $\downarrow$	RMSE $\downarrow$	LPIPS $\downarrow$	FID $\downarrow$
Pix2pix [24]	0.9565	19.0794	0.0130	0.1121	0.1379	105.0758
PG <sup>2</sup> [42]	0.9631	19.6109	0.0115	0.1055	0.1297	102.6674
ReshapeGAN	<b>0.9641</b>	<b>20.1008</b>	<b>0.0103</b>	<b>0.0998</b>	<b>0.1030</b>	<b>49.2096</b>



Figure 6: Visual comparison of facial expression generation on KDEF dataset.



Figure 7: Visual comparison of facial expression generation on RaFD dataset.

First, to evaluate the image generation performance for reshaping, we conduct experiments on facial expression

datasets using KDEF [5] and RaFD [31] datasets. For these two facial datasets, we use paired data to train all the models. KDEF contains 70 different identities with 7 different emotional representations and 5 different poses, while RaFD contains 67 identities with 8 different emotional expressions, 5 different poses and 3 eye gaze directions. For both of the two emotional datasets, we only use the frontal images with neutral faces as input images and generate images with other facial expressions. So we reorganize the two datasets, thus make KDEF (420 totally, 336 for training and 84 for evaluating) and RaFD (469 totally, 392 for training and 77 for evaluating) for facial expression reshaping.

Here we compare our ReshapeGAN with two state-of-the-art supervised methods, i.e., Pix2pix [24] and PG<sup>2</sup> [42] for evaluation, and the quantitative comparison results can be found in Table 1 and Table 2, in terms of SSIM, PSNR, MSE, RMSE, LPIPS and FID, where our ReshapeGAN gets the best performance among the three methods. For ReshapeGAN and PG<sup>2</sup>, we use dlib [26] to obtain the geometric information as guidance. While for Pix2pix, to compare fairly, we encode the emotional or viewpoint expression as one-hot code and inject it into the bottleneck of generator as guidance. More visual results can be found in Fig. 6 and Fig. 7. As it can be seen, Pix2pix can not generate acceptable results, while the outputs of PG<sup>2</sup> have blur boundary and some dirty color blocks. Our ReshapeGAN generates reasonably clear outputs which preserve the identity information of inputs.

To investigate the efficiency of different components in our approach, we design additional experiments on RaFD for ablation study, and the quantitative results are listed in Table 3. It can be seen that the perturbed loss  $\mathcal{L}_{perturb}$  improves the PSNR score dramatically, and the geometric information helps to reduce the LPIPS distance, while the perceptual loss  $\mathcal{L}_{percep}$  reduces the LPIPS and FID, since the perceptual loss actually provides multi-scale constraints to the generator by using a cascade architecture.

#### 4.2.2 Viewpoint transfer

Then, we devise viewpoint transfer task, aiming to generate different pose (viewpoint) faces on FEI dataset [61], which contains 200 identities with 11 different pose directions. Since we can not extract the accurate geometric in-

Table 3: Quantitative comparison for ablation study of our ReshapeGAN on facial expression generation from RaFD dataset. The symbol  $\uparrow$  ( $\downarrow$ ) indicates that the larger (smaller) the value, the better the performance.

Method	SSIM $\uparrow$	PSNR $\uparrow$	MSE $\downarrow$	RMSE $\downarrow$	LPIPS $\downarrow$	FID $\downarrow$
Backbone	0.9632	19.8962	0.0106	0.1017	0.1365	106.0801
Geometry	<b>0.9655</b>	19.7940	0.0109	0.1030	0.1242	89.5522
Geometry	0.9636	<b>20.1368</b>	<b>0.0101</b>	<b>0.0991</b>	0.1241	81.6016
+ $\mathcal{L}_{perturb}$						
Geometry	0.9628	19.5127	0.0116	0.1064	0.1184	61.8619
+ $\mathcal{L}_{percep}$						
ReshapeGAN	<b>0.9641</b>	<b>20.1008</b>	<b>0.0103</b>	<b>0.0998</b>	<b>0.1030</b>	<b>49.2096</b>

Table 4: Quantitative comparison of viewpoint transfer (facial pose translation) on FEI dataset. The symbol  $\uparrow$  ( $\downarrow$ ) indicates that the larger (smaller) the value, the better the performance.

Method	SSIM $\uparrow$	PSNR $\uparrow$	MSE $\downarrow$	RMSE $\downarrow$	LPIPS $\downarrow$	FID $\downarrow$
Pix2pix [24]	0.9274	16.9358	0.0161	0.1228	0.2236	79.3540
DR-GAN [62]	0.9384	17.4995	0.0138	0.1138	0.2001	83.5483
PG <sup>2</sup> [42]	0.9354	17.3366	0.0144	0.1160	0.1887	<b>70.7044</b>
ReshapeGAN	<b>0.9578</b>	<b>19.2299</b>	<b>0.0097</b>	<b>0.0941</b>	<b>0.1721</b>	73.5864

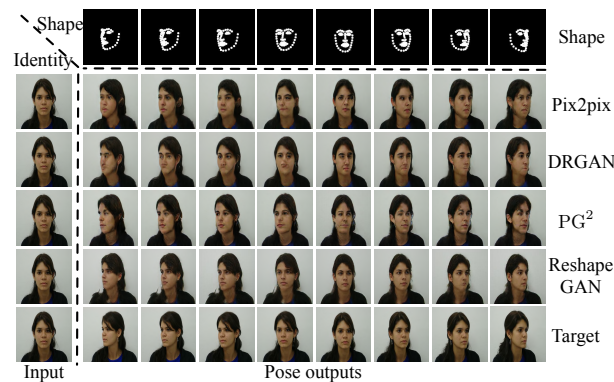


Figure 8: Visual comparison of viewpoint transfer (facial pose translation) on FEI dataset.

formation from the full left or full right images, we don't use them for training. We use frontal images as inputs to

generate other 8 pose images. We compare our ReshapeGAN with Pix2pix [24], PG<sup>2</sup> [42] and DR-GAN [62]. For ReshapeGAN, Pix2pix and PG<sup>2</sup>, we also use the same settings as facial expression generation (Section 4.2.1) for adding guidance to these three methods, and for DR-GAN, we use the default setting of facial pose generation from the authors but using our training data. Besides, we use 160 identities for training and other 40 for testing. The quantitative comparison results are listed in Table 4. Our ReshapeGAN gets the lowest LPIPS and FID values among the four methods. Visual synthesized results are shown in Fig. 8 for comparison. Pix2pix generates images with blur boundary and some dirty color blocks, DR-GAN and PG<sup>2</sup> distort the faces, while our ReshapeGAN synthesizes better results with detailed parts (such as eyes, mouth and nose) as well as reasonable faces. By combining the geometric information, our ReshapeGAN can easily capture the relationships between different parts.

Note that, the geometric information does not limit to landmarks, instance maps and pose skeletons. We can also extract any other accurate geometric information and regard them as additional guidance to obtain better results.

### 4.3 Reshaping by within-domain guidance with unpaired data

#### 4.3.1 Controllable face translation on CelebA

For this task, we aim to achieve the controllable facial reshaping using unpaired data. We conduct the experiments on CelebA dataset [38], which contains approximately 200 thousand images with high diversity. To get precise facial expression, we reorganize this dataset for getting a sub-dataset using dlib [26] to achieve face detection and crop the detected face regions then resize them to  $256 \times 256$ . We use 157,619 images for training and 39,404 images for testing. In theory, our method could generate  $39404 \times 39404$  images on testing stage. Considering the huge consumption of computing resources, we only generate 1,000 images with different geometric representation for one input sample.

Fig. 9 shows that our ReshapeGAN can synthesize 25 plausible images, via 5 different input images, where each row has the appearance and content consistency, while each column has geometric consistency. It can also be seen that, ReshapeGAN can generate images with re-

quired geometric representation by providing a geometric guidance, even if there are no paired samples for training, indicating that ReshapeGAN learns the facial expressions from the entire dataset rather than some specific sample. Moreover, Fig. 10 exhibits the random 96 generated images using different geometric guidance. Note the smaller image framed by red box in the lower right corner of every generated image is the given reference image (for providing geometric information). The synthesized images preserve the appearance and identity of input images, and our ReshapeGAN is able to generate corresponding outputs with emotional and layout/pose information according to the given reference images.

