# StyleCariGAN: Caricature Generation via StyleGAN Feature Map Modulation

WONJONG JANG, POSTECH, Republic of Korea
GWANGJIN JU, POSTECH, Republic of Korea
YUCHEOL JUNG, POSTECH, Republic of Korea
JIAOLONG YANG, Microsoft Research Asia, China
XIN TONG, Microsoft Research Asia, China
SEUNGYONG LEE, POSTECH, Republic of Korea

|  |  |  |  |  |  |
|---|---|---|---|---|---|
| (a) Photo | (b) CariGANs | (c) WarpGAN | (d) AutoToon | (e) Ours | (f) Ours + InterFaceGAN |

Fig. 1. *Comparison of previous caricature generation methods with ours.* Different state-of-the-art caricature generation methods are compared to our proposed StyleCariGAN; CariGANs [Cao et al. 2018] (b); WarpGAN [Shi et al. 2019] (c); AutoToon [Gong et al. 2020] (d). Our method (e) can be used with other StyleGAN-based image editing methods, e.g., InterFaceGAN [Shen et al. 2020] (f). Photos: ©John Roaux/AP Photo, [Chen et al. 2002], ©Kevin Winter/Getty Images.

Authors' addresses: Wonjong Jang, POSTECH, Pohang, Republic of Korea, wonjong@postech.ac.kr; Gwangjin Ju, POSTECH, Pohang, Republic of Korea, gwangjin@postech.ac.kr; Yucheol Jung, POSTECH, Pohang, Republic of Korea, ycjung@postech.ac.kr; Jiaolong Yang, Microsoft Research Asia, Beijing, China, jiaoyan@microsoft.com; Xin Tong, Microsoft Research Asia, Beijing, China, xtong@microsoft.com; Seungyong Lee, POSTECH, Pohang, Republic of Korea, leesy@postech.ac.kr.

We present a caricature generation framework based on shape and style manipulation using StyleGAN. Our framework, dubbed *StyleCariGAN*, automatically creates a realistic and detailed caricature from an input photo with optional controls on shape exaggeration degree and color stylization type. The key component of our method is shape exaggeration blocks that are used for modulating coarse layer feature maps of StyleGAN to produce desirable caricature shape exaggerations. We first build a layer-mixed StyleGAN for photo-to-caricature style conversion by swapping fine layers of the StyleGAN for photos to the corresponding layers of the StyleGAN trained to generate caricatures. Given an input photo, the layer-mixed model produces detailed color stylization for a caricature but without shape exaggerations. We then append shape exaggeration blocks to the coarse layers of the layer-mixed model and train the blocks to create shape exaggerations while preserving the characteristic appearances of the input. Experimental results show that our StyleCariGAN generates realistic and detailed caricatures compared to the current state-of-the-art methods. We demonstrate

StyleCariGAN also supports other StyleGAN-based image manipulations, such as facial expression control.

CCS Concepts: • **Computing methodologies** → *Image processing*.

Additional Key Words and Phrases: 2D caricature, StyleGAN, shape exaggeration block, layer-swapping

## 1 INTRODUCTION

A caricature is a type of portrait wherein artists exaggerate the most recognizable characteristics of their subjects while oversimplifying other characteristics. To draw a caricature, artists have to learn how to not only paint in cartoon style but also capture and exaggerate the most salient facial features, which requires long time of training for development. Even after developing the skills, artists can take hours or even days to draw a single piece of caricature. Automatic computational strategies for photo-to-caricature translation can reduce the training and production burdens and make caricatures readily available for the general public.

Early studies in caricature generation [Akleman 1997; Akleman et al. 2000; Gooch et al. 2004] suggested methods to capture and deform facial geometric features, but they relied on user interaction to produce impressive results. Several automatic caricature generation methods were proposed, but their artistic styles were limited to pre-defined rules [Brennan 1985; Le et al. 2011; Mo et al. 2004].

Recent deep learning algorithms for image-to-image translation can find a mapping from an input domain to an output domain given training examples, but giving a proper supervision to learn automatic photo-to-caricature translation is not trivial. A straight-forward way would be to collect photo-caricature pairs, but the challenge is that photos and caricatures only have weak correspondences. For example, the caricatures corresponding to one person can vary in pose, expression, and exaggeration style. Constructing pairs using only the identity matching would make a sparse dataset for training photo-to-caricature translation.

Unpaired image-to-image translation algorithms such as Cycle-GAN [Zhu et al. 2017] can be trained without explicit pairing between input and output training examples. However, these unpaired image-to-image translation approaches still have a challenge when crossing domains using a sparse dataset such as photo and caricature dataset. Simply using a generic unpaired image-to-image translation approach does not result in realistic and visually pleasing results, as demonstrated in [Cao et al. 2018]. It seems the supervision between the photo domain and the caricature domain must be specifically crafted because of large shape variations between the two domains.

Recent automatic caricature algorithms design specialized network architectures for photo-to-caricature translation by separating geometric deformation and texture stylization [Cao et al. 2018; Gong et al. 2020; Shi et al. 2019]. However, the geometric deformation modules are based on 2D image warping and usually rely on interpolation of sparse control points [Cao et al. 2018; Shi et al. 2019],

resulting in loss of detailed shape deformations present in real caricatures. Collecting a small amount of dense and detailed caricature shape deformation examples was studied [Gong et al. 2020], but collecting a large amount of such dense data would be laborious.

To create realistic and detailed deformations, we take a generative approach that produces shape variations at multiple scales by modulating the corresponding feature maps in the network. By moving away from 2D image warping, we can handle a wide range of deformations from changes of overall facial shape to addition of wrinkles. Specifically, in this paper, we use StyleGAN [Karras et al. 2019] to manipulate a rich space of shape deformations in an unsupervised way. StyleGAN generates intermediate feature maps of multiple scales, where each scale represents a different spatial scale of the modeled image. It was showcased that StyleGAN can separately control over coarse, medium, and fine scale facial attributes.

One challenge for using StyleGAN as a photo-to-caricature translator is that it models the distribution of only one domain. We can train two StyleGANs separately for facial photos and caricature images. Then, as our translator must take a photo as the starting point, we need a way to connect manipulated photo feature maps to caricature styles in the final image. We resolve the problem by simply swapping layers of two StyleGANs, as in [Pinkney and Adler 2020]. Our layer-mixed StyleGAN uses the first four layers of the photo StyleGAN and the last three layers of the caricature Style-GAN. Another challenge for using StyleGAN to build a system that takes a photo as the input is the encoding of the input photo into a StyleGAN latent vector. We combine two previous ideas on GAN inversion [Karras et al. 2020b; Tewari et al. 2020] to obtain a latent code optimized for the input photo.

Given the latent code of an input photo, our layer-mixed Style-GAN can create an output image with caricature color styles. However, the output would not contain enough shape exaggerations as the facial shape is dominated by the first four layers copied from the photo StyleGAN. We then append *shape exaggeration blocks* to the four layers of the layer-mixed StyleGAN to handle shape deformations with feature map modulation. Our shape exaggeration blocks compute desirable feature map changes to be added to the original feature maps, and provides multi-scale control for shape exaggerations.

There are desirable properties for shape exaggeration blocks: 1) The resulting deformations should generate an image that looks like a realistic caricature. 2) The deformations should preserve unique visual features of the input photo. The first property can be achieved by adversarial training using a caricature dataset. However, the second property is not trivial to achieve. To design a constraint on preservation of important features, we need a tool to extract a set of visual features shared between photos and caricatures. The extraction is challenging because of large shape and texture gaps between the two domains.

Previous photo-to-caricature translators based on deep learning solved the problem of preservation of visual features in different ways. Cao *et al.* [2018] designed a *characteristic loss* to constrain their 2D image warping using the cosine similarity between input and output landmark displacements, and the loss was used along with a cycle consistency loss between a photo and a caricature. Shi *et al.* [2019] designed an *identity-preserving adversarial loss* that

makes the discriminator for the adversarial training also perform face identification between a photo and a caricature.

Our approach to enable the shape exaggeration blocks to preserve the visual features of an input image is based on two key ideas. 1) We establish weak correspondence between photos and caricatures in an unsupervised way with *cycle consistency*, as in [Zhu et al. 2017] and [Cao et al. 2018]. For the cycle consistency, we create two models for photo-to-caricature translation and caricature-to-photo translation. 2) We guide the shape exaggeration blocks with a direct supervision for preservation of facial attributes using a recently released dataset, WebCariA [Ji et al. 2020], which provides a shared set of annotations between photos and caricatures. WebCariA provides 50 label annotations based on intrinsic facial attributes. The annotations include shape attribute labels of overall facial shapes, nose shapes, eye shapes, and mouth shapes, i.e., *Wide Nose, Mustache, Big Eyes, Arched Eyebrows, etc*, and are provided both for photo and caricature images. Using the dataset, we train two attribute classifiers for photos and caricatures separately. We then design an *attribute matching loss* that constrains the attributes detected in the input photo to be preserved in the output caricature.

To summarize, our main contributions are as follows:

- We propose a novel *StyleCariGAN* framework for photo-to-caricature translation. Based on StyleGAN, it manipulates feature maps at multiple scales to produce desired shape exaggeration and color stylization.
- To create desirable shape exaggerations in photo-to-caricature translation, we propose shape exaggeration blocks that are appended to the coarse layers of a layer-mixed StyleGAN.
- We design an attribute matching loss that is used along with cycle consistency to constrain the shape exaggeration blocks to preserve important visual features of the input.
- Our StyleCariGAN framework generates realistic and detailed caricatures compared to state-of-the-art methods, supporting other image manipulations such as facial expression control.

## 2 RELATED WORK

### 2.1 Rule-based automatic caricature generation

Rule-based caricature generation methods create caricatures by transforming an input image with a fixed sequence of operations. The seminal work of Brennan [1985] suggested an automatic algorithm based on exaggeration of relative feature point differences between an input face and a template mean face. This idea was further explored in many other algorithms. Mo *et al.* [2004] extended it by considering the variances of different feature point positions. Liao *et al.* [2004] exaggerated facial feature points constructed in a tree structure. Le *et al.* [2011] warped an input face based on exaggeration of anthropometric ratios between facial components.

