

# SofGAN: A Portrait Image Generator with Dynamic Styling

ANPEI CHEN\* and RUIYANG LIU\*, ShanghaiTech University

LING XIE, ShanghaiTech University

ZHANG CHEN, ShanghaiTech University

HAO SU, University of California San Diego

JINGYI YU, ShanghaiTech University

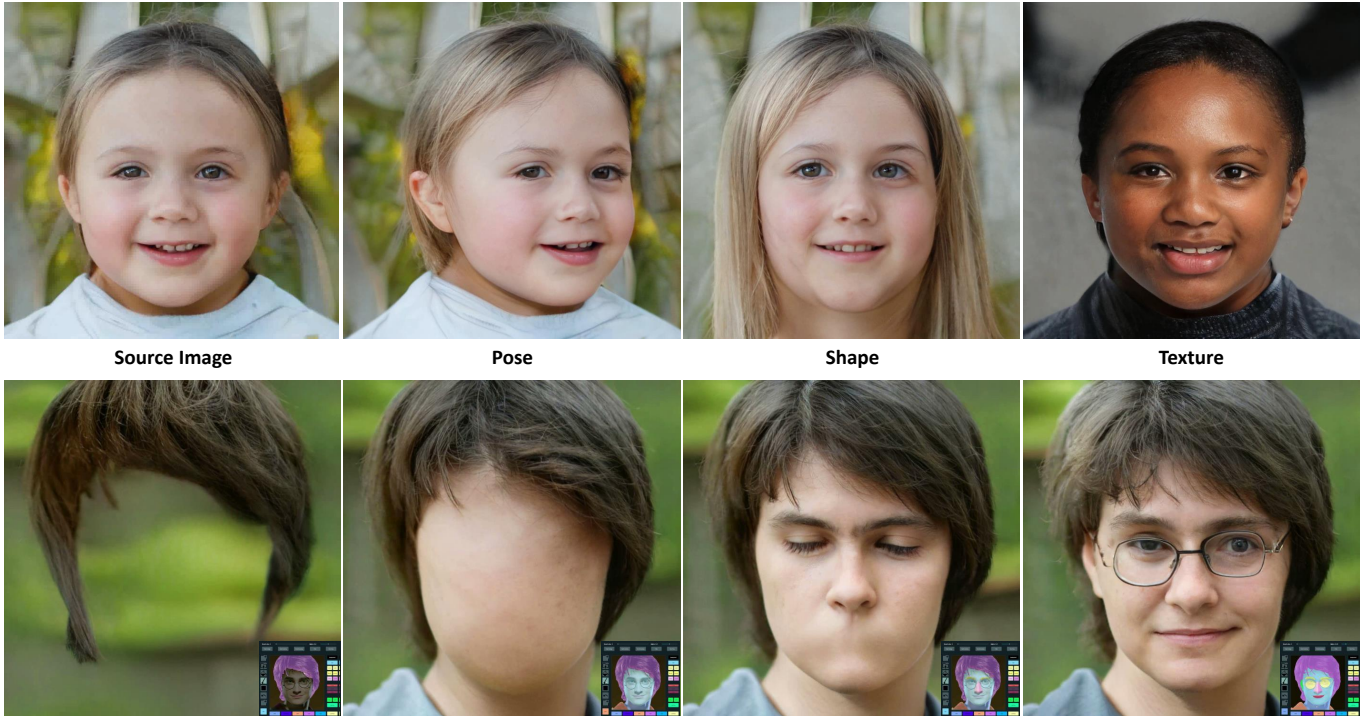


Fig. 1. First row: our portrait image generator allows explicit control over pose, shape and texture styles. Starting from the source image, we explicitly change it's head pose (2nd image), facial/hair contour (3rd image) and texture styles. Second row: interactive image generation from incomplete segmaps. Our method allow users to gradually add parts to the segmap and generate colorful images on-the-fly.

Recently, *Generative Adversarial Networks* (GANs) have been widely used for portrait image generation. However, in the latent space learned by GANs, different attributes, such as pose, shape, and texture style, are generally entangled, making the explicit control of specific attributes difficult. To address this issue, we propose a *SofGAN* image generator to decouple the latent space of portraits into two subspaces: a geometry space and a

\*Authors contributed equally to this work.

Author's addresses: Anpei Chen, Ruiyang Liu, Zhang Chen, Ling Xie, Jingyi Yu, Shanghai Engineering Research Center of Intelligent Vision and Imaging, School of Information Science and Technology, ShanghaiTech University, Shanghai, China; Anpei Chen, Ruiyang Liu, and Zhang Chen are also affiliated with the Shanghai Institute of Microsystem and Information Technology (SIMIT) and the University of the Chinese Academy of Science (UCAS); Hao Su, Department of Computer Science and Engineering, University of California San Diego, La Jolla, CA, USA..

© 2021 Association for Computing Machinery.

This is the author's version of the work. It is posted here for your personal use. Not for redistribution. The definitive Version of Record was published in *ACM Transactions on Graphics*, <https://doi.org/10.1145/3470848>.

texture space. The latent codes sampled from the two subspaces are fed to two network branches separately, one to generate the 3D geometry of portraits with canonical pose, and the other to generate textures. The aligned 3D geometries also come with semantic part segmentation, encoded as a *semantic occupancy field* (SOF). The SOF allows the rendering of consistent 2D semantic segmentation maps at arbitrary views, which are then fused with the generated texture maps and stylized to a portrait photo using our *semantic instance-wise* (SIW) module. Through extensive experiments, we show that our system can generate high quality portrait images with independently controllable geometry and texture attributes. The method also generalizes well in various applications such as appearance-consistent facial animation and dynamic styling. The code is available at [sofgan.github.io](https://sofgan.github.io).

CCS Concepts: • **Computing methodologies** → **Rendering**; *Computational photography*; Shape representations.

Additional Key Words and Phrases: Image synthesis, 3D modeling, Generative Adversarial Networks

**ACM Reference Format:**

Anpei Chen, Ruiyang Liu, Ling Xie, Zhang Chen, Hao Su, and Jingyi Yu. 2021. SofGAN: A Portrait Image Generator with Dynamic Styling. *ACM Trans. Graph.* 41, 1, Article 1 (July 2021), 27 pages. <https://doi.org/10.1145/3470848>

## 1 INTRODUCTION

Portrait modeling and rendering find broad applications in computer graphics, ranging from virtual reality (VR) and augmented reality (AR), to photorealistic face swaps, and to avatar-based tele-immersion. Effective portrait generators should be capable of synthesizing a diverse set of styles at high realism with explicit controls over attributes including illumination, shape (e.g., expression, poses, wearings), texture styles (e.g., age, gender, hair, and clothes), etc.

Physically-based rendering approaches have sought to explicitly model shape, materials, lighting, and textures and can potentially provide users complete controls over these attributes for synthesizing various styles. In reality, most existing techniques require carefully processed models as well as extensive manual tuning to achieve desired photorealism. More recent approaches have adopted learning-based schemes that aim to model styles as specific image distributions from a large-scale image dataset. Their major advantage is the avoidance of explicitly formulating geometry, lighting, or material that generally requires extensive experience and knowledge. For example, methods based on Generative Adversarial Networks (GAN) [Karras et al. 2019b,a] can now achieve extremely high realism and variety; however, there is still ample room to improve controllability over styles.

The seminal attribute-based approaches [Abdal et al. 2020; Chen et al. 2016; Shen et al. 2020; Tewari et al. 2020] explore to find attribute-specific paths in the latent spaces, e.g., to use direction codes to drag random style code towards each attribute direction. Such methods assume that attributes are linearly interpolative and independent from each other in the pre-trained generation space. However, this assumption is generally violated because attributes are entangled in the latent space, leading to flickering and undesirable attribute change.

It is also possible to employ 3D priors [Deng et al. 2020; Tewari et al. 2020] in portrait synthesis. These techniques can produce convincing facial animations and relighting results under user control. However, limited by the representation ability of explicit 3D facial priors (most commonly, the 3D morphable model), such approaches fail to model decorative attributes such as hairstyles, clothing, etc. The recent *Semantic Region-Adaptive Normalization (SEAN)* technique [Zhu et al. 2020] provides regional style adjustment conditioned on a segmentation map for more flexible image synthesis and editing. However, it also requires boundary-aligned segmentation maps and images (i.e., paired) for training that are difficult to acquire in practice.

We propose a disentangled portrait generation method by drawing inspirations from the practices of highly successful image rendering systems. In the image rendering literature, it is a basic practice to decompose the modeling of scenes as constructing the geometry component and the texture component. Likewise, we learn two individual latent spaces, i.e., a geometry space and a texturing space. Sparked by recent works on implicit geometry modeling [Chen and Zhang 2019; Mescheder et al. 2019; Park et al. 2019a;

Sitzmann et al. 2019b], we extend the traditional occupancy field to *semantic occupancy field (SOF)* to model portrait geometry. *SOF* describes the probabilistic distribution over  $k$  semantic classes (including hair, face, neck, cloth, etc.) for each spatial point. We train this *SOF* from calibrated multi-view semantic segmentation maps without groundtruth 3D supervision. To synthesize images, we first raytrace the *SOF* to obtain 2D segmentation maps from a given user-specified viewpoint then adopt GAN generator to texture each semantic region with a style code sampled from the texturing space. We propose a *semantic instance-wise (SIW)* texturing module to support dynamic and regional style control. Specifically, we design a novel semantic-wise “demodulation” and a novel training scheme that spatially mixes two random style codes for each semantic region during training. We further encode the semantic segmentation maps in a low dimensional space with a three-layer encoder to encourage continuity during view changing.

We have evaluated our method on the FFHQ dataset [Karras et al. 2019b] and the CelebA dataset [Lee et al. 2020]. Our generator achieves a lower Fréchet Inception Distance (FID score) and higher Learned Perceptual Image Patch Similarity (LPIPS) metric than the SOTA image synthesis methods. We will release our code, pre-trained models, and results upon acceptance.

To summarize, we propose a photorealistic portrait generator (*SofGAN*) with the following characteristics:

- **Attribute-consistent style adjustment.** Our new representation provides explicit controls over individual attributes with the rest unchanged and hence can support respective rendering effects such as free-viewpoint rendering, global and regional style adjustment, face morphing, expression editing, and artificial aging.
- **Training on unpaired data.** Unlike previous approaches that require using paired/aligned RGB and segmentation images for training, our SIW module can be directly trained using unpaired real-world images and synthesized semantic segmentation maps.
- **On-demand and Interactive Generation.** The tailored architecture of our generator supports photorealistic portrait synthesis from inaccurate or even incomplete segmentation maps. We have hence built a user-friendly tool that allows users to hand draw semantic contours for interactive portrait design (see the supplementary video).

## 2 RELATED WORK

**Unconditional image generation.** In the context of photo-realistic image generation, the pioneering ProgressiveGAN [Karras et al. 2017] conducts progressive training on both generator and discriminator and achieves significant improvements on rendering quality, training efficiency, and stability. Subsequent StyleGAN and extensions [Karras et al. 2019b,a] re-design the generator architecture to provide scale-specific control: the generator starts from a learned constant block and adjusts the “style” at each convolution layer. Such a structure enables direct controls over image features at different scales, i.e., unsupervised separation of high-level attributes from stochastic variations (e.g., freckles, hair) to support intuitive scale-specific mixing and interpolation. By providing a less entangled



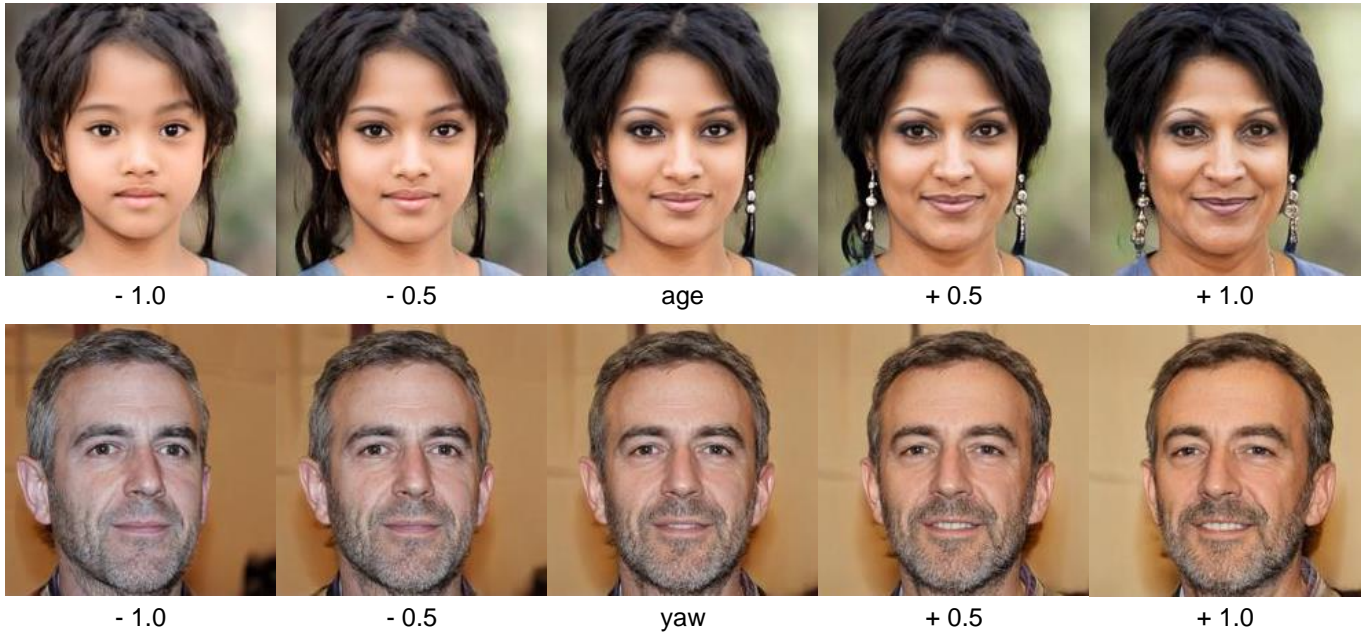


Fig. 2. Visualization of entangled attributes in the generation space of StyleGANs. In the first row, when we change the age attribute of the face, other attributes like hair style and ear rings are also changed. In the second row, complexion and facial emotions are changed along with yaw pose.

representation, these approaches produce state-of-the-art (SOTA) results.