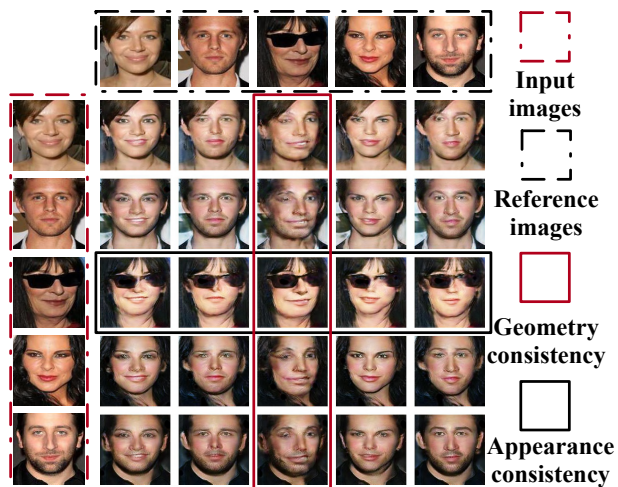


Figure 9: The  $5 \times 5$  outputs by our ReshapeGAN using random 5 images as both inputs and references on CelebA dataset. Images framed by dotted lines are input images and reference images, and each row shows images with the identical appearance (e.g., the row framed by black line) while each column exhibits images with same geometric representation (e.g., the column framed by red line).

To our knowledge, we are the first to achieve arbitrary identity reshaping using unpaired data, with only limited images from each identity. Some traditional methods could perform face swapping, while Deepfakes<sup>2</sup> only translates one identity to another identity using abundant

<sup>2</sup><https://github.com/deepfakes/Faceswap>

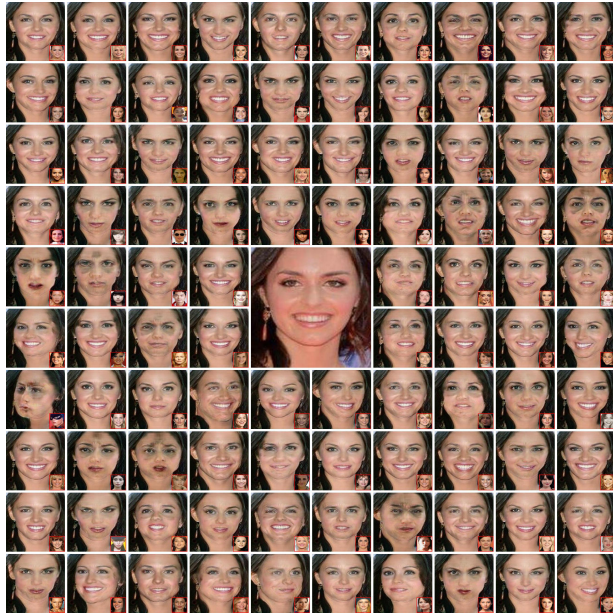


Figure 10: Random 96 synthesized images by our ReshapeGAN from one input sample on CelebA dataset. The middle enlarged image shows the source input image. The smaller images framed by red box in the lower right corner of every generated image is the given reference image. The synthesized images have the same emotional and pose information with given reference images while preserving the appearance and identity information of input images. Please zoom in to see more details.

images from both source identity and target identity, but it fails to translate between multiple identities. Thus, we choose Faceswap from OpenCV<sup>3</sup>) as traditional method for comparison. Besides, we also make comparison using some popular image-to-image translation methods. Due to the lack of paired data, we conduct experiments using Pix2pix and PG<sup>2</sup> by removing the pixel-level loss between outputs and ground truth (denote as Pix2pix-- and PG<sup>2</sup>--) respectively. Here we regard the geometric information as the conditional information and provide the geometric information to generator by concatenating the raw input and geometry, through this way, we conduct experiments using StarGAN [10], and in consideration of

<sup>3</sup><https://opencv.org>

that there is only one domain in CelebA dataset, we remove the classification loss of StarGAN (denote as StarGAN--).

The quantitative comparison is given in Table 5. SSIM and LPIPS scores computed by landmarks show the geometric consistency, higher SSIM and lower LPIPS scores indicate better geometric consistency with reference images. FID computes the distance between generated sample distribution and real reference images, and lower FID score represents better image generation quality. Our ReshapeGAN performs best in terms of SSIM, LPIPS and FID. We also use average identification distance to evaluate the identity similarity, Faceswap fails to preserve the identity with distance larger than 0.6, and the other methods can keep the identity information while our ReshapeGAN performs better on reshaping with larger distance (but below 0.6). The user study of identity/shape judgement can also arrive at similar conclusion that our ReshapeGAN works well considering both identity preservation and object reshaping.

Visual comparison of different methods can be found in Fig. 11. As shown, Faceswap is limited to only borrow the face appearance from target reference and merging it with the corresponding region without considering the relationship between parts, and the output images seem implausible and don't preserve the identity information of the inputs. Pix2pix-- and PG<sup>2</sup>-- fail to reshape the input images and generate images with artifacts. And StarGAN-- also fails to synthesize images with required shape. Compared to these methods, our ReshapeGAN achieves face reshaping by providing one single reference image. The visual comparison of Fig. 11 provides the same clues as quantitative comparison shown in Table. 5. To show the powerful performance of our ReshapeGAN, we exhibit more synthesized results in Fig. 12. We organize these generated results sorted by progressive increase of identification distance with source input image in Fig. 12(a), and the right image show the visualization results of identification distance, which shows the relationship between the distance and the identity preservation with object reshaping.

Based on the above synthesized results, we see that not all of the generated images can borrow the effective information reasonably, some generated images have extremely exaggerated regions. More interestingly, if the geometric information of given reference image is from

Table 5: Quantitative comparison of controllable face translation on CelebA dataset. Higher SSIM and lower LPIPS scores represent better geometric matching with guided geometric information. Lower FID score means better image quality. Identification distance larger than 0.6 shows different identities, while larger distance below 0.6 indicates better reshaping. Higher votes to both identity and shape of user study indicate better reshaping performance. Please refer to Section 4.1 for details about evaluation metrics.

Method	SSIM (Landmark)	LPIPS (Landmark)	FID	Identification Distance	Identity / Shape (User Study)
Faceswap	0.6949	0.1775	111.7921	<u>0.7230</u>	0.025 / 0.605
Pix2pix--	0.6049	0.2561	338.4090	0.4419	0.771 / 0.093
PG <sup>2</sup> --	0.6605	0.2097	194.1435	0.3369	1.0 / 0.102
StarGAN--	0.6544	0.2156	157.4540	0.2576	1.0 / 0.075
ReshapeGAN	0.8332	0.0932	110.7313	0.5700	0.385 / 0.719

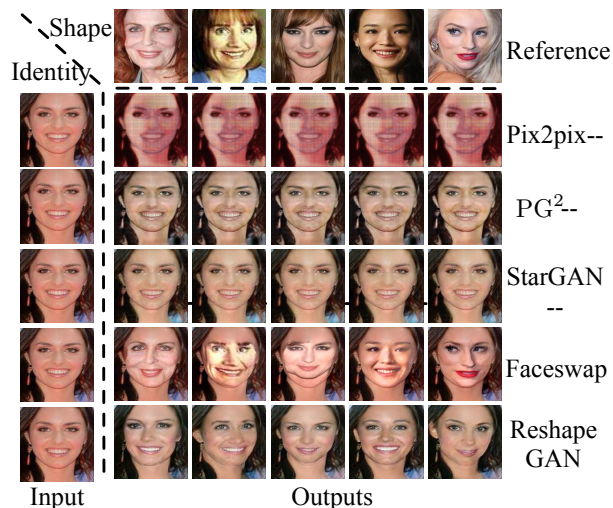


Figure 11: Visual comparison of controllable face translation using different methods on CelebA dataset.

a man while the input image depicts a woman, there is a chance that some sort of masculine pattern will be generated. We guess that the provided geometric information may somehow carry the gender and age characteristics. If the pose distance between input image and given reference image is too large, our model fails to generate plausible outputs. To obtain more intuitive results, we compute the LPIPS distance between each reference image and the input image, and visualize the results in Fig. 12(b). Here we reorganize the generated images sorted by progressive increase of LPIPS distance between target reference im-

age and source input image. As it can be seen, ReshapeGAN performs well if the LPIPS is less than an implicit threshold. However, when the LPIPS distance is too large, ReshapeGAN is not able to translate the input image to the given target geometric space with limited prior information.

### 4.3.2 Controllable face translation on UTKFace

Then, to evaluate whether the geometric information carries underlying information such as gender and age characteristics, we conduct the face translation experiments on another large dataset UTKFace [77]. This dataset contains over 20,000 face images with long age span (range from 0 to 116 years old). In this task, we also take experiments to generate images with arbitrary geometric information. According to our above assumption, the geometric information may contain underlying scale, position and facial expression information. By providing information to the generator and discriminator, they can build a mapping function between input images and geometric information. Visual results are shown in Fig. 13, we can see that the generated images have a large range of ages. We also visualize the identification distance, from which, we can see that the given geometric information does covers intrinsic characteristics. In order to fool the discriminator, the generator should dig in the implicit information and express it, reasonably.

