# WarpGAN: Automatic Caricature Generation

Yichun Shi*    Debayan Deb*    Anil K. Jain
Michigan State University, East Lansing MI 48824

{shiyichu, debdebay}@msu.edu, jain@cse.msu.edu

Figure 1: Example photos and caricatures of two subjects in our dataset. Column (a) shows each identity's real face photo, while two generated caricatures of the same subjects by WarpGAN are shown in column (b) and (c). Caricatures drawn by artists are shown in the column (d) and (e).

## Abstract

*We propose, WarpGAN, a fully automatic network that can generate caricatures given an input face photo. Besides transferring rich texture styles, WarpGAN learns to automatically predict a set of control points that can warp the photo into a caricature, while preserving identity. We introduce an identity-preserving adversarial loss that aids the discriminator to distinguish between different subjects. Moreover, WarpGAN allows customization of the generated caricatures by controlling the exaggeration extent and the visual styles. Experimental results on a public domain dataset, WebCaricature, show that WarpGAN is capable of generating caricatures that not only preserve the identities but also outputs a diverse set of caricatures for each input photo. Five caricature experts suggest that caricatures generated by WarpGAN are visually similar to hand-drawn ones and only prominent facial features are exaggerated.*

## 1. Introduction

A *caricature* is defined as "*a picture, description, or imitation of a person or a thing in which certain striking char-acteristics are exaggerated in order to create a comic or grotesque effect*" [1]. Paradoxically, caricatures are images with facial features that represent the face more than the face itself. Compared to cartoons, which are 2D visual art that try to re-render an object or even a scene in a usually simplified artistic style, caricatures are portraits that have exaggerated features of a certain persons or things. Some example caricatures of two individuals are shown in Figure 1. The fascinating quality of caricatures is that even with large amounts of distortion, the identity of person in the caricature can still be easily recognized by humans. In fact, studies have found that we can recognize caricatures even more accurately than the original face images [2].

Caricature artists capture the most important facial features, including the face and eye shapes, hair styles, etc. Once an artist sketches a rough draft of the face, they will start to exaggerate person-specific facial features towards a larger deviation from an average face. Nowadays, artists can create realistic caricatures through computer softwares through: (1) warping the face photo to exaggerate the shape

---

* indicates equal contribution

and (2) re-rendering the texture style [3]. By mimicking this process, researchers have been working on automatic caricature generation [4, 5]. A majority of the studies focus on designing a good structural representation to warp the image and change the face shape. However, neither the identity information nor the texture differences between a caricature and a face photo are taken into consideration. In contrast, numerous works have made progress with deep neural networks to transfer image styles [6, 7]. Still these approaches merely focus on translating the texture style forgoing any changes in the facial features.

In this work, we aim to build a completely automated system that can create new caricatures from photos by utilizing Convolutional Neural Networks (CNNs) and Generative Adversarial Networks (GANs). Different from previous works on caricature generation and style transfer, we emphasize the following challenges in our paper:

- The caricature generation involves both texture changes and shape deformation.

- The faces need to be exaggerated in a manner such that they can still be recognized.

- Caricature samples exist in various visual and artistic styles (see Figure 1).

In order to tackle these challenges, we propose a new type of style transfer network, named WarpGAN, which decouples the shape deformation and texture rendering into two tasks. Akin to a human operating an image processing software, the generator in our system automatically predicts a set of control points that warp the input face photo into the closest resemblance to a caricature and also transfers the texture style through non-linear filtering. The discriminator is trained via an identity-preserving adversarial loss to distinguish between different identities and styles, and encourages the generator to synthesize diverse caricatures while automatically exaggerating facial features specific to the identity. Experimental results show that compared to state-of-the-art generation methods, WarpGAN allows for texture update along with face deformation in the image space, while preserving the identity. Compared to other style transfer GANs [13, 7], our method not only permits a transfer in texture style, but also deformation in shape. The contributions of the paper can be summarized as follows:

- A domain transfer network that decouples the texture style and geometric shape by automatically estimating a set of sparse control points to warp the images.

- A joint learning of texture style transfer and image warping for domain transfer with adversarial loss.

- A quantitative evaluation through face recognition performance shows that the proposed method retains identity information after transferring texture style and



(a) Global Parameters [14] [15] [16]    (b) Dense Deformation Field [17]

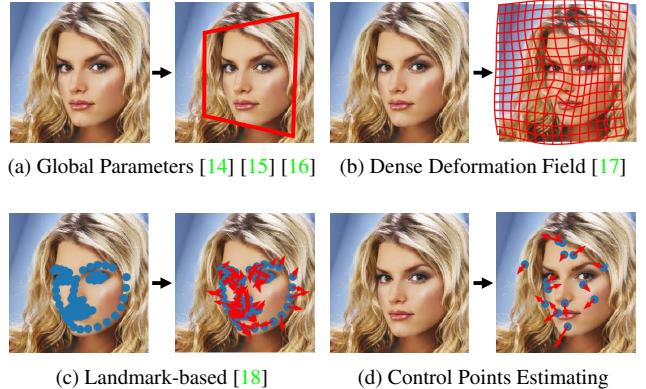(c) Landmark-based [18]    (d) Control Points Estimating

Figure 2: Inputs and outputs of different types of warping modules in neural networks. Given an image, WarpGAN can automatically predict both control points and their displacements based on local features.

warping. In addition, we conducted two perceptual studies where five caricature experts suggest that WarpGAN generates caricatures that are (1) visually appealing, (2) realistic; where only the appropriate facial features are exaggerated, and (3) our method outperforms the state-of-the-art.

- An open-source[1] automatic caricature generator where users can customize both the texture style and exaggeration degree.

## 2. Related Work

### 2.1. Automatic Image Warping

Many works have been proposed to enhance the spatial variability of neural networks via automatic warping. Most of them warp images by predicting a set of global transformation parameters [14, 16] or a dense deformation field [17]. Parametric methods estimate a small number of global transformation parameters and therefore cannot handle fine-grained local warping while dense deformation needs to predict all the vertices in a deformation grid, most of which are useless and hard to estimate. Cole *et al.* [18] first proposed to use spline interpolation in neural networks to allow control point-based warping, but their method requires pre-detected landmarks as input. Several recent works have attempted to combine image warping with GANs to improve the spatial variability of the generator, however these methods either train the warping module separately [15, 12], or need paired data as supervision [15, 19]. In comparison, our warping module can be inserted as an enhancement of a normal generator and can be trained as part of an end-to-end system without further

[1] https://github.com/seasonSH/WarpGAN

| Approach | Methodology | | | Examples |
|---|---|---|---|---|
| | **Study** | **Exaggeration Space** | **Warping** | |
| Shape Deformation | Brennan *et al.* [8] | Drawing Line | User-interactive | |
| | Liang *et al.* [4] | 2D Landmarks | User-interactive | |
| | CaricatureShop [9] | 3D Mesh | Automatic | [8]  [4]  [9] |
| Texture Transfer | Zheng *et al.* [10] | Image to Image | None | |
| | CariGAN [11] | Image + Landmark Mask | None | [10]  [11] |
| Texture + Shape | CariGANs [12] | PCA Landmarks | Automatic | |
| | WarpGAN | Image to Image | Automatic | [12]  Ours |

Table 1: Comparison of various studies on caricature generation. Majority of the published studies focus on either deforming the faces or transferring caricature styles, unlike the proposed WarpGAN which focuses on both. On the other hand, WarpGAN deforms the face in the image space thereby, truly capturing the transformations from a real face photo to a caricature. Moreover, WarpGAN does not require facial landmarks for generating caricatures.

modification. To the best of our knowledge, this study is the first work on automatic image warping with self-predicted control points using deep neural networks. An overview of different warping methods are shown in Figure 2.