Although these rule-based methods create plausible caricature deformations based on carefully designed procedures, the procedurally deformed facial shape may not resemble a realistic and artist-drawn caricature. The designed procedures cannot cover a wide range of caricature deformations performed by human artists. A natural solution to this problem would be modeling the mapping from a photo to a caricature in an end-to-end fashion, as done in recent deep learning based image translation methods.

### 2.2 Deep image-to-image translation

With the advance of deep convolutional neural networks, image-to-image translation methods have been developed to convert an image in the source domain to the target domain. Given paired images from the two domains as training data, Isola *et al.* [2017] and Wang *et al.* [2018] showcase end-to-end training schemes for image translation based on adversarial training [Goodfellow et al. 2014]. These algorithms work well when input-output pairs are well-defined. Using unpaired data for image-to-image translation learning has also been actively studied. CycleGAN [Zhu et al. 2017] learns the mapping between two image domains without paired correspondence using a cycle consistency constraint. U-GAT-IT [Kim et al. 2020] learns the mapping using attention and a new feature normalization method. MUNIT [Huang et al. 2018] considers the fact that there can be multiple outputs possible for a given input and proposes a multimodal translation framework.

The difficulty of applying image-to-image translation to caricature generation comes from lack of abundant identity-matched photo-caricature pairs. Moreover, the two domains of real photos and caricatures differ at not only image style, but also face shapes. Naively applying unpaired translation would result in inferior caricature quality and less identity preservation, as shown in [Cao et al. 2018] . In this paper, we do not explicitly pair photos and caricatures, but supervise the correspondence between the two with constraint on facial attribute preservation. In addition, we make the mapping be learned more easily with a decomposition in caricature image formation: shape manipulation and appearance stylization.

### 2.3 Deep caricature generation

Automatic caricature generation algorithms based on deep learning has been studied recently. Cao *et al.* [2018] proposed CariGANs for photo-to-caricature translation using two networks for geometric deformation and style transfer. They constrain the output caricature to preserve visual features of the input shape by applying a cycle consistency and a cosine similarity constraint to landmark deformations. For more flexible shape variation, WarpGAN [Shi et al. 2019] trained an image warping module with dynamically decided control points using deep neural networks. WarpGAN created caricatures that resemble the input by making the discriminator used for the adversarial learning perform identity classification. AutoToon [Gong et al. 2020] adopted dense deformation fields, instead of using coarse control points, and trained the caricature generator using a small amount of paired dataset comprising photos and corresponding caricatures generated by artists using 2D deformation tools. CariGAN [Li et al. 2020] proposed an image fusion mechanism to encourage the caricature generator to focus on important local regions around sparse facial landmarks. 3D caricature algorithms can be used to create a 2D caricature by warping the input image based on the created 3D shape. Han *et al.* [2018] create 3D caricatures with deep neural networks by generating and manipulating a rough sketch of an exaggerated 3D facial shape. The final 2D caricature image is generated by blending a warped image based on the 3D warping and the original image.

Our main difference to previous deep-learning-based caricature generation methods is that we take a generative approach for caricature exaggerations by modeling spatially dense and detailed shape

Table 1. *Comparison of state-of-the-art deep learning based 2D automatic caricature generation methods and ours.*

| | CariGANs [Cao et al. 2018] | WarpGAN [Shi et al. 2019] | AutoToon [Gong et al. 2020] | Ours |
|---|---|---|---|---|
| Shape manipulation type | 2D spatial interpolation of landmark deformations | 2D spatial interpolation of control point deformations | 2D spatial dense displacement map | StyleGAN feature map modulation |
| Shape manipulation spatial density | Pre-defined 63 landmarks detected from input image | Dynamically decided 16 control points | Dense | Dense |
| Training dataset type | Unpaired large (10K photos + <10K caricatures) | Unpaired large (6K photos + <10K caricatures) | Paired small (101 photos + 101 caricatures) | Unpaired large (70K photos + <10K caricatures) |
| User control | single exaggeration magnitude control | single exaggeration magnitude control | single exaggeration magnitude control | 4-scale exaggeration magnitude control |

deformations using a large amount of unpaired photo-caricature data. We create realistic and detailed caricatures by obtaining desirable shape deformations through modulations of StyleGAN feature maps. We also provide a multi-scale control over shape exaggeration magnitude in the caricature output. Table 1 summarizes the differences among caricature generation methods.

## 2.4 Image generation and editing using StyleGAN

StyleGAN [Karras et al. 2019, 2020b] is a powerful generative image model, which can synthesize high-resolution face images that are even hard to distinguish from real photos. It is shown that the feature maps at different resolutions of the StyleGAN architecture characterize the image styles at different scales and control the generation in a coarse-to-fine fashion. With StyleGAN, image manipulation can be achieved by embedding a real photo into the StyleGAN latent space and editing the embedded latent code for re-generation [Abdal et al. 2019, 2020; Deng et al. 2020; Karras et al. 2020b; Tewari et al. 2020; Zhu et al. 2020]. Semantically-meaningful editing in GAN latent space has also been studied [Bau et al. 2019; Brock et al. 2016; Creswell and Bharath 2018; Richardson et al. 2020; Yeh et al. 2017; Zhu et al. 2016]. Recent studies based on StyleGAN achieve such editing via either supervised [Abdal et al. 2021; Shen et al. 2020] or unsupervised [Shen and Zhou 2020] learning.

*Toonification* [Pinkney and Adler 2020] creates an image with cartoon structure and photo-realistic rendering using StyleGAN. It swaps layers of two StyleGANs trained for photos and caricatures. The StyleGAN for photos is trained from scratch and the StyleGAN for caricatures is fine-tuned from the photo model. It uses the fine-tuned caricature StyleGAN for coarse layers to create cartoon structures, and the photo StyleGAN for fine layers to create photo-realistic textures. To generate a *toonified* image from an input photo, separately generated cartoon structures are blended into the input photo feature maps. FreezeG [bryandlee 2020] achieves similar pseudo image translation through training instead of layer-swapping. It freezes a certain range of layers in a StyleGAN for

one domain, and fine-tunes the StyleGAN to generate images of another domain. However, in both Toonification and FreezeG, it is unclear how to explicitly supervise the translations to preserve visual features of input images.

Our automatic caricature system generates an image with exaggerated photo structure and cartoon-style rendering. The exaggerated photo structure is created from the input, preserving important visual features. We handle the problem of crossing the domains from photo to caricature by swapping layers of two StyleGANs, similarly to Toonification [Pinkney and Adler 2020]. However, differently from Toonification, the coarse layers are copied from the photo StyleGAN and the fine layers are from the caricature StyleGAN in our method. Moreover, we append learnable shape exaggeration blocks to the copied coarse layers and design a supervision for the blocks to preserve characteristic attributes of the input. Comparison between our method and Toonification is presented in Sec. 4.5.

## 3 PHOTO-TO-CARICATURE TRANSLATION FRAMEWORK

Our StyleCariGAN framework contains two fixed StyleGANs: one is trained to generate photos, and the other is fine-tuned from the photo model to generate caricature images (Fig. 2). After pre-training, the two models are fixed and not updated during our process. We refer to the StyleGAN trained to generate regular face images as the *photo StyleGAN* and another trained for caricature images as the *caricature StyleGAN*. By swapping layers of the two StyleGANs, we make a layer-mixed model that generates a caricature image starting from a latent code.

In the layer-mixed StyleGAN model, the coarse (low-resolution) layers are from photo StyleGAN to control the overall structural and identity information of the output caricature. The fine (high-resolution) layers are from caricature StyleGAN to create detailed color styles for a caricature. However, naively copying the coarse layers from photo StyleGAN would produce ordinary face shapes without exaggerations.

Fig. 2. *Image samples generated using the fine-tuned StyleGAN for caricatures.*[1] To sample the images, we apply the truncation trick [Karras et al. 2019, 2020b] with $\psi = 0.7$.



Fig. 3. *Overview of our approach for caricature generation using a layer-mixed StyleGAN with shape exaggeration blocks.* A $\mathcal{W}+$ code is first obtained by embedding an input to the latent space of photo StyleGAN. Then, the first four coarse feature maps generated by the layers copied from photo StyleGAN using the $\mathcal{W}+$ code are modulated with shape exaggeration blocks. Finally, the modulated features are fed to the fine layers copied from caricature StyleGAN to produce the output caricature.

To impose shape exaggeration capability on the layer-mixed model, we append *shape exaggeration blocks* to the coarse layers. Our key idea is to train the shape exaggeration blocks with a specially crafted supervision for photo-to-caricature translation. We train the blocks to introduce feature map modulations that create realistic and identity-preserving caricature exaggerations. Shape exaggeration blocks implicitly increase the deviation from the average facial geometry in feature space because our model is trained to generate highly diverse caricature geometry with GAN loss. We design cycle consistency and attribute matching losses to guide the deviation to be distinctively reflecting the features of the input.

The starting point of photo-to-caricature translation, *i.e.*, the input to the layer-mixed model, is obtained from the input photo using photo-to-latent embedding. We obtain the $\mathcal{W}+$ latent code [Abdal et al. 2019] representing the input image through iterative optimization, akin to [Tewari et al. 2020] and [Karras et al. 2020b]. From the $\mathcal{W}+$ code, our layer-mixed model generates multi-scale feature maps and eventually the output caricature image.

Fig. 3 shows the overall process of our photo-to-caricature translation. The following subsections explain each step of the process in detail.

### 3.1 Layer swapping between two StyleGANs

For the layer-mixed model, we copy the parameters of coarse and fine layers from photo and caricature StyleGANs, respectively. When training caricature StyleGAN with fine-tuning starting from photo StyleGAN, we align caricatures and photos based on facial landmarks. The positions of important visual features are shared between photos and caricatures after the alignment. This spatial alignment enables layer swapping to produce a plausible mixture of a photo and a caricature (Fig. 4).

