

Total Generate: Cycle in Cycle Generative Adversarial Networks for Generating Human Faces, Hands, Bodies, and Natural Scenes

Hao Tang and Nicu Sebe

Abstract—We propose a novel and unified Cycle in Cycle Generative Adversarial Network (C2GAN) for generating human faces, hands, bodies, and natural scenes. Our proposed C2GAN is a cross-modal model exploring the joint exploitation of the input image data and guidance data in an interactive manner. C2GAN contains two different generators, i.e., an image-generation generator and a guidance-generation generator. Both generators are mutually connected and trained in an end-to-end fashion and explicitly form three cycled subnets, i.e., one image generation cycle and two guidance generation cycles. Each cycle aims at reconstructing the input domain and simultaneously produces a useful output involved in the generation of another cycle. In this way, the cycles constrain each other implicitly providing complementary information from both image and guidance modalities and bringing an extra supervision gradient across the cycles, facilitating a more robust optimization of the whole model. Extensive results on four guided image-to-image translation subtasks demonstrate that the proposed C2GAN is effective in generating more realistic images compared with state-of-the-art models. The code is available at <https://github.com/Ha0Tang/C2GAN>.

Index Terms—GANs, Cycle in Cycle, Cycle Consistency, Guided Image-to-Image Translation

I. INTRODUCTION

In this work, we focus on how to generate a target image given an input image. This has many application scenarios such as human-computer interaction, entertainment, virtual reality, and data augmentation. However, this task is challenging since it needs a high-level semantic understanding of the image mapping between the input and the output domains. Recently, Generative Adversarial Networks (GANs) [1] have shown the potential to solve this challenging task. GANs have produced promising results in many tasks such as image generation [2], image inpainting [3], and cross-modal translation [4].

Recent works have developed powerful image-to-image translation systems, e.g., Pix2pix [5] and GauGAN [6] in supervised settings, and CycleGAN [7] and DualGAN [8] in unsupervised settings. However, these methods are tailored to merely two domains at a time and scaling them to more requires a quadratic number of models to be trained. For instance, with m different image domains, CycleGAN and Pix2pix need to train $m(m-1)/2$ and $m(m-1)$ models, respectively. To overcome this, Choi et al. propose StarGAN [9], in which a single generator/discriminator performs image-to-image translation for multiple domains. However, StarGAN

is not effective in handling some specific image-to-image translation tasks such as human pose generation [10], [11], hand gesture generation [12], and cross-view image translation [13], in which image generation could involve infinite image domains since human body, hand gesture, and natural scene in the wild can have arbitrary poses, sizes, appearances, locations, and viewpoints.

To address these limitations, many methods are proposed to generate images based on extra semantic guidance, such as object keypoints [14], [10], human skeletons [11], [12], or segmentation maps [15], [16], [6], [17], [18]. For instance, Song et al. [14] propose a G2GAN framework for facial expression synthesis based on facial landmarks. Siarohin et al. [11] introduce a PoseGAN model for pose-based human image generation conditioned on human body skeletons. Regmi and Borji [13] propose both X-Fork and X-Seq for cross-view image translation conditioned on segmentation maps. However, the current state-of-the-art guided image-to-image translation methods such as PG2 [10], PoseGAN [11], X-Fork [13], and X-Seq [13] have two main issues: 1) they directly transfer a source image and the target guidance to the target domain (i.e., $[I_x, L_y] \xrightarrow{G_i} I'_y$ in Fig. 1), without considering the mutual translation between each other, while the translation across different image and guidance modalities in a unified framework would bring rich cross-modal information; 2) they simply employ the guidance data as input to guide the generation process, without involving the generated guidance as supervisory signals to further improve the network optimization. Both issues lead to unsatisfactory results.

To address both, we propose a novel and unified Cycle In Cycle Generative Adversarial Network (C2GAN), in which three cycled sub-nets are explicitly formed to learn both image and guidance modalities in a joint model. The framework of the proposed C2GAN is shown in Fig. 1. Specifically, to address the first limitation, C2GAN contains an image cycle, i.e., I2I2I ($[I_x, L_y] \xrightarrow{G_i} [I'_y, L_x] \xrightarrow{G_i} I'_x$), which aims at reconstructing the input and further refines the generated images I'_y . To address the second limitation, the guidance information (such as the human body skeleton) in C2GAN is not only utilized as input but also acts as output, meaning that the guidance is also a generative objective. The input and output of the guidance are connected by two novel guidance cycles, i.e., G2I2G ($[I_x, L_y] \xrightarrow{G_i} I'_y \xrightarrow{G_g} L'_y$) and G2R2G ($[I'_y, L_x] \xrightarrow{G_i} I'_x \xrightarrow{G_g} L'_x$), where G_i and G_g denote an image and a guidance generator, respectively. In this

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.

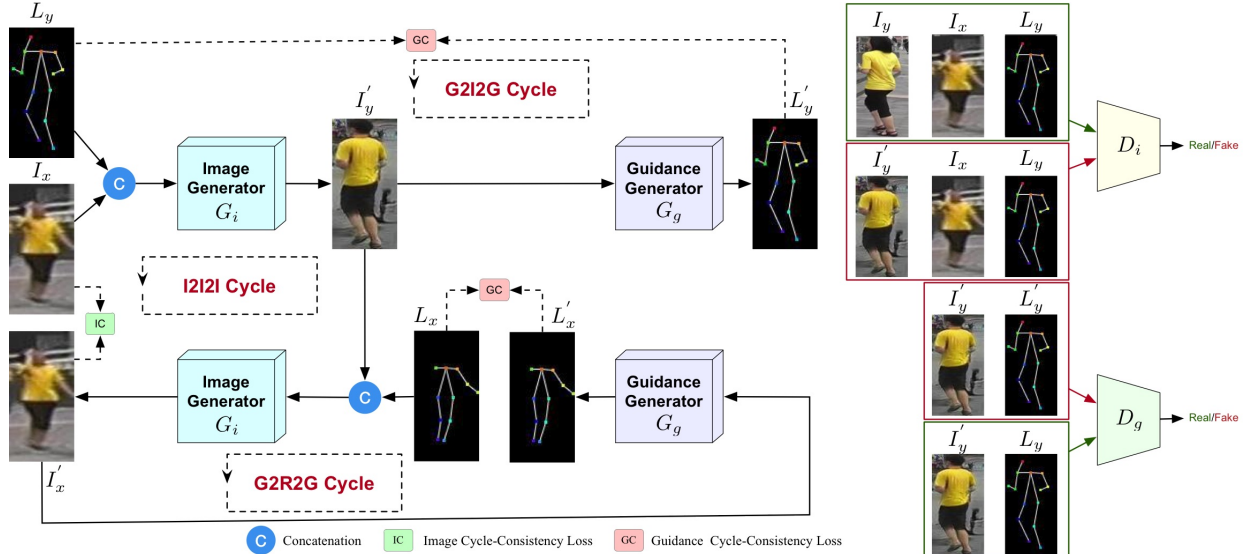


Fig. 1: Overview of the proposed C2GAN, which consists of two types of generators, i.e., image generator G_i and guidance generator G_g . Parameter-sharing strategies can be used in between the image or the guidance generators to reduce the model capacity. During the training stage, two generators G_i and G_g are explicitly connected and trained by three cycles, i.e., the image cycle I2I2I: $[I_x, L_y] \xrightarrow{G_i} [I'_y, L_x] \xrightarrow{G_i} I'_x$ and two guidance cycles G2I2G: $[I_x, L_y] \xrightarrow{G_i} I'_y \xrightarrow{G_g} L'_y$, G2R2G: $[I'_y, L_x] \xrightarrow{G_i} I'_x \xrightarrow{G_g} L'_x$. The right side of the figure shows the cross-modal discriminators (i.e., D_i and D_g) for a better network optimization.

way, guidance cycles can provide weak supervision to the generated images I'_y . The intuition behind the guidance cycles is that if the generated guidance is very close to the real guidance, then the corresponding images should be similar (see Fig. 2). In other words, a better guidance generation will boost the performance of image generation, and conversely the improved image generation will further facilitate the guidance generation. The proposed three cycles inherently constraint each other in an end-to-end training fashion. Moreover, for a better optimization of the proposed three cycles, we further propose two novel cycle losses, i.e., Image Cycle-consistency loss (IC) and Guidance Cycle-consistency loss (GC). With both cycle losses, each cycle can benefit from each other in a joint learning way. We also propose two cross-modal discriminators corresponding to the generators.

