# MangaGAN: Unpaired Photo-to-Manga Translation Based on The Methodology of Manga Drawing

Hao Su[1], Jianwei Niu[1,2,3*], Xuefeng Liu[1], Qingfeng Li[1], Jiahe Cui[1], and Ji Wan[1]

[1]State Key Lab of VR Technology and System, School of Computer Science and Engineering, Beihang University
[2]Industrial Technology Research Institute, School of Information Engineering, Zhengzhou University
[3]Hangzhou Innovation Institute, Beihang University
{bhsuhao, niujianwei, liu_xuefeng, liqingfeng, cuijiahe, wanji}@buaa.edu.cn

**This paper has been accepted by AAAI 2021.**



**Figure 1:** *Left:* the combination of manga faces we generated and body components we collected from a popular manga work *Bleach* [27], which shows a unified style with a strong attractiom. *Right:* the input frontal face photos and our results, where our method can effectively endow output results with both the facial similarity and the target manga style.

## ABSTRACT

Manga is a world popular comic form originated in Japan, which typically employs black-and-white stroke lines and geometric exaggeration to describe humans' appearances, poses, and actions. In this paper, we propose MangaGAN, the first method based on Generative Adversarial Network (GAN) for unpaired photo-to-manga translation. Inspired by how experienced manga artists draw manga, MangaGAN generates the geometric features of manga face by a designed GAN model and delicately translates each facial region into the manga domain by a tailored multi-GANs architecture. For training MangaGAN, we construct a new dataset collected from a popular manga work, containing manga facial features, landmarks, bodies and so on. Moreover, to produce high-quality manga faces, we further propose a structural smoothing loss to smooth stroke-lines and avoid noisy pixels, and a similarity preserving module to improve the similarity between domains of photo and manga. Extensive experiments show that MangaGAN can produce high-quality manga faces which preserve both the facial similarity and a popular manga style, and outperforms other related state-of-the-art methods.

## 1  INTRODUCTION

*Manga*, originated in Japan, is a worldwide popular comic form of drawing on serialized pages to present long stories. Typical manga is printed in black-and-white (as shown in Fig. 1 left), which employs abstract stroke lines and geometric exaggeration to describe humans' appearances, poses, and actions. Professional manga artists usually build up personalized drawing styles during their careers, and their styles are hard to be imitated by other peers. Meanwhile, drawing manga is a time-consuming process, and even a professional manga artist requires several hours to finish one page of high-quality work.

As an efficient approach to assist with manga drawing, automatically translating a face photo to manga with an attractive style is much desired. This task can be described as the image translation that is a hot topic in the computer vision field. In recent years, deep learning based image translation has made significant progress and derived a series of systematic methods. Among the examples are the *Neural Style Transfer* (NST) methods (e.g.,[4, 11, 21, 29]) which use tailored CNNs and objective functions to stylize images, the *Generative Adversarial Network* (GAN)[12] based methods (e.g.,

*The corresponding author.

[19, 36, 65]) which work well for mapping paired or unpaired images from the original domain to the stylized domain.

Although these excellent works have achieved good performances in their applications, they have difficulties to generate a high-quality manga due to the following four *challenges*. First, in the manga domain, humans' faces are abstract, colorless, geometrically exaggerated, and far from that in the photo domain. The facial correspondences between the two domains are hard to be matched by networks. Second, the style of manga is more represented by the structure of stroke lines, face shape, and facial features' sizes and locations. Meanwhile, for different facial features, manga artists always use different drawing styles and locate them with another personalized skill. These independent features (i.e., appearance, location, size, style) are almost unable to be extracted and concluded by a network simultaneously. Third, a generated manga has to faithfully resemble the input photo to keep the identity of a user without comprising the abstract manga style. It is a challenge to keep both of them with high performances. Forth, the training data of manga is difficult to collect. Manga artists often use local storyboards to show stories, which makes it difficult to find clear and complete manga faces with factors such as covered by hair or shadow, segmented by storyboards, low-resolution and so on. Therefore, related state-of-the-art methods of image stylization (e.g., [11, 19, 29, 35, 53, 60, 65]) are not able to produce desired results of manga[1].

To address these challenges, we present MangaGAN, the first GAN-based method for translating frontal face photos to the manga domain with preserving the attractive style of a popular manga work *Bleach* [27]. We observed that an experienced manga artist generally takes the following steps when drawing manga: first outlining the exaggerated face and locating the geometric distributions of facial features, and then fine-drawing each of them. MangaGAN follows the above process and employs a multi-GANs architecture to translate different facial features, and to map their geometric features by another designed GAN model. Moreover, to obtain high-quality results in an unsupervised manner, we present a Similarity Preserving (SP) module to improve the similarity between domains of photo and manga, and leverage a structural smoothing loss to avoid artifacts.

To summarize, our main contributions are three-fold:

- We propose MangaGAN, the first GAN-based method for unpaired photo-to-manga translation. It can produce attractive manga faces with preserving both the facial similarity and a popular manga style. MangaGAN uses a novel network architecture by simulating the drawing process of manga artists, which generates the exaggerated geometric features of faces by a designed GAN model, and delicately translates each facial region by a tailored multi-GANs architecture.
- We propose a similarity preserving module that effectively improves the performances on preserving both the facial similarity and manga style. We also propose a structural smoothing loss to encourage producing results with smooth stroke-lines and less messy pixels.
- We construct a new dataset called MangaGAN-BL (containing manga facial features, landmarks, bodies, etc.), collected from a world popular manga work *Bleach*. Each sample has been manually processed by cropping, angle-correction, and

repairing of disturbing elements (e.g, hair covering, shadows). MangaGAN-BL will be released for academic use.

## 2 RELATED WORK

Recent literature suggests two main directions with the ability to generate manga-like results: neural style transfer, and GAN-based cross-domain translation.

### 2.1 Neural style transfer

The goal of neural style transfer (NST) is to transfer the style from an art image to another content target image. Inspired by the progress of CNN, Gatys et al. [11] propose the pioneering NST work by utilizing CNN's power of extracting abstract features, and the style capture ability of Gram matrices [10]. Then, Li and Wand [29] use the Markov Random Field (MRF) to encode styles, and present an MRF-based method (CNNMRF) for image stylization. Afterward, various follow-up works have been presented to improve their performances on visual quality [13, 20, 35, 39, 48, 63], generating speed [4, 17, 21, 33, 50], and multimedia extension [3, 5, 15, 54, 59].

Although these methods work well on translating images into some typical artistic styles, e.g., oil painting, watercolor, they are not good at producing black-and-white manga with exaggerated geometry and discrete stroke lines, since they tend to translate textures and colors features of a target style and preserve the structure of the content image.

### 2.2 GAN-based cross-domain translation

Many GAN-based cross-domain translation methods work well on image stylization, whose goal is to learn a mapping from a source domain to a stylized domain. There are a series of works based on GAN [12] presented and applied for image stylization. Pix2Pix [19] first presents a unified framework for image-to-image translation based on conditional GANs [40]. BicycleGAN [66] extends it to multi-modal translation. Some methods including CycleGAN [65], DualGAN [61], DiscoGAN [23], UNIT [36], DTN [55] etc. are presented for unpaired one-to-one translation. MNUIT [18], startGAN [7] etc. are presented for unpaired many-to-many translation.

The methods mentioned above succeed in translation tasks that are mainly characterized by color or texture changes only (e.g., summer to winter, and apples to oranges). For photo-to-manga translation, they fail to capture the correspondences between two domains due to the abstract structure, colorless appearance, and geometric deformation of manga drawing.

Besides the above two main directions, there are also some works specially designed for creating artistic facial images. They employ techniques of Non-photorealistic rendering (NPR), data-driven synthesizing, computer graphics, etc., and have achieved much progress in many typical art forms, e.g., caricature and cartoon [2, 6, 16, 31, 47, 52, 58, 64], portrait and sketching [1, 8, 34, 43, 45, 46, 49, 56, 57, 60, 62]. However, none of them involve the generation of manga face.

---

[1]Comparison results as shown in Figure 11 and 12 of experiments.
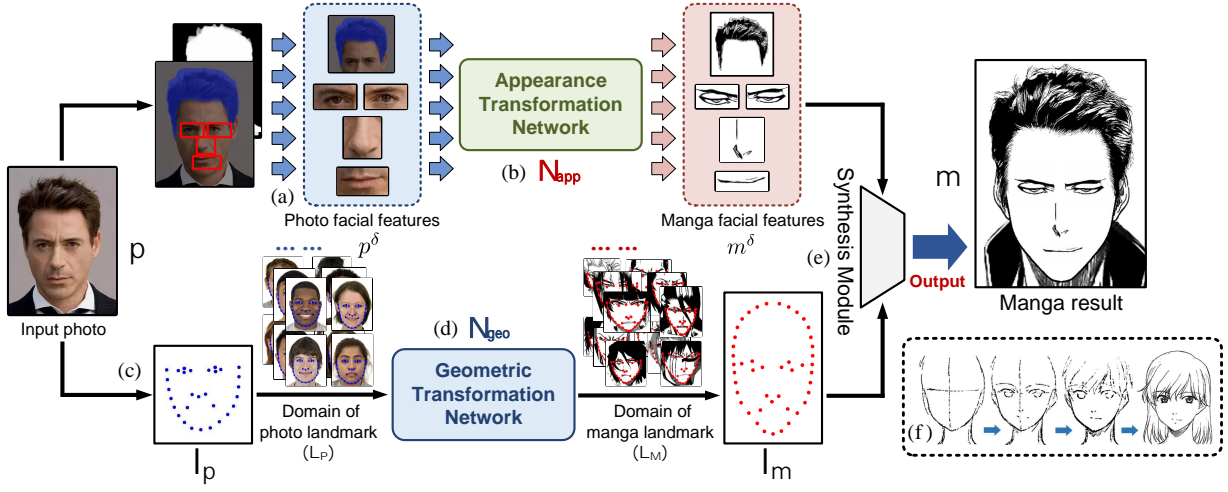
**Figure 2: Overall pipeline of MangaGAN. Inspired by the prior knowledge of manga drawing, MangaGAN consists of two branches: one branch learns the geometric mapping by a Geometric Transformation Network (GTN); the other branch learns the appearance mapping by an Appearance Transformation Network (ATN). On the end, a Synthesis Module is designed to fuse them and end up with the manga face.**

## 3 METHOD

### 3.1 Overview

Let $P$ and $M$ indicate the face photo domain and the manga domain respectively, where no pairing exists between them. Given an input photo $p \in P$, our MangaGAN learns a mapping $\Psi : P \to M$ that can transfer $p$ to a sample $m = \Psi(p)$, $m \in M$, while endowing $m$ with manga style and facial similarity.

