# Asymmetric Generative Adversarial Networks for Image-to-Image Translation

Hao Tang, Dan Xu, Hong Liu and Nicu Sebe

**State-of-the-art models for unpaired image-to-image translation with Generative Adversarial Networks (GANs) can learn the mapping from the source domain to the target domain using a cycle-consistency loss. The intuition behind these models is that if we translate from one domain to the other and back again we should arrive at where we started. However, existing methods always adopt a symmetric network architecture to learn both forward and backward cycles. Because of the task complexity and cycle input difference between the source and target image domains, the inequality in bidirectional forward-backward cycle translations is significant and the amount of information between two domains is different. In this paper, we analyze the limitation of the existing symmetric GAN models in asymmetric translation tasks, and propose an AsymmetricGAN model with both translation and reconstruction generators of unequal sizes and different parameter-sharing strategy to adapt to the asymmetric need in both unsupervised and supervised image-to-image translation tasks. Moreover, the training stage of existing methods has the common problem of model collapse that degrades the quality of the generated images, thus we explore different optimization losses for better training of AsymmetricGAN, and thus make image-to-image translation with higher consistency and better stability. Extensive experiments on both supervised and unsupervised generative tasks with several publicly available datasets demonstrate that the proposed AsymmetricGAN achieves superior model capacity and better generation performance compared with existing GAN models. To the best of our knowledge, we are the first to investigate the asymmetric GAN framework on both unsupervised and supervised image-to-image translation tasks. The source code, data and trained models are available at https://github.com/Ha0Tang/AsymmetricGAN.**

***Index Terms*—Generative Adversarial Networks (GANs), Asymmetric Networks, Image-to-Image Translation, Style Transfer.**

## I. INTRODUCTION

Recently, Generative Adversarial Networks (GANs) [1] have received considerable attention in computer vision community. GANs are generative models which are particularly designed for generation tasks. Recent works have been able to yield promising image translation performance, e.g., Pix2pix [2], in a supervised setting given carefully annotated image pairs. However, pairing the training data is usually difficult and costly. To tackle this problem, several GAN approaches, such as CycleGAN [3], DualGAN [4] and ComboGAN [5], target to effectively learn a mapping from the source domain to the target domain without paired training data. Some progress has been made by these cross-modal translation frameworks on the unpaired image translation task. However, these are not efficient for the multi-domain image translation. For example, for $m$ different domains, BicycleGAN and Pix2pix need the training of $m(m-1)$ models; CycleGAN and DualGAN require $\frac{m(m-1)}{2}$ models; ComboGAN needs to train $m$ models for different $m$ image domains.

To fix the aforementioned limitation, Choi et al. propose StarGAN [6], which performs multi-domain image translation using only one generator/discriminator pair and an extra domain classifier [7]. Mathematically, we assume $X$ and $Y$ represent the source and target domains, and $x \in X$ and $y \in Y$ denote images in domain $X$ and domain $Y$, respectively; we define $z_x$ and $z_y$ denote category labels of domain $X$ and $Y$,
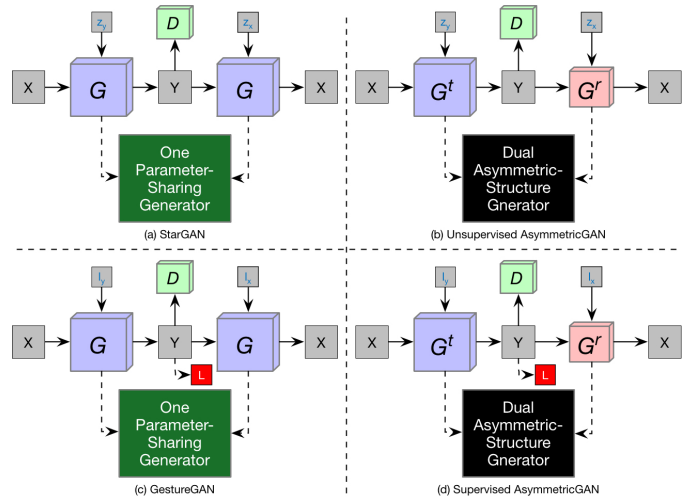


Fig. 1: An illustration of the unsupervised and supervised frameworks of the proposed AsymmetricGAN (b and d) compared with StarGAN [6] (a) and GestureGAN [10] (c).

respectively. StarGAN utilizes a symmetric GAN model and uses the same generator $G$ twice to translate $X$ into $Y$ with the target label $z_y$, i.e., $G(x, z_y) \approx y$, and reconstructs the input image $x$ from the translated output $G(x, z_y)$ and the label $z_x$, i.e., $G(G(x, z_y), z_x) \approx x$. In this way, the generator $G$ shares a common mapping and data structures for two different tasks, i.e., image translation and image reconstruction. Moreover, we note that StarGAN cannot handle with some specific image translation tasks such as person image generation [8], [9] and hand gesture-to-gesture translation [10], since both tasks have infinite image domain $m$ as illustrated in [10].

To solve the limitation, Tang et al. [10] propose GestureGAN, which can produce hand gestures with different

Hao Tang and Nicu Sebe are with the Department of Information Engineering and Computer Science (DISI), University of Trento, Trento 38123, Italy. E-mail: hao.tang@unitn.it, sebe@disi.unitn.it.

Dan Xu is with the Department of Engineering Science, University of Oxford, Oxford OX1 2JD, United Kingdom. E-mail: danxu@robots.ox.ac.uk.

Hong Liu is with the Engineering Laboratory on Intelligent Perception for Internet of Things, Shenzhen Graduate School, Peking University, Shenzhen 518055, China. E-mail: hongliu@pku.edu.cn.

poses, sizes and locations by using hand skeletons $l_x$ and $l_y$. Note that GestureGAN also uses a symmetric structure, i.e., GestureGAN utilizes the same generator $G$ twice for both image translation and image reconstruction, which can be defined as $G(x, l_y) \approx y$ and $G(G(x, l_y), l_x) \approx x$, respectively. In summary, both StarGAN and GestureGAN use a symmetric structure of GANs. We argue that since each task has unique information and distinct targets, it is harder to optimize the generator and to make it gain a good generalization ability on both tasks.

In this paper, we analyze the limitation of both StarGAN and GestureGAN, and propose a novel Asymmetric Generative Adversarial Network (AsymmetricGAN) for both unsupervised and supervised image-to-image translation tasks. Unlike StarGAN and GestureGAN, AsymmetricGAN consists of two different asymmetric generators of unequal sizes to adapt to the asymmetric need in both image translation and image construction.

There are three reasons for designing the asymmetric-structured generators. Firstly, the translation generator $G^t$ transforms images from $X$ to $Y$, and the reconstruction generator $G^r$ uses the translated images from $G^t$ and the original domain guidance $z_x/l_x$ to reconstruct the original $x$. Generators $G_t$ and $G_r$ cope with different tasks. Usually, we mainly focus on the image translation rather that the image reconstruction. Secondly, the input data distribution for them is different. The inputs of the translation generator $G_t$ are a real image and a target domain guidance. The goal of $G_t$ is to generate the target domain image. While $G_r$ accepts a translated image and an original domain guidance as input, and tries to reconstructs the original input image. For generator $G_t$ and $G_r$, the input images are a real image and a generated/fake image respectively, and thus the data distribution is difference between them. Thirdly, because of the complexity difference between the source and target image domains, the complexity inequality in a bidirectional image-to-image translation is significant. Therefore, it is intuitive to design different network structures for the two different generators. These two generators are allowed to use different network architecture designs and different levels of parameter sharing strategy according to the diverse difficulty of the tasks. By doing so, each generator can have its own network parts which usually helps to learn better each task-specific mapping in a multi-task setting [11]. A motivation illustration of the proposed AsymmetricGAN compared with the most related two works, i.e., StarGAN and GestureGAN, is presented in Fig. 1.

Moreover, to avoid the model collapse issue in training AsymmetricGAN for both unsupervised and supervised image-to-image translation tasks, we further explore different objective functions for better optimization. (i) The color cycle-consistency loss which targets solving the "channel pollution" problem proposed in [10] by separately generating red, green and blue color channels instead of generating all three at one time; (ii) The multi-scale SSIM loss, which preserves the information of luminance, contrast and structure between reconstructed images and input images across different scales; and (iii) The conditional identity preserving loss, which helps

retaining the identity information of the input images. These loss functions are jointly embedded in the proposed AsymmetricGAN for training and help to generate results with higher consistency and better stability. The main contributions of this paper are:

- We propose a novel Asymmetric Generative Adversarial Network (AsymmetricGAN) for both unsupervised and supervised image-to-image translation tasks. The asymmetric dual generators, allowing different network structures and different-level parameter sharing, are designed to specifically cope with image translation and image reconstruction tasks, which facilitates obtaining a better generalization ability of the proposed model to improve the generation performance.
- We explore jointly utilizing different objectives for a better optimization of the proposed AsymmetricGAN, and thus obtaining both unsupervised and supervised image-to-image translation with higher consistency and better stability.
- We extensively evaluate AsymmetricGAN on both unpaired multi-domain image-to-image translation and paired hand gesture-to-gesture translation tasks with several different datasets, demonstrating its superiority in model capacity and its better generation performance compared with state-of-the-art methods.

This paper is organized as follows. Sec. II surveys the evolution of image-to-image translation related methods. Sec. III presents both unsupervised and supervised frameworks of the proposed AsymmetricGAN. In Sec. IV, experimental evaluations and detailed discussions on two popular generative tasks, i.e., multi-domain image-to-image translation and hand gesture-to-gesture translation, are presented. Finally, Sec. V concludes this paper.

## II. RELATED WORK

**Generative Adversarial Networks (GANs)** [1] are generative models, which have achieved promising results on different generative tasks, e.g., image generation [12], [13], [14]. Moreover, to generate images controlled by users, Mirza et al. [15] propose Conditional GANs (CGANs), which use a conditional information to guide the image generation process. Extra conditional guidance information can be category labels [6], [16], object keypoints [8], [17], [18], human skeletons [10], [9], text descriptions [19], [20], semantic maps [21], [22] and conditional images [2]. CGANs have been successfully used in various applications, such as text-to-image [19], audio-to-image [23], image-to-image [2], [24] and video-to-video [25] translation tasks. In this paper, we mainly focus on the image-to-image translation task, methods of this task can be divided into two categories, i.e., supervised/paired and unsupervised/unpaired.