Our contributions can be summarized as follows:

- We propose C2GAN, a novel and unified cross-modal generative adversarial network for guided image-to-image translation tasks, which organizes the guidance and image data in an interactive manner, instead of using as input only the guidance information.
- The proposed cycle in cycle network structure is a new design which explores the effective use of cross-modal information for guided image-to-image translation tasks. The designed cycled subnetworks connect different modalities and implicitly constraint each other, leading to extra supervision signals for better image generation. We also investigate cross-modal discriminators and cycle losses for a more robust network optimization.
- Extensive results on four challenging guided image-to-image translation tasks, i.e., person image generation, facial expression generation, hand gesture-to-gesture translation, and cross-view image translation demonstrate the effectiveness of the proposed C2GAN and show more photorealistic

images compared with state-of-the-art models.

Part of this work has been published in [19]. We extend it in numerous ways: 1) A more detailed analysis is presented in “Introduction” section by giving a deeper analysis on the motivation and the difference from relevant works. 2) We extend the model proposed in [19] to a unified GAN framework for handling different guided image-to-image translation tasks. 3) We present an in-depth description of the proposed approach, providing all the architectural and implementation details of the method, with special emphasis on guaranteeing the reproducibility of the experiments. 4) We substantially extend the experimental evaluation.

II. RELATED WORK

Generative Adversarial Networks (GANs) have shown the capability of generating high-quality images. To generate images that meet user requirements, Conditional GAN (CGAN) [20] is employing the conditioned guidance information to guide the image generation process. A CGAN model always combines a vanilla GAN and an external source, such as segmentation maps [21], [15], [17], conditional images [5], and attention maps [22], [23]. However, synthesizing images based on global constraints does not provide control over pose, object location, or shape.

Image-to-Image Translation models use input-output data to learn a mapping between the source domain and the target domain. Isola et al. propose Pix2pix [5], which employs a CGAN to learn a image mapping from the input domain to the output domain. Moreover, unpaired image-to-image translation approaches [24], [8], [23], [22] have been proposed to learn the mapping without paired training data. However, these existing image translation models are inefficient and ineffective as indicated in the introduction section. Most importantly, these aforementioned approaches cannot handle some specific

guided image-to-image translation tasks, such as person image generation [11], and hand gesture-to-gesture translation [12]. **Guided Image-to-Image Translation.** To address these aforementioned limitations, several works have been proposed to generate images based on object keypoints [10], [25], human hand/body skeleton [26], [27], [28], and scene segmentation map [15], [29], [17]. For instance, Wang et al. [25] propose a Conditional MultiMode Network (CMM-Net) for facial landmark guided smile generation. Tang et al. [12] propose a novel GestureGAN to perform the hand gesture-to-gesture translation task conditioned on an input image and several novel hand skeletons. Park et al. [6] propose a novel spatially-adaptive normalization for the semantic image synthesis task based on semantic labels. These methods only focus on a single image generation task.

We propose a multitask framework aiming at handling two tasks using a single network, i.e., image generation and guidance generation. During the training stage, the two generation tasks are restricted mutually by the proposed three cycles and then benefit from each other. To the best of our knowledge, the proposed C2GAN is the first attempt to generate both the image and the guidance domain in an interactive generation strategy within a unified cycle in cycle GAN model, for guided image-to-image translation tasks.

III. CYCLE IN CYCLE GAN (C2GAN)

A. Model Overview

The proposed C2GAN learns two different generators in a single network, i.e., image generator and guidance generator. The two generators are mutually connected through three novel generative adversarial cycles, i.e., one image-oriented cycle and two guidance-oriented cycles. In the training stage, these three cycles are jointly optimized in an end-to-end way and each generator can benefit from the others due to the richer cross-modal information and the crossing cycle supervision. The core framework of C2GAN is illustrated in Fig. 1.

B. Image-Domain Generative Adversarial Cycle

I2I2I Cycle. The image cycle I2I2I aims to generate the image I'_y by using the combination of the input image I_x and the target guidance L_y , and then reconstruct the input image I_x by using the combination of the generated image I'_y and the guidance L_x of image I_x :

$$[I_x, L_y] \xrightarrow{G_i} [I'_y, L_x] \xrightarrow{G_i} I'_x, \quad (1)$$

where G_i is the image generator.

Different from the previous guided image-to-image translation methods such as PG2 [10], X-Fork [13], and PoseGAN [11] employing only one mapping $[I_x, L_y] \xrightarrow{G_i} I'_y$, StarGAN [9] employs the target and the original domain labels l_y and l_x as extra guidance information to reconstruct the input image. However, StarGAN can only handle tasks which have a fixed number of the image categories. To solve this limitation, we replace the domain labels l_y and l_x in StarGAN by using the guidances L_y and L_x . The guidance can be object keypoints, human skeletons, or scene segmentation

maps. Specifically, I_x and L_y are first fed into the image generator G_i to generate the desired image I'_y . Next, the generated image I'_y and the guidance L_x are concatenated as the input of G_i to reconstruct the original image I_x . In this way, the forward and backward consistency can be guaranteed.

Image Generator. The U-Net architecture [30] is adopted for our image generator G_i . U-Net consists of an encoder and a decoder with skip connections between them. The generator G_i is used two times for generating image I'_y and reconstructing the original image I_x . To reduce the model capacity, the image generator G_i shares parameters between image generation and reconstruction. For image generation, the target of G_i is generating an image $I'_y = G_i(I_x, L_y)$ conditioned on the target guidance L_y which is similar to the real image I_y . For image reconstruction, the goal of G_i is recovering an image $I'_x = G_i(I'_y, L_x)$ that looks similar to the input image I_x . The image generator G_i learns a combined data distribution between the image generation and the image reconstruction by sharing parameters, meaning that G_i receives double data during the network optimization compared to those generators without using the parameter-sharing strategy.

Cross-Modal Image Discriminator. Different from previous works such as PG2 [10] employing a single-modal discriminator, we propose a novel cross-modal discriminator which receives both image and guidance data as input (Fig. 1). The image discriminator D_i receives two images and one guidance data as input. More specifically, D_i aims to distinguish between the generated triplet $[I_x, L_y, G_i(I_x, L_y)]$ and the real triplet $[I_x, L_y, I_y]$ during the image generation stage.

We further propose an image adversarial loss based on the vanilla adversarial loss [1], which can be expressed as:

$$\begin{aligned} & \mathcal{L}_{GAN}^i(G_i, D_i, I_x, I_y, L_y) \\ &= \mathbb{E}_{I_x, L_y, I_y \sim p_{\text{data}}(I_x, L_y, I_y)} [\log D_i([I_x, L_y, I_y])] \\ &+ \mathbb{E}_{I_x, L_y \sim p_{\text{data}}(I_x, L_y)} [\log(1 - D_i([I_x, L_y, G_i(I_x, L_y)]))], \end{aligned} \quad (2)$$

where the image generator G_i tries to minimize the image adversarial loss $\mathcal{L}_{GAN}^i(G_i, D_i, I_x, I_y, L_y)$ while the image discriminator D_i tries to maximize it.

Another image adversarial loss for the image reconstruction mapping $G_i : [I'_y, L_x] \rightarrow I'_x$ is defined as:

$$\begin{aligned} & \mathcal{L}_{GAN}^i(G_i, D_i, I_x, I_y, L_x) \\ &= \mathbb{E}_{I_x, L_x, I_y \sim p_{\text{data}}(I_x, L_x, I_y)} [\log D_i([I_y, L_x, I_x])] \\ &+ \mathbb{E}_{I'_y, L_x, I_y \sim p_{\text{data}}(I'_y, L_x, I_y)} [\log(1 - D_i([I_y, L_x, G_i(I'_y, L_x)]))], \end{aligned} \quad (3)$$

where the image discriminator D_i aims at distinguishing between the fake triplet $[I_y, L_x, G_i(I'_y, L_x)]$ and the real triplet $[I_y, L_x, I_x]$. Therefore, the overall image adversarial loss is the sum of Eq. (2) and Eq. (3):

$$\begin{aligned} & \mathcal{L}_{GAN}^i(G_i, D_i, I_x, I_y, L_x, L_y) \\ &= \mathcal{L}_{GAN}^i(G_i, D_i, I_x, I_y, L_y) + \mathcal{L}_{GAN}^i(G_i, D_i, I_x, I_y, L_x). \end{aligned} \quad (4)$$

Image Cycle-Consistency (IC) Loss. We also propose the IC

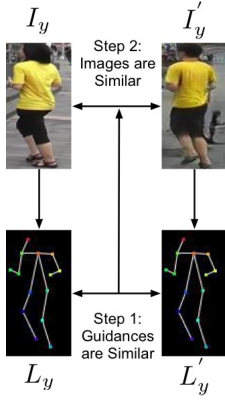


Fig. 2: The motivation of the guidance cycle. If the generated guidance L'_y is close to the real guidance L_y , then the corresponding images (i.e., I'_y and I_y) should be similar.

loss to better learn the image cycle I2I2I:

$$\begin{aligned} \mathcal{L}_{CYC}^i(G_i, I_x, L_x, L_y) \\ = \mathbb{E}_{I_x, L_x, L_y \sim p_{\text{data}}(I_x, L_x, L_y)} [\|G_i(G_i(I_x, L_y), L_x) - I_x\|_1]. \end{aligned} \quad (5)$$

The reconstructed image $I'_x = G_i(G_i(I_x, L_y), L_x)$ should closely match with the input image I_x . Notably, the image generator G_i is used two times with the parameter-sharing strategy and the $L1$ distance is adopted in Eq. (5) to compute a pixel-to-pixel difference between the recovered image I'_x and the real input image I_x .