## 2.2. Style Transfer Networks

Stylizing images by transferring art characteristics has been extensively studied in literature. Given the effective ability of CNNs to extract semantic features [20, 21, 22, 23], powerful style transfer networks have been developed. Gatys *et al.* [24] first proposed a neural style transfer method that uses a CNN to transfer the style content from the style image to the content image. A limitation of this method is that both the style and content images are required to be similar in nature which is not the case for caricatures. Using Generative Adversarial Networks (GANs) [25, 26] for image synthesis has been a promising field of study, where state-of-the-art results have been demonstrated in applications ranging from text to image translation [27], image inpainting [28], to image super-resolution [23]. Domain Transfer Network [29], Cycle-GAN [13], StarGAN [30], UNIT [6], and MUNIT [7] attempt image translation with unpaired image sets. All of these methods only use a de-convolutional network to construct images from the latent space and perform poorly on caricature generation due to the large spatial variation [11, 12].

## 2.3. Caricature Generation

Studies on caricature generation can be mainly classified into three categories: deformation-based, texture-based and methods with both. Traditional works mainly focused on exaggerating face shapes by enlarging the deviation of the given shape representation from average, such as 2D landmarks or 3D meshes [8, 4, 5, 9], whose deformation capability is usually limited as shape modeling can only happen in the representation space. Recently, with the success of GANs, a few works have attempted to apply style transfer networks to image-to-image caricature generation [10, 11]. However, their results suffer from poor visual quality because these networks are not suitable for problems with large spatial variation. Cao et al. [12] recently proposed to decouple texture rendering and geometric deformation with two CycleGANs trained on image and landmark space, respectively. But with their face shape modeled in the PCA subspace of landmarks, they suffer from the same problem of the traditional deformation-based methods. In this work, we propose an end-to-end system with a joint learning of texture rendering and geometric warping. Compared with previous works, WarpGAN can model both shapes and textures in the image space with flexible spatial variability, leading to better visual quality and more artistic shape exaggeration. The differences between caricature generation methods are summarized in Table 1.

## 3. Methodology

Let $\mathbf{x}_p \in \mathcal{X}_p$ be images from the domain of face photos, $\mathbf{x}_c \in \mathcal{X}_c$ be images from the caricature domain and $\mathbf{s} \in \mathcal{S}$ be the latent codes of texture styles. We aim to build a network that transforms a photo image into a caricature by both transferring its texture style and exaggerating its geometric shape. Our system includes one deformable generator (see
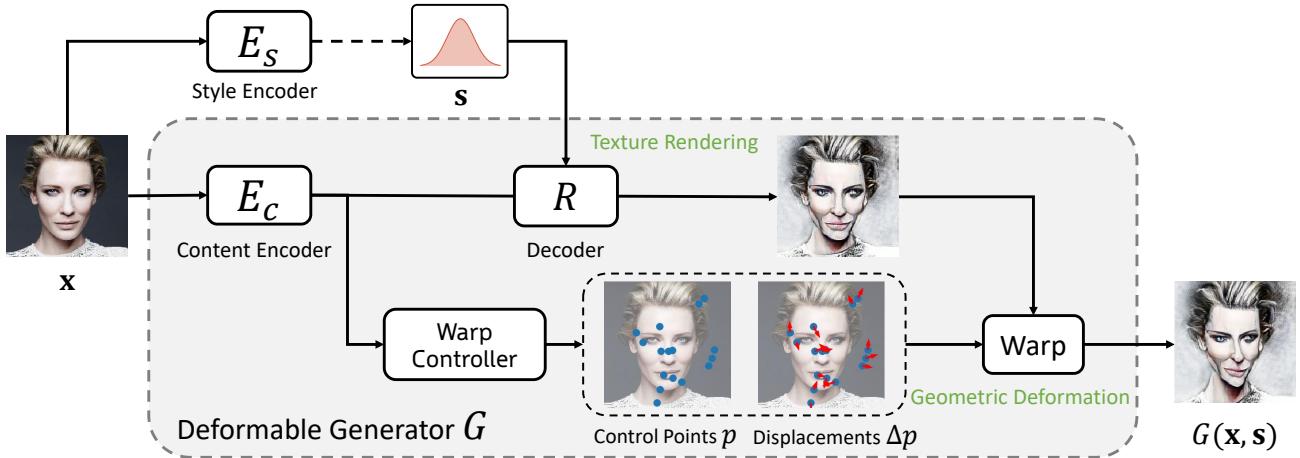
Figure 3: The generator module of WarpGAN. Given a face image, the generator outputs an image with a different texture style and a set of control points along with their displacements. A differentiable module takes the control points and warps the transferred image to generate a caricature.

| Name | Meaning | Name | Meaning |
|------|---------|------|---------|
| $\mathbf{x}_p$ | real photo image | $y^p$ | label of photo image |
| $\mathbf{x}_c$ | real caricature image | $y^c$ | label of caricature image |
| $E_c$ | content encoder | $R$ | decoder |
| $E_s$ | style encoder | $D$ | discriminator |
| $p$ | estimated control points | $\Delta p$ | displacements of $p$ |
| $M$ | number of identities | $k$ | number of control points |

Table 2: Important notations used in this paper.

Figure 3) $G$, one style encoder $E_s$ and one discriminator $D$ (see Figure 4). The important notations used in this paper are summarized in Table 2.

## 3.1. Generator

The proposed deformable generator in WarpGAN is composed of three sub-networks: a content encoder $E_c$, a decoder $R$ and a warp controller. Given any image $\mathbf{x} \in \mathbb{R}^{H \times W \times C}$, the encoder outputs a feature map $E_c(\mathbf{x})$. Here $H$, $W$ and $C$ are height, width and number of channels respectively. The content decoder takes $E_c(\mathbf{x})$ and a random latent style code $\mathbf{s} \sim \mathcal{N}(0, \mathbf{I})$ to render the given image into an image $R(E_c(\mathbf{x}), \mathbf{s})$ of a certain style. The warp controller estimates the control points and their displacements to warp the rendered images. An overview of the deformable generator is shown in Figure 3.

**Texture Style Transfer**  Since there is a large variation in the texture styles of caricatures images (See Figure 1), we adopt an unsupervised method [7] to disentangle the style representation from the feature map $E_c(\mathbf{x})$ so that we can transfer the input photo into different texture styles

present in the caricature domain. During the training, the latent style code $s \sim \mathcal{N}(0, \mathbf{I})$ is sampled randomly from a normal distribution and passed as an input into the decoder $R$. A multi-layer perceptron in $R$ decodes $\mathbf{s}$ to generate the parameters of the Adaptive Instance Normalization (AdaIN) layers in $R$, which have been shown to be effective in controlling visual styles [31]. The generated images $R(E_c(\mathbf{x}), \mathbf{s})$ with random styles are then warped and passed to the discriminator. Various styles obtained from Warp-GAN can be seen in Figure 5.

To prevent $E_c$ and $R$ from losing semantic information during texture rendering, we combine the identity mapping loss [29] and reconstruction loss [7] to regularize $E_c$ and $R$. In particular, a style encoder $E_s$ is used to learn the mapping from the image space to the style space $S$. Given its own style code, both photos and caricatures should be reconstructed from the latent feature map:

$$\mathcal{L}_{idt}^p = \mathbb{E}_{\mathbf{x}_p \in \mathcal{X}_p}[\|R(E_c(\mathbf{x}_p), E_s(\mathbf{x}_p)) - \mathbf{x}_p\|_1] \quad (1)$$

$$\mathcal{L}_{idt}^c = \mathbb{E}_{\mathbf{x}_c \in \mathcal{X}_c}[\|R(E_c(\mathbf{x}_c), E_s(\mathbf{x}_c)) - \mathbf{x}_c\|_1] \quad (2)$$