As shown in Figure 2(f), our method is inspired by the prior knowledge that how experienced manga artists doing drawing manga: first outline the exaggerated face and locate the geometric distributions of facial features, and finally do the fine-drawing. Accordingly, MangaGAN consists of two branches: one branch learns a geometric mapping $\Psi_{geo}$ by a Geometric Transformation Network (GTN) $N_{geo}$ which adopted to translate the facial geometry from $P$ to $M$ [Figure 2(d)]; the other branch learns an appearance mapping $\Psi_{app}$ by an Appearance Transformation Network (ATN) $N_{app}$ [Figure 2(b)] which used to produce components of all facial features. At the end, a Synthesis Module is designed to fuse facial geometry and all components, and end up with the output manga $m \in M$ [Figure 2(e)]. Then, we will detail the ATN, the GTN, and the Synthesis Module in Section 3.2, Section 3.3, and Section 3.4 respectively.

### 3.2 Appearance transformation network

As shown in Figure 3, ATN $N_{app}$ is a network with multi-GAN architecture includes a set of four locals GANs, $N_{app}=\{N^{eye}, N^{nose}, N^{mouth}, N^{hair}\}$, where $N^{eye}$, $N^{nose}$, $N^{mouth}$, and $N^{hair}$ are respectively trained for translating facial regions of eye, nose, mouth, and hair, from the input $p \in P$ to the output $m \in M$.

*3.2.1 Translating regions of eyes and mouths.* Eyes and mouths are the critical components of manga faces but are the hardest parts to translate, since they are most noticed, error sensitive, and vary with different facial expressions. For $N^{eye}$ and $N^{mouth}$, for better



**Figure 3: ATN is a network with multi-GANs architecture, consists of four local GANs, designed to translate each facial region respectively. Moreover, we tailor different training strategies and encoders to improve their performances.**

mapping the unpaired data, we couple it with a reverse mapping, inspired by the network architecture of CycleGAN [65]. Accordingly, the baseline architecture of $N^{\delta}(\delta \in \{eye, mouth\})$ includes the forward / backward generator $G_M^{\delta} / G_P^{\delta}$ and the corresponding discriminator $D_P^{\delta} / D_M^{\delta}$. $G_M^{\delta}$ learns the mapping $\Psi_{app}^{\delta} : p^{\delta} \to \widehat{m}^{\delta}$, and $G_P^{\delta}$ learns the reverse mapping $\Psi_{app}^{\delta'} : m^{\delta} \to \widehat{p}^{\delta}$, where $\widehat{m}_i^{\delta}$ and $\widehat{p}_i^{\delta}$ are the generated fake samples; the discriminator $D_P^{\delta} / D_M^{\delta}$ learn to distinguish real samples $p^{\delta} / m^{\delta}$ and fake samples $\widehat{p}^{\delta} / \widehat{m}^{\delta}$. Our generators $G_P^{\delta}$, $G_M^{\delta}$ use the Resnet 6 blocks [14], and $D_P^{\delta}$, $D_M^{\delta}$ use the Markovian discriminator of $70 \times 70$ patchGANs [19, 28, 30].

We adopt the stable least-squares losses [38] instead of negative log-likelihood objective [12] as our adversarial losses $L_{adv}$, defined

as

$$\mathcal{L}_{adv}^{\delta}(G_M^{\delta}, D_M^{\delta}) = \mathbb{E}_{m^{\delta} \sim M^{\delta}}[(D_M^{\delta}(m^{\delta}) - 1)^2] + \mathbb{E}_{p^{\delta} \sim P^{\delta}}[D_M^{\delta}(G_M^{\delta}(p^{\delta}))^2] \quad , \tag{1}$$

while $\mathcal{L}_{adv}^{\delta}(G_P^{\delta}, D_P^{\delta})$ is defined in a similar manner.

$\mathcal{L}_{cyc}$ is the cycle-consistency loss [65] that is used to constrain the mapping solution between the input and the output domain, defined as

$$\mathcal{L}_{cyc}^{\delta}(G_P^{\delta}, G_M^{\delta}) = \mathbb{E}_{p^{\delta} \sim P^{\delta}}[\|G_P^{\delta}(G_M^{\delta}(p^{\delta})) - p^{\delta}\|_1] + \mathbb{E}_{m^{\delta} \sim M^{\delta}}[\|G_M^{\delta}(G_P^{\delta}(m^{\delta})) - m^{\delta}\|_1] \quad . \tag{2}$$

However, we find that the baseline architectures of $N^{eye}$ and $N^{mouth}$ with $\mathcal{L}_{adv}$ and $\mathcal{L}_{cyc}$ still fail to preserve the similarity between two domains. Specifically, for regions of eye and mouth, it always produces messy results since the networks almost unable to match colored photos and discrete black lines of mangas. Therefore, we further make three following improvements to optimize their performances.

First, we design a *Similarity Preserving (SP)* module with an SP loss $\mathcal{L}_{SP}$ to enhance the similarity. Second, we train an encoder $E^{eye}$ that can extract the main backbone of $p^{eye}$ to binary results, as the input of $N^{eye}$, and an encoder $E^{mouth}$ that encodes $p^{mouth}$ to binary edge-lines, used to guide the shape of manga mouth[2]. Third, a structural smoothing loss $\mathcal{L}_{SS}$ is designed for encouraging networks to produce manga with smooth stroke-lines, defined as

$$\mathcal{L}_{SS}(G_P^{\delta}, G_M^{\delta}) = \frac{1}{\sqrt{2\pi}\sigma} \Big[ \sum_{j \in \{1,2,\dots,N\}} \exp\Big( \frac{-(G_P^{\delta}(m^{\delta})_j - \mu)^2}{2\sigma^2} \Big) + \sum_{k \in \{1,2,\dots,N\}} \exp\Big( -\frac{(G_M^{\delta}(p^{\delta})_k - \mu)^2}{2\sigma^2} \Big) \Big] \quad , \tag{3}$$

where $\mathcal{L}_{SS}$ based on a Gaussian model with $\mu = \frac{255}{2}$, $G_P^{\delta}(m^{\delta})_j$ or $G_M^{\delta}(p^{\delta})_k$ is the $j$-th or $k$-th pixel of $G_P^{\delta}(m^{\delta})$ or $G_M^{\delta}(p^{\delta})$. The underlying idea is that producing unnecessary gray areas will distract and mess the manga results since manga mainly consists of black and white stroke lines. Thus, we give a pixel smaller loss when its gray value closer to black (0) or white (255), to smooth the gradient edges of black stroke lines and produce clean results.

**Similarity Preserving Module.** The main idea of SP module is that keeping the similarity between two images at a lower resolution can give them similar spatial distributions and different pixel details when they are up-sampled to a higher resolution. As shown in Figure 4(a), we append two SP modules on both forward and backward mappings of $N^{\delta}$. SP module leverages a pre-trained network $\phi$ that we designed to extract feature maps in different latent spaces and resolutions. The architecture of $\phi$ as shown in Figure 4(b), it only uses few convolutional layers since we consider the correspondences of encoded features are relatively clear. For the forward mapping $\Psi_{app}^{\delta} : \widehat{m}^{\delta} = G_M^{\delta}(p^{\delta})$, we input $p^{\delta}$ and $G_M^{\delta}(p^{\delta})$ to SP module, and optimize $G_M^{\delta}$ by minimizing the loss functions $\mathcal{L}_{SP}(G_M^{\delta}, p^{\delta})$ defined as

$$\mathcal{L}_{SP}(G_M^{\delta}, p^{\delta}) = \sum_{i \in \phi} \lambda_i \mathcal{L}_{feat}^{\phi,i}[f_i^{\phi}(p^{\delta}), f_i^{\phi}(G_M^{\delta}(p^{\delta}))] + \lambda_I \mathcal{L}_{pixel}^I [p^{\delta}, G_M^{\delta}(p^{\delta})] \quad , \tag{4}$$

---

[2]Training details of the two encoders are described in Section 7.5.1.



**Figure 4: (a) We append two SP modules on both forward and backward mappings. (b) SP module extracts feature maps with different resolutions and measures the similarities between two inputs in different latent spaces.**

where $\lambda_i$, $\lambda_I$ controls the relative importance of each objective, $\mathcal{L}_{pixel}^I$ and $\mathcal{L}_{feat}^{\phi,i}$ are used to keep the similarity on pixel-wise and different feature-wise respectively. $\mathcal{L}_{pixel}^I$ and $\mathcal{L}_{feat}^{\phi,i}$ defined as

$$\mathcal{L}_{feat}^{\phi,i}[f_i^{\phi}(p^{\delta}), f_i^{\phi}(G_M^{\delta}(p^{\delta}))] = \|f_i^{\phi}(p^{\delta}) - f_i^{\phi}(G_M^{\delta}(p^{\delta}))\|_2^2,$$
$$\mathcal{L}_{pixel}^I[p^{\delta}, G_M^{\delta}(p^{\delta})] = \|p^{\delta} - G_M^{\delta}(p^{\delta})\|_2^2, \tag{5}$$

where $f_i^{\phi}(x)$ is a feature map extracted from $i$-th layer of network $\phi$ when $x$ as the input. Note that we only extract feature maps after pooling layers.