C. Guidance-Domain Generative Adversarial Cycle

The motivation of the guidance cycle is that, if the generated guidance is similar to the real guidance, then the corresponding two images should be very close (see Fig. 2). The proposed C2GAN has two guidance cycles, i.e., G2I2G and G2R2G, as shown in Fig. 1. Both cycles can provide extra supervision information for better optimizing the image cycle I2I2I.

G2I2G Cycle. For the G2I2G cycle, $[I_x, L_y]$ is first fed into the image generator G_i to produce the target image I'_y . Then the guidance generator G_g tries to produce the guidance L'_y from the generated image I'_y . The generated guidance L'_y should be very close to the real guidance L_y . The formulation of the G2I2G cycle can be expressed as:

$$[I_x, L_y] \xrightarrow{G_i} I'_y \xrightarrow{G_g} L'_y. \quad (6)$$

G2R2G Cycle. For the G2R2G cycle, the generated image I'_y and guidance L_x are first concatenated, and then fed into G_i to produce the recovered image I'_x . Next, the guidance generator G_g generates the guidance L'_x from the recovered image I'_x . We assume that the generated guidance L'_x is very similar to the real guidance L_x . The G2R2G cycle can be formulated as:

$$[I'_y, L_x] \xrightarrow{G_i} I'_x \xrightarrow{G_g} L'_x. \quad (7)$$

Both generated guidances $L'_y = G_g(G_i(I_x, L_y))$ and $L'_x = G_g(G_i(I'_y, L_x))$ should have a close match to the real guidance L_y and L_x . Note that the guidance generator G_g could share parameters between these two cycles.

Guidance Generator. The U-Net structure [30] is employed for our guidance generator G_g . The input of G_g is an image and the output is a guidance. The guidance generator respectively produces $L'_y = G_g(I'_y)$ and $L'_x = G_g(I'_x)$ from the

generated images I'_y and I'_x , which further provide more supervision gradient to guide the image generator G_i to produce more realistic images.

Cross-Modal Guidance Discriminator. As shown in Fig. 1, the proposed guidance discriminator D_g is a cross-modal discriminator receiving both image and guidance data as inputs. Thus, the guidance adversarial loss for D_g can be defined as:

$$\begin{aligned} \mathcal{L}_{GAN}^g(G_g, D_g, I'_y, L_y) \\ = \mathbb{E}_{I'_y, L_y \sim p_{\text{data}}(I'_y, L_y)} [\log D_g([I'_y, L_y])] \\ + \mathbb{E}_{I'_y \sim p_{\text{data}}(I'_y)} [\log(1 - D_g([I'_y, G_g(I'_y)]))], \end{aligned} \quad (8)$$

where the guidance generator G_g aims to minimize the guidance adversarial loss $\mathcal{L}_{GAN}^g(G_g, D_g, I'_y, L_y)$ while the guidance discriminator D_g tries to maximize it. The discriminator D_g aims to distinguish between the fake pair $[I'_y, L'_y]$ and the real pair $[I'_y, L_y]$.

A similar guidance adversarial loss for the mapping function $G_g : I'_x \rightarrow L'_x$ is defined as:

$$\begin{aligned} \mathcal{L}_{GAN}^g(G_g, D_g, I'_x, L_x) \\ = \mathbb{E}_{I'_x, L_x \sim p_{\text{data}}(I'_x, L_x)} [\log D_g([I'_x, L_x])] \\ + \mathbb{E}_{I'_x \sim p_{\text{data}}(I'_x)} [\log(1 - D_g([I'_x, G_g(I'_x)]))], \end{aligned} \quad (9)$$

where the guidance discriminator D_g aims to distinguish between the fake pair $[I'_x, L'_x]$ and the real pair $[I'_x, L_x]$.

Thus, the total guidance adversarial loss is the sum of Eq. (8) and Eq. (9):

$$\begin{aligned} \mathcal{L}_{GAN}^g(G_g, D_g, I'_x, I'_y, L_x, L_y) \\ = \mathcal{L}_{GAN}^g(G_g, D_g, I'_y, L_y) + \mathcal{L}_{GAN}^g(G_g, D_g, I'_x, L_x). \end{aligned} \quad (10)$$

Guidance Cycle-Consistency (GC) Loss. A novel GC loss is further proposed to better learn both the guidance cycles (i.e., G2I2G and G2R2G), which can be expressed as:

$$\begin{aligned} \mathcal{L}_{CYC}^g(G_g, G_i, I_x, I'_y, L_x, L_y) \\ = \mathbb{E}_{I_x, L_y \sim p_{\text{data}}(I_x, L_y)} [\|G_g(G_i(I_x, L_y)) - L_y\|_1] \\ + \mathbb{E}_{I'_y, L_x \sim p_{\text{data}}(I'_y, L_x)} [\|G_g(G_i(I'_y, L_x)) - L_x\|_1], \end{aligned} \quad (11)$$

where the $L1$ distance is used to compute the pixel-to-pixel difference between the generated guidance (i.e., L'_y and L_y) and the corresponding real guidance (i.e., L_x and L_y).

During the training stage, the proposed guidance cycle-consistency loss can back-propagate errors from the guidance generator G_g to the image generator G_i facilitating the optimization of the image generator and then boosting the image generation performance.

D. Joint Optimization Objective

We follow existing methods [11], [16] and use the image pixel loss to reduce the changes between the generated image $I'_y = G_i(I_x, L_y)$ and the corresponding real one I_y :

$$\begin{aligned} \mathcal{L}_{PIXEL}^i(G_i, I_x, L_y, I_y) \\ = \mathbb{E}_{I_x, L_y, I_y \sim p_{\text{data}}(I_x, L_y, I_y)} [\|G_i(I_x, L_y) - I_y\|_1], \end{aligned} \quad (12)$$

where the $L1$ distance is adopted as the loss measurement of the image pixel loss. By doing so, more constrains can be added on both the image generator G_i .

Consequently, the complete objective loss of the proposed C2GAN framework is:

$$\begin{aligned} \mathcal{L}(G_i, G_g, D_i, D_g) &= \lambda_{gan}^i * \mathcal{L}_{GAN}^i + \lambda_{cyc}^i * \mathcal{L}_{CYC}^i + \lambda_{pixel}^i * \mathcal{L}_{PIXEL}^i \quad (13) \\ &+ \lambda_{gan}^g * \mathcal{L}_{GAN}^g + \lambda_{cyc}^g * \mathcal{L}_{CYC}^g, \end{aligned}$$

where λ_{gan}^i , λ_{cyc}^i , λ_{pixel}^i , λ_{gan}^g , and λ_{cyc}^g are parameters controlling the relative relation of objectives terms.

E. Implementation Details

Network Architecture. We adopt the U-Net architecture [30] consisting of an encoder and a decoder for our generators G_i and G_g . Moreover, we employ the PatchGAN discriminator [5] for our discriminators D_i and D_g , which has shown effectiveness in previous image-to-image translation works [5], [7]. The difference between a PatchGAN and a regular GAN discriminator is that the regular GAN maps from an image to a single scalar output, ‘real’ or ‘fake’, whereas the PatchGAN tries to classify if each $N \times N$ patch in an image is real or fake. By doing so, PatchGAN can alleviate the generation of visual artifacts and achieve better performance.

The amount of network parameters in the proposed method is twice that of Pix2pix [5] because our method contains two generators and two discriminators, while Pix2pix has only one generator and one discriminator. Although Pix2pix [5] has fewer network parameters, it is only able to translate between two domains and cannot handle some specific guided image-to-image translation tasks, such as person image generation, facial expression generation, hand gesture-to-gesture translation, and cross-view image translation. In contrast, our method is a universal method that can handle all these tasks.