---

[1] Since we added zero-padding on borders of training images after aligning the caricature dataset with facial landmarks, black borders sometimes appear in the results.
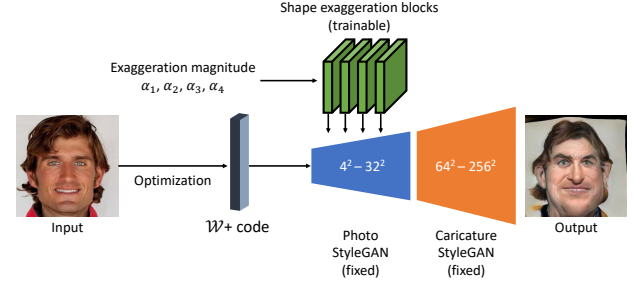
In layer swapping, the choice of the boundary between coarse and fine layers is important. Different choices will lead to different levels of input structure preservation and detail stylization, as shown in Fig. 4. When selecting the boundary, our goal is to enable the layer-mixed model to generate caricature images faithful to the input photo in terms of face structure and identity. We empirically choose the first four scales ($4^2 \sim 32^2$) as the coarse layers and the rest as the fine layers. Fig. 4 shows that taking more than three fine layers from caricature StyleGAN would introduce excessive shape deformations that may change the visual appearance of the input photo.

The fine layers of our StyleCariGAN take style latent codes that can be selected by a user to control detail styles. For convenient control, we construct a curated set of style latent codes by selecting a number of good examples from randomly generated styles.

The copied parameters of the layers are not updated further after layer swapping. That is, we do not apply additional end-to-end fine-tuning to the layer-mixed model, as the coarse and fine layers of the model already have the desired properties of handling facial shapes and detailed styles, respectively. The only trained components in the layer-mixed model are shape exaggeration blocks appended to the coarse layers.

To train shape exaggeration blocks, we define another layer-mixed model for caricature-to-photo translation to enable cycled training. Specifically, we construct this new layer-mixed model by taking the first four coarse layers from caricature StyleGAN and the remaining fine layers from photo StyleGAN. We denote this new layer-mixed model as the *c2p* model, while the original one, which is the backbone for our StyleCariGAN, is denoted as the *p2c* model.

### 3.2 Shape exaggeration blocks

Our key component for photo-to-caricature translation is the shape exaggeration blocks. Four shape exaggeration blocks are attached to the p2c model, and we call the modified p2c model *p2c-StyleCariGAN*, or simply *StyleCariGAN*. We create *c2p-StyleCariGAN* similarly, then use the model for enforcing cycle consistency.

Our shape exaggeration blocks are convolutional layers that produce additive feature modulation maps for the coarse layers of StyleCariGAN. A shape exaggeration block takes a feature map of
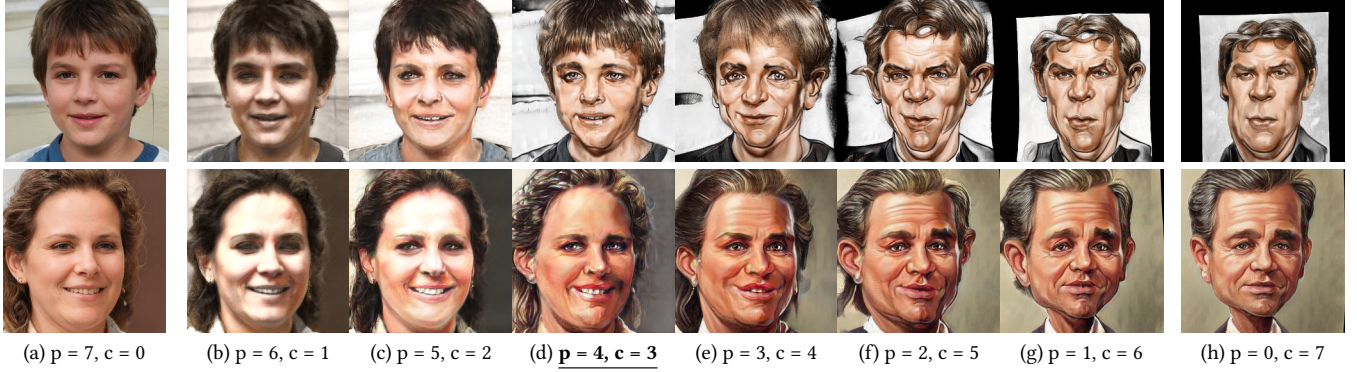
| (a) p = 7, c = 0 | (b) p = 6, c = 1 | (c) p = 5, c = 2 | (d) **p = 4, c = 3** | (e) p = 3, c = 4 | (f) p = 2, c = 5 | (g) p = 1, c = 6 | (h) p = 0, c = 7 |

Fig. 4. *Layer swapping examples.* $p$ denotes the number of coarse photo-StyleGAN layers used in the layer-mixed model, and $c$ denotes the number of fine caricature-StyleGAN layers. For instance, $p = 4, c = 3$ means that the first four coarse layers are from the photo StyleGAN and the last three fine layers are from the caricature StyleGAN. As $p$ decreases, the face structure of the input is deformed and the identity is lost. We choose $p = 4, c = 3$ for our layer-mixed model that retains the shape information in the input while producing enough stylization of details.

size $n^2$ and creates a $n^2$ feature modulation map which is added back to the input feature map, akin to residual learning [He et al. 2016a]. The modulated feature map is then fed to the next layer that creates a higher resolution feature map to handle finer-scale structures. The architecture of a shape exaggeration block of size $n^2$ is the same as that of a convolution layer of size $n^2$ in the original StyleGAN [Karras et al. 2019], except that we exclude the upsampling block.

We train the shape exaggeration blocks to achieve two goals. The first goal is to introduce shape deformations that resemble real caricatures, which is achieved by two types of generative adversarial losses: $\mathcal{L}_{feat}$ and $\mathcal{L}_{GAN}$. The second goal is to preserve important visual features from the input photo, which is achieved by three losses: feature map cycle consistency loss $\mathcal{L}_{fcyc}$, identity cycle consistency loss $\mathcal{L}_{icyc}$, and attribute matching loss $\mathcal{L}_{attr}$ (Fig. 5).

*Generative Adversarial Loss.* We use the same generative adversarial loss and supervision method as StyleGAN2 [Karras et al. 2020b]. For the *real* data needed for the supervision of discriminators, we use generated examples from the caricature StyleGAN, instead of real caricatures. This strategy helps the training process by providing a diverse set of samples and avoiding image loading overhead. One consequence is that our StyleCariGAN would produce caricatures similar to caricature StyleGAN. However, the caricature StyleGAN can create high-quality caricatures as shown in Fig. 2, supporting the strategy. We apply adversarial learning to both the modulated feature maps and final images by introducing two discriminators. We use the non-saturating logistic loss [Goodfellow et al. 2014] with $R_1$ regularization [Mescheder et al. 2018] for both, and calculate the sum of the two losses as our total adversarial loss $\mathcal{L}_{adv}$:

$$\mathcal{L}_{adv} = \mathcal{L}_{feat} + \lambda_{GAN}\mathcal{L}_{GAN}, \tag{1}$$

where $\mathcal{L}_{feat}$ and $\mathcal{L}_{GAN}$ are the losses on feature maps and final images, respectively.

*Cycle consistency Loss.* To guide the shape exaggeration blocks to build a correspondence between photos and caricatures, we constrain the blocks to be cycle-consistent. To implement cycle consistency, we utilize the symmetric set of StyleCariGANs: p2c-StyleCariGAN and c2p-StyleCariGAN. We impose cycle consistency

both on intermediate modulated feature maps and generated images. We refer to the cycle consistencies for the modulated feature maps and images as feature map cycle consistency $\mathcal{L}_{fcyc}$ and identity cycle consistency $\mathcal{L}_{icyc}$, respectively.

The feature map consistency forces the effect of shape exaggeration blocks to be invertible with a cycle at each feature map scale. The shape exaggeration blocks for p2c-StyleCariGAN define photo-to-caricature feature modulation $\mathcal{S}_{p \to c}$. The corresponding blocks in c2p-StyleCariGAN define caricature-to-photo feature modulation $\mathcal{S}_{c \to p}$. That is, the blocks in c2p-StyleCariGAN does inverse of exaggeration, reverting an exaggerated shape into a regular one. We call these blocks *shape de-exaggeration blocks*. Using the two feature modulations, the cycle consistency loss $\mathcal{L}_{fcyc}$ is defined as:

$$\mathcal{L}_{fcyc}^{p \to c} = \sum_{i=1}^{4} (\mathbb{E}_{\mathcal{F}_p^i \sim G_p^i(w)} [\|\mathcal{S}_{c \to p}^i(\mathcal{S}_{p \to c}^i(\mathcal{F}_p^i)) - \mathcal{F}_p^i\|_2]),$$

$$\mathcal{L}_{fcyc}^{c \to p} = \sum_{i=1}^{4} (\mathbb{E}_{\mathcal{F}_c^i \sim G_c^i(w)} [\|\mathcal{S}_{p \to c}^i(\mathcal{S}_{c \to p}^i(\mathcal{F}_c^i)) - \mathcal{F}_c^i\|_2]),$$

$$\mathcal{L}_{fcyc} = \mathcal{L}_{fcyc}^{p \to c} + \mathcal{L}_{fcyc}^{c \to p}, \tag{2}$$

where $F_c^i$ is a caricature feature map generated by the $i$-th block $G_c^i$ of the caricature StyleGAN, $F_p^i$ is a photo feature map of the $i$-th block $G_p^i$ of the photo StyleGAN, $S_{p \to c}^i$ is the $i$-th photo-to-caricature shape exaggeration block, and $S_{c \to p}^i$ is the $i$-th caricature-to-photo shape de-exaggeration block. Fig. 6 illustrates the losses.