**Supervised/Paired Image-to-Image Translation.** CGANs learn a mapping between image inputs and image outputs using convolutional neural networks. For example, Isola et al. propose Pix2pix [2], which is a conditional framework using a CGAN to learn the mapping function. Based on Pix2pix, Wang et al. present Pix2pixHD [21], which can be used for high-resolution photo-realistic image translation.
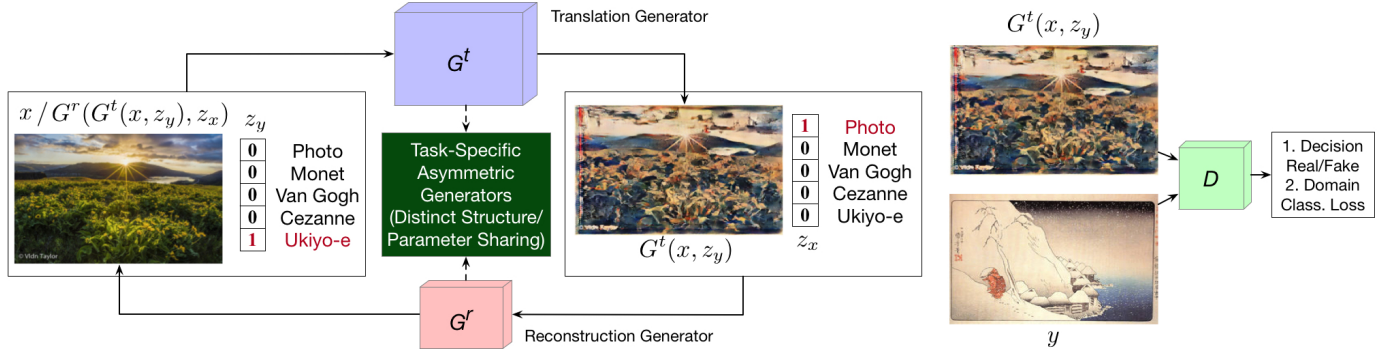
Fig. 2: The unsupervised framework of AsymmetricGAN for the multi-domain image-to-image translation task. $z_x$ and $z_y$ indicate the category labels of domain $X$ and $Y$, respectively. $G^t$ and $G^r$ are task-specific asymmetric generators. The translation generator $G^t$ translates images from domain $X$ into domain $Y$ and the reconstruction generator $G^r$ receives the generated image $G^t(x, z_y)$ and the original domain label $z_x$ and tries to recover the input image $x$ during the training stage with the proposed objective functions.

Similar ideas have also been applied to many other tasks, e.g., pose-guided person image generation [8], [26], [9] and hand gesture-to-gesture translation [10]. However, all of these models require paired training data, which are usually costly to obtain.

**Unsupervised/Unpaired Image-to-Image Translation.** To alleviate the issue of pairing training data, Zhu et al. introduce CycleGAN [3], which learns the mappings between two unpaired image domains without supervision with the aid of a cycle-consistency loss. Apart from CycleGAN, there are other variants proposed to tackle the same problem such as [4], [27], [28], [29], [30], [31]. However, these models are only suitable in cross-domain translation tasks.

**Multi-Domain Image-to-Image Translation.** There are only very few recent methods attempting to implement multi-modal image-to-image translation in an efficient way. Anoosheh et al. propose ComboGAN [5], which only requires to train $m$ generator/discriminator pairs for $m$ different image domains. Choi et al. present the StarGAN framework [6], which equips a single symmetric-structured generator and is able to perform the task with a complexity of $\Theta(1)$. Although the model complexity is low, jointly learning both image translation and image reconstruction tasks with the same generator require the sharing of all parameters, which increases the optimization complexity and reduces the generalization ability, thus leading to unsatisfactory generation performance. The proposed framework targets at obtaining a good balance between the network capacity and image generation quality. Along with this research line, we propose a novel Asymmetric Generative Adversarial Network (AsymmetricGAN), which achieves this target via using two task-specific and asymmetric generators. Moreover, we explore various optimization objectives to train better the model to produce more consistent and more stable results.

## III. ASYMMETRIC GENERATIVE ADVERSARIAL NETWORKS

We first start with the unsupervised framework of the proposed AsymmetricGAN for the multi-domain image-to-image translation task, and then introduce the supervised framework

of the proposed AsymmetricGAN for the hand gesture-to-gesture translation task.

### A. Unsupervised Framework

We focus on the multi-domain image-to-image translation task with unpaired training data. The overview of the proposed AsymmetricGAN is illustrated in Fig. 2. Existing cross-domain generation models, such as CycleGAN [3], Disco-GAN [28] and DualGAN [4], which need to separately train $\frac{m(m-1)}{2}$ models for $m$ different image domains. However, the proposed AsymmetricGAN is specifically designed for tackling the multi-domain image translation problem with significant advantages in the model complexity and in the training overhead, which only needs to train a single model. To directly compare with StarGAN [6], which simply adopts the same generator for both image reconstruction and image translation tasks. We argue that the training of a single generator model for multiple domains is a challenging problem as mentioned in the introduction section, thus we propose a more effective asymmetric generator network structure and more robust optimization objectives to stabilize the training process. In summary, our work focuses on exploring different strategies to improve the optimization of the multi-domain translation model, which we aim to give useful insights into the design of more effective multi-domain translation generators.

To achieve this goal, the translation generator $G^t$ is learned to translate an input image $x$ into an output image $y$ which is conditioned on the target domain label $z_y$, this process can be expressed as $G^t(x, z_y) \rightarrow y$. Then the reconstruction generator $G^r$ receives the translated image $G^t(x, z_y)$ and the original domain label $z_x$ as input, and learns to recover the input image $x$, this process can be formulated as $G^r(G^t(x, z_y), z_x) \rightarrow x$. The asymmetric generators are task-specific generators which allow for different network designs and different levels of parameter sharing for learning better the generators. The discriminator $D$ tries to distinguish between the real image $y$ and the generated image $G^t(x, z_y)$, and also to classify the translated image $G^t(x, z_y)$ to the corresponding domain label $z_y$ via the domain classification loss. Moreover, we investigate how the distinct network designs and different network

sharing schemes for the asymmetric generators dealing with different sub-tasks could balance the generation performance and network complexity. These aspects are not covered and considered in the multi-domain model StarGAN [6].

**Model Optimization.** The optimization objective of the proposed AsymmetricGAN contains several different losses. These optimization losses are jointly embedded into the proposed AsymmetricGAN during the training stage.

*Color Cycle-Consistency Loss.* The cycle-consistency loss can be regarded as "pseudo" pairs in training data even though we do not have corresponding samples in the target domain. This loss function can be defined as:

$$\mathcal{L}_{cyc}(x) = \mathbb{E}_{x \sim p_{\text{data}}(x)} \left[ ||G^r(G^t(x, z_y), z_x) - x||_1 \right]. \quad (1)$$

The $L1$ norm is adopted for the reconstruction loss. The goal of this optimization loss is to make the reconstructed image $G^r(G^t(x, z_y), z_x)$ as close as possible to the input image $x$. However, the generation of a whole image at one time makes the different color channels influence each other, thus leading to artifacts in the generation results [10]. To overcome this limitation, we propose a novel color cycle-consistency loss, which constructs the consistence loss for each channel, separately. This loss can be expressed as,

$$\mathcal{L}_{colorcyc} = \sum_{i \in \{r,g,b\}} \mathcal{L}_{cyc}(x^i), \quad (2)$$

where $x^b, x^g, x^r$ denote the blue, yellow and red channels of the image $x$. We calculate the pixel loss for the red, green, blue channels separately and then sum up these three color losses as the final loss. In this way, the generator can be enforced to generate each channel independently to avoid the "channel pollution" problem.

*Multi-Scale SSIM Loss.* The structural similarity index (SSIM) has been originally proposed in [32] to measure the similarity of two images. We introduce it here to help to preserve the information of luminance, contrast and structure across scales. For the reconstructed image $\widehat{x}=G^r(G^t(x, z_y), z_x)$ and the input image $x$, the SSIM loss is written as:

$$\mathcal{L}_{ssim}(\widehat{x}, x) = [l(\widehat{x}, x)]^\alpha [c(\widehat{x}, x)]^\beta [s(\widehat{x}, x)]^\gamma, \quad (3)$$

where $l(\widehat{x}, x) = \frac{2\mu_{\widehat{x}}\mu_x + C_1}{\mu_{\widehat{x}}^2 + \mu_x^2 + C_1}$, $c(\widehat{x}, x) = \frac{2\sigma_{\widehat{x}}\sigma_x + C_2}{\sigma_{\widehat{x}}^2 + \sigma_x^2 + C_2}$ and $s(\widehat{x}, x) = \frac{\sigma_{\widehat{x}x} + C_3}{\sigma_{\widehat{x}}\sigma_x + C_3}$. These three terms compare the luminance, contrast and structure information between $\widehat{x}$ and $x$. $\alpha, \beta$ and $\gamma$ are hyper-parameters to control the relative weight of $l(\widehat{x}, x)$, $c(\widehat{x}, x)$ and $s(\widehat{x}, x)$, respectively; $\mu_{\widehat{x}}$ and $\mu_x$ are the means of $\widehat{x}$ and $x$; $\sigma_{\widehat{x}}$ and $\sigma_x$ are the standard deviations of $\widehat{x}$ and $x$; $\sigma_{\widehat{x}x}$ is the covariance of $\widehat{x}$ and $x$; $C_1$, $C_2$ and $C_3$ are predefined parameters. To make the model benefit from multi-scale deep information, we refer to a multi-scale implementation of SSIM [33] which constrains SSIM over scales. The Multi-Scale SSIM loss can be written as:

$$\mathcal{L}_{msssim}(\widehat{x}, x) = [l_M(\widehat{x}, x)]^{\alpha_M} \prod_{j=1}^{M} [c_j(\widehat{x}, x)]^{\beta_j} [s_j(\widehat{x}, x)]^{\gamma_j}. \quad (4)$$

Through using this loss, the luminance, contrast and structure information of the input images is expected to be preserved.