**Automatic Image Warping**  The warp controller is a sub-network of two fully connected layers. With latent feature map $E_c(\mathbf{x})$ as input, the controller learns to estimate $k$ control points $p = \{\mathbf{p}_1, \mathbf{p}_2, ...., \mathbf{p}_k\}$ and their displacement vectors $\Delta p = \{\Delta \mathbf{p}_1, \Delta \mathbf{p}_2, ... \Delta \mathbf{p}_k\}$, where each $\mathbf{p}_i$ and $\Delta \mathbf{p}_i$ is a 2D vector in the u-v space. The points are then fed into a differentiable warping module [18]. Let $p' = \{\mathbf{p}'_1, \mathbf{p}'_2, ..., \mathbf{p}'_k\}$ be the destination points, where $\mathbf{p}'_i = \mathbf{p}_i + \Delta \mathbf{p}_i$. A grid sampler of size $H \times W$ can then be computed via thin-plate spline interpolation:

$$f(\mathbf{q}) = \sum_{i=1}^{k} w_i \phi(\|\mathbf{q} - \mathbf{p}'_i\|) + \mathbf{v}^T \mathbf{q} + \mathbf{b} \quad (3)$$
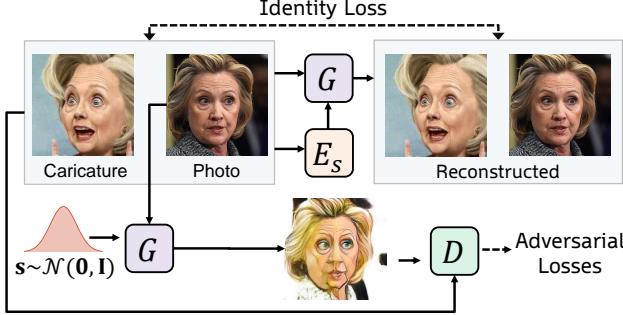
Figure 4: Overview of the proposed WarpGAN.

where the vector $\mathbf{q}$ denotes the u-v location of a pixel in the target image, and $f(\mathbf{q})$ gives the inverse mapping of the pixel $\mathbf{q}$ in the original image, and $\phi(r) = r^2 log(r)$ is the kernel function. The parameters $\mathbf{w}, \mathbf{v}, \mathbf{b}$ are fitted to minimize $\sum_j^k \left\| f(\mathbf{p}'_j) - \mathbf{p}_j \right\|^2$ and a curvature constraint, which can be solved in closed form [32]. With the grid sampler constructed via inverse mapping function $f(\mathbf{q})$, the warped image

$$G(\mathbf{x}, \mathbf{s}) = \text{Warp}\left(R(E_c(\mathbf{x}), \mathbf{s}), p, \Delta p\right) \quad (4)$$

can then be generated through bi-linear sampling [14]. The entire warping module is differentiable and can be trained as part of an end-to-end system.

### 3.2. Discriminator

**Patch Adversarial Loss** We first used a fully convolutional network as a patch discriminator [7, 13]. The patch discriminator is trained as a 3-class classifier to enlarge the difference between the styles of generated images and real photos [29]. Let $D_1$, $D_2$ and $D_3$ denote the logits for the three classes of caricatures, photos and generated images, respectively. The patch adversarial loss is as follows:

$$\mathcal{L}_p^G = - \mathbb{E}_{\mathbf{x}_p \in \mathcal{X}_p, \mathbf{s} \in S}[\log D_1(G(\mathbf{x}_p, \mathbf{s}))] \quad (5)$$

$$\mathcal{L}_p^D = - \mathbb{E}_{\mathbf{x}_c \in \mathcal{X}_c}[\log D_1(\mathbf{x}_c)] - \mathbb{E}_{\mathbf{x}_p \in \mathcal{X}_p}[\log D_2(\mathbf{x}_p)]$$
$$- \mathbb{E}_{\mathbf{x}_p \in \mathcal{X}_p, \mathbf{s} \in S}[\log D_3(G(\mathbf{x}_p, \mathbf{s}))] \quad (6)$$

**Identity-Preservation Adversarial Loss** Although patch discriminator is suitable for learning visual style transfer, it fails to capture the distinguishing features of different identities. The exaggeration styles for different people could actually be different based on their facial features (See Section 4.3). To combine the identity-preservation and identity-specific style learning, we propose to train the discriminator as a $3M$-class classifier, where $M$ is the number of identities. The first, second, and third $M$ classes correspond to different identities of real photos, real caricatures and fake caricatures, respectively. Let $y^p, y^c \in \{1, 2, 3, ...M\}$ be the

identity labels of the photos and caricatures, respectively. The identity-preservation adversarial losses for $G$ and $D$ are as follows:

$$\mathcal{L}_g^G = - \mathbb{E}_{\mathbf{x}_p \in \mathcal{X}_p, \mathbf{s} \in S}[\log D(y_p; G(\mathbf{x}_p, \mathbf{s}))] \quad (7)$$

$$\mathcal{L}_g^D = - \mathbb{E}_{\mathbf{x}_c \in \mathcal{X}_c}[\log D(y_c; \mathbf{x}_c)]$$
$$- \mathbb{E}_{\mathbf{x}_p \in \mathcal{X}_p}[\log D(y_p + M; \mathbf{x}_p)] \quad (8)$$
$$- \mathbb{E}_{\mathbf{x}_p \in \mathcal{X}_p, \mathbf{s} \in S}[\log D(y_p + 2M; G(\mathbf{x}_p, \mathbf{s}))]$$

Here, $D(y; x)$ denotes the logits of class $y$ given an image $x$. The discriminator is trained to tell the differences between real photos, real caricatures, generated caricatures as well as the identities in the image. The generator is trained to fool the discriminator in recognizing the generated image as a real caricature of the corresponding identity. Finally, the system is optimized in an end-to-end way with the following objective functions:

$$\min_G \mathcal{L}_G = \lambda_p \mathcal{L}_p^G + \lambda_g \mathcal{L}_g^G + \lambda_{idt}(\mathcal{L}_{idt}^c + \mathcal{L}_{idt}^p) \quad (9)$$

$$\min_D \mathcal{L}_D = \lambda_p \mathcal{L}_p^D + \lambda_g \mathcal{L}_g^D \quad (10)$$

## 4. Experiments

**Dataset** We use the images from a public domain dataset, WebCaricature [33][2], to conduct the experiments. The dataset consists of $6,042$ caricatures and $5,974$ photos from 252 identities. We align all the images with five landmarks. Then, the images are aligned through similarity transformation using the five landmarks and are resized to $256 \times 256$. We randomly split the dataset into a training set of 126 identities ($3,016$ photos and $3,112$ caricatures) and a testing set of 126 identities ($2,958$ photos and $2,930$ caricatures). *All the testing images in this paper are from identities in the testing set.*

**Training Details** We use ADAM optimizers in Tensorflow with $\beta_1 = 0.5$ and $\beta_2 = 0.9$ for the whole network. Each mini-batch consists of a random pair of photo and caricature. We train the network for $100,000$ steps. The learning rate starts with $0.0001$ and is decreased linearly to $0$ after $50,000$ steps. We empirically set $\lambda_g = 1.0$, $\lambda_p = 2.0$, $\lambda_{idt} = 10.0$ and number of control points $k = 16$. We conduct all experiments using Tensorflow r1.9 and one Geforce GTX 1080 Ti GPU. The average speed for generating one caricature image on this GPU is 0.082s. The details of the architecture are provided in the supplementary material.

### 4.1. Comparison to State-of-the-Art

We qualitatively compare our caricature generation method with **CycleGAN** [13], **StarGAN** [30], **Unsupervised Image-to-Image Translation (UNIT)** [6],

---

[2]https://cs.nju.edu.cn/rl/WebCaricature.htm

| Input | CycleGAN [13] | StarGAN [30] | UNIT [6] | MUNIT [7] | WarpGAN-1 | WarpGAN-2 | WarpGAN-3 |
|-------|---------------|--------------|----------|-----------|-----------|-----------|-----------|



Figure 5: Comparison of 3 different caricature styles from WarpGAN and four other state-of-the-art style transfer networks. WarpGAN is able to deform the faces unlike the baselines.