**Conditional image generation**, also commonly referred to as style transfer, leverages GANs to learn a mapping from the source domain to the target domain ([Isola et al. 2017; Park et al. 2019b; Wang et al. 2018; Zhu et al. 2020]). Pix2Pix [Isola et al. 2017] firstly adopts a U-Net decoder in the generator to directly share the low-level information between input and output through skip-connections. It then trains the network with both perceptual loss and GAN loss. To tackle instability in adversarial training as well as to handle high-resolution image generation, the subsequent Pix2PixHD [Wang et al. 2018] uses a multi-scale generator and discriminator architectures. Although effective, these approaches face vanishing/exploding gradients when using a deep network. SPADE [Park et al. 2019b] proposes to use spatially adaptive normalization on all decoding layers rather than the first few layers of the network. The recent SEAN [Zhu et al. 2020] aims to synthesize style specified images by combining style latent vector and semantic maps, and achieves SOTA *Frechet Inception Distance* (FID) score. Since both SPADE and SEAN use the perceptual loss [Simonyan and Zisserman 2014] based on paired training data, their generated images still present certain blurriness.

**Controllable image synthesis.** Unconditional GANs unanimously synthesize images from randomly sampled latent vectors and thus are incapable of providing attribute-specific control (e.g., pose, eye, age). Inspired by the disentangled representation in the VAE and GAN literature, several recent approaches explore to achieve controllable image synthesis by disentangling the generation space into semantically meaningful attributes, either based on vector arithmetic phenomenon [Collins et al. 2020; Shen et al. 2019] or shape

priors [Tewari et al. 2020] such as 3DMM [Blanz and Vetter 1999]. Other controllable neural modeling methods [Chen and Zhang 2019; Mescheder et al. 2019; Nguyen-Phuoc et al. 2019, 2020; Park et al. 2019a] employ traditional implicit geometry representations such as the signed distance field (SDF) and occupancy field (OF) but under neural network approximations. In a similar vein, controllable neural rendering techniques extend physical-based models to conduct dynamic image generation tasks such as free-viewpoint relighting [Chen et al. 2020], reenactment [Cao et al. 2013; Geng et al. 2018; Siarohin et al. 2019; Tewari et al. 2020; Thies et al. 2016], novel views synthesis [Chen et al. 2018; Lombardi et al. 2019; Mildenhall et al. 2020; Sitzmann et al. 2019a,b; Thies et al. 2019; Zhou et al. 2020], etc, even with sparse inputs.

Instead of using explicit SDF and occupancy-based neural scene representations, our approach models the shared coarse geometry across different portraits via an embedded semantic occupancy implicit fields. Specifically, we set out to model the SOF using a set of semi-calibrated segmentation maps and subsequently encoding for the distribution of each 3D point over  $k$  semantic classes. This allows us to obtain free-viewpoint segmentation maps by projecting SOF onto a 2D space under a specific view camera. We further use a novel SIW module for modeling texture/appearance on each region in the semantic map. The combination of SOF and SIW modules enables better disentanglement between geometry attributes and texture styles, thus providing a new flexible control over image generation while maintaining generation quality.

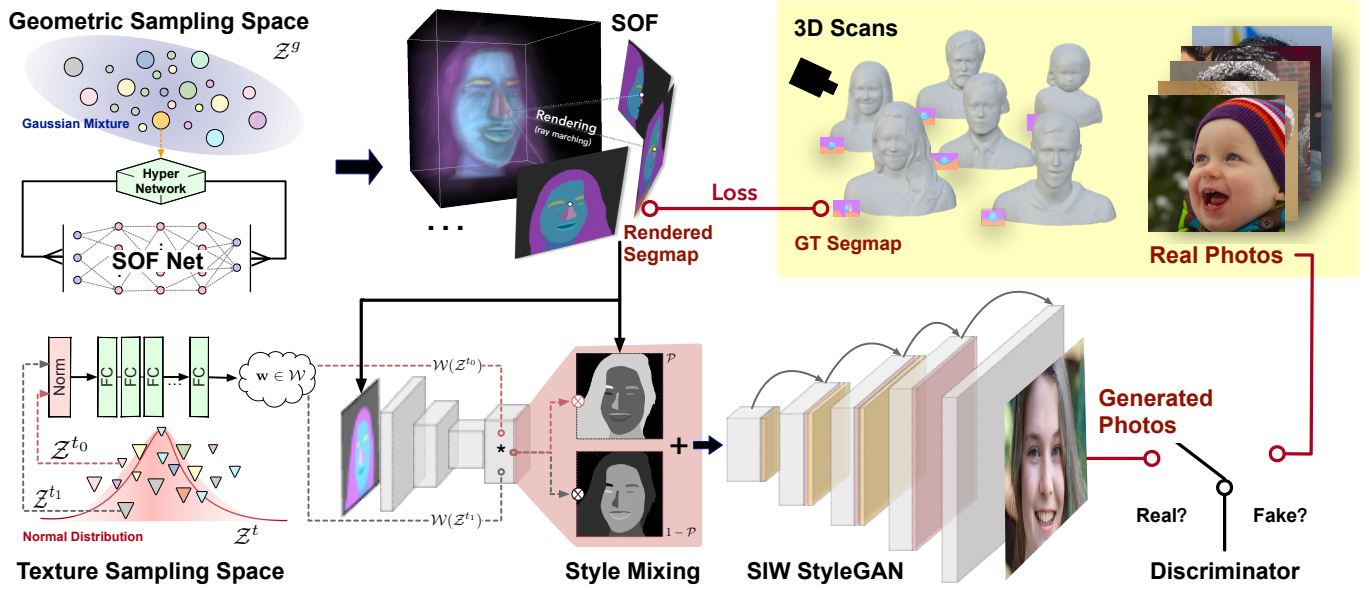


Fig. 3. An overview of our generation pipeline. Geometric latent code  $z^g$  in a learned Gaussian mixture. We adopt a hyper network to decode  $z^g$  into the weights of a MLP (*SOF Net*) which represents a contiguous semantic occupancy field (*SOF*) among 3D space. Then we render free-viewpoint segmentation maps from *SOF* with a ray-casting-marching scheme and generate two region-wise distance maps  $\mathcal{P}$ ,  $1 - \mathcal{P}$  for the following texturing. During texturing, we sample random vectors ( $z_0^t, z_1^t$ ) from the texture space and mod/demod into per-layer style vectors, then we adopt a Semantic Instance Wise (*SIW*) StyleGAN module to regionally stylize the generated segmaps. We use 5 style mixing layers (Pink) and 3 SPADE layer (Brown) in the *SIW*-StyleGAN. Our *SOF* and *SIW* are trained separately, where we render multi-view segmentation maps from synthetic portrait scans to train *SOF* and use real photos for *SIW*-StyleGAN training.

### 3 OVERVIEW

Our portrait generator synthesizes photo-realistic portrait images with controls over attributes including shape, pose, and texture styles. In StyleGAN, styles are controlled via features at each Conv layer based on the input latent vectors at different levels: coarse (contour), medium (expressions, hairstyle) and fine levels (color distribution, freckles). Although effective, such a control mechanism does not provide independent controls over individual attributes, largely due to the entanglement of various attributes.

To address this issue, we decompose the generation space into two sub-spaces: a geometry space and a texture space (Fig. 3 left). Each sample in the geometry space can be decoded into the weights of a *SOF Net* that represents a 3D-continuous occupancy field (*SOF*) with companion semantic labelings. At the rendering stage, given an arbitrary query viewpoint, we use a ray marching framework to map the *SOF* to a 2D segmentation map. The use of *SOF* ensures view consistency. Next we follow the semantic image synthesis framework [Isola et al. 2017; Park et al. 2019b; Wang et al. 2018; Zhu et al. 2017, 2020] and present a semantic-based, instance-wise (*SIW*) generation module to generate photorealistic images.

In a nutshell, our complete image generation process can be formulated as:

$$I = G(\mathcal{W}(z^t), \mathcal{R}(\text{SOF}(z^g), C)), \quad (1)$$

where  $z^g \in \mathbb{R}^n$  denotes the geometric code,  $z^t \in \mathbb{R}^m$  the texture code, and  $m$  and  $n$  the dimensions of the geometric and texture latent spaces respectively. *SOF* corresponds to a neural geometric representation defined by  $z^g$  whereas  $\mathcal{R}$  a differentiable renderer to render *SOF* onto a segmentation map at the query viewpoint  $C$ .  $G$  refers to the *SIW*-StyleGAN image generator and  $\mathcal{W}$  is a basis transformation operator on the texture code.

It is worth mentioning that a major limitation of previous condition-GAN style transfer methods is the requirement of using pairwise content/style images for training. In practice, collecting such image pairs is not only time-consuming but also requires elaborate alignment. Our technique does not impose such constraints. Rather, we set out to separately train *SOF* and *SIW*. In particular, training a *SOF* only requires using a relatively small set of multi-view segmentation maps that can even originate from scanned 3D models (Section 4). In sec.5, we present *SIW*-StyleGAN to support regional texturing from segmentation maps.

### 4 GEOMETRY MODELING

Previous approaches use the signed distance function (*SDF*) or the occupancy field (*OF*) as an alternative geometric representation to 3D points or meshes. For example, the *SDF* function  $\mathcal{F}(x)$  defined on  $\mathbb{R}^3$  models the signed distance value or occupancy probability. Such functions only characterize surface geometry (spatial location, normal, etc..) but do not consider the semantic meanings of the

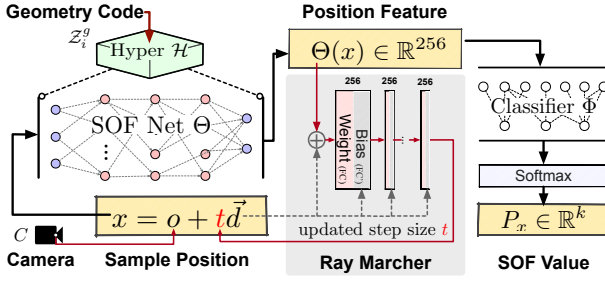


Fig. 4. The free viewpoint segmaps rendering process. Given a query camera  $C$  and geometric latent code  $z_i^g$ , we first use a hypernet  $\mathcal{H}$  to compute the weight of SOF network ( $\Theta$ ). Next, we conduct ray casting on SOF. To render each pixel, we cast its corresponding ray towards the SOF, starting from the given camera’s Center-of-Projection (CoP)  $o$ . We conduct iterative depth refinement using the SOF where the refinement step is also computed from an MLP-based ray marcher. Finally, we feed the estimated surface point and its feature vector to the classifier  $\Phi$  to decode them into a  $k$ -class semantic probability  $P_x \in \mathbb{R}^k$ .

surface. We observe that such semantic components are critical for photo-realistic image generation: they provide crucial guidance on generating convincing texture and shading details.

We instead use a novel semantic occupancy field (SOF) representation that extends the SDF by mapping an input point  $x \in \mathbb{R}^3$  to a  $k$ -dimensional vector to describe the occupancy probability of it belonging to  $k$  different semantic classes (e.g., eyes, mouth, hair, hat, clothing, etc.). The combination of semantic labels with geometry in SOF have several key benefits.

- (1) Regions defined by semantic classes serve as the basic units for regional style synthesis and adjustment. Such regional controls over individual components of the portrait (e.g. hair vs. eyes vs. mouth) greatly outperform the brute-force approaches that treat an image as a single entity.
- (2) Free-viewpoint image generation is a natural extension of SOF that can be seamlessly integrated into the deep networks: SOF describes the geometry attributes where new views can be directly synthesized via subsequent neural rendering.

An SOF models both geometry and semantic labels as:

$$\mathcal{F} : \mathbb{R}^3 \rightarrow \mathbb{R}^k, \quad \mathcal{F}(x) = P_x \quad (2)$$

where the semantic occupancy function  $\mathcal{F}(x)$  assigns to every point  $x \in \mathbb{R}^3$  in 3D space a  $k$ -dimensional probability ( $P_x$ ) over  $k$  semantic classes.

#### 4.1 Neural SOF Representation

Inspired by recent advances on neural scene representations [Oechsle et al. 2020; Park et al. 2019a; Peng et al. 2020; Saito et al. 2019; Sitzmann et al. 2019b], we approximate  $\mathcal{F}$  using two multi-layer perceptrons (MLPs), an SOF net  $\Theta$ , and a classifier  $\Phi$ . As shown in Fig. 4,  $\Theta$  maps each spatial location  $x \in \mathbb{R}^3$  to a spatial feature vector  $f \in \mathbb{R}^n$  that encodes the semantic property whereas  $\Phi$  corresponds to a semantic classifier with softmax activation to decode  $f$  to a  $k$ -dimensional semantic probability vector  $P_x$  as:

$$P_x \approx \Phi(\Theta(x)) \quad \text{where} \quad \Theta : \mathbb{R}^3 \rightarrow \mathbb{R}^n, \Phi : \mathbb{R}^n \rightarrow \mathbb{R}^k \quad (3)$$

We train the SOF by rendering known 3D portrait models into multi-view segmentation maps (segmaps) and compute the binary cross entropy loss between rendered segmaps and ground truth segmaps. The ground truth segmentation maps can be obtained by either semantically parsing real multi-view portrait images or by rendering synthetic 3D models with semantic labels (e.g., on textures), as shown in Fig. 19 of the Appendix.

To tackle diversity across 3D portraits, the brute-force approach would be to train a separate SOF for every portrait instance. Such approaches are highly inefficient and thus impractical for real-world applications. We instead train a geometric sampling space to support various portrait instances in the latent space using latent codes. We observe that 3D human portraits exhibit similar layouts, i.e., facial components such as eyes, hair, nose, and lips follow consistent spatial layouts. This enables us to train a shared canonical latent geometry space. Specifically, we adopt a shared hyper-network  $\mathcal{H}$  (Fig. 4) to learn, for each instance’s network, the weights of SOF net  $\Theta$ . Under the shared hyper-network, we can represent each instance in the training dataset with a latent vector  $z^g$  as:

$$\mathcal{H} : \mathbb{R}^m \rightarrow \mathbb{R}^{|\Theta|}, \quad \mathcal{H}(z_i^g) = \Theta_i, \quad i \in \{1, \dots, \|\mathcal{D}\|\}, \quad (4)$$

where  $i$  refers to the  $i$ th instance in training dataset  $\mathcal{D}$ . We use all  $z^g$ s to form our geometry sampling space  $\mathcal{G}$  with the Gaussian mixture model. In our implementation, we build a dataset  $\mathcal{D}$  with 122 segmented portrait scans and for each model we render 64 segmaps from randomly sampled viewpoints.