Training Strategy. We follow the standard optimization method from [1] to optimize the proposed C2GAN, i.e., we alternate between one gradient descent step on G_i , D_i , G_g , and D_g . The Adam solver [31], with a learning rate of 0.0002, and momentum terms $\beta_1=0.5$, $\beta_2=0.999$, is adopted as our optimizer. For each task, we keep the same learning rate for the first half of the number of epochs and linearly decay the rate to zero over the next half of the epochs. Take facial expression generation as an example, we keep the same learning rate for the first 100 epochs and linearly decay the rate to zero over the next 100 epochs. The proposed C2GAN is trained end-to-end and can generate image and guidance simultaneously, then the generated guidances will benefit the quality of the generated images. Moreover, to slow down the rate of discriminators D_i and D_g relative to generators G_i and G_g , we divide the objectives by 2 while optimizing the discriminators.

The public software OpenFace [32] is employed to extract facial landmarks on the Radboud Faces dataset for facial expression generation. While OpenPose [33] is used to extract human hand and body skeleton on the Creative Senz3D and Market-1501 datasets for hand gesture-to-gesture translation and person image generation, respectively. Next, RefineNet

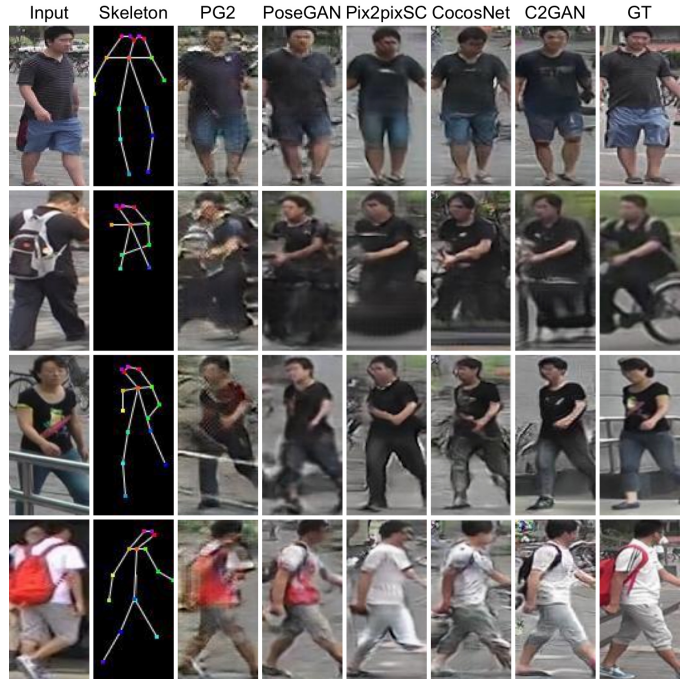


Fig. 3: Qualitative comparison of person image generation on the Market-1501 dataset. From left to right: Input, Body Skeleton, PG2 [10], PoseGAN [11], Pix2pixSC [36], CocosNet [37], C2GAN (Ours), and Ground Truth (GT).

[34] is employed to extract segmentation maps from the Dayton dataset for cross-view image translation.

Inference Strategy. During the inference stage, the proposed G2GAN receives an image I_x and a guidance L_y into the image generator G_i , and outputs a target image I'_y . At the same time, the guidance generator G_g receives the image I_x as input and outputs the corresponding guidance L'_x .

Parameter Setting. For a fair comparison, all competing models are trained for 200 epochs on the Radboud Faces dataset for facial expression generation. All models are trained around 90 epochs on person image generation. For hand gesture-to-gesture translation, all models are trained with 20 epochs. For cross-view image translation, we train the models for 35 epochs. The mask loss proposed in PG2 [10] is also used for person image generation. Our C2GAN is implemented using the public deep learning software PyTorch.

IV. EXPERIMENTS

A. Person Image Generation

Datasets. We employ the Market-1501 dataset [35] for the person image generation task. This dataset [35] is a challenging person-reID dataset containing 32,668 images of 1,501 persons collected from six surveillance cameras. We adopt the training and testing splits used in [11] and obtain 263,631 and 12,000 pairs for the training and testing subset.

Evaluation Metrics. We follow [11], [10] and adopt Inception Score (IS) [40], Structural Similarity (SSIM) [41], and their masked versions Mask-SSIM and Mask-IS as our evaluation metrics. Moreover, the AMT perceptual user study is adopted to evaluate the generated images by different models.

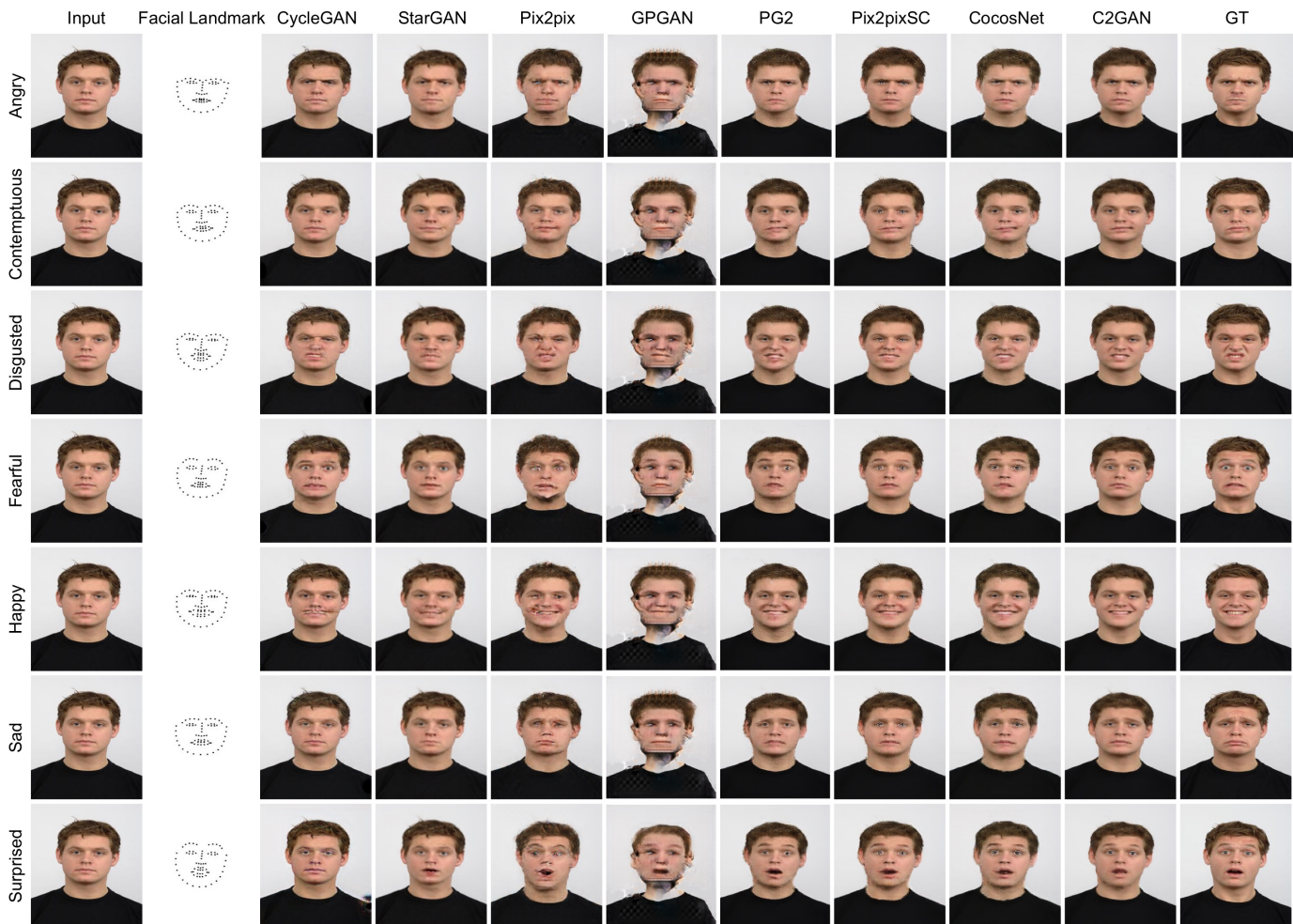


Fig. 4: Qualitative comparison of facial expression generation on the Radboud Faces dataset. From left to right: Input, Facial Landmark, CycleGAN [7], StarGAN [9], Pix2pix [5], GPGAN [38], PG2 [10], Pix2pixSC [36], CocosNet [37], C2GAN (Ours), and Ground Truth (GT).