For this task, we also conduct some comparative experiments using aforementioned methods. The comparative visual results are shown in Fig. 14. We can see that the

Table 6: Quantitative comparison of controllable face translation on UTKFace dataset. Higher SSIM and lower LPIPS scores represent better geometric matching with guided geometric information. Lower FID score means higher image quality. Identification distance larger than 0.6 shows different identities, while larger distance below 0.6 indicates better reshaping. Higher votes to both identity and shape of user study indicate better reshaping performance. Please refer to Section 4.1 for details about evaluation metrics.

Method	SSIM (Landmark)	LPIPS (Landmark)	FID	Identification Distance	Identity / Shape (User Study)
Faceswap	0.5595	0.2404	120.09001	<u>0.7354</u>	0.045 / 0.807
Pix2pix--	0.6063	0.2197	251.4699	0.4144	0.786 / 0.188
PG <sup>2</sup> --	0.5197	0.3127	173.5819	0.2628	1.0 / 0.024
StarGAN--	0.5501	0.2793	186.4931	0.1850	0.991 / 0.048
ReshapeGAN	0.8103	0.0849	103.8411	0.5902	0.494 / 0.850

synthesized images using Pix2pix--, PG<sup>2</sup>-- and StarGAN-- do not achieve the facial reshaping. Comparing with Faceswap, our ReshapeGAN can preserve the identity information better. Quantitative comparison of different methods are given in Table 6, as can be seen, our ReshapeGAN performs best with highest SSIM, lowest LPIPS and FID, indicating best geometric matching with highest image quality. Besides, ReshapeGAN also gets largest identification distance compared to all other methods below 0.6. Meanwhile, the reshaping ability with identity preservation of ReshapeGAN can also be validated by user study.

### 4.3.3 Controllable cat reshaping

Moreover, we implement our ReshapeGAN on an interesting application that is to reshape a cat by providing geometric information. We use data and annotations from the cat head dataset [76], and each annotation file describes the coordinates of 9 defined points (6 for ears, 2 for eyes and 1 for mouth). Here we use 7,287 cropped images (according to the coordinates) for training and 1,000 images for testing. Experimental results are shown in Fig. 15, our ReshapeGAN can reshape the cat image with a given geometric information, while keeping the detailed information of input image. For this task, we only provide the limited 9 pointed annotation, which describes the location of ears, eyes and mouth. So, without sufficient guidance, our method can also work well. For this task, we only compute the FID between generated cat images and real images for reference, and the FID score is 76.7446.

### 4.3.4 Controllable pose reshaping

In addition, the skeleton guidance of human body provides more detailed pose information. Thus, we use pose estimation model OpenPose [6] to extract pose information of human images and regard them as geometric guidance for pose reshaping application. Here we perform our ReshapeGAN on Panama dataset<sup>4</sup>. Note that, we don't use the paired training for this task, and only add the cycle-consistency constrains to our model. By fusing pose information and appearance information from given reference images, our method can generate reasonable results as shown in Fig. 16. Here we only generate images by frames, we leave the pose sequence generation by adding spatial and temporal constrains as our future work. For this task, we also only compute the FID between generated pose images and real images for reference, and the FID score is 145.3365.

### 4.3.5 Ablation study

On the one hand, to investigate the architecture of our geometry-guided method, we compare three cases for generator and discriminator without or with geometric information: 1) only add geometric information to generator, 2) only add geometric information to discriminator, 3) add geometric information to both generator and discriminator. We conduct experiments on CelebA dataset for this ablation study in the three cases, and the visual results are shown in Fig. 17(a). Discriminator without geometric

<sup>4</sup>[https://github.com/llSourcecell/Everybody\\_Dance\\_Now](https://github.com/llSourcecell/Everybody_Dance_Now)

Table 7: Quantitative comparison of different cases for our ReshapeGAN without or with geometric information on CelebA dataset. Only adding geometric information to generator  $G$  (Only G) and discriminator  $D$  (Only D) get very low identification distance, showing that they fail to achieve geometric reshaping with given geometric constrain, that is, they fall into the trivial solution to generate the very similar output to the input image (very lower identification distance). The user study also validates the same point (higher votes to both identity and shape of user study indicate better reshaping performance). Our ReshapeGAN, by adding geometric information to both  $G$  and  $D$ , gets higher SSIM and lower LPIPS scores, indicating to reshape the input image with given geometric information better, and also gets the lower FID score representing higher image quality.

Method	SSIM (Landmark)	LPIPS (Landmark)	FID	Identification Distance	Identity / Shape (User Study)
Only G	0.6763	0.1855	168.2893	0.0438	1.0 / 0.1739
Only D	0.6543	0.1924	176.5256	0.0352	1.0 / 0.073
G+D (ReshapeGAN)	0.8332	0.0932	110.7313	0.5700	0.385 / 0.719

Table 8: Quantitative comparison of different cases for our ReshapeGAN without or with perceptual loss on CelebA dataset. Although the model without the perceptual loss can get better performance in terms of FID, LPIPS and SSIM, the identification distance is larger than 0.6. Thus the perceptual loss does helps to preserve the identity and appearance information from source input image.

Method	SSIM (Landmark)	LPIPS (Landmark)	FID	Identification Distance	Identity / Shape (User Study)
w/o $\mathcal{L}_{percep}$	0.8683	0.0706	84.6713	<u>0.6075</u>	0.145 / 0.879
w/ $\mathcal{L}_{percep}^{gr}$	0.8509	0.0774	99.9562	<u>0.6188</u>	0.186 / 0.934
w/ $\mathcal{L}_{percep}^{gi}$	0.8332	0.0932	110.7313	0.5700	0.385 / 0.719

Table 9: Quantitative comparison of different cases for our ReshapeGAN without or with perturbed loss on CelebA dataset. The perturbed loss can reduce the identification distance and reserve the source identity information while helping improve image quality.

Method	SSIM (Landmark)	LPIPS (Landmark)	FID	Identification Distance	Identity / Shape (User Study)
w/o $\mathcal{L}_{perturb}$	0.8473	0.0827	115.4396	0.5790	0.256 / 0.896
W/ $\mathcal{L}_{perturb}$	0.8332	0.0932	110.7313	0.5700	0.385 / 0.719

information guidance captures the pose information but fails to provide effective gradient direction to generator, so that the model can not generate required controllable images. If we don't provide the given geometric information to generator, the model can not achieve controllable reshaping. Thus the guidance is necessary for both generator and discriminator. Table 7 with quantitative comparison also concludes this point.

On the other hand, we explore the efficiency of perceptual loss by considering three cases: 1) without perceptual loss, 2) with perceptual loss computed between synthe-

sized fake image and real reference image ( $\mathcal{L}_{percep}^{gr}$ ), 3) with perceptual loss computed between synthesized fake image and real input image ( $\mathcal{L}_{percep}^{gi}$ ). Visual results can be found in Fig. 17(b), and quantitative results are listed in Table 8. From the results, we can see that the model falls in trivial solution and encounters over-fitting problem by adding constrains between generated images and given target reference images, that is, the generated images are very similar to the given reference images. Meanwhile, computing the perceptual loss between synthesized images and input images can guarantee content and iden-

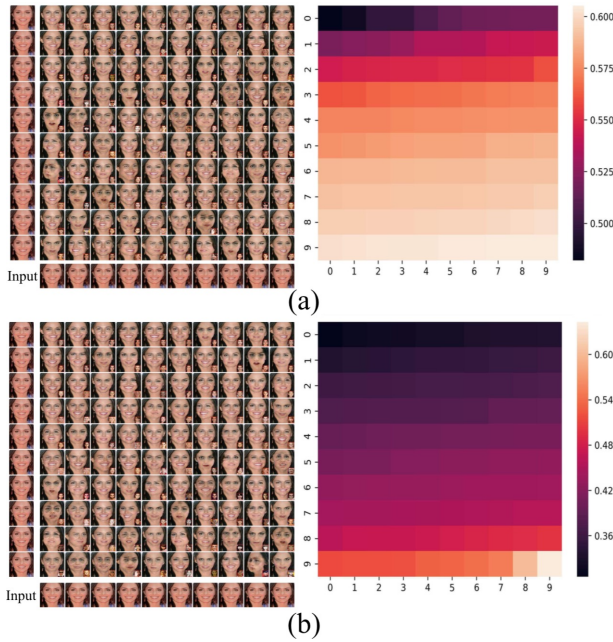


Figure 12: (a) The random 100 generated results (left) sorted by progressive increase of identification distance (right) between the generated image and the source input image on CelebA dataset. (b) The random 100 generated results (left) sorted by progressive increase of LPIPS distance (right) between the corresponding target reference image and the source input image on CelebA dataset. We can see that the generated image with lower identification distance has higher similarity, while the generated image with lower LPIPS distance has higher quality. Please zoom in to see more details.

tity consistency. Comparing with L1 loss, perceptual loss focuses on high-level semantic matching between images rather than pixel-level matching.