#### 4.2 Free-Viewpoint Segmap Rendering

Given the trained SOF, we can subsequently render free-viewpoint segmaps. Fig. 4 shows our detailed rendering procedure: we extract the portrait surface from the SOF via a ray marching scheme and estimate the per-ray depth  $t$  (from the camera center to the surface) as the sum of  $N$  marching step-sizes  $t_i, i \in [1, N]$ . The portrait surface  $\mathcal{S}$  corresponds to points:

$$\mathcal{S} = \{x \in \mathbb{R}^3 \mid x = o + \vec{d} \sum_{i=1}^N t_i\} \quad (5)$$

For each forward rendering pass, we first cast rays to the scene, and then use a ray marching scheme to obtain the surface point  $x$  (middle of Fig. 5) by iteratively estimating the marching step-size  $t_i$ , and finally predict its semantic property with the Eq.3. The multiply steps in the ray marching can be simulated with a LSTM network[Sitzmann et al. 2019b], and achieves pleasing scene representation performance. However, we observe that the LSTM structure is quite unstable and sensitive to novel views during viewpoint changing, as shown in Fig. 15 of the Appendix.

We thus propose a more stable ray marcher (Fig. 4, middle) that predicts the step sizes from the current position feature and the ray direction. Each ray marching MLP layer contains two sub-fully-connected layers  $FC_i$  and  $FC'_i$  where  $FC_i$  is a linear mapping without additive bias (i.e. the  $b$  in linear layer  $y = xA^T + b$ ) whereas  $FC'_i$  maps the ray direction to be the layer bias.



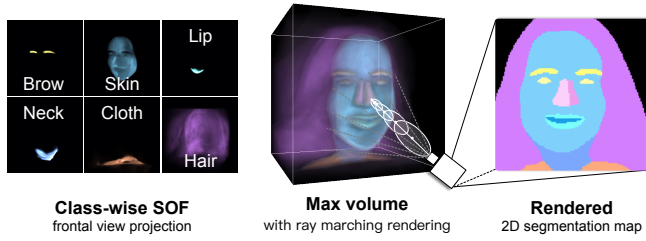


Fig. 5. *SOF* is a 3-dimensional volume, with a  $k$ -class semantic probability for each spatial point. The probability for each semantic class also forms a contiguous 3D volume (left). For each position, we take the semantic value with max probability and adopt a differentiable ray marching scheme (middle) to render it into 2D segmentation maps (right).

$$f_i = \text{ReLU}(FC_i(f_{i-1}) + FC'_i(\vec{d})), \quad (6)$$

$$i \in [1, \dots, 6], \quad f_0 = \Theta(x)$$

Fig. 5 shows a sample SOF with the predicted semantic labels (we only show the label that corresponds to the highest probability) for each point within the volume. We also demonstrate the improvement in the rendering quality of our proposed MLP ray marching module over the LSTM one in Fig. 15 of the Appendix. To further ensure depth consistency across views at the inference stage, we extract a coarse portrait proxy with marching cube algorithm (MC) and use the reprojected depth as the initial depth of ray marching.

## 5 TEXTURE SYNTHESIS

Next, we show how to use the 2D semantic segmaps from an arbitrary viewpoint to synthesize photorealistic images. Our approach is to conduct image-to-image translation [Isola et al. 2017; Park et al. 2019b; Wang et al. 2018; Zhu et al. 2017, 2020] and we present a *SIW-StyleGAN* technique that support high image quality and flexible style adjustments at both global and local scales.

To achieve attribute-specific generation, existing StyleGANs [Karras et al. 2019b,a] start synthesizing from a learnable  $4 \times 4 \times 512$  constant block and upsample it to higher resolutions ( $4^2 - 1024^2$ ) conditioned on the style vector  $z^t$ . Specifically, we transform the style vector to its corresponding kernel scales and conv image features  $F_{in}$  at each convolution layer where the output of each layer  $F_o$  can be formulated as:

$$F_o = F_{in} * w' \quad \text{with} \quad w' = \alpha \cdot w \quad \text{and} \quad \alpha = \Psi(\mathcal{W}(z^t)), \quad (7)$$

where  $w$  and  $w'$  are the original and the modulated convolution kernels, respectively,  $*$  is the convolution operator, and  $\alpha \in \mathbb{R}^{1 \times CH}$  is a channel-wise feature scaling vector obtained by feeding the style vector  $\mathcal{W}(z^t)$  to the modulation/demodulation function  $\Psi$ .  $z \sim \mathcal{N}(0, 1)$  is the input of StyleGANs sampled from a 512 dimensional normal distribution with mean 0 and variance 1.

To fully exploit the semantic segmaps from our scheme, we further formulate the texturing process as to regionally stylize the one-hot semantic segmap  $\mathcal{M} \in \mathbb{R}^{K \times 32 \times 32}$  with a vector  $z^t \in \mathbb{R}^{512}$  of normal distribution. We first use three convolutional layers to downscale  $\mathcal{M}$

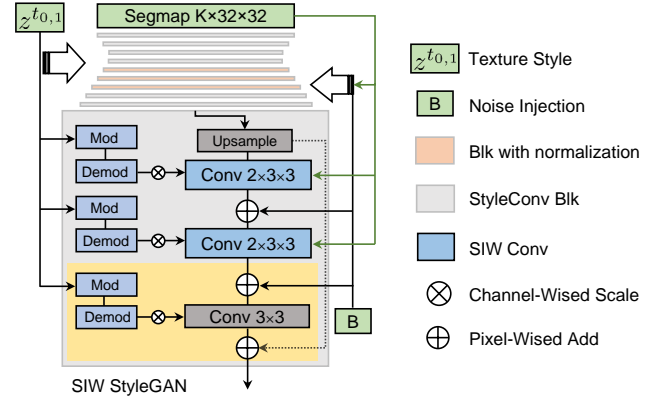


Fig. 6. Our SIW generation module. SIW extends StyleGANs2 to generate images from one-hot segmap and two normal distributed  $z^{t_{0,1}}$ , where  $z^{t_i}$  stands for texture style code. "Mod/Demod" refers to the  $\Psi$  term in Eq. 7 that transforms the style vectors to the style code  $\mathcal{W}(z_{0,1}^t)$  after Mod operation. We formulate the SIW Conv operation as Eq.7 and adopt SPADE normalization layer [Park et al. 2019b] at  $64^2 - 256^2$  resolutions (i.e., the orange blocks). Details of each layer can be found in Appendix A.

to a low-resolution feature map. Next we use a semantic instance-wise (SIW) convolution layer to conduct region-based upscale of the feature map to a high-resolution image (as the SIW StyleConv block in Fig. 6). We achieve the SIW style convolution operation by applying a region-specific convolution. Specifically, we first extend the modulation into  $K$  (i.e., the number of semantic labels, we set to 17 in all our examples) style vectors and then convert it into a pixel- and channel-wise scaling map, thus the output feature maps of each layer are formulated as the sum of all regional feature maps:

$$F_o = \sum_{i=1}^K (F_{in} * w'_i) \cdot \mathcal{M}_i \quad (8)$$

$$w'_i = \alpha_i \cdot w_i; \quad \alpha_i = \Psi(\mathcal{W}(z_i^t))$$

where  $\alpha \in \mathbb{R}^{K \times CH}$  includes the scaling factors for each feature channel and semantic region.

To strengthen the influence of the semantic segmaps in the generation process, we add an SPADE layer [Park et al. 2019b; Zhu et al. 2020] in the intermediate scaling layers (i.e.,  $64^2 - 256^2$  resolution shown as the orange blocks in Figure 6) and we have:

$$F_o = \gamma \sum_{i=1}^K (F_{in} * w'_i) \cdot \mathcal{M}_i + \beta \quad (9)$$

where  $\beta$  and  $\gamma$  are the mean and variance of the spatially adaptive normalization parameters:  $\gamma, \beta = \text{SPADE}(\mathcal{M})$ .

Recall that the operation above requires convolving the feature block with  $K$  different style kernels before fusing the intermediate feature blocks into one feature block (e.g., the summing in Eq.9). Brute-force implementation is neither efficient nor necessary: the one-hot semantic segmap contains mostly zeros so that most intermediate features do not contribute to the fused feature block.

Note that, the modulated regional texture styles  $w'$  are generated from random vectors in the training phase, we hence approximate

the generation by spatially mixing two style vectors  $z^{l_{0,1}}$  with a region-wise distance map  $\mathcal{P} \in [0, 1]$  (i.e., the mixed style training scheme), where the distance represents the similarity between the two styles. Each semantic region maps to a single style distance so that the output features can be simplified to a regional linear blending, the linear combination of two modulated features. Eq. 9 is then simplified to:

$$F_o = \gamma \cdot (F_{in} * \mathcal{W}(z^{l_0}) \cdot \mathcal{P} + F_{in} * \mathcal{W}(z^{l_1}) \cdot (1.0 - \mathcal{P})) + \beta \quad (10)$$

At the training stage, we obtain the distance map  $\mathcal{P}$  by assigning a random value to each semantic region. This mixed-style scheme simplifies  $K$  Conv and addition operations to 2 and therefore supports highly efficient training.

## 6 EXPERIMENTS

Next, we discuss our implementation details, conduct quantitative evaluations, and demonstrate our technique in a variety of applications.

### 6.1 Datasets

We use the following datasets in our experiments:

1) **CelebAMask-HQ** [Lee et al. 2020] that consists of 30k facial images and corresponding segmentation maps with attributes from 19 semantic part classes. We merge left/right pairs of parts into the same label (e.g., *left eye* and *right eye* are merged into a single *eye* class) and divide the nose region into left/right parts to emphasize the geometry structure in nose region.

2) **FFHQ** [Karras et al. 2019b] that consists of 70k high-quality images. We label the semantic classes with a pre-trained face parser.

3) **Self-captured**. Portrait scans exhibit relatively small geometric variations after segmented into semantic regions. This enables us to use a small 3D dataset to form the sample space of the *SOF*. Consequently, we capture 122 portrait scans with manually labeled semantic properties on their texture maps (as shown in the top right of Fig. 3). We render 64 random views for each scan to produce 7,808 segmentation maps in total for the *SOF* training. These 3D scans together with their corresponding semantic labels are ready to be released to the community.

### 6.2 Implementation

**Training the SOF.** Our *SOF* consists of three trainable sub-modules: hyper-net, ray marcher, and classifier, as shown in Fig. 4. The hyper-network  $\mathcal{H}$  in Eq. 4 contains 3 layers with 256 channels, and generates the weight of  $\Theta$  for each  $z^g$  via 4 fully connected (FC) layers. The ray marcher takes spatial features and ray directions as input and predicts the marching step size, as shown the middle of Fig. 4. The classifier is an FC layer with 256/ $k$  dimensional input/output channels respectively, as in Eq. 3. We use Adam optimizer [Kingma and Ba 2014], linear warm-up and cosine decay learning rate scheduler with  $lr = 1e - 4$ . For 122 scans, it takes about 10 hours to train *SOF* till convergence (with  $mIOU > 0.95$ ).

**Training the SIW-StyleGAN.** We adopt the same settings for training as StyleGAN2 [Karras et al. 2019a] in the following aspects: the dimensions of  $z^f$ , the architecture of the basis transformation

network  $\mathcal{W}$  (contains 8 fully connected layers), the leaky ReLU activation with  $\alpha = 0.2$ , the exponential moving average of generator weights [Karras et al. 2019a], the non-saturating logistic loss [Goodfellow et al. 2014] with R1 regularization [Mescheder et al. 2018], and the Adam optimizer with  $\beta_1 = 0, \beta_2 = 0.99, \epsilon = 10^{-8}$ . We refer the readers to StyleGAN2 [Karras et al. 2019a] for more details on the training process. It is worth mentioning that we only use the discriminator to minimize the distribution distance between our SIW-StyleGAN outputs and real captured photos, which shows that the SIW Conv layer implicitly learns the alignment between segmaps and texture styles, thus relaxing the pair-wise constraint between segmaps and portrait photos in training.

We perform all training with path regularization every 4 steps, style mixing with  $p = 0.9$ , and data augmentation via random scaling (1.0–1.7) and cropping. Our 1024×1024 model takes about 22 days to train on 4 RTX 2080 Ti GPUs and CUDA 10.1 (10000*king* iterations<sup>1</sup>). We observe that our module is already able to obtain visually pleasing results with only 1500*king* iterations (by image, which takes about three days), and the rest training iterations mostly target on improving high-frequency details, such as lighting, hair, pores, etc.

### 6.3 Quantitative Evaluation.

We quantitatively evaluate our method on various metrics including the Fréchet Inception Distance (FID) [Heusel et al. 2017], Learned Perceptual Image Patch Similarity (LPIPS) [Zhang et al. 2018], mIOUs and facial identity metrics. We compare our method with several recent image synthesis methods: Pix2PixHD [Wang et al. 2018], SPADE [Park et al. 2019b], SEAN [Zhu et al. 2020], and the baseline StyleGAN2 [Karras et al. 2019a].

Table 1 shows the requirements and achievable effects of each method. In contrast to the competing methods, our SIW-StyleGAN enables both global and local style adjustments without requiring pairwise data for training. However, we require an additional small dataset with 3D segmaps to form the basis of our geometric generation space.

To evaluate the above methods under the same settings of dataset, iterations, and resolution, we retrain their modules on the FFHQ [Karras et al. 2019b] and CelebA [Karras et al. 2017] datasets with 800*king* iterations at 512<sup>2</sup> resolution, and swap their segmentation maps as the evaluation set (i.e., if trained with dataset A, then use the segmentation maps of dataset B for testing). We set *truncation* = 1.0 and calculate FID value on 50k randomly sampled images. As shown in Fig. 7, our method achieves the best FID scores on both FFHQ and CelebA datasets<sup>2</sup>.