*Conditional Least Square Loss.* We use a least square loss [34], [3] to stabilize our model during the training stage. The least square loss is more stable than the negative log likelihood objective $\mathcal{L}_{cgan}(G^t, D_s, z_y) = \mathbb{E}_{y \sim p_{\text{data}}(y)} [\log D_s(y)] + \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log(1 - D_s(G^t(x, z_y)))]$, and is converging faster than Wasserstein GAN [35]. This loss can be expressed as:

$$\mathcal{L}_{lsgan} = \mathbb{E}_{y \sim p_{\text{data}}(y)} \left[ (D_s(y) - 1)^2 \right] + \\ \mathbb{E}_{x \sim p_{\text{data}}(x)} \left[ D_s(G^t(x, z_y))^2 \right], \quad (5)$$

where $z_y$ are the category labels of domain $y$, $D_s$ is the probability distribution over sources produced by discriminator $D$. The target of $G^t$ is to generate an image $G^t(x, z_y)$ that is expected to be similar to the images from domain $Y$, while $D$ aims to distinguish the generated images $G^t(x, z_y)$ from the real ones $y$.

*Domain Classification Loss.* To perform multi-domain image translation with a single discriminator, previous works employ an auxiliary classifier [7], [6] on the top of the discriminator, and impose the domain classification loss when updating both the generator and discriminator. We also consider this loss in our optimization:

$$\mathcal{L}_c = \mathbb{E}_{x \sim p_{\text{data}}(x)} \left\{ - \left[ \log D_c(z_x|x) + \log D_c(z_y|G^t(x, z_y)) \right] \right\}, \quad (6)$$

where $D_c(z_x|x)$ represents the probability distribution over the domain labels given by discriminator $D$. $D$ learns to classify $x$ to its corresponding domain $z_x$. $D_c(z_y|G^t(x, z_y)$ denotes the domain classification for fake images. We minimize the domain classification loss to produce the image $G^t(x, z_y)$ that can be classified to the corresponding domain $z_y$.

*Conditional Identity Preserving Loss.* To reinforce the identity of the input image during the translation, we use a conditional identity preserving loss [3]. This loss encourages the mapping to preserve identity information such as color information between the input and the output, which can be formulated as,

$$\mathcal{L}_{id}(G^t, G^r, z_x) = \mathbb{E}_{x \sim p_{\text{data}}(x)} \left[ ||G^r(x, z_x) - x||_1 \right]. \quad (7)$$

By doing so, the generator also tries to preserve the identity via the back-propagation of this loss.

*Full Objective.* Given the loss functions presented above, the complete optimization objective of AsymmetricGAN for the multi-domain image-to-image translation task can be written as:

$$\mathcal{L} = \mathcal{L}_{lsgan} + \lambda_c \mathcal{L}_c + \lambda_{cyc} \mathcal{L}_{colorcyc} + \lambda_m \mathcal{L}_{msssim} + \lambda_{id} \mathcal{L}_{id}, \quad (8)$$

where $\lambda_c$, $\lambda_{cyc}$, $\lambda_m$ and $\lambda_{id}$ are parameters controlling the relative importance of the corresponding objective. All objectives are jointly optimized in an end-to-end fashion. We set $\lambda_c=1, \lambda_{cyc}=10, \lambda_m=1, \lambda_{id}=0.5$ in our experiments.

**Network Architecture.** The proposed AsymmetricGAN consists of an asymmetric dual-generator and a discriminator. The asymmetric dual generators are designed to specifically deal with different tasks in GANs, i.e., the translation and the reconstruction tasks, which have different targets for training the network. We can design different network structures for the different generators to make them learn better task-specific
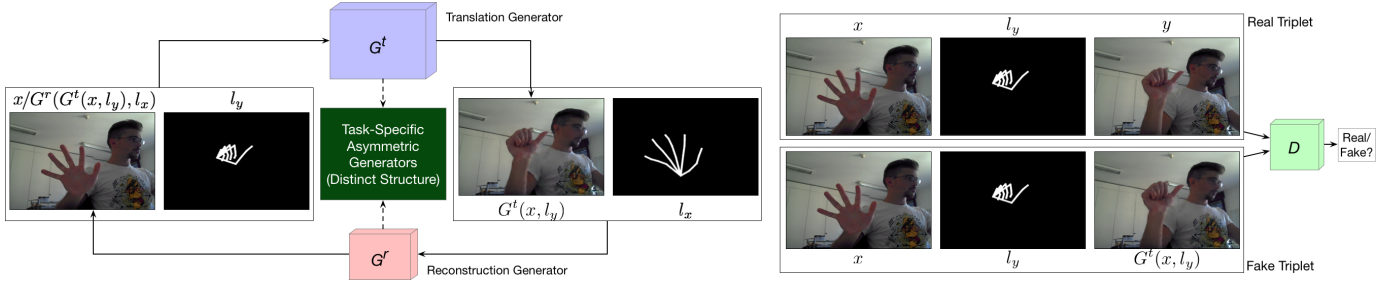
Fig. 3: The supervised framework of AsymmetricGAN for the hand gesture-to-gesture translation task. $l_x$ and $l_y$ indicate the hand skeletons of images $x$ and $y$, respectively. $G^t$ and $G^r$ are task-specific asymmetric generators. The translation generator $G^t$ converts images from domain $X$ into domain $Y$ and the construction generator $G^r$ receives the generated image $G^t(x, l_y)$ and the original hand skeleton $l_x$ and attempts to reconstruct the original image $x$ during the optimization with the proposed different objective losses. We have two cycles, i.e., $x \mapsto y' \mapsto \widehat{x} \approx x$ and $y \mapsto x' \mapsto \widehat{y} \approx y$, but we only show one here, i.e., $x \mapsto y' \mapsto \widehat{x} \approx x$.

objectives, which allows us to share parameters between the generators to further reduce the model capacity since the shallow image representations are shareable for both generators. The parameter sharing facilitates the achievement of good balance between the model complexity and the generation quality. Our model generalizes the model of StarGAN [6]. When the parameters are fully shared with the usage of the same network structure for both generators, our framework becomes a StarGAN. Moreover, We represent the class labels $z_x$ and $z_y$ using a one-hot vector, and then the vector is passed through a linear layer to obtain a label embedding with 64 dimensions. This embedding is replicated to form feature maps that are further concatenated with the image feature maps for follow-up convolution operations with residual blocks and several deconvolution layers to obtain the target images. For the discriminator $D$, we use PatchGAN [3], [6]. After the discriminator, a convolution layer is applied to produce a final one-dimensional output which indicates whether local image patches are real or fake.

**Network Training.** For reducing model oscillation, we adopt a cache of generated images to update the discriminator as in [36]. In the experiments, we set the number of image buffer to 50. The batch size is set to 1 for all the experiments and all the models are trained with 200 epochs.

### B. Supervised Framework

In this part, we start to introduce the supervised framework of the proposed AsymmetricGAN for the hand gesture-to-gesture translation task. The framework is shown in Fig. 3, which consists of asymmetric dual generators (i.e., $G^t$ and $G^r$) and a discriminator $D$. Specifically, we concatenate the input image $x$ and the target hand skeleton $l_y$, and input them into the translation generator $G^t$ and synthesize the target image $y'=G^t(x, l_y)$. Different from GestureGAN [10], which adopts the same generator to reconstruct the original input image, we propose an asymmetric reconstruction generator $G^r$ to benefit more from the image translation process. The conditional hand skeleton $l_x$ together with the generated image $y'$ are input into the reconstruction generator $G^r$, and produce the reconstructed input image $\widehat{x}$. We formalize the process as $\widehat{x}=G^r(y', l_x)=G^r(G^t(x, l_y), l_x)$. Then the optimization objective is to make $\widehat{x}$ as close as possible to $x$.

**Model Optimization.** For better optimizing the proposed AsymmetricGAN on the hand gesture-to-gesture translation task, we adopt six loss functions, i.e., cycle-consistency loss, color loss, adversarial loss, identity preserving loss, perceptual loss and total variation loss. These optimization losses and the proposed framework are jointly learned in an end-to-end fashion during the training stage.

*Cycle-Consistency Loss.* This loss ensures the consistency between the source image $x$ and the reconstructed image $\widehat{x}$, and it can be defined as,

$$\mathcal{L}_{cyc}(x) = \mathbb{E}_{x \sim p_{\text{data}}(x)} \left[ ||G^r(G^t(x, l_y), l_x) - x||_1 \right]. \quad (9)$$

Note that different from PG$^2$ [8], DPIG [26] and PoseGAN [9], which only use the target keypoints or skeletons to guide the image generation process, the proposed AsymmetricGAN not only uses the target keypoints or skeletons to guide the image translation process, but also uses them to guide the image reconstruction process. In this way, the cycle consistency can further be guaranteed.

*Color Loss.* We also adopt an improved pixel loss, i.e., channel-wise color loss, to reduce the "channel pollution" issue [10]. The loss can be expressed as,

$$\mathcal{L}_{color} = \sum_{i \in \{r,g,b\}} \mathbb{E}_{x \sim p_{\text{data}}(x)} \left[ ||G^t(x^i, l_y) - y^i||_1 \right]. \quad (10)$$

where $y^r$, $y^g$ and $y^b$ denote the red, green and blue color channels of image $y$. The intuition is that the generation of a 3-channel image is much more complex than the generation of a 1-channel image. Moreover, by calculating the loss of each channel independently, the error from each channel will not influence other channels.

*Conditional Adversarial Loss.* The goal of vanilla GAN loss is to train the generator $G$ which learns the mapping from random noise $z$ to the image $y$. The mapping $G(z) \to y$ can be learned through the following function,

$$\mathcal{L}_{gan}(G^t, D) = \mathbb{E}_{y \sim p_{\text{data}}(y)} \left[ \log D(y) \right] + \\ \mathbb{E}_{z \sim p_{\text{data}}(z)} \left[ \log(1 - D(G^t(z))) \right]. \quad (11)$$

Base on this, CGANs try to learn the mapping from a conditional image $x$ to the target image $y$. The generator $G^t$ tries to generate image $y'=G^t(x)$ which cannot be distinguished from the real image $y$, while the discriminator $D$ tries to detect the

TABLE I: Description of the datasets used in the multi-domain image-to-image translation task.

| Dataset | Type | #Domain | #Mapping | Resolution | Unpaired/Paired | #Train | #Test | #Total |
|---|---|---|---|---|---|---|---|---|
| Facades [37] | Architectures | 2 | 2 | 256×256 | Paired | 800 | 212 | 1,012 |
| AR Face [38] | Faces | 4 | 12 | 768×576 | Paired | 920 | 100 | 1,020 |
| Bu3dfe [39] | Faces | 7 | 42 | 512×512 | Paired | 2,520 | 280 | 2,800 |
| Alps [5] | Natural Seasons | 4 | 12 | - | Unpaired | 6,053 | 400 | 6,453 |
| RaFD [40] | Faces | 8 | 56 | 1024×681 | Unpaired | 5,360 | 2,680 | 8,040 |
| Collection Style [3] | Painting Style | 5 | 20 | 256×256 | Unpaired | 7,837 | 1,593 | 9,430 |

fake images produced by $G^t$. Thus, the objective function of a CGAN can be expressed as,

$$\mathcal{L}_{cgan}(G^t, D) = \mathbb{E}_{y \sim p_{\text{data}}(y)} \left[ \log D(x, y) \right] + \\ \mathbb{E}_{x \sim p_{\text{data}}(x)} \left[ \log(1 - D(x, G^t(x, l_y))) \right], \quad (12)$$

where $D$ tries to distinguish the fake image pair $(x, G^t(x, l_y))$ from the real image pair $(x, y)$. To jointly learn the images and the hand skeletons, we make a modification base on Eq. (12),

$$\mathcal{L}_{cgan}(G^t, D) = \mathbb{E}_{y \sim p_{\text{data}}(y)} \left[ \log D(x, l_y, y) \right] + \\ \mathbb{E}_{x \sim p_{\text{data}}(x)} \left[ \log(1 - D(x, l_y, G^t(x, l_y))) \right], \quad (13)$$

where $D$ tries to distinguish the fake image triplet $(x, l_y, G^t(x, l_y))$ from the real image triplet $(x, l_y, y)$. In this way, $D$ takes consideration of both images and hand skeletons during optimization.