Furthermore, we also devise similar experiments to explore the effectiveness of the perturbed loss, the results are shown in Table 9 and Fig. 17(c). As shown, the perturbed loss helps to preserve the identity information and improves the generation quality. Although the model without perturbed loss can get a little higher SSIM and lower LPIPS scores (computed between geometric information), we should pay more attention to preserve the appearance consistency and image generation quality.

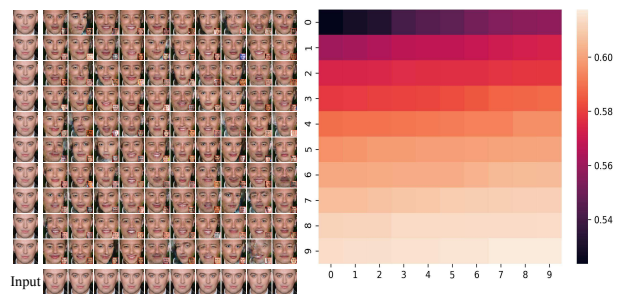


Figure 13: The random 100 generated results (left) sorted by progressive increase of identification distance (right) between the generated image and the source input image on UTKFace dataset. Please zoom in to see more details.

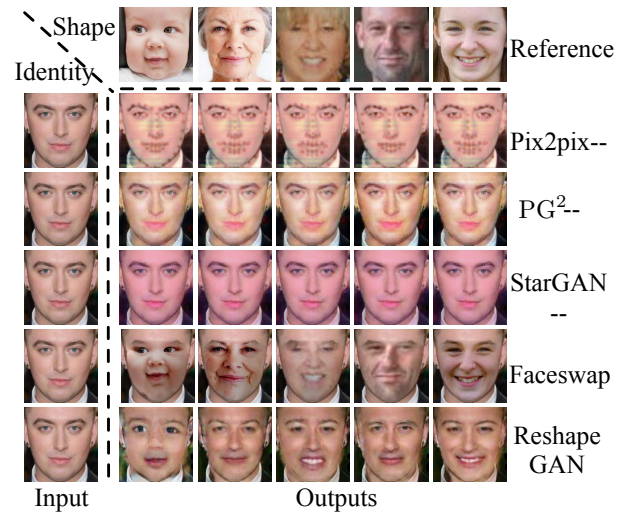


Figure 14: Visual comparison of controllable face translation using different methods on UTKFace dataset.

## 4.4 Reshaping by cross-domain guidance

### 4.4.1 Facial reshaping across multiple datasets

In this task, we use 5 facial image datasets including Yale [18], WSEFEP [49], ADFES [63], KDEF [5] and RaFD [31] for facial reshaping.

First, considering domain gaps between different datasets, we devise three different experiments for evaluating the efficiency of our ReshapeGAN. First, we conduct experiments on RaFD and KDEF datasets, and we



Table 10: Quantitative comparison of controllable face translation on KDEF-RaFD dataset. Higher SSIM and lower LPIPS scores represent better geometric matching with guided geometric information. Lower FID score means higher image quality. Identification distance larger than 0.6 shows different identities, while larger distance below 0.6 indicates better reshaping. Please refer to Section 4.1 for details about evaluation metrics. The results above the dashed line are computed in the first case: the input image is from KDEF and the reference image is from RaFD, while the results below the dashed line are computed in the second case: the input image is from RaFD and the reference image is from KDEF.

Method	SSIM (Landmark)	LPIPS (Landmark)	FID	Identification Distance
CycleGAN [80]	0.6181	0.2608	172.8081	0.4591
MUNIT [23]	0.6083	0.2754	182.4521	0.4721
DRIT [35]	0.6949	0.2692	115.8099	0.5427
StarGAN [10]	0.6623	0.2493	130.7385	0.0612
ReshapeGAN	0.7981	0.1184	100.8653	0.5909
CycleGAN [80]	0.6323	0.2372	140.2661	0.4602
MUNIT [23]	0.6315	0.2482	154.2534	0.4654
DRIT [35]	0.6521	0.1932	84.2044	0.6361
StarGAN [10]	0.6256	0.2503	100.9369	0.0489
ReshapeGAN	0.7895	0.1222	122.9286	0.5930

choose CycleGAN, MUNIT, DRIT and StarGAN for comparison. In order to compare fairly, we provide the geometric information to all the generators as the additional constrains for all compared methods. Additionally, for StarGAN, we use the one-hot encoding to perform the facial translation on the two different datasets. By combining all the emotional categories appearing on both two datasets, we get 7 different emotional labels (contemptuous, disgusted, fearful, happy, sad, and surprised) and conduct experiments following [10]. For the testing stage, we randomly select 91 input-reference pairs from the testing images of the two domains. Please note that there are two cases for the cross-domain reshaping on KDEF-RaFD datasets: 1) the input image is from KDEF and the reference is from RaFD, 2) the input image is from RaFD and the reference image is from KDEF. The visual synthesized images are exhibited in Fig. 18 and Fig. 19(a). CycleGAN and MUNIT only achieve the domain adaption according to the given reference image, and they fail to perform reshaping by combining the geometric information. DRIT generates images with underlying position changes, but it also fails to reshape input images with given reference images. Above all, recent unpaired cross-domain methods among two different domains (CycleGAN, MUNIT

and DRIT) can achieve domain adaption (including the low-level texture and background translation), but they pay more attention to the low-level matching rather than the high-level geometrical matching, that is, they fail to generate images with required shape by providing the geometric guidance. StarGAN fails to achieve the facial translation by providing additional emotional label, we guess the reason is that the same two datasets contain similar emotional expression but different background information and data distribution, so that the model is hard to focus on the emotional representation. Compared to above methods, our ReshapeGAN can achieve the facial reshaping by only providing a single reference image. For the quantitative comparison, Table 10 lists all the results. As shown, our method gets highest SSIM and lowest LPIPS scores, which indicates that ReshapeGAN achieves geometric consistency with reference images. At the same time, our method can preserve the appearance information (including the background and domain information) well in terms of identification distance.

Second, we perform our method on WSEFEP, ADFES and Yale datasets, where the three datasets have small intra-domain distances, and the visual results are shown in Fig. 19(b). For the above two cases, our method can per-



Figure 15: Random 96 synthesized images by our ReshapeGAN from one input sample on cat head dataset. The middle enlarged image shows the source input image. And the smaller image framed by red box in the lower right corner of every generated image is the given target reference image. Please zoom in to see more details.

form well on both object reshaping and generation quality.

Finally, we conduct experiments on all above five emotional datasets, and the visual results can be seen in Fig. 19(c). We observe that the synthesized images have dirty color blocks if the domain gap is large (such as KDEF and WSEFEP), since the samples from the two datasets have extremely different location as well as pose information, such that our method can not generate reasonable outputs without sufficient prior information.