TABLE I: Quantitative comparison of person image generation on the Market-1501 dataset. For all the metrics, higher is better.

Model	AMT (R2G) \uparrow	AMT (G2R) \uparrow	SSIM \uparrow	IS \uparrow	Mask-SSIM \uparrow	Mask-IS \uparrow
PG2 [10]	11.2	5.5	0.253	3.460	0.792	3.435
DPIG [39]	-	-	0.099	3.483	0.614	3.491
PoseGAN [11]	22.7	50.2	0.290	3.185	0.805	3.502
Pix2pixSC [36]	18.6	41.5	0.275	3.141	0.790	3.468
CocosNet [37]	20.1	45.7	0.280	3.275	0.801	3.514
C2GAN (Ours)	23.8	47.3	0.285	3.362	0.813	3.526
Real Data	-	-	1.000	3.860	1.000	3.360

State-of-the-Art Comparison. We compare C2GAN with PG2 [10], DPIG [39], PoseGAN [11], Pix2pixSC [36], and CocosNet [37]. Different from these models which focus on person image generation, our method is a general framework and learns image and guidance generation simultaneously in a joint network. Quantitative results are shown in Table I. C2GAN achieves better results than PG2, DPIG, Pix2pixSC, and CocosNet. Moreover, compared to PoseGAN [11], C2GAN yields better results on most metrics, i.e., AMT (R2G), IS, mask-SSIM, and mask-IS. Qualitative comparison results compared with PG2 and PoseGAN are shown in Fig. 3. C2GAN can generate more clear and visually plausible person images than both leading methods, validating the effectiveness of C2GAN. Moreover, our generated images are more similar to

the ground truth.

B. Facial Expression Generation

Datasets. We employ the Radboud Faces dataset [42] for the facial expression generation task. This dataset contains over 8,000 color face images with eight different facial expressions. We randomly select 67% of the images for training and the rest 33% images for testing. We remove the images in which the face is not correctly detected by OpenFace [32], then combine two different facial expression images of the same person to form an image pair for training. Therefore, 5,628 and 1,407 image pairs are obtained for training and testing, respectively. **Evaluation Metrics.** We first adopt the AMT user study to evaluate the generated images. Moreover, we employ

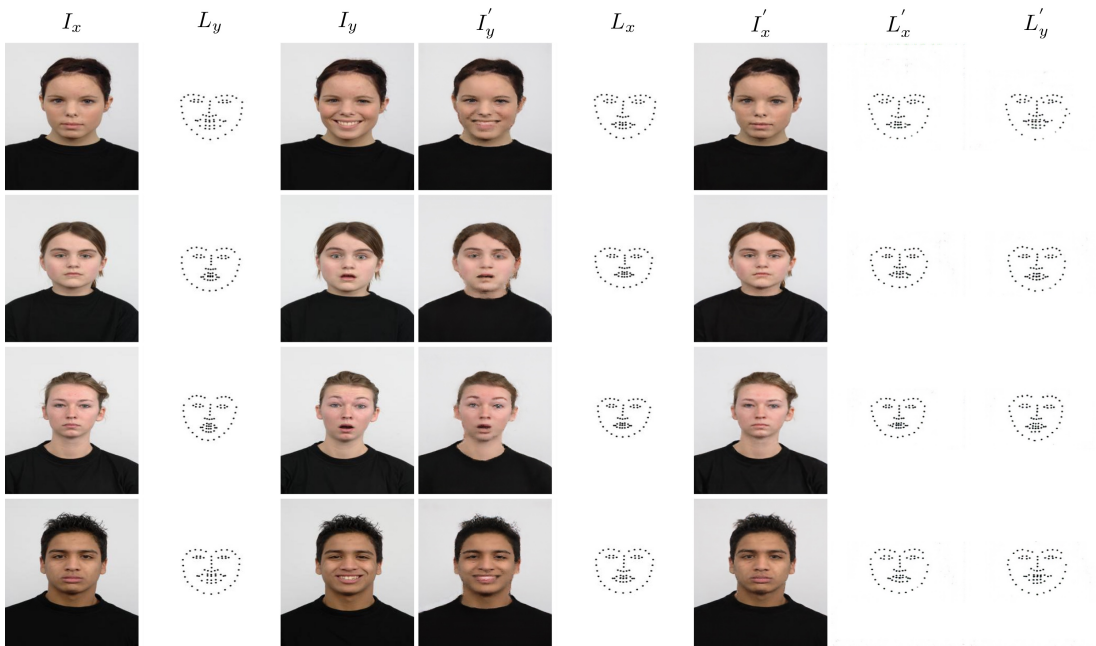


Fig. 5: Visualization of facial landmark generation on the facial expression generation task.

TABLE II: Quantitative comparison of facial expression generation on the Radboud Faces dataset. For all the metrics except LPIPS, higher is better.

Model	AMT \uparrow	SSIM \uparrow	PSNR \uparrow	LPIPS \downarrow
CycleGAN [7]	19.5	0.8307	18.8067	-
StarGAN [9]	24.7	0.8345	19.6451	-
Pix2pix [5]	13.4	0.8217	19.9971	0.1334
GPGAN [38]	0.3	0.8185	18.7211	0.2531
PG2 [10]	28.4	0.8462	20.1462	0.1130
Pix2pixSC [36]	30.8	0.8433	20.3584	0.1042
CocosNet [37]	31.3	0.8524	20.7915	0.0985
C2GAN (Ours)	34.2	0.8618	21.9192	0.0934

SSIM [41], Peak Signal-to-Noise Ratio (PSNR), and LPIPS [43] for quantitative measurements. SSIM and PSNR measure the image quality from a pixel level, while LPIPS evaluates the generated images from a deep feature level.

State-of-the-Art Comparisons. The proposed method is compared with several facial image generation models, i.e., CycleGAN [7], StarGAN [9], Pix2pix [5], GPGAN [38], PG2 [10], Pix2pixSC [36], and CocosNet [37]. Note that CycleGAN and StarGAN are unsupervised generation methods, while the others are supervised generation models. The comparisons with CycleGAN and StarGAN are just to see how big the gap between supervised and unsupervised methods is for this task. The results are shown in Table II. We observe that the proposed C2GAN achieves the best results on all four evaluation metrics, validating the effectiveness of our method.

Qualitative comparison results compared with the leading methods are shown in Fig. 4. Clearly, GPGAN performs worse among all comparison models. Pix2pix can generate the correct expression, but the faces are distorted. StarGAN can generate sharper faces, but the details of these generated faces are not convincing. For instance, the mouths in StarGAN are blurred or not correct. Moreover, the results of PG2 tend to be blurry. Compared with the existing leading methods, the re-

TABLE III: Quantitative comparison of hand gesture-to-gesture translation on the Senz3D dataset. For all metrics except FRD, higher is better.

Method	PSNR \uparrow	AMT \uparrow	FRD \downarrow
PG2 [10]	26.5138	2.8	3.0933
SAMG [27]	26.9545	2.3	3.1006
DPIG [39]	26.9451	6.9	3.0846
Pix2pixSC [36]	27.0569	7.2	3.0814
CocosNet [37]	27.1532	7.9	3.0741
PoseGAN [11]	27.3014	8.6	3.0467
GestureGAN [12]	27.9749	22.6	2.9836
C2GAN (Ours)	27.2531	12.7	3.0573

sults generated by the proposed C2GAN are smoother, sharper and contain more details. We also show some generated facial landmarks in Fig. 5. We see that the proposed method not only produces realistic images but also generates reasonable facial landmarks. This is not provided by any existing facial expression generation works.

C. Hand Gesture-to-Gesture Translation

Datasets. We follow GestureGAN [12] and adopt the Creative Senz3D dataset [44] for the hand gesture-to-gesture translation task. This dataset contains 11 different hand gestures performed by four people, each performing gesture is repeated 30 times, resulting in 4 subjects \times 11 gestures \times 30 times = 1320 images in total. We follow [12] and select 12,800 and 135,504 pairs as testing and training data, respectively.

Evaluation Metrics. We follow [12] and adopt Peak Signal-to-Noise Ratio (PSNR) and FRD [12] as evaluation metrics. PSNR measures the similarity between the real image and the generated image at a pixel level. FRD measures the distance between the real image and the fake image from a deep feature level. Moreover, we follow [12] and conduct a user study to evaluate the generated image by different models.