Combining Eq.(1)-(5), the full objective for learning the appearance mappings of $N^{\delta}(\delta \in \{eye, mouth\})$ is:

$$\mathcal{L}_{app}^{\delta} = \mathcal{L}_{adv}^{\delta}(G_M^{\delta}, D_M^{\delta}) + \mathcal{L}_{adv}^{\delta}(G_P^{\delta}, D_P^{\delta}) + \alpha_1 \mathcal{L}_{cyc}^{\delta}(G_P^{\delta}, G_M^{\delta}) + \alpha_2 \mathcal{L}_{SP}^{\delta}(G_M^{\delta}, p^{\delta}) + \alpha_3 \mathcal{L}_{SP}^{\delta}(G_P^{\delta}, m^{\delta}) + \alpha_4 \mathcal{L}_{SS}^{\delta}(G_M^{\delta}, G_P^{\delta}) \tag{6}$$

where $\alpha_1$ to $\alpha_4$ used to balance the multiple objectives.

*3.2.2 Translating regions of nose and hair.* Noses are insignificant to manga faces since almost all characters have a similar nose in the target manga style. Therefore, $N^{nose}$ adopts a generating method instead of a translating one, which follows the architecture of progressive growing GANs [22] that can produce a large number of high-quality results similar to training data. As shown in Figure 3(d), we first train a variational autoencoder [26] to encode the nose region of the input photo into a feature vector, then make the vector as a seed to generate a default manga nose, and we also allow users to change it according to their preferences.

$N^{hair}$ employs a pre-trained generator of APDdrawingGAN [60] that can produce binary portrait hair with the style similar to manga.
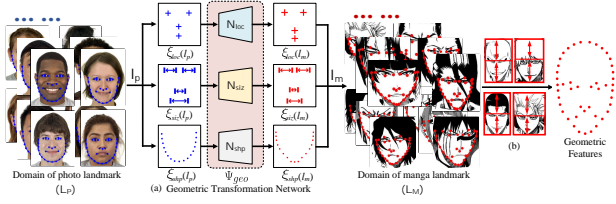
**Figure 5: The pipeline of GTN. (a)** To improve the variety of facial collocation mode, GTN divides geometric information into three independent attributes, i.e., facial features' locations, sizes, and face shape, while including three sub-GANs $N_{loc}$, $N_{siz}$, $N_{sha}$ to targetedly translate them. **(b)** According to the pre-computed proportion of cheek and forehead, we produce the geometric features of a whole manga face.

In addition, the coordinates of generated portraits can accurately correspond to the input photos. As shown in Figure 3(e), we first extract the rough hair region by a hair segmentation method [41] with a fine-tune of expanding the segmented area, and then remove the extra background area by a portrait segmentation method [51].

### 3.3 Geometric transformation network

The goal of GTN is to translate the geometric features of faces from the photo domain to the manga domain, where we represent these features with facial landmarks. Let $L_P$ and $L_M$ express the domain of landmarks corresponding to photo and manga. GTN learns a geometric mapping $\Psi_{geo}: l_p \in L_P \rightarrow l_m \in L_M$, where $l_m$ must be similar to $l_p$ and follow manga's geometric style. For training data, each landmark $l_p$ can be extracted by an existing face landmark detector [24], and 106 facial landmarks of manga data $l_m$ are manually marked by us.

When translating facial landmarks, an issue is that the collocation mode of facial features constrains the variety of results. For example, people with the same face shape may have different sizes or locations of eyes, nose, or mouth. However, GAN may generate them in a fixed or similar collocation mode when it is trained by the landmarks of global faces. Accordingly, as shown in Figure 5, we divide the geometric features into three attributions (face shape, facial features' locations and sizes) and employ three sub-GANs $N_{sha}$, $N_{loc}$, $N_{siz}$ to translate them respectively.

**Input of sub-GANs.** For $N_{loc}$, we employ relative locations instead of absolute coordinates, since directly generating coordinates may incur few facial features beyond the face profile. As shown in Figure 6(b), for $l_p$, relative locations are represented as a vector $\xi_{loc}(l_p)$. $\xi_{loc}(l_p) = \{l_p^{el}, l_p^{er}, l_p^n, l_p^m\}$ and $\xi_{loc}(l_m)$ is represented similarly, where $l_p^{el}, l_p^{er}, l_p^n, l_p^m$ represent regions of left eye, right eye, nose, and mouth respectively. Take $l_p^n$ as an example, its relative location is represented as three scalars $l_{p(d)}^{n\_cl}, l_{p(d)}^{n\_cr}, l_{p(d)}^{n\_cb}$, corresponding to distances of nose's center to cheek's left edge, right edge, and bottom edge respectively, and $l_p^{el}, l_p^{er}, l_p^m$ are defined similarly. $N_{siz}$ only learns the mapping of facial features' widths, since the length-width ratio of the generated manga facial regions are fixed. Then, the size features of $l_p$ is represent as $\xi_{siz}(l_p) = \{l_{p(w)}^{el}, l_{p(w)}^{er}, l_{p(w)}^n, l_{p(w)}^m\}$, where $l_{p(w)}^{el}, l_{p(w)}^{er}, l_{p(w)}^n, l_{p(w)}^m$ represent the width of left eye, right eye, nose, and mouth respectively. $N_{sha}$ learns to translate



**Figure 6: (a)** Architectures of $N_{loc}$, $N_{siz}$, and $N_{sha}$. **(b)** Definitions of relative locations in $\xi_{loc}(l_p)$ and $\xi_{loc}(l_m)$.

the face shape, where the face shape is represented as the landmark of cheek region containing 17 points.

**Network architecture.** As shown in Figure 6(a), $N_{loc}$, $N_{siz}$, and $N_{shp}$ roughly follow the structure of CycleGAN [65] with adversarial loss $\mathcal{L}_{adv}$ as eq(1) and cycle loss $\mathcal{L}_{cyc}$ as eq(2). Moreover, we replace all convolutional layer in generators with the fully connected layers, and add the characteristic loss $\mathcal{L}_{cha}$ [2] that leverages the differences between a face and the mean face to measure the distinctive features after exaggeration. Let $\mathcal{L}_{cha}^{L_M}(G_{L_M})$ indicates the characteristic loss on the forward mapping, defined as

$$\mathcal{L}_{cha}^{L_M}(G_{L_M}) = \mathbb{E}_{\xi_*(l_p) \sim \xi_*(L_P)} \{ 1 - \cos[\xi_*(l_p) \\ - \xi_*(\overline{L_P}), G_{L_M}(\xi_*(l_p)) - \xi_*(\overline{L_M})] \} \quad (7)$$

where $\xi_*(\overline{L_P})$ or $\xi_*(\overline{L_M})$ denotes the averages of vector $\xi_*(L_P)$ or $\xi_*(L_M)$ whose format defined by network $N_*$, $* \in \{loc, siz, shp\}$, while the reverse loss $\mathcal{L}_{cha}^{L_P}$ is defined similarly. We let $\underset{loc}{\mathcal{L}}$ denotes the loss of $N_{loc}$, and losses of $N_{siz}$ and $N_{sha}$ are represented in a similar manner. The objective function $\mathcal{L}_{geo}$ to optimize GTN is

$$\mathcal{L}_{geo} = \underset{loc}{\mathcal{L}}_{adv}^{L_P} + \underset{loc}{\mathcal{L}}_{adv}^{L_M} + \beta_1 \underset{loc}{\mathcal{L}}_{cyc} + \beta_2(\underset{loc}{\mathcal{L}}_{cha}^{L_P} + \underset{loc}{\mathcal{L}}_{cha}^{L_M})$$
$$+ \underset{siz}{\mathcal{L}}_{adv}^{L_P} + \underset{siz}{\mathcal{L}}_{adv}^{L_M} + \beta_3 \underset{siz}{\mathcal{L}}_{cyc} + \beta_4(\underset{siz}{\mathcal{L}}_{cha}^{L_P} + \underset{siz}{\mathcal{L}}_{cha}^{L_M}), \quad (8)$$
$$+ \underset{shd}{\mathcal{L}}_{adv}^{L_P} + \underset{shd}{\mathcal{L}}_{adv}^{L_M} + \beta_5 \underset{sha}{\mathcal{L}}_{cyc} + \beta_6(\underset{shd}{\mathcal{L}}_{cha}^{L_P} + \underset{shp}{\mathcal{L}}_{cha}^{L_M})$$



**Figure 7: In synthesis module, we generate manga by fusing all facial components and their geometric features.**

where $\beta_1$ to $\beta_6$ used to balance the multiple objectives.

Finally, as shown in Figure 5(b), according to the pre-defined proportion of cheek and forehead, we produce the geometric features of the whole manga face.

## 3.4 Synthesis Module

The goal of this module is to synthesis an attractive manga face by combining facial components and their geometric features. As mentioned above, facial components of eyes, nose, mouth, and hair are generated by ATN in Section 3.2, and the geometric features of them are generated by GTN in Section 3.3.

The pipeline of fusing components is shown in Figure 7. First, we resize and locate facial components following the geometric features [Figure 7(a)]. Second, the face shape is drawn by the fitting curve of generated landmarks, based on the method of *Piecewise Cubic Hermite Interpolating Polynomial (PCHIP)* [9], where PCHIP can obtain a smooth curve and effectively preserving the face shape [Figure 7 (b)]. Then, for ear regions, we provide 10 components of manga ears instead of generating them, since they are stereotyped and unimportant for facial expression. Moreover, we collect 8 manga bodies in our dataset, 5 for male, and 3 for female, that mainly used for decorating faces. In the end, we output a default manga result, and provide a toolkit that allows users to fast fine-tune the size and location of each manga component, and to switch components that insignificant for facial expression (i.e., noses, ears, and bodies) following their preferences [Figure 7 (c)].

