

# StyleSDF: High-Resolution 3D-Consistent Image and Geometry Generation

Roy Or-Ei<sup>1</sup>   Xuan Luo<sup>1</sup>   Mengyi Shan<sup>1</sup>   Eli Shechtman<sup>2</sup>  
 Jeong Joon Park<sup>3</sup>   Ira Kemelmacher-Shlizerman<sup>1</sup>  
<sup>1</sup>University of Washington   <sup>2</sup>Adobe Research   <sup>3</sup>Stanford University

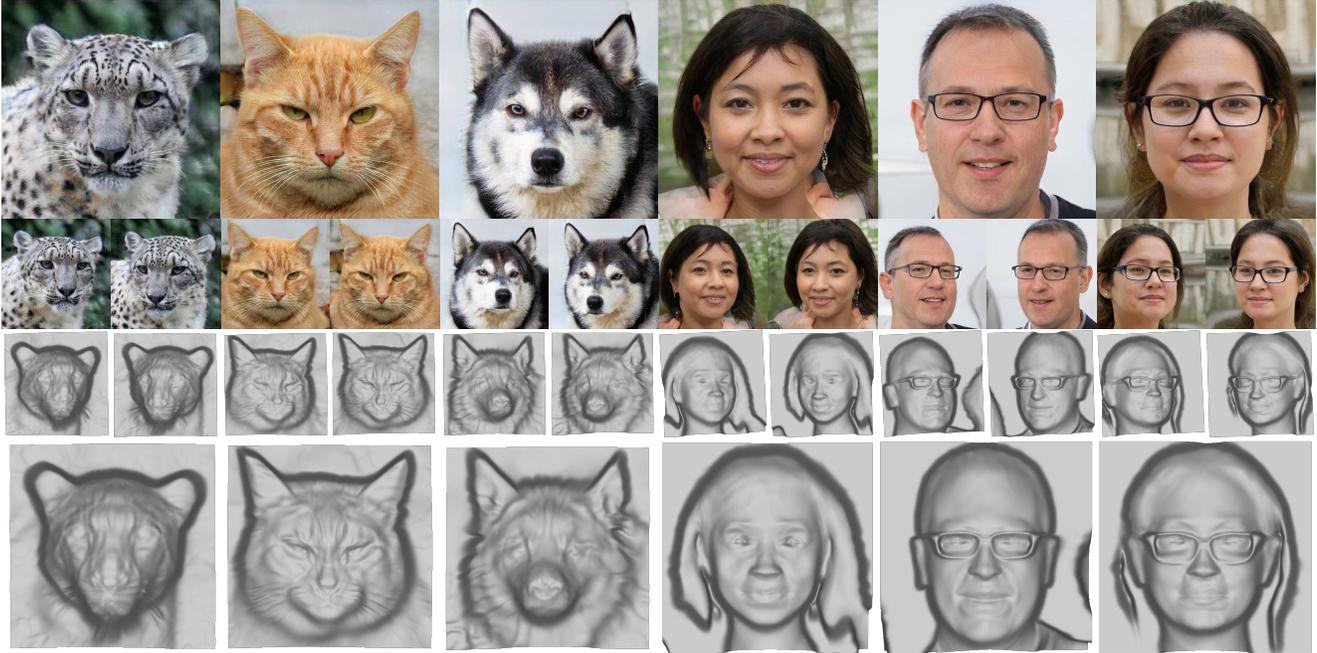


Figure 1. Our proposed framework—StyleSDF—learns to jointly generate high resolution, 3D-consistent images (top rows) along with their detailed view-consistent geometry represented with SDFs (depth maps in bottom rows), while being trained on single view RGB images.

## Abstract

We introduce a high resolution, 3D-consistent image and shape generation technique which we call StyleSDF. Our method is trained on single-view RGB data only, and stands on the shoulders of StyleGAN2 for image generation, while solving two main challenges in 3D-aware GANs: 1) high-resolution, view-consistent generation of the RGB images, and 2) detailed 3D shape. We achieve this by merging a SDF-based 3D representation with a style-based 2D generator. Our 3D implicit network renders low-resolution feature maps, from which the style-based network generates view-consistent,  $1024 \times 1024$  images. Notably, our SDF-based 3D modeling defines detailed 3D surfaces, leading to consistent volume rendering. Our method shows higher quality results compared to state of the art in terms of visual and geometric quality.

Project Page: <https://stylesdf.github.io/>

## 1. Introduction

StyleGAN architectures [37–39] have shown an unprecedented quality of RGB image generation. They are, however, designed to generate single RGB views rather than 3D content. In this paper, we introduce StyleSDF, a method for generating 3D-consistent  $1024 \times 1024$  RGB images and geometry, trained only on single-view RGB images.

Related 3D generative models [9, 52, 57, 61, 66] present shape and appearance synthesis via coordinate-based multi-layer-perceptrons (MLP). These works, however, often require 3D or multi-view data for supervision, which are difficult to collect, or are limited to low-resolution rendering outputs as they rely on expensive volumetric field sampling. Without multi-view supervision, 3D-aware GANs [9, 52, 61] typically use opacity fields as geometric proxy, forgoing well-defined surfaces, which results in low-quality depth maps that are inconsistent across views.

At the core of our architecture lies the SDF-based 3D volume renderer and the 2D StyleGAN generator. We use a coordinate-based MLP to model Signed Distance Fields (SDF) and radiance fields which render low resolution feature maps. These feature maps are then efficiently transformed into high-resolution images using the StyleGAN generator. Our model is trained with an adversarial loss that encourages the networks to generate realistic images from all sampled viewpoints, and an Eikonal loss that ensures proper SDF modeling. These losses automatically induce view-consistent, detailed 3D scenes, without 3D or multi-view supervision. The proposed framework effectively addresses the resolution and the view-inconsistency issues of existing 3D-aware GAN approaches that base on volume rendering. Our system design opens the door for interesting future research in vision and graphics that involves a latent space of high quality shape and appearance.

Our approach is evaluated on the FFHQ [38] and AFHQ [13] datasets. We demonstrate through extensive experiments that our system outperforms the state-of-the-art 3D-aware methods, measured by the quality of the generated images and surfaces, and their view-consistencies.

## 2. Related Work

In this section, we review related approaches in 2D image synthesis, 3D generative modeling, and 3D-aware image synthesis.

**Generative Adversarial Networks:** State-of-the-art Generative Adversarial Networks [21] (GANs) can synthesize high-resolution RGB images that are practically indistinguishable from real images [36–39]. Substantial work has been done in order to manipulate the generated images, by exploring meaningful latent space directions [1–3, 14, 26, 30, 62, 63, 67, 68], introducing contrastive learning [64], inverse graphics [79], exemplar images [33] or multiple input views [42]. While 2D latent space manipulation produces realistic results, these methods tend to lack explicit camera control, have no 3D understanding, require shape priors from 3DMM models [67, 68], or reconstruct the surface as a preprocessing step [42].

**Coordinate-based 3D Models:** While multiple 3D representations have been proposed for generative modeling [24, 73, 75], recent coordinate-based neural implicit models [10, 46, 57] stand out as an efficient, expressive, and differentiable representation.

Neural implicit representations (NIR) have been widely adopted for learning shape and appearance of objects [4, 11, 15, 22, 47, 53, 55, 59, 60], local parts [19, 20], and full 3D scenes [7, 12, 31, 58] from explicit 3D supervisions. Moreover, NIR approaches have been shown to be a powerful tool for reconstructing 3D structure from multi-view 2D supervision via fitting their 3D models to the multi-view images using differentiable rendering [48, 54, 66, 77].

Two recent seminal breakthroughs are NeRF [48] and SIREN [65]. NeRF introduced the use of volume rendering [34] for reconstructing a 3D scene as a combination of neural radiance and density fields to synthesize novel views. SIREN replaced the popular ReLU activation function with sine functions with modulated frequencies, showing great single scene fitting results. We refer readers to [70] for more comprehensive review.

**Single-View Supervised 3D-Aware GANs:** Rather than relying on 3D or multi-view supervisions, recent approaches aim at learning a 3D generative model from a set of unconstrained single-view images. These methods [9, 18, 27, 32, 44, 49–52, 61] typically optimize their 3D representations to render realistic 2D images from all randomly sampled viewpoints using adversarial loss.

Most inline with our work are methods that use implicit neural radiance fields for 3D-aware image and geometry generation (GRAF [61] and Pi-GAN [9]). However, these methods are limited to low-resolution outputs due to the high computational costs of the volume rendering. In addition, the use of density fields as proxy for geometry provides ample amount of leeway for the networks to produce realistic images while violating 3D consistency, leading to inconsistent volume rendering w.r.t. the camera viewpoints (the rendered RGB or depth images are not 3D-consistent).

To improve the surface quality, ShadeGAN [56] introduces a shading-guided pipeline, and GOF [74] gradually shrink the sampling region of each camera ray. However, the image output resolution ( $128 \times 128$ ) is still bounded by the computational burden of the volume rendering. GIRAFFE [52] proposed a dual stage rendering process. A backbone volume renderer generates low resolution feature maps ( $16 \times 16$ ) that are passed to a 2D CNN to generate outputs at  $256 \times 256$  resolution. Despite improved image quality, GIRAFFE outputs lack view consistency. The hairstyle, facial expression, and sometimes the object’s identity, are entangled with the camera viewpoint inputs, likely because 3D outputs at  $16 \times 16$  are not descriptive enough.

Concurrent works [8, 17, 25, 80] adopt two-stage rendering process or smart sampling procedures for high-resolution image generation, yet these works still do not model well-defined, view-consistent 3D geometry.

## 3. Algorithm

### 3.1. Overview

Our framework consists of two main components. A backbone conditional SDF volume renderer, and a 2D style-based generator [39]. Each component also has an accompanied mapping network [38] to map the input latent vector into modulation signals for each layer. An overview of our architecture can be seen in Figure 2.