The identity consistency forces the effect of shape exaggeration blocks to be invertible with a cycle by inspecting photos generated from feature maps. We design the identity consistency loss $\mathcal{L}_{icyc}$ based on a face embedding network trained on photos [Schroff et al. 2015], so the identity consistency is calculated only for the cycle starting from photo feature maps. Using the face embedding network, we design the identity cycle consistency loss $L_{icyc}$ as the $L_2$ distance between two embeddings: 1) the face embedding of a starting photo and 2) the face embedding of the cycle-reconstructed photo, which is created using coarse feature maps modulated by $S_{p \to c}$ and $S_{c \to p}$ sequentially.
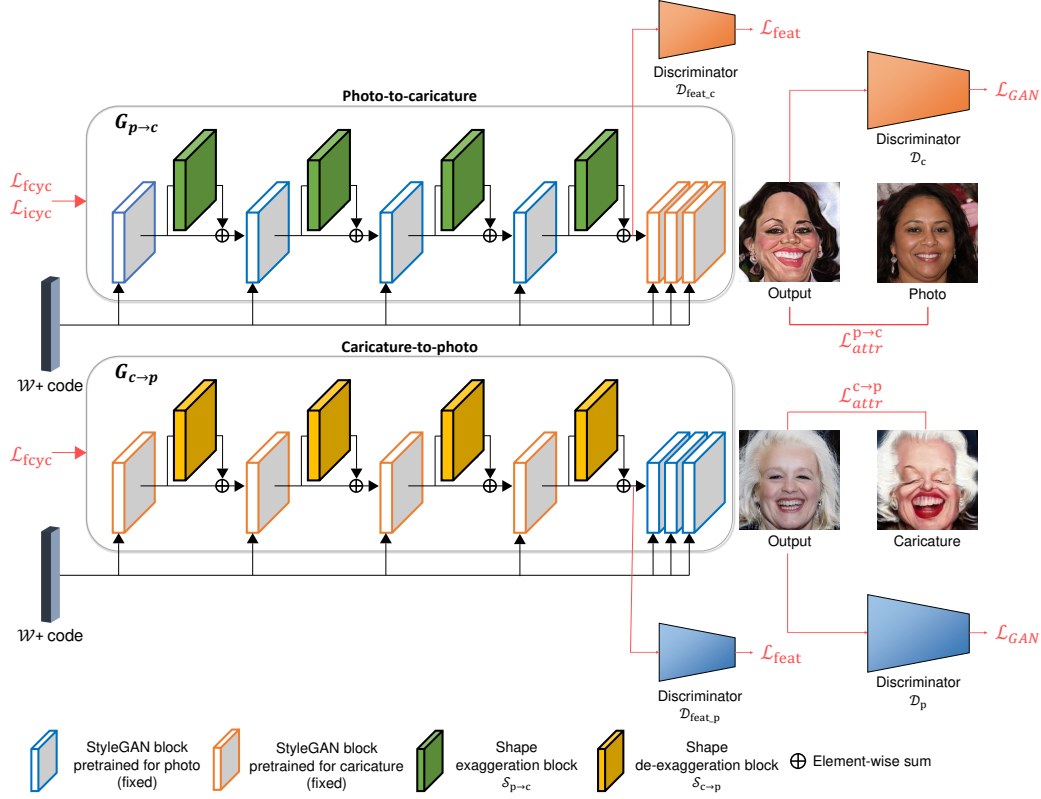
Fig. 5. *Training shape exaggeration blocks.* With fixed layers copied from photo StyleGAN and caricature StyleGAN, we train two sets of shape exaggeration blocks: $\delta_{p \to c}$ in p2c-StyleCariGAN and $\delta_{c \to p}$ in c2p-StyleCariGAN. To achieve realistic and facial-attribute-preserving exaggerations, shape exaggeration blocks are supervised with the adversarial losses ($\mathcal{L}_{feat}$, $\mathcal{L}_{GAN}$), the cycle consistency losses ($\mathcal{L}_{fcyc}$, $\mathcal{L}_{icyc}$), and the attribute matching loss ($\mathcal{L}_{attr}$).
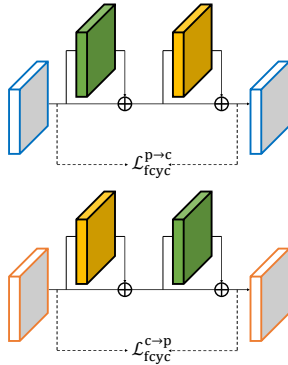


Fig. 6. *Features map cycle consistency.* We enforce cycle consistency between the feature maps of coarse layers in p2c-StyleCariGAN and c2p-StyleCariGAN.

Note that our cycle consistency loss does not require a caricature-to-feature-map encoder because our photo-to-caricature and caricature-to-photo mappings are feature-map-to-feature-map mappings. The mappings do not take images as the input, and we can directly use feature maps to implement the cycle consistency loss.

Finally, the cycle consistency loss $\mathcal{L}_{cyc}$ is computed as

$$\mathcal{L}_{cyc} = \mathcal{L}_{fcyc} + \lambda_{icyc} \mathcal{L}_{icyc}. \tag{3}$$

*Attribute Matching Loss.* Even with the cycle consistency, the resultant exaggeration is not guaranteed to preserve important features of the input photo. To constrain the shape exaggeration blocks to produce valid caricature deformations, we use facial attribute classifiers for photos and caricatures. The facial attribute classifiers are trained on an attribute-labeled dataset, WebCariA [Ji et al. 2020]. The dataset annotates 50 attributes that describe intrinsic facial shape features such as the sizes of facial components or the overall shape of the face. These attributes are not affected by extrinsic factors such as cosmetics or poses. The annotation is done both for photos and caricatures using the same set of annotation labels.

We constrain p2c-StyleCariGAN to create a caricature with the same attributes as the input photo. We apply similar attribute matching to c2p-StyleCariGAN. The attribute matching loss $\mathcal{L}_{attr}$ is defined using binary cross entropy losses between photo attributes and caricature attributes:

$$\mathcal{L}_{attr}^{p \to c} = -\mathbb{E}_{w \sim \mathcal{W}}[\phi_p(G_p(w)) \log \phi_c(G_{p \to c}(w))$$
$$+ (1 - \phi_p(G_p(w))) \log(1 - \phi_c(G_{p \to c}(w)))],$$

$$\mathcal{L}_{attr}^{c \to p} = -\mathbb{E}_{w \sim \mathcal{W}}[\phi_c(G_c(w)) \log \phi_p(G_{c \to p}(w))$$
$$+ (1 - \phi_c(G_c(w))) \log(1 - \phi_p(G_{c \to p}(w)))],$$

$$\mathcal{L}_{attr} = \mathcal{L}_{attr}^{p \to c} + \mathcal{L}_{attr}^{c \to p}, \tag{4}$$

(a) input      (b) $\mathcal{L}_{adv}$      (c) $\mathcal{L}_{adv} + \mathcal{L}_{cyc}$      (d) $\mathcal{L}_{adv} + \mathcal{L}_{attr}$      (e) $\mathcal{L}_{adv} + \mathcal{L}_{cyc} + \mathcal{L}_{attr}$
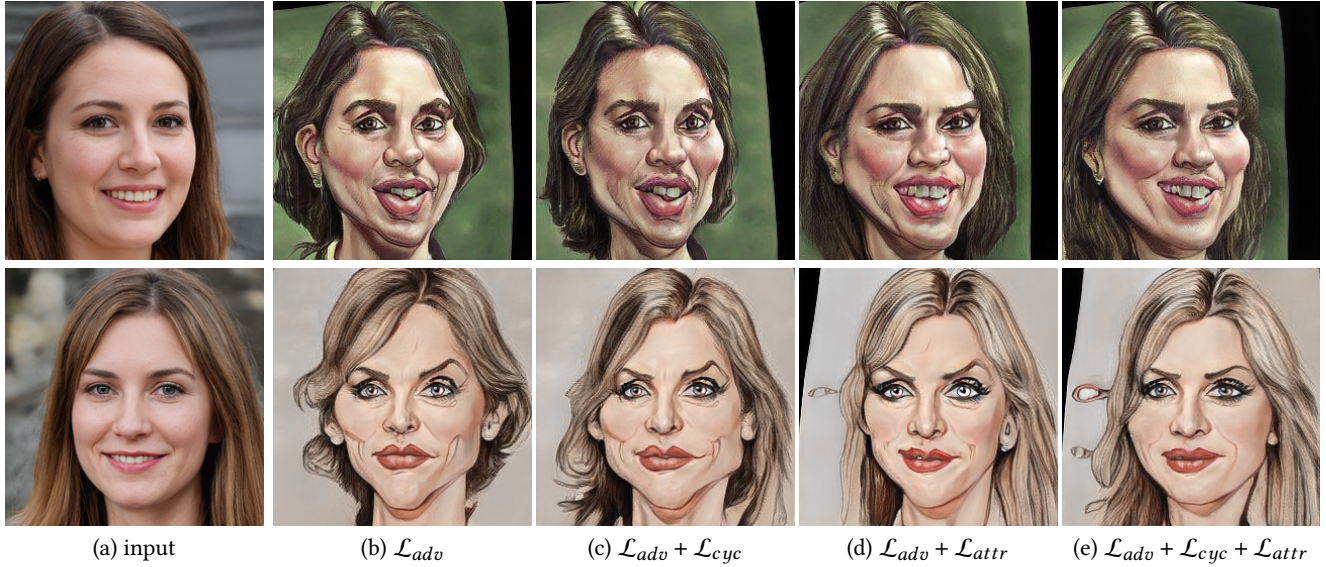
Fig. 7. *Effect of each loss used for training shape exaggeration blocks.* The input image (a) is transformed to a cartoon-style shape with the effect of adversarial loss $\mathcal{L}_{adv}$ (b). The cartoon-style deformation introduced by $\mathcal{L}_{adv}$ is random and does not respect the input shape. Adding only one of cycle consistency loss $\mathcal{L}_{cyc}$ (c) or attribute matching loss $\mathcal{L}_{attr}$ (d) changes the output to have a more similar facial contour and hair to the input, but the results often show excessive wrinkles or deformed eye shapes. Combining the attribute matching loss $\mathcal{L}_{attr}$ together with the cycle consistency loss $\mathcal{L}_{cyc}$ (e), we obtain caricatures that preserve important features of the input without artifacts.

where $\phi_p$ is a photo attribute classifier, $\phi_c$ is a caricature attribute classifier, $G_p$ is the photo StyleGAN, $G_c$ is the caricature StyleGAN, $G_{p \to c}$ is p2c-StyleCariGAN, and $G_{c \to p}$ is c2p-StyleCariGAN.

To summarize, our full objective function for training is as follows:

$$\mathcal{L}_G = \lambda_{adv}\mathcal{L}_{adv} + \lambda_{cyc}\mathcal{L}_{cyc} + \lambda_{attr}\mathcal{L}_{attr}, \qquad (5)$$

where $\lambda_{adv}$, $\lambda_{cyc}$, and $\lambda_{attr}$ are constants defining the loss weights. The effects of different loss components are visualized in Fig. 7, and more examples can be found in the supplementary material.