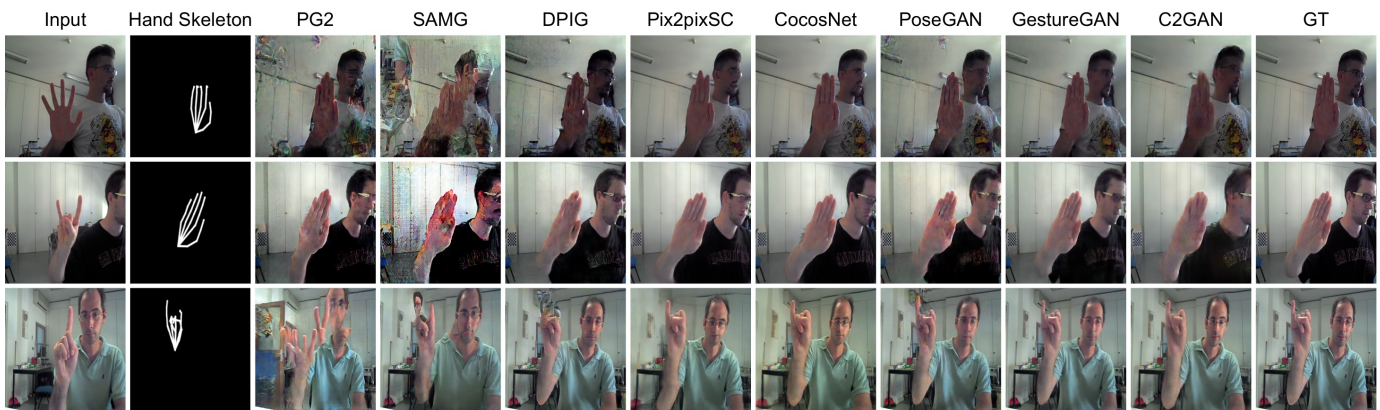


Fig. 6: Qualitative comparison of hand gesture-to-gesture translation on the Senz3d dataset. From left to right: Input, Hand Skeleton, PG2 [10], SAMG [27], DFIG [39], Pix2pixSC [36], CocosNet [37], PoseGAN [11], GestureGAN [12], C2GAN (Ours), and Ground Truth (GT).



Fig. 7: Visualization of hand skeleton generation on the hand gesture-to-gesture translation task.

State-of-the-Art Comparisons. We adopt the most related several works, i.e., PG2 [10], DFIG [39], PoseGAN [11], GestureGAN [12], SAMG [27], Pix2pixSC [36], and CocosNet [37], as our baselines for the facial expression generation task. Comparison results are shown in Table III. The proposed method achieves very competitive results compared with the leading methods. Specifically, the proposed C2GAN achieves significantly better results than PG2, SAMG, DFIG, Pix2pixSC, and CocosNet on all metrics. Moreover, PoseGAN obtains slightly better results than our C2GAN on both PSNR and FRD metrics, however, the proposed C2GAN achieves better AMT than PoseGAN. Moreover, the proposed C2GAN achieves visually better results than PoseGAN, as shown in Fig. 6. Lastly, GestureGAN achieves better results than C2GAN on all metrics. The reason is that GestureGAN is carefully tailored and designed for the specific hand gesture-to-gesture translation task, meaning that GestureGAN is fine-tuned to this task with the network structure, loss objective, and hyper-parameter selection. However, the proposed C2GAN is a novel and unified GAN model, which can be used to handle all kinds of settings of guided image-to-

image translation without modifying the network structure, the loss objective, and hyper-parameters. Furthermore, our C2GAN can generate both images and guidances, which is not considered in GestureGAN.

Qualitative comparison results compared with PG2, DFIG, PoseGAN, GestureGAN, SAMG are shown in Fig. 6. The proposed method generates much better images than PG2, DFIG, SAMG, and PoseGAN. Moreover, our results are very close to those generated by GestureGAN. Our C2GAN is a joint learning framework and it is not only able to generate the target images but is also able to produce the hand skeleton of the input image, which will benefit other computer vision tasks such as hand pose estimation. The results of the generated hand skeletons are shown in Fig. 7. The generated hand skeleton L'_x is very similar to the real hand skeleton L_x , which verifies the effectiveness of the guidance generator G_g and our joint learning strategy.

D. Cross-View Image Translation

Datasets. We follow [13] and adopt the Dayton dataset [45] to evaluate the cross-view image translation task. This



Fig. 8: Qualitative comparison of cross-view image translation on the Dayton dataset. From left to right: Input, Pix2pix [5], X-SO [15], X-Fork [13], X-Seq [13], Pix2pix++ [5], X-Fork++ [13], X-Seq++ [13], SelectionGAN [16], C2GAN (Ours), and Ground Truth (GT).

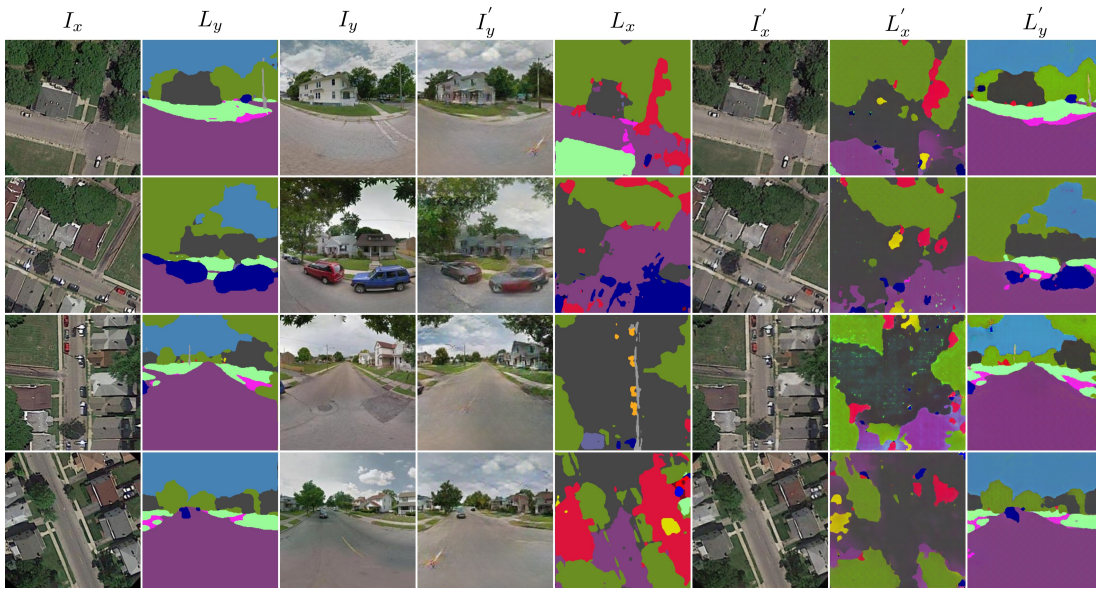


Fig. 9: Visualization of segmentation map generation on the cross-view image translation task.

TABLE IV: Quantitative comparison of cross-view image translation on the Dayton dataset in a2g direction. For all metrics except KL, higher is better.

Method	Accuracy (%)				Inception Score			KL ↓
	Top-1 ↑		Top-5 ↑		All ↑	Top-1 ↑	Top-5 ↑	
Pix2pix [5]	6.80	9.15	23.55	27.00	2.8515	1.9342	2.9083	38.26 ± 1.88
X-SO [15]	27.56	41.15	57.96	73.20	2.9459	2.0963	2.9980	7.20 ± 1.37
X-Fork [13]	30.00	48.68	61.57	78.84	3.0720	2.2402	3.0932	6.00 ± 1.28
X-Seq [13]	30.16	49.85	62.59	80.70	2.7384	2.1304	2.7674	5.93 ± 1.32
Pix2pix++ [5]	32.06	54.70	63.19	81.01	3.1709	2.1200	3.2001	5.49 ± 1.25
X-Fork++ [13]	34.67	59.14	66.37	84.70	3.0737	2.1508	3.0893	4.59 ± 1.16
X-Seq++ [13]	31.58	51.67	65.21	82.48	3.1703	2.2185	3.2444	4.94 ± 1.18
SelectionGAN [16]	42.11	68.12	77.74	92.89	3.0613	2.2707	3.1336	2.74 ± 0.86
C2GAN (Ours)	45.80	75.28	76.03	90.67	2.9603	2.1225	2.9435	2.70 ± 1.02

dataset contains 76,048 images and a training/testing split of 55,000/21,048. The original size of the image is 354×354 resolution. The images are resized to 256×256 .

Evaluation Metrics. Following [13], Inception Score (IS), top-k prediction accuracy, and KL score are employed for the quantitative analysis. These three metrics evaluate the generated images from a high-level feature space.