To generate an image, we sample a latent vector  $\mathbf{z}$  from

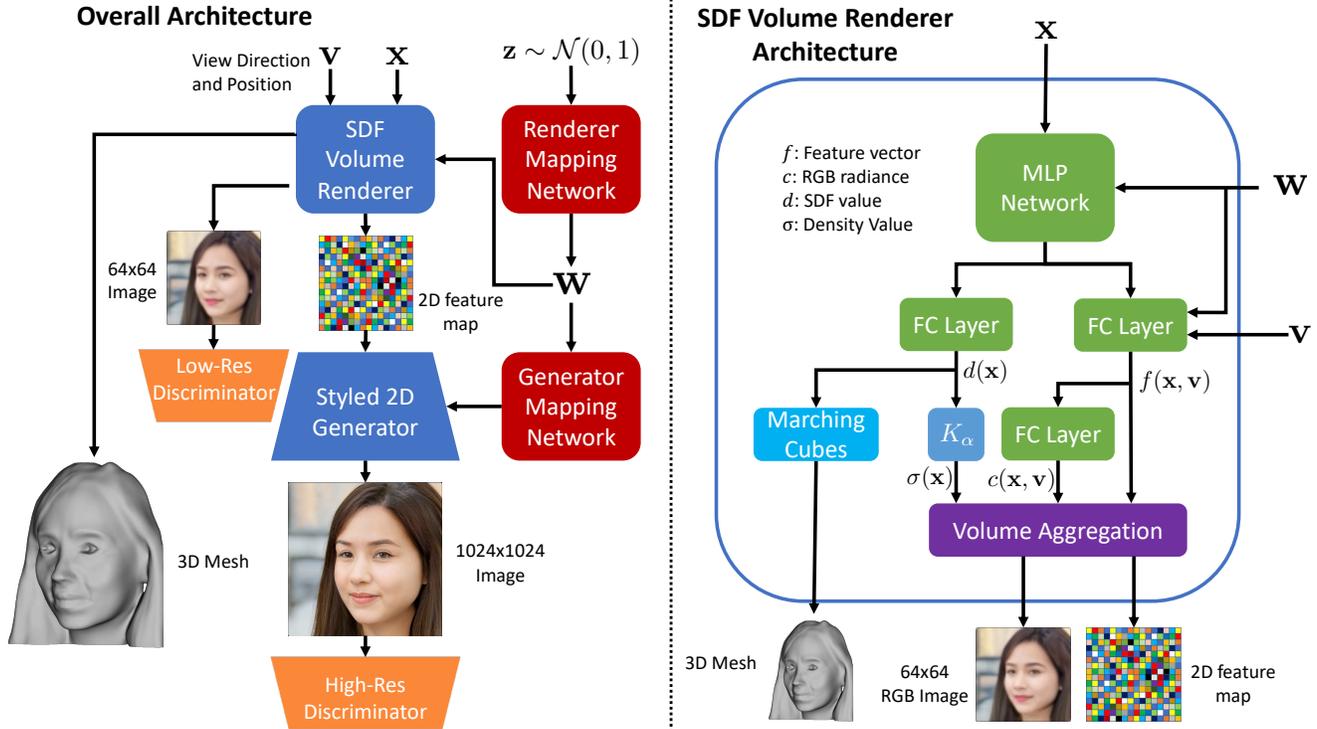


Figure 2. StyleSDF Architecture: (Left) Overall architecture: SDF volume renderer takes in a latent code and camera parameters, queries points and view directions in the volume, and projects the 3D surface features into the 2D view. The projected features are fed to the Styled 2D generator that creates the high resolution image. (Right) our SDF volume renderer jointly models volumetric SDF and radiance field, providing a well defined and view consistent geometry.

the unit normal distribution, and camera azimuth and elevation angles  $(\phi, \theta)$  from the dataset’s estimated object pose distribution. For simplicity, we assume that the camera is positioned on the unit sphere and directed towards the origin. Next, our volume renderer outputs the signed distance value, RGB color, and a 256 element feature vector for all the sampled volume points along the camera rays. We calculate the surface density for each sampled point from its SDF value and apply volume rendering [48] to project the 3D surface features into 2D feature map. The 2D generator then takes the feature map and generates the output image from the desired viewpoint. The 3D surface can be visualized with volume-rendered depths or with the mesh from marching-cubes algorithm [43].

### 3.2. SDF-based Volume Rendering

Our backbone volume renderer takes a 3D query point,  $\mathbf{x}$  and a viewing direction  $\mathbf{v}$ . Conditioned by the latent vector  $\mathbf{z}$ , it outputs an SDF value  $d(\mathbf{x}, \mathbf{z})$ , a view dependent color value  $\mathbf{c}(\mathbf{x}, \mathbf{v}, \mathbf{z})$ , and feature vector  $\mathbf{f}(\mathbf{x}, \mathbf{v}, \mathbf{z})$ . For clarity, we omit  $\mathbf{z}$  from hereon forward.

The SDF value indicates the distance of the queried point from the surface boundary, and the sign indicates whether the point is inside or outside of a watertight surface. As shown in VolSDF [76], the SDF can be serve as a proxy for the density function used for the traditional volume render-

ing [48]. Assuming a non-hollow surface, we convert the SDF value into the 3D density fields  $\sigma$ ,

$$\sigma(\mathbf{x}) = K_{\alpha}(d(\mathbf{x})) = \frac{1}{\alpha} \cdot \text{Sigmoid}\left(\frac{-d(\mathbf{x})}{\alpha}\right), \quad (1)$$

where  $\alpha$  is a learned parameter that controls the tightness of the density around the surface boundary.  $\alpha$  values that approach 0 represent a solid, sharp, object boundary, whereas larger  $\alpha$  values indicate a more “fluffy” object boundary. A large positive SDF value would drive the sigmoid function towards 0, meaning no density outside of the surface, and a high-magnitude negative SDF value would push the sigmoid towards 1, which means maximal density inside the surface.

We render low resolution  $64 \times 64$  feature maps and color images with volume rendering. For each pixel, we query points on a ray that originates at the camera position  $\mathbf{o}$ , and points at the camera direction  $\mathbf{r}(t) = \mathbf{o} + t\mathbf{v}$ . and calculate the RGB color and feature map as follows:

$$\begin{aligned} \mathbf{C}(\mathbf{r}) &= \int_{t_n}^{t_f} T(t)\sigma(\mathbf{r}(t))\mathbf{c}(\mathbf{r}(t), \mathbf{v})dt, \\ \mathbf{F}(\mathbf{r}) &= \int_{t_n}^{t_f} T(t)\sigma(\mathbf{r}(t))\mathbf{f}(\mathbf{r}(t), \mathbf{v})dt, \end{aligned} \quad (2)$$

$$\text{where } T(t) = \exp\left(-\int_{t_n}^t \sigma(\mathbf{r}(s))ds\right),$$

which we approximate with discrete sampling along rays.

Unlike NeRF [48] and other 3D-aware GANs such as Pi-GAN [9] and StyleNeRF [25] we do not use stratified sampling. Instead, we split  $[t_n, t_f]$  into  $N$  evenly-sized bins, draw a single offset term uniformly  $\delta \sim \mathcal{U}[0, \frac{t_f - t_n}{N}]$ , and sample  $N$  evenly-spaced points,

$$t_i = \frac{t_f - t_n}{N} \cdot i + \delta, \quad \text{where } i \in \{0, \dots, N - 1\}. \quad (3)$$

In addition, we forgo hierarchical sampling altogether, thereby reducing the number of samples by 50%. We discuss the merits of our sampling strategy in the supplementary material.

The incorporation of SDFs provides clear definition of the surface, allowing us to extract the mesh via Marching Cubes [43]. Moreover, the use of SDFs along with the related losses (Sec. 3.4.1) leads to higher quality geometry in terms of expressiveness and view-consistency (as shown in Sec. 4.4), even with a simplified volume sampling strategy.

The architecture of our volume renderer mostly matches that of Pi-GAN [9]. The mapping network consists of a 3 layer MLP with LeakyReLU activation and maps an input latent code  $\mathbf{z}$  into  $\mathbf{w}$  space and then generates frequency modulation,  $\gamma_i$ , and phase shift,  $\beta_i$ , for each layer of the volume renderer. The volume rendering network contains eight shared modulated FC layers with SIREN [65] activation:

$$\phi_i(x) = \sin(\gamma_i(W_i \cdot x + b_i) + \beta_i), \quad i \in \{0, \dots, 7\} \quad (4)$$

where  $W_i$  and  $b_i$  are the weight matrix and bias vector of the fully connected layers. The volume renderer then splits into two paths, the SDF path and the color path. The SDF path is implemented using a single FC layer denoted  $\phi_d$ . In the color path, the output of the last shared layer  $\phi_7$  is concatenated with the view direction input and passed into one additional FiLM siren layer [9]  $\phi_f$  followed by a single FC layer  $\phi_c$  that generates the color output. To summarize:

$$\begin{aligned} \sigma(\mathbf{x}) &= K_\alpha \circ \phi_d \circ \phi_7 \circ \dots \circ \phi_0(\mathbf{x}), \\ f(\mathbf{x}, \mathbf{v}) &= \phi_f(\phi_7 \circ \dots \circ \phi_0(\mathbf{x}), \mathbf{v}) \\ c(\mathbf{x}, \mathbf{v}) &= \phi_c \circ \phi_f. \end{aligned} \quad (5)$$

The output features of  $\phi_f$  are passed to the 2D style-based generator, and the generated low resolution color image is fed to a discriminator for supervision. The discriminator is identical to the Pi-GAN [9] discriminator.

We observed that using view-dependent color  $c(\mathbf{x}, \mathbf{v})$  tends to make the networks overfit to biases in the dataset. For instance, people in FFHQ [38] tend to smile more when facing the camera. This makes the facial expression change with the viewpoint although the geometry remains consistent. However, when we removed view-dependent color, the model did not converge. Therefore, to get view consistent images, we train our model with view dependent color, but fix the view direction  $\mathbf{v}$  to the frontal view during inference.

### 3.3. High-Resolution Image Generation

Unlike NeRF [48], where the reconstruction loss is computed individually for each ray, adversarial training needs a full image to be present. Therefore, scaling a pure volume renderer to high-resolution quickly becomes untractable, as we need to sample over  $10^7$  queries to render a single  $1024 \times 1024$  image. As such, we seek to fuse a volume renderer with the StyleGAN2 network that has a proven capabilities of synthesizing high-resolution 2D images.

To combine the two architectures, we truncate the early layers of the StyleGAN2 generator up until the  $64 \times 64$  layer and feed the generator with the  $64 \times 64$  feature maps generated by the backbone volume renderer. In addition, we cut StyleGAN2’s mapping network from eight layers to five layers, and feed it with the  $\mathbf{w}$  latent code from the volume renderer’s mapping network, instead of the original latent vector  $\mathbf{z}$ . The discriminator is left unchanged.

This design choice allows us to enjoy the best of both worlds. The volume renderer learns the underline geometry, explicitly disentangles the object’s pose from it’s appearance, and enables full control of the camera position during inference. The StyleGAN2 generator upsamples the low resolution feature maps, adds high frequency details, and mimics complex light transport effects such as sub-surface scattering and inter-reflections that are difficult to model with the low-resolution volume renderer.

### 3.4. Training

We employ a two-stage training procedure. First we train only the SDF-based volume renderer, then we freeze the volume renderer weights, and train the StyleGAN generator.

#### 3.4.1 Volume Renderer training

We use the non-saturating GAN loss with R1 regularization [45], denoted  $\mathcal{L}_{adv}$ , to train our volume renderer. On top of that, we use 3 additional regularization terms.

**Pose Alignment Loss:** This loss is designed to make sure that all the generated objects are globally aligned. On top of predicting whether the image is real or fake, the discriminator also tries to predict the two input camera angles  $(\phi, \theta)$ . We penalize the prediction error using a smoothed L1 loss:

$$\mathcal{L}_{view} = \begin{cases} (\hat{\theta} - \theta)^2 & \text{if } |\hat{\theta} - \theta| \leq 1 \\ |\hat{\theta} - \theta| & \text{otherwise} \end{cases}. \quad (6)$$

This loss is applied on both view angles for the generator and the discriminator, however, since we don’t have ground truth pose data for the original dataset, this loss is only applied to the fake images in the discriminator pass.