### 3.3 Exaggeration degree control

Our shape exaggeration blocks enable multi-scale deformation control over the generated images. Since these blocks perform additive modulation on four feature maps representing different spatial scales, we can freely change the amount of modulation by simply attaching a scaling factor to the output of each block, which alters the degree of exaggeration. For example, to reduce the deformation of the overall face shape, we can multiply a weight less than 1 to the output of the first shape exaggeration block before it is added to the original feature map. To remove wrinkles incurred by extreme deformations, we can multiply zero or a small scaling factor to the output of the third or the fourth shape exaggeration block.

### 3.4 Photo-to-latent embedding

To use a photo image for the input of our StyleCariGAN, we optimize the $\mathcal{W}+$ latent code of StyleGAN to reproduce the photo. Since there can be multiple latent codes that result in similar photos, the mapping from the input photo to the latent code is not unique.

However, not all of the possible latent codes are adequate for generating feature maps suitable for image editing. Therefore, we need to choose a good latent code that is meaningful for editing.

The most intuitive solution to obtain the desired latent code is simply applying an existing GAN inversion method. Some GAN inversion methods are good at reconstructing an input image [Abdal et al. 2019, 2020], but their embedded latent codes are often out-of-distribution [Zhu et al. 2020]. Therefore, we considered GAN inversion methods that work for image editing in the latent space [Karras et al. 2020b; Tewari et al. 2020].

Hierarchical optimization [Tewari et al. 2020] and StyleGAN2 inversion [Karras et al. 2020b] showed the most promising results in our task. However, hierarchical optimization, which first finds the latent code in $\mathcal{W}$ space and then updates it in $\mathcal{W}+$ space, is restricted to restoring high-frequency features in images since it does not consider noise optimization. StyleGAN2 inversion optimizes noise as well as latent code, but it may move away from the prior distribution of StyleGAN since it simply initializes the latent code with the mean latent code. To improve existing methods, we combine hierarchical optimization [Tewari et al. 2020] and noise optimization techniques [Karras et al. 2020b] for latent code embedding.

Specifically, we first find the $\mathcal{W}$ latent code that reconstructs an input facial image via optimization. Then, we optimize $\mathcal{W}+$ latent code and noise simultaneously by initializing the latent code with the obtained $\mathcal{W}$ latent code. We also apply latent code perturbation and noise regularization [Karras et al. 2020b]. As a result, we can obtain the desired latent code suitable for our framework, as illustrated in Fig. 8.
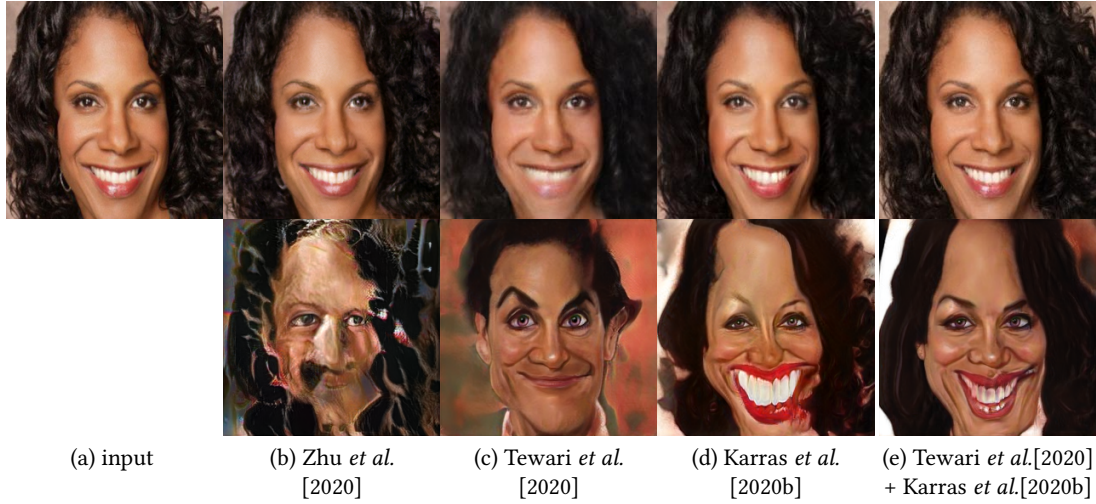
|                  |                  |                  |                  |                  |
|:---------------:|:---------------:|:---------------:|:---------------:|:-----------------:|
| (a) input | (b) Zhu *et al.* [2020] | (c) Tewari *et al.* [2020] | (d) Karras *et al.* [2020b] | (e) Tewari *et al.*[2020] + Karras *et al.*[2020b] |

Fig. 8. *Comparison of different methods for photo-to-latent embedding.* The input image (a) is embedded into a $\mathcal{W}+$ latent vector using different methods. The reproduced images using the embedded latent vectors are visualized in the top row. The caricatures generated from the embedded vectors using our StyleCariGAN are in the bottom row. All embedded vectors using different methods produce plausible reconstructions, but those from previous methods (b-d) do not result in satisfactory caricatures. In contrast, our approach (e) combining [Tewari et al. 2020] and [Karras et al. 2020b] obtains an embedded vector that achieves a more pleasing caricature. Input: ©Craig Sjodin/Getty Images.

## 4 EXPERIMENTS

In this section, we first describe the settings of our experiments. We then evaluate our StyleCariGAN[2] qualitatively and quantitatively and compare it with state-of-the-art methods.

*StyleGANs.* Our framework requires two StyleGANs; one generates photos, and the other generates caricatures. Both StyleGANs use the architecture and the training algorithm of StyleGAN2 [Karras et al. 2020b]. The StyleGAN for photos was trained with FFHQ dataset [Karras et al. 2019] resized to $256 \times 256$ resolution. The StyleGAN for caricatures was fine-tuned from the photo model with caricatures in WebCaricature [Huo et al. 2017] using FreezeD [Mo et al. 2020] and ADA [Karras et al. 2020a]. The WebCaricature dataset was aligned using five landmarks (the centers of the eyes, the tip of the nose, and the two corners of the mouth) and resized to $256 \times 256$ as well. We used a PyTorch implementation [rosinality 2020] for training both models.

*Face embedding network.* To employ the identity cycle consistency loss $\mathcal{L}_{icyc}$, a face embedding network is required. We use a pretrained FaceNet [Schroff et al. 2015] as the face embedding network.

*Attribute classifiers.* To implement $\mathcal{L}_{attr}$, we need two face attribute classifiers for photos and caricatures, respectively. We train the attribute classifiers using WebCariA dataset [Ji et al. 2020], which provides labels for both photos and caricatures. We found that the label distributions for photos and caricatures are reasonably similar in the dataset. The backbone architecture of the attribute classifiers is ResNet-18 [He et al. 2016b] with the only change in the output channel size of the last fully connected layer. Our output channel size is 50, which is the number of attributes in WebCariA. We finetuned the pre-trained ResNet-18 provided by PyTorch [Paszke et al.

---

[2]Our code is available at https://github.com/wonjongg/StyleCariGAN

2019]. The test accuracy on the test set of the WebCariA dataset was $85\%$ for photos and $82\%$ for caricatures.

*Training shape exaggeration blocks.* A shape exaggeration block consists of two convolutional layers, each with a leaky ReLU. Each convolutional layer has kernel size = 3, stride = 1, and padding = 1. The leaky ReLU layers have a negative slope 0.2.

We use the Adam optimizer [Kingma and Ba 2014] in PyTorch with $\beta_1 = 0$, $\beta_2 = 0.99$, and learning rate 0.002. Each mini-batch with a batch size 32 consists of a randomly generated photos and caricatures. We stopped training after $1,000$ iterations. We empirically set the weights for the losses as $\lambda_{adv} = 1$, $\lambda_{GAN} = 10$, $\lambda_{cyc} = 10$, $\lambda_{icyc} = 1000$, and $\lambda_{attribute} = 10$. The training time with 4 NVIDIA Quadro RTX 8000 (48 GB) GPUs was about 8 hours.

*Running time.* In test time, the GAN inversion from an input photo to the latent code takes three to four minutes. After the inversion, generating a caricature image takes about 40 ms. The measurement was done on a server with NVIDIA Quadro RTX 8000 and Intel Xeon Gold 6226R. We use $256 \times 256$ size for both the input photo and the output caricature.

### 4.1 Comparison to state-of-the-art methods

We qualitatively evaluate our method by comparing with two classes of methods: generic image-to-image translation methods and deep caricature generation methods (Fig. 9).

We compare our results with U-GAT-IT [Kim et al. 2020] and StarGAN v2 [Choi et al. 2020]. Among various image-to-image translation models, the two models have showcased visually pleasing results in their work when translating two domains with large shape changes. We trained U-GAT-IT and StarGAN v2 from scratch on the WebCaricature dataset using their official implementations.

| Input | Ours | U-GAT-IT | StarGAN v2 | WarpGAN | AutoToon |
|---|---|---|---|---|---|

Fig. 9. *Comparison with other state-of-the-art methods for image-to-image translation and caricature generation.* Our method can deform and stylize the input faces to produce more realistic and detailed caricatures in comparison to other methods. More visual comparisons can be found in the supplementary material. Inputs: ©Jordan Strauss/Invision/AP, ©AF archive/Alamy Stock Photo, ©Stuart C. Wilson/Getty Images, ©WENN Rights Ltd/Alamy Stock Photo, ©Danny Moloshok/Invision/AP, ©Slaven Vlasic/Getty Images.

U-GAT-IT does not successfully generate caricatures from the input, creating blurry images. Sometimes the changes in the output compared to the input are minuscule. StarGAN v2 generates more stable results, but sometimes artifacts are found in textures. Besides, the degree of exaggeration is smaller than ours.

