

LinesToFacePhoto: Face Photo Generation From Lines With Conditional Self-Attention Generative Adversarial Networks

YUHANG LI, University of Science and Technology of China, China

XUEJIN CHEN, University of Science and Technology of China, China

SIYU HU, University of Science and Technology of China, China

ZHENG-JUN ZHA, University of Science and Technology of China, China

SING BING KANG, Microsoft Research, USA

We explore the task of generating photo-realistic face images from lines in this paper. Previous methods based on conditional generative adversarial networks (cGANs) have shown their power in generating visually plausible images when the input conditional image and the output image share well-aligned structures. However, this model is not able to synthesize face images with the whole set of well-defined structures, e.g. eyes, noses, mouths, and etc., especially when the conditional line maps lack of one or several parts of the structure. To address this problem, we propose a conditional self-attention generative adversarial network (CSA-GAN) to capture the global structure of faces and long-range dependencies across the entire image. We demonstrate the effectiveness of the proposed model with experiments on translating face images of CelebA dataset from the corresponding edge maps. We evaluate our model by two kinds of perceptual user studies and MS-SSIM. The experimental results demonstrate that our proposed method generates higher-quality realistic face images from rough lines, and well preserve face structures, comparing to the state-of-the-art methods.

CCS Concepts: • Computing methodologies → Neural networks;

Additional Key Words and Phrases: Self-attention, conditional generative adversarial nets, face, line map, realistic images

ACM Reference Format:

Yuhang Li, Xuejin Chen, Siyu Hu, Zheng-Jun Zha, and Sing Bing Kang. 2018. LinesToFacePhoto: Face Photo Generation From Lines With Conditional Self-Attention Generative Adversarial Networks. 1, 1 (July 2018), 16 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

Realistic image synthesis has been a hot topic in computer vision and computer graphics for years. Traditional methods [10, 11, 16] establish databases of existing images, and generate images by matching and fusing images in the database patch-wisely. With the emergence of deep neural networks (DNN), several promising DNN-based approaches for image synthesis have been proposed. Variational autoencoders (VAEs) [23], which maximize a variational lower bound on the log-likelihood of the training data, have brought some progress in generating visually plausible imagesl.

Authors' addresses: Yuhang Li, University of Science and Technology of China, No. 96, Jinzhai Road, Hefei, Anhui, 230027, China, lyh9001@mail.ustc.edu.cn; Xuejin Chen, University of Science and Technology of China, No. 96, Jinzhai Road, Hefei, Anhui, 230027, China, xjchen99@ustc.edu.cn; Siyu Hu, University of Science and Technology of China, No. 96, Jinzhai Road, Hefei, Anhui, 230027, China, sy891228@mail.ustc.edu.cn; Zheng-Jun Zha, University of Science and Technology of China, No. 96, Jinzhai Road, Hefei, Anhui, 230027, China, zhazj@ustc.edu.cn; Sing Bing Kang, Microsoft Research, One Microsoft Way, Redmond, WA, 98052, USA, SingBing.Kang@microsoft.com.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2018 Association for Computing Machinery.

XXXX-XXXX/2018/7-ART \$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

However, the generated samples suffer from being blurry. The PixelCNN decoders [42] generate convincing images using autoregressive models. However, the pixel-wise sampling procedure make this model computationally expensive.

Among these DNN-based methods, generative adversarial networks (GANs) [13] offer a new and promising mechanism to generate images. GANs take noise vectors as input, and train two networks playing a minmax game to guide the generated samples to be indistinguishable from the real ones. Conditional GANs, which generate image from assigned conditional information instead of noise vector, are conditional versions of GANs. Conditional GANs are trained in a supervised manner and shown to be powerful in modeling the conditional distributions with respect to the assigned conditions. A variety of conditions have been applied to conditional GANs, such as discrete class labels [32], texts [50, 51], and images. As a specific application of CGANs, image-to-image translation has drawn a lot of attention recently. It aims to apply a conditional image in one domain to generate the corresponding target image in another, while preserving shared concepts, objects or scenes in these two images. Since the first image-to-image model (pix2pix) [19] was proposed, there have been many variants of this approach in both supervised and unsupervised manner [28, 48, 53, 54]. These models successfully synthesize realistic textures when fine structures are given in the conditional image, or inferring the global structure in a specific object class. When the geometric structure is partially provided with different level of details, which is exactly the case of line drawings or edge maps, these models fails to complete the missing structure while preserving the given structure details.

The underlying reason is mainly two-fold. First, the existing GANs are mainly built based on convolutional layers. Since the convolution operator has a local receptive field depending on the size of its kernels, a large receptive field is achieved by cooperation of several convolution layers. It is hard for the network optimizer to discover proper parameter values that model the long-range dependencies through several convolutional layers [49]. Secondly, the existing discriminator used in GANs focuses on examining local patches instead of capturing the global information, and therefore fails to guide the generator to synthesize the global structure of the conditional image.

Considering the first reason, we introduce a conditional self-attention mechanism to the generator of image-to-image models to address the problem. Self-attention [7, 43, 45, 49], which computes the response at a position as a weighted sum of the features at all positions, is able to capture the long-range dependencies across different regions of images and feature maps. In order to adapting the conditional setting of image-to-image translation and encouraging the model to leverage the information of the conditional image directly, we propose a conditional self-attention module (CSAM) which enables the higher layers to sense the conditional image. For the second reason, we consider to establish multiple discriminators to capture information of different levels, both patch-wisely and globally. We note that similar idea of multiple discriminators has been raised by [4, 8, 18, 51] who resizes the real/fake samples and applies multiple discriminators to these multi-scale samples.

On the whole, our contributions are summarized as follow:

- (1) We firstly introduce the self-attention mechanism to image-to-image translation and propose a novel conditional self-attention generative adversarial networks for the image-to-image translation task. Unlike convolution-based methods, the proposed model is able to model the long-range dependencies and global structure across images.
- (2) We propose a multi-level discriminator to the image-to-image translation. The proposed discriminator is able to capture the global structure information as well as the local realism.
- (3) We show the effectiveness of the proposed model by experiments. Two kinds of user studies are investigated to show the perceptual evaluation of the results generated by the proposed

method. Quantitative evaluation is conducted by calculating the MS-SSIM of the pix2pix model and the proposed model.

The rest of this article is organized as follow. Related work is presented in Section 2. The method we proposed is introduced in Section 3. We demonstrate the effectiveness of method by a series of experiments in Section 4. Section 5 summarizes our conclusions and outlines possible future work.

2 RELATED WORK

Our work is based on image-to-image translation frameworks, which are variants of GANs in a conditional setting. In this section, we present related research in GANs, conditional GANs, and image-to-image translation models. We also give a brief review on recently proposed attention models.