## 4 EXPERIMENT

In the following experiments, we first introduce our dataset and training details in Section 7.8 and then evaluate the effectiveness of our improvements in Section 7.5.1. Finally, in Section 7.5.2, we compare our MangaGAN with other state-of-the-art works. We implemented MangaGAN in PyTorch [42] and all experiments are performed on a computer with an NVIDIA Tesla V100 GPU.

### 4.1 Training

**Dataset.** The datasets we used in experiments are divided into three parts, i.e., the manga dataset $\mathcal{D}_m$, the photo dataset $\mathcal{D}_p$, and the portrait dataset $\mathcal{D}_b$. $\mathcal{D}_m$, called MangaGAN-BL, is a novel dataset constructed by us and is collected from a world popular manga work *Bleach* [27]. It contains manga facial features of 448 eyes, 109 noses, 179 mouths, and 106 frontal view of manga faces whose landmarks have been marked manually. Moreover, each sample of $\mathcal{D}_m$ is normalized to 256×256 and is optimized by cropping, angle-correction, and repairing of disturbing elements (e.g, covering of hairs, glasses, shadows); $\mathcal{D}_p$ contains 1197 front view of face photos collected from CFD [37], and $\mathcal{D}_b$ contains 1197 black-and-white portraits generated by APDrawingGAN [60] when $\mathcal{D}_p$ as input.

**Training details.** For training MangaGAN, each training data of the photo domain and the manga domain is converted to grayscale with 1 channel, and each landmark of manga face is pre-processed by symmetric processing to generate more symmetrical faces. For all experiments, we set $\alpha_1$=10, $\alpha_{\{2,3\}}$=5, $\alpha_4$=1 in Eq.(6); $\beta_{\{1,3,5\}}$=10, $\beta_{\{2,4,6\}}$=1 in Eq.(8); the parameters of $\mathcal{L}_{SP}$ in Eq.(4) are fixed at $\lambda_I$=1, $\lambda_{pool5}$=1, $\lambda_i$=0, $i \in \{pool1, pool2, pool3, pool4\}$ with the output resolution of 256×256. Moreover, we employ the Adam solver [25]

with a batch size of 5. All networks use the learning rate of 0.0002 for the first 100 epochs, where the rate is linearly decayed to 0 over the next 100 epochs.

### 4.2 Ablation experiment of our improvements

In Section 3.2.1, encoders $E^{eye}$ and $E^{mouth}$ help GANs to capture the abstract correspondences of eye and mouth regions, respectively. $E^{eye}$ is a conditional GAN model basically following [19], and is pretrained by paired eye regions of photos from dataset $\mathcal{D}_p$ and their binary result from dataset $\mathcal{D}_b$; $E^{mouth}$ includes a landmark detector [24] and a pre-processed program that smoothly connects landmarks of mouth to the black edge-lines to guide the shape of a manga mouth.

With the help of $E^{eye}$ and $E^{mou}$, as shown in Figure 10, our method can effectively preserve the shape of eyebrows (red arrows), eyes, and mouths, and further abstract them into manga style. Without $E^{eye}$ or $E^{mou}$, the network cannot capture the correspondences or generated messy results, as shown in the $6^{th}$ and $12^{th}$ columns in Figure 8.

SP module is essential to keep the similarity between the photo domain and the manga domain. As shown in the $5^{th}$ and $11^{th}$ columns in Figure 8, without the SP module, neither the manga style nor the similarity between input and output can be well preserved.

Structural Smoothing (SS) loss is also a key to produce mangas with clean appearances and smooth stroke-lines. As shown in the $4^{th}$ and $10^{th}$ columns in Figure 8, for both eyes and mouth, when training with SS loss, the structure of black stroke lines are effectively smoothed and the gray messy pixels are reduced as well.

### 4.3 Comparison with state-of-the-art methods

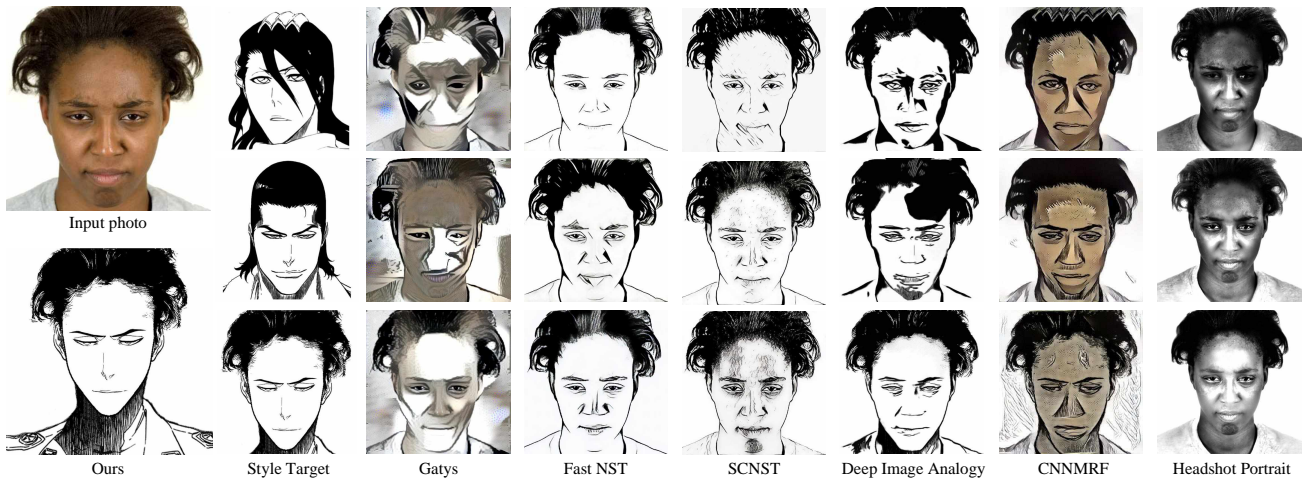We compare MangaGAN with nine state-of-the-art methods that have potentials to produce manga-like results: the first class is NST methods, containing Gatys [11], Fast NST [21], SCNST [20], Deep Image Analogy [35], CNNMRF [29], and Headshot Portrait [53]. For fair comparison, as shown in Figure 11, we employ three different manga faces (one of which is our result) as the style targets to stylize each input photo respectively. The results show t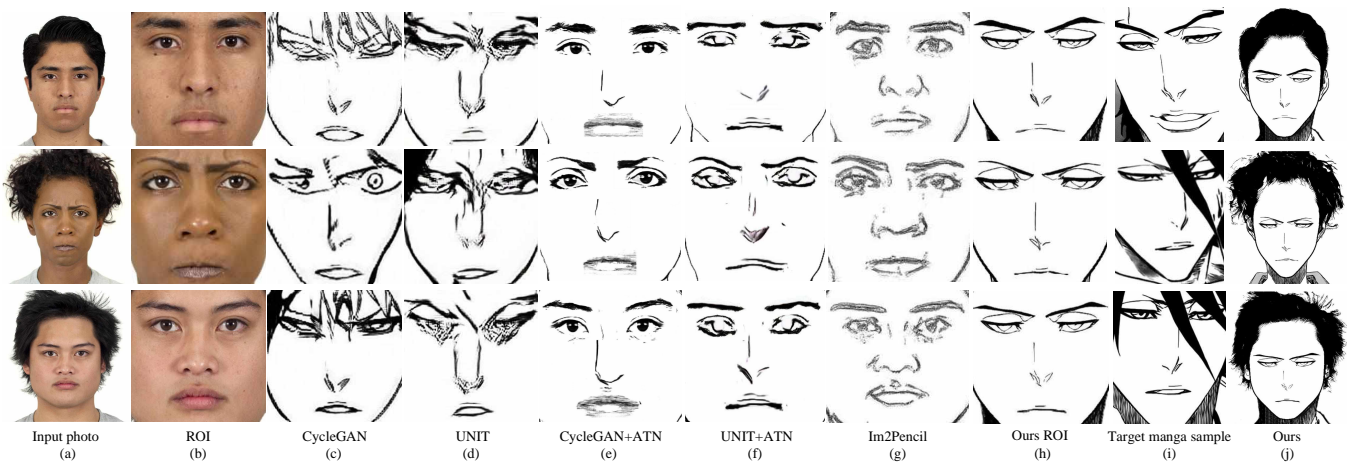hat these methods generally produce warping stroke lines and fail to produce clean manga face, since they focus on transferring the texture and color from the style target. They roughly follow the structure of the photo content, and ignore the transformation of geometric features.

The second class we compared with is cross-domain translation methods, containing CycleGAN [65], UNIT [36], and Im2pencil [32] as shown in Figure 12. For fair comparison, we train CycleGAN and UNIT to translate the whole face region, and translate each facial feature, respectively. For the whole face region translation as shown in Figure 12(c)(d), we only train the ROI [Figure 12(b)] to make these methods easier to find the correspondences between photo and manga, where the photo domain trained by 1197 frontal facial photos' ROIs in $\mathcal{D}_b$, and the manga domain trained by 83 frontal manga faces' ROIs in $\mathcal{D}_m$. For each facial feature translation as shown in Figure 12(e)(f), we append CycleGAN and UNIT on the ATN structure, and train each facial region by the same data as we use. Comparison results in Figure 12 show that the other methods get trouble in matching the poor correspondences between photo and manga, i.e., they focus on matching the dark region of photos and manga, and do not translate the face shape and stroke-line structures.