State-of-the-Art Comparisons. Several leading cross-view image translation methods are adopted as our baselines, i.e.,

Pix2pix [5], X-SO [15], X-Fork [13] and X-Seq [13]. These methods aim to generate images based on a given image. To further evaluate the proposed C2GAN, we introduce four strong baselines, i.e., Pix2pix++ [5], X-Fork++ [13], X-Seq++ [13] and SelectionGAN [16]. These four models aim to generate images based on a given image and several novel segmentation maps. Note that we implement Pix2pix++, X-Fork++ and X-Seq++ using their released public code.

Comparison results are shown in Table IV. The proposed

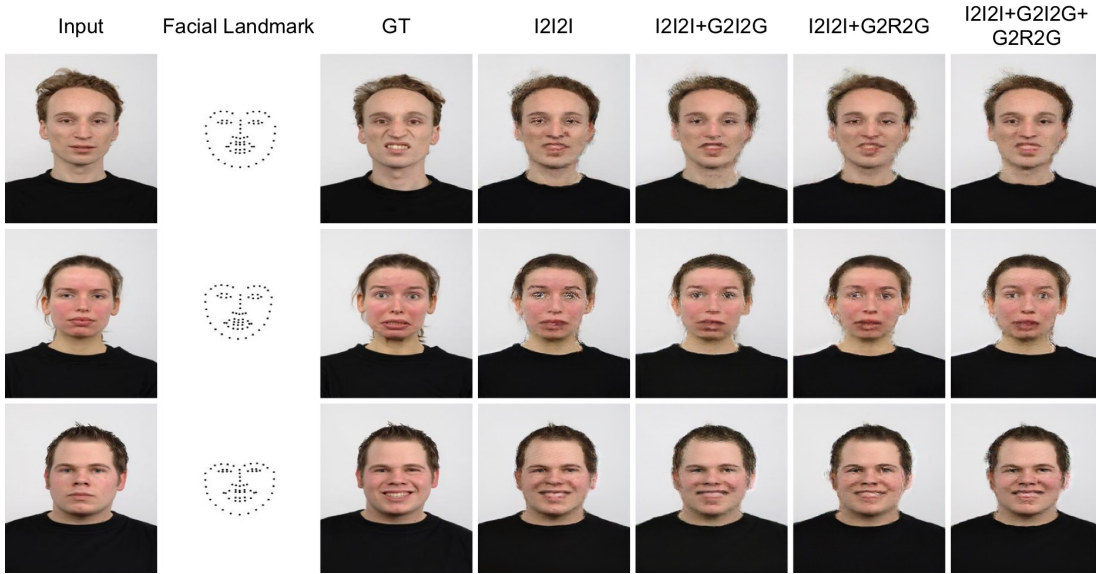


Fig. 10: Influence of individual generation cycle on the Radbound Faces dataset.

TABLE V: Quantitative comparison of ablation studies on the Radbound Faces dataset. For all metrics, higher is better.

Baseline	AMT \uparrow	PSNR \uparrow	SSIM \uparrow
C2GAN w/ I2I2I	25.3	21.2030	0.8449
C2GAN w/ I2I2I + G2I2G	28.2	20.8708	0.8419
C2GAN w/ I2I2I + G2R2G	28.7	21.0156	0.8437
C2GAN w/ I2I2I + G2I2G + G2R2G	30.8	21.6262	0.8540
C2GAN w/ Single-Modal D	26.4	21.2794	0.8426
C2GAN w/ Non-Sharing G	32.9	21.6353	0.8611

C2GAN achieves the best results on several metrics such as KL and Top-1 Accuracy. For other metrics, the proposed method still achieves very competitive results, which validates the effectiveness of the proposed C2GAN.

Several qualitative comparison results are also provided in Fig. 8. Our C2GAN generates much better realistic images than other baselines. Moreover, we show the generated segmentation maps by our method in Fig. 9, the proposed C2GAN can generate reasonable segmentation maps, which we believe our method can be used to improve the performance of semantic segmentation tasks.

E. Ablation Study

We conduct extensive experiments on the Radbound Faces dataset to verify the effectiveness of each component of the proposed C2GAN. All experiments are trained with 50 epochs and Table V shows the quantitative comparison results.

Influence of Individual Generation Cycle. To evaluate the influence of individual generation cycle, we test with four different combinations of the cycles, i.e., ‘I2I2I’, ‘I2I2I+G2I2G’, ‘I2I2I+G2R2G’, and ‘I2I2I+G2I2G+G2R2G’. All four combinations use the same training strategies and hyper-parameters. Comparison results are shown in Table V. Clearly, ‘I2I2I’, ‘G2I2G’, and ‘G2R2G’ are all critical to the final result and the removal of one of them degrades the generation performance, demonstrating that by using cross-modal data in a joint framework and by making the cycles constraint on each other improve the final generation performance.

TABLE VI: Influence of λ_{pixel}^i .

SSIM \uparrow	1	10	100
λ_{pixel}^i	0.8143	0.8540	0.8349

Moreover, ‘I2I2I+G2I2G+G2R2G’ obtains the best performance among all four combination settings. Meanwhile, ‘I2I2I+G2I2G+G2R2G’ achieves remarkably better results than I2I2I on all metrics, demonstrating the effectiveness of constraining both image and guidance cycles facilitating thus a more robust optimization of the whole model. Moreover, some visualization results are provided in Fig. 10 to show the influence of each generation cycle. We can obtain the similar conclusion as the one from Table V, further validating our network design.

Cross-Modal vs. Single-Modal Discriminator. We then evaluate the influence of the proposed cross-modal discriminator (‘C2GAN w/ I2I2I+G2I2G+G2R2G’). Our baseline is the traditional single-modal discriminator (‘C2GAN w/ Single-Modal D ’). The results are listed in Table V. The proposed cross-modal discriminator achieves much better results than the single-modal discriminator on all metrics, meaning that the rich cross-modal information helps to learn a better discriminator and thus facilitates the optimization of the generator.

Parameter Sharing between Generators. The parameter sharing could remarkably reduce parameters of the whole network. We then evaluate how the parameter-sharing strategy would affect the generation results. Specifically, two different baselines are tested: one is ‘C2GAN w/ I2I2I+G2I2G+G2R2G’, which shares the network parameters between the two image generators, and between the two guidance generators, respectively. While ‘C2GAN w/ Non-Sharing G ’ learns four different generators separately. As can be seen in Table V, the non-sharing one achieves slightly better performance than the sharing one. However, the number of parameters of non-sharing one is 217.6M, which is twice as much as that of the sharing one. This means that the parameter-sharing strategy is a good way to balance both

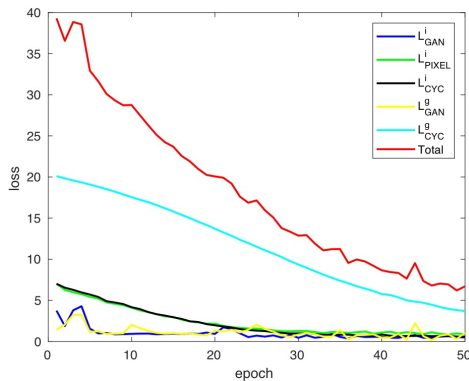


Fig. 11: Model convergence loss in Eq. (13).

image performance and network overhead.

Influence of Hyper-Parameters. For the hyper-parameters in Eq. (13). We first follow Pix2pix [5] and set the hyper-parameters of adversarial losses (i.e., λ_{gan}^i and λ_{gan}^g) to 1. Next, we follow CycleGAN [7] and set the hyper-parameters of cycle-consistency losses (i.e., λ_{cyc}^i and λ_{cyc}^g) to 10. Lastly, we investigate the influence of the hyper-parameter of pixel reconstruction loss (i.e., λ_{pixel}^i) on the performance of our model. Comparison results are shown in Table VI, the proposed method achieves the best performance when $\lambda_{pixel}^i=10$. Therefore, the hyper-parameters λ_{gan}^i , λ_{gan}^g , λ_{cyc}^i , λ_{cyc}^g , and λ_{pixel}^i in Eq. (13) are set to 1, 1, 10, 10, and 10, respectively, in all experiments.

Model Convergence and Training Time. Fig. 11 illustrates the convergence loss of the proposed method in Eq. (13). Note that the proposed model ensures a very fast yet stable convergence. Moreover, our proposed method takes about 10 hours to finish the training of the ablation study on a single TITAN Xp GUP, while CocosNet [37] and PG2 [10] task around 18 and 14 hours, respectively.

V. CONCLUSION

We propose a novel and unified Cycle In Cycle Generative Adversarial Network (C2GAN) for guided image-to-image translation tasks. The proposed C2GAN contains two different types of generators, i.e., image-oriented generator and guidance-oriented generator. Both generators are connected in three generation cycles and can be optimized in an end-to-end fashion. Extensive qualitative and quantitative experimental results on four challenging generative tasks demonstrate that the proposed C2GAN is effective to generate photorealistic images with convincing details.