**Eikonal Loss:** This term ensures that the learned SDF is physically valid [23]:

$$\mathcal{L}_{eik} = \mathbb{E}_{\mathbf{x}}(\|\nabla d(\mathbf{x})\|_2 - 1)^2. \quad (7)$$

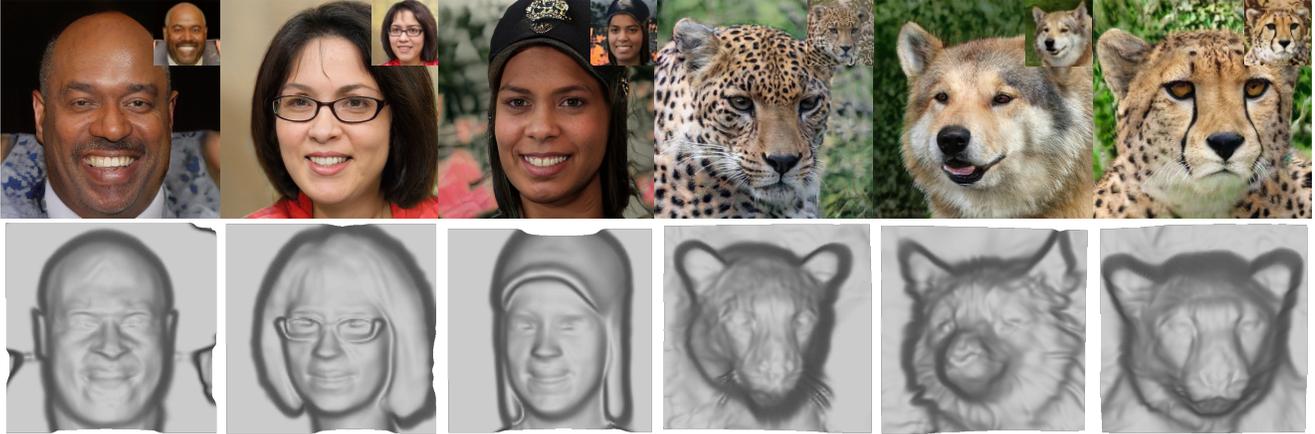


Figure 3. Generated high-res RGB images (top), low-res volume rendered images (inset) and depth maps (bottom) for the same view. The  $64 \times 64$  volume rendering output features are passed to the StyleGAN generator for high-resolution RGBs. Note that the object identities and structures are preserved between the image pairs. Furthermore, as can be seen in the jaguar and cheetah examples, the StyleGAN generator occasionally corrects badly modeled background signal from the volume renderer.

**Minimal Surface Loss:** We encourage the 3D network to describe the scenes with minimal volume of zero-crossings to prevent spurious and non-visible surfaces from being formed within the scenes. That is, we penalize the SDF values that are close to zero:

$$\mathcal{L}_{surf} = \mathbb{E}_{\mathbf{x}} (\exp(-100|d(x)|)). \quad (8)$$

The overall loss function is then,

$$\mathcal{L}_{vol} = \mathcal{L}_{adv} + \lambda_{view} \mathcal{L}_{view} + \lambda_{eik} \mathcal{L}_{eik} + \lambda_{surf} \mathcal{L}_{surf}, \quad (9)$$

where  $\lambda_{view} = 15$ ,  $\lambda_{eik} = 0.1$ , and  $\lambda_{surf} = 0.05$ . The weight of the R1 loss is set according to the dataset.

### 3.4.2 Styled Generator Training

We train our Styled generator with the same losses and optimizer parameters as the original implementation, a non saturating adversarial loss, R1 regularization, and path regularization. As in the volume renderer training, we set the weight of the R1 regularization according to the dataset.

While it is possible to have a reconstruction loss between the low-resolution and high-resolution output images, we find that the inductive bias of the 2D convolutional architecture and the sharing of style codes is strong enough to preserve important structures and identities between the images (Fig. 3).

## 4. Experiments

### 4.1. Datasets & Baselines

We train and evaluate our model on the FFHQ [38] and AFHQ [13] datasets. FFHQ contains 70,000 images of diverse human faces at  $1024 \times 1024$  resolution, which are

centered and aligned according to the procedure introduced in Karras *et al.* [36]. The AFHQ dataset consists of 15,630 images of cats, dogs and wild animals at  $512 \times 512$  resolution. Note that the AFHQ images are not aligned and contain diverse animal species, posing a significant challenge to StyleSDF.

We compare our method against the state-of-the-art 3D-aware GAN baselines, GIRAFFE [52], PiGAN [9], GRAF [61] and HoloGAN [49], on the above datasets by measuring the quality of the generated images, shapes, and rendering consistency.

### 4.2. Qualitative Evaluations

**Comparison to Baseline Methods:** We compare the visual quality of our images to the baseline methods by rendering the same identity (latent code) from 4 different viewpoints, results are shown in Figure 4. To compare the quality of the underlying geometry, we also show the surfaces extracted by marching cubes from StyleSDF, Pi-GAN, and GRAF (Note that GIRAFFE and HoloGAN pipelines do not generate shapes). Our method generates superior images as well as more detailed 3D shapes. Additional generation results from our method can be seen in Figures 1 and 3.

**Novel View Synthesis:** Since our method learns strong 3D shape priors, it can generate images from viewpoints that are not well represented in the dataset distribution. Examples of out-of-distribution view synthesis are displayed in Figure 5.

**Video Results:** We urge readers to view our [project's website](#) that includes a larger set of results and videos to better appreciate the multi-view capabilities of StyleSDF.

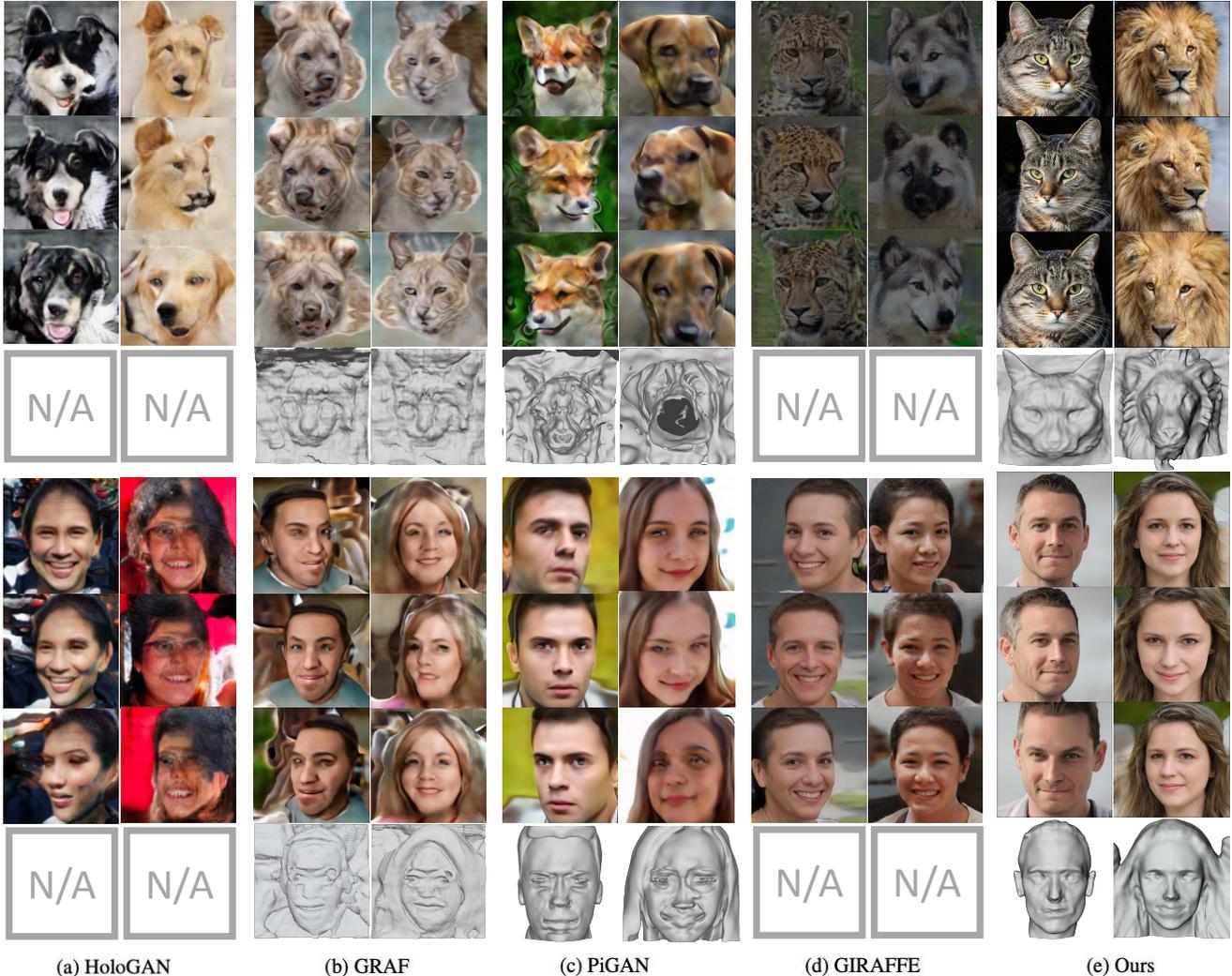


Figure 4. Qualitative image and geometry comparisons. We compare our sample renderings and corresponding 3D meshes against the state-of-the-art 3D-aware GAN approaches ([9, 49, 52, 61]). Note that HoloGAN and GIRAFFE are unable to create 3D mesh from their representations. Both HoloGAN (a) and GRAF (b) produce renderings that are of lower quality. The 3D mesh reconstructed from PiGAN’s learned opacity fields reveal noticeable artifacts (c). While GIRAFFE (d) produces realistic low-resolution images, the identity of the person often changes with the viewpoints. StyleSDF (d) produces  $1024 \times 1024$  realistic view consistent RGB, while also generating high quality 3D. Best viewed digitally.

### 4.3. Quantitative Image Evaluations

We evaluate the visual quality and the diversity of the generated images using the Frchet Inception Distance (FID) [28] and Kernel Inception Distance (KID) [6]. We compare our scores against the aforementioned baseline models on the FFHQ and AFHQ datasets.

All the baseline models are trained following their given pipelines to generate  $256 \times 256$  images, with the exception of Pi-GAN, which is trained on  $128 \times 128$  images and renders  $256 \times 256$  images at inference time. The results, summarized in Table 1, show that StyleSDF performs consistently better than all the baselines in terms of visual quality. It is also on par with reported scores from concurrent

works such as StyleNerf [25] and CIPS-3D [80].

### 4.4. Volume Rendering Consistency

Volume rendering has emerged as an essential technique to differentially optimize a volumetric field from 2D images, as its wide-coverage point sampling leads to stable gradient-flow during training. Notably, volume rendering excels at modeling thin surfaces or transparent objects, e.g., human hairs, which are difficult to model with explicit surfaces, e.g., 3D meshes.

However, we notice that the volume rendering of existing 3D-aware GANs [9, 61] using unregularized opacity fields severely lacks view-consistency due to the absence of multi-

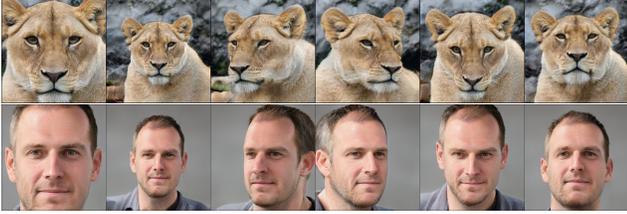


Figure 5. Out-of-distribution view synthesis (field of view and camera angles). Although StyleSDF was trained with a fixed field of view, increasing and decreasing FOV by 25% (columns 1-2) still looks realistic. Similarly with 1.5 standard deviations of the camera angles distribution used for training (columns 3-6).

Dataset:	FFHQ		AFHQ	
	FID	KID	FID	KID
HoloGAN	90.9	75.5	95.6	77.5
GRAF	79.2	55.0	129.5	85.1
PiGAN	83.0	85.8	52.4	30.7
GIRAFFE	31.2	20.1	33.5	15.1
Ours	<b>11.5</b>	<b>2.65</b>	<b>12.8</b>	<b>4.47</b>

Table 1. FID and KID evaluations. All datasets were evaluated at a resolution of  $256 \times 256$ . Our method demonstrates the best performance. Note that we report  $KID \times 1000$  for simplicity.