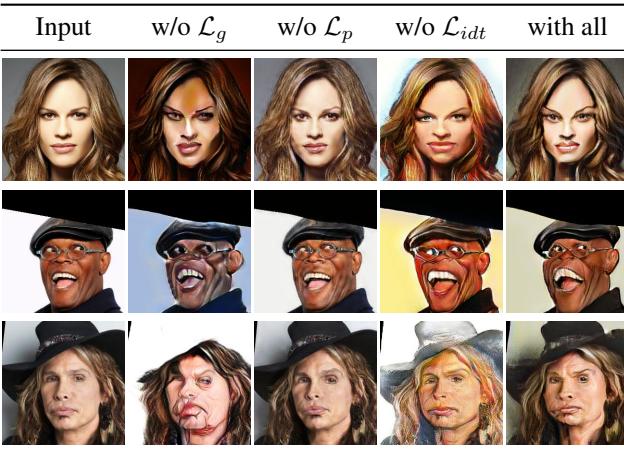| Input | w/o $\mathcal{L}_g$ | w/o $\mathcal{L}_p$ | w/o $\mathcal{L}_{idt}$ | with all |
|-------|---------------------|---------------------|-------------------------|----------|



Figure 6: Different variants of the WarpGAN without certain loss functions.

and Multimodal UNsupervised Image-to-image Translation (**MUNIT**) [7] for style transfer approaches[3]. We find that among all the three baseline style transfer networks, Cycle-GAN and MUNIT demonstrate the most visually appealing texture styles (see Figure 5). StarGAN and UNIT produce very photo-like images with minimal or erroneous changes in texture. Since all these networks focus only on transferring the texture styles, they fail to deform the faces into caricatures, unlike WarpGAN. The other issue with the baselines methods is that they do not have a module for warping the images and therefore, they try to compensate for deformations in the face using only texture. Due to the com-

plexity of this task, it becomes increasingly difficult to train them and they usually result in generating collapsed images.

## 4.2. Ablation Study

To analyze the function of different modules in our system, we train three variants of WarpGAN for comparison by removing $\mathcal{L}_g$, $\mathcal{L}_p$ and $\mathcal{L}_{idt}$, respectively. Figure 6 shows a comparison of WarpGAN variants that include all the loss functions. Without the proposed identity-preservation adversarial loss, the discriminator only focuses on local texture styles and therefore the geometric warping fails to capture personal features and is close to randomness. Without the patch adversarial loss, the discriminator mainly focuses on facial shape and the model fails to learn diverse texture styles. The model without identity mapping loss still performs well in terms of texture rendering and shape exaggeration. We keep the identity loss to improve the visual quality of the generated images.

## 4.3. Shape Exaggeration Styles

Caricaturists usually define a set of prototypes of face parts and have certain modes on how to exaggerate them [34]. In WarpGAN we do not adopt any method to exaggerate the facial regions explicitly, but instead we introduce the identity preservation constraint as part of the adversarial loss. This forces the network to exaggerate the faces to be more distinctive from other identities and implicitly encourages the network to learn different exaggeration styles for people with different salient features. Some example exaggeration styles learned by the network are shown in Figure 7.

---

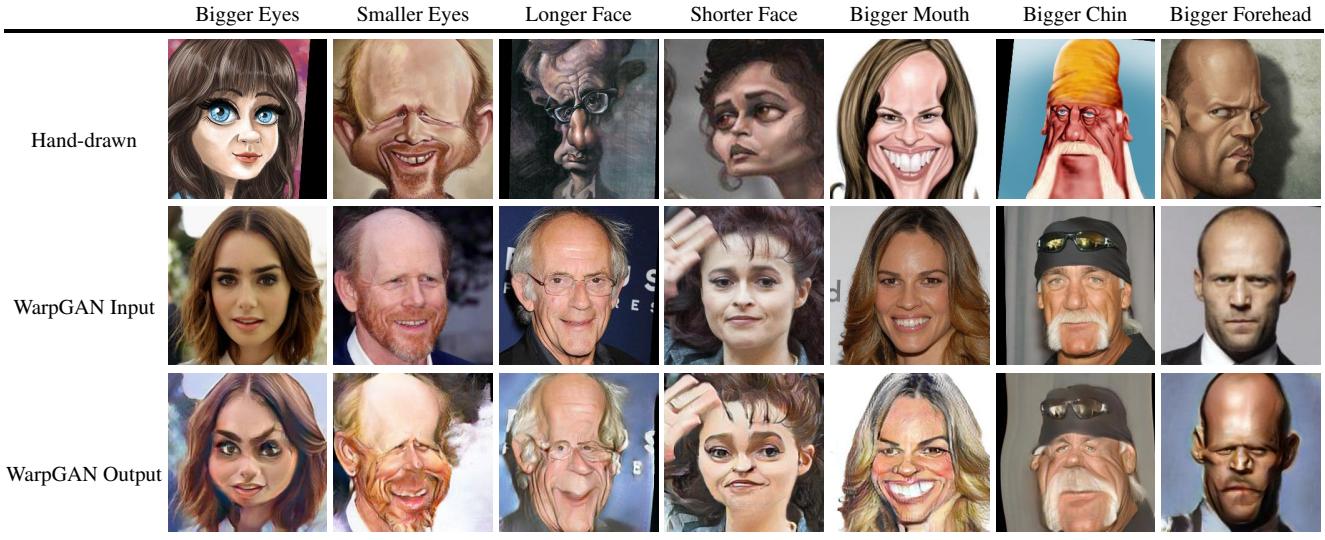[3]We train the baselines using their official implementations.

Figure 7: A few typical exaggeration styles learned by WarpGAN. First row shows hand-drawn caricatures that have certain exaggeration styles. The second and third row show the input images and the generated images of WarpGAN with the corresponding exaggeration styles. All the identities are from the testing set.
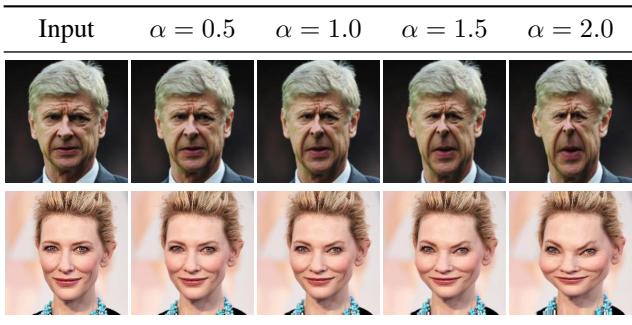


Figure 8: The result of changing the amount of exaggeration by scaling the $\Delta p$ with an input parameter $\alpha$.

| Method | COTS | SphereFace [35] |
|---|---|---|
| Photo-to-Photo | $94.81 \pm 1.22\%$ | $90.78 \pm 0.64\%$ |
| Hand-drawn-to-Photo | $41.26 \pm 1.16\%$ | $45.80 \pm 1.56\%$ |
| WarpGAN-to-Photo | $79.00 \pm 1.46\%$ | $72.65 \pm 0.84\%$ |

Table 3: Rank-1 identification accuracy for three different matching protocols using two state-of-the-art face matchers, COTS and SphereFace [35].

## 4.5. Quantitative Analysis

**Face Recognition** In order to quantify identity preservation accuracy for caricatures generated by WarpGAN, we evaluate automatic face recognition performance using two state-of-the-art face matchers: (1) a Commercial-Off-The-Shelf (COTS) matcher[4] and (2) an open source SphereFace [35] matcher.

An identification experiment is conducted where one photo of the identity is kept in the gallery while all remaining photos, or all hand-drawn caricatures, or all synthesized caricatures for the same identity are used as probes. We evaluate the Rank-1 identification accuracy using 10-fold cross validation and report the mean and standard deviation across the folds in Table 3. We find that the generated caricatures can be matched to real face images with a higher accuracy than hand-drawn caricatures. We also observe the same trend for both the matchers, which suggests that recognition on synthesized caricatures is consistent and matcher-independent.

## 4.4. Customizing the exaggeration

Although the WarpGAN is trained as a deterministic model, we introduce a parameter $\alpha$ during deployment to allow customization of the exaggeration extent. Before warping, the displacement of control points $\Delta p$ will be scaled by $\alpha$ to control how much the face shape will be exaggerated. The results are shown in Figure 8. When $\alpha = 0.0$, only the texture is changed and $\alpha = 1.0$ leads to the original output of the WarpGAN. Even when changing $\alpha$ to 2.0, the resulting images appear as caricatures, but only the distinguishing facial features are exaggerated. Since the texture styles are learned in a disentangled way, WarpGAN can generate various texture styles. Figure 5 shows results from WarpGAN with three randomly sampled styles.