We attribute our major improvements to:

- (1) Different from SPADE and SEAN, our SIW-StyleGAN is trained with a non-pair-wise and unsupervised scheme instead of a perception loss on the pair-wise target photos, resulting in more freedom in the texture style space and thus increasing the variety of the generated images.

<sup>1</sup>*king* evaluates how many images used in the discriminator, which is independent of batchSize and GPU number

<sup>2</sup>The results for the comparison methods might be slightly different from their original papers, as we use random styles during evaluation.

	INPUTS		TRAINING		EFFECTS		
	latent code	segmentation	pairwise	3D segmaps	global style	local style	free view
<b>Pix2PixHD</b> [Wang et al. 2018]	√	√	√				
<b>SPADE</b> [Park et al. 2019b]	√	√	√		√		
<b>SEAN</b> [Zhu et al. 2020]	√	√	√		√	√	
<b>StyleGAN2</b> [Karras et al. 2019a]	√				√		
<b>SofGAN</b>	√	√		√	√	√	√

Table 1. Comparisons on training data requirement and controls over styles of our vs. SOTA methods. Pix2pix [Wang et al. 2018], SPADE [Park et al. 2019b], and SEAN [Zhu et al. 2020] unanimously require a large number of pair-wise (image and its registered segmap) training data whereas ours requires multi-view segmaps obtained from a small number of 3D scans. On controls, our technique provides more explicit and more varieties of controls than StyleGANs.

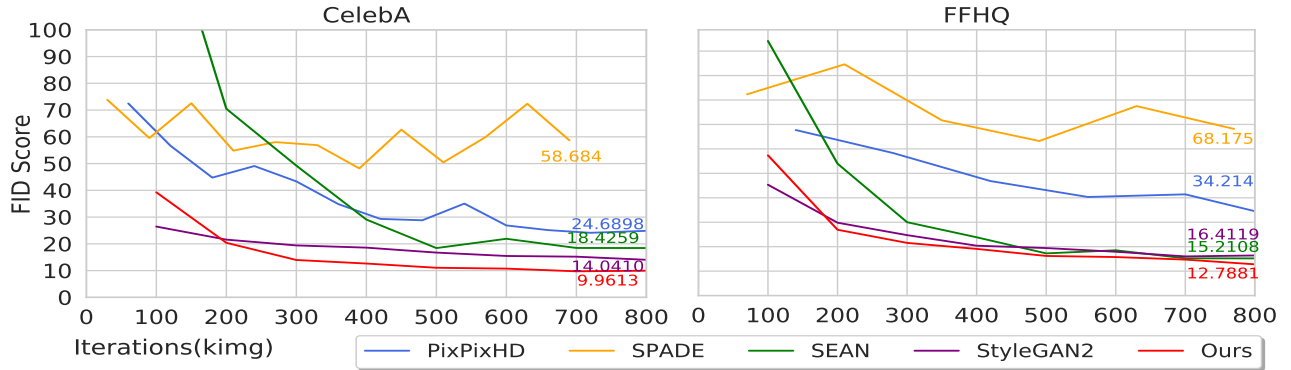


Fig. 7. Quantitative comparisons between our method vs. SOTA: we compare FID scores on CelebA [Lee et al. 2020] and FFHQ [Karras et al. 2019b]. Specifically, we first re-train all networks on one dataset and then generate images with random style vectors on other datasets for evaluation.

	CelebA			FFHQ		
	AlexNet ↑	VGG 16 ↑	SqueezeNet ↑	AlexNet ↑	VGG 16 ↑	SqueezeNet ↑
<b>Pix2PixHD</b> [Wang et al. 2018]	0.63±7.38	0.60±4.92	0.45±6.93	0.53±8.27	0.55±6.04	0.35±6.71
<b>SPADE</b> [Park et al. 2019b]	0.52±8.27	0.53±6.48	0.36±6.93	0.49±10.51	0.51±8.50	0.34±9.17
<b>SEAN</b> [Zhu et al. 2020]	0.51±12.52	0.57±11.40	0.41±12.30	0.58±6.26	0.66±6.04	0.50±6.71
<b>StyleGAN2</b> [Karras et al. 2019a]	0.62±7.38	0.62±7.38	0.47±7.16	0.64±6.04	0.66±5.59	0.51±6.48
<b>SofGAN</b>	<b>0.65±6.71</b>	<b>0.66±6.04</b>	<b>0.50±7.16</b>	<b>0.66±5.81</b>	<b>0.69±5.59</b>	<b>0.53±6.26</b>

Table 2. **Comparisons on diversity.** We compare the average perceptual divergence (with scaled std errors in unit of  $10^{-2}$  as subscript) of 50k image pairs generated by our SofGAN vs. SOTAs. SofGAN can produce a more diverse class of styles on both CelebAMaskHD [Lee et al. 2020] and FFHQ [Karras et al. 2019b] dataset.

- (2) Since we divide the image generation into semantic regions which have relatively similar texture patterns (i.e., semantic regions with the same class share a similar data distribution), our generator is able to focus on each small region separately instead of handling the whole image at the same time, and therefore, it converges faster than StyleGAN2.

Fig. 8, Fig. 28 and Fig. 29 in the Appendix give several visual comparisons of our method v.s. recent conditional image synthesis methods: Pix2PixHD [Wang et al. 2018], SPADE [Park et al. 2019b], SEAN [Zhu et al. 2020]<sup>3</sup>, and StyleGAN2 [Karras et al. 2019a] on both

<sup>3</sup>SEAN extracts style codes from reference style images, while other methods use randomly sampled style codes as input. Thus, we project the images in the evaluation

CelebAMask-HQ [Lee et al. 2020] and FFHQ [Karras et al. 2019b] datasets. From the comparisons, we can see that our framework performs better than Pix2PixHD, SPADE, and SEAN in both style diversity and realism of the generated images with comparable quality as StyleGAN2, which only supports global stylization.

Besides the improved realism, SofGAN further improves the diversity of texture styles. We generally expect the GANs to produce

set to the texture sampling space and acquire style codes for each image. To satisfy the pairwise constraint on semantic classes in SEAN, we only use the style images that have the same semantic classes with the conditioning semantic maps for the generation.



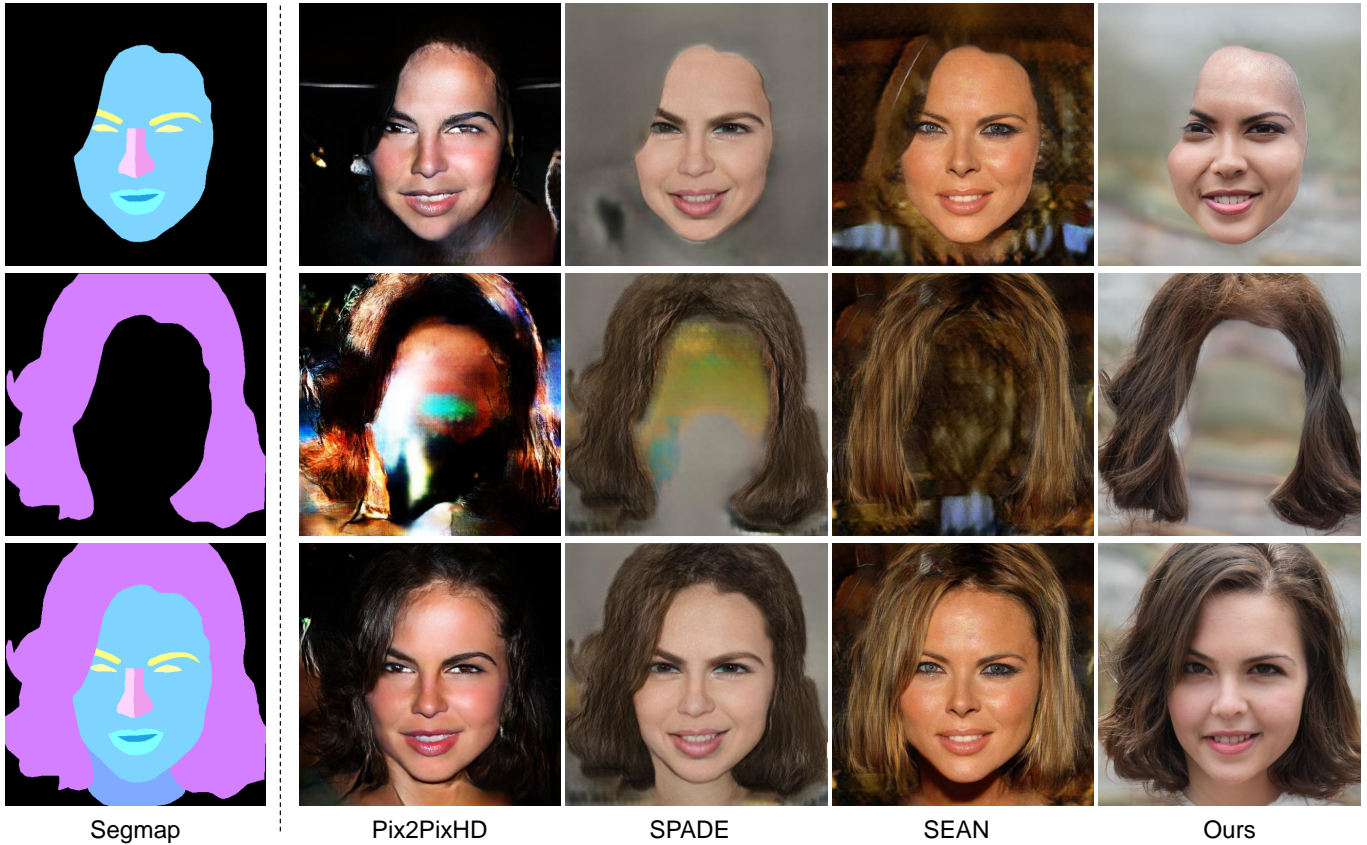


Fig. 8. Image synthesis from partial segment maps using SofGAN vs. SOTA. The SIW-StyleGAN in our SofGAN further improves visual realism.

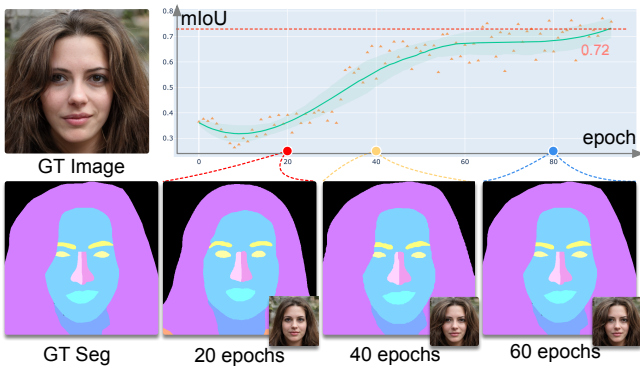


Fig. 9. Representability evaluation of the *SOF* space. We reproject 4000 segmaps from CelebAMask-HQ to our *SOF* space and evaluate their *mIOU* scores. (a) a semantic segmentation sample and its textured image using the SIW-StyleGAN. (b) Generated image and segmentation map after 20 epochs. (c) Generated image and segmentation map after 40 epochs ( $mIoU > 0.5$ ). (d) Final regressed segmentation map and generated image.

images with a wide variety while training GANs, a commonly encountered issue is that the discriminator gets stuck in a local minimum and the generator starts to produce similar outputs to favor the

specific discriminator, i.e., the *Mode Collapse* problem. During the experiment, we found that decomposing the generation space into sub-spaces and regionally generation could also help for dealing with the *Mode Collapse* problem and improve diversity.

To evaluate the representability of the trained *SOF* space, we therefore first randomly select 4,000 segmap samples from the CelebAMask-HQ dataset [Lee et al. 2020] and then search for the corresponding geometry latent code in our geometric sampling space with the *mIOU* metric for similarity evaluation. Fig. 9 shows that, for most samples, our scheme manages to find the corresponding geometry code with around 6k iterations. This indicates that the space formed with 122 3D scans can already cover a large variety of portrait shapes. We then jointly optimize  $\Phi$  and  $\Theta$  by casting rays from randomly sampled views for each instance  $D_i$  in  $\mathcal{D}$  and calculate the cross-entropy loss between the predicted segmentation and the ground truth segmentation rendered from  $D_i$ .

To evaluate the diversity among the generated images, we randomly sample 50k image pairs<sup>4</sup> from the same checkpoint with the FID evaluation (800kimg iterations), and calculate the mean LPIPS [Zhang et al. 2018] value with three backbone architectures

<sup>4</sup>With 50k image pairs we can constraint the standard derivative of LPIPS for all pairs at around  $\pm 0.0003$ , which we think is small enough to represent the general statics of the whole generation space.



Fig. 10. Free-viewpoint generation results on a same *SOF*. The first three rows show the free-viewpoint results with a same texture style and the last row shows results of applying different texture styles to the *SOF*. Our free-viewpoint stylizing effects conform the perspective imaging rules and also is able to preserve the facial identity.

(VGG16, AlexNet, and SqueezeNet). Table 2 compares our result with several recent GAN-based image generators. A higher LPIPS value means the generated images are more diverse, while a lower score means more similar to each other. Considering that artifacts and image noise could also pull up the LPIPS score, we compare the FID and LPIPS score simultaneously, and observe that our method achieves better image quality (FID score) with more diversity (LPIPS score).

Next, we evaluate the diversity and identity of conditional image generation by controlling the segmap and poses. For the diversity evaluation, we first randomly sample  $1k$  segmaps from CelebAMask dataset and stylize the above each random sampled segmap with 6 random texture styles and calculate the LPIPS [Zhang et al. 2018]

metric between the generated images (10 random image pairs for each). We show qualitative and quantitative results in Fig. 22 and the average LPIPS scores for single segmap generation (i.e., the SSG Diversity) of Table 3. Fig. 22 shows that our SofGAN can generate more realistic and less artifact results.