Dataset:	FFHQ	AFHQ
PiGAN	11.04	8.66
Ours	<b>0.40</b>	<b>0.63</b>

Table 2. Depth consistency results. We measure the average modified Chamfer distance (Eq. (10)) over 1,000 random pairs of depth maps for each dataset. Each pair contains one frontal view depth map and one side view depth map. Our method demonstrates significantly stronger consistency (see Fig. 6).

view supervision. That is, depth values, computed as the expected termination distance of each camera ray [16, 48], from different viewpoints do not consistently overlap in the global coordinate. This means that neural implicit features are evaluated at inconsistent 3D locations, undermining the inductive bias of the implicit 3D representation for view-consistent renderings. As such, we measure and compare the depth map consistency across viewpoints to gauge the quality of volume rendering for each system.

We sample 1,000 identities, render their  $128 \times 128$  depth maps from the frontal view and a fixed side view, and compute the alignment between the two views. The depth value is defined as the expected termination distance of 128 uniformly sampled points along each ray. Note that we remove non-terminating rays whose accumulated opacity is below 0.5. We set the side viewpoint to be  $1.5 \times$  the standard de-

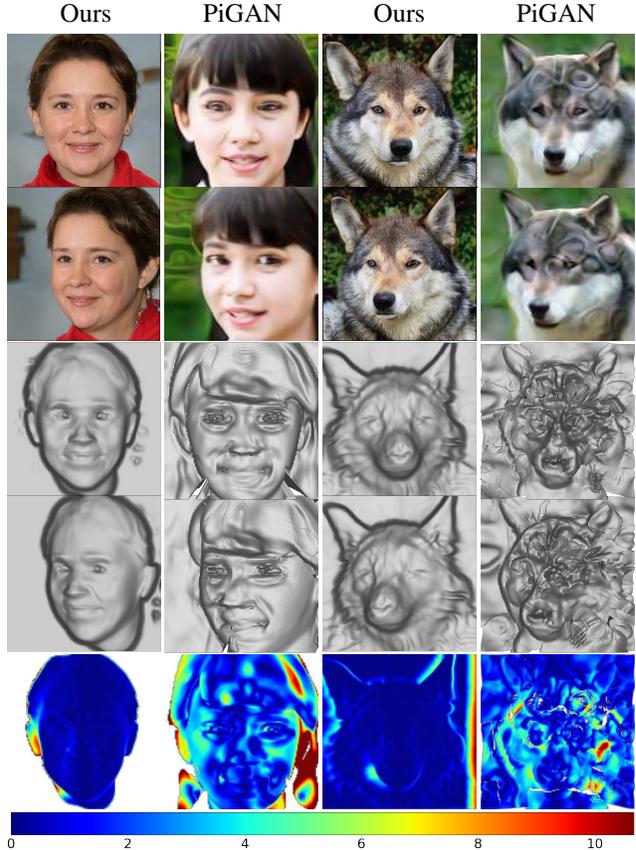


Figure 6. Visual comparison of depth consistency. We visualize the nearest neighbour distances (in sample bin units) from the frontal depth maps to side-view depth maps. Our SDF-based technique significantly improves depth consistency compared to the baseline.

viation of the azimuth distribution in training. See supplementary for more experiment details.

To measure the alignment errors between the depth points, we adopt a modified Chamfer distance metric. I.e., we replace the usual mean distance definition with the median of the distances to nearest points,

$$CD(S_1, S_2) = \text{med} \min_{x \in S_1} \min_{y \in S_2} \|x - y\|_2^2 + \text{med} \min_{y \in S_2} \min_{x \in S_1} \|x - y\|_2^2, \quad (10)$$

for some point sets  $S_1$  and  $S_2$ . This metric is more robust to outliers that come from occlusion and background mismatch that we are not interested in measuring. To put the metric at scale, we normalize the distances by the volume sampling bin size.

As shown in Table 2, our use of SDF representation dramatically improves depth consistency compared to the strongest current baseline PiGAN [9]. Figure 6 shows the sample depth map pairs used for the evaluation and the error visualizations (in terms of distance to the closest point). The color map shows that our depth maps align well except for the occluded regions and backgrounds. In contrast, PiGAN

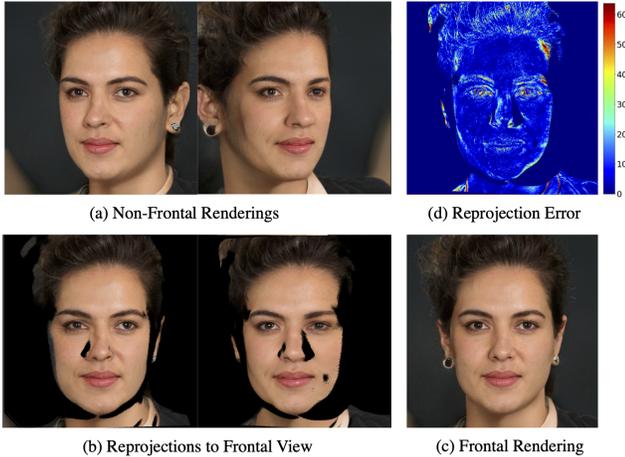


Figure 7. RGB rendering view-consistency. We render two side views (a) and project them to the frontal view (b) using the depth maps rendered from each views, ignoring occluded pixels. Note the high similarity between the reprojected images and the rendered frontal view (c), as can be seen from the error map (d). The error map shows mean absolute pixel difference for RGB channels (0-255) for the right side-view image. The errors are mostly from regions with high frequency textures and geometry (e.g., ear, hair), or occlusion boundaries (right forehead).

depth maps show significant noise and spurious concave regions (e.g., nose of the dog).

Moreover, we show that our consistent volume rendering naturally leads to high view-consistency for our RGB renderings. As shown in Fig. 7, we visualize the reprojected side-view renderings to the frontal view, using the depth values from volume rendering. The reprojected pixels closely match those of the original frontal view, indicating that our high-res multi-view RGB renderings and depth maps are all consistent to each other. Refer to supplementary for more detailed experiments.

## 5. Limitations & Future Work

StyleSDF might exhibit minor aliasing and flickering, e.g., in teeth area. We leave it for future work since we expect those two to be corrected similarly to Mip-NeRF [5] and Alias-free StyleGAN [37]. See example at left two columns of Figure 8. Specularities or other strong lighting effects currently introduce depth dents since StyleSDF might find it hard to disambiguate with no multi-view data (Figure 8 third column from the left). Adjusting the losses to include those effects is left for future work. Similarly, we do not currently separate foreground from background and use a single SDF for the whole image. Figure 8 (right column) shows how the cat’s face is rendered properly, but the transition to the background is too abrupt, potentially diminishing photorealism. A potential solution could be adding an additional volume renderer to model the back-

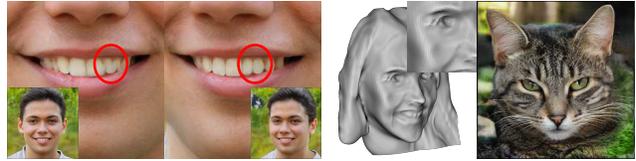


Figure 8. Limitations: potential aliasing artifacts, e.g., in teeth (left two columns). Specularities and shadows may create artifacts (3rd column from the left, cheek and eyes area), high curvatures are enhanced with radiance scaling filter [71]. Inconsistencies in background might decrease photorealism (right column).

ground as suggested in NeRF++ [78].

Finally, one may consider two improvements to the algorithm. First one is training the two parts as a single end-to-end framework, instead of the current two networks. In such case the StyleGAN2 discriminator would send proper gradients back to the volume renderer to produce optimal feature maps, which might lead to even more refined geometry. However, end-to-end training poses a trade-off. The increased GPU memory consumption of this setup would require either a decreased batch size, which might hurt the overall performance, or increased training time if we keep the batch size and accumulate gradients. Second improvement could be to create a volume sampling strategy tied to SDF’s surface boundary (to reduce the number of query points at each forward pass) and eliminate the need for a 2D CNN that upsamples feature maps. That would tie 3D geometry directly to the high resolution image.

## 6. Conclusions

We introduced StyleSDF, a method that can render 1024x1024 view-consistent images along with the detailed underlying geometry. The proposed architecture combines SDF-based volume renderer and a 2D StyleGAN network and is trained to generate realistic images for all sampled viewpoints via adversarial loss, naturally inducing view-consistent 3D scenes. StyleSDF represents and learns complex 3D shape and appearance without multi-view or 3D supervision, requiring only a dataset of single-view images, suggesting a new route ahead for neural 3D content generation, editing, and reconstruction.

## Acknowledgements

We wish to thank Aleksander Holynski for his valuable advice. This work was supported by the UW Reality Lab, Meta, Google, Futurewei, Amazon and Adobe. J.J. Park was supported by the Apple fellowship.

## References

- [1] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan: How to embed images into the stylegan latent