---

[4]Uses a convolutional neural network for face recognition.

| Method | Visual Quality | Exaggeration |
|--------|----------------|--------------|
| Hand-Drawn | 7.70 | 7.16 |
| CycleGAN [13] | 2.43 | 2.27 |
| MUNIT [7] | 1.82 | 1.83 |
| **WarpGAN** | **5.61** | **4.87** |

Table 4: Average perceptual scores from 5 caricature experts for visual quality and exaggeration extent. Scores range from 1 to 10.



Input    Warping Only    Texture Only    Both

Figure 9: Example result images generated by the WarpGAN trained without texture/warping and with both.

**Perceptual Study** We conducted two perceptual studies by recruiting 5 caricature artists who are experts in their field to compare hand-drawn caricatures with images synthesized by our baselines along with our WarpGAN. A caricature is generated from a random image for each 126 subjects in the WebCaricature testing set. The first perceptual study uses 30 of them and 96 are used for the second. Experts do not have any knowledge of the source of the caricatures and they rely solely on their perceptual judgment.

The first study assesses the overall similarity of the generated caricatures to the hand-drawn ones. Each caricature expert was shown a face photograph of a subject along with three corresponding caricatures generated by Cycle-GAN, MUNIT, and WarpGAN, respectively. The experts then rank each of the three generated caricatures from "most visually closer to a hand-drawn caricature" to "least similar to a hand-drawn caricature". We find that caricatures generated by WarpGAN is ranked as the most similar to a real caricature $99\%$ of the time, compared to $0.5\%$ and $0.5\%$ for CycleGAN and MUNIT, respectively.

In the second study, experts scored the generated caricatures according to two criteria: (i) visual quality, and (ii) whether the caricatures are exaggerated in proper manner where only prominent facial features are deformed. Experts are shown three photographs of a subject along with a caricature image that can either be (i) a real hand-drawn caricature, or (ii) generated using one of the three automatic style transfer methods. From Table 4 we find that Warp-GAN receives the best perceptual scores out of the three methods. Even though hand-drawn caricatures rate higher, our approach, WarpGAN, has made a tremendous leap in automatically generating caricatures, especially when compared to state-of-the-art.

## 5. Discussion

**Joint Rendering and Warping Learning** Unlike other visual style transfer tasks [29, 13, 7], transforming photos into caricatures involves both texture difference and geometric transition. Texture is import in exaggerating local fine-grained features such as depth of the wrinkles while geometric deformation allows exaggeration of global features such as face shape. Conventional style transfer networks [29, 13, 7] aims to reconstruct an image from feature space using a decoder network. Because the decoder is a stack of nonlinear local filters, they are intrinsically inflexible in terms of spatial variation and the decoded images usually suffer from poor quality and severe information loss when there is a large geometric discrepancy between the input and output domain. On the other hand, warping-based methods are limited by nature to not being able to change the content and fine-grained details. Therefore, both style transfer and warping module are necessary parts for our adversarial learning framework. As shown in Figure 6, without either module, the generator will not be able to close the gap between photos and caricatures and the balance of competition between generator and discriminator will be broken, leading to collapsed results.

**Identity-preservation Adversarial Loss** The discriminator in conventional GANs are usually trained as a binary [13] or ternary classifiers [29], with each class representing a visual style. However, we found that because of the large variation of shape exaggeration in the caricatures, treating all the caricatures as one class in the discriminator would lead to the confusion of the generator, as shown in Figure 6. However, we observe that caricaturists tend to give similar exaggeration styles to the same person. Therefore, we treat each identity-domain pair as a separate class to reduce the difficulty of learning and also encourage the identity-preservation after the shape exaggeration.

## 6. Conclusion

In this paper, we proposed a new method of caricature generation, namely WarpGAN, that addresses both style transfer and face deformation in a joint learning framework. Without explicitly requiring any facial landmarks, the identity-preserving adversarial loss introduced in this work appropriately learns to capture caricature artists' style while preserving the identity in the generated caricatures. We evaluated the generated caricatures by matching synthesized caricatures to real photos and observed that the recognition accuracy is higher than caricatures drawn by artists. Moreover, five caricature experts suggest that caricatures synthesized by WarpGAN are not only pleasing to the eye, but are also realistic where only the appropriate facial features are exaggerated and that our WarpGAN indeed outperforms the state-of-the-art networks.

# References

[1] Oxford English Dictionary. Caricature Definition. https://en.oxforddictionaries.com/definition/caricature, 2018. [Online; accessed 31-October-2018]. 1

[2] Gillian Rhodes, Susan Brennan, and Susan Carey. Identification and ratings of caricatures: Implications for mental representations of faces. *Cognitive Psychology*, 1987. 1

[3] Benny Qibal. Photoshop caricature tutorial. https://www.youtube.com/watch?v=EeL2F4cgyPs, 2015. [Online; accessed 04-November-2018]. 2

[4] Lin Liang, Hong Chen, Ying-Qing Xu, and Heung-Yeung Shum. Example-based caricature generation with exaggeration. In *Pacific Conf. on Computer Graphics and Applications*, 2002. 2, 3, 10, 11

[5] Thomas Lewiner, Thales Vieira, Dimas Martínez, Adelailson Peixoto, Vinícius Mello, and Luiz Velho. Interactive 3d caricature from harmonic exaggeration. *Computers & Graphics*, 2011. 2, 3

[6] Ming-Yu Liu, Thomas Breuel, and Jan Kautz. Unsupervised image-to-image translation networks. In *NIPS*, 2017. 2, 3, 6

[7] Xun Huang, Ming-Yu Liu, Serge Belongie, and Jan Kautz. Multimodal unsupervised image-to-image translation. *arXiv:1804.04732*, 2018. 2, 3, 4, 5, 6, 8, 10

[8] Susan E Brennan. Caricature generator: The dynamic exaggeration of faces by computer. *Leonardo*, 1985. 3

[9] Xiaoguang Han, Kangcheng Hou, Dong Du, Yuda Qiu, Yizhou Yu, Kun Zhou, and Shuguang Cui. Caricatureshop: Personalized and photorealistic caricature sketching. *arXiv:1807.09064*, 2018. 3, 10, 11

[10] Ziqiang Zheng, Haiyong Zheng, Zhibin Yu, Zhaorui Gu, and Bing Zheng. Photo-to-caricature translation on faces in the wild. *arXiv:1711.10735*, 2017. 3, 10, 11

[11] Wenbin Li, Wei Xiong, Haofu Liao, Jing Huo, Yang Gao, and Jiebo Luo. CariGAN: Caricature Generation through Weakly Paired Adversarial Learning. *arXiv:1811.00445*, 2018. 3, 10, 11

[12] Kaidi Cao, Jing Liao, and Lu Yuan. CariGANs: Unpaired Photo-to-Caricature Translation. *arXiv:1811.00222*, 2018. 2, 3, 10, 11

[13] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *ICCV*, 2017. 2, 3, 5, 8

[14] Max Jaderberg, Karen Simonyan, Andrew Zisserman, et al. Spatial transformer networks. In *NIPS*, 2015. 2, 5

[15] Haoye Dong, Xiaodan Liang, Ke Gong, Hanjiang Lai, Jia Zhu, and Jian Yin. Soft-gated warping-gan for pose-guided person image synthesis. *arXiv:1810.11610*, 2018. 2

[16] Chen-Hsuan Lin, Ersin Yumer, Oliver Wang, Eli Shechtman, and Simon Lucey. St-gan: Spatial transformer generative adversarial networks for image compositing. In *CVPR*, 2018. 2

[17] Yaroslav Ganin, Daniil Kononenko, Diana Sungatullina, and Victor Lempitsky. Deepwarp: Photorealistic image resynthesis for gaze manipulation. In *ECCV*, 2016. 2