2.1 Generative Adversarial Networks (GANs)

Generative adversarial networks (GANs) [13] have obtained a great success in recent years. Based on the minmax game theory, a classical architecture of GANs contains a generator network and a discriminator network. The generator takes a noise vector as input and generates samples indistinguishable from the real ones, while the discriminator, in opposite, attempts to find out whether its input is real or synthesized. The minmax game played by these two networks guides the generated distribution to be similar to the real data distribution. Compared to other DNN frameworks of image generation [23, 42], GANs are able to synthesize images with less blurriness, and provide a more efficient process to generate samples. However, GANs suffer from several problems in the early stage, such as the instability of training and the mode collapse problem. To stabilize the training process of GANs and enable GANs to generate images with high quality and large diversity, many efforts have been made. Deep convolutional GANs (DCGANs) [38] first introduced a convolutional architecture which led to improved visual quality. An approach is proposed to train discriminator in a semi-supervised fashion [40], granting the discriminator's internal representations knowledge of the class structure of (some fraction of) the training data is presented. Energy-based GANs (EBGANs) [52] were proposed as a class of GANs that aim to model the discriminator as an energy function. This variant converges more stably and is not only easy to train but also robust to hyper-parameter variations. Wasserstein distance, which acts as a loss as well as a measure of convergence in training process, is brought to GANs [1, 14] to benefit both the stability and mode coverage. Several other methods [3, 24, 31] also make progress in stabilizing the training and increasing the diversity of the results of GANs. A recently proposed technique [49] introduced self-attention mechanism to unconditional GANs and achieved state-of-the-art results. Inspired by this self-attention mechanism, we introduce a conditional self-attention module to capture long-range dependences between parts of faces when transferring lines to realistic images.

2.2 Conditional Generative Adversarial Networks

Conditional GANs are generalized versions of GANs in a conditional setting. Instead of taking a noise vector as input, conditional GANs generate images based on the assigned conditions and model the conditional distributions of the samples with respect to the conditions. Conditional GANs were firstly introduced [32] as treating the conditional generation problem as the inverse processing of image classification and using discrete labels as condition to generate images. Previous work has explored GANs for generating images based on a wide variety of conditions. [6, 36] took both noise vectors and discrete class label as input and added a classification task to the discriminator in two different architectures in order to generate images with high recognizability. [9] trained convolutional networks to generate images of objects given object style, viewpoint and color.

With the experiments of interpolating viewpoints, they showed that networks learn a meaningful representation of 3D models. [50, 51] generated high-resolution photo-realistic images conditioned on text descriptions in two stages, where GANs sketched the basic shapes and colors in the first stage and added details to the generated images in the second stage. A recently proposed method [34] leverages the conditional information in a novel way, where the discriminator involves an inner product term between the condition vector and the feature vector in a middle level layer. This formulation is based on the observation that the loss function of GANs can be decomposed into two log likelihood ratios.

Our framework utilizes GANs in a conditional setting to generate images from rough lines, and it carries the condition information directly even in the high-level layers.

2.3 Image-to-Image Translation with GANs

Given an image in one domain, image-to-image translation methods generate a corresponding image in another domain. These two images are possible representations of the same scene or object. Image-to-image translation with GANs is a special case of conditional GANs where images are applied to be conditions.

The pix2pix method [19] firstly introduced the concept of image-to-image translation. Pix2pix is trained in a supervised manner, where the training dataset is a set of paired images. To deliver the fine-level structure of the conditional image to the generated image, pix2pix applies skip connections [39] between mirrored layers in the generator to make sure low-level information pass through its encoder-decoder architecture and uses patch discriminators to increase the performance of the generator. However, the convolution-based architecture makes it difficult to discover the long-range dependencies across the images and feature maps. Moreover, the patch-wise discriminator is not able to ensure the global structure information to be well captured by the model. In addition to pix2pix, many image-to-image tasks are trained in a supervised manner. [4, 44] used coarse-to-fine refinement frameworks to synthesis photographic images from semantic label maps. [20] studied to generate images of outdoor scenes from semantic label maps coupled with attributes. [54] presented a framework that is able to model the multi-modal distribution of possible outputs. Image-to-image translation has also been well-studied in an unsupervised setting. [27] studied on unpaired image-to-image translation by training a two-branch GAN. Each branch is composed with a encoder, a generator and a discriminator. With the idea that high-level representation of a pair of corresponding images in two domains should be the same, high-level layers share weights between two branches in encoders, generators and discriminators. CycleGAN [53], DiscoGAN [21] and DualGAN [48] developed similar architectures to translate unpaired images which contain, for each, two generators and two discriminators. These methods learn two mappings in an adversarial training process such that an input image in one domain is mapped to a generated image in another, and then the generated image is mapped to a reconstructed image which is closed to the input image in some measures. These methods shared the same idea that the generated image should contain the content of the input image in order reconstruct the input image from it. In comparison, our work focuses on translating a rough lines map to a realistic face photo, whose key challenge is to learn the long-range dependencies and global structure across different regions in a face image.

2.4 Attention Mechanism

Since the convolution operation has a local receptive field, several layers and large kernel sizes are required to sense the global structure in a large receptive field. However, simply combining several convolutional layers loses the computational and statistical efficiency. Recently, attention mechanisms have been introduced to capture global dependencies [2, 47]. Self-attention [45, 49] has been shown to be powerful in a variety of tasks. [43] applied self-attention to machine translation

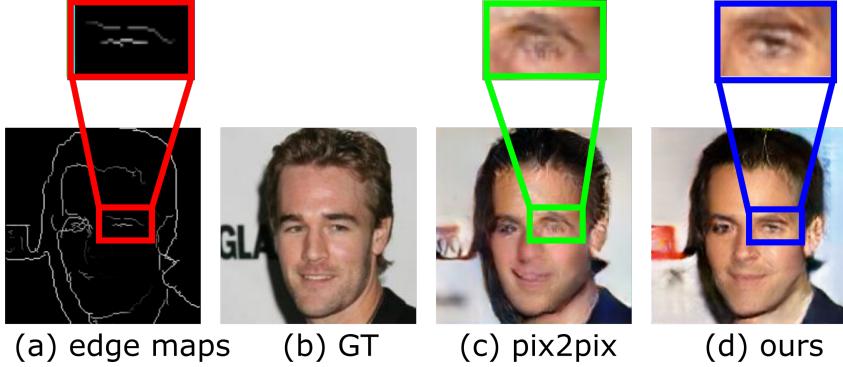


Fig. 1. An example of translating face images (b) from the corresponding edge maps (a). In this example, the edge map does not contain the complete edges of the left eye (on the right hand side of the reader). It is obvious that there should be a left eye in the red square according to the global structural information. The pix2pix model (c) fails to render a recognizable eye in the corresponding location (green square) while the proposed method (d) is able to generate the entire structure even when the conditional edge map lacks several parts of the global structure.

models, and demonstrated the plausible effectiveness of self-attention mechanism. [37] studied on combining the self-attention mechanism with autoregressive models, and proposed an image transformer model in image generation. Inspired by non-local operations in computer vision, [45] utilized self-attention mechanism as a non-local operation to model long-range spatial-temporal dependencies for video processing. [49] introduced self-attention to unconditional GANs and achieved state-of-the-art results in generating natural images from noise vectors. Inspired by previous works, we explore the self-attention mechanisms in the context of image-to-image translation.