We also compare our results with WarpGAN [Shi et al. 2019] and AutoToon [Gong et al. 2020], which are state-of-the-art methods for caricature generation based on deep learning. We used the pre-trained models of WarpGAN and AutoToon released by the authors. Similar to U-GAT-IT and StarGAN v2, WarpGAN was trained on
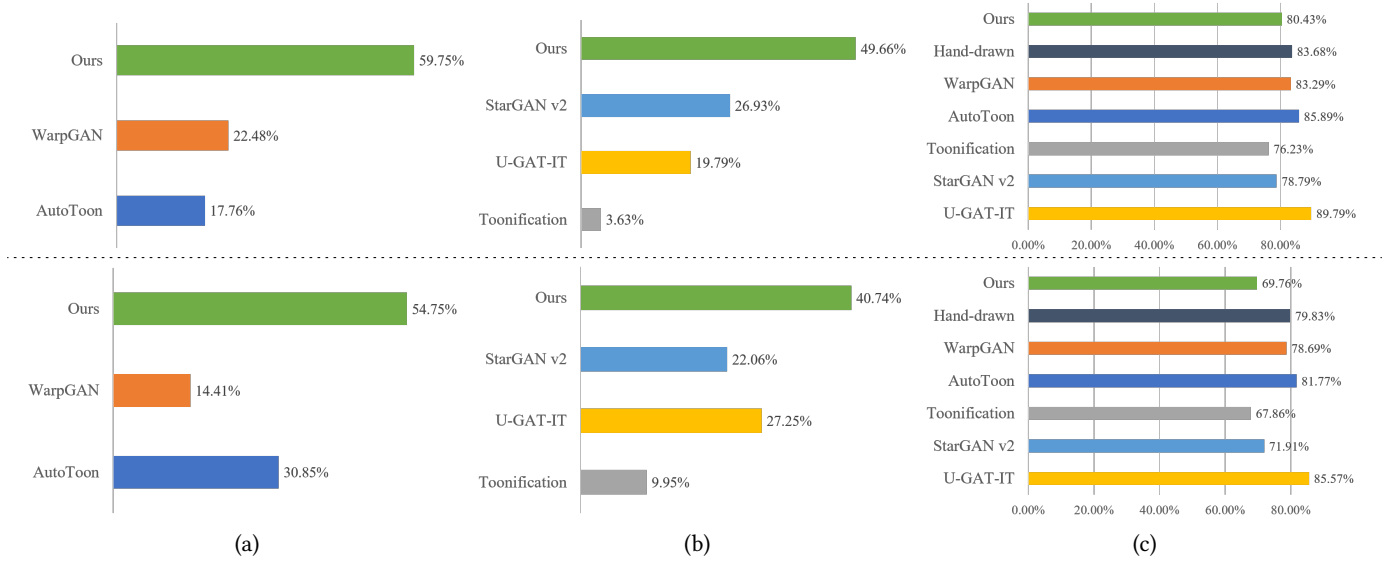
Fig. 10. *User study results.* Top and bottom parts show the results when our exaggeration factors are 0.5 and 1.0, respectively. (a, b) Faithfulness to hand-drawn caricature styles. Our method outperforms other methods with a large margin in terms of the quality of caricature styles. (c) Identity preservation. Our method shows comparable identification rates to other methods.

the WebCaricature dataset. In contrast, AutoToon was trained on its own dataset containing photo and caricature pairs for supervised learning. We also present comparisons with CariGANs [Cao et al. 2018] and CariGAN [Li et al. 2020] in the supplementary material.

As can be seen from Fig. 9, these deep learning based caricature generation methods afford large shape changes using explicit 2D image deformations. However, the generated shape deformations do not provide enough details to create realistic caricatures. WarpGAN uses only 16 sparse control points to handle deformations, and AutoToon is trained with only 101 photo-caricature pairs. Because of the limited density of control points and the sparsity of training samples, the generalization of these methods to realistic caricatures would be hard. Visually inspected, the results of our method are not only more pleasing but also closer to artist-created contents.

## 4.2 Quantitative analysis

We quantitatively measured the faithfulness to caricature image distribution with an evaluation of FID [Heusel et al. 2017]. Compared to other state-of-the-art methods (AutoToon, U-GAT-IT, WarpGAN, StarGAN v2), our method showed the best FID score (Table 2). This shows that our generated caricatures are closest to the distribution of caricatures. We trained U-GAT-IT, StarGAN v2, and ours with FFHQ [Karras et al. 2019] for photo and WebCaricature [Huo et al. 2017] for caricature dataset. We used officially released codes and models for AutoToon and WarpGAN. The reported values are FIDs between the generated caricatures and all caricatures in the WebCaricature dataset. The generated caricatures are created from all images in the CelebA [Liu et al. 2015] dataset, which was not involved in training.

## 4.3 Perceptual study

We perceptually evaluated the faithfulness to hand-drawn caricature styles through a user study. Given a photo, we asked users to

Table 2. *FIDs between generated and hand-drawn caricatures.* Lower is better.

| Method | FID |
|--------|-----|
| AutoToon | 114.84 |
| U-GAT-IT | 92.79 |
| WarpGAN | 74.60 |
| StarGAN v2 | 57.94 |
| Ours | **52.35** |

select the best caricature that looks like a hand-drawn caricature among different caricatures generated from state-of-the-arts methods and ours. We compared the responses for WarpGAN, AutoToon, U-GAT-IT, StarGANv2, and Toonification with ours. We evaluate our method using two different exaggeration factors (0.5 and 1.0). We run the experiment twice, one for each exaggeration factor of our method. Different users were involved in each experiment to prevent a learning effect. There are several methods to compare, and it could be hard for the user to answer consistently if the caricature results from all methods are shown together for each question. Therefore, we run the user study by splitting the methods into two groups: specifically designed caricature generator methods (WarpGAN, AutoToon) and general image-to-image translation methods (U-GAT-IT, StarGANv2, Toonification).

For each group, each user was asked 30 questions randomly sampled from the pool of 71 questions. The question pool was constructed by generating caricatures for the input photos previously used for the compared methods and the input photos in Fig. 9. We added five duplicate questions to filter out random selections. We excluded users that show inconsistent answers in more than three duplicate questions to obtain valid responses from 60 users for each group. Before starting the experiment, we asked each user five training questions using hand-drawn caricatures to expose the user to

realistic caricature styles. In each training question, a user was asked to pick the best hand-drawn caricature that matches the input photo. The user study was done using Amazon Mechanical Turk. In both of the groups, regardless of the exaggeration factor, our method significantly outperformed the previous methods (Figs. 10a-b). The questions and results are included in the supplementary material.

We also perceptually evaluated the degree of identity preservation of our method through another user study. Similarly to [Cao et al. 2018], we asked users to pick the photo that matches the identity of an input caricature. In each question, the input caricature was one of Hand-drawn, WarpGAN, AutoToon, Toonification, U-GAT-IT, StarGANv2, and ours. As in the other user study, we evaluated our method using two different exaggeration factors (0.5 and 1.0). For a fair comparison, we run the experiment twice with different exaggeration factors of our method, instead of running the experiment once with both results of ours, to balance the number of exposures per method. Consequently, in each experiment, there are seven caricature generators in total.

We provided five photo choices per question: one containing the answer with the same identity as the input for the caricature but with a different pose, and the rest containing similar faces to the answer. The similar faces for the wrong choices were selected by first encoding the answer into a FaceNet [Schroff et al. 2015] feature and estimating pose using [Deng et al. 2019], then selecting images that show close feature distances and similar poses to the answer among the union of FFHQ and CelebA datasets.

We run each experiment on caricatures of 34 randomly selected identities from WebCaricature [Huo et al. 2017]. With the seven methods to evaluate, there are $34 \times 7$ caricatures. We created seven sets of questions, each set containing 39 questions that consist of 34 questions and five duplicate questions for filtering out randomly-picking users. The caricatures of the sets are mutually exclusive and the union of the sets covers all the $34 \times 7$ caricatures. For each set, we asked 50 users to answer the questions. Before starting the session, we asked each user five training questions, where the inputs are hand-drawn caricatures and the correct answers are shown to the user. The user study was done using Amazon Mechanical Turk. The results show that ours (0.5) shows comparable identity preservation to other methods including hand-drawn caricatures, while ours (1.0) shows lower identity preservation (Fig. 10c). We included the questions and responses in the supplementary material.

Some baselines such as AutoToon, WarpGAN, and U-GAT-IT show better identity preservation than hand-drawn caricatures. The implementation of AutoToon available on the internet and we used for experiments only deforms the input image without changing the textures, leaving clues for inferring the identity. WarpGAN is trained with a reconstruction loss that forces the output to be similar to the input, preserving hints for inferring the input identity. U-GAT-IT often generates results that are not very different from the input. Our method generates caricatures that are close to hand-drawn caricatures, and with a proper exaggeration control, can generate caricatures with comparable identity preservation to hand-drawings. In Fig. 10c, the results for Hand-drawn and five previous methods differ among the two experiments on top and bottom, because the users participated in the Amazon Mechanical Turk experiments were not the same although the same sets of caricatures were used.

Table 3. *Quantitative analysis on attribute preservation.* The average accuracy of all identity-matching photo and hand-drawn caricature pairs is the upper bound, and the average accuracy of randomly selected pairs is the lower bound. For the identity-matching hand-drawn case, ground-truth labels are annotations in the dataset, and predicted labels are from our attribute classifiers. StyleCariGAN shows a reasonably high average accuracy compared to hand-drawn caricatures with predicted labels.

| Caricature source | Accuracy |
|---|---|
| identity-matching hand-drawn (using GT labels) | 85.26% |
| identity-matching hand-drawn (using predicted labels) | 73.73% |
| StyleCariGAN (using predicted labels) | 70.38% |
| random hand-drawn (using GT labels) | 66.16% |

### 4.4 Attribute preservation accuracy

To evaluate the attribute preservation capability of our StyleCari-GAN, we analyze the detected attributes of our caricature results. Given a pair of photo attributes and caricature attributes, we calculate the matching accuracy as $m/M$, where $m$ is the number of matched elements between photo and caricature attributes, and $M$ is the total number of attributes. We first compute the upper and lower bounds of the accuracy using hand-drawn caricatures, then show our level of attribute preservation by comparing our accuracy to the bounds.

We use WebCariA [Ji et al. 2020] dataset for the evaluation. The dataset provides labels for photos and hand-drawn caricatures. An attribute classifier is pre-trained on the caricatures of the WebCariA dataset to detect attributes from caricature images. For our evaluation, we filter the attribute labels of each photo by marking only the labels that are shared at least by 50% of all identity-matching photos. The filtering is done to remove noisy labels coming from extrinsic factors such as pose.

In this experiment, the upper bound is the average accuracy of all identity-matching photo-caricature pairs in the dataset, while the lower bound is the average accuracy of randomly-selected photo-caricature pairs. Our accuracy is calculated using StyleCariGAN results generated from all photos in the dataset. The evaluation result shows that our method preserves facial attributes reasonably well in the generated caricatures (Table 3). Note that WebCariA dataset has many mutually exclusive labels, which lead to many zeros in the attribute vectors. Consequently, attribute vectors could be similar even for random pairs, and the accuracy for random hand-drawn becomes relatively high.