*Conditional Identity Preserving Loss.* To further preserve the identity information, we propose an conditional identity preserving loss, which can be defined as,

$$\mathcal{L}_{id} = \mathbb{E}_{x \sim p_{\text{data}}(x)} \left[ ||G^t(x, l_x) - x||_1 \right] + \\ \mathbb{E}_{y \sim p_{\text{data}}(y)} \left[ ||G^t(y, l_y) - y||_1 \right]. \quad (14)$$

We adopt the $L1$ loss to minimize the difference between the images ($x$ and $y$) and the self-guided results, i.e., $G^t(x, l_x)$ and $G^t(y, l_y)$.

*Perceptual Loss and Total Variation Loss.* We also use a perceptual loss and a total variation loss between the generated image $y'$ and the real image $y$ to better optimize our model. Both losses have been shown to be useful in Pix2pixHD [21] and SelectionGAN [22], respectively.

*Full Objective.* The final objective function of the proposed AsymmetricGAN on the hand gesture-to-gesture translation task can be expressed as,

$$\mathcal{L} = \mathcal{L}_{cgan} + \lambda_c \mathcal{L}_{color} + \lambda_{cyc} \mathcal{L}_{cyc} + \\ \lambda_{id} \mathcal{L}_{id} + \lambda_{vgg} \mathcal{L}_{vgg} + \lambda_{tv} \mathcal{L}_{tv}, \quad (15)$$

where hyper-parameters $\lambda_c$, $\lambda_{cyc}$, $\lambda_{id}$, $\lambda_{vgg}$ and $\lambda_{tv}$ are controlling the relative importance of each loss. All objectives are jointly optimized in an end-to-end fashion. In our experiments, we empirically set $\lambda_c$=800, $\lambda_{cyc}$=0.1, $\lambda_{id}$=0.01, $\lambda_{vgg}$=1000 and $\lambda_{tv}$=1e−6.

**Network Architecture.** We adopt the architecture from [41] as our generators. Since our focus is on the translation direction, it means that $G^t$ is more important than $G^r$. Thus we use a deeper network for $G^t$ and a shallow network for $G^r$. More specific, we use nine residual blocks for both generators. However, the filters in first convolutional layer of $G^t$ and $G^r$ are 64 and 4, respectively. For the discriminator $D$, we adopt 70×70 PatchGAN proposed in [2].

**Network Training.** The batch size is set to 4 for both datasets and all the models are trained with 20 epochs. Moreover, we follow [8], [10] and employ OpenPose [42] to extract hand skeletons as training data.

## IV. EXPERIMENTS

We conduct experiments on two different unsupervised and supervised generative tasks, i.e., multi-domain image-to-image translation and hand gesture-to-gesture translation, to evaluate the effectiveness of the proposed AsymmetricGAN. We employ the Adam optimizer [43] with $\beta_1$=0.5 and $\beta_2$=0.999 to optimize the whole model. We sequentially update the translation generator $G^t$ and the reconstruction generator $G^r$ after the discriminator $D$ updates at each iteration.

### A. Multi-Domain Image-to-Image Translation Task

**Datasets.** Six publicly available datasets including Facades [37], AR Face [38], Alps Season [5], Bu3dfe [39], RaFD [40] and Collection Style [3], are used to validate the proposed AsymmetricGAN on the unpaired multi-domain image-to-image translation task. Table I shows the details of these datasets.

**Baselines.** We employ several image translation models as our competing baselines, i.e., CycleGAN [3], DualGAN [4], ComboGAN [5], DistanceGAN [30], Dist.+Cycle [30], Self Dist. [30], BicycleGAN [44] and Pix2pix [2]. We train the models multiple times for every pair of two different image domains except for ComboGAN [5], which only needs to train $m$ models for $m$ different domains. We also adopt StarGAN [6] as a baseline which performs multi-domain image translation using one generator/discriminator pair. The fully supervised Pix2pix and BicycleGAN are trained with paired data, the other baselines and the proposed AsymmetricGAN are trained with unpaired data. Note that since BicycleGAN can produce several different outputs with one single input image, then we randomly pick one output from them for comparison. For a fair comparison, results of all baselines are produced by using the authors' publicly available codes with the same training strategy as our approach.

**Comparison Against Baselines.** The proposed AsymmetricGAN is evaluated on four different tasks, i.e., label↔photo translation, facial expression-to-expression translation, season translation and painting style transfer. We will describe the comparison with the state-of-the-art approaches in the following.

*Task 1: Label↔Photo Translation.* We use the Facades dataset to evaluate the label↔photo translation task. The results on the Facades dataset are shown in Fig. 4, which are only to show that the proposed AsymmetricGAN is also applicable
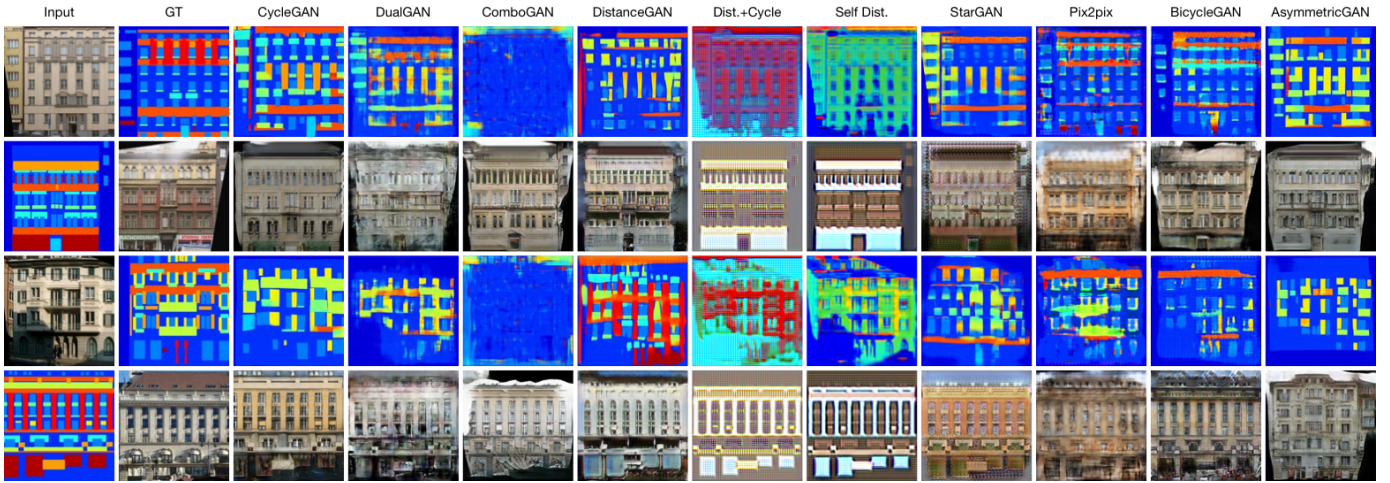
Fig. 4: Different methods for label↔photo translation trained on the Facades dataset.
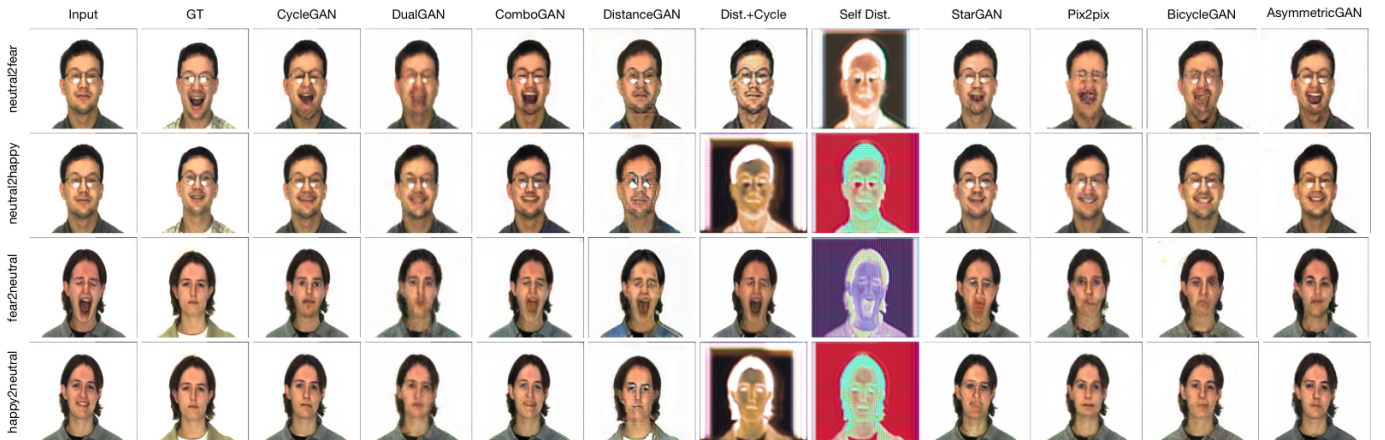


Fig. 5: Different methods for multi-domain facial expression-to-expression translation trained on the AR Face dataset.

on the translation on two domains only and could produce competitive performance. We can see that Dist.+Cycle, Self Dist., ComboGAN fail to generate reasonable results on the photo→label task. For the opposite mapping, i.e., label→photo, Dist.+Cycle, Self Dist., DualGAN, StarGAN and Pix2pix suffer from model collapse, leading reasonable but blurry results. However, the proposed AsymmetricGAN achieves compelling results in both directions compared with baselines.