#### 4.4.2 Caricature reshaping across multiple datasets

For this task, we aim to achieve caricature translation between multiple datasets. The datasets include CUHK Face Sketch database [65], KDEF [5], IIIT-CFW [45] and PHOTO-SKETCH [66]. Among the four datasets, we could get 5 different image style domains (there are 2 domains from PHOTO-SKETCH dataset). For each dataset,

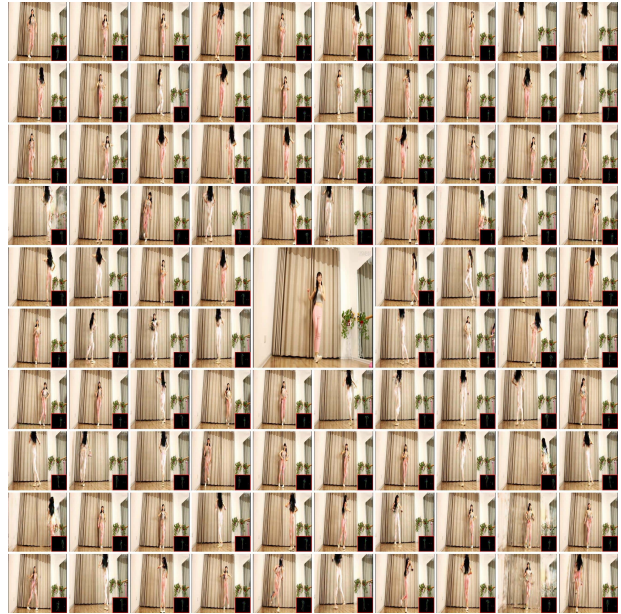


Figure 16: Random 96 synthesized images by our ReshapeGAN from one input sample on Panama dataset. The middle enlarged image shows the source input image. And the smaller image framed by red box in the lower right corner of every generated image is the given pose reference image extracted using OpenPose. We could generate a different video sequence with poses from target identity while using appearance and identity information from source inputs. Please zoom in to see more details.

we randomly select approximately four fifths of samples for training and others for testing. First, we conduct comparative experiments on PHOTO-SKETCH dataset. Due to the lack of conditional label, here we don't use StarGAN for comparison and follow the training/testing splitting in [71]. Using the same testing criteria, we randomly get 199 input-reference images. Note that the input-reference paired images are not from the same identity. The visual synthesized results can be seen in Fig. 20. As shown, our ReshapeGAN can reshape the input images according to the shape of reference images, while CycleGAN, MUNIT and DRIT only achieve the domain transfer between two different image manifolds, failing to find the geometric mapping function. The quantitative comparison is listed in Table 11, from which, our Reshape-

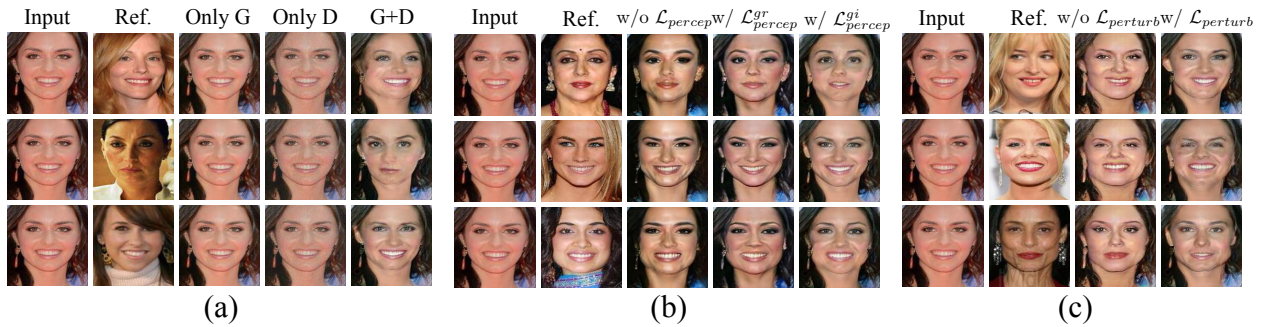


Figure 17: Visual synthesized results of different cases without or with geometric information (a), of different cases without or with perceptual loss (b), of different cases without or with perturbed loss (c), on CelebA dataset.

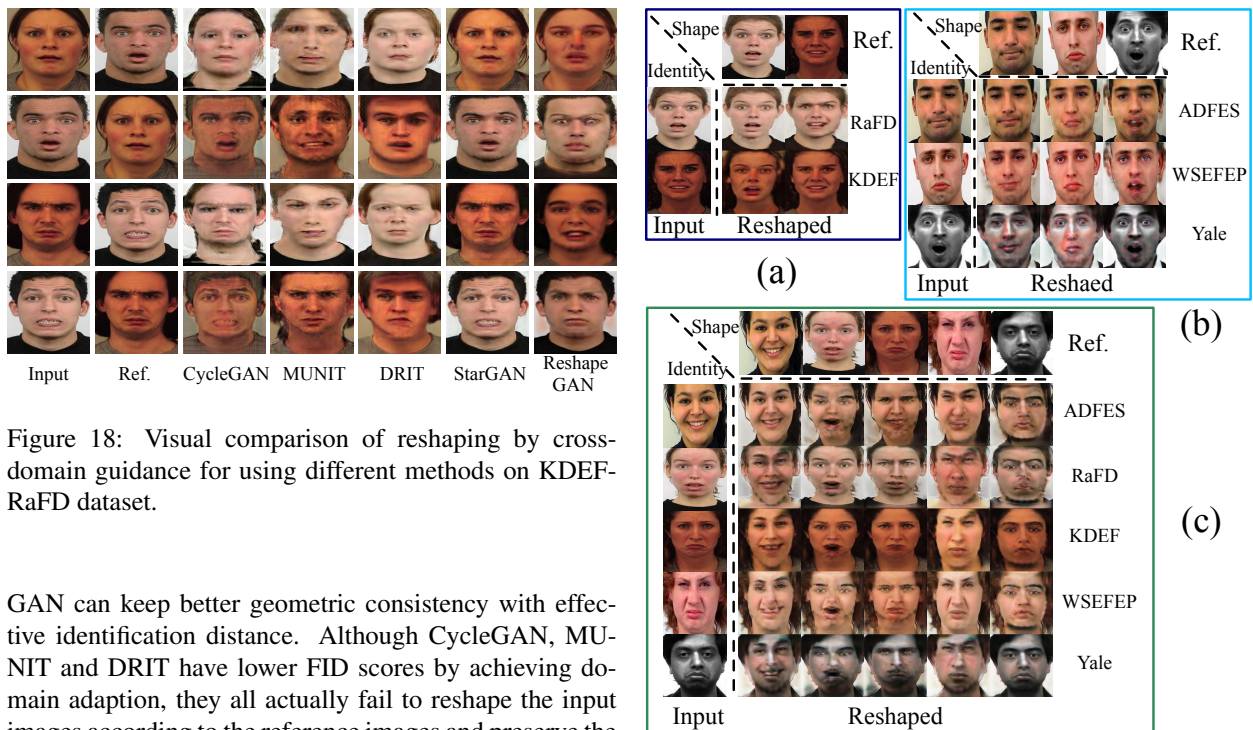


Figure 18: Visual comparison of reshaping by cross-domain guidance for using different methods on KDEF-RaFD dataset.

GAN can keep better geometric consistency with effective identification distance. Although CycleGAN, MUNIT and DRIT have lower FID scores by achieving domain adaption, they all actually fail to reshape the input images according to the reference images and preserve the appearance and style information of input images.

Actually, in the testing stage, we can generate arbitrary domain image with arbitrary geometric guidance, and we can also generate  $n \times n$  images based on  $n$  images from 5 image domains for this task. To reduce the computational cost, we randomly assemble 5 images from 5 image domains to obtain a combination. So we can get  $n = \max n_i$  combinations, where  $n_i$  represents the num-

Figure 19: Visual results synthesized by our ReshapeGAN on RaFD and KDEF dataset (a), on WSEFEP, ADFES, and Yale (b), on above five datasets (c), for facial reshaping.

ber of testing images from the 5 image domains separately. In order to evaluate the performance of our Re-

Table 11: Quantitative comparison of controllable face translation on PHOTO-SKETCH dataset. Higher SSIM and lower LPIPS scores represent better geometric matching with guided geometric information. Lower FID score means higher image quality. Identification distance larger than 0.6 shows different identities, while larger distance below 0.6 indicates better reshaping. Please refer to Section 4.1 for details about evaluation metrics. The results above the dashed line are computed in the first case: the input image is from PHOTO domain and the reference image is from SKETCH domain, while the results below the dashed line are computed in the second case: the input image is from SKETCH domain and the reference image is from PHOTO domain.