To evaluate the view consistency metric among the generated free viewpoint images, we randomly sample  $1k$  shape instances from the *SOF* spaces and render 15 views pure instance (e.g., as shown in the first three rows of Fig. 10). Then, we randomly sample 10 image pairs of the 15 views and evaluate their similarity score with the dlib face recognition algorithm<sup>5</sup>. Note that we center crop

<sup>5</sup>[https://github.com/ageitgey/face\\_recognition](https://github.com/ageitgey/face_recognition)





Fig. 11. Interactive image generation. Our method enables interactively synthesizing photo-realistic portrait images by drawing semantic segmentation maps (b columns). Our region-based texturing scheme is able to preserve shape and generalizes well even for extremely unnatural segmaps (as the "TOG" and "2020" patterns in the first row). In the second row, we add glasses and an earring to the left figure while we add glasses and a hat to the right figure. We can observe that our method can dynamically adjust the Appearance (even the lighting) according to the manipulation. In the last row, we manipulate the facial shapes and expressions.

face region with a face detection scheme before the comparison to favorite the input setting of the face parser. Table. 3 gives the quantitative evaluation for the free viewpoint video generation. We observe that our SIW module is able to better preserve the facial identity compared with the existing method [Park et al. 2019b; Zhu et al. 2020].

	FVV Identity $\uparrow$		SSG Diversity $\uparrow$	
	CelebA	FFHQ	CelebA	FFHQ
<b>SPADE</b>	0.460 $\pm$ 12.3	0.414 $\pm$ 10.8	0.116 $\pm$ 4.87	0.085 $\pm$ 7.41
<b>SEAN</b>	0.438 $\pm$ 11.8	0.399 $\pm$ 10.1	0.463 $\pm$ 0.18	<b>0.534<math>\pm</math>0.14</b>
<b>SofGAN</b>	<b>0.471<math>\pm</math>13.3</b>	<b>0.448<math>\pm</math>12.8</b>	<b>0.463<math>\pm</math>6.17</b>	0.480 $\pm$ 5.42

Table 3. **Quantitative evaluation of controlled generation.** We compare the preservation of facial identity under free viewpoint video generation (FVV Identity), and style variance generated with single segmap generation (SSG Diversity) with SPACE [Park et al. 2019b] and SEAN [Zhu et al. 2020]. Our architecture is able to better preserve identity during free-viewpoint generation and achieve compatible image diversity with SEAN when conditional on a same segmap. We show the mean FFV/SSG values of 1k images in each cell with std error as subscript in the unit of  $10^{-2}$ .

## 6.4 Applications

Recall that our generator is controlled by three individual components: camera poses  $C[R, T, K]$ , shape space latent code  $z^g$ , and texture space latent code  $z^t$ . In the following section, we will demonstrate various applications by exploring the above three components.

**Free-Viewpoint Video Synthesis.** Under our generation framework, we can generate free-viewpoint portrait images from geometric samples or real images by changing the camera pose.

As shown in Fig. 10, the last row of Fig. 12 and Fig. 23, since our *SOF* is trained with multi-view semantic segmentation maps, the geometric projection constraint between views is encoded in the *SOF*, which enables our method to keep shape and expression consistent when changing the viewpoint. Besides, our model is able to preserve overall consistency (regardless of some local surface details) of texture style including facial appearance and hairstyle even under significant view angle variation. For real images, we first parse a monocular segmap from the image, and reproject the segmap back to the geometric space. Then we generate *SOF* for the segmap and render free-view segmaps. The last row of Fig. 12 shows an free-viewpoint Obama result and Fig. 21 gives several examples for photos from CelebA dataset and internet.





Fig. 12. Our method preserves geometric and textural consistency in several applications. In expression editing (1st row), the hairstyle and background remain unchanged when the character changes from a neutral expression to a smiling). In gender/age morphing (2nd row), we maintain photo-realism of the central shape. In style mixing (3rd and 4th rows), we adopt two different styles and our results maintain consistency across frames (see the supplemental materials for comparison between ours vs. SOTA). In free-viewpoint rendering (bottom row), we first parse the segmap of President Obama’s photo, then back project the results to the *SOF* latent space, and finally stylize the segmaps with a random texture style.

**Shape Space Exploration.** As shown in Fig. 12 and Fig. 11, our method also supports generation and attribute editing in both 2D image space and the 3D geometric space. To generate new shapes, we fit a Bayesian Gaussian Mixture Model (GMM) with Dirichlet Process Prior to all samples  $z^g$  in the trained dispersed geometric sampling space  $\mathcal{G}$ , and then sample from the continuous GMM to generate new shapes outside the training dataset. For attribute editing, we borrow the idea of attribute decoupling [Shen et al. 2019] and extract shape-related eigenvectors (e.g., expression, age) from the *SOF*  $z^g \in \mathbb{R}^{256}$  via principal component analysis. The editing process is contiguous and allows us, for example, in the 1st row of Fig. 12, to gradually turn the expression of the little girl from neutral to smile.

The regional-specific property of our method also helps in dealing with the entanglement issue in most existing decoupling-based methods for the controllable generation ([Deng et al. 2020; Shen

et al. 2019, 2020]). The first row of Fig. 12 shows a smile expression interpolation example, with the general knowledge that smiling is mostly related to facial regions (e.g., eye, eyebrow and mouse regions) except other regions (e.g., hair, cloth, etc.), thus we only change the related region while fixing none-related region that enables to better preserve geometric and textural consistency during the generation. As our SIW-StyleGAN is region-aware, when applying shape morphing, the generated intermediate shape between two random shapes is also photorealistic (the 2nd row of Fig. 12).

Apart from sampling shapes from the *SOF* latent space, our method can also generate portrait images from existing photos and videos or hand-painted 2D segmentation maps, as demonstrated in Fig. 11 and the bottom rows of Fig. 12. In Fig. 11, we interactively modify segmentation maps to obtain novel user-specified shapes. Such an effect can benefit artistic creation or special effect generation. In Fig. 12, we collect a video clip from Internet and generate

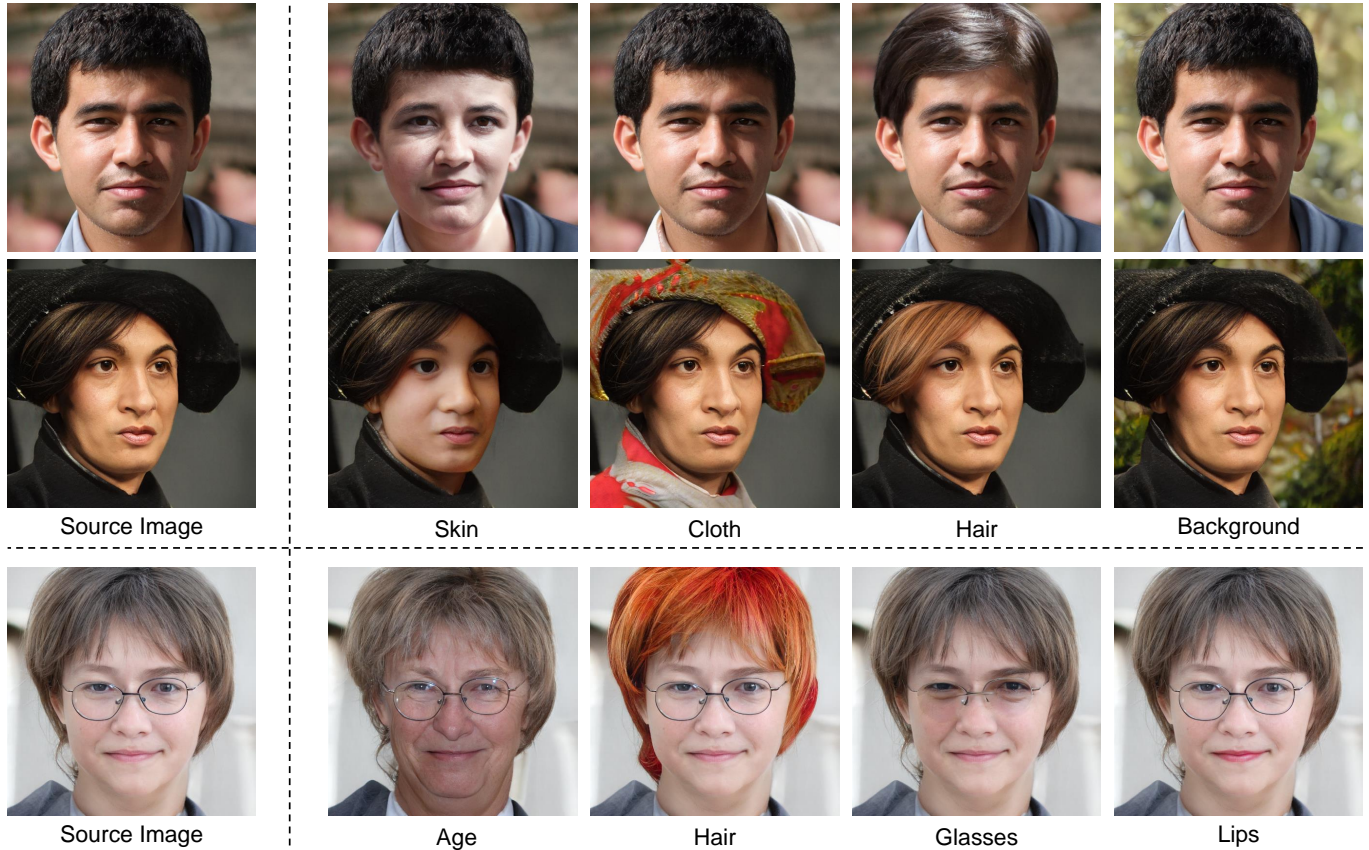


Fig. 13. Semantic-level style adjustment. Our algorithm enables semantic level style adjustment on 17 classes: skin, hair, lips, eyes & eyebrows, wearings, background and etc..

segmentation maps for each frame with a pre-trained face parser [Yu et al. 2018]. Our method can preserve texture style and shape consistency on various poses and expressions without any temporal regularization. We refer readers to our supplementary video for the animated sequences.

**Texture Space Exploration.** One of the key features of our SIW-StyleGAN is semantic-level style controlling. Benefiting from the StyleConv blocks and style mixing training strategy, we could separately control the style for each semantic region by adjusting the composed  $z^l$  through the similarity map  $\mathcal{P}$ . Fig. 13, 27 demonstrates the dynamic styling effects by changing texture styles in background, skin, hair, facial-region, hair, lip and wearing (e.g., cloth, hat and glasses etc.) regions of source images (more results could be found in the Appendix). Still, our approach could keep the texture styles unchanged except for the target region and adaptively adjust the semantic boundaries to ensure that the output looks natural. Moreover, our method could preserve global lighting when adjusting regional styles. Please refer to our supplemental video for more effects.

To better demonstrate the performance of our method, we further train a SIW-StyleGAN model with 10000 *king* iterations (one image per batch) under  $1024^2$  resolution. Fig. 24, 25 shows more results

on global style adjustment, Fig. 26 shows regional style controlling during generation, while Fig. 27 shows real-captured photo editing via our method.

## 7 CONCLUSION

In this paper, we have presented a novel two-stage portrait image generation framework that enables dynamic styling. We employed a semantic occupancy field as the geometric representation and a semantic instancewise StyleGAN for regional texturing. With our decomposed generation framework, we not only enable attribute-specific control over the generation but also allow users to generate from existing segmaps and editing attributes of existing images. In particular, our approach implicitly learns 3D geometric priors from 2D semantic maps without 3D supervision, eliminating the need for high-quality 3D portrait scans. Besides, our unsupervised training scheme for SIW-StyleGAN does not require paired segmaps and photos for training. Comprehensive experiments have shown that our approach achieves SOTA FID and LPISP scores on both CelebA and FFHQ datasets and can be used in a broad class of synthesis tasks.



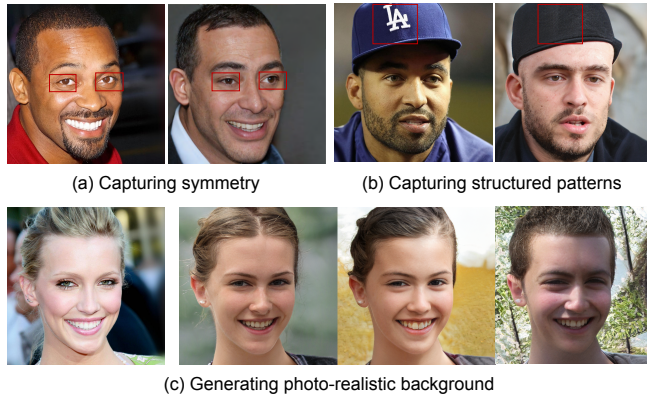


Fig. 14. We explore the limited power of SofGAN: (a) Capturing symmetry. The eyes may look towards different directions when generating extreme side faces (this phenomenon can be also seen in the right side of Fig. 10). (b) Capturing well-structured patterns. The SofGAN tends to generate flattened textures in decorative regions, and failed to reconstruct structured patterns like the “LA” in the hat. (c) Generating photo-realistic background, compared with real-captured photo (left), the generated background region tends to be blurry or noisy. Note that the left sides of each sub-figure are real-captured photos while the right sides are our generated images.

By using disentangled geometric and texture space, we could guarantee multi-view consistency at geometry level with *SOF*. However, as *SIW-StyleGAN* texture each semantic region independently, we do not guarantee pixel-level multi-view consistency, and significant shape modification would also leads to noticeable texture changes (Fig. 8, last column). Moreover, our method still exhibits limitations in capturing symmetry, structured patterns and photo-realistic generation in the background region. For example, the gazing directions sometimes are inconsistency (e.g., Fig. 14 (a)). Also, our method tends to output either flattened or regionally repetitive texture in the same semantic region, thus perform poorly on synthesizing complex patterns within the same semantic region, like the “LA” pattern in Fig. 14 (b). As shown in Fig. 14 (c), though our method succeeds in generating photo-realistic textures for the facial region, the background region is generally noisy or blurry. This may originate from the fact that our discriminator is not region-aware and only discriminates at the global distribution. As a result, such a design would flatten semantically specified styles. We plan to redesign the discriminator to be region-based and improve the semantically specified structures and styles in future work.