*Task 2: Facial Expression Synthesis.* We employ three face datasets, i.e., AR Face, Bu3dfe and RaFD, to evaluate the facial expression synthesis task. The results are shown in Fig. 5, we can see that Dist.+Cycle and Self Dist. fail to produce faces similar to the target domain. DualGAN generates reasonable but blurry faces. DistanceGAN, StarGAN, BicycleGAN and Pix2pix produces much sharper results, but still contain some artifacts in the translated faces, e.g., twisted mouths of StarGAN, Pix2pix and BicycleGAN on the "neutral2fear" direction. ComboGAN, CycleGAN and the proposed AsymmetricGAN work better than other baselines on this dataset. Similar results can be seen on the Bu3dfe dataset in Fig. 6. We also present results on the RaFD dataset compared with the most related two works, i.e., CycleGAN and StarGAN, in Fig. 7. We observe that our method achieves visually better results than both CycleGAN and StarGAN.

*Task 3: Season Translation.* We evaluate the proposed AsymmetricGAN on the season translation task. Fig. 8 shows the results. Clearly, DistanceGAN, Dist.+Cycle, Self Dist., DualGAN fail to produce reasonable results. StarGAN can generate reasonable but blurry results, and there are some visual artifacts in the translated results. ComboGAN, CycleGAN and the proposed AsymmetricGAN are able to produce better results than other methods. However, ComboGAN yields some visual artifacts in some cases, such as the "summer2autumn" direction. We also show one failure case of the proposed method on this dataset as shown in the last row of Fig. 8. Our method generates images similar to the input domain, while CycleGAN and DualGAN generate visually better results compared with the proposed AsymmetricGAN on the "winter2spring" direction. Note that both DualGAN and CycleGAN require to train twelve generators for this task on the dataset, while the proposed AsymmetricGAN only needs to train two generators, and thus our model complexity is significantly lower.

*Task 4: Painting Style Transfer.* The comparison results on the painting style dataset compared with two state-of-the-art methods, i.e., CycleGAN and StarGAN, are shown in Fig. 9. We can see that StarGAN generates less diverse generations crossing different styles compared with CycleGAN and AsymmetricGAN. The proposed AsymmetricGAN has
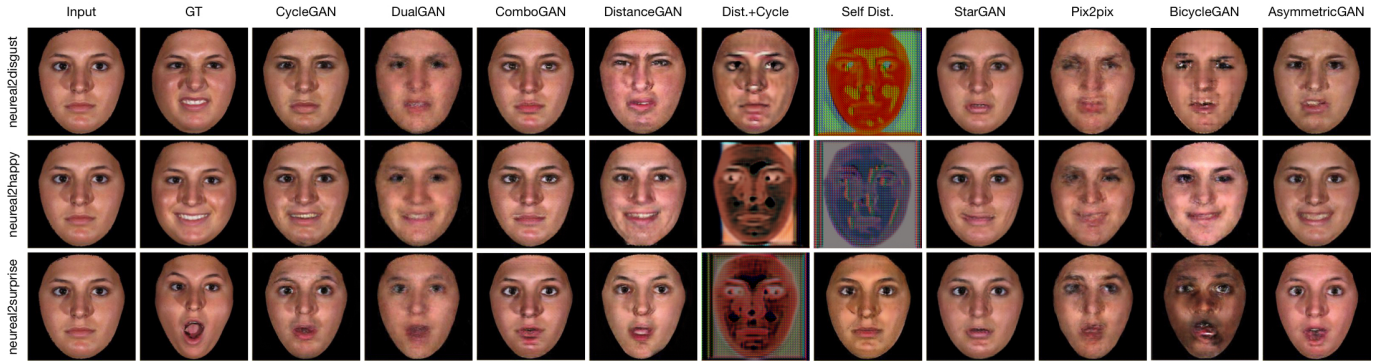
Fig. 6: Different methods for multi-domain facial expression-to-expression translation trained on the Bu3dfe dataset.
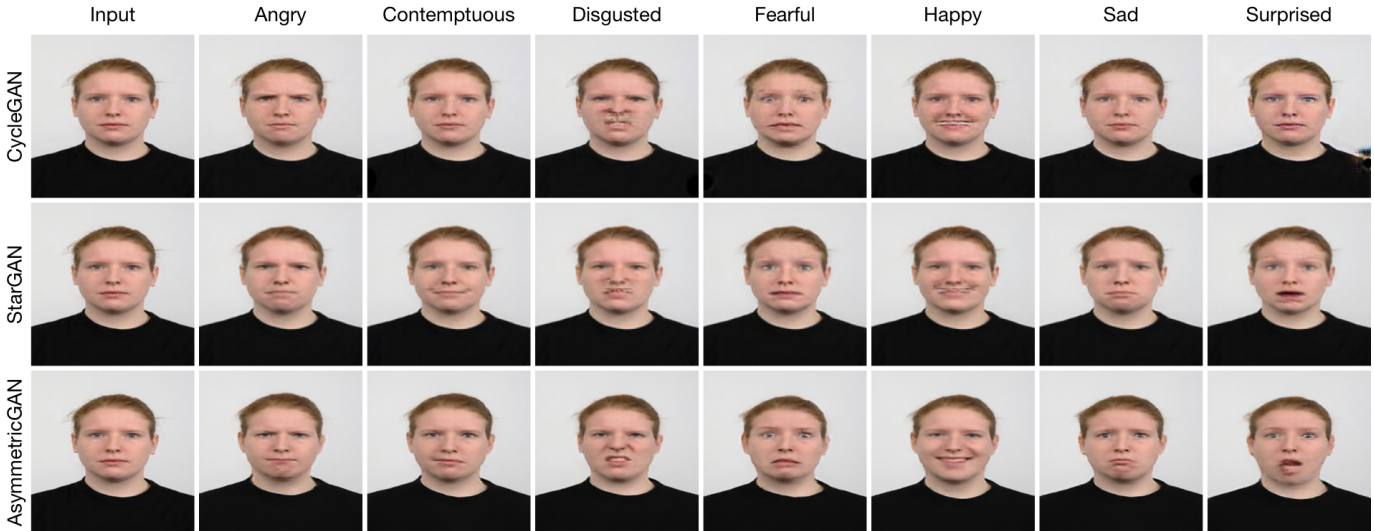


Fig. 7: Different methods for multi-domain facial expression-to-expression translation trained on the RaFD dataset.

comparable performance with CycleGAN, requiring only one single model for all the styles, and thus the network complexity is remarkably lower compared with CycleGAN which trains an individual model for each pair of styles.

*Quantitative Comparison on All Tasks.* Moreover, we provide quantitative performance on the four tasks. Different evaluation metrics are considered, i.e., AMT perceptual studies [3], [2], IS [45], FID [46] and Classification Accuracy (CA) [6].

We follow both CycleGAN and StarGAN, and use the same perceptual study protocol to evaluate the generated images, which is a "real vs fake" perceptual metric to assess the realism from a holistic level. Tables V, IV and II report the performance of the AMT perceptual test. As we can see from Tables V, IV and II, the proposed AsymmetricGAN achieves very competitive results compared with baselines. We observe that the proposed AsymmetricGAN significantly outperforms StarGAN trained using one generator on most of the metrics and on all the datasets. We also note that supervised Pix2pix shows worse results than unpaired methods in Table II, which can be also observed in DualGAN [4].

We then adopt IS [45] to measure the quality of synthesized images. The results are shown in Tables V and III. AsymmetricGAN generates sharper and more photo-realistic results than Dist.+Cycle, Self Dist. and StarGAN, while the latter models present slightly higher IS. However, higher IS does not necessarily mean higher quality. Higher quality

images may have smaller IS as demonstrated in other image generation [8] and super-resolution tasks [41]. Moreover, we employ FID [46] to evaluate the performance on both RaFD and painting style datasets. Results are shown in Tables V and IV, we can see that AsymmetricGAN achieves the best results compared with both StarGAN and CycleGAN.

We finally compute Classification Accuracy (CA) on the generated images as in [6]. We train different classifiers on the AR Face, Alps, Bu3dfe, Collection datasets, respectively. For each dataset, we take the real image as training data and the generated images of different models as testing data. For the AR Face, Alps and Collection datasets, we report top 1 classification accuracy. For the Bu3dfe dataset, we present both top 1 and top 5 classification accuracies. Tables IV and III show the results. Note that the proposed AsymmetricGAN outperforms the baselines on the AR Face, Bu3dfe and Collection datasets. On the Alps dataset, StarGAN achieves slightly better performance than ours but the translated images by our model contain fewer visual artifacts than StarGAN as shown in Fig. 8.