Method	SSIM (Landmark)	LPIPS (Landmark)	FID	Identification Distance
CycleGAN [80]	0.4839	0.3567	45.6227	<u>0.6320</u>
MUNIT [23]	0.4623	0.3462	89.6343	<u>0.6103</u>
DRIT [35]	0.4823	0.3514	67.8397	<u>0.6420</u>
ReshapeGAN	0.8088	0.0871	125.4597	0.5749
-----				
CycleGAN [80]	0.4786	0.3703	77.2058	<u>0.6038</u>
MUNIT [23]	0.4823	0.3643	84.2434	<u>0.6243</u>
DRIT [35]	0.4927	0.3513	65.3231	<u>0.6182</u>
ReshapeGAN	0.7931	0.1123	101.5948	0.5818

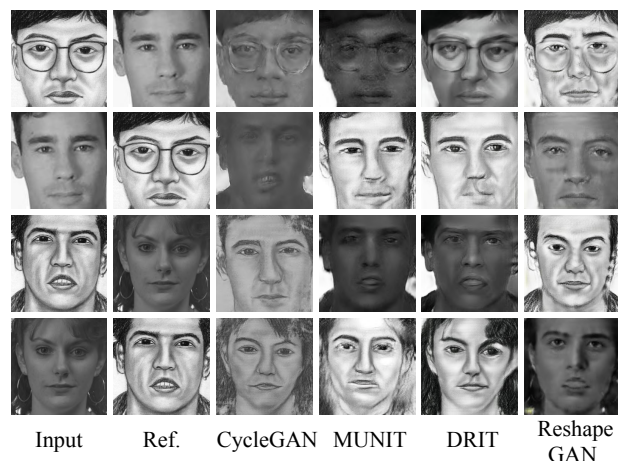


Figure 20: Visual comparison of reshaping by cross-domain guidance for using different methods on PHOTO-SKETCH dataset.

shapeGAN on cross-domain reshaping, we compute the FID scores between generated images and real images. Here we compute the FID scores of 4 different cases: 1) between synthesized images without the domain classification loss  $\mathcal{L}_{cls}(D_d)$  and reference images (Without CLS for abbreviation), 2) between synthesized images opti-

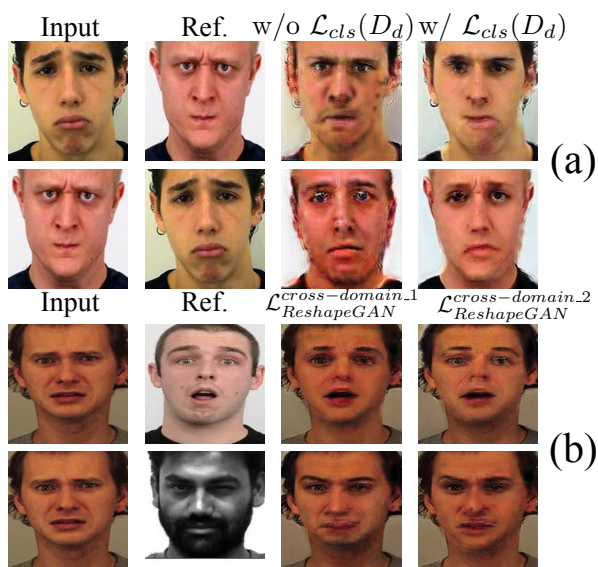


Figure 21: Visual comparison results generated by our ReshapeGAN of different cases without or with  $\mathcal{L}_{cls}(D_d)$  (a), of different cases with  $\mathcal{L}_{ReshapeGAN}^{cross-domain.1}(G, D, D_d)$  or with  $\mathcal{L}_{ReshapeGAN}^{cross-domain.2}(G, D, D_d)$  (b) on selected 5 facial datasets.

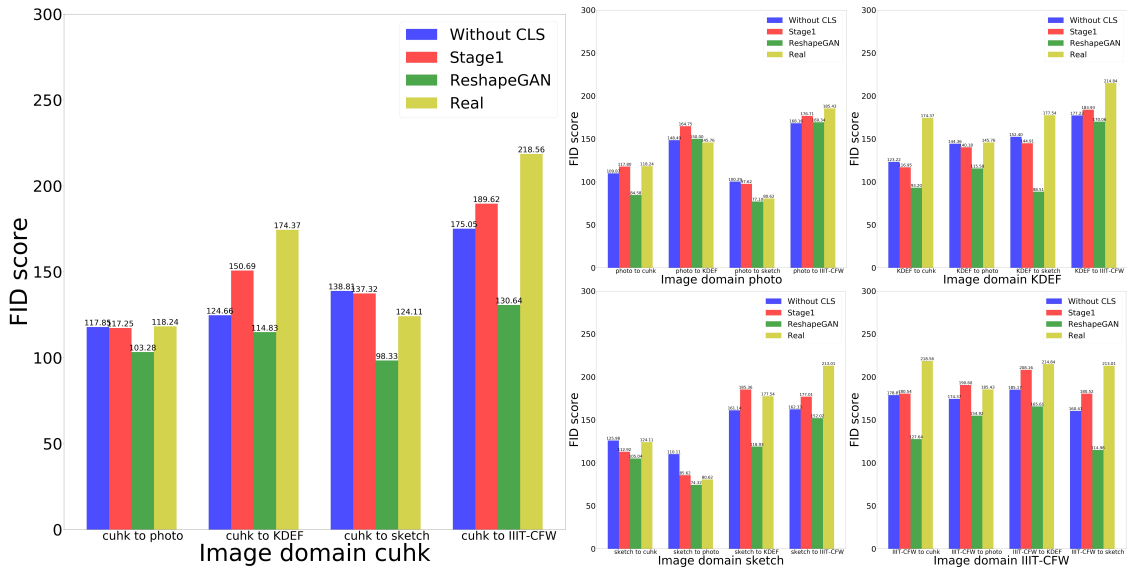


Figure 22: The visualization results of of FID distance using different methods on four datasets for caricature reshaping.

mized by  $\mathcal{L}_{ReshapeGAN}^{cross-domain.1}(G, D, D_d)$  and reference images (Stage1 for abbreviation), 3) between synthesized images optimized by  $\mathcal{L}_{ReshapeGAN}^{cross-domain.2}(G, D, D_d)$  and reference images (ReshapeGAN for abbreviation), 4) between source input images and reference images (Real for abbreviation). We visualize all the FID results in Fig. 22. Each sub figure shows the translation from a source domain to another 4 target domains. We see that the FID score of the model without domain classification loss is higher than the model with domain classification loss (ReshapeGAN). So the domain classification loss does helps to improve the image generation quality. Besides, if we split a difficult cross-domain object reshaping task to two easier tasks, we can get better generation performance. Almost all the FID scores of proposed model are lower than those of other models using only one stage. That is, the second refining stage does improve the image generation quality. More ablation study can be found in Section 4.5. With geometric information from target domain, we see that the FID between generated images and reference images is lower than FID between source input images and reference images. Thus our ReshapeGAN can reduce domain distances by combining the geometric in-

formation.

## 4.5 Ablation study

In order to show the effectiveness of our pipeline architecture, we design another ablation study with one-stage generation architecture by combining domain classification loss  $\mathcal{L}_{cls}(D_d)$  and geometric information to a strong conditional input. Fig. 21(a) shows the visual comparison results, from which, we can see that that the one-stage generation method without  $\mathcal{L}_{cls}(D_d)$  can not preserve the identity and domain information well, and the generation is not sharp while the boundary is not clear. And with  $\mathcal{L}_{cls}(D_d)$ , the model generates images without dirty color blocks. In addition, we also explore the effectiveness of the second refining stage, the visual comparison is exhibited in Fig. 21(b). As shown, the generated images with refining stage have better geometric changes according to given guidance.

## 5 Failure and limitation

In this section, we will discuss the limitation of our ReshapeGAN. Our model requires extra geometric information as an input, e.g., we use dlib [26] to obtain face landmark as geometric information. So the face landmark should be exist and accurate, and sometimes we have to abandon the face images that dlib cannot extract landmarks (e.g., the method fails to obtain geometric information for the full left and right images on FEI dataset). Obviously, it will affect the performance of our model if the geometric information is wrong. For using the unsupervised manner in our experiments to achieve the cross-dataset image reshaping, we still lack prior knowledge when the semantic domain gap is huge between different image domains. Fig. 23 shows some failure cases of ReshapeGAN, where our model fails to preserve the appearance of input images well, and the outputs look like just the simple combination of source image and target image. Visually, our ReshapeGAN fails to reshape the input image to the reference geometric shape if the distance between the inputs and the references is too large, since our model has insufficient corresponding prior information to achieve reasonable translation. In this case, we will try to disentangle the content representation to make the content information preserved, and we leave this as our future work.

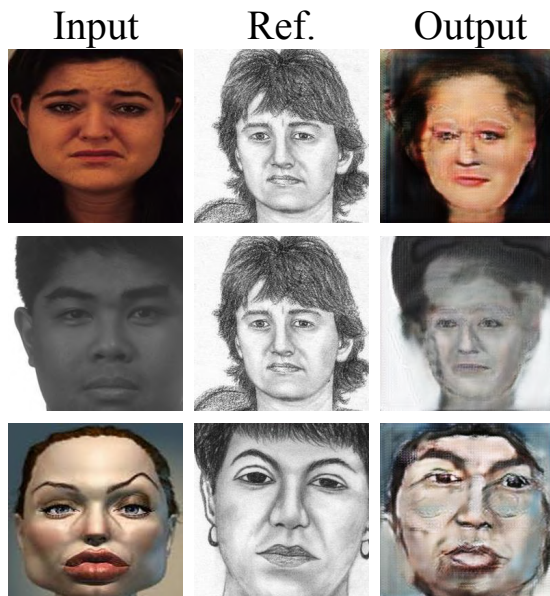


Figure 23: The failure cases on object reshaping by our ReshapeGAN across multiple datasets.