3 METHOD

The proposed Conditional Self-attention Generative Adversarial Networks (CSAGANs), which translate images from one domain to another and are able to capture long-range dependencies and reserve the global structures across image. We first review the pix2pix model as our baseline (Sec.3.1). And then we introduce the Conditional Self-attention Module (SCAM) (Sec. 3.2). Finally, we describe the idea of multiple level patch discriminator (Sec.3.3 and the architecture we proposed (Sec. 3.4).

3.1 The Pix2pix Model

Since our model is based on the pix2pix model [19], we review this model in this sub-section. The pix2pix model is an image-to-image translation framework based on conditional GANs, which trains a generator network G and a discriminator network D alternatively. The generator G takes a conditional image as input and outputs the corresponding target image, while the discriminator D distinguishes real images from the synthesized ones. To train these two networks in a supervise manner, a set of corresponding image pairs $\{(\mathbf{x}_i, \mathbf{y}_i)\}$ is required as training set, where \mathbf{x}_i is a conditional image and \mathbf{y}_i is the corresponding target image. These two networks play a minmax game to guide the generator to model the conditional distribution of real images given the conditional images. The objective is given by:

$$\min_G \max_D \mathcal{L}_{adv}(G, D) + \lambda \mathcal{L}_{L1}(G), \quad (1)$$

where G aims to minimize this objective while D tries to maximize it inversely. The adversarial loss function is generally given by

$$\mathcal{L}_{adv}(G, D) = E_{(\mathbf{x}, \mathbf{y}) \sim p_{data}(\mathbf{x}, \mathbf{y})} [\log D(\mathbf{x}, \mathbf{y})] + E_{\mathbf{x} \sim p_{data}(\mathbf{x})} [\log(1 - D(\mathbf{x}, G(\mathbf{x})))], \quad (2)$$

and the L_1 loss is given by

$$\mathcal{L}_{L1}(G) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim p_{data}(\mathbf{x}, \mathbf{y})} [\|\mathbf{y} - G(\mathbf{x})\|_1] \quad (3)$$

The generator of the pix2pix model is a fully-convolution-based U-Net [39]. The input of the generator is only applied to the first layer. The discriminator of the pix2pix model is a patch-wise discriminator, which examines only a patch of its input image and uses the average of outputs from all patches of the input image as the ultimate output. The size of each patch is set to 70×70 . The conditional image is concatenated channel-wisely to the synthesized image or real image as the input of the discriminator.

However, in the task of translating a face image from the corresponding edge map, the pix2pix model has troubles to generate realistic face images in some cases with structural constrains. Since faces have well-defined structural parts, e.g. noses, mouths, eyes and etc., the synthesized face images should contain the whole set of these structural part to be realistic, even when the conditional edge maps lack of edges on the supposed locations of these parts. The pix2pix fails to generate realistic structural part in this circumstance. An example is displayed in Figure 1. In this example, an edge map of a face, shown on the left of the figure, only captures a part of edges of the left eye (in the green square) rather than the whole set of edges of the entire left eye. The face image generated by the pix2pix model on the condition of this edge map is shown on the right in the figure. According to the global structural information, there should be a left eye in the red square obviously. However, we can observe that the pix2pix model fails to render a recognizable left eye in the synthesized face image.

This phenomenon might be caused by two reasons. 1) The pix2pix model is a convolution-based model which relies on convolutional operations to model the dependencies across different regions of images and feature maps. Convolutional operations have local receptive fields depending on the kernel sizes and are not able to balance between ability to model long-range dependencies and efficiency of computation and statistics in some cases [49]. 2) The discriminator used by the pix2pix model is patch-wised, based on the assumption that pixels separated by a distance more than a patch diameter are independent to each other. This assumption is true in some cases like texture generation and style transfer, and has been applied in previous work [12, 26]. However, this assumption fails in the case with global structural constrains. Therefore the patch-wise discriminator fails to grasp the global structure information and is not able to guide the generator to be aware of the structure of the faces.

To address the problem caused by the first reason, we introduce self-attention to the generator of image-to-image models to address the problem. Self-attention [7, 43, 45, 49], which computes the response at a position as a weighted sum of the features at all positions, is able to capture the long-range dependencies across different regions of images and feature maps. In order to adapting the conditional setting of image-to-image translation and encouraging the model to leverage the information of the conditional image directly, we propose a conditional self-attention module (CSAM), which enables the higher layers to sense the conditional image, as a general module of networks. For the second reason, we consider to establish a multi-level discriminator to capture the information of its input image both patch-wisely and globally. We note that similar ideas of multiple discriminators have been raised by [4, 8, 18, 51] with different architectures. We describe CSAM and the multi-level discriminator in next sections.

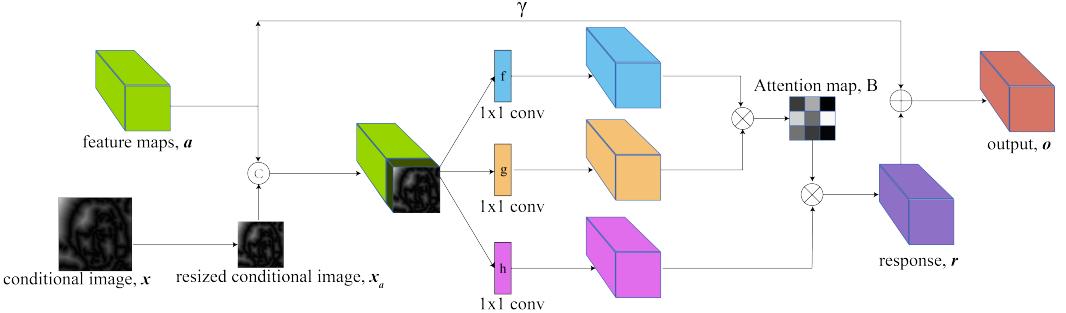


Fig. 2. The proposed CSAM. Given the conditional image and feature maps from the previous layer, the output feature maps are calculated in a self-attention manner. This module is designed to be added after any assigned layers.

3.2 Conditional Self-Attention Module (CSAM)

We improve the pix2pix model by utilizing self-attention mechanism to capture the long-range dependencies of images and feature maps. A recently proposed method [49] has introduced self-attention to unconditional GANs and achieved state-of-the-art results. Inspired by this method, we propose a conditional self-attention module (CSAM) which is suitable for image-to-image translation framework and able to leverage the conditional information directly. This module is designed as a general module of conditional frameworks and can be added after any existing modules. We will provide details of our architecture in Subsection 3.4. The formulation of CSAM is described below.