### 4.5 Comparison to *Toonification* based on StyleGAN

Pinkney and Adler [2020] proposed the *Toonification* method that adds a cartoon structure into the input photo using StyleGAN layer swapping. The method generates a photo-realistic rendering on top of the cartoon structure by mixing coarse feature maps of cartoons and fine feature maps of photos. Although Toonification and our framework share the same idea of layer swapping, the goals and the
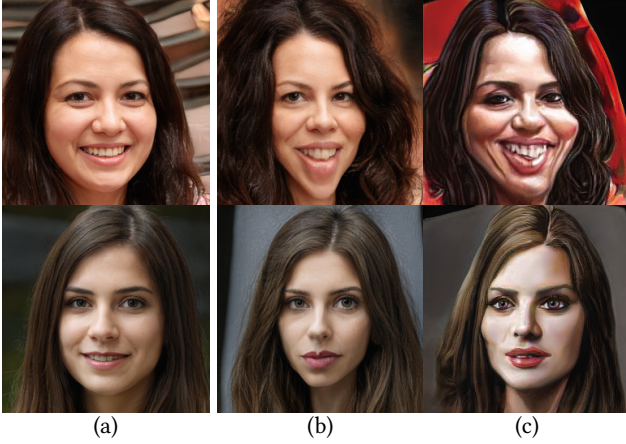
(a)          (b)          (c)

Fig. 11. *Comparison with Toonification [Pinkney and Adler 2020].* Toonification (b) blends an arbitrary caricature structure that largely changes the characteristics of the input (a). For example, Toonification does not preserve the overall facial shape of the input. Our framework creates shape deformations from the input image using *shape exaggeration blocks* trained explicitly to generate caricatures that preserves important visual features. As a result, our framework generates caricatures (c) that possesses exaggerated input structures and cartoon-style rendering.

results of the two frameworks are different. Toonification simply performs layer swapping using two StyleGANs trained for photos and cartoon images to create a cartoon image from an input photo. Our framework also performs layer swapping but our shape exaggeration blocks additionally create meaningful and facial-attribute-preserving shape deformations from the input photo. Fig. 11 shows a visual comparison, where our caricature StyleGAN is used for Toonification, instead of a StyleGAN trained to generate cartoon images. In the comparison, Toonification simply blends arbitrary caricature-style structures into the input photos. In contrast, our framework generates caricatures that not only have pleasing caricature styles but also preserve important features of the input photos.

Pinkney and Adler [2020] also mentioned that an image with cartoon-like textures can be created by using photo-StyleGAN features as the coarse layers and caricature-StyleGAN features as the fine layers. In that case, the method becomes identical to our setting except there are no shape exaggeration blocks. The results of this case using different mixing boundaries can be observed in Fig. 4. However, the deformations obtained by simple layer swapping are hard to control for achieving desirable caricature styles. For example, a result generated with more caricature layers can map a photo of a woman to a caricature image of an old man as seen in Fig. 4. On the other hand, if we take few caricature layers, we cannot obtain large enough deformations. Instead of inducing deformations with layer swapping only, we create caricature deformations by using shape exaggeration blocks that are directly supervised with losses designed for training a caricature generator.

### 4.6 Caricature-to-photo translation

Since StyleCariGAN is trained with cycle consistency, it has an inverse mapping for caricature-to-photo translation. To translate a caricature to a photo, we first encode a caricature image into the



Fig. 12. *Caricature-to-photo translation.* Given caricatures, we can generate photos using the inverse mapping of StyleCariGAN trained for cycle consistency. The generated photos share important visual features with the input caricatures. Input caricatures: ©MCT/Getty Images.

$\mathcal{W}+$ space of the caricature StyleGAN. Given the encoded vector as the input, we simply use *c2p-StyleCariGAN* trained for cycle consistency to generate a photo that resembles the given caricature. As illustrated in Fig. 12, our method can generate convincing caricature-to-photo translation results.

### 4.7 Exaggeration control

The first row of Fig. 13 demonstrates the contribution of shape exaggeration blocks, where the leftmost and rightmost images show the results without/with the blocks, respectively. It clearly shows our shape exaggeration blocks play a critical role for generating demanded shape changes in caricatures while preserving overall color stylization. Fig. 13 also shows some typical examples of controlling the degree of exaggerations in the caricature results by attaching scaling factors to shape exaggeration blocks. More examples without/with shape exaggeration blocks can be found in the supplementary material.

### 4.8 StyleGAN manipulation

Since our framework uses StyleGAN as the backbone network, various kinds of StyleGAN manipulations can be applied, including caricature color style customization and caricature expression control. Note that our framework does not require additional training process for such manipulations.

*Caricature appearance style selection.* In caricature generation, providing user control on appearance style is desirable. We can build such a control with style mixing [Karras et al. 2019], which can apply the tone and texture of a reference image onto the structure of a source image by injecting fine-scale latent codes of the reference image to the latent code of the source image. We use $\mathcal{W}+$ latent codes to perform the style mixing. A set of reference appearances are curated from randomly generated caricature samples using the caricature StyleGAN. As shown in Fig. 14, the appearance style can be flexibly changed with different reference styles, while the original structural contents can be well preserved.

*Caricature expression manipulation.* Our framework can be readily used with other StyleGAN-based semantic image editing methods.

Fig. 13. *Multi-scale control over shape exaggerations.* StyleCariGAN provides 4-scale control for exaggeration degrees. $\alpha_1$ controls the exaggeration of the coarsest scale, and $\alpha_4$ controls the finest scale. Each row shows different scales of edits. At each row, the set of parameters in the label is modulated, while unspecified parameters are fixed to 1. The first row shows the effects of shape exaggeration blocks, from no exaggeration to full exaggeration. Removing the blocks results in a caricature with photo-like structures (Row 1, Column 1). Applying all the blocks results in a caricature with desirable and realistic deformations (Row 1, Column 4). More examples with/without shape exaggeration blocks can be found in the supplementary material. The other rows show multi-scale control of exaggerations. Adding feature modulations of different scales from $\alpha_4$ to $\alpha_2$ produces structure deformations of increasing scales (Rows 2-4, Column 1). Parameter $\alpha_1$ for the coarsest feature modulation affects overall facial shape (Row 4, Columns 1-4). Note that the caricature results in the rightmost column are all the same since their parameter values are all 1.

Recent studies [Shen et al. 2020; Shen and Zhou 2020] have proposed semantic face editing based on a pretrained GAN. The basic idea of those methods is to find latent-space directions corresponding to the semantics, *e.g.*, smile, and then move the latent code of an input image along the desired directions. In Fig. 15, we use InterFaceGAN [Shen et al. 2020] to search for the direction of editing images towards smiling expression. The direction search is performed using the plain StyleGAN for photos. Then, we edit the latent code of the input image along the direction to modulate the magnitude of smile. StyleCariGAN produces a caricature image with the modulated smile magnitude by taking the edited latent code as input. As shown in the examples, this simple approach can add strong cartoon-like smiles without loss of identity information.

## 5 CONCLUSION

In this paper, we presented StyleCariGAN, a novel caricature generation framework based on StyleGAN. Our framework handles

the problem of caricature generation by modulating coarse-level feature maps with shape exaggeration blocks and swapping fine-level layers to the corresponding layers of the StyleGAN trained for caricature images. The shape exaggeration blocks are supervised to produce feature modulations for realistic and facial-attribute-preserving deformations. The modulated feature maps are rendered to a cartoon-style image with swapped fine layers. Our framework provides a multi-scale exaggeration control, and works with other StyleGAN-based image manipulation techniques. It creates realistic and detailed caricatures compared to other state-of-the-art methods.

*Limitations.* In this paper, we did not explicitly handle data bias in the training dataset. Due to the bias, our system does not successfully preserve some traits of the input in certain cases (Fig. 16). For example, the output caricature sometimes fails to preserve glasses contained in the input photo. In addition, some inputs that are largely different from the photo dataset, *e.g.*, gray-scale images or

Fig. 14. *Caricature appearance style selection.* Our framework can apply various appearance styles to the generated caricature shapes using a set of curated reference caricatures.

photos under a dim light, may result in invalid shapes and cause visual artifacts. Handling the bias in the dataset in a systematic way is necessary for a more stabilized real-world automatic caricature system. Besides, our caricature generator may disregard the hairstyle of the input subject under large deformations. Since our attribute matching loss only constrains the attributes related to facial shape, the preservation of hairstyle is not explicitly handled by supervision. This problem may be handled by adding new constraints on hairstyles, which is left as our future work.

## ACKNOWLEDGMENTS

## REFERENCES

Rameen Abdal, Yipeng Qin, and Peter Wonka. 2019. Image2stylegan: How to embed images into the stylegan latent space?. In *Proc. ICCV*.
Rameen Abdal, Yipeng Qin, and Peter Wonka. 2020. Image2StyleGAN++: How to edit the embedded images?. In *Proc. CVPR*.
Rameen Abdal, Peihao Zhu, Niloy Mitra, and Peter Wonka. 2021. Styleflow: Attribute-conditioned exploration of stylegan-generated images using conditional continuous normalizing flows. *arXiv* (2021).
Ergun Akleman. 1997. Making caricatures with morphing. In *Proc. ACM SIGGRAPH*.
Ergun Akleman, James Palmer, and Ryan Logan. 2000. Making extreme caricatures with a new interactive 2D deformation technique with simplicial complexes. In *Proc. Visual*.
David Bau, Hendrik Strobelt, William Peebles, Jonas Wulff, Bolei Zhou, Jun-Yan Zhu, and Antonio Torralba. 2019. Semantic photo manipulation with a generative image prior. *ACM Trans. Graph.* 38, 4 (2019), 1–11.