**Figure 8: Comparison results of eye and mouth regions based on different improvements. Obviously, without our improvements, the network produces poor manga results with messy regions and artifacts, and even cannot capture the correspondences between inputs and outputs.**



**Figure 9: (a) Samples of eye regions in target manga work. (b) Samples of mouth regions in target manga work. Comparison with the generated results in Figure 8 and 10, we observe our method effectively preserve the style of the target manga work.**

Unlike them, our method can effectively make the output similar to the appearance of the target manga (e.g., exaggerated eyelids, smooth eyebrows, simplified mouths) as shown in Figure 12(h)(i).

## 5 DISCUSSION

**The performance on preserving manga style**. Most of the state-of-the-art methods prone to translate the color or texture of the artistic image, and ignore the translation of geometric abstraction. As shown in Figure 11 and 12, the stylized faces they generated are similar to the input photos with only color or texture changing, which makes them more like the realistic sketches or portraits than the abstract mangas. Unlike them, we extend the translation to the structure of stroke lines and the geometric abstraction of facial features (e.g., simplified eyes and mouths, beautified facial proportions), which makes our results more like the works drawn by the manga artist.

**The performance on preserving user identity**. We generate manga face guided by the input photo, however, manga characters are typically fictitious, simplified, idealized and much unlike real people. Specifically, manga faces are usually designed to own optimum proportions, and the facial features are simplified to several black lines [Figure 12(i)]. Therefore, the excessive similarity between the output and input will make the output unlike a manga. To generate typical and clean manga faces, we even remove the detail textures and beautify the proportions of facial features, which



**Figure 10: (a) Samples of eye regions. (b) Samples of mouth regions (red lines indicate landmarks). Our method can effectively preserve the shape of eyebrows (red arrows), eyes, and mouths, and further abstracts them into manga style.**

compromise the performance on preserving the user identity. Accordingly, it is reasonable that there are some dissimilarities between the output manga face and the input facial photo.

**More evaluations**. To subjectively evaluate the performances of our methods on preserving manga style, user identity, and visual attractiveness, we conduct a series of user studies in Section 2 of the supplementary materials. Moreover, we also show more experimental results and generated manga faces in Section 5 of our supplementary materials.

**Figure 11: Comparison results with NST methods, containing Gatys [11], Fast NST [21], SCNST [20], Deep Image Analogy [35], CNNMRF [29], and Headshot Portrait [53]. For fair comparison, we employ three different manga faces (one of which is our result) as the style targets to stylize each input photo respectively.**



**Figure 12: Comparison results with cross-domain translation methods. (a) Input photo. (b) ROI of the input photo. (c)-(h) Results of CycleGAN [65], UNIT [36], Im2Pencil [32], APDrawingGAN [60], and our method, respectively. (i) Some typical face samples in target manga work [27]. We obverse that our method can effectively preserve the manga style of (i), e.g., exaggerated eyelid, smooth eyebrow, and simplified mouth. More generated samples as shown in Figure 8 and 9 in our Supplemental Material.**

# 6 CONCLUSION

In this paper, we propose the first GAN-based method for unpaired photo-to-manga translation, called MangaGAN. It is inspired by the prior-knowledge of drawing manga, and can translate a frontal face photo into the manga domain with preserving the style of a popular manga work. Extensive experiments and user studies show that MangaGAN can produce high-quality manga faces and outperforms other state-of-the-art methods.

# 7 SUPPLEMENTAL MATERIAL

## 7.1 Overview

In this document we provide the following supplementary contents:

- a series of user studies to subjectively evaluate our method and related state-of-the-art works (Section 7.2);
- more details about the ablation experiment of our improvements (Section 7.5.1);
- more qualitative results of comparison with state-of-the-art style methods (Section 7.5.2);
- details about our network architectures (Section 7.6);
- more generated samples of our MangaGAN (Section 7.7);
- our dataset and download link (Section 7.8);

- some failure cases (Section 7.9).

## 7.2 User Study

To subjectively evaluate our performances on preserving manga style, user identity and visual attractiveness, we conduct two user studies in Section 7.3 and Section 7.4 respectively.

## 7.3 Preserve manga style and user identity

**Method.** We design an online questionnaire, which first shows some samples of input photo and our corresponding results, and then appends two questions, "*How much do you think our results are similar to the target manga style?*" and "*How much do you think our results are similar to the input photos?*". All users are required to vote one of five selections (very dissimilar, dissimilar, common, similar, and very similar) according to their observation. To evaluate our work professionally, we anonymously open the questionnaire to a professional manga forum, and ask the experienced manga readers to attend this user study.

**Result.** In a two-week period, 157 participants attended this user study. The summarized results as shown in Figure 14. We observe that 86.62 % participants believe our results preserve the style of target manga, and 78.98 % participants believe our results are similar to the input photos, which indicates that our method has good performances on both two aspects.

## 7.4 Visual attractiveness

**Method.** We invited 20 volunteers (10 males and 10 females) irrelevant to this work to conduct a user study. Preparing for the experiment, we firstly select ten face photos from $\mathcal{D}_p$ randomly. Then, each photo is expanded to two group of images. The first group containing: one input photo and stylized results that are produced by six NST methods (Gatys [11], Fast NST [21], SCNST [20], Deep Image Analogy [35], CNNMRF [29], Headshot Portrait [53]), and our MangaGAN, respectively; another group containing: one input



**Figure 13: Results of our online user study.** *Upper:* **The user study on how much the similarity between our results and target manga style.** *Bottom:* **The user study on how much the similarity between our results and input photos.**



**Figure 14: User studies on the visual attractiveness of NST methods and ours.** *Left*: **Voting results of the method that has the most attractive results.** *Right*: **Boxplot of scoring results for visual attractiveness.**



**Figure 15: User studies on the visual attractiveness of cross-domain translation methods and ours.** *Left*: **Voting results of the method that has the most attractive results.** *Right*: **Boxplot of scoring results for visual attractiveness.**

photo and stylized results that are produced by five cross-domain translation methods (CycleGAN [65], UNIT [36], CycleGAN+ATN, UNIT+ATN, and Im2pencil [32]), and our MangaGAN, respectively. Finally, each volunteer is asked to complete two tasks for each image group: the first task is scoring 1 to 5 for each method's result, where a higher score indicates a higher attractiveness; another task is to vote for the method with the most attractive results.

**Result.** As shown in Figure 14, compared with NST methods, our method scored the highest on visual-quality, and over 70% volunteers believe our results are the most attractive ones. As shown in Figure 15, compared with cross-domain translation methods, our method still gets the highest score and the most number of votes. The above user studies show that our MangaGAN has reached the state-of-the-art level on visual attractiveness.

## 7.5 Supplemental Experiment

*7.5.1 Ablation experiment of our improvements.* In Figure 16, we show more comparison results corresponding to the ablation experiments in Section 4.2 of the main paper. We can observe that: the structural smoothing loss $L_{SS}$ can make the structure of stroke lines smooth, and constrain the generation of mess gray areas; the SP module successfully preserves the similarity between the input photos and the output mangas; the encoder $E^{eye}$ effectively helps the network extract the main structure of the eye region and capture the poor correspondences between photos and mangas. Without the above improvements, the model cannot generate high-quality results with clean stroke lines and an attractive manga style.

*7.5.2 More qualitative results of comparison.* According to Section 4.3 of the main paper, for more fair comparisons, we leverage related state-of-the-art methods and our methods to translate the same local facial regions (e.g., eye and mouth) respectively. For NST methods, we use three different manga eyes and mouths (one

of which is our result) as the style targets to stylize the input photo respectively. For cross-domain translation methods, we train them to translate the same local facial region, using the same dataset as us.

Comparison results as shown in Figure 17. We observe that neither the NST methods nor the cross-domain methods can generate clean and attractive manga eyes and mouthes, due to the reasons we concluded in Section 4.3 of the main paper.

## 7.6 Network Architecture

In Section 3.2 of the main paper, $N^{eye}$, $N^{mouth}$, and $N_{nose}$ are respectively trained for translating facial regions of eye, mouth, and nose, from the input photo $p \in P$ to the output manga $m \in M$. The generators of $N^{eye}$ and $N^{mouth}$ use the Resnet 6 blocks [14, 65], and the discriminators use the Markovian discriminator of $70 \times 70$ patchGANs [19, 28, 30]. We also tested using U-Net [44] or Resnet 9 blocks [14] as the generators of $N^{eye}$ and $N^{mouth}$, but they often produce messy results. Table 1 illustrates the network architectures used for the generators of $N^{eye}$ and $N^{mouth}$.

**Table 1: Network architecture used for the generators of $N^{eye}$ and $N^{mouth}$.**

| Type | Kernal Size | Output Channels | Output Size |
|---|---|---|---|
| Input | N/A | 1 | 256 |
| Conv | 7 | 64 | 256 |
| ReLu+Conv+IN | 3 | 128 | 128 |
| Residual block | 3 | 256 | 64 |
| Residual block | 3 | 256 | 64 |
| Residual block | 3 | 256 | 64 |
| Residual block | 3 | 256 | 64 |
| Residual block | 3 | 256 | 64 |
| Residual block | 3 | 256 | 64 |
| ReLu+DeConv+IN | 3 | 128 | 128 |
| ReLu+DeConv+IN | 3 | 64 | 256 |
| ReLu+Conv+IN | 7 | 1 | 256 |

$N^{nose}$ employs a generating method instead of a translating one, which follows the architecture of progressive growing GANs [22]. The network architectures of $N^{nose}$ as illustrated in Table 2.

*MangaGAN-BL* can be downloaded by the ***Google Drive*** link: https://drive.google.com/drive/folders/1viLG8fbT4lVXAwrYBO xVLoJrS2ZTUC3o?usp=sharing.