[18] Forrester Cole, David Belanger, Dilip Krishnan, Aaron Sarna, Inbar Mosseri, and William T Freeman. Synthesizing normalized faces from facial identity features. In *CVPR*, 2017. 2, 4

[19] Aliaksandr Siarohin, Enver Sangineto, Stéphane Lathuilière, and Nicu Sebe. Deformable gans for pose-based human image generation. In *CVPR*, 2018. 2

[20] Joan Bruna, Pablo Sprechmann, and Yann LeCun. Super-resolution with deep convolutional sufficient statistics. *arXiv:1511.05666*, 2015. 3

[21] Leon Gatys, Alexander S Ecker, and Matthias Bethge. Texture synthesis using convolutional neural networks. In *NIPS*, pages 262–270, 2015. 3

[22] Leon A Gatys, Alexander S Ecker, Matthias Bethge, Aaron Hertzmann, and Eli Shechtman. Controlling perceptual factors in neural style transfer. In *CVPR*, 2017. 3

[23] Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew P Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, et al. Photo-realistic single image super-resolution using a generative adversarial network. In *CVPR*, 2017. 3

[24] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. Image style transfer using convolutional neural networks. In *CVPR*, 2016. 3

[25] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NIPS*, 2014. 3

[26] Jiajun Wu, Chengkai Zhang, Tianfan Xue, Bill Freeman, and Josh Tenenbaum. Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling. In *NIPS*, 2016. 3

[27] Scott Reed, Zeynep Akata, Xinchen Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak Lee. Generative adversarial text to image synthesis. *arXiv:1605.05396*, 2016. 3

[28] Raymond Yeh, Chen Chen, Teck Yian Lim, Mark Hasegawa-Johnson, and Minh N Do. Semantic image inpainting with perceptual and contextual losses. arxiv preprint. *arXiv:1607.07539*, 2016. 3

[29] Yaniv Taigman, Adam Polyak, and Lior Wolf. Unsupervised cross-domain image generation. In *ICLR*, 2017. 3, 4, 5, 8

[30] 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. 2018. 3, 5

[31] Xun Huang and Serge J Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *ICCV*, 2017. 4, 10

[32] Chris A Glasbey and Kantilal Vardichand Mardia. A review of image-warping methods. *Journal of applied statistics*, 1998. 5

[33] Jing Huo, Wenbin Li, Yinghuan Shi, Yang Gao, and Hujun Yin. Webcaricature: a benchmark for caricature face recognition. *arXiv:1703.03230*, 2017. 5, 10

[34] Brendan F Klare, Serhat S Bucak, Anil K Jain, and Tayfun Akgul. Towards automated caricature recognition. In *ICB*, 2012. 6

[35] Weiyang Liu, Yandong Wen, Zhiding Yu, Ming Li, Bhiksha Raj, and Le Song. Sphereface: Deep hypersphere embedding for face recognition. In *CVPR*, 2017. 7

[36] Dmitry Ulyanov, Andrea Vedaldi, and Victor S Lempitsky. Improved texture networks: Maximizing quality and diversity in feed-forward stylization and texture synthesis. In *CVPR*, 2017. 10

[37] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016. 10

[38] Zhenyao Mo, John P Lewis, and Ulrich Neumann. Improved automatic caricature by feature normalization and exaggeration. In *SIGGRAPH*, 2004. 10, 11

[39] Mahdi M Kalayeh, Misrak Seifu, Wesna LaLanne, and Mubarak Shah. How to take a good selfie? In *ACM MM*, 2015. 11

## A. Implementation Details

**Preprocessing** We align all the images with five landmarks (left eye, right eye, nose, mouth left, mouth right) using the ones provided in the WebCaricature dataset [33] protocol. Since the protocol does not provide the locations of eye centers, we estimate them by taking the average of the corresponding eye corners. Then, a similarity transformation is applied for all the images using the five landmarks. The aligned images are resized to $256 \times 256$. The whole dataset consists of $6,042$ caricatures and $5,974$ photos from $252$ identities. We randomly split the dataset into a training set of $126$ identities ($3,016$ photos and $3,112$ caricatures) and a testing set of $126$ identities ($2,958$ photos and $2,930$ caricatures). **All the testing images in the main paper and this supplementary material are from the identities in the testing split**.

**Experiment Settings** We conduct all experiments using Tensorflow r1.9 and one Geforce GTX 1080 Ti GPU. The average speed for generating one caricature image on this GPU is 0.082s.

**Architecture** Our network architecture is modified based on MUNIT [7]. Let `c7s1-k` be a $7 \times 7$ convolutional layer with $k$ filters and stride 1. `dk` denotes a $4 \times 4$ convolutional layer with $k$ filters and stride 2. `Rk` denotes a residual block that contains two $3 \times 3$ convolutional layers. `uk`

denotes a $2\times$ upsampling layer followed by a $5 \times 5$ convolutional layer with $k$ filters and stride 1. `fck` denotes a fully connected layer with $k$ filters. `avgpool` denotes a global average pooling layer. We apply Instance Normalization (IN) [36] to the content encoder and Adaptive Instance Normalization (AdaIN) [31] to the decoder. No normalization is used in the style encoder. We use Leaky ReLU with slope 0.2 in the discriminator and ReLU activation everywhere else. The architectures of different modules are as follows:

- Style Encoder:
  `c7s1-64,d128,d256,avgpool,fc8`
- Content Encoder:
  `c7s1-64,d128,d256,R256,R256,R256`
- Decoder:
  `R256,R256,R256,u128,u64,c7s1-3`
- Discriminator:
  `d32,d64,d128,d256,d512,fc512,fc3M`

A separate branch of $1 \times 1$ convolutional layer with 3 filters and stride 1 is attached to the last convolutional layer of the discriminator to output $D_1, D_2, D_3$ for patch adversarial losses. The style decoder (the multi-layer perceptron) has two hidden fully connected layers of 128 filters without normalization and the warp controller has only one hidden fully connected layer of 128 filters with Layer Normalization [37]. The length of the latent style code is set to 8.

## B. Additional Baselines

In the main paper, we compared WarpGAN with state-of-the-art style transfer networks as baselines. Here, we compare WarpGAN with other caricature generation works [4, 38, 9, 10, 11, 12]. Since these methods do not release their code and use different testing images, we crop the images from their papers and compare with them one by one. All the baseline results are also taken from their original papers. The results are shown in Figure 10.

## C. Transformation Methods

To see the advantage of the proposed control-points estimation for automatic warping, we train three variants of our model by replacing the warping method with (1) projective transformation, (2) dense deformation and (3) landmark-based warping. In projective transformation, the warp controller outputs 8 parameters for the transformation matrix. In dense deformation, the warp controller outputs a $16 \times 16$ deformation grid, which is further interpolated into $256 \times 256$ for grid sampling. In landmark-based warping, we use the landmarks provided by Dlib[5] and the warp