ACKNOWLEDGMENTS

This work was supported by the EU H2020 AI4Media No. 951911 project, by the Italy-China collaboration project TALENT:2018YFE0118400, and by the PRIN project PREVUE.

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 *NeurIPS*, 2014. 1, 3, 5
- [2] H. Tang, H. Liu, and N. Sebe, "Unified generative adversarial networks for controllable image-to-image translation," *IEEE TIP*, vol. 29, pp. 8916–8929, 2020. 1
- [3] J. Zhang, J. Chen, H. Tang, W. Wang, Y. Yan, E. Sangineto, and N. Sebe, "Dual in-painting model for unsupervised gaze correction and animation in the wild," in *ACM MM*, 2020. 1
- [4] B. Duan, W. Wang, H. Tang, H. Latapie, and Y. Yan, "Cascade attention guided residue learning gan for cross-modal translation," in *ICPR*, 2021. 1
- [5] P. Isola, J. Zhu, T. Zhou, and A. Efros, "Image-to-image translation with conditional adversarial networks," in *CVPR*, 2017. 1, 2, 5, 6, 7, 9, 11
- [6] T. Park, M.-Y. Liu, T.-C. Wang, and J.-Y. Zhu, "Semantic image synthesis with spatially-adaptive normalization," in *CVPR*, 2019. 1, 3
- [7] J. Zhu, T. Park, P. Isola, and A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *ICCV*, 2017. 1, 5, 6, 7, 11
- [8] Z. Yi, H. Zhang, P. T. Gong *et al.*, "Dualgan: Unsupervised dual learning for image-to-image translation," in *ICCV*, 2017. 1, 2
- [9] 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, 3, 6, 7
- [10] L. Ma, X. Jia, Q. Sun, B. Schiele, T. Tuytelaars, and L. Van Gool, "Pose guided person image generation," in *NeurIPS*, 2017. 1, 3, 5, 6, 7, 8, 11
- [11] A. Siarohin, E. Sangineto, S. Lathuiliere, and N. Sebe, "Deformable gans for pose-based human image generation," in *CVPR*, 2018. 1, 3, 4, 5, 6, 7, 8
- [12] 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, 3, 7, 8
- [13] K. Regmi and A. Borji, "Cross-view image synthesis using conditional gans," in *CVPR*, 2018. 1, 3, 8, 9
- [14] L. Song, Z. Lu, R. He, Z. Sun, and T. Tan, "Geometry guided adversarial facial expression synthesis," in *ACM MM*, 2018. 1
- [15] K. Regmi and A. Borji, "Cross-view image synthesis using geometry-guided conditional gans," *Elsevier CVIU*, vol. 187, p. 102788, 2019. 1, 2, 3, 9
- [16] 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. 1, 4, 9
- [17] H. Tang, D. Xu, Y. Yan, P. H. Torr, and N. Sebe, "Local class-specific and global image-level generative adversarial networks for semantic-guided scene generation," in *CVPR*, 2020. 1, 2, 3
- [18] H. Tang, X. Qi, D. Xu, P. H. Torr, and N. Sebe, "Edge guided gans with semantic preserving for semantic image synthesis," *arXiv preprint arXiv:2003.13898*, 2020. 1
- [19] 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
- [20] M. Mirza and S. Osindero, "Conditional generative adversarial nets," *arXiv preprint:1411.1784*, 2014. 2
- [21] H. Tang, S. Bai, and N. Sebe, "Dual attention gans for semantic image synthesis," in *ACM MM*, 2020. 2
- [22] H. Tang, H. Liu, D. Xu, P. H. Torr, and N. Sebe, "Attentiongan: Unpaired image-to-image translation using attention-guided generative adversarial networks," *arXiv preprint arXiv:1911.11897*, 2019. 2
- [23] H. Tang, D. Xu, N. Sebe, and Y. Yan, "Attention-guided generative adversarial networks for unsupervised image-to-image translation," in *IJCNN*, 2019. 2
- [24] H. Tang, D. Xu, W. Wang, Y. Yan, and N. Sebe, "Dual generator generative adversarial networks for multi-domain image-to-image translation," in *ACCV*, 2018. 2
- [25] W. Wang, X. Alameda-Pineda, D. Xu, P. Fua, E. Ricci, and N. Sebe, "Every smile is unique: Landmark-guided diverse smile generation," in *CVPR*, 2018. 3
- [26] H. Tang, S. Bai, P. H. Torr, and N. Sebe, "Bipartite graph reasoning gans for person image generation," in *BMVC*, 2020. 3
- [27] Y. Yan, J. Xu, B. Ni, W. Zhang, and X. Yang, "Skeleton-aided articulated motion generation," in *ACM MM*, 2017. 3, 7, 8
- [28] H. Tang, S. Bai, L. Zhang, P. H. Torr, and N. Sebe, "Xinggan for person image generation," in *ECCV*, 2020. 3
- [29] G. Liu, H. Tang, H. Latapie, and Y. Yan, "Exocentric to egocentric image generation via parallel generative adversarial network," in *ICASSP*, 2020. 3
- [30] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *MICCAI*, 2015. 3, 4, 5
- [31] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *ICLR*, 2015. 5
- [32] B. Amos, B. Ludwiczuk, and M. Satyanarayanan, "Openface: A general-purpose face recognition library with mobile applications," CMU-CS-16-118, CMU School of Computer Science, Tech. Rep., 2016. 5, 6

- [33] Z. Cao, T. Simon, S.-E. Wei, and Y. Sheikh, "Realtime multi-person 2d pose estimation using part affinity fields," in *CVPR*, 2017. 5
- [34] G. Lin, A. Milan, C. Shen, and I. Reid, "Refinenet: Multi-path refinement networks for high-resolution semantic segmentation," in *CVPR*, 2017. 5
- [35] L. Zheng, L. Shen, L. Tian, S. Wang, J. Wang, and Q. Tian, "Scalable person re-identification: A benchmark," in *ICCV*, 2015. 5
- [36] M. Wang, G.-Y. Yang, R. Li, R.-Z. Liang, S.-H. Zhang, P. M. Hall, and S.-M. Hu, "Example-guided style-consistent image synthesis from semantic labeling," in *CVPR*, 2019. 5, 6, 7, 8
- [37] P. Zhang, B. Zhang, D. Chen, L. Yuan, and F. Wen, "Cross-domain correspondence learning for exemplar-based image translation," in *CVPR*, 2020. 5, 6, 7, 8, 11
- [38] X. Di, V. A. Sindagi, and V. M. Patel, "Gp-gan: gender preserving gan for synthesizing faces from landmarks," in *ICPR*, 2018. 6, 7
- [39] L. Ma, Q. Sun, S. Georgoulis, L. Van Gool, B. Schiele, and M. Fritz, "Disentangled person image generation," in *CVPR*, 2018. 6, 7, 8
- [40] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen, "Improved techniques for training gans," in *NeurIPS*, 2016. 5
- [41] 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. 5, 7
- [42] O. Langner, R. Dotsch, G. Bijlstra, D. H. Wigboldus, S. T. Hawk, and A. Van Knippenberg, "Presentation and validation of the radboud faces database," *Taylor & Francis Cognition and emotion*, 2010. 6
- [43] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in *CVPR*, 2018. 7
- [44] A. Memo and P. Zanuttigh, "Head-mounted gesture controlled interface for human-computer interaction," *Springer Multimedia Tools and Applications*, vol. 77, no. 1, pp. 27–53, 2018. 7
- [45] N. N. Vo and J. Hays, "Localizing and orienting street views using overhead imagery," in *ECCV*, 2016. 8



Hao Tang is a Ph.D. candidate in the Department of Information Engineering and Computer Science at the University of Trento. He received the Master degree in computer application technology in 2016 at the School of Electronics and Computer Engineering, Peking University. He was a visiting scholar in the Department of Engineering Science at the University of Oxford, from 2019 to 2020. His research interests are deep learning, machine learning and their applications to computer vision.



Nicu Sebe is Professor in the University of Trento, Italy, where he is leading the research in the areas of multimedia analysis and human behavior understanding. He was the General Co-Chair of the IEEE FG 2008 and ACM Multimedia 2013. He was a program chair of ACM Multimedia 2011 and 2007, ECCV 2016, ICCV 2017 and ICPR 2020. He is a general chair of ACM Multimedia 2022 and a program chair of ECCV 2024. He is a fellow of IAPR.