- space? In *ICCV*, pages 4432–4441, 2019. 2
- [2] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan++: How to edit the embedded images? In *CVPR*, pages 8296–8305, 2020. 2
- [3] Rameen Abdal, Peihao Zhu, Niloy J Mitra, and Peter Wonka. Styleflow: Attribute-conditioned exploration of stylegan-generated images using conditional continuous normalizing flows. *ACM TOG*, 40(3):1–21, 2021. 2
- [4] Matan Atzmon and Yaron Lipman. Sal: Sign agnostic learning of shapes from raw data. In *CVPR*, pages 2565–2574, 2020. 2
- [5] Jonathan T. Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P. Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *ICCV*, pages 5855–5864, October 2021. 8, 17
- [6] Mikołaj Bińkowski, Danica J. Sutherland, Michael Arbel, and Arthur Gretton. Demystifying mmd gans. *arXiv preprint arXiv:1801.01401*, 2018. 6
- [7] Rohan Chabra, Jan E Lenssen, Eddy Ilg, Tanner Schmidt, Julian Straub, Steven Lovegrove, and Richard Newcombe. Deep local shapes: Learning local sdf priors for detailed 3d reconstruction. In *ECCV*, pages 608–625. Springer, 2020. 2
- [8] Eric R. Chan, Connor Z. Lin, Matthew A. Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas Guibas, Jonathan Tremblay, Sameh Khamis, Tero Karras, and Gordon Wetzstein. Efficient geometry-aware 3D generative adversarial networks. In *CVPR*, 2022. 2
- [9] Eric R. Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *CVPR*, pages 5799–5809, 2021. 1, 2, 4, 5, 6, 7, 12, 16
- [10] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *CVPR*, pages 5939–5948, 2019. 2
- [11] Julian Chibane, Thiemo Alldieck, and Gerard Pons-Moll. Implicit functions in feature space for 3d shape reconstruction and completion. In *CVPR*. IEEE, jun 2020. 2
- [12] Julian Chibane, Aymen Mir, and Gerard Pons-Moll. Neural unsigned distance fields for implicit function learning. In *NeurIPS*, December 2020. 2
- [13] Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. Stargan v2: Diverse image synthesis for multiple domains. In *CVPR*, pages 8188–8197, 2020. 2, 5
- [14] Edo Collins, Raja Bala, Bob Price, and Sabine Susstrunk. Editing in style: Uncovering the local semantics of gans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5771–5780, 2020. 2
- [15] Thomas Davies, Derek Nowrouzezahrai, and Alec Jacobson. Overfit neural networks as a compact shape representation. *arXiv preprint arXiv:2009.09808*, 2020. 2
- [16] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised nerf: Fewer views and faster training for free. *arXiv preprint arXiv:2107.02791*, 2021. 7
- [17] Yu Deng, Jiaolong Yang, Jianfeng Xiang, and Xin Tong. Gram: Generative radiance manifolds for 3d-aware image generation. In *CVPR*, 2022. 2
- [18] Matheus Gadelha, Subhansu Maji, and Rui Wang. 3d shape induction from 2d views of multiple objects. In *3DV*, pages 402–411. IEEE, 2017. 2
- [19] Kyle Genova, Forrester Cole, Avneesh Sud, Aaron Sarna, and Thomas Funkhouser. Local deep implicit functions for 3d shape. In *CVPR*, pages 4857–4866, 2020. 2
- [20] Kyle Genova, Forrester Cole, Daniel Vlasic, Aaron Sarna, William T Freeman, and Thomas Funkhouser. Learning shape templates with structured implicit functions. In *CVPR*, pages 7154–7164, 2019. 2
- [21] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NeurIPS*, pages 2672–2680, 2014. 2
- [22] Amos Gropp, Lior Yariv, Niv Haim, Matan Atzmon, and Yaron Lipman. Implicit geometric regularization for learning shapes. In *ICML*, pages 3569–3579, 2020. 2
- [23] Amos Gropp, Lior Yariv, Niv Haim, Matan Atzmon, and Yaron Lipman. Implicit geometric regularization for learning shapes. *arXiv preprint arXiv:2002.10099*, 2020. 4
- [24] Thibault Groueix, Matthew Fisher, Vladimir G Kim, Bryan C Russell, and Mathieu Aubry. A papier-mâché approach to learning 3d surface generation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 216–224, 2018. 2
- [25] Jiatao Gu, Lingjie Liu, Peng Wang, and Christian Theobalt. Stylenerf: A style-based 3d-aware generator for high-resolution image synthesis. *ICLR*, 2022. 2, 4, 6
- [26] Erik Härkönen, Aaron Hertzmann, Jaakko Lehtinen, and Sylvain Paris. Ganspace: Discovering interpretable gan controls. In *NeurIPS*, 2020. 2
- [27] Paul Henderson, Vagia Tsiminaki, and Christoph H Lampert. Leveraging 2d data to learn textured 3d mesh generation. In *CVPR*, pages 7498–7507, 2020. 2
- [28] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *NeurIPS*, page 6629–6640, 2017. 6
- [29] K Hill and J White. Designed to deceive: Do these people look real to you. *New York Times*, 2020. 12
- [30] Ali Jahanian, Lucy Chai, and Phillip Isola. On the “steerability” of generative adversarial networks. In *ICLR*, 2020. 2
- [31] Chiyu Jiang, Avneesh Sud, Ameesh Makadia, Jingwei Huang, Matthias Nießner, Thomas Funkhouser, et al. Local implicit grid representations for 3d scenes. In *CVPR*, pages 6001–6010, 2020. 2
- [32] Danilo Jimenez Rezende, SM Eslami, Shakir Mohamed, Peter Battaglia, Max Jaderberg, and Nicolas Heess. Unsupervised learning of 3d structure from images. *NeurIPS*, 29:4996–5004, 2016. 2
- [33] Omer Kafri, Or Patashnik, Yuval Alaluf, and Daniel Cohen-Or. Stylefusion: A generative model for disentangling spatial segments. *arXiv preprint arXiv:2107.07437*, 2021. 2
- [34] James T Kajiya and Brian P Von Herzen. Ray tracing volume densities. *ACM SIGGRAPH computer graphics*, 18(3):165–174, 1984. 2

- [35] Kimmo Karkkainen and Jungseock Joo. Fairface: Face attribute dataset for balanced race, gender, and age for bias measurement and mitigation. In *WACV*, pages 1548–1558, 2021. [1](#), [2](#)
- [36] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of GANs for improved quality, stability, and variation. In *ICLR*, 2018. [2](#), [5](#)
- [37] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Alias-free generative adversarial networks. *arXiv preprint arXiv:2106.12423*, 2021. [1](#), [2](#), [8](#), [17](#)
- [38] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *CVPR*, pages 4401–4410, 2019. [1](#), [2](#), [4](#), [5](#), [13](#), [17](#)
- [39] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *CVPR*, pages 8110–8119, 2020. [1](#), [2](#), [16](#)
- [40] Hyeonwoo Kim, Pablo Garrido, Ayush Tewari, Weipeng Xu, Justus Thies, Matthias Niessner, Patrick Pérez, Christian Richardt, Michael Zollhöfer, and Christian Theobalt. Deep video portraits. *ACM Transactions on Graphics (TOG)*, 37(4):1–14, 2018. [12](#)
- [41] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. [15](#)
- [42] Thomas Leimkühler and George Drettakis. Freestylegan: Free-view editable portrait rendering with the camera manifold. *arXiv preprint arXiv:2109.09378*, 2021. [2](#)
- [43] William E. Lorensen and Harvey E. Cline. Marching cubes: A high resolution 3d surface construction algorithm. *ACM TOG*, 21(4):163–169, 1987. [3](#), [4](#)
- [44] Sebastian Lunz, Yingzhen Li, Andrew Fitzgibbon, and Nate Kushman. Inverse graphics gan: Learning to generate 3d shapes from unstructured 2d data. *arXiv preprint arXiv:2002.12674*, 2020. [2](#)
- [45] Lars Mescheder, Andreas Geiger, and Sebastian Nowozin. Which training methods for gans do actually converge? In *ICML*, pages 3481–3490, 2018. [4](#)
- [46] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *CVPR*, pages 4460–4470, 2019. [2](#)
- [47] Mateusz Michalkiewicz, Jhony K Pontes, Dominic Jack, Mahsa Baktashmotlagh, and Anders Eriksson. Implicit surface representations as layers in neural networks. In *ICCV*, pages 4743–4752, 2019. [2](#)
- [48] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, pages 405–421. Springer, 2020. [2](#), [3](#), [4](#), [7](#), [16](#)
- [49] Thu Nguyen-Phuoc, Chuan Li, Lucas Theis, Christian Richardt, and Yong-Liang Yang. Hologan: Unsupervised learning of 3d representations from natural images. In *ICCV*, pages 7588–7597, 2019. [2](#), [5](#), [6](#)
- [50] Thu Nguyen-Phuoc, Christian Richardt, Long Mai, Yong-Liang Yang, and Niloy Mitra. Blockgan: Learning 3d object-aware scene representations from unlabelled images. In *NeurIPS*, 2020. [2](#)
- [51] Michael Niemeyer and Andreas Geiger. Campari: Camera-aware decomposed generative neural radiance fields. In *3DV*, pages 951–961. IEEE, 2021. [2](#)
- [52] Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. In *CVPR*, pages 11453–11464, 2021. [1](#), [2](#), [5](#), [6](#)
- [53] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Occupancy flow: 4d reconstruction by learning particle dynamics. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5379–5389, 2019. [2](#)
- [54] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. In *CVPR*, pages 3504–3515, 2020. [2](#)
- [55] Michael Oechsle, Lars Mescheder, Michael Niemeyer, Thilo Strauss, and Andreas Geiger. Texture fields: Learning texture representations in function space. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4531–4540, 2019. [2](#)
- [56] Xingang Pan, Xudong Xu, Chen Change Loy, Christian Theobalt, and Bo Dai. A shading-guided generative implicit model for shape-accurate 3d-aware image synthesis. In *NeurIPS*, 2021. [2](#)
- [57] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *CVPR*, pages 165–174, 2019. [1](#), [2](#)
- [58] Songyou Peng, Michael Niemeyer, Lars Mescheder, Marc Pollefeys, and Andreas Geiger. Convolutional occupancy networks. In *ECCV*, pages 523–540. Springer, 2020. [2](#)
- [59] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In *ICCV*, pages 2304–2314, 2019. [2](#)
- [60] Shunsuke Saito, Tomas Simon, Jason Saragih, and Hanbyul Joo. PifuHD: Multi-level pixel-aligned implicit function for high-resolution 3d human digitization. In *CVPR*, pages 84–93, 2020. [2](#)
- [61] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. Graf: Generative radiance fields for 3d-aware image synthesis. In *NeurIPS*, 2020. [1](#), [2](#), [5](#), [6](#)
- [62] Yujun Shen, Jinjin Gu, Xiaoou Tang, and Bolei Zhou. Interpreting the latent space of gans for semantic face editing. In *CVPR*, pages 9243–9252, 2020. [2](#)
- [63] Yujun Shen and Bolei Zhou. Closed-form factorization of latent semantics in gans. In *CVPR*, pages 1532–1540, 2021. [2](#)
- [64] Alon Shoshan, Nadav Bhonker, Igor Kviatkovsky, and Gérard Medioni. Gan-control: Explicitly controllable gans. In *ICCV*, pages 14083–14093, 2021. [2](#)
- [65] Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. Implicit neural representa-

- tions with periodic activation functions. *Advances in Neural Information Processing Systems*, 33, 2020. 2, 4
- [66] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. *NeurIPS*, 2019. 1, 2
- [67] Ayush Tewari, Mohamed Elgharib, Florian Bernard, Hans-Peter Seidel, Patrick Pérez, Michael Zollhöfer, and Christian Theobalt. Pie: Portrait image embedding for semantic control. *ACM TOG*, 39(6):1–14, 2020. 2
- [68] Ayush Tewari, Mohamed Elgharib, Gaurav Bharaj, Florian Bernard, Hans-Peter Seidel, Patrick Pérez, Michael Zollhofer, and Christian Theobalt. Stylerig: Rigging stylegan for 3d control over portrait images. In *CVPR*, pages 6142–6151, 2020. 2
- [69] Ayush Tewari, Ohad Fried, Justus Thies, Vincent Sitzmann, Stephen Lombardi, Kalyan Sunkavalli, Ricardo Martin-Brualla, Tomas Simon, Jason Saragih, Matthias Nießner, et al. State of the art on neural rendering. In *Computer Graphics Forum*, volume 39, pages 701–727. Wiley Online Library, 2020. 12
- [70] Ayush Tewari, Justus Thies, Ben Mildenhall, Pratul Srinivasan, Edgar Tretschk, Yifan Wang, Christoph Lassner, Vincent Sitzmann, Ricardo Martin-Brualla, Stephen Lombardi, Tomas Simon, Christian Theobalt, Matthias Niessner, Jonathan T. Barron, Gordon Wetzstein, Michael Zollhofer, and Vladislav Golyanik. Advances in neural rendering. *arXiv preprint arXiv:2111.05849*, 2021. 2
- [71] Romain Vergne, Romain Pacanowski, Pascal Barla, Xavier Granier, and Christopher M. Schlick. Radiance scaling for versatile surface enhancement. In *I3D '10*, 2010. 8
- [72] Sheng-Yu Wang, Oliver Wang, Richard Zhang, Andrew Owens, and Alexei A Efros. Cnn-generated images are surprisingly easy to spot... for now. In *CVPR*, pages 8695–8704, 2020. 12
- [73] Jiajun Wu, Chengkai Zhang, Tianfan Xue, William T Freeman, and Joshua B Tenenbaum. Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, pages 82–90, 2016. 2
- [74] Xudong Xu, Xingang Pan, Dahua Lin, and Bo Dai. Generative occupancy fields for 3d surface-aware image synthesis. *NeurIPS*, 34, 2021. 2
- [75] Guandao Yang, Xun Huang, Zekun Hao, Ming-Yu Liu, Serge Belongie, and Bharath Hariharan. Pointflow: 3d point cloud generation with continuous normalizing flows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4541–4550, 2019. 2
- [76] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces. *arXiv preprint arXiv:2106.12052*, 2021. 3
- [77] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Basri Ronen, and Yaron Lipman. Multiview neural surface reconstruction by disentangling geometry and appearance. *Advances in Neural Information Processing Systems*, 33, 2020. 2
- [78] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*, 2020. 8
- [79] Yuxuan Zhang, Wenzheng Chen, Huan Ling, Jun Gao, Yinan Zhang, Antonio Torralba, and Sanja Fidler. Image gans meet differentiable rendering for inverse graphics and interpretable 3d neural rendering. In *ICLR*, 2020. 2
- [80] Peng Zhou, Lingxi Xie, Bingbing Ni, and Qi Tian. Cips-3d: A 3d-aware generator of gans based on conditionally-independent pixel synthesis. *arXiv preprint arXiv:2110.09788*, 2021. 2, 6

## Appendix

In this appendix we provide additional qualitative and quantitative results on our approach, along with the technical details that supplement the main paper. In Sec. Appendix A, we discuss possible societal impacts of our technology. Then, we present additional experiments on the view-consistency of our RGB renderings via image reprojection (Sec. Appendix B). We further demonstrate the quality of our 3D shapes in Sec. Appendix C. In Sec. Appendix D, we describe the content of the supplementary videos and introduce a geometry-aware noise injection procedure to reduce flickering. Next, we discuss implementation details and the merits of our proposed sampling strategy in (Sec. Appendix E and Appendix F respectively). We conduct ablation studies on our approach in Appendix G). Finally, we continue our discussion on our method’s limitations in Sec. Appendix H.