## 6 Conclusion

In this paper, we proposed a new architecture that can generate domain-specific or cross-domain images with geometric guidance, which is the first general framework for a wide range of object reshaping by providing a single reference. Comprehensive experiments on both ablation study and comparisons with comparable state-of-the-art models, in terms of all kinds of applicable quantitative and qualitative as well as human subjective evaluation metrics, show that our model performs better for identity preserved object reshaping.

## Acknowledgment

This work was supported in part by the Royal Society under IEC\R3\ 170013 - International Exchanges 2017 Cost Share (Japan and Taiwan only), in part by a Microsoft Research Asia (MSRA) Collaborative Research 2019 Grant, and in part by the National Natural Science Foundation of China under grants 61771440 and 41776113.

## References

- [1] Antipov, G., Baccouche, M., Dugelay, J.L.: Face aging with conditional generative adversarial networks. arXiv preprint arXiv:1702.01983 (2017)
- [2] Arjovsky, M., Chintala, S., Bottou, L.: Wasserstein generative adversarial networks. In: International Conference on Machine Learning, pp. 214–223 (2017)
- [3] Bau, D., Zhu, J.Y., Strobel, H., Zhou, B., Tenenbaum, J.B., Freeman, W.T., Torralba, A.: GAN dissection: Visualizing and understanding generative adversarial networks. In: International Conference on Learning Representation (2019)
- [4] Borji, A.: Pros and cons of gan evaluation measures. *Computer Vision and Image Understanding* **179**, 41–65 (2019)
- [5] Calvo, M.G., Lundqvist, D.: Facial expressions of emotion (kdef): Identification under different display-duration conditions. *Behavior Research Methods* **40**(1), 109–115 (2008)
- [6] Cao, Z., Simon, T., Wei, S.E., Sheikh, Y.: Realtime multi-person 2d pose estimation using part affinity fields. In: IEEE Conference on Computer Vision and Pattern Recognition, pp. 7291–7299 (2017)
- [7] Chen, Q., Koltun, V.: Photographic image synthesis with cascaded refinement networks. In: IEEE International Conference on Computer Vision, pp. 1511–1520 (2017)
- [8] Chen, X., Duan, Y., Houthoofd, R., Schulman, J., Sutskever, I., Abbeel, P.: InfoGAN: Interpretable representation learning by information maximizing generative adversarial nets. In: Advances in Neural Information Processing Systems, pp. 2172–2180 (2016)
- [9] Chen, Z., Tong, Y.: Face super-resolution through Wasserstein GANs. arXiv preprint arXiv:1705.02438 (2017)
- [10] Choi, Y., Choi, M., Kim, M., Ha, J.W., Kim, S., Choo, J.: StarGAN: Unified generative adversarial networks for multi-domain image-to-image translation. In: IEEE Conference on Computer Vision and Pattern Recognition, pp. 8789–8797 (2018)
- [11] Creswell, A., White, T., Dumoulin, V., Arulkumaran, K., Sengupta, B., Bharath, A.A.: Generative adversarial networks: An overview. *IEEE Signal Processing Magazine* **35**(1), 53–65 (2018)
- [12] Denton, E.L., Chintala, S., Fergus, R., et al.: Deep generative image models using a laplacian pyramid of adversarial networks. In: Advances in Neural Information Processing Systems, pp. 1486–1494 (2015)
- [13] Devon Hjelm, R., Jacob, A.P., Che, T., Trischler, A., Cho, K., Bengio, Y.: Boundary-seeking generative adversarial networks. arXiv preprint arXiv:1702.08431 (2017)
- [14] Donahue, J., Krähenbühl, P., Darrell, T.: Adversarial feature learning. arXiv preprint arXiv:1605.09782 (2016)
- [15] Dosovitskiy, A., Brox, T.: Generating images with perceptual similarity metrics based on deep networks. In: Advances in Neural Information Processing Systems, pp. 658–666 (2016)
- [16] Dumoulin, V., Belghazi, I., Poole, B., Mastropietro, O., Lamb, A., Arjovsky, M., Courville, A.: Adversarially learned inference. arXiv preprint arXiv:1606.00704 (2016)
- [17] Gatys, L.A., Ecker, A.S., Bethge, M.: A neural algorithm of artistic style. arXiv preprint arXiv:1508.06576 (2015)
- [18] Georghiades, A., Belhumeur, P., Kriegman, D.: Yale face database. Center for computational Vision and Control at Yale University, <http://cvc.yale.edu/projects/yalefaces/yalefa> 2 (1997)
- [19] Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. In: Advances in Neural Information Processing Systems, pp. 2672–2680 (2014)
- [20] Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., Courville, A.C.: Improved training of Wasserstein GANs. In: Advances in Neural Information Processing Systems, pp. 5767–5777 (2017)
- [21] Hore, A., Ziou, D.: Image quality metrics: Psnr vs. ssim. In: International Conference on Pattern Recognition, pp. 2366–2369. IEEE (2010)
- [22] Huang, H., Yu, P.S., Wang, C.: An introduction to image synthesis with generative adversarial nets. arXiv preprint arXiv:1803.04469 (2018)
- [23] Huang, X., Liu, M.Y., Belongie, S., Kautz, J.: Multimodal unsupervised image-to-image translation. In: European Conference on Computer Vision, pp. 172–189 (2018)
- [24] Isola, P., Zhu, J.Y., Zhou, T., Efros, A.A.: Image-to-image translation with conditional adversarial networks. In: IEEE conference on Computer Vision and Pattern Recognition, pp. 1125–1134 (2017)
- [25] Kim, T., Cha, M., Kim, H., Lee, J.K., Kim, J.: Learning to discover cross-domain relations with generative adversarial networks. In: International Conference on Machine Learning, pp. 1857–1865. JMLR.org (2017)
- [26] King, D.E.: Dlib-ml: A machine learning toolkit. *Journal of Machine Learning Research* **10**(Jul), 1755–1758 (2009)
- [27] Kingma, D.P., Welling, M.: Auto-encoding variational Bayes. arXiv preprint arXiv:1312.6114 (2013)
- [28] Kodali, N., Abernethy, J., Hays, J., Kira, Z.: On convergence and stability of gans. arXiv preprint arXiv:1705.07215 (2017)
- [29] Kurach, K., Lucic, M., Zhai, X., Michalski, M., Gelly, S.: The GAN landscape: Losses, architectures, regularization, and normalization. arXiv preprint arXiv:1807.04720 (2018)
- [30] Lample, G., Zeghidour, N., Usunier, N., Bordes, A., Denoyer, L., et al.: Fader networks: Manipulating images by sliding attributes. In: Advances in Neural Information Processing Systems, pp. 5967–5976 (2017)
- [31] Langner, O., Dotsch, R., Bijlstra, G., Wigboldus, D., Hawk, S., van Knippenberg, A.: Presentation and validation of the radboud faces database. *Cognition and Emotion* **24**, 1377–1388 (2010)
- [32] Larsen, A.B.L., Sønderby, S.K., Larochelle, H., Winther, O.: Autoencoding beyond pixels using a learned similarity metric. In: International Conference on Machine Learning, pp. 1558–1566 (2016)
- [33] Lassner, C., Pons-Moll, G., Gehler, P.V.: A generative model of people in clothing. In: IEEE International Conference on Computer Vision, pp. 853–862 (2017)