Given the conditional image $x \in \mathbb{R}^{3 \times N_x}$ and feature maps from the previous layer $a \in \mathbb{R}^{C \times N_a}$, we first resize the conditional image x to match the size of a and get $x_a \in \mathbb{R}^{3 \times N_a}$. Here $N_x = H_x \times W_x$, where H_x, W_x are the height and width of the conditional image x . N_a is defined similarly for a . Then we concatenate the resized conditional image x_a to the feature maps a to get $[a, x_a]$ as conditioned features, where $[\cdot, \cdot]$ is the concatenation operation. This allows the information of conditional image to convey to every attention module and guide the network to form the attention directly based on the conditional image.

In order to calculate the attention, we map the conditional features $[a, x_a]$ to two feature spaces by:

$$f([a, x_a]) = W_f [a, x_a], \quad (4)$$

$$g([a, x_a]) = W_g [a, x_a], \quad (5)$$

where $W_f, W_g \in \mathbb{R}^{\hat{C} \times (C+3)}$ are trainable weights and are implemented by 1×1 convolutions. Here, we use $\hat{C} = C/8$ in our experiments following the setting of previous work [49]. Let $\mathbf{B} \in \mathcal{R}^{N_a \times N_a}$ be the attention map. Every element in \mathbf{B} is denoted as $b_{j,i}$ which indicates the extent to which the model attends to the i^{th} location when synthesizing the j^{th} region and is calculated by

$$b_{j,i} = \frac{\exp(s_{ij})}{\sum_{i=1}^{N_a} \exp(s_{ij})} \quad (6)$$

where $s_{ij} = f([a, x_a])^T g([a, x_a])$. Next, we use $b_{j,i}$ as the attention weights and compute the response $\mathbf{r} = (r_1, r_2, \dots, r_{N_a}) \in \mathbb{R}^{C \times N_a}$ at every position as a weighted sum of the features at all

positions, where

$$\mathbf{r}_j = \sum_{i=1}^{N_a} b_{j,i} h([\mathbf{a}, \mathbf{x}_a]), \quad (7)$$

where $h([\mathbf{a}, \mathbf{x}_a]) = \mathbf{W}_h [\mathbf{a}, \mathbf{x}_a]$ and $\mathbf{W}_h \in \mathbb{R}^{(C+3) \times (C+3)}$. As suggested in [49], we further multiply the response of the attention layer by a scale parameter γ and add back to the input feature maps. The final output is calculated by

$$\mathbf{o}_i = \gamma \mathbf{r}_i + \mathbf{a}_i, \quad (8)$$

where γ is trainable value and is set to 0 at the beginning of the training process. This is because at the early stage of training process, the networks are able to learn the local dependencies, and then learn the long-range dependencies by assign more weight to the non-local evidence progressively.

3.3 Multi-Level Patch Discriminator

The discriminator of the pix2pix model is patch-wised, which distinguishes the real/synthesized images patch by patch convolutionally with in a local receptive field much smaller than the size of the input images. The average value of all responses is provided as the ultimate output of d . This is based on the assumption of independence between pixels separated by more than a patch diameter. However, since the structural constrain is global information across the entire image, the patch-wise discriminator lacks ability to capture this global information. We add another global discriminator D_g with a receptive field as large as the entire image to capture the global structure information. The patch discriminator D_p and the global discriminator D_g share weights in first few layers since the lower features of these discriminators should be the same, as shown in Figure 3. The objective of the minmax game therefore is modified from Equation 1 to

$$\min_G \max_{D_g, D_p} \mathcal{L}_{adv}(G; D_g, D_p) + \lambda \mathcal{L}_{L1}(G), \quad (9)$$

where the adversarial loss is given by

$$\begin{aligned} \mathcal{L}_{adv}(G; D_g, D_p) &= E_{(\mathbf{x}, \mathbf{y}) \sim p_{data}(\mathbf{x}, \mathbf{y})} [\log D_g(\mathbf{x}, \mathbf{y}) + \log D_p(\mathbf{x}, \mathbf{y})] \\ &\quad + E_{\mathbf{x} \sim p_{data}(\mathbf{x})} [\log(1 - D_g(\mathbf{x}, G(\mathbf{x}))) + \log(1 - D_p(\mathbf{x}, G(\mathbf{x})))], \end{aligned} \quad (10)$$

and L_1 loss is still the same as Equation 3.

3.4 Architecture

Our architecture, shown in Figure 3, is based on that of the pix2pix method which uses a convolution-based U-Net [39] as its generator and a patch-wise discriminator. We add a proposed CSAM after every convolutional layer of the generator except the first and last ones. The conditional image is resized to specific size and concatenate to the previous feature maps as the input of every CSAM. CSAMs are able to access the information of the conditional image directly and model the long-range dependencies across images and feature maps. Also, we switch the patch-wise discriminator into the proposed multiple level patch discriminator to enable the discriminator network to capture both global and local information and therefore guide the generator to generate images with more structural layout. More details of the architecture are discussed below.

Noise vector. Some past conditional GANs add a noise vector to the generator as input to avoid it producing a deterministic output. However, the pix2pix model has shown that the noise vector is just ignored by the generator network and hardly change the output samples. We observe the same phenomenon in our experiments and do not apply the noise vector in our model.

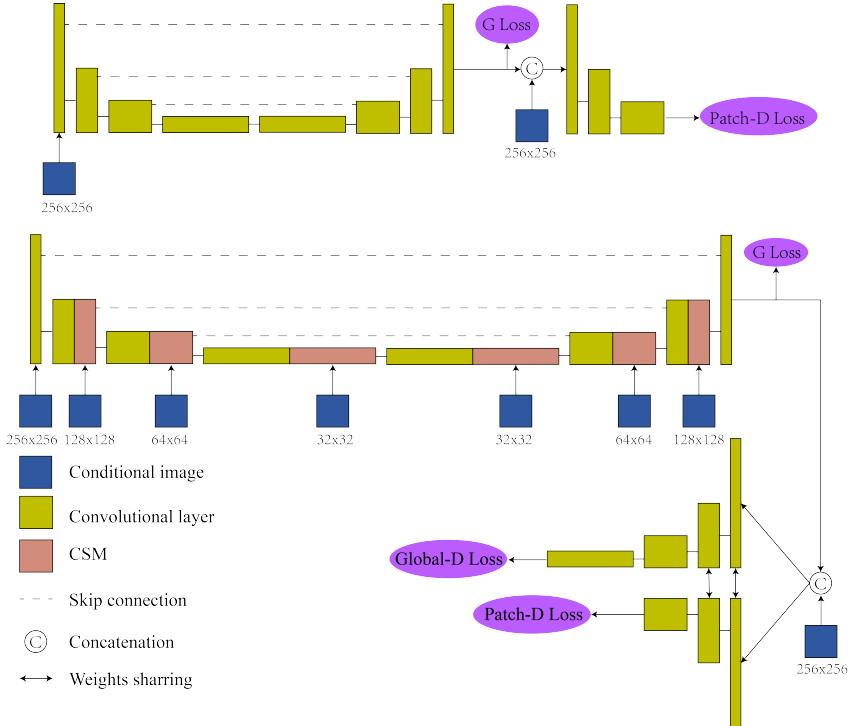


Fig. 3. The pix2pix framework (the upper architecture) and the proposed framework (the lower architecture). Compared to the pix2pix model, in the generator of CSAGAN CSAMs are added after every convolutional layers except the first and the last ones. The discriminator is changed into a multi-level discriminator. The conditional image is resized and connect to every CSAM in the generator.