### A. Societal Impacts

Image and 3D model-generating technologies (e.g., deepfakes) could be used for spreading misinformation about existing or non-existing people [29,40]. Our proposed technology allows generating multi-view renderings of a person, and might be used for creating more realistic fake videos. These problems could potentially be addressed by developing algorithms to detect neural network-generated content [72]. We refer readers to [69] for strategies of mitigating negative social impacts of neural rendering. Moreover, image generative models are optimized to follow the training distribution, and thus could inherit the ethnic, gender, or other biases present in the training dataset. A possible solution is creating a more balanced dataset, e.g., as in [35].

### B. View Consistency of RGB Renderings

#### B.1. Volume Rendering Consistency

The consistent volume rendering from our SDF-based technique naturally leads to high view consistency of our RGB renderings. To show the superior 3D-consistency of our SDF-based volume rendering, we measure the reprojection error when a side view pixels are warped to the frontal view. We randomly sample 1,000 identities and render the depth and RGB images at  $256 \times 256$  and set the side view to be  $1.5 \times$  the standard deviation of the azimuth distribution in training (which is 0.45 radians for FFHQ and 0.225 radians for AFHQ). We reproject the side-view RGB renderings to the frontal view using the side-view depth, and we do not ignore occluded pixels. We measure color inconsistency with the median of pixel-wise L1 error in RGB (0 - 255), averaged over the 1,000 samples. The use of median effectively removes the large errors coming from occlusions. Note that

Dataset:	FFHQ	AFHQ
PiGAN [9]	14.7	16.5
Ours (volume renderer)	<b>2.9</b>	<b>2.6</b>

Table 3. Quantitative view-consistency comparison of the RGB renderings. We evaluate the color error of the RGB volume renderings between the frontal view and the reprojection from a fixed side view. The error is measured as the median of the per-pixel mean absolute difference (0 - 255). We average the color inconsistency over 1,000 samples for each dataset. Our underlying SDF geometry representation promotes superior 3D consistency. (also see Fig. 9).

since PiGAN is trained with center-cropped FFHQ images (resized to  $320 \times 320$  and center-cropped to  $256 \times 256$ ), we apply the same transformation on our results before computing the median.

As shown in Tab. 3, StyleSDF presents significantly improved color consistency compared to the strongest current baseline, PiGAN [9]. Fig. 9 shows the sample depth and color rendering pairs used for the evaluation, along with the pixelwise error maps. The error maps demonstrate that our volume RGB renderings have high view consistency, as the large reprojection errors are mostly in the occluded regions. On the other hand, PiGAN’s reprojections do not align well with the frontal view, showing big errors also near the eyes, mouth, in presence of specular highlights, etc.

#### B.2. High-Resolution RGB Consistency

In Fig. 10, we present the reprojection experiment results using our high-resolution RGB outputs. As in the volume rendering consistency experiment, we reproject the RGB pixels from non-frontal views (with varying azimuth and elevation) to the frontal views. The results demonstrate the strong 3D-consistency of our high-resolution images, as the reprojected non-frontal images are similar to the frontal renderings. However, as mentioned in the limitation section of the main paper, the current implementation of StyleGAN2 comes with significant aliasing of the high-frequency components, resulting in noticeable pixel errors on regions with high-frequency details, e.g., hair, ears, eyes, etc. To identify the errors in the high-frequency details, we visualize the mean reprojection images. I.e. we project non-frontal views and average the pixel values across views. As can be seen in Fig. 11, the mean reprojection images closely replicate the identities and important structures of the frontal view, demonstrating strong view-consistencies. The error map confirms that most of the errors are concentrated on the high-frequency noise of the StyleGAN generator.

### C. Qualitative 3D results

We demonstrate the consistency of our 3D representation by overlaying the point clouds from the frontal and

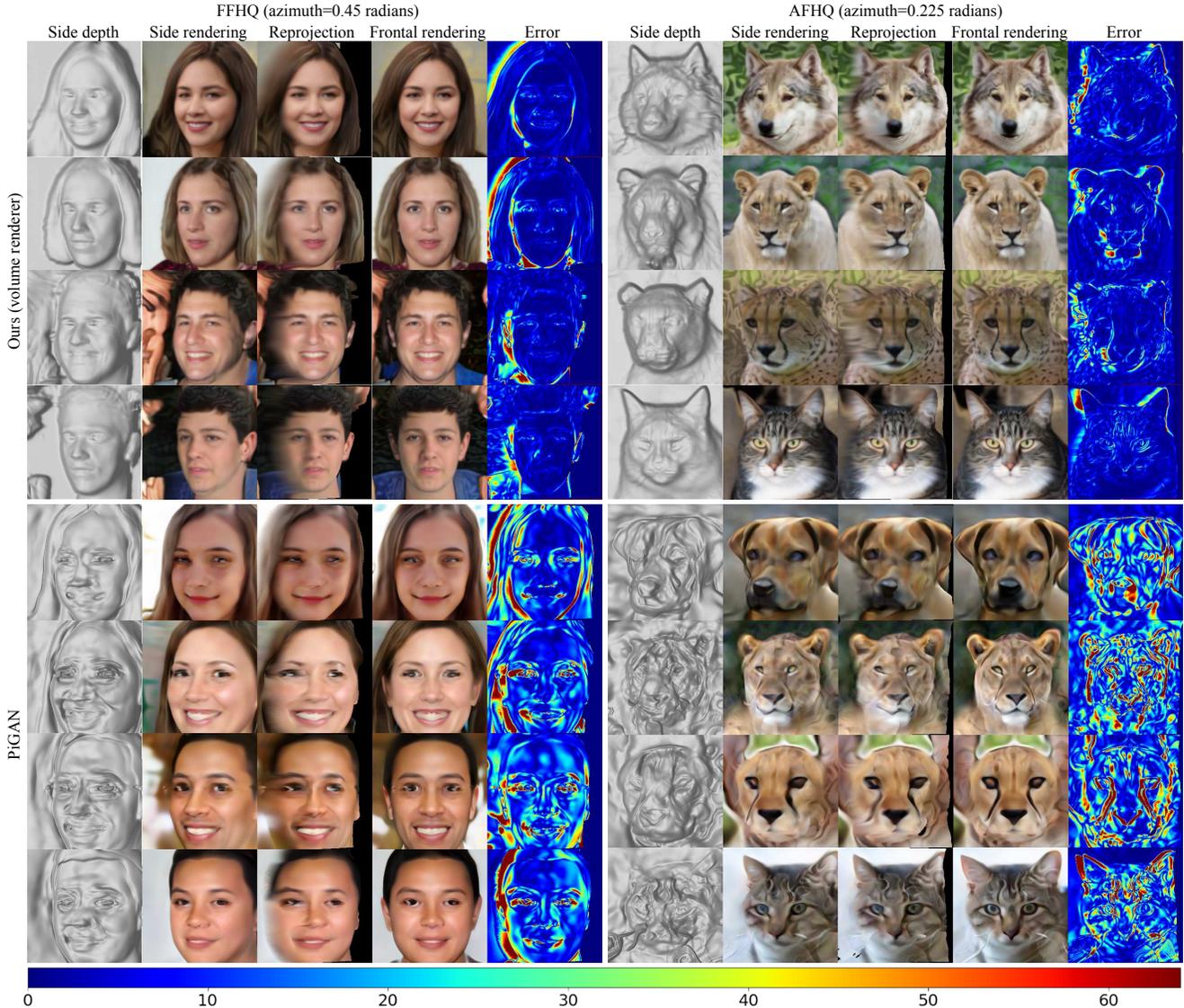


Figure 9. Qualitative view consistency comparison of RGB renderings. We project the rendering from a side view using its corresponding depth map to the frontal view. We compare the reprojection to the frontal-view rendering and compute the error map showing mean absolute difference in RGB channels (0 - 255). Our SDF-based technique generate superior depth quality and significantly improves the view-consistency of the RGB renderings. Most of our errors concentrate on the occlusion boundaries whereas PiGAN’s errors spread across the whole subject (e.g., eyes, mouth, specular highlights, fur patterns).

side view depth maps (Fig. 12b). The visualization, shown in two different colors, clearly shows high consistency between the depth maps. To show the quality and plausibility of our 3D models, we extract meshes on our SDFs via marching cubes and visualize them in extreme angles (Fig. 12c).

## D. Video Results

Since our 3D-consistent high-resolution image generation can be better appreciated with videos, we have attached 24 sequences in the supplementary material, featur-

ing view-generation results on the two datasets using two different camera trajectories. For each identity, we provide two videos, one for RGB and another for depth rendering. The videos are presented in the [project’s website](#).

### D.1. Geometry-Aware StyleGAN Noise

Even though the images shown in the main paper on multi-view RGB generation look highly realistic, we note that for generating a video sequence, the random noise of StyleGAN2 [38], when naïvely applied to 2D images, could result in severe flickering of high-frequency details between frames. The flickering artifacts are especially prominent for