**Model Analysis.** We investigate four aspects of Asymmetric-GAN for the multi-domain image-to-image translation task.

*(1) Importance of Distinct Network Designs for Different Generators.* We design three different experimental settings (i.e., S1, S2, S3) with three different generator architectures ranging from light-weight to heavy-weight: (i) Architecture I has the
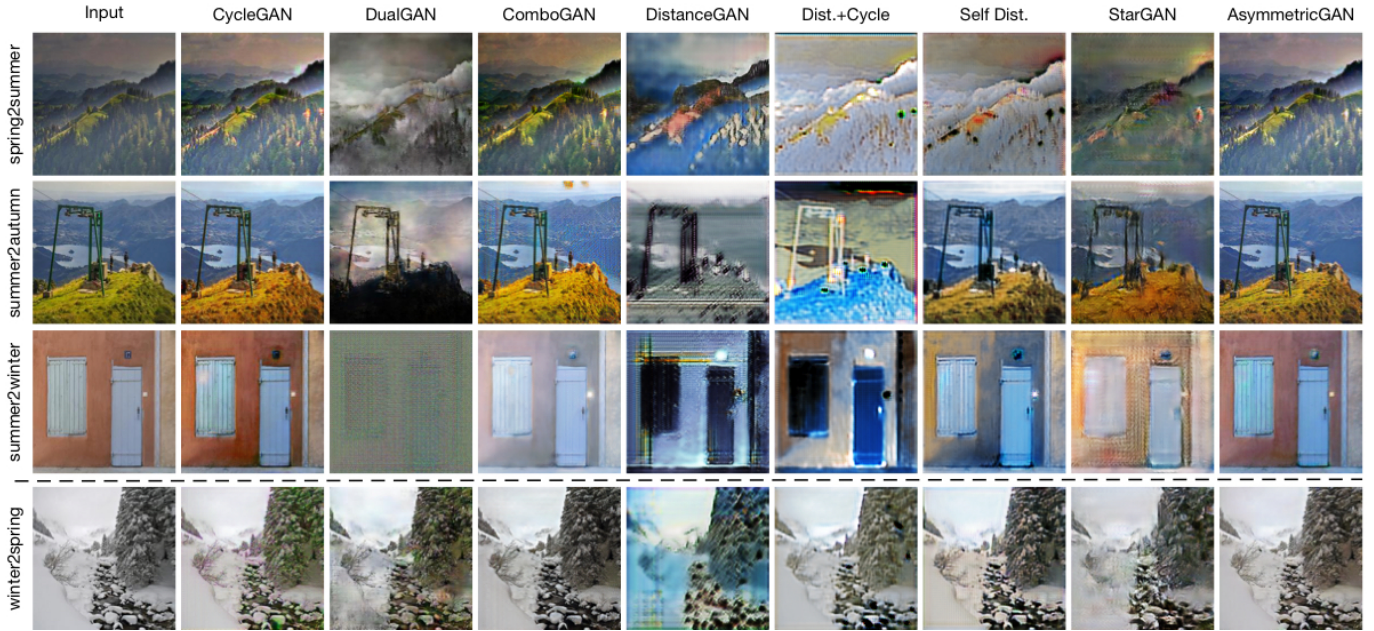
Fig. 8: Different methods for multi-domain season translation trained on the Alps dataset.

TABLE II: AMT "real vs fake" study of multi-domain image-to-image translation on the Facades, AR Face, Alps, Bu3dfe datasets.

| Model | label→photo | photo→label | AR Face | Alps | Bu3dfe |
|---|---|---|---|---|---|
| CycleGAN [3] | 8.8±1.5 | 4.8±0.8 | **24.3±1.7** | 39.6±1.4 | 16.9±1.2 |
| DualGAN [4] | 0.6±0.2 | 0.8±0.3 | 1.9±0.6 | 18.2±1.8 | 3.2±0.4 |
| ComboGAN [5] | 4.1±0.5 | 0.2±0.1 | 4.7±0.9 | 34.3±2.2 | **25.3±1.6** |
| DistanceGAN [30] | 5.7±1.1 | 1.2±0.5 | 2.7±0.7 | 4.4±0.3 | 6.5±0.7 |
| Dist.+Cycle [30] | 0.3±0.2 | 0.2±0.1 | 1.3±0.5 | 3.8±0.6 | 0.3±0.1 |
| Self Dist. [30] | 0.3±0.1 | 0.1±0.1 | 0.1±0.1 | 5.7±0.5 | 1.1±0.3 |
| StarGAN [6] | 3.5±0.7 | 1.3±0.3 | 4.1±1.3 | 8.6±0.7 | 9.3±0.9 |
| Pix2pix [2] | 4.6±0.5 | 1.5±0.4 | 2.8±0.6 | - | 3.6±0.5 |
| BicycleGAN [44] | 5.4±0.6 | 1.1±0.3 | 2.1±0.5 | - | 2.7±0.4 |
| AsymmetricGAN, Fully-Sharing | 4.6±0.9 | 2.4±0.4 | 6.8±0.6 | 15.4±1.9 | 13.1±1.3 |
| AsymmetricGAN, Partially-Sharing | 8.2±1.2 | 3.6±0.7 | 16.8±1.2 | 36.7±2.3 | 18.9±1.1 |
| AsymmetricGAN, No-Sharing | **10.3±1.6** | **5.6±0.9** | 22.8±1.9 | **47.7±2.8** | 23.6±1.7 |

TABLE III: IS and CA of multi-domain image-to-image translation on the Facades, AR Face, Alps and Bu3dfe datasets.

| Model | Facades | AR Face | | Alps | | Bu3dfe | |
|---|---|---|---|---|---|---|---|
| | IS ↑ | IS ↑ | CA (%) ↑ | IS ↑ | CA (%) ↑ | IS ↑ | CA (%) ↑ |
| CycleGAN [3] | 3.6098 | 2.8321 | @1:27.333 | 4.1734 | @1:42.250 | 1.8173 | @1:48.292, @5:94.167 |
| DualGAN [4] | 3.7495 | 1.9148 | @1:28.667 | 4.2661 | @1:53.488 | 1.7176 | @1:40.000, @5:90.833 |
| ComboGAN [5] | 3.1289 | 2.4750 | @1:28.250 | 4.2438 | @1:62.750 | 1.7887 | @1:40.459, @5:90.714 |
| DistanceGAN [30] | 3.9988 | 2.3455 | @1:26.000 | 4.8047 | @1:31.083 | 1.8974 | @1:46.458, @5:90.000 |
| Dist.+Cycle [30] | 2.6897 | **3.5554** | @1:14.667 | **5.9531** | @1:29.000 | 3.4618 | @1:26.042, @5:79.167 |
| Self Dist. [30] | 3.8155 | 2.1350 | @1:21.333 | 5.0584 | @1:34.917 | **3.4620** | @1:10.625, @5:74.167 |
| StarGAN [6] | **4.3182** | 2.0290 | @1:26.250 | 3.3670 | @1:**65.375** | 1.5640 | @1:52.704, @5:94.898 |
| Pix2pix [2] | 3.6664 | 2.2849 | @1:22.667 | - | - | 1.4575 | @1:44.667, @5:91.750 |
| BicycleGAN [44] | 3.2217 | 2.0859 | @1:28.000 | - | - | 1.7373 | @1:45.125, @5:93.125 |
| AsymmetricGAN, Fully-Sharing | 4.2615 | 2.3875 | @1:28.396 | 3.6597 | @1:61.125 | 1.9728 | @1:52.985, @5:95.165 |
| AsymmetricGAN, Partially-Sharing | 4.1689 | 2.4846 | @1:28.835 | 4.0158 | @1:62.325 | 1.5896 | @1:53.456, @5:95.846 |
| AsymmetricGAN, No-Sharing | 4.0819 | 2.6522 | @1:**29.667** | 4.3773 | @1:63.667 | 1.8714 | @1:**55.625**, @5:**96.250** |

simplest network structure, only consisting of 7 non-linear transformation operations with each using a convolution and a ReLU layer. The number of parameters of this architecture is 2.9K. (ii) Architecture II uses an encoder-decoder network with a symmetric structure, which has 1.3M parameters. (iii) Architecture III employs the same encoder-decoder network as architecture II while adding extra 6 residual blocks. It has the largest network capability (8.4M parameters) in the considered three. In the multi-domain image-to-image translation task, the final target is to make the network have a good generation ability. Thus the translation generator $G^t$

is expected to use a more powerful architecture, while the reconstruction generator $G^r$ can employ a lighter structure. We consider the following combinations for the translation and the reconstruction generators: (i) In S1, $G^t$ uses the generator architecture III, and $G^r$ uses the generator architecture I. (ii) In S2, $G^t$ uses the architecture III, and $G^r$ uses the generator architecture II. (iii) In S3, $G^t$ and $G^r$ use the same generator architecture III. We report the results of the running-time for training one epoch, the total number of generator parameters and the quantitative performance on the Bu3dfe dataset. Our results are compared with the most related model StarGAN
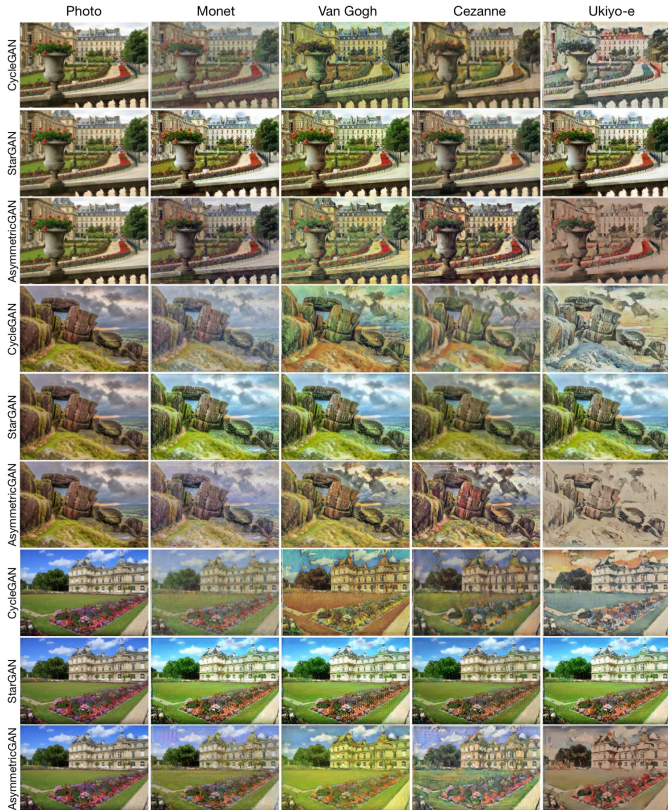
Fig. 9: Different methods for multi-domain painting style transfer trained on the Collection Style dataset.

TABLE IV: Comparison results of the multi-domain image-to-image translation task on the Collection Style dataset.

| Model | AMT ↑ | FID ↓ | CA (%) ↑ |
|---|---|---|---|
| CycleGAN [3] | 16.8±1.9 | 47.4823 | 73.72 |
| StarGAN [6] | 13.9±1.4 | 58.1562 | 44.63 |
| AsymmetricGAN | **19.8±2.4** | **43.7473** | **78.84** |
| Real Data | - | - | 91.34 |

TABLE V: Comparison results of the multi-domain image-to-image translation task on the RaFD dataset.

| Model | AMT ↑ | IS ↑ | FID ↓ |
|---|---|---|---|
| CycleGAN [3] | 19.5 | 1.6942 | 52.8230 |
| StarGAN [6] | 24.7 | 1.6695 | 51.6929 |
| AsymmetricGAN | **29.1** | **1.7187** | **51.2765** |

as shown in Table VI. It is clear to observe that the proposed AsymmetricGAN achieves much better performance than Star-GAN on all metrics when we consider only a light-weight generator structure for the reconstruction generator (S1). By so doing, the number of parameters for ours is only 2.9k more than StarGAN, while the performance is significantly boosted, which shows that the distinct network designs for different generators are very important for learning better both the generators, demonstrating our initial motivation.