Spectral Normalization. Spectral normalization [33] is a recently proposed normalization technique, which restricts the spectral norm of each layer of the discriminator to constrain its Lipschitz constant. Spectral normalization is computationally efficient and require no extra hyper-parameter. It has shown that spectral normalization also benefit the training of generator by avoiding unusual gradients. We add spectral normalization to the discriminator and CSAMs in the generator.

4 EXPERIMENT

We propose the CSGANs framework, which translate images from one domain to another, being able to capture long-range dependencies and reserve the global structures. To demonstrate the effectiveness of our framework, we have performed several experiments. In this section, we discuss

4.1 Implementation Details

In order to comparing with the pix2pix model, we basically follow its implementation details. We use minibatch SGD and Adam [22] optimizer with learning rate $lr = 0.002$ and momentum parameters $\beta_1 = 0.5, \beta_2 = 0.999$. We update one step for either of G and D alternatively. Batch normalization is used in convolutional layers of the generator. Batch size is set to 8.

4.2 Dataset

We evaluate our method with the task of translating edge maps to natural images, e.g. the target images are face images while the conditional images are the corresponding edge maps. The face images of the dataset we used are face images in CelebA dataset [29], a large-scale face attributes dataset with more than 200K celebrity images. Faces have well-defined structure of eyes, noses, mouths, and etc., and therefore the artifacts are visually sensitive for observers. This makes face images suitable for evaluating the proposed method. We utilize the cropped and aligned version of dataset with the size of every images being 218×178 . In order to meet the original setting of the pix2pix method, we center-crop the images and resize the image to 128×128 in both experiments of the pix2pix model and the proposed model. The face attributes are attached in the dataset but not included in our experiments.

The edge maps we use are generated in a pipeline similar to that used in pix2pix paper. Specifically, the edges are firstly extracted using a deep edge detector named holistically-nested edge detector (HED) [46]. We keep the values of each edge pixels calculated by HED in the edge maps. Each of these values is supposed to indicate the probability of being edge in the positions of pixel. And then several steps of post-processing are conducted to obtain simpler and clearer edge maps with fewer edge fragments, including thinning, short edge removal, and erosion. In addition, since the edge maps are very sparse, we add one more step to the process to decrease the sparsity of edge maps. We calculate an unsigned euclidean distance field for each edge map to obtain a dense representation. We note that similar idea of distance filed representations can be found in some recent works [5, 15, 35]. In Section 4.3, we will prove the advantages of distance fields by experiments. 56

4.3 Evaluation Metrics

The evaluation of generative models is an open and complicated task, because a model with good performance with respect to one criterion need not imply good performances with respect to the other criteria [30, 41]. Traditional metrics, such as pixel-wise mean-squared error do not present the joint statistics of the synthesized samples and therefore is not able to evaluate the performance to a conditional generated model. Inception Score (IS) [40] is a widely-used criterion. However, IS has been pointed out to have serious limitations that it focuses more on the recognizability of the generated images rather than realism of details or intra-class diversity [41]. Moreover, IS is an evaluation metric for class-aware task which is not suitable for our experiments.

Since the goal of image-to-image translation is to generate from the conditional image an corresponding image visually plausible to human, we mainly compare the results between different models by perceptual user studies. Several related works have proposed similar perceptual experiments [4, 8, 25, 40, 44]. Following the similar procedure as described in [4], we conduct two different kinds of experiments: unlimited time user study and limited time user study. In addition, we use another popular criterion, multi-scale structure similar index [17], to prove the effectiveness of proposed method quantitatively. More details are explained below.

Unlimited Time User Study. We utilize perceptual user study experiments to compare the generated samples between different models. In every trial, we randomly select a conditional image from the testing dataset and generate two synthesized images from pix2pix and our model that are going to be compared with each other. These three images are displayed side by side, and the user is asked to pick one from the two synthesized images within unlimited time based on "which is more realistic and matches the conditional image better". The options offered to users are two of the synthesized images. No feedback is provided after every trial to avoid misguiding the user's perceptual judgment and preference. 00 users participate this experiments, and 00 trials are provided to every user.

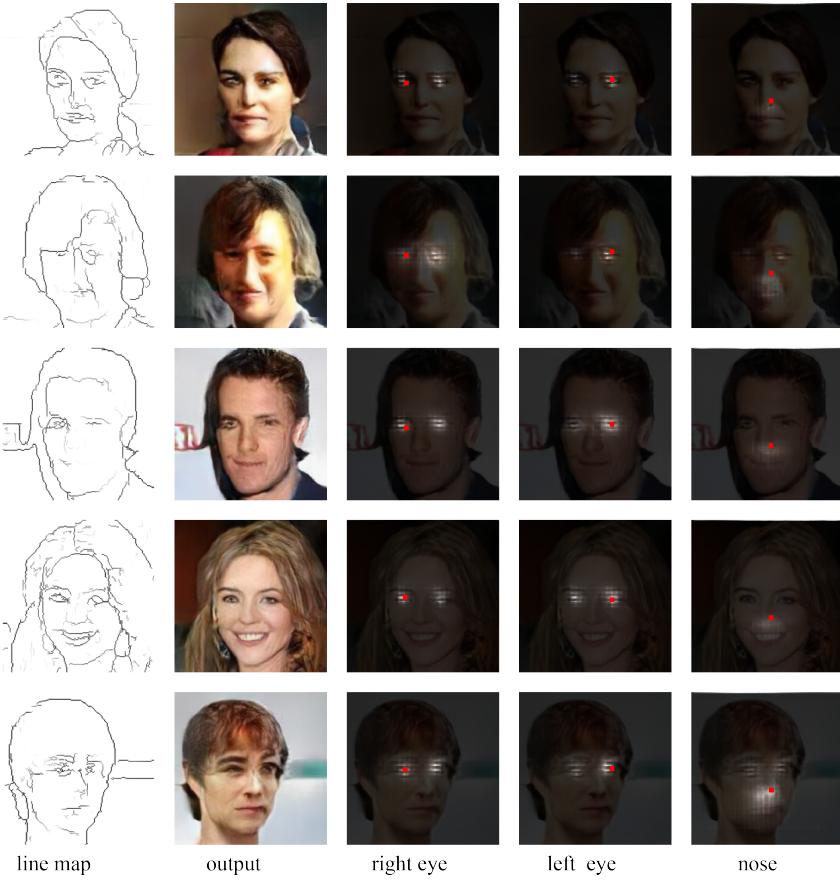


Fig. 4. The attention maps are shown with the conditional edge maps and the generated images. The attention maps are drawn from the last CSAM in the generator, since it is closest to the generated images. Three locations of the nose and two eyes are marked in red while the attention maps with respect to these locations are shown. The larger values in the attention maps are brighter in the figure. We can observe that the long-range dependencies are captured by the CSAM.