---

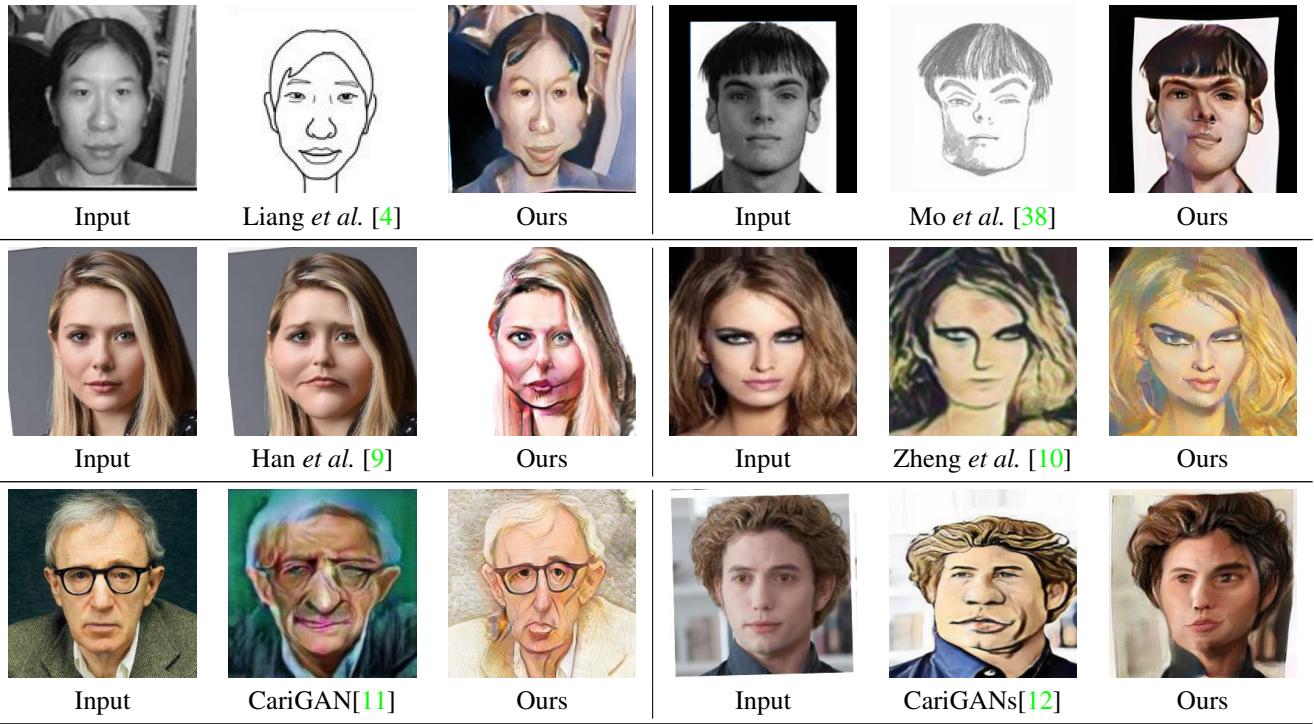[5]http://dlib.net/face_landmark_detection.html

Figure 10: Comparison with previous works on caricature generation. In each cell, the left and middle images are the input and result images taken from the baseline paper, respectively. The right images are the results of WarpGAN.

controller only outputs the displacements. As shown in Figure 12, the warping is too limited in projective transformation for generating artistic caricatures and too unconstrained in dense deformation that it is difficult to train. Landmark-based warping yields reasonable results, but it is limited by the landmark detector. In comparison, our methods does not require any domain knowledge, has little limitation and leads to visually satisfying warping results.

## D. More Results

**Ablation Study**  We show more results of the ablation study in Figure 11. The results are consistent with those in the main paper: (1) the joint learning of texture rendering and warping are crucial for generating realistic caricature images and (2) without patch adversarial loss or identity-preservation adversarial loss, the model cannot learn to generate caricatures with various texture styles and shape exaggeration styles.

**Different Texture Styles**  More results of texture style controlling are shown in Figure 13. Five latent style codes are randomly sampled from the normal distribution $\mathcal{N}(0, \mathbf{I})$. Images in the same column in Figure 13 are generated with the same style code.

**Selfie Dataset**  To test the performance of our model in more application scenarios, we download the public Selfie dataset[6] [39] for cross-dataset evaluation. The dataset includes $46,836$ public selfies crawled from Internet. Unlike our training dataset (WebCaricature), the identities in this dataset are not restricted to celebrities and there is a difference between the visual styles of these images and the ones in our training dataset. The results are shown in Figure 14.

---

[6]http://crcv.ucf.edu/data/Selfie/

| Input | w/o texture | w/o warping | w/o $\mathcal{L}_g$ | w/o $\mathcal{L}_p$ | w/o $\mathcal{L}_{idt}$ | with all |
|-------|-------------|-------------|---------------------|---------------------|-------------------------|----------|



Figure 11: More results on ablation study. Input images are shown in the first column. The subsequent columns show the results of different models trained without a certain module or loss. The texture style codes are randomly sampled from the normal distribution.

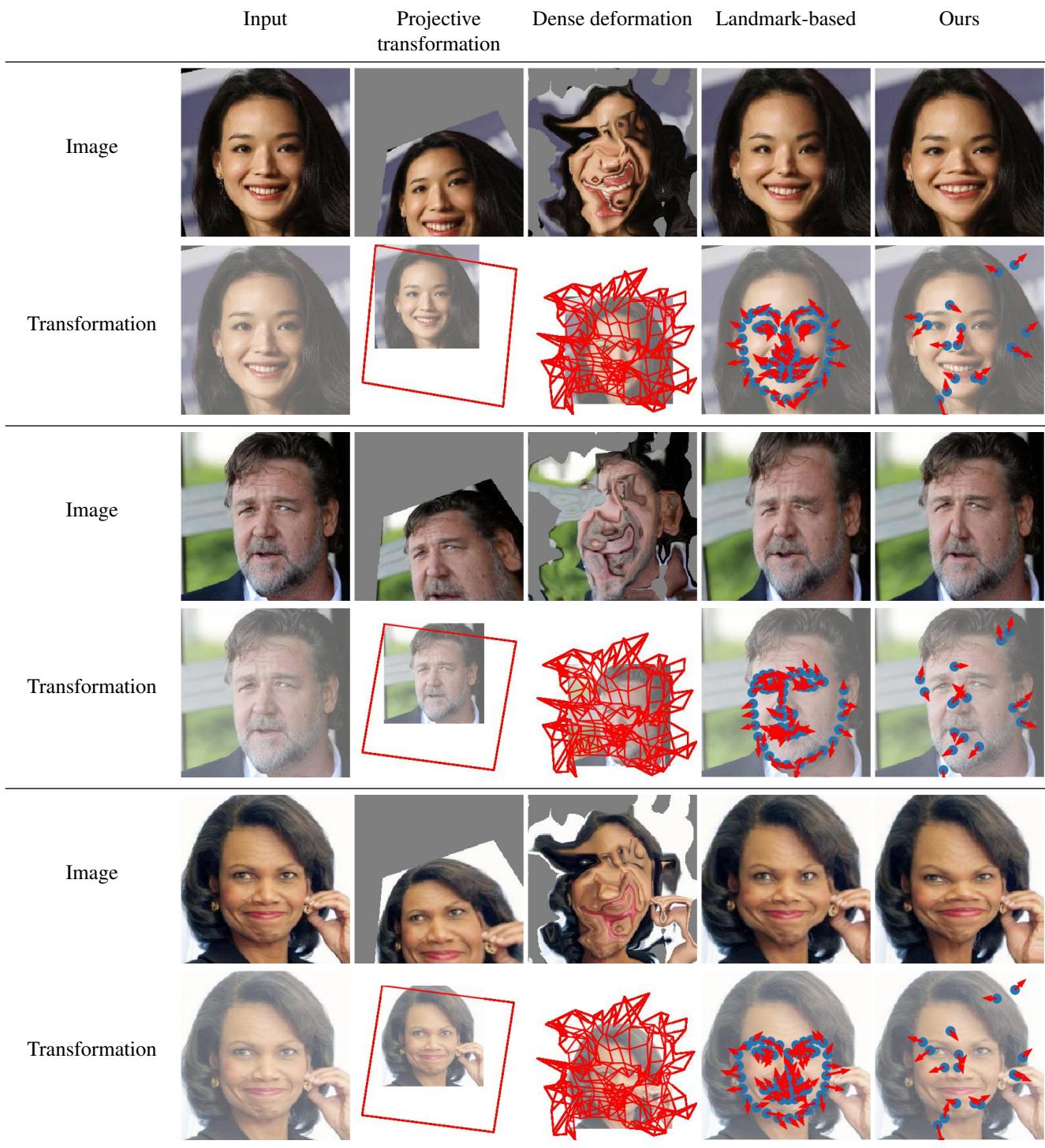|  | Input | Projective transformation | Dense deformation | Landmark-based | Ours |