## 8 ACKNOWLEDGEMENTS

We thank Xinwei Li, Qiuyue Wang for dubbing the video, Zhixin Piao for comments and discussions, as well as Kim Seonghyeon and Adam Geitgey for sharing their *StyleGAN2* implementation and face recognition code for our comparisons and quantity evaluation. This work was supported by NSFC programs (61976138, 61977047), the National Key Research and Development Program (2018YFB2100500), STCSM (2015F0203-000-06) and SHMEC (2019-01-07-00-01-E00003).

## REFERENCES

- Rameen Abdal, Peihao Zhu, Niloy Mitra, and Peter Wonka. 2020. StyleFlow: Attribute-conditioned Exploration of StyleGAN-Generated Images using Conditional Continuous Normalizing Flows. *arXiv:2008.02401* [cs.CV]
- Volker Blanz and Thomas Vetter. 1999. A morphable model for the synthesis of 3D faces. In *Proceedings of the 26th annual conference on Computer graphics and interactive techniques*. ACM Press/Addison-Wesley Publishing Co., 187–194.
- Chen Cao, Yanlin Weng, Shun Zhou, Yiyang Tong, and Kun Zhou. 2013. Facewarehouse: A 3d facial expression database for visual computing. *IEEE Transactions on Visualization and Computer Graphics* 20, 3 (2013), 413–425.
- Anpei Chen, Minye Wu, Yingliang Zhang, Nianyi Li, Jie Lu, Shenghua Gao, and Jingyi Yu. 2018. Deep surface light fields. *Proceedings of the ACM on Computer Graphics and Interactive Techniques* 1, 1 (2018), 1–17.
- Xi Chen, Yan Duan, Rein Houthoofd, John Schulman, Ilya Sutskever, and Pieter Abbeel. 2016. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. In *Advances in neural information processing systems*. 2172–2180.
- Zhang Chen, Anpei Chen, Guli Zhang, Chengyuan Wang, Yu Ji, Kiriakos N Kutulakos, and Jingyi Yu. 2020. A Neural Rendering Framework for Free-Viewpoint Relighting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5599–5610.
- Zhiqin Chen and Hao Zhang. 2019. Learning implicit fields for generative shape modeling. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 5939–5948.
- Edo Collins, Raja Bala, Bob Price, and Sabine Süsstrunk. 2020. Editing in Style: Uncovering the Local Semantics of GANs. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Yu Deng, Jialong Yang, Dong Chen, Fang Wen, and Xin Tong. 2020. Disentangled and Controllable Face Image Generation via 3D Imitative-Contrastive Learning. *ArXiv abs/2004.11660* (2020).
- Jiahao Geng, Tianjia Shao, Youyi Zheng, Yanlin Weng, and Kun Zhou. 2018. Warp-guided GANs for single-photo facial animation. *ACM Transactions on Graphics (TOG)* 37, 6 (2018), 1–12.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. In *Advances in neural information processing systems*. 2672–2680.
- 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 *Advances in neural information processing systems*. 6626–6637.
- Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. 2017. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1125–1134.
- Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. 2017. Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196* (2017).
- Tero Karras, Samuli Laine, and Timo Aila. 2019b. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 4401–4410.
- Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. 2019a. Analyzing and improving the image quality of stylegan. *arXiv preprint arXiv:1912.04958* (2019).
- Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- Cheng-Han Lee, Ziwei Liu, Lingyun Wu, and Ping Luo. 2020. MaskGAN: Towards Diverse and Interactive Facial Image Manipulation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. 2019. Neural volumes: Learning dynamic renderable volumes from images. *arXiv preprint arXiv:1906.07751* (2019).
- Lars Mescheder, Andreas Geiger, and Sebastian Nowozin. 2018. Which training methods for GANs do actually converge? *arXiv preprint arXiv:1801.04406* (2018).
- Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. 2019. Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 4460–4470.
- Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. 2020. NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis. In *ECCV*.
- Thu Nguyen-Phuoc, Chuan Li, Lucas Theis, Christian Richardt, and Yong-Liang Yang. 2019. Hologan: Unsupervised learning of 3d representations from natural images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 7588–7597.
- Thu Nguyen-Phuoc, Christian Richardt, Long Mai, Yong-Liang Yang, and Niloy Mitra. 2020. BlockGAN: Learning 3D object-aware scene representations from unlabelled images. *arXiv preprint arXiv:2002.08988* (2020).
- Michael Oechsle, Michael Niemeyer, Lars M. Mescheder, Thilo Strauss, and Andreas Geiger. 2020. Learning Implicit Surface Light Fields. *ArXiv abs/2003.12406* (2020).



Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. 2019a. Deepsdf: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 165–174.

Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. 2019b. Semantic Image Synthesis with Spatially-Adaptive Normalization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.

Songyou Peng, Michael Niemeyer, Lars M. Mescheder, Marc Pollefeys, and Andreas Geiger. 2020. Convolutional Occupancy Networks. *ArXiv abs/2003.04618* (2020).

Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. 2019. PIFu: Pixel-Aligned Implicit Function for High-Resolution Clothed Human Digitization. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)* (2019), 2304–2314.

Yujun Shen, Jinjin Gu, Xiaoou Tang, and Bolei Zhou. 2019. Interpreting the latent space of gans for semantic face editing. *arXiv preprint arXiv:1907.10786* (2019).

Yujun Shen, Ceyuan Yang, Xiaoou Tang, and Bolei Zhou. 2020. InterFaceGAN: Interpreting the Disentangled Face Representation Learned by GANs. *arXiv preprint arXiv:2005.09635* (2020).

Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. 2019. First Order Motion Model for Image Animation. In *Advances in Neural Information Processing Systems*. 7135–7145.

Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014).

Vincent Sitzmann, Justus Thies, Felix Heide, Matthias Nießner, Gordon Wetzstein, and Michael Zollhöfer. 2019a. DeepVoxels: Learning Persistent 3D Feature Embeddings. In *Proc. CVPR*.

Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. 2019b. Scene Representation Networks: Continuous 3D-Structure-Aware Neural Scene Representations. In *NeurIPS*.

Ayush Tewari, Mohamed Elgharib, Gaurav Bharaj, Florian Bernard, Hans-Peter Seidel, Patrick Pérez, Michael Zollhöfer, and Christian Theobalt. 2020. StyleRig: Rigging StyleGAN for 3D Control over Portrait Images. *arXiv preprint arXiv:2004.00121* (2020).

Justus Thies, Michael Zollhöfer, and Matthias Nießner. 2019. Deferred neural rendering: Image synthesis using neural textures. *ACM Transactions on Graphics (TOG)* 38, 4 (2019), 1–12.

Justus Thies, Michael Zollhofer, Marc Stamminger, Christian Theobalt, and Matthias Nießner. 2016. Face2face: Real-time face capture and reenactment of rgb videos. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2387–2395.

Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. 2018. High-Resolution Image Synthesis and Semantic Manipulation with Conditional GANs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.

Changqian Yu, Jingbo Wang, Chao Peng, Changxin Gao, Gang Yu, and Nong Sang. 2018. Bisenet: Bilateral segmentation network for real-time semantic segmentation. In *Proceedings of the European conference on computer vision (ECCV)*. 325–341.

Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. 2018. The Unreasonable Effectiveness of Deep Features as a Perceptual Metric. In *CVPR*.

Hang Zhou, Jihao Liu, Ziwei Liu, Yu Liu, and Xiaogang Wang. 2020. Rotate-and-Render: Unsupervised Photorealistic Face Rotation from Single-View Images. *arXiv preprint arXiv:2003.08124* (2020).

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

Peihao Zhu, Rameen Abdal, Yipeng Qin, and Peter Wonka. 2020. SEAN: Image Synthesis with Semantic Region-Adaptive Normalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5104–5113.

## A NETWORK ARCHITECTURES

Table 4 shows detailed specifications of each sub-module in the neural *SOF* representation described in Sec. 4.1 of the main manuscript. Table 5 shows detailed network specifications for the texture modeling described in Sec. 5 of the main manuscript.

## B ABLATION STUDY

To better analyze the contribution of each component in our SofGAN, we conduct ablation studies on:

- 1) ray marching with LSTM vs. our proposed architecture.
- 2) constant input vs. with semantic map encoder.
- 3) with vs. w/o SIW style mixing block.

	Layer	Channels	Input	
$\mathbf{T}$	$fc_{0\_0}$	3/256	point location $x$	
	$fc_{0\_1}$	256/256	$fc_{0\_0}$	
	$fc_{0\_2}$	256/256	$fc_{0\_1}$	
<b>RayMarcher</b>	$fc_{1\_0\_l}$	256/256	$fc_{0\_2}$	
	$fc_{1\_0\_b}$	3/256	ray_dir $\vec{d}$	
	$fc_{1\_i\_l}$	256/256	$fc_{1_{(i-1)}\_l} + fc_{1_{(i-1)}\_b}$	
	$i \in [1, 5]$	$fc_{1\_i\_b}$	3/256	ray_dir $\vec{d}$
		$fc_{1\_6\_l}$	256/1	$fc_{1\_5\_l} + fc_{1\_5\_b}$
		$fc_{1\_6\_b}$	3/1	ray_dir $\vec{d}$
<b>Classifier</b>	$fc_{2\_0}$	256/256	$fc_{1\_6\_l} + fc_{1\_6\_b}$	
	$fc_{2\_1}$	256/256	$fc_{2\_0}$	
	$\Phi$	$fc_{2\_2}$	256/256	$fc_{2\_1}$
		$fc_{2\_3}$	256/20	$fc_{2\_2}$

Table 4. **Layer specifications for the neural *SOF* representation.**  $fc$  denotes a fully-connected layer with **Layer Normalization** and **ReLU** activation.

	Layer	Channels	s	Input
$i \in [1, 3]$	$Conv_{0\_0}$	16	2	$\mathcal{M}$
	$Conv_{0\_1}$	16	1	$Conv_{0\_0}$
	$Conv_{i\_0}$	$16 \times 2^i$	2	$Conv_{(i-1)\_0}$
$j \in [1, 3]$	$Conv_{i\_1}$	$16 \times 2^i$	1	$Conv_{i\_0}$
	$\Psi_0$	512	1/2	$z^t, Conv_{3\_1}$
	$\Psi_j$	512	1/2	$z^t, \Psi_{j-1}$
<b>SIW-StyleGAN</b>	$\Psi_4, \mathcal{N}$	512	1/2	$z^t, \Psi_3, \mathcal{M}, \mathcal{P}$
	$\Psi_5, \mathcal{N}$	256	1/2	$z^t, \Psi_4, \mathcal{M}, \mathcal{P}$
	$\Psi_6, \mathcal{N}$	128	1/2	$z^t, \Psi_5, \mathcal{M}, \mathcal{P}$
	$\Psi_7$	64	1/2	$z^t, \Psi_6, \mathcal{P}$
	$\Psi_8$	32	1/2	$z^t, \Psi_7, \mathcal{P}$
	$\Psi_9$	3	1/2	$z^t, \Psi_8, \mathcal{P}$

Table 5. **Network architecture for texture modeling (SIW-StyleGAN).**  $s$  denotes stride size and  $\mathcal{M}$  denotes the semantic segmentation map.  $z^t \in \mathcal{R}^{2 \times 512}$  and  $\Psi$  is the texture style modulation function as mentioned in Sec. 5.  $\mathcal{N}$  is the spatially adaptive normalization layer while  $\mathcal{P}$  is the style distance map for the mixed style training.

The modules in these ablation studies are **only** trained with  $800kimg$  iterations and  $512^2$  resolution as it is challenging to train each module with  $10000kimg$  and  $1024^2$  due to limited computing resources. Hence the image quality demonstrated in this section is much lower than the application section (Sec. 6.4).

### Ray Marching Architecture

We observed that the LSTM ray marcher [Sitzmann et al. 2019b] is sensitive to the initial camera position and can easily fail to render novel views when only trained on sparse views. Thus, we propose a new ray marching architecture, which estimates the step size for marching only based on the current position feature and ray direction without temporal status. This section conducts an ablation study on the ray marcher by feeding with various camera positions.

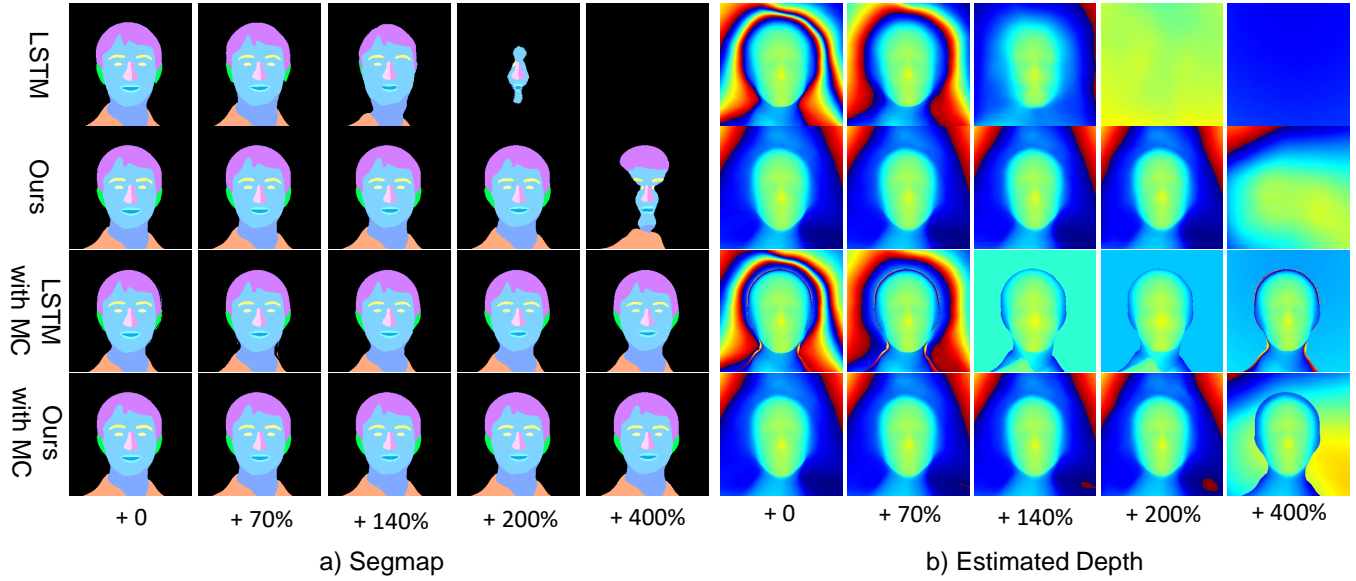


Fig. 15. Ablation analysis on the proposed ray marching network. It shows segmentation maps (left 5 columns) and depth maps (right 5 columns) estimated using different ray marching schemes. We gradually increase the distance rate of view camera, i.e., (novel view - training sets) / (training sets - world center). Each column shows the corresponding results under a distance rate.