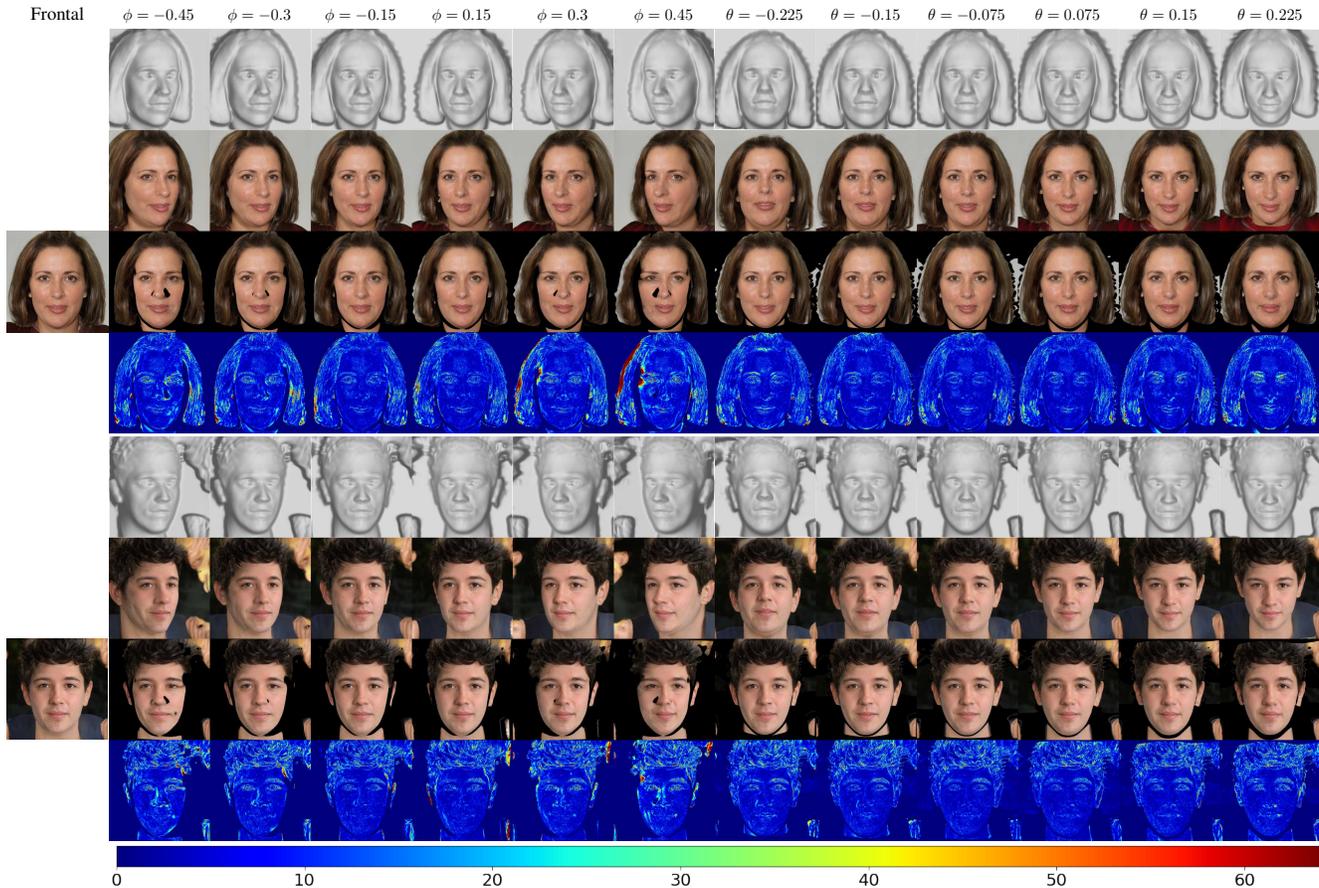


Figure 10. View-consistency visualization of high-resolution renderings. We use the side-view depth maps (first rows) to warp the side-view RGB renderings (second rows) to the frontal view (first column). The reprojected pixels that pass the occlusion testing are shown in the third row. We compare the rejections with the frontal-view renderings and show the per-pixel error maps (fourth rows). Our rejections well align with the frontal view with errors mostly in the occlusion boundaries and high-frequency details.

the AFHQ dataset due to high-frequency textures from the fur patterns.

Therefore, we aim at reducing this flickering by adding the Gaussian noise in a 3D-consistent manner, i.e., we want to attach the noise on the 3D surface. We achieve this by extracting a mesh (at 128 resolution grid) for each sequence from our SDF representation and attach a unit Gaussian noise to each vertex, and render the mesh using vertex coloring. Since higher resolution intermediate features require up to  $1024 \times 1024$  noise map, we subdivide triangle faces of the extracted mesh once every layer, starting from  $128 \times 128$  feature layers. The video results show that the geometry-aware noise injection reduce the flickering problem on the AFHQ dataset, but noticeable flickering still exist. Furthermore, we observe that the geometry-aware noise slightly sacrifices individual frame’s image quality, presenting less pronounced high-frequency details, likely due to the change of the Gaussian noise distribution during the rendering process. The videos rendered with geometry-aware noise can be viewed at the [project’s website](#).

## E. Implementation Details

### E.1. Dataset Details

**FFHQ:** We trained FFHQ with R1 regularization loss of 10. The camera field of view was fixed to  $12^\circ$  and its azimuth and elevation angles are sampled from Normal distributions with zero mean and standard deviations of 0.3 and 0.15 respectively. We set the near and far fields to  $[0.88, 1.12]$  and sample 24 points per ray during training. We trained our volume renderer for  $200k$  iterations and the 2D-Styled generator for  $300k$  iterations.

**AFHQ:** The AFHQ dataset contains training and validation sets for 3 classes, cats, dogs and wild animals. We merged all the training data into a single training set. We apply R1 regularization loss of 50. Both azimuth and elevation angles are sampled from a Gaussian distribution with zero mean and standard deviation of 0.15 and a camera field of view of  $12^\circ$ . The near and far fields as well as the number of samples per ray are identical to the FFHQ setup. Our

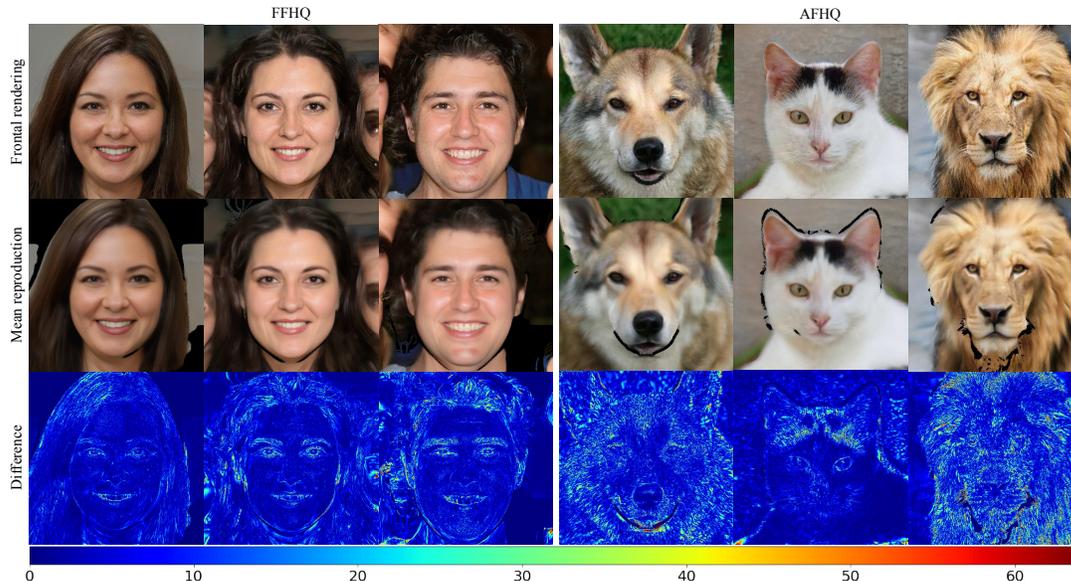


Figure 11. Color consistency visualization with mean faces. We reproject the RGB renderings from the side views to the frontal view (as in Fig. 10). We show the mean reprojections that pass the occlusion testing and their differences to the frontal-view renderings. The mean reprojections are well aligned with the frontal rendering. The majority of the errors are in the high-frequency details, generated from the random noise maps in the StyleGAN component. This demonstrates the strong view consistency of our high-resolution renderings.

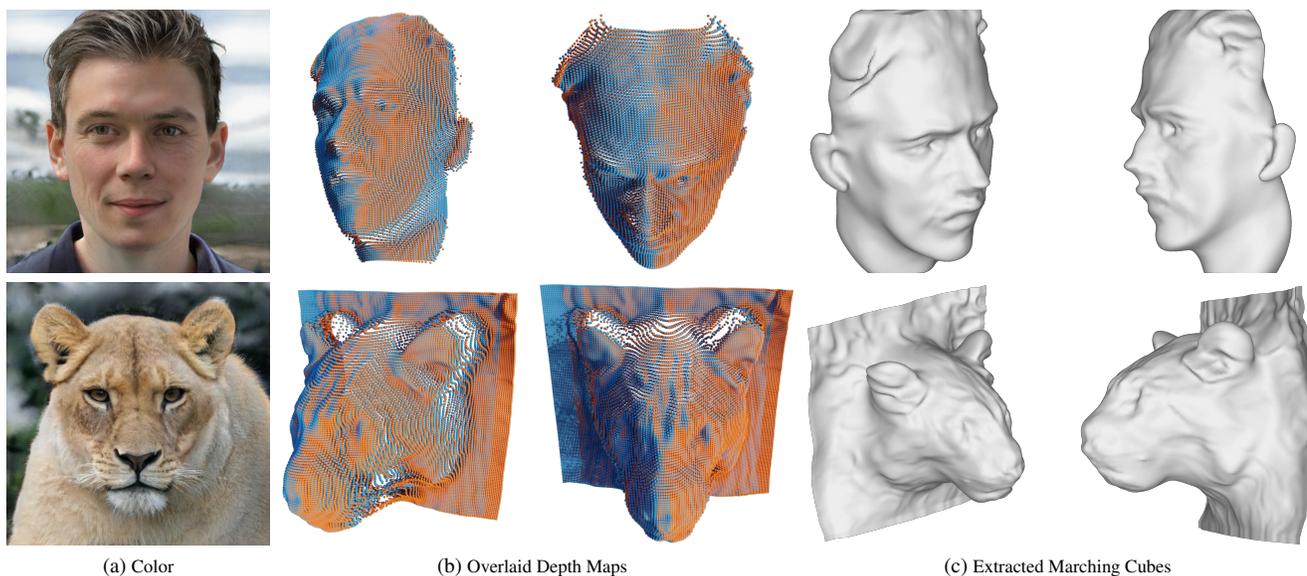


Figure 12. Consistent and plausible 3D shapes. (a) Color images. (b) Overlaid point clouds extracted from frontal and side view depth maps. (c) Marching cubes meshes, rendered from extreme angles.

volume renderer as well as the 2D-Styled generator were trained for  $200k$  iterations.

## E.2. Training Details

**Sphere Initialization:** During our experiments we have noticed that our SDF volume renderer can get stuck at a local minimum, which generates concave surfaces. To avoid this optimization failure, we first initialize the MLP to generate an SDF of a sphere centered at the origin with a fixed radius.

We analytically compute the signed distance of the sampled points from the sphere and fit the MLP to match these distances. We run this procedure for  $10k$  iterations before the main training. The importance of sphere initialization is discussed in Appendix G.