- [34] Ledig, C., Theis, L., Huszár, F., Caballero, J., Cunningham, A., Acosta, A., Aitken, A., Tejani, A., Totz, J., Wang, Z., et al.: Photo-realistic single image super-resolution using a generative adversarial network. In: *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4681–4690 (2017)
- [35] Lee, H.Y., Tseng, H.Y., Huang, J.B., Singh, M., Yang, M.H.: Diverse image-to-image translation via disentangled representations. In: *European Conference on Computer Vision*, pp. 35–51 (2018)
- [36] Li, J., Liang, X., Wei, Y., Xu, T., Feng, J., Yan, S.: Perceptual generative adversarial networks for small object detection. In: *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1222–1230 (2017)
- [37] Liu, M.Y., Breuel, T., Kautz, J.: Unsupervised image-to-image translation networks. In: *Advances in Neural Information Processing Systems*, pp. 700–708 (2017)
- [38] Liu, Z., Luo, P., Wang, X., Tang, X.: Deep learning face attributes in the wild. In: *IEEE International Conference on Computer Vision*, pp. 3730–3738 (2015)
- [39] Lucic, M., Kurach, K., Michalski, M., Gelly, S., Bousquet, O.: Are GANs created equal? a large-scale study. In: *Advances in Neural Information Processing Systems*, pp. 700–709 (2018)
- [40] Luo, Y., Xu, Y., Ji, H.: Removing rain from a single image via discriminative sparse coding. In: *ICCV*, pp. 3397–3405 (2015)
- [41] Ma, L., Jia, X., Georgoulis, S., Tuytelaars, T., Van Gool, L.: Exemplar guided unsupervised image-to-image translation. In: *International Conference on Learning Representation* (2019)
- [42] Ma, L., Jia, X., Sun, Q., Schiele, B., Tuytelaars, T., Van Gool, L.: Pose guided person image generation. In: *Advances in Neural Information Processing Systems*, pp. 406–416 (2017)
- [43] Mao, X., Li, Q., Xie, H., Lau, R.Y., Wang, Z., Smolley, S.P.: Least squares generative adversarial networks. In: *IEEE international conference on computer vision (ICCV)*, pp. 2813–2821. *IEEE* (2017)
- [44] Mirza, M., Osindero, S.: Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784* (2014)
- [45] Mishra, A., Rai, S.N., Mishra, A., Jawahar, C.: IIIT-CFW: A benchmark database of cartoon faces in the wild. In: *European Conference on Computer Vision*, pp. 35–47 (2016)
- [46] Nguyen, A., Clune, J., Bengio, Y., Dosovitskiy, A., Yosinski, J.: Plug & play generative networks: Conditional iterative generation of images in latent space. In: *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4467–4477 (2017)
- [47] Nogue, F.C., Huie, A., Dasgupta, S.: Object detection using domain randomization and generative adversarial refinement of synthetic images. *arXiv preprint arXiv:1805.11778* (2018)
- [48] Odena, A., Olah, C., Shlens, J.: Conditional image synthesis with auxiliary classifier gans. In: *International Conference on Machine Learning*, pp. 2642–2651. *JMLR.org* (2017)
- [49] Olszanowski, M., Pochwatko, G., Kuklinski, K., Scibor-Rylski, M., Lewinski, P., Ohne, R.K.: Warsaw Set of Emotional Facial Expression Pictures: a validation study of facial display photographs. *Frontiers in Psychology* **5**, 1516 (2015)
- [50] Park, T., Liu, M.Y., Wang, T.C., Zhu, J.Y.: Semantic image synthesis with spatially-adaptive normalization. In: *IEEE Conference on Computer Vision and Pattern Recognition* (2019)
- [51] Pathak, D., Krahenbuhl, P., Donahue, J., Darrell, T., Efros, A.A.: Context encoders: Feature learning by inpainting. In: *CVPR*, pp. 2536–2544 (2016)
- [52] Press, O., Bar, A., Bogin, B., Berant, J., Wolf, L.: Language generation with recurrent generative adversarial networks without pre-training. *arXiv preprint arXiv:1706.01399* (2017)
- [53] Qi, G.J.: Loss-sensitive generative adversarial networks on lipschitz densities. *arXiv preprint arXiv:1701.06264* (2017)
- [54] Radford, A., Metz, L., Chintala, S.: Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434* (2015)
- [55] Reed, S., Akata, Z., Yan, X., Logeswaran, L., Schiele, B., Lee, H.: Generative adversarial text to image synthesis. *arXiv preprint arXiv:1605.05396* (2016)
- [56] Reed, S., van den Oord, A., Kalchbrenner, N., Bapst, V., Botvinick, M., de Freitas, N.: Generating interpretable images with controllable structure. In: *International Conference on Learning Representation Workshop* (2017)
- [57] Reed, S.E., Akata, Z., Mohan, S., Tenka, S., Schiele, B., Lee, H.: Learning what and where to draw. In: *Advances in Neural Information Processing Systems*, pp. 217–225 (2016)
- [58] Shrivastava, A., Pfister, T., Tuzel, O., Susskind, J., Wang, W., Webb, R.: Learning from simulated and unsupervised images through adversarial training. In: *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2107–2116 (2017)
- [59] Sønderby, C.K., Caballero, J., Theis, L., Shi, W., Huszár, F.: Amortised map inference for image super-resolution. In: *International Conference on Learning Representation*, pp. 1–17 (2017)
- [60] Taigman, Y., Polyak, A., Wolf, L.: Unsupervised cross-domain image generation. In: *International Conference on Learning Representation* (2017)
- [61] Thomaz, C.E., Giraldo, G.A.: A new ranking method for principal components analysis and its application to face image analysis. *Image and Vision Computing* **28**(6), 902–913 (2010)
- [62] Tran, L.Q., Yin, X., Liu, X.: Representation learning by rotating your faces. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2018)
- [63] Van Der Schalk, J., Hawk, S.T., Fischer, A.H., Doosje, B.: Moving faces, looking places: validation of the Amsterdam Dynamic Facial Expression Set (ADFES). *Emotion* **11**(4), 907 (2011)
- [64] Wang, T.C., Liu, M.Y., Zhu, J.Y., Tao, A., Kautz, J., Catanzaro, B.: High-resolution image synthesis and semantic manipulation with conditional gans. *arXiv preprint arXiv:1711.11585* (2017)
- [65] Wang, X., Tang, X.: Face photo-sketch synthesis and recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **31**(11), 1955–1967 (2008)
- [66] Wang, X., Tang, X.: Face photo-sketch synthesis and recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **31**(11), 1955–1967 (2009)



- [67] Willmott, C.J., Matsuura, K.: Advantages of the mean absolute error (mae) over the root mean square error (rmse) in assessing average model performance. *Climate research* **30**(1), 79–82 (2005)
- [68] Yan, X., Yang, J., Sohn, K., Lee, H.: Attribute2image: Conditional image generation from visual attributes. In: *European Conference on Computer Vision* (2016)
- [69] Yang, C., Lu, X., Lin, Z., Shechtman, E., Wang, O., Li, H.: High-resolution image inpainting using multi-scale neural patch synthesis. In: *CVPR*, pp. 6721–6729 (2017)
- [70] Yang, C., Wang, Z., Zhu, X., Huang, C., Shi, J., Lin, D.: Pose guided human video generation. In: *European Conference on Computer Vision*, pp. 201–216 (2018)
- [71] Yi, Z., Zhang, H., Tan, P., Gong, M.: DualGAN: Unsupervised dual learning for image-to-image translation. In: *IEEE International Conference on Computer Vision*, pp. 2849–2857 (2017)
- [72] Zhang, H., Sindagi, V., Patel, V.M.: Image de-raining using a conditional generative adversarial network. *arXiv preprint arXiv:1701.05957* (2017)
- [73] Zhang, H., Xu, T., Li, H., Zhang, S., Wang, X., Huang, S.X., Metaxas, D.N.: StackGAN++: Realistic image synthesis with stacked generative adversarial networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2018)
- [74] Zhang, H., Xu, T., Li, H., Zhang, S., Wang, X., Huang, X., Metaxas, D.N.: StackGAN: Text to photo-realistic image synthesis with stacked generative adversarial networks. In: *IEEE International Conference on Computer Vision*, pp. 5907–5915 (2017)
- [75] Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable effectiveness of deep features as a perceptual metric. *arXiv preprint* (2018)
- [76] Zhang, W., Sun, J., Tang, X.: Cat head detection-how to effectively exploit shape and texture features. In: *European Conference on Computer Vision*, pp. 802–816. Springer (2008)
- [77] Zhang, Z., Song, Y., Qi, H.: Age progression/regression by conditional adversarial autoencoder. In: *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5810–5818 (2017)
- [78] Zhao, B., Wu, X., Cheng, Z.Q., Liu, H., Jie, Z., Feng, J.: Multi-view image generation from a single-view. In: *2018 ACM Multimedia Conference on Multimedia Conference*, pp. 383–391. ACM (2018)
- [79] Zhao, J., Mathieu, M., LeCun, Y.: Energy-based generative adversarial network. *arXiv preprint arXiv:1609.03126* (2016)
- [80] Zhu, J.Y., Park, T., Isola, P., Efros, A.A.: Unpaired image-to-image translation using cycle-consistent adversarial networks. In: *IEEE International Conference on Computer Vision*, pp. 2223–2232 (2017)
- [81] Zhu, J.Y., Zhang, R., Pathak, D., Darrell, T., Efros, A.A., Wang, O., Shechtman, E.: Toward multimodal image-to-image translation. In: *Advances in Neural Information Processing Systems*, pp. 465–476 (2017)