As demonstrated in Fig. 15, we first compare LSTM approximation vs. our proposed architecture in the first two rows. The first column of Fig. 15 (a) shows a novel camera that is close to the training cameras. For the (a) column of Fig. 15, we move the camera away from the training camera positions and decrease the Field of View (FOV) synchronously to ensure the face is roughly at the same location in the image plane. Fig. 15 (b) shows the estimated depth  $t$  for each camera position in Fig. 15 (a). We can observe that our proposed architecture is much less sensitive to the camera position.

The last two rows of Fig. 15 further analyze the results when giving good depth initials, i.e., we initialize  $t$  for each ray with the Marching Cube (MC) Algorithm (as shown in Fig. ??). In practically, we first uniformly sample the “background” probability and to obtain the portrait surface via 0.5 probability threshold, and then project the surface to the rendering camera to obtain the initial ray marching depth  $t$ . As shown in the last two rows of Fig. 15 (b), LSTM can obtain fine predictions inside the object but performs poorly on the boundaries between the portrait and the background (the boundaries are eroded in this case due to wrong depth estimation). On the contrary, our architecture is able to completely recover both interior and boundary regions.

#### Constant Input vs. With Encoder

Unlike the StyleGAN2, which focuses on synthesizing realistic static images, our generator attempts to enable controllable generation including local and global styling, free-viewpoint generation, and stylizing image sequences. We observe that using constant input would cause two artifacts. 1) The “phase” artifacts: as shown in Fig. 17, refer to a strong localized appearance which is especially noticeable when viewed as an image sequence (keep static during animation). 2) The style “entanglement” artifacts: our generation

process is controlled by both the style codes and the semantic maps, we observe that they are usually incomparable to each other, e.g., feeding a woman’s semantic map and a man’s texture styles would lead to significant entanglement (shown in Fig. 16 (a)).

More specifically, we observe that the “phase” artifacts are caused by the translation invariance property of CNN layer, especially deconvolution from the constant feature blocks. This results in same feature values. At the same time, since the normalization layers are very shallow (two convolutional layers) and have a small receptive field, the design with constant input could not capture the global shape features of the semantic maps and leads to the style “entanglement” artifact. To address these two artifacts, we use additional resBlocks to enhance the connection between style and semantic maps, as shown in Fig. 16 (b). The encoder can efficiently reduce the artifacts by dynamically changing the generator’s conditional features. In this way, the network is able to capture global shape features (e.g., contour, position in the image, etc.) and strengthen semantic control.

#### With vs. w/o SIW Style Mixing

To enable the regional stylizing effect, we proposed a mixed style training strategy in Sec. 5. To demonstrate the effectiveness of this strategy, we train another model that removes this design during training and generates an image with two random styles only during evaluation (shown in the second row of Fig. 18). In the case of Fig. 18, the skin region is generated with style  $z^0$ , while other regions are synthesised with  $z^1$ . We can observe that significant artifacts exist on the boundaries and the images look unnatural for the second row, while the results trained with the mixed style strategy (the first row) produce soft boundaries.

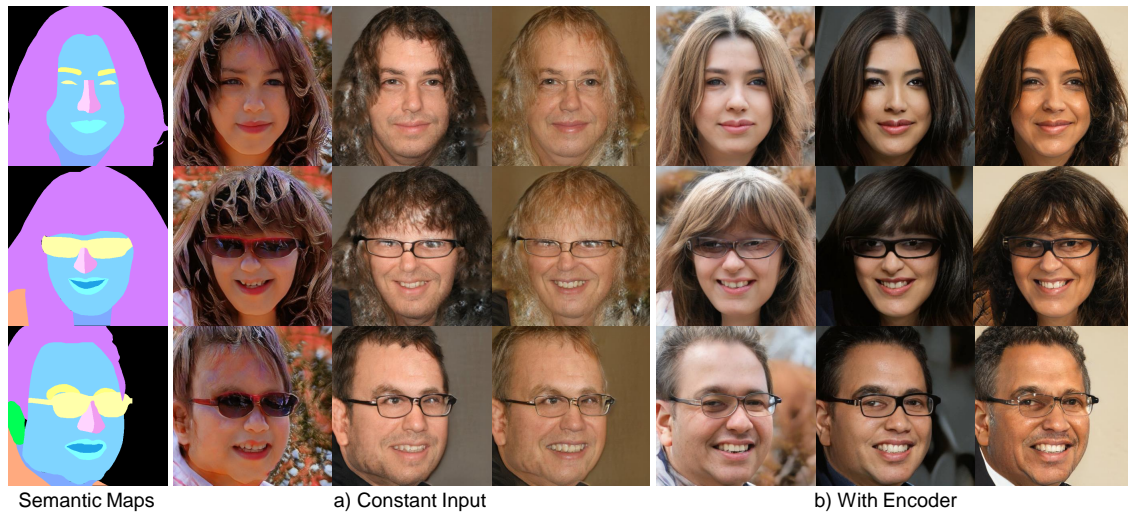


Fig. 16. Ablation analysis on w/ and w/o the segmentation map encoder. Each column shares the same style code while each rows denotes different segmap encoder: a)  $512 \times 4 \times 4$  constant input, b)  $17 \times 128 \times 128$  one hot semantic maps downscaled to  $128 \times 16 \times 16$  with SIW StyleConv block, which can dynamically adjust gender styles according to different input segmentation maps and reduce gender entanglement.



Fig. 17. Visualization of the "phase" artifact: strong localized texture appearance keeps static during animation.

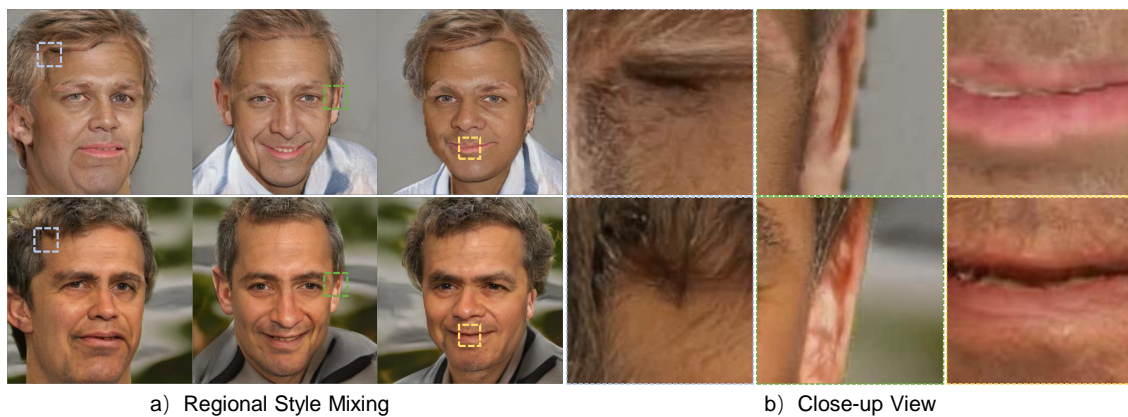


Fig. 18. Ablation analysis on w/ and w/o the mixed style training scheme. These results are generated after  $800k_{img}$  iterations. Top row: mixing style with SIW mix style blocks. Bottom row: mixing style results without SIW mix style blocks. Hair, eyebrows and lips share the same style, while skin uses another one. Mixed style training scheme can significantly improve the smoothness around the semantic boundaries.

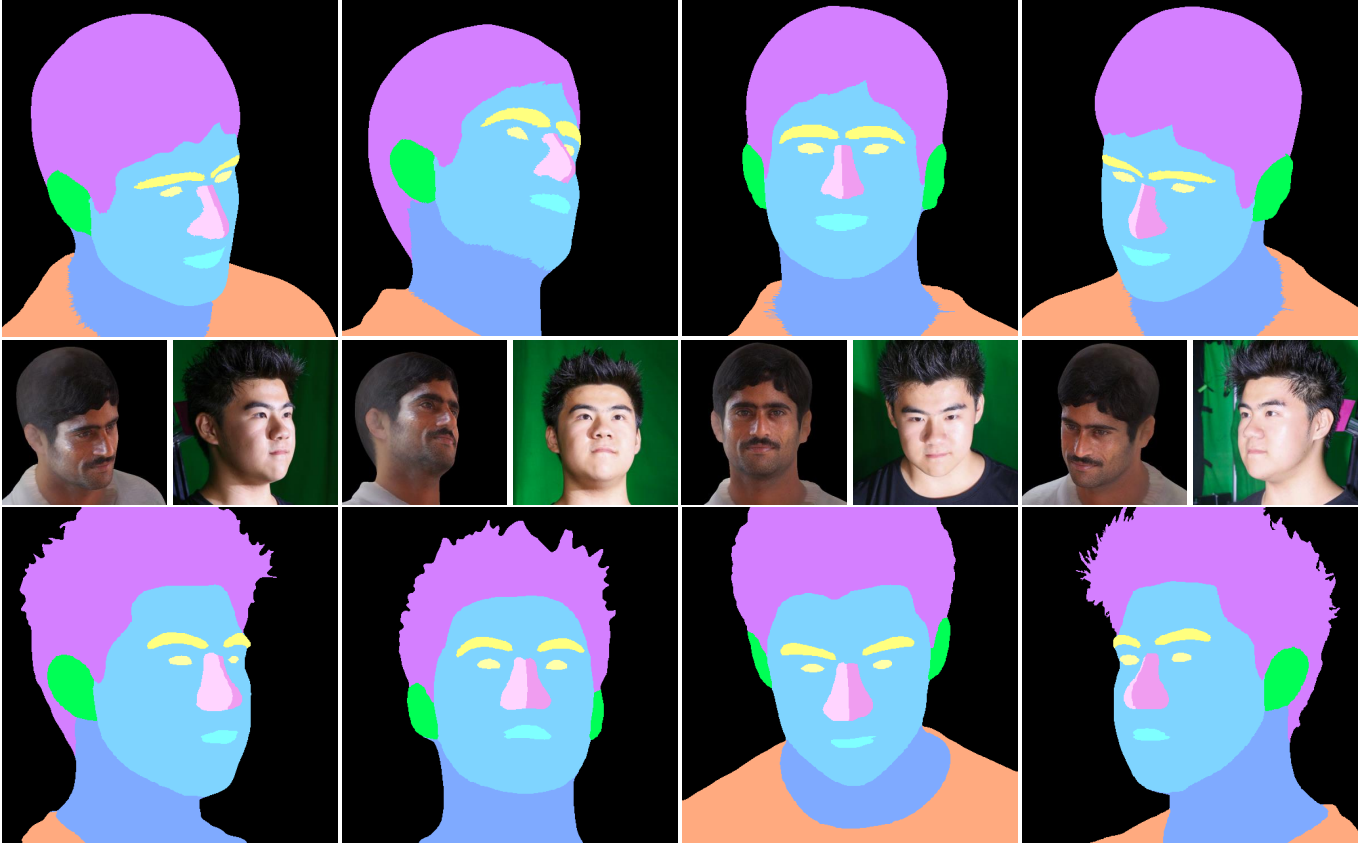


Fig. 19. Data samples from *SOF* training set. Top row: multi-view segmaps rendered from synthetic portraits (middle left). Bottom row: segmaps parsed from multi-view photos (middle right).

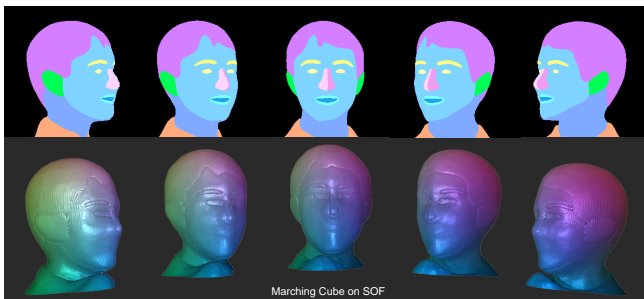


Fig. 20. Visualization of the *SOF*. To visualize the trained semantic occupancy field, we uniformly sample the occupancy of the “background” class in the *SOF* and use the Marching Cubes algorithm with  $level = 0.5$  to obtain 3D surface (foreground). We can observe explicit boundaries around each semantic class on the surface as the *SOF* encodes geometry along with its semantic properties.

## C SUPPORT FIGURES

Here we provide additional supporting figures mentioned in the main paper. Fig. 19 shows two multi-view segmap examples of our *SOF* training set. Fig. ?? gives visualization of a trained *SOF*. Fig. 22 compares the controlled image generation conditioned on same segmaps.

## D ADDITIONAL RESULTS

We provide additional results to better demonstrate the robustness of our SofGAN. Fig. 23 shows free-viewpoint generation results conditioned on different *SOF*. Fig. 24, 25 and 26 show style adjustment on both global and local semantic regions. Fig. 28 and 29 demonstrate the visual comparisons between Pix2PixHD [Wang et al. 2018], SPADE [Park et al. 2019b], SEAN [Zhu et al. 2020], StyleGAN2 [Karras et al. 2019a] and our SofGAN. For the corresponding quantitative evaluation and implementation details, please refer to Sec. 6.3 and 6.2.