Limited Time User Study. For this task, we evaluate how quickly the users can perceive the differences between images. In every comparison, we select three images corresponding to one randomly drawn edge maps (two generated by pix2pix and our model, and the ground truth). Similarly, two of these three images are displayed to the user with the edge map side by side for a short period of time. The user is asked to pick one of two displayed face images also based on "which is more realistic and matches the conditional image better". The duration is randomly selected between 1/8 seconds and 8 seconds.

Multi-scale Structural Similarity (MS-SSIM). Multi-scale structural similarity (MS-SSIM) is a popular similarity metric between two sets of images. MS-SSIM is based on the assumption that the human visual system is highly adapted for extracting structural information from the scene, and therefore a measure of structural similarity can provide a good approximation to perceived image quality. We calculate the MS-SSIM between the target images and the images generated by edges maps in the testing set.

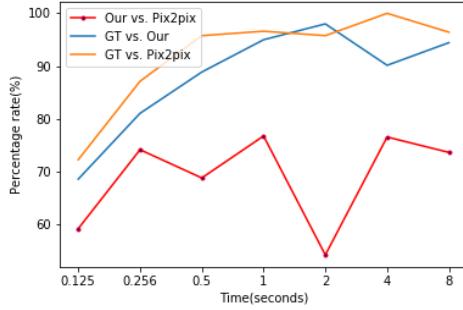


Fig. 5. The results of limited time user studies.

4.4 Comparison with pix2pix

In this section, experiments are conducted to compare the images generated by the pix2pix model and the proposed model. The comparisons are described below.

Two kinds of user studies are performed. The unlimited time user study is designed to evaluate the perceptual quality of the generated image, the results of which are shown in Table 1. We can observe that given unlimited time, users are able to discover the visual differences between these two methods. The advantages of our method are obvious in this user study. On the other hand, the limited time user studies are designed to evaluate how quickly the users can perceive the differences between images. Figure 5 shows the results. When images are shown in a very short time (1/8 seconds), users are not able to sense the differences among these two methods and the ground truth. With the increase of the time, more differences are perceived by users. The advantages of the proposed method is also obvious in this study.

Table 1. Unlimited time user study.

	pix2pix [19]	Ours
user preference	0.718%	0.282%

Table 2. Evaluation metrics

Images	MS-SSIM
pix2pix [19]	0.2538
Ours w/o distance fields	0.2878
Ours w/o global Discriminator	0.3196
Ours w/o conditional Connection	0.2778
Ours, full model	0.3629

4.5 Ablation Study

We examine the importance of every part of our model by MS-SSIM, shown in Table 2. Experiments are conducted by removing specific parts from the full model and then generating images without this part. Specifically, we remove 1) the calculation of the unsigned distance fields of edge maps

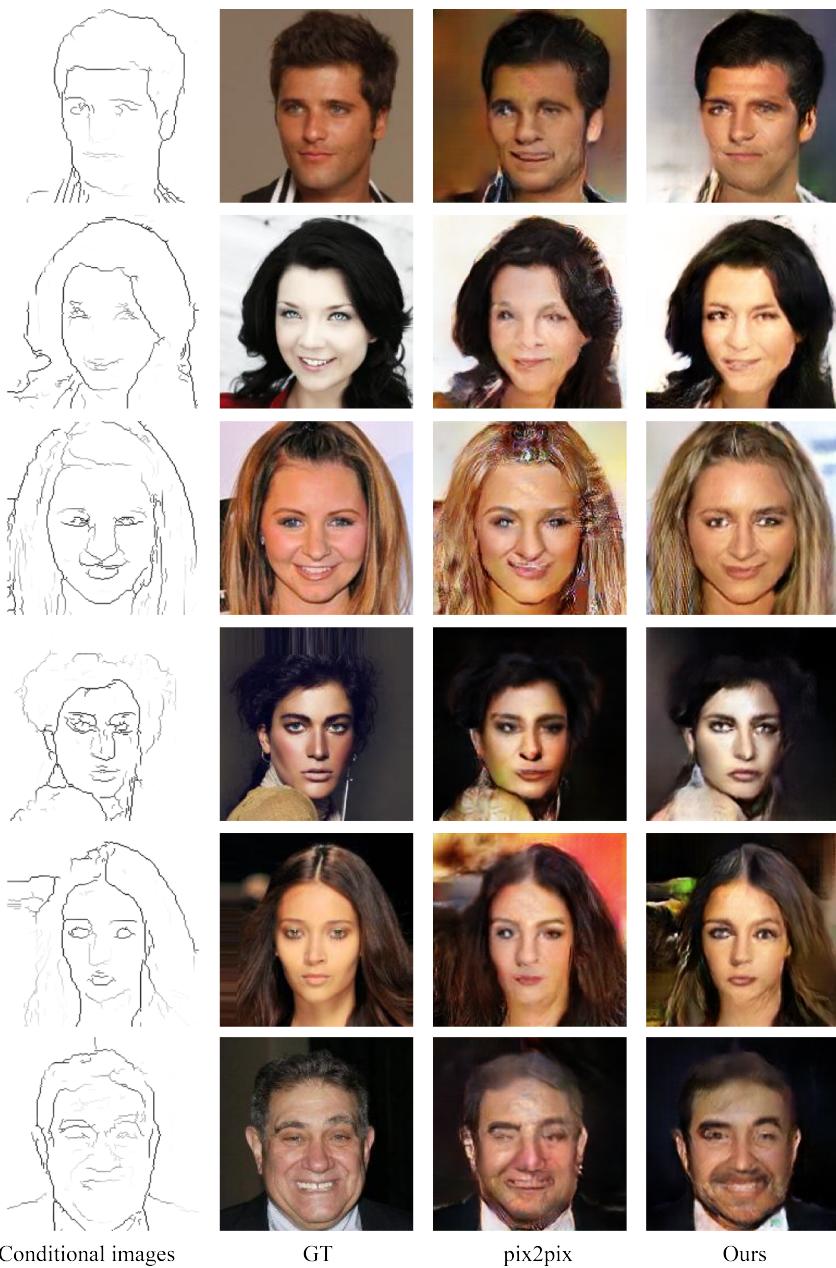


Fig. 6. The face images generated by the pix2pix and the proposed method on the condition of edge maps. Ground truth (GT) images that we use to obtain the edge maps are shown in the second column.

before input to the generator and discriminator(-distance fields), 2) the spectral normalization in the CSAMs of generator and convolutional layers of discriminator (-spectral normalization), 3) the global discriminator in the multi-level discriminator (-global discriminator), and 4) the connection in each CSAM of resized images (-conditional connection).