## 7.7 Generated Samples

In Table 3, we show some generated samples of paired eye regions. For one people with different facial expressions, our method successfully preserves the similarities of manga eyes, and the appearances of manga eyes are adaptively changed with the facial expressions as well; for different people, our method can effectively preserve the shape of eyebrows and eyes, and further abstract them into manga style.

Some generated samples of manga noses as shown in Figure 20. Moreover, in Figure 18 and Figure 19, we show some manga faces with high resolution, generated for males and females.

**Table 2: Network architecture used for the generators of $N^{nose}$.**

| Generator | | | |
|---|---|---|---|
| Type | Kernal Size | Output Channels | Output Size |
| Latent vector | N/A | 512 | 1 |
| Conv+LReLU | 4 | 512 | 4 |
| Conv+LReLU | 3 | 512 | 4 |
| Upsample | N/A | 512 | 8 |
| Conv+LReLU | 3 | 512 | 8 |
| Conv+LReLU | 3 | 512 | 8 |
| Upsample | N/A | 512 | 16 |
| Conv+LReLU | 3 | 512 | 16 |
| Conv+LReLU | 3 | 512 | 16 |
| Upsample | N/A | 512 | 32 |
| Conv+LReLU | 3 | 512 | 32 |
| Conv+LReLU | 3 | 512 | 32 |
| Upsample | N/A | 512 | 64 |
| Conv+LReLU | 3 | 256 | 64 |
| Conv+LReLU | 3 | 256 | 64 |
| Upsample | N/A | 256 | 128 |
| Conv+LReLU | 3 | 128 | 128 |
| Conv+LReLU | 3 | 128 | 128 |
| Upsample | N/A | 64 | 256 |
| Conv+LReLU | 3 | 64 | 256 |
| Conv+LReLU | 3 | 64 | 256 |
| Conv+liner | 1 | 3 | 256 |
| Discriminator | | | |
| Type | Kernal Size | Output Channels | Output Size |
| Input image | N/A | 3 | 256 |
| Conv+LReLU | 1 | 64 | 256 |
| Conv+LReLU | 3 | 64 | 256 |
| Conv+LReLU | 3 | 128 | 256 |
| Downsample | N/A | 128 | 128 |
| Conv+LReLU | 3 | 128 | 128 |
| Conv+LReLU | 3 | 256 | 128 |
| Downsample | N/A | 256 | 64 |
| Conv+LReLU | 3 | 256 | 64 |
| Conv+LReLU | 3 | 512 | 64 |
| Downsample | N/A | 512 | 32 |
| Conv+LReLU | 3 | 512 | 32 |
| Conv+LReLU | 3 | 512 | 32 |
| Downsample | N/A | 512 | 16 |
| Conv+LReLU | 3 | 512 | 16 |
| Conv+LReLU | 3 | 512 | 16 |
| Downsample | N/A | 512 | 8 |
| Conv+LReLU | 3 | 512 | 8 |
| Conv+LReLU | 3 | 512 | 8 |
| Downsample | N/A | 512 | 4 |
| Conv+LReLU | 3 | 512 | 4 |
| Conv+LReLU | 3 | 512 | 4 |
| Conv+LReLU | 4 | 512 | 1 |
| Fully-connected+linear | N/A | 1 | 1 |

## 7.8 Dataset

Our dataset *MangaGAN-BL* is collected from a world popular manga work *Bleach* [? ]. It contains manga facial features of 448 eyes, 109 noses, 179 mouths, and 106 frontal view of manga faces whose landmarks have been marked manually. Moreover, each sample of MangaGAN-BL is normalized to 256×256 and optimized by cropping, angle-correction, and repairing of disturbing elements (e.g, covering of hairs, glasses, shadows).

## 7.9 Failure Cases

Although our method can generate attractive manga faces in many cases, the network still produces some typical failure cases. As shown in Figure 21, when the input eyes are close to the hair, part of the hair area will be selected into the input image, which results in some artifacts in the generated manga. These failure cases are caused by the incomplete content of our dataset. For example, our data for training manga eyes only include clean eye regions, thus the model cannot be adaptive to some serious interference elements (e.g., hair, glasses).

**Figure 16: Ablation experiment of our improvements on eye regions.** *From left to right*: input face photos, results of encoder $E^{eye}$, our results, results of removing structural smoothing loss $L_{SS}$, results of removing SP module, and results of removing $E^{eye}$.

**Table 3: Some samples of eye regions in input photos and generated mangas.**

| Left Eye | | Right Eye | | Left Eye | | Right Eye | |
|---|---|---|---|---|---|---|---|
| Input | Result | Input | Result | Input | Result | Input | Result |
| People with different facial expressions | | | | | | | |



| Results of different people | | | | | | | |

**Figure 17:** *Upper:* **comparison results with NST methods, containing Gatys [11], Fast NST [21], Deep Image Analogy [35], and CNNMRF [29].** *Bottom:* **comparison results with GAN-based one-to-one translation methods, containing CycleGAN [65] and UNIT [36].**

## REFERENCES

[1] Itamar Berger, Ariel Shamir, Moshe Mahler, Elizabeth Carter, and Jessica Hodgins. Style and abstraction in portrait sketching. *ACM Transactions on Graphics (TOG)*, 32(4):55, 2013.

[2] Kaidi Cao, Jing Liao, and Lu Yuan. Carigans: unpaired photo-to-caricature translation. In *SIGGRAPH Asia 2018 Technical Papers*, page 244. ACM, 2018.

[3] Dongdong Chen, Jing Liao, Lu Yuan, Nenghai Yu, and Gang Hua. Coherent online video style transfer. In *Proc. Intl. Conf. Computer Vis.*, 2017.

[4] Dongdong Chen, Lu Yuan, Jing Liao, Nenghai Yu, and Gang Hua. Stylebank: An explicit representation for neural image style transfer. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1897–1906, 2017.

[5] Dongdong Chen, Lu Yuan, Jing Liao, Nenghai Yu, and Gang Hua. Stereoscopic neural style transfer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6654–6663, 2018.

[6] Yang Chen, Yu-Kun Lai, and Yong-Jin Liu. Cartoongan: Generative adversarial networks for photo cartoonization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9465–9474, 2018.

[7] Yunjey Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In *Proceedings of the IEEE Conference on Computer*