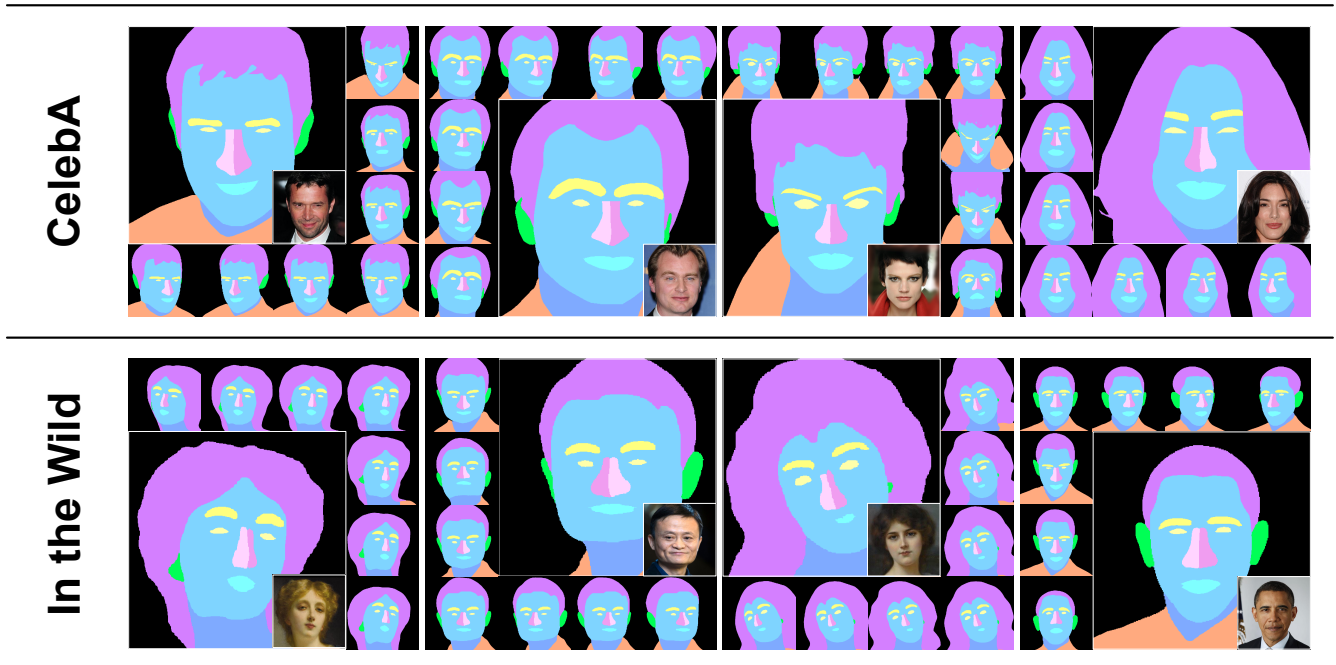


Fig. 21. Generating multi-view segmaps for real images. We first parse a monocular segmap from a given image, then project the segmap (central segmap) to SOF and generate multi-view segmaps (surrounding segmaps).

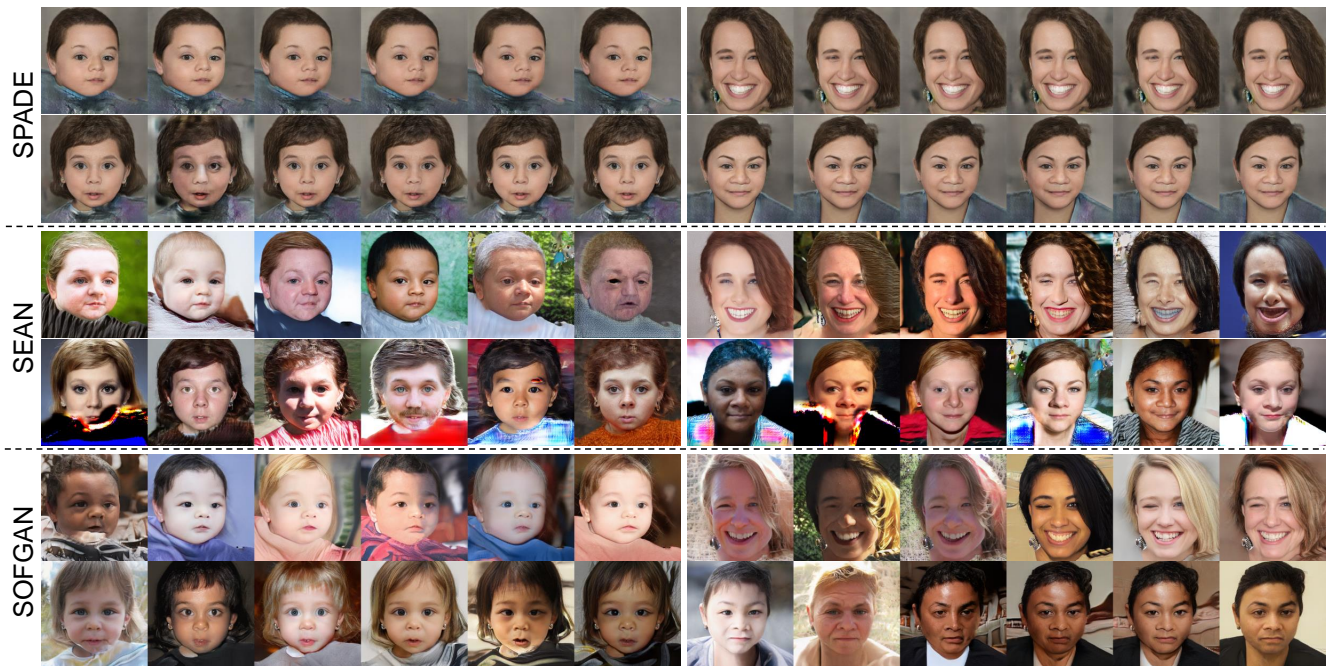


Fig. 22. Visual comparison of generation conditioned on same segmaps. Here we use the first four segmaps in the FFHQ dataset. Please refer to Sec. 6.3 for the experiment setting. The corresponding quantitative results are shown in Tab. 3. Note that the models in this comparison are trained for **only 800k** iterations and we set truncation to 1, thus we can observe some artifacts in the results.





Fig. 23. Free-viewpoint generation results. Images in the same row share the same geometry and texture style vectors, but use different camera poses. We can observe that our method is able to preserve shape and texture consistency even under large view angle variation.





Fig. 24. Results for global style adjustment.





Fig. 25. Results for global style adjustment.





Fig. 26. Results for regional style adjustment.



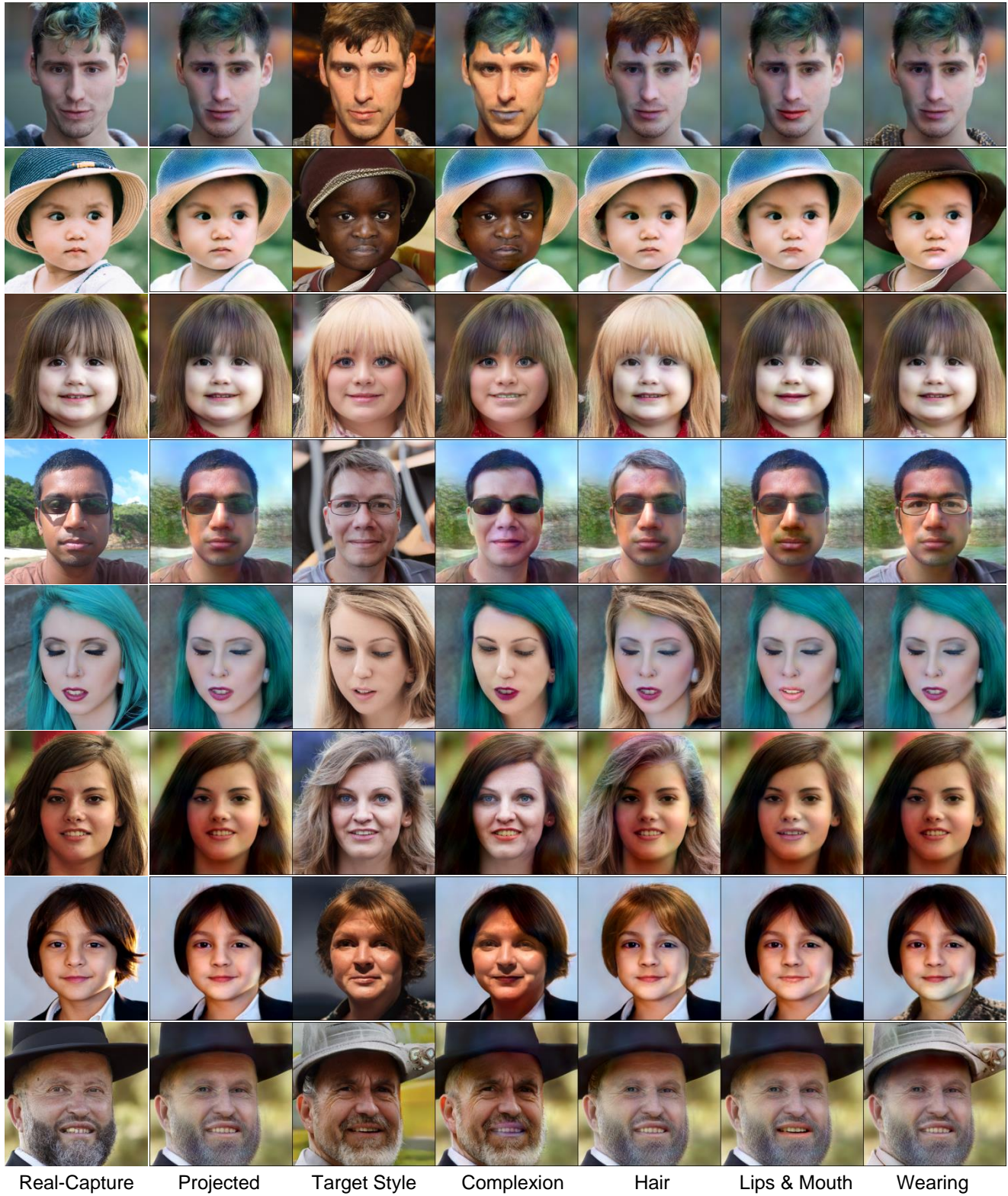


Fig. 27. Results for real-photo editing. We first project the real-captured photos into our texture space  $z^0$ , then we randomly sample a target texture style  $z^t_1$  from our texture space, furthermore, we regionally edit the project image via the regional style mixing scheme mentioned in Sec. 6.4.



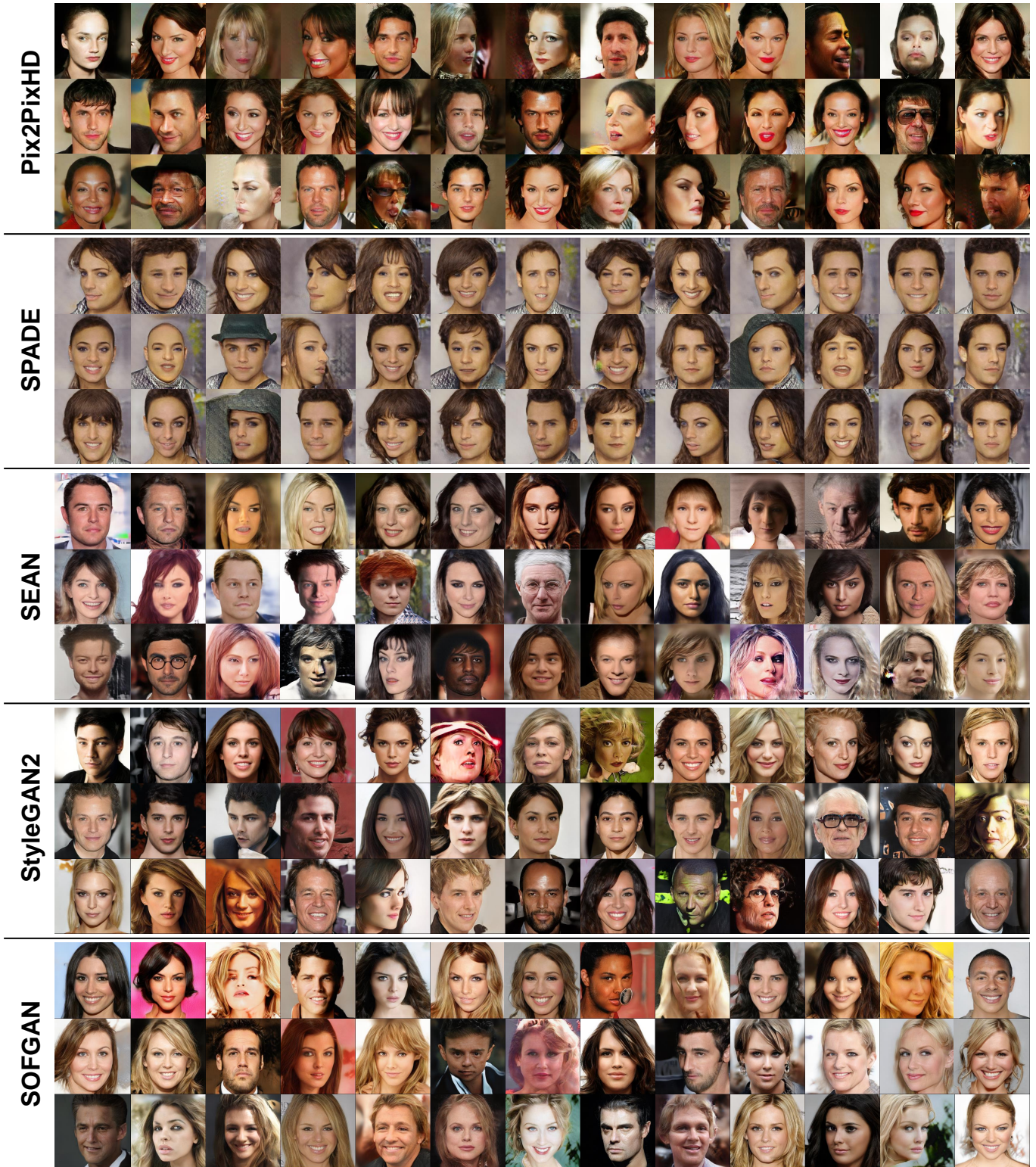


Fig. 28. Visual comparison on CelebAMask-HQ dataset [Lee et al. 2020]. Each model is trained with 800k images in resolution 512<sup>2</sup>.





Fig. 29. Visual comparison on FFHQ dataset [Karras et al. 2019b]. Each model is trained with 800k images in resolution  $512^2$ .