4.6 Self-Attention

We visualize the attention maps to find out how pixels of all locations in the images and feature maps are learned attend to one specific pixel. We select the attention maps of the last CSAM in the generator since this layer is closest to the generated image. Figure 4 shows some examples. In this figure, three locations of the nose and two eyes are marked in red while the attention maps with respect to these locations are shown. The larger values in the attention maps are brighter in the figure. We can observe that the long-range dependencies are captured by the CSAM. For example, to generate the pixels in one eye, the regions of both eyes are assigned high attentions. In another words, the information of generating a specific pixel comes from not only its local area but also related regions far away from this pixel.

5 CONCLUSION

In this work, we introduced self-attention mechanism to the conditional GANs and proposed conditional self-attention GANs (CSAGANs). This framework is able to capture the long-range dependencies across different regions of images and global structural information. With the help of the proposed conditional self-attention module, the proposed model is able to leverage the information of conditional images directly. With a series of experiments We evaluate the effectiveness of the proposed model and advantages compared to the pix2pix model.

REFERENCES

- [1] Martin Arjovsky, Soumith Chintala, and Léon Bottou. 2017. Wasserstein gan. *arXiv preprint arXiv:1701.07875* (2017).
- [2] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473* (2014).
- [3] David Berthelot, Thomas Schumm, and Luke Metz. 2017. BEGAN: boundary equilibrium generative adversarial networks. *arXiv preprint arXiv:1703.10717* (2017).
- [4] Qifeng Chen and Vladlen Koltun. 2017. Photographic image synthesis with cascaded refinement networks. In *IEEE International Conference on Computer Vision (ICCV)*, Vol. 1. 3.
- [5] Wengling Chen and James Hays. 2018. SketchyGAN: Towards Diverse and Realistic Sketch to Image Synthesis. *CoRR abs/1801.02753* (2018). arXiv:1801.02753 <http://arxiv.org/abs/1801.02753>
- [6] Xi Chen, Yan Duan, Rein Houthooft, John Schulman, Ilya Sutskever, and Pieter Abbeel. 2016. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. In *Advances in Neural Information Processing Systems*. 2172–2180.
- [7] Jianpeng Cheng, Li Dong, and Mirella Lapata. 2016. Long Short-Term Memory-Networks for Machine Reading. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*. 551–561.
- [8] Emily L Denton, Soumith Chintala, Rob Fergus, et al. 2015. Deep generative image models using aæij laplacian pyramid of adversarial networks. In *Advances in neural information processing systems*. 1486–1494.
- [9] Alexey Dosovitskiy, Jost Tobias Springenberg, Maxim Tatarchenko, and Thomas Brox. 2017. Learning to generate chairs, tables and cars with convolutional networks. *IEEE transactions on pattern analysis and machine intelligence* 39, 4 (2017), 692–705.
- [10] A. A. Efros and T. K. Leung. 1999. Texture synthesis by non-parametric sampling. In *Proceedings of the Seventh IEEE International Conference on Computer Vision*, Vol. 2. 1033–1038 vol.2. <https://doi.org/10.1109/ICCV.1999.790383>
- [11] W. T. Freeman, T. R. Jones, and E. C. Pasztor. 2002. Example-based super-resolution. *IEEE Computer Graphics and Applications* 22, 2 (Mar 2002), 56–65. <https://doi.org/10.1109/38.988747>
- [12] L. A. Gatys, A. S. Ecker, and M. Bethge. 2016. Image Style Transfer Using Convolutional Neural Networks. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2414–2423. <https://doi.org/10.1109/CVPR.2016.265>
- [13] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative Adversarial Nets. In *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 2 (NIPS'14)*. MIT Press, Cambridge, MA, USA, 2672–2680. <http://dl.acm.org/citation.cfm?id=2969033.2969125>
- [14] Ishaaq Gulrajani, Faruk Ahmed, Martín Arjovsky, Vincent Dumoulin, and Aaron C. Courville. 2017. Improved Training of Wasserstein GANs. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*. 5769–5779. <http://papers.nips.cc/paper/7159-improved-training-of-wasserstein-gans>