**Training setup:** Our system is trained in a two-stage strategy. First, we train the backbone SDF volume renderer on  $64 \times 64$  images with a batch size of 24 using the ADAM [41] optimizer with learning rates of  $2 \cdot 10^{-5}$  and  $2 \cdot 10^{-4}$  for the

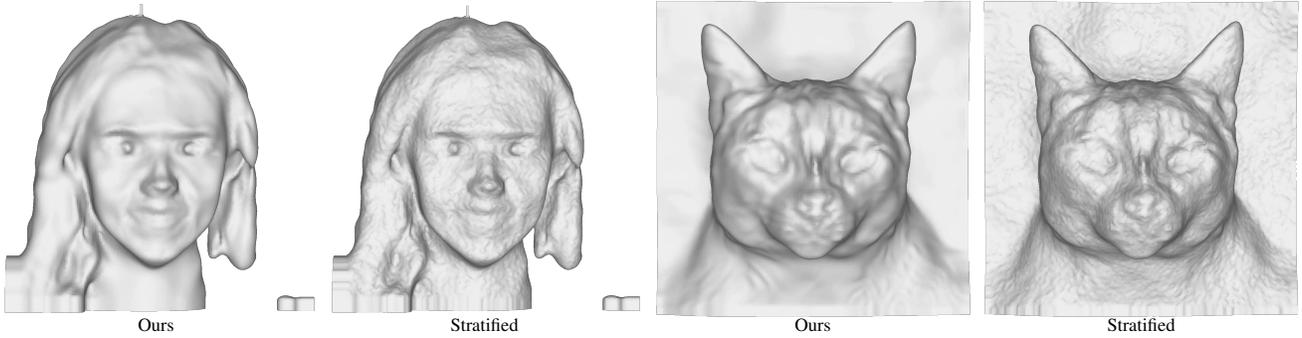


Figure 13. We compare extracted meshes using our sampling strategy vs. stratified sampling. Note the noise induced by stratified sampling. (zoom in for details)

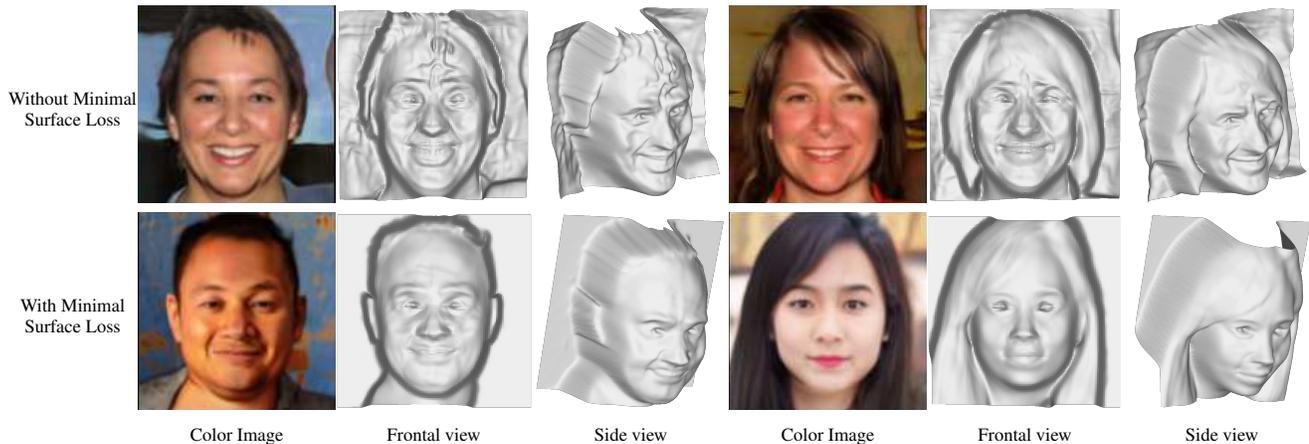


Figure 14. Minimal surface loss ablation study. We visualize the volume rendered RGB and depth images from volume renderers trained with and without the minimal surface loss. The Depth map meshes are visualized from the front and side views. Note how a model trained with the minimal surface loss generates smoother surfaces and is less prone to shape-radiance ambiguities, e.g., specular highlights are baked into the geometry.

generator and discriminator respectively and  $\beta_1 = 0, \beta_2 = 0.9$ . We accumulate gradients in order to fit to the GPU memory constraints. For instance, a setup of 2 NVIDIA A6000 GPUs (a batch of 12 images per GPU) requires the accumulation of two forward passes (6 images per forward pass) and takes roughly 3.5 days to train. We use an exponential moving average model during inference.

In the second phase, we freeze the volume renderer weights and train the 2D styled generator with identical setup to StyleGAN2 [39]. This includes ADAM optimizer with 0.002 learning rate and  $\beta_1 = 0, \beta_2 = 0.99$ , equalized learning rate, lazy R1 and path regularization, batch size of 32, and exponential moving average. We trained the styled generator on 8 NVIDIA TeslaV100 GPUs for 7 days.

## F. Sampling Strategy

NeRF [48], along with existing 3D-aware GANs like PiGAN [9], rely on hierarchical sampling strategy for obtaining more samples near the surface. Our use of SDFs al-

lows sampling the volume with smaller number of samples without sacrificing the surface quality, thereby reducing the memory footprints and simplifying the implementation.

Stratified sampling randomizes the distance between adjacent samples along each ray, adding undesired noise to the volume rendering (Fig. 13). The randomness also amplifies flickering in RGB videos. Our sampling strategy ensures that the integration intervals are of the same length, which eliminates the noise and results in smoother volume rendering outputs.

## G. Ablation studies

We perform two ablation studies to show the necessity of the minimal surface loss (see main paper) and the sphere initialization. As can be seen in Figure 14, on top of preventing spurious and non-visible surfaces from being formed, the minimal surface loss also helps to disambiguate between shape and radiance. Penalizing values that are close to zero essentially minimizes the surface area and

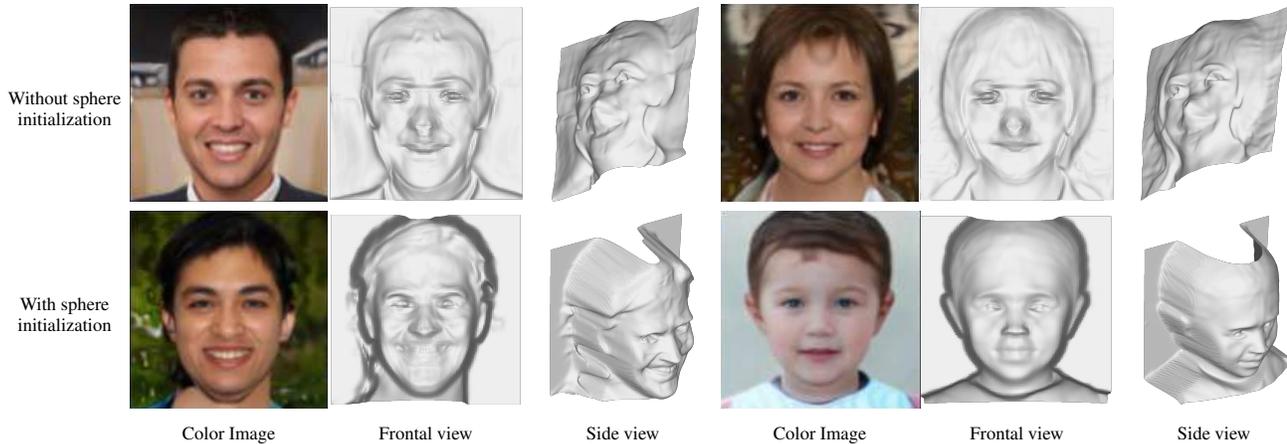


Figure 15. Sphere initialization ablation study. We visualize volume-rendered RGB and depth images from volume renderers trained with and without sphere initialization. The Depth map meshes are visualized from the front and side views. Note how a model trained without model sphere initialization generates concave surfaces.

makes the network prefer smooth SDFs.

In Figure 15, we show the importance of the sphere initialization in breaking the concave/convex ambiguity. Without properly initializing the weights, the network gets stuck at a local minimum that generates concave surfaces. Although concave surfaces are physically incorrect, they can perfectly explain multi-view images, as they are essentially the "mirror" surface. Concave surfaces cause the images to be rendered in the opposite azimuth angle, an augmentation that the discriminator cannot detect as fake. Therefore, the generator cannot recover from this local minima.

## H. Limitations (continued)

As mentioned in the main paper, our high-resolution generation network is based on the implementation of StyleGAN2 [38], and thus might experience the same aliasing and flickering at regions with high-frequency details (e.g., hair), which are recently addressed in Alias-free GAN [37] or Mip-NeRF [5]. Moreover, we observe that the reconstructed geometry for human eyes contain artifacts, characterized by concave, instead of convex, eye balls. We believe that these artifacts often lead to slight gaze changes along with the camera views. As stated in the main paper, our current implementation of volume rendering during inference uses fixed frontal view directions for RGB queries  $c(x, v)$ , and thus cannot express moving specular highlights along with the camera.

## I. Additional Results

We show uncurated set of images generated by our networks (Fig. 16).

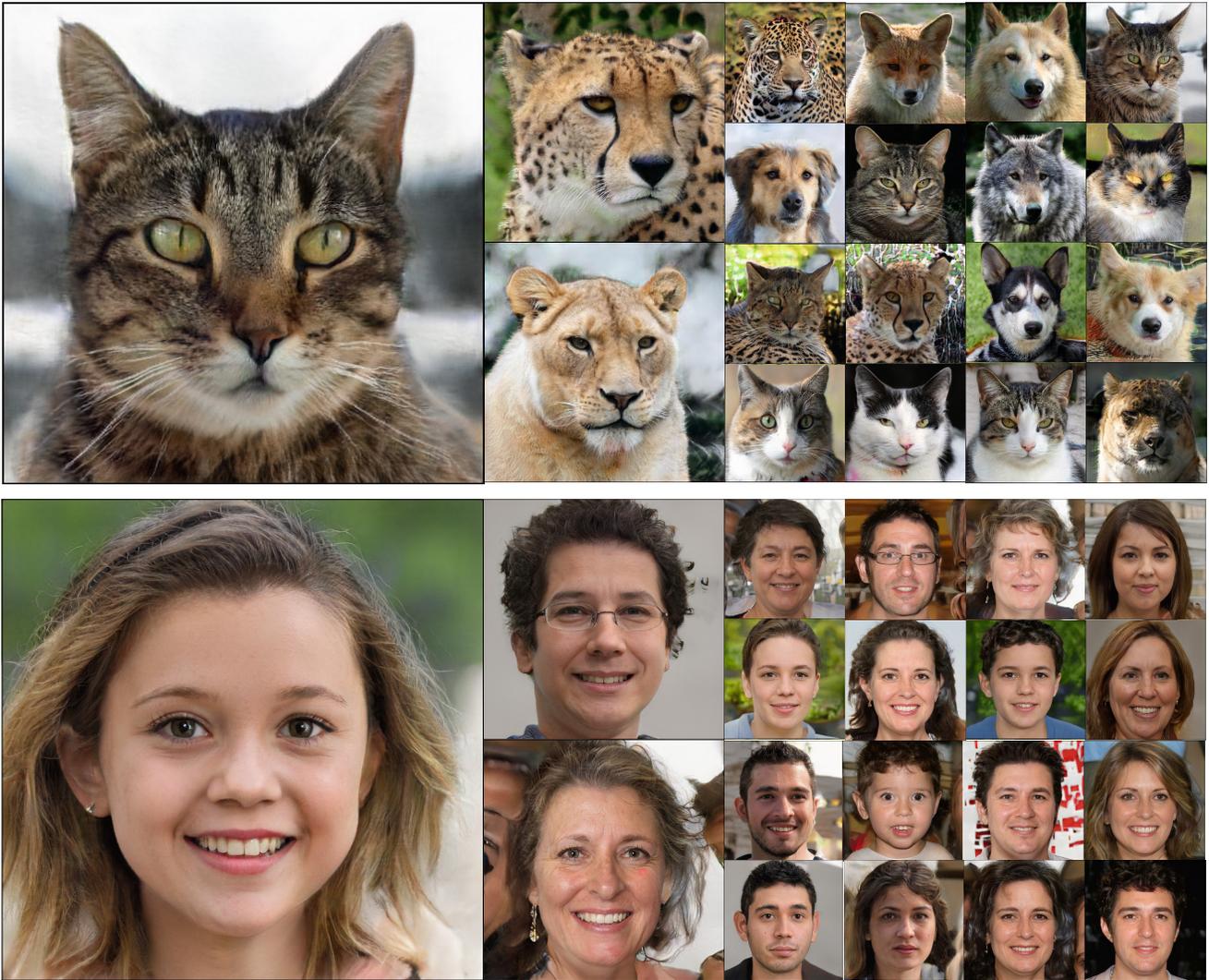


Figure 16. Uncurated high-resolution RGB images that are randomly generated by StyleSDF.