Susan E Brennan. 1985. Caricature generator: The dynamic exaggeration of faces by computer. *Leonardo* 18, 3 (1985), 170–178.
Andrew Brock, Theodore Lim, James M Ritchie, and Nick Weston. 2016. Neural photo editing with introspective adversarial networks. *arXiv* (2016).
bryandlee. 2020. FreezeG. https://github.com/bryandlee/FreezeG.
Kaidi Cao, Jing Liao, and Lu Yuan. 2018. CariGANs: Unpaired photo-to-caricature translation. *ACM Trans. Graph.* 37, 6, Article 244 (2018).
Hong Chen, Nan-Ning Zheng, Lin Liang, Yan Li, Ying-Qing Xu, and Heung-Yeung Shum. 2002. PicToon: a personalized image-based cartoon system. In *Proc. MM*.
Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. 2020. Stargan v2: Diverse image synthesis for multiple domains. In *Proc. CVPR*.
Antonia Creswell and Anil Anthony Bharath. 2018. Inverting the generator of a generative adversarial network. *IEEE Trans. Neural Networks and Learning Systems* 30, 7 (2018), 1967–1974.
Yu Deng, Jiaolong Yang, Dong Chen, Fang Wen, and Tong Xin. 2020. Disentangled and controllable face image generation via 3D imitative-contrastive learning. In *Proc. CVPR*.
Yu Deng, Jiaolong Yang, Sicheng Xu, Dong Chen, Yunde Jia, and Xin Tong. 2019. Accurate 3d face reconstruction with weakly-supervised learning: From single image to image set. In *Proc. CVPR Workshops*.
Julia Gong, Yannick Hold-Geoffroy, and Jingwan Lu. 2020. AutoToon: Automatic Geometric Warping for Face Cartoon Generation. In *Proc. WACV*.
Bruce Gooch, Erik Reinhard, and Amy Gooch. 2004. Human facial illustrations: Creation and psychophysical evaluation. *ACM Trans. Graph.* 23, 1 (2004), 27–44.
Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial networks. In *Proc. NeurIPS*.
Xiaoguang Han, Kangcheng Hou, Dong Du, Yuda Qiu, Shuguang Cui, Kun Zhou, and Yizhou Yu. 2018. Caricatureshop: Personalized and photorealistic caricature sketching. *IEEE Trans. Visualization & Computer Graphics* 26, 7 (2018), 2349–2361.
Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016a. Deep residual learning for image recognition. In *Proc. CVPR*.
Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016b. Identity mappings in deep residual networks. In *Proc. ECCV*.
Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. 2017. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Proc. NeurIPS*.
Xun Huang, Ming-Yu Liu, Serge Belongie, and Jan Kautz. 2018. Multimodal unsupervised image-to-image translation. In *Proc. ECCV*.
Jing Huo, Wenbin Li, Yinghuan Shi, Yang Gao, and Hujun Yin. 2017. Webcaricature: a benchmark for caricature recognition. *arXiv* (2017).
Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. 2017. Image-to-image translation with conditional adversarial networks. In *Proc. CVPR*.
Wen Ji, Kelei He, Jing Huo, Zheng Gu, and Yang Gao. 2020. Unsupervised domain attention adaptation network for caricature attribute recognition. In *Proc. ECCV*.
Tero Karras, Miika Aittala, Janne Hellsten, Samuli Laine, Jaakko Lehtinen, and Timo Aila. 2020a. Training generative adversarial networks with limited data. In *Proc. NeurIPS*.
Tero Karras, Samuli Laine, and Timo Aila. 2019. A style-based generator architecture for generative adversarial networks. In *Proc. CVPR*.
Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. 2020b. Analyzing and improving the image quality of stylegan. In *Proc. ICCV*.
Junho Kim, Minjae Kim, Hyeonwoo Kang, and Kwanghee Lee. 2020. U-gat-it: unsupervised generative attentional networks with adaptive layer-instance normalization for image-to-image translation. In *Proc. ICLR*.
Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv* (2014).
Nguyen Kim Hai Le, Yong Peng Why, and Golam Ashraf. 2011. Shape stylized face caricatures. In *Proc. MM*.
Wenbin Li, Wei Xiong, Haofu Liao, Jing Huo, Yang Gao, and Jiebo Luo. 2020. Carigan: Caricature generation through weakly paired adversarial learning. *Neural Networks* 132 (2020), 66–74.
Pei-Ying Chiang Wen-Hung Liao and Tsai-Yen Li. 2004. Automatic caricature generation by analyzing facial features. In *Proc. ACCV*.
Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. 2015. Deep Learning Face Attributes in the Wild. In *Proc. ICCV*.
Lars Mescheder, Andreas Geiger, and Sebastian Nowozin. 2018. Which training methods for GANs do actually converge?. In *Proc. ICML*.
Sangwoo Mo, Minsu Cho, and Jinwoo Shin. 2020. Freeze Discriminator: A Simple Baseline for Fine-tuning GANs. *arXiv* (2020).
Zhenyao Mo, John P Lewis, and Ulrich Neumann. 2004. Improved automatic caricature by feature normalization and exaggeration. In *Proc. ACM SIGGRAPH*.
Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan
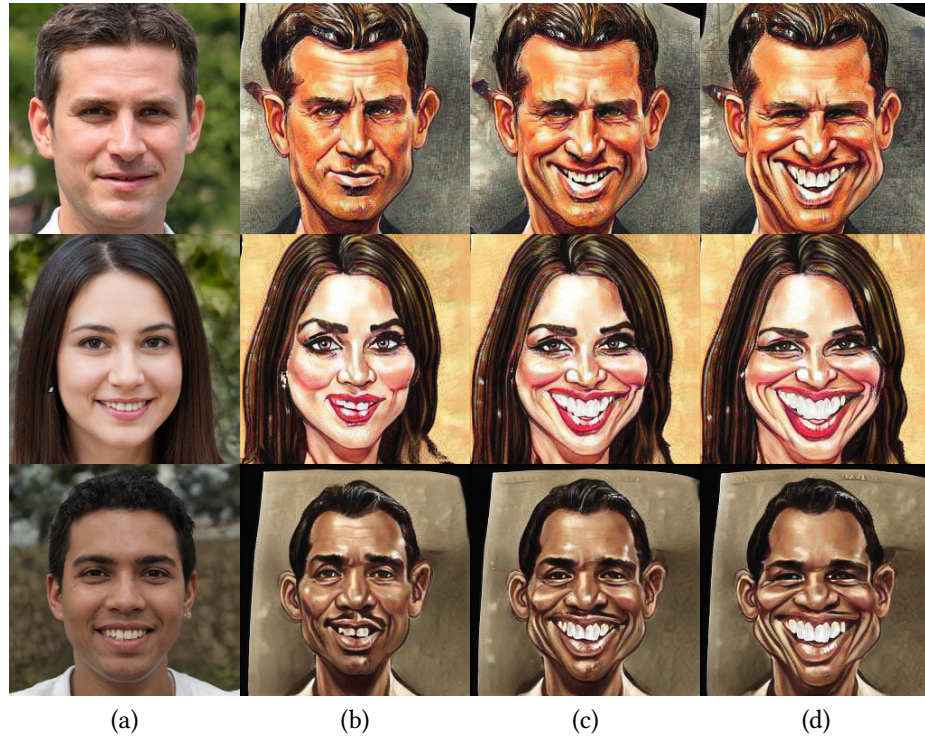
(a)       (b)       (c)       (d)

Fig. 15. *Caricature expression manipulation.* Using a StyleGAN-based semantic image editing method, we can manipulate our caricature results. Starting from the input (a), caricatures are generated (b). We can add desired amount of smile increasingly (c, d). Note that the latent code editing direction to create smiles was searched using the plain StyleGAN trained for photos, but the edits create natural and realistic caricatures even with strong smiles.



(a)       (b)

Fig. 16. *Limitations.* Some types of input photos (a) lead to failed caricatures (b). The input on the top row is a gray-scale image, which results in unsuccessful photo-to-latent embedding because the training examples for the photo StyleGAN did not contain many gray-scale images. The failure in the embedding incurs an invalid output. The person on the bottom row wears orange sunglasses, but the output does not preserve this unique look. Besides, the hairstyle is changed as well. As we did not apply explicit supervision for glasses or hairstyles, such visual features may get removed in the generated caricatures. Inputs: ©Masheter Movie Archive/Alamy Stock Photo, ©Angela Weiss/Getty Images.

Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *Advances in Neural Information Processing Systems 32*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.). Curran Associates, Inc., 8024–8035.

Justin NM Pinkney and Doron Adler. 2020. Resolution Dependant GAN Interpolation for Controllable Image Synthesis Between Domains. *arXiv* (2020).

Elad Richardson, Yuval Alaluf, Or Patashnik, Yotam Nitzan, Yaniv Azar, Stav Shapiro, and Daniel Cohen-Or. 2020. Encoding in style: a stylegan encoder for image-to-image translation. *arXiv* (2020).

rosinality. 2020. stylegan2-pytorch. https://github.com/rosinality/stylegan2-pytorch.

Florian Schroff, Dmitry Kalenichenko, and James Philbin. 2015. Facenet: A unified embedding for face recognition and clustering. In *Proc. CVPR*.

Yujun Shen, Ceyuan Yang, Xiaoou Tang, and Bolei Zhou. 2020. InterFaceGAN: Interpreting the disentangled face representation learned by GANs. In *Proc. PAMI*.

Yujun Shen and Bolei Zhou. 2020. Closed-form factorization of latent semantics in gans. *arXiv* (2020).

Yichun Shi, Debayan Deb, and Anil K Jain. 2019. Warpgan: Automatic caricature generation. In *Proc. CVPR*.

Ayush Tewari, Mohamed Elgharib, Florian Bernard, Hans-Peter Seidel, Patrick Pérez, Michael Zollhöfer, and Christian Theobalt. 2020. Pie: Portrait image embedding for semantic control. *ACM Trans. Graph.* 39, 6 (2020), 1–14.

Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. 2018. High-resolution image synthesis and semantic manipulation with conditional gans. In *Proc. CVPR*.

Raymond A Yeh, Chen Chen, Teck Yian Lim, Alexander G Schwing, Mark Hasegawa-Johnson, and Minh N Do. 2017. Semantic image inpainting with deep generative models. In *Proc. CVPR*.

Jiapeng Zhu, Yujun Shen, Deli Zhao, and Bolei Zhou. 2020. In-domain gan inversion for real image editing. In *Proc. ECCV*.

Jun-Yan Zhu, Philipp Krähenbühl, Eli Shechtman, and Alexei A Efros. 2016. Generative visual manipulation on the natural image manifold. In *Proc. ECCV*.

Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. 2017. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proc. ICCV*.