- [15] X. Han, Z. Li, H. Huang, E. Kalogerakis, and Y. Yu. 2017. High-Resolution Shape Completion Using Deep Neural Networks for Global Structure and Local Geometry Inference. In *2017 IEEE International Conference on Computer Vision (ICCV)*. 85–93. <https://doi.org/10.1109/ICCV.2017.19>
- [16] James Hays and Alexei A. Efros. 2008. Scene Completion Using Millions of Photographs. *Commun. ACM* 51, 10 (Oct. 2008), 87–94. <https://doi.org/10.1145/1400181.1400202>
- [17] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, Günter Klambauer, and Sepp Hochreiter. 2017. GANs Trained by a Two Time-Scale Update Rule Converge to a Nash Equilibrium. *CoRR* abs/1706.08500 (2017). arXiv:1706.08500 <http://arxiv.org/abs/1706.08500>
- [18] Xun Huang, Yixuan Li, Omid Poursaeed, John Hopcroft, and Serge Belongie. 2017. Stacked generative adversarial networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Vol. 2.
- [19] P. Isola, J. Y. Zhu, T. Zhou, and A. A. Efros. 2017. Image-to-Image Translation with Conditional Adversarial Networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 5967–5976. <https://doi.org/10.1109/CVPR.2017.632>
- [20] Levent Karacan, Zeynep Akata, Aykut Erdem, and Erkut Erdem. 2016. Learning to generate images of outdoor scenes from attributes and semantic layouts. *arXiv preprint arXiv:1612.00215* (2016).
- [21] Taeksoo Kim, Moonsu Cha, Hyunsoo Kim, Jung Kwon Lee, and Jiwon Kim. 2017. Learning to discover cross-domain relations with generative adversarial networks. *arXiv preprint arXiv:1703.05192* (2017).
- [22] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [23] Diederik P Kingma and Max Welling. 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114* (2013).
- [24] Naveen Kodali, Jacob Abernethy, James Hays, and Zsolt Kira. 2017. On convergence and stability of GANs. *arXiv preprint arXiv:1705.07215* (2017).
- [25] C. Ledig, L. Theis, F. Huszár, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang, and W. Shi. 2017. Photo-Realistic Single Image Super-Resolution Using a Generative Adversarial Network. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 105–114. <https://doi.org/10.1109/CVPR.2017.19>
- [26] Chuan Li and Michael Wand. 2016. Precomputed real-time texture synthesis with markovian generative adversarial networks. In *European Conference on Computer Vision*. Springer, 702–716.
- [27] Ming-Yu Liu, Thomas Breuel, and Jan Kautz. 2017. Unsupervised Image-to-Image Translation Networks. In *Advances in Neural Information Processing Systems 30*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.). Curran Associates, Inc., 700–708. <http://papers.nips.cc/paper/6672-unsupervised-image-to-image-translation-networks.pdf>
- [28] Ming-Yu Liu and Oncel Tuzel. 2016. Coupled Generative Adversarial Networks. In *Proceedings of the 30th International Conference on Neural Information Processing Systems (NIPS’16)*. Curran Associates Inc., USA, 469–477. <http://dl.acm.org/citation.cfm?id=3157096.3157149>
- [29] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaou Tang. 2015. Deep Learning Face Attributes in the Wild. In *Proceedings of International Conference on Computer Vision (ICCV)*.
- [30] Mario Lucic, Karol Kurach, Marcin Michalski, Sylvain Gelly, and Olivier Bousquet. 2017. Are gans created equal? a large-scale study. *arXiv preprint arXiv:1711.10337* (2017).
- [31] X. Mao, Q. Li, H. Xie, R. Y. K. Lau, Z. Wang, and S. P. Smolley. 2017. Least Squares Generative Adversarial Networks. In *2017 IEEE International Conference on Computer Vision (ICCV)*. 2813–2821. <https://doi.org/10.1109/ICCV.2017.304>
- [32] Mehdi Mirza and Simon Osindero. 2014. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784* (2014).
- [33] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. 2018. Spectral normalization for generative adversarial networks. *arXiv preprint arXiv:1802.05957* (2018).
- [34] Takeru Miyato and Masanori Koyama. 2018. cGANs with Projection Discriminator. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=ByS1VpgRZ>
- [35] D. T. Nguyen, B. S. Hua, M. K. Tran, Q. H. Pham, and S. K. Yeung. 2016. A Field Model for Repairing 3D Shapes. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 5676–5684. <https://doi.org/10.1109/CVPR.2016.612>
- [36] Augustus Odena, Christopher Olah, and Jonathon Shlens. 2017. Conditional Image Synthesis with Auxiliary Classifier GANs. In *Proceedings of the 34th International Conference on Machine Learning (Proceedings of Machine Learning Research)*, Doina Precup and Yee Whye Teh (Eds.), Vol. 70. PMLR, International Convention Centre, Sydney, Australia, 2642–2651. <http://proceedings.mlr.press/v70/odena17a.html>
- [37] Niki Parmar, Ashish Vaswani, Jakob Uszkoreit, Lukasz Kaiser, Noam Shazeer, and Alexander Ku. 2018. Image Transformer. *CoRR* abs/1802.05751 (2018). arXiv:1802.05751 <http://arxiv.org/abs/1802.05751>
- [38] Alec Radford, Luke Metz, and Soumith Chintala. 2015. Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks. *arXiv preprint arXiv:1511.06434* (2015).
- [39] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. 2015. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*. Springer,

- 234–241.
- [40] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, Xi Chen, and Xi Chen. 2016. Improved Techniques for Training GANs. In *Advances in Neural Information Processing Systems 29*, D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett (Eds.). Curran Associates, Inc., 2234–2242. <http://papers.nips.cc/paper/6125-improved-techniques-for-training-gans.pdf>
 - [41] Lucas Theis, Aäron van den Oord, and Matthias Bethge. 2015. A note on the evaluation of generative models. *arXiv preprint arXiv:1511.01844* (2015).
 - [42] Aaron van den Oord, Nal Kalchbrenner, Oriol Vinyals, Lasse Espeholt, Alex Graves, and Koray Kavukcuoglu. 2016. Conditional Image Generation with PixelCNN Decoders. *arXiv preprint arXiv:1606.05328* (2016).
 - [43] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*. 5998–6008.
 - [44] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. 2018. High-Resolution Image Synthesis and Semantic Manipulation with Conditional GANs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
 - [45] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. 2017. Non-local neural networks. *arXiv preprint arXiv:1711.07971* 10 (2017).
 - [46] Saining "Xie and Zhuowen" Tu. 2015. Holistically-Nested Edge Detection. In *Proceedings of IEEE International Conference on Computer Vision*.
 - [47] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. 2015. Show, attend and tell: Neural image caption generation with visual attention. In *International conference on machine learning*. 2048–2057.
 - [48] Z. Yi, H. Zhang, P. Tan, and M. Gong. 2017. DualGAN: Unsupervised Dual Learning for Image-to-Image Translation. In *2017 IEEE International Conference on Computer Vision (ICCV)*. 2868–2876. <https://doi.org/10.1109/ICCV.2017.310>
 - [49] Han Zhang, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. 2018. Self-Attention Generative Adversarial Networks. *arXiv preprint arXiv:1411.1784* (2018).
 - [50] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiaolei Huang, and Dimitris N. Metaxas. 2017. StackGAN++: Realistic Image Synthesis with Stacked Generative Adversarial Networks. *CoRR* abs/1710.10916 (2017). [arXiv:1710.10916](https://arxiv.org/abs/1710.10916) <http://arxiv.org/abs/1710.10916>
 - [51] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiaolei Huang, and Dimitris N Metaxas. 2017. StackGAN: Text to Photo-Realistic Image Synthesis With Stacked Generative Adversarial Networks. In *Proceedings of the IEEE International Conference on Computer Vision*. 5907–5915.
 - [52] Junbo Zhao, Michael Mathieu, and Yann LeCun. 2016. Energy-based generative adversarial network. *arXiv preprint arXiv:1609.03126* (2016).
 - [53] J. Y. Zhu, T. Park, P. Isola, and A. A. Efros. 2017. Unpaired Image-to-Image Translation Using Cycle-Consistent Adversarial Networks. In *2017 IEEE International Conference on Computer Vision (ICCV)*. 2242–2251. <https://doi.org/10.1109/ICCV.2017.244>
 - [54] Jun-Yan Zhu, Richard Zhang, Deepak Pathak, Trevor Darrell, Alexei A Efros, Oliver Wang, and Eli Shechtman. 2017. Toward Multimodal Image-to-Image Translation. In *Advances in Neural Information Processing Systems 30*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.). Curran Associates, Inc., 465–476. <http://papers.nips.cc/paper/6650-toward-multimodal-image-to-image-translation.pdf>