Figure 12: Different transformation methods. Input images are shown in the first column. The next four columns show the results and the transformation visualizations of four different models trained with different transformation methods. The landmark-based model uses 68 landmarks detected by Dlib. Texture rendering is hidden here for clarity.
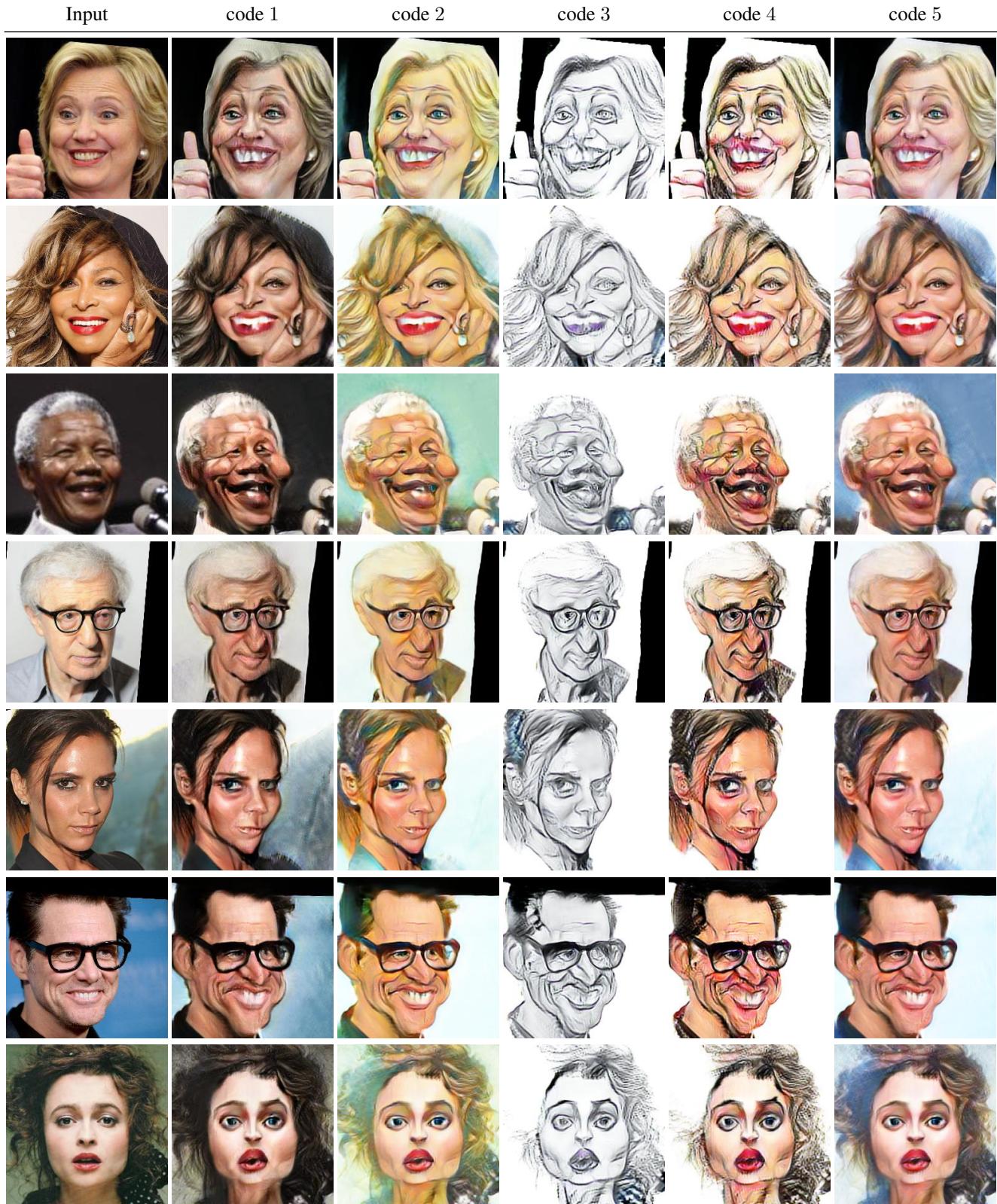
Figure 13: Results of five different texture styles. Input images are shown in the first column. Subsequent five columns show the results of WarpGAN using five style codes sampled randomly from the normal distribution. All the images in the same column are generated with the same latent style code.
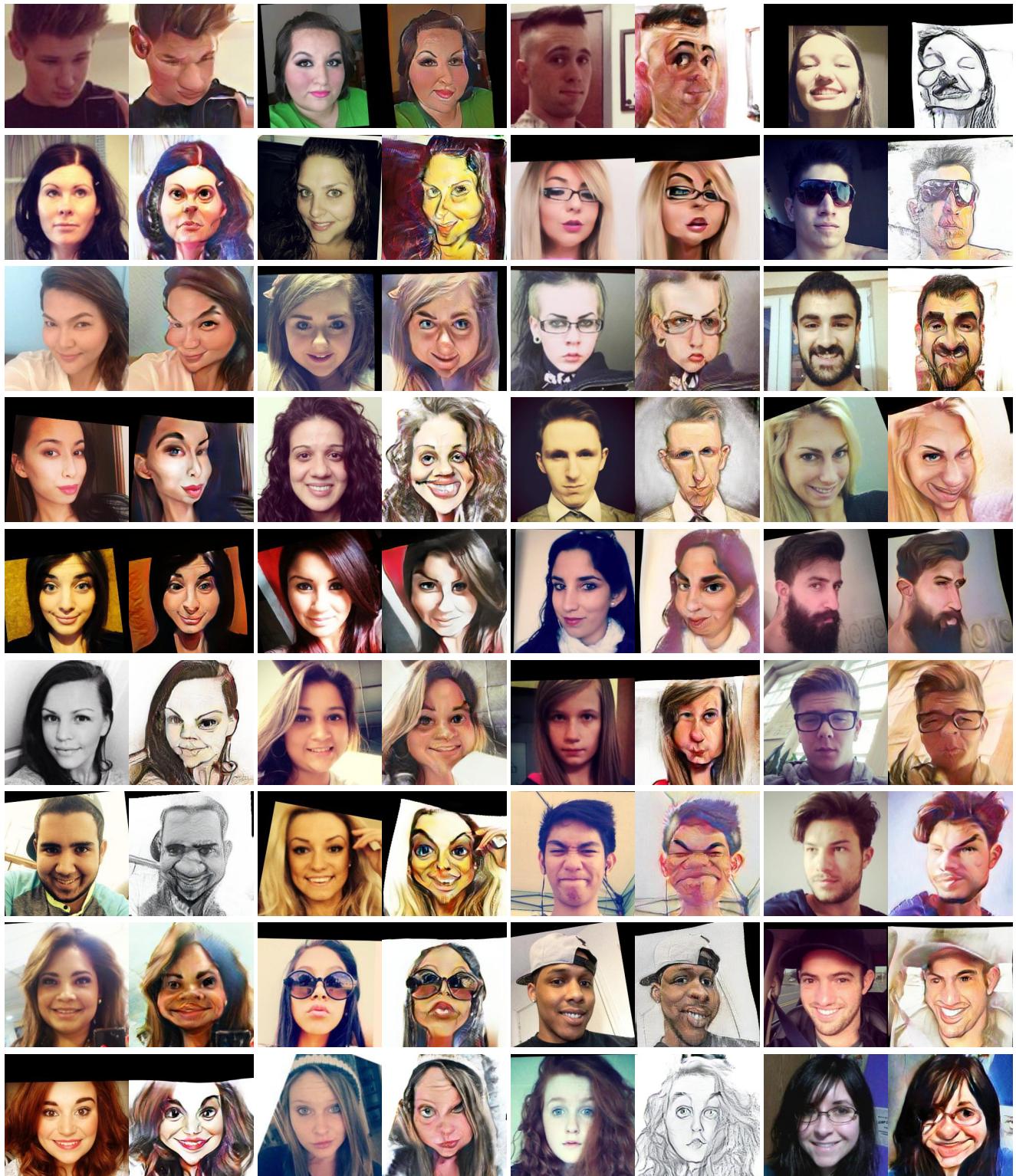
Figure 14: Example results on the Selfie dataset. This is a cross-dataset evaluation and no training is involved. In each pair, the left image is the input and the right image is the output of WarpGAN with a random texture style.