*(2) Generation Performance v.s. Network Complexity.* Through a comparison of the performance among the setting S1, S2, S3 in Table VI, we also observe that using a more complex generator indeed improves the generation performance, while the network capacity is consequently increased. Specifically, from S2 to S3, the number of parameter changes from 8.4M+1.3M to 8.4M+8.4M. Although the parameters remarkably increase, the generation performance has slight improvement (IS and

CA metrics are even worse), meaning that the balance between the network complexity and the generation performance should be also an important consideration in designing a good GAN.
*(3) Model Component Analysis.* We conduct an ablation study of the proposed AsymmetricGAN on several datasets, i.e., Facades, AR Face and Bu3dfe. We report the generation performance without using the conditional identity preserving loss (I), multi-scale SSIM loss (S), color cycle-consistency loss (C) and double discriminators strategy (D), respectively. We also employ two different discriminators as in [47], [10] to further improve our generation performance. In order to investigate the parameter-sharing strategy of the asymmetric generator, we perform experiments on different schemes including: (i) Fully-sharing, i.e., the two generators share the same parameters. (ii) Partially-sharing, i.e., only the encoder part shares the same parameters. (iii) No-sharing, i.e., two independent generators. The basic generator structure follows StarGAN [6]. Quantitative results of both AMT score and CA are reported in Table VII. We observe that without using double discriminators slightly degrades performance, meaning that the proposed model can achieve good results trained using the proposed asymmetric dual generators and one discriminator. However, removing the conditional identity preserving loss, multi-scale SSIM loss and color cycle-consistency loss substantially degrades the performance, meaning that the proposed joint optimization objectives are particularly important to stabilize the training process and thus produce much better generation performance. The results of different parameter sharing strategies are shown in Tables II, III and VIII, we observe that different-level parameter sharing influences both the generation performance and the model capacity, demonstrating our initial motivation.
*(4) Overall Model Capacity Analysis.* We also compare the overall model capacity with several state-of-the-art methods. The number of trained models and the number of model parameters on the Bu3dfe dataset for $m$ domains are presented in Table VIII. We note that BicycleGAN and Pix2pix are supervised models, thus they need to train $A_m^2$ models for $m$ domains. CycleGAN, DiscoGAN, DualGAN, DistanceGAN are unsupervised methods, and they require $C_m^2$ models to learn $m$ domains, but each model of them contains two generators and two discriminators. ComboGAN needs only $m$ models to learn all the mappings of $m$ domains, while StarGAN and the proposed AsymmetricGAN only need to train one model to learn all the mappings of $m$ domains. In addition, we report the number of parameters on the Bu3dfe dataset in Table VIII. This dataset contains seven different facial expressions, which means $m$=7. Note that DualGAN employs the fully connected layers in its generators, which brings a significantly larger number of parameters. CycleGAN and DistanceGAN adopt the same generator and discriminator architectures, which means they have the same number of parameters. The proposed AsymmetricGAN uses fewer parameters compared with the other baselines except for StarGAN, but we achieve significantly better generation results in most metrics as shown in Tables V, IV, II and III. When we adopt a parameter-sharing strategy, our generation performance is only slightly lower (but still outperforming StarGAN) while

TABLE VI: Comparison results with different generator settings of AsymmetricGAN on the Bu3dfe dataset.

| Model | #Time | #Parameters | AMT ↑ | IS ↑ | CA (%) ↑ |
|---|---|---|---|---|---|
| StarGAN [6] | 2m23s | 8.4M | 9.3±0.9 | 1.5640 | @1:52.704, @5:94.898 |
| S1: $G^t$ (Architecture III), $G^r$ (Architecture I) | 2m27s | 8.4M+2.9K | 18.9±1.4 | 1.8790 | @1:55.575, @5:96.014 |
| S2: $G^t$ (Architecture III), $G^r$ (Architecture II) | 2m29s | 8.4M+1.3M | 20.1±1.4 | 1.9293 | @1:56.173, @5:97.112 |
| S3: $G^t$ (Architecture III), $G^r$ (Architecture III) | 2m33s | 8.4M+8.4M | 23.6±1.7 | 1.8714 | @1:55.625, @5:96.250 |

TABLE VII: Ablation study of AsymmetricGAN on the Facades, AR Face and Bu3dfe datasets for the multi-domain image-to-image translation task. All: full version of AsymmetricGAN, I: Identity preserving loss, S: multi-scale SSIM loss, C: Color cycle-consistency loss, D: Double discriminators strategy.

| Model | Label→Photo | Photo→Label | AR Face | | Bu3dfe | |
|---|---|---|---|---|---|---|
| | AMT ↑ | AMT ↑ | AMT ↑ | CA (%) ↑ | AMT ↑ | CA (%) ↑ |
| All | **10.3±1.6** | **5.6±0.9** | **22.8±1.9** | **@1:29.667** | **23.6±1.7** | **@1:55.625, @5:96.250** |
| All - I | 2.6±0.4 | 4.2±1.1 | 4.7±0.8 | @1:29.333 | 16.3±1.1 | @1:53.739, @5:95.625 |
| All - S - C | 4.4±0.6 | 4.8±1.3 | 8.7±0.6 | @1:28.000 | 14.4±1.2 | @1:42.500, @5:95.417 |
| All - S - C - I | 2.2±0.3 | 3.9±0.8 | 2.1±0.4 | @1:24.667 | 13.6±1.2 | @1:41.458, @5:95.208 |
| All - D | 9.0±1.5 | 5.3±1.1 | 21.7±1.7 | @1:28.367 | 22.3±1.6 | @1:53.375, @5:95.292 |
| All - D - S | 3.3±0.7 | 4.5±1.1 | 14.7±1.7 | @1:27.333 | 20.1±1.4 | @1:42.917, @5:91.250 |
| All - D - C | 8.7±1.3 | 5.1±0.9 | 19.4±1.5 | @1:28.000 | 21.6±1.4 | @1:45.833, @5:93.875 |

TABLE VIII: Comparison of the overall model capacity of different models with the number of image domain $m=7$ for the multi-domain image-to-image translation task.

| Model | #Models | #Parameters |
|---|---|---|
| Pix2pix [2] | $A_m^2 = m(m-1)$ | 57.2M×42 |
| BicycleGAN [44] | | 64.3M×42 |
| CycleGAN [3] | | 52.6M×21 |
| DiscoGAN [28] | $C_m^2 = \frac{m(m-1)}{2}$ | 16.6M×21 |
| DualGAN [4] | | 178.7M×21 |
| DistanceGAN [30] | | 52.6M×21 |
| ComboGAN [5] | $m$ | 14.4M×7 |
| StarGAN [6] | 1 | 53.2M×1 |
| AsymmetricGAN, Fully-Sharing | 1 | 53.2M×1 |
| AsymmetricGAN, Partial-Sharing | 1 | 53.8M×1 |
| AsymmetricGAN, No-Sharing | 1 | 61.6M×1 |

the number of parameters is comparable with StarGAN.

## B. Hand Gesture-to-Gesture Translation Task

Besides the unsupervised image-to-image translation task, we also conduct more experiments on supervised image-to-image translation, i.e., hand gesture-to-gesture translation, to validate the effectiveness of the proposed AsymmetricGAN.
**Datasets.** We follow GestureGAN [10] and employ the NTU Hand Digit [48] and Creative Senz3D [49] datasets to evaluate the proposed AsymmetricGAN. The number of train/test image pair for the NTU Hand Digit and Creative Senz3D datasets are 75,036/9,600 and 135,504/12,800, respectively.
**Baselines.** We compare the proposed AsymmetricGAN with the most related five works, i.e., GestureGAN [10], PG$^2$ [8], DPIG [26], PoseGAN [9] and SAMG [50]. All these five methods and the proposed AsymmetricGAN are paired image-to-image translation models.
**Metrics.** Following GestureGAN [10], we employ Peak Signal-to-Noise Ratio (PSNR), Fréchet Inception Distance (FID) [46], Fréchet ResNet Distance (FRD) [10] and user study (AMT) to evaluate the quality of generated images.
**Ablation Study.** We conduct ablation studies between SymmetricGAN and AsymmetricGAN on the NTU Hand Digit dataset to validate our motivation of the asymmetric network design. (i) SymmetricGAN has two separate generators with the identity structure for both $G^t$ and $G^r$, which

TABLE IX: Comparison between SymmetricGAN and AsymmetricGAN for the hand gesture-to-gesture translation task on the NTU Hand Digit dataset.

| Model | PSNR ↑ | AMT ↑ | FID ↓ | FRD ↓ | #Parameters |
|---|---|---|---|---|---|
| SymmetricGAN | 32.5740 | 27.9 | 6.8711 | 1.7519 | 11.388M*2 |
| AsymmetricGAN | **32.6686** | **29.7** | **6.7132** | **1.7341** | **11.388M+0.046M** |

has 11.388M*2=22.776M parameters totally. (ii) AsymmetricGAN also has two separate generators for both $G^t$ and $G^r$. However, the filters in first convolutional layer of $G^t$ and $G^r$ are 64 and 4, respectively. It means $G^t$ and $G^r$ have 11.388M and 0.046M parameters, respectively. Comparison results of both generation performance and network parameters are reports in Table IX. We observe that although the total number of parameters of AsymmetricGAN is much less than SymmetricGAN, AsymmetricGAN still achieves slightly better results than SymmetricGAN on all metrics, which validate our motivation of the asymmetric generator design.
**Comparison Against Baselines.** Qualitative results compared with several state-of-the-art approaches are shown in Fig. 10. We can see that the proposed AsymmetricGAN produces much more photo-realistic results with convincing details compared with other approaches, i.e., GestureGAN [10], PG$^2$ [8], DPIG [26], PoseGAN [9] and SAMG [50]. Moreover, we provide quantitative comparison with those methods. Results are shown in Table X. We note that our results are significantly much better than baseline models on both datasets.

## V. CONCLUSION

In this paper, we present a new Asymmetric Generative Adversarial Network (AsymmetricGAN), a robust and efficient GAN model that can perform both paired and unpaired image-to-image translations. The proposed asymmetric dual generators, allowing for different network architectures and different-level parameter sharing strategy, are designed for the image translation and image reconstruction tasks. Moreover, we explore jointly using different objective functions to optimize our AsymmetricGAN, and thus generating images with better fidelity and high quality. Extensive experimental results on different scenarios demonstrate that our AsymmetricGAN achieves more photo-realistic results and less model capacity

Fig. 10: Different methods for hand gesture-to-gesture translation trained on the NTU Hand Digit (Top) and Senz3D (Bottom) datasets.

