

# 3D-GIF: 3D-Controllable Object Generation via Implicit Factorized Representations

Minsoo Lee      Chaeyeon Chung      Hojun Cho      Minjung Kim  
 Sanghun Jung      Jaegul Choo      Minhyuk Sung  
 KAIST, Daejeon, South Korea

{alstn2022, cy\_chung, hojun.cho, emjay73, shjung13, jchoo, mhsung}@kaist.ac.kr