**Figure 18: Samples of input photos and generated manga faces**

**Figure 19: Samples of input photos and generated manga faces**
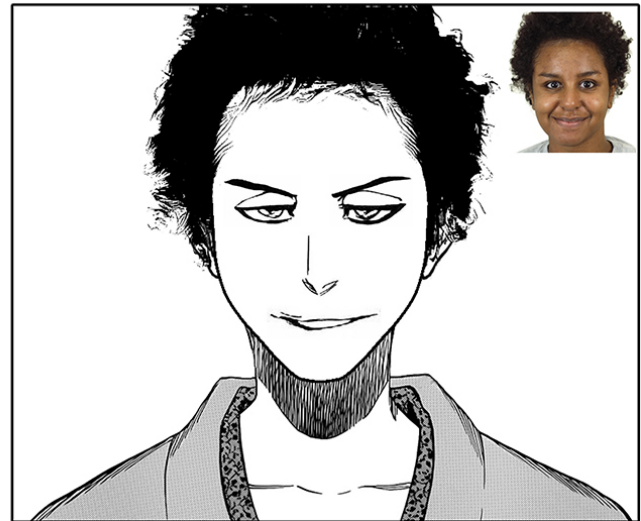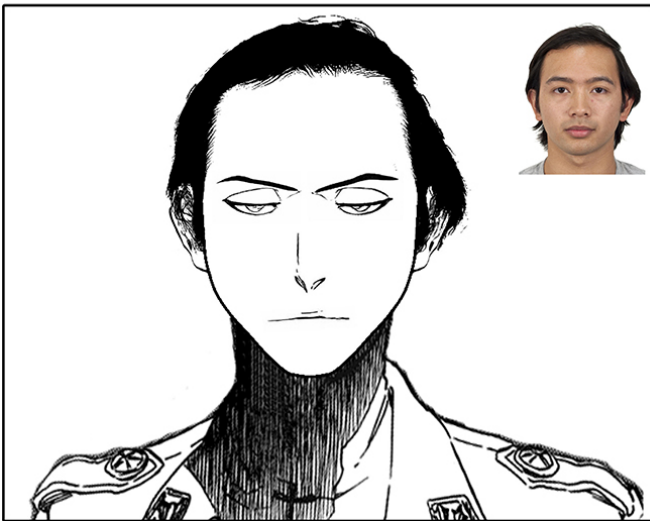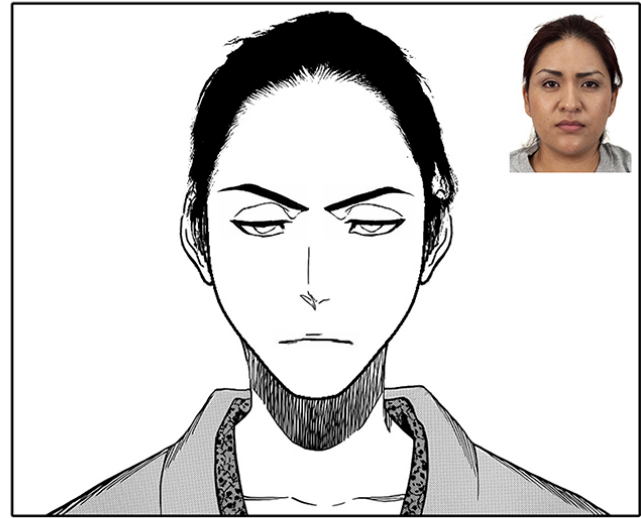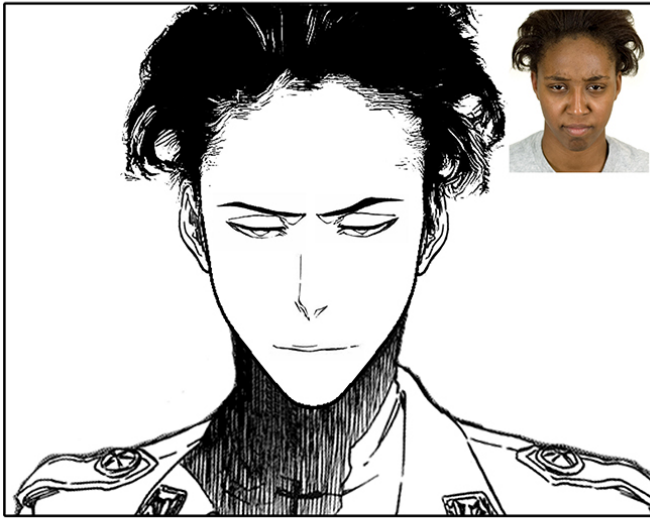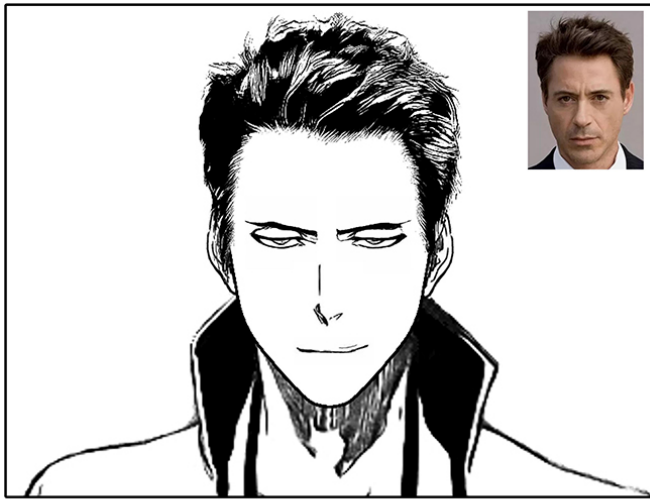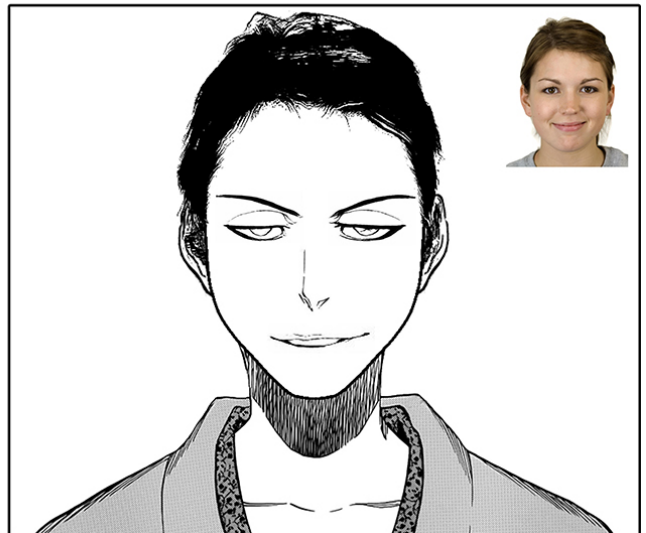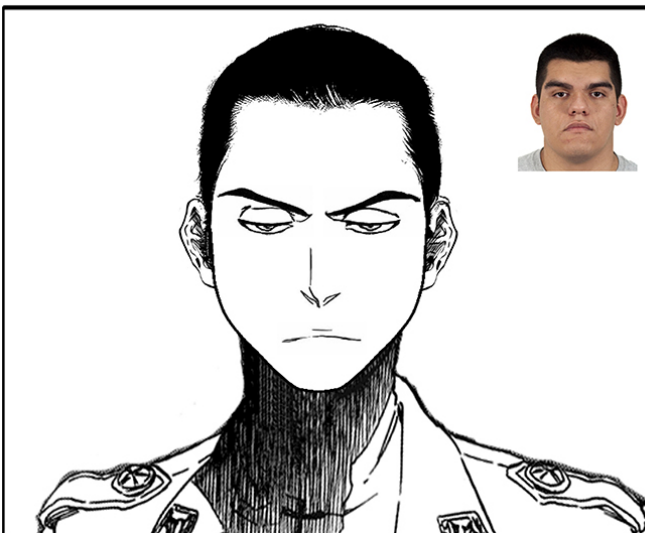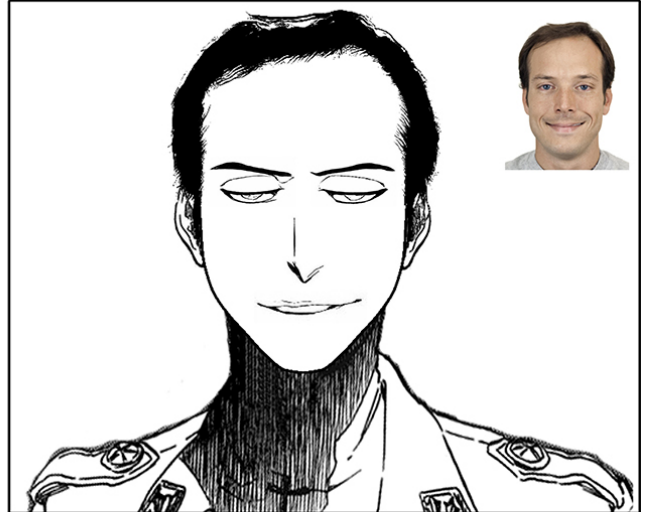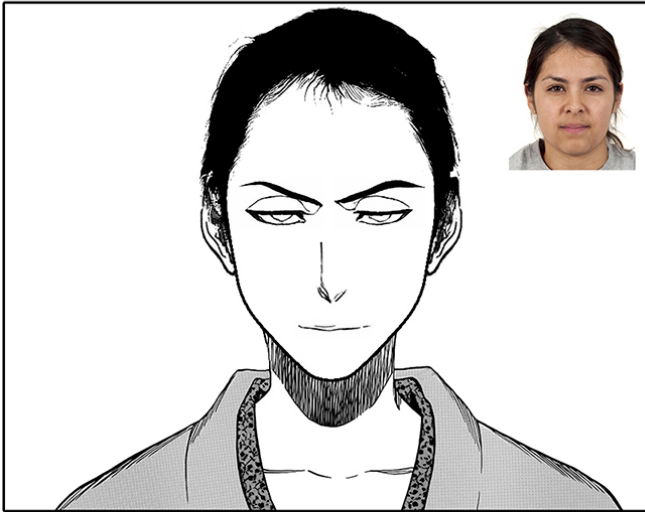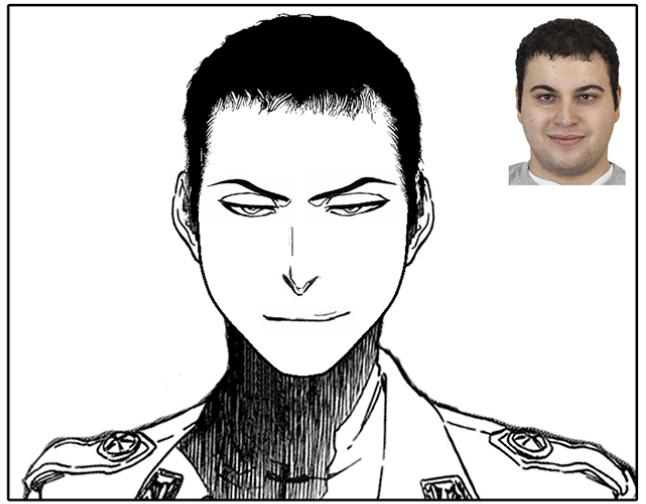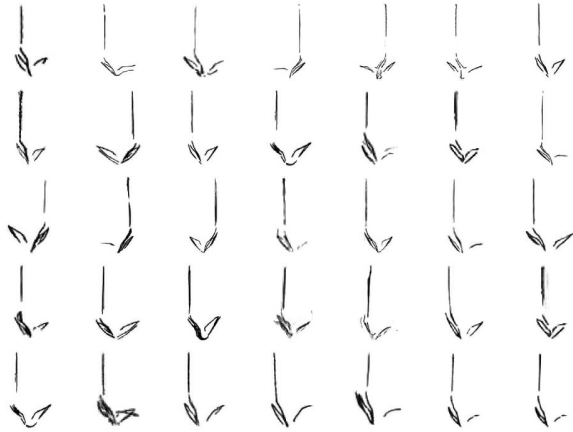
**Figure 20: Samples of generated manga noses.**



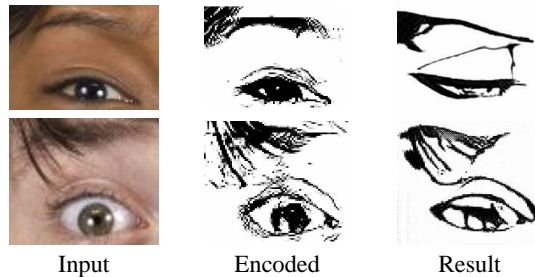Input                    Encoded                    Result

**Figure 21: Typical failure cases of our method. When the input eyes are close to the hair, part of the hair area may be selected into the input image, which results in some artifacts in the generated manga.**

*Vision and Pattern Recognition*, pages 8789–8797, 2018.

[8] Jakub Fišer, Ondřej Jamriška, Michal Lukáč, Eli Shechtman, Paul Asente, Jingwan Lu, and Daniel Sỳkora. Stylit: illumination-guided example-based stylization of 3d renderings. *ACM Transactions on Graphics (TOG)*, 35(4):92, 2016.

[9] Frederick N Fritsch and Ralph E Carlson. Monotone piecewise cubic interpolation. *SIAM Journal on Numerical Analysis*, 17(2):238–246, 1980.