TABLE X: Comparison with different models for hand gesture-to-gesture translation on the NTU Hand Digit and Senz3D datasets.

| Model | NTU Hand Digit | | | | Senz3D | | | |
|---|---|---|---|---|---|---|---|---|
| | PSNR ↑ | AMT ↑ | FID ↓ | FRD ↓ | PSNR ↑ | AMT ↑ | FID ↓ | FRD ↓ |
| PG² [8] | 28.2403 | 3.5 | 24.2093 | 2.6319 | 26.5138 | 2.8 | 31.7333 | 3.0933 |
| SAMG [50] | 28.0185 | 2.6 | 31.2841 | 2.7453 | 26.9545 | 2.3 | 38.1758 | 3.1006 |
| DPIG [26] | 30.6487 | 7.1 | 6.7661 | 2.6184 | 26.9451 | 6.9 | 26.2713 | 3.0846 |
| PoseGAN [9] | 29.5471 | 9.3 | 9.6725 | 2.5846 | 27.3014 | 8.6 | 24.6712 | 3.0467 |
| GestureGAN [10] | 32.6091 | 26.1 | 7.5860 | 2.5223 | 27.9749 | 22.6 | 18.4595 | 2.9836 |
| AsymmetricGAN | **32.6686** | **29.7** | **6.7132** | **1.7341** | **31.5624** | **28.1** | **12.4326** | **2.2011** |

than existing methods for both unsupervised and supervised image-to-image translation tasks. Finally, the proposed GAN model and training skills can be easily injected into other GAN frameworks.

## REFERENCES

[1] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *NIPS*, 2014. 1, 2

[2] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-image translation with conditional adversarial networks," in *CVPR*, 2017. 1, 2, 6, 8, 9, 11

[3] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *ICCV*, 2017. 1, 3, 4, 5, 6, 8, 9, 10, 11

[4] Z. Yi, H. Zhang, P. T. Gong *et al.*, "Dualgan: Unsupervised dual learning for image-to-image translation," in *ICCV*, 2017. 1, 3, 6, 8, 9, 11

[5] A. Anoosheh, E. Agustsson, R. Timofte, and L. Van Gool, "Combogan: Unrestrained scalability for image domain translation," in *CVPR Workshop*, 2018. 1, 3, 6, 9, 11

[6] Y. Choi, M. Choi, M. Kim, J.-W. Ha, S. Kim, and J. Choo, "Stargan: Unified generative adversarial networks for multi-domain image-to-image translation," in *CVPR*, 2018. 1, 2, 3, 4, 5, 6, 8, 9, 10, 11

[7] A. Odena, C. Olah, and J. Shlens, "Conditional image synthesis with auxiliary classifier gans," in *ICML*, 2017. 1, 4

[8] L. Ma, X. Jia, Q. Sun, B. Schiele, T. Tuytelaars, and L. Van Gool, "Pose guided person image generation," in *NIPS*, 2017. 1, 2, 3, 5, 6, 8, 11, 12

[9] A. Siarohin, E. Sangineto, S. Lathuilière, and N. Sebe, "Deformable gans for pose-based human image generation," in *CVPR*, 2018. 1, 2, 3, 5, 11, 12

[10] H. Tang, W. Wang, D. Xu, Y. Yan, and N. Sebe, "Gesturegan for hand gesture-to-gesture translation in the wild," in *ACM MM*, 2018. 1, 2, 3, 4, 5, 6, 10, 11, 12

[11] S. Ruder, "An overview of multi-task learning in deep neural networks," *arXiv preprint arXiv:1706.05098*, 2017. 2

[12] T. Karras, S. Laine, and T. Aila, "A style-based generator architecture for generative adversarial networks," in *CVPR*, 2019. 2

[13] A. Brock, J. Donahue, and K. Simonyan, "Large scale gan training for high fidelity natural image synthesis," in *ICLR*, 2019. 2

[14] Y. Guo, Q. Chen, J. Chen, Q. Wu, Q. Shi, and M. Tan, "Auto-embedding generative adversarial networks for high resolution image synthesis," *IEEE TMM*, 2019. 2

[15] M. Mirza and S. Osindero, "Conditional generative adversarial nets," *arXiv preprint arXiv:1411.1784*, 2014. 2

[16] H. Tang, W. Wang, S. Wu, X. Chen, D. Xu, N. Sebe, and Y. Yan, "Expression conditional gan for facial expression-to-expression translation," in *ICIP*, 2019. 2

[17] H. Tang, D. Xu, G. Liu, W. Wang, N. Sebe, and Y. Yan, "Cycle in cycle generative adversarial networks for keypoint-guided image generation," in *ACM MM*, 2019. 2

[18] Y. Yan, B. Ni, W. Zhang, J. Xu, and X. Yang, "Structure-constrained motion sequence generation," *IEEE TMM*, 2018. 2

[19] H. Zhang, T. Xu, H. Li, S. Zhang, X. Wang, X. Huang, and D. N. Metaxas, "Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks," in *ICCV*, 2017. 2

[20] S. Reed, Z. Akata, X. Yan, L. Logeswaran, B. Schiele, and H. Lee, "Generative adversarial text-to-image synthesis," in *ICML*, 2016. 2

[21] T.-C. Wang, M.-Y. Liu, J.-Y. Zhu, A. Tao, J. Kautz, and B. Catanzaro, "High-resolution image synthesis and semantic manipulation with conditional gans," in *CVPR*, 2018. 2, 6

[22] H. Tang, D. Xu, N. Sebe, Y. Wang, J. J. Corso, and Y. Yan, "Multi-channel attention selection gan with cascaded semantic guidance for cross-view image translation," in *CVPR*, 2019. 2, 6

[23] B. Duan, W. Wang, H. Tang, H. Latapie, and Y. Yan, "Cascade attention guided residue learning gan for cross-modal translation," *arXiv preprint arXiv:1907.01826*, 2019. 2

[24] X. Guo, R. Nie, J. Cao, D. Zhou, L. Mei, and K. He, "Fusegan: Learning to fuse multi-focus image via conditional generative adversarial network," *IEEE TMM*, 2019. 2

[25] T.-C. Wang, M.-Y. Liu, J.-Y. Zhu, G. Liu, A. Tao, J. Kautz, and B. Catanzaro, "Video-to-video synthesis," in *NeurIPS*, 2018. 2

[26] L. Ma, Q. Sun, S. Georgoulis, L. Van Gool, B. Schiele, and M. Fritz, "Disentangled person image generation," in *CVPR*, 2018. 3, 5, 11, 12

[27] L. Chen, L. Wu, Z. Hu, and M. Wang, "Quality-aware unpaired image-to-image translation," *IEEE TMM*, 2019. 3

[28] T. Kim, M. Cha, H. Kim, J. Lee, and J. Kim, "Learning to discover cross-domain relations with generative adversarial networks," in *ICML*, 2017. 3, 11

[29] H. Tang, D. Xu, N. Sebe, and Y. Yan, "Attention-guided generative adversarial networks for unsupervised image-to-image translation," in *IJCNN*, 2019. 3

[30] S. Benaim and L. Wolf, "One-sided unsupervised domain mapping," in *NIPS*, 2017. 3, 6, 9, 11

[31] Y.-F. Zhou, R.-H. Jiang, X. Wu, J.-Y. He, S. Weng, and Q. Peng, "Branchgan: Unsupervised mutual image-to-image transfer with a single encoder and dual decoders," *IEEE TMM*, 2019. 3

[32] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE TIP*, vol. 13, no. 4, pp. 600–612, 2004. 4

[33] Z. Wang, E. P. Simoncelli, and A. C. Bovik, "Multiscale structural similarity for image quality assessment," in *Asilomar Conference on Signals, Systems and Computers*, 2003. 4

[34] X. Mao, Q. Li, H. Xie, R. Y. Lau, Z. Wang, and S. P. Smolley, "Least squares generative adversarial networks," in *ICCV*, 2017. 4

[35] M. Arjovsky, S. Chintala, and L. Bottou, "Wasserstein gan," in *ICML*, 2017. 4

[36] A. Shrivastava, T. Pfister, O. Tuzel, J. Susskind, W. Wang, and R. Webb, "Learning from simulated and unsupervised images through adversarial training," in *CVPR*, 2017. 5

[37] R. Tyleček and R. Šára, "Spatial pattern templates for recognition of objects with regular structure," in *GCPR*, 2013. 6

[38] A. M. Martinez, "The ar face database," *CVC TR*, 1998. 6

[39] L. Yin, X. Wei, Y. Sun, J. Wang, and M. J. Rosato, "A 3d facial expression database for facial behavior research," in *FGR*, 2006. 6

[40] O. Langner, R. Dotsch, G. Bijlstra, D. H. Wigboldus, S. T. Hawk, and A. Van Knippenberg, "Presentation and validation of the radboud faces database," *Cognition and emotion*, vol. 24, no. 8, pp. 1377–1388, 2010. 6

[41] J. Johnson, A. Alahi, and L. Fei-Fei, "Perceptual losses for real-time style transfer and super-resolution," in *ECCV*, 2016. 6, 8

[42] T. Simon, H. Joo, I. Matthews, and Y. Sheikh, "Hand keypoint detection in single images using multiview bootstrapping," in *CVPR*, 2017. 6

[43] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *ICLR*, 2015. 6

[44] J.-Y. Zhu, R. Zhang, D. Pathak, T. Darrell, A. A. Efros, O. Wang, and E. Shechtman, "Toward multimodal image-to-image translation," in *NIPS*, 2017. 6, 9, 11

[45] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen, "Improved techniques for training gans," in *NIPS*, 2016. 8

[46] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter, "Gans trained by a two time-scale update rule converge to a local nash equilibrium," in *NIPS*, 2017. 8, 11

[47] T. Nguyen, T. Le, H. Vu, and D. Phung, "Dual discriminator generative adversarial nets," in *NIPS*, 2017. 10

[48] Z. Ren, J. Yuan, J. Meng, and Z. Zhang, "Robust part-based hand gesture recognition using kinect sensor," *IEEE TMM*, vol. 15, no. 5, pp. 1110–1120, 2013. 11

[49] A. Memo and P. Zanuttigh, "Head-mounted gesture controlled interface for human-computer interaction," *Springer MTA*, pp. 1–27, 2016. 11

[50] Y. Yan, J. Xu, B. Ni, W. Zhang, and X. Yang, "Skeleton-aided articulated motion generation," in *ACM MM*, 2017. 11, 12