[10] Leon Gatys, Alexander S Ecker, and Matthias Bethge. Texture synthesis using convolutional neural networks. In *Advances in neural information processing systems*, pages 262–270, 2015.

[11] Leon A. Gatys, Alexander S. Ecker, and Matthias Bethge. Image style transfer using convolutional neural networks. In *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, pages 2414–2423, 2016.

[12] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.

[13] Shuyang Gu, Congliang Chen, Jing Liao, and Lu Yuan. Arbitrary style transfer with deep feature reshuffle. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8222–8231, 2018.

[14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

[15] Haozhi Huang, Hao Wang, Wenhan Luo, Lin Ma, Wenhao Jiang, Xiaolong Zhu, Zhifeng Li, and Wei Liu. Real-time neural style transfer for videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 783–791, 2017.

[16] Junhong Huang, Mingkui Tan, Yuguang Yan, Chunmei Qing, Qingyao Wu, and Zhuliang Yu. Cartoon-to-photo facial translation with generative adversarial networks. In *Asian Conference on Machine Learning*, pages 566–581, 2018.

[17] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1501–1510, 2017.

[18] Xun Huang, Ming-Yu Liu, Serge Belongie, and Jan Kautz. Multimodal unsupervised image-to-image translation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 172–189, 2018.

[19] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5967–5976. IEEE, 2017.

[20] Yongcheng Jing, Yang Liu, Yezhou Yang, Zunlei Feng, Yizhou Yu, Dacheng Tao, and Mingli Song. Stroke controllable fast style transfer with adaptive receptive fields. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 238–254, 2018.

[21] Justin Johnson, Alexandre Alahi, and F.-F. Li. Perceptual losses for real-time style transfer and super-resolution. In *Proc. Eur. Conf. Comput. Vis.*, pages 694–711, 2016.

[22] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*, 2017.

[23] Taeksoo Kim, Moonsu Cha, Hyunsoo Kim, Jung Kwon Lee, and Jiwon Kim. Learning to discover cross-domain relations with generative adversarial networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1857–1865. JMLR. org, 2017.

[24] Davis E King. Dlib-ml: A machine learning toolkit. *Journal of Machine Learning Research*, 10(Jul):1755–1758, 2009.

[25] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

[26] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.

[27] Tite Kubo. Bleach. *Weekly Jump*, 2001-2016.

[28] Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, et al. Photo-realistic single image super-resolution using a generative adversarial network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4681–4690, 2017.

[29] Chuan Li and Michael Wand. Combining markov random fields and convolutional neural networks for image synthesis. In *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, pages 2479–2486, 2016.

[30] Chuan Li and Michael Wand. Precomputed real-time texture synthesis with markovian generative adversarial networks. In *European Conference on Computer Vision*, pages 702–716. Springer, 2016.

[31] Wenbin Li, Wei Xiong, Haofu Liao, Jing Huo, Yang Gao, and Jiebo Luo. Carigan: Caricature generation through weakly paired adversarial learning. *arXiv preprint arXiv:1811.00445*, 2018.

[32] Yijun Li, Chen Fang, Aaron Hertzmann, Eli Shechtman, and Ming-Hsuan Yang. Im2pencil: Controllable pencil illustration from photographs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1525–1534, 2019.

[33] Yijun Li, Chen Fang, Jimei Yang, Zhaowen Wang, Xin Lu, and Ming-Hsuan Yang. Universal style transfer via feature transforms. In *Advances in neural information processing systems*, pages 386–396, 2017.

[34] Dongxue Liang, Kyoungju Park, and Przemyslaw Krompiec. Facial feature model for a portrait video stylization. *Symmetry*, 10(10):442, 2018.

[35] Jing Liao, Yuan Yao, Lu Yuan, Gang Hua, and Sing Bing Kang. Visual attribute transfer through deep image analogy. *ACM Transactions on Graphics (TOG)*, 36(4):120, 2017.

[36] Ming-Yu Liu, Thomas Breuel, and Jan Kautz. Unsupervised image-to-image translation networks. In *Advances in neural information processing systems*, pages 700–708, 2017.

[37] Debbie S Ma, Joshua Correll, and Bernd Wittenbrink. The chicago face database: A free stimulus set of faces and norming data. *Behavior research methods*, 47(4):1122–1135, 2015.

[38] Xudong Mao, Qing Li, Haoran Xie, Raymond YK Lau, Zhen Wang, and Stephen Paul Smolley. Least squares generative adversarial networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2794–2802, 2017.

[39] Yifang Men, Zhouhui Lian, Yingmin Tang, and Jianguo Xiao. A common framework for interactive texture transfer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6353–6362, 2018.

[40] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014.

[41] Umar Riaz Muhammad, Michele Svanera, Riccardo Leonardi, and Sergio Benini. Hair detection, segmentation, and hairstyle classification in the wild. *Image and Vision Computing*, 2018.

[42] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017.

[43] Chunlei Peng, Xinbo Gao, Nannan Wang, Dacheng Tao, Xuelong Li, and Jie Li. Multiple representations-based face sketch–photo synthesis. *IEEE transactions on neural networks and learning systems*, 27(11):2201–2215, 2015.

[44] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015.

[45] Paul L Rosin and Yu-Kun Lai. Non-photorealistic rendering of portraits. In *Proceedings of the workshop on Computational Aesthetics*, pages 159–170. Eurographics Association, 2015.

[46] Paul L Rosin, David Mould, Itamar Berger, John P Collomosse, Yu-Kun Lai, Chuan Li, Hua Li, Ariel Shamir, Michael Wand, Tinghuai Wang, et al. Benchmarking non-photorealistic rendering of portraits. In *NPAR*, pages 11–1, 2017.

[47] Takafumi Saito and Tokiichiro Takahashi. Comprehensible rendering of 3-d shapes. In *ACM SIGGRAPH Computer Graphics*, volume 24, pages 197–206. ACM, 1990.

[48] Ahmed Selim, Mohamed Elgharib, and Linda Doyle. Painting style transfer for head portraits using convolutional neural networks. pages 129:1–129:18, 2016.

[49] Ahmed Selim, Mohamed Elgharib, and Linda Doyle. Painting style transfer for head portraits using convolutional neural networks. *ACM Transactions on Graphics (ToG)*, 35(4):129, 2016.

[50] Falong Shen, Shuicheng Yan, and Gang Zeng. Neural style transfer via meta networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8061–8069, 2018.

[51] Xiaoyong Shen, Aaron Hertzmann, Jiaya Jia, Sylvain Paris, Brian Price, Eli Shechtman, and Ian Sachs. Automatic portrait segmentation for image stylization. In *Computer Graphics Forum*, volume 35, pages 93–102. Wiley Online Library, 2016.

[52] Yichun Shi, Debayan Deb, and Anil K Jain. Warpgan: Automatic caricature generation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 10762–10771, 2019.

[53] YiChang Shih, Sylvain Paris, Connelly Barnes, William T Freeman, and Frédo Durand. Style transfer for headshot portraits. *ACM Transactions on Graphics (TOG)*, 33(4):148, 2014.

[54] Hao Su, Jianwei Niu, Xuefeng Liu, Qingfeng Li, Ji Wan, Mingliang Xu, and Tao Ren. An end-to-end method for producing scanning-robust stylized qr codes. *arXiv preprint arXiv:2011.07815*, 2020.

[55] Yaniv Taigman, Adam Polyak, and Lior Wolf. Unsupervised cross-domain image generation. *arXiv preprint arXiv:1611.02200*, 2016.

[56] Lidan Wang, Vishwanath Sindagi, and Vishal Patel. High-quality facial photo-sketch synthesis using multi-adversarial networks. In *2018 13th IEEE international conference on automatic face & gesture recognition (FG 2018)*, pages 83–90. IEEE, 2018.

[57] Nannan Wang, Xinbo Gao, Leiyu Sun, and Jie Li. Bayesian face sketch synthesis. *IEEE Transactions on Image Processing*, 26(3):1264–1274, 2017.

[58] Holger Winnemöller, Sven C Olsen, and Bruce Gooch. Real-time video abstraction. In *ACM Transactions On Graphics (TOG)*, volume 25, pages 1221–1226. ACM, 2006.

[59] Mingliang Xu, Hao Su, Yafei Li, Xi Li, Jing Liao, Jianwei Niu, Pei Lv, and Bing Zhou. Stylized aesthetic qr code. *IEEE Transactions on Multimedia*, 2019.

[60] Ran Yi, Yong-Jin Liu, Yu-Kun Lai, and Paul L Rosin. Apdrawinggan: Generating artistic portrait drawings from face photos with hierarchical gans. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 10743–10752, 2019.

[61] Zili Yi, Hao Zhang, Ping Tan, and Minglun Gong. Dualgan: Unsupervised dual learning for image-to-image translation. In *Proceedings of the IEEE international conference on computer vision*, pages 2849–2857, 2017.

[62] Shengchuan Zhang, Xinbo Gao, Nannan Wang, Jie Li, and Mingjin Zhang. Face sketch synthesis via sparse representation-based greedy search. *IEEE transactions on image processing*, 24(8):2466–2477, 2015.

[63] Yexun Zhang, Ya Zhang, and Wenbin Cai. Separating style and content for generalized style transfer. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8447–8455, 2018.

[64] Yong Zhang, Weiming Dong, Chongyang Ma, Xing Mei, Ke Li, Feiyue Huang, Bao-Gang Hu, and Oliver Deussen. Data-driven synthesis of cartoon faces using different styles. *IEEE Transactions on image processing*, 26(1):464–478, 2016.

[65] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232, 2017.

[66] Jun-Yan Zhu, Richard Zhang, Deepak Pathak, Trevor Darrell, Alexei A Efros, Oliver Wang, and Eli Shechtman. Toward multimodal image-to-image translation. In *Advances in Neural Information Processing Systems*, pages 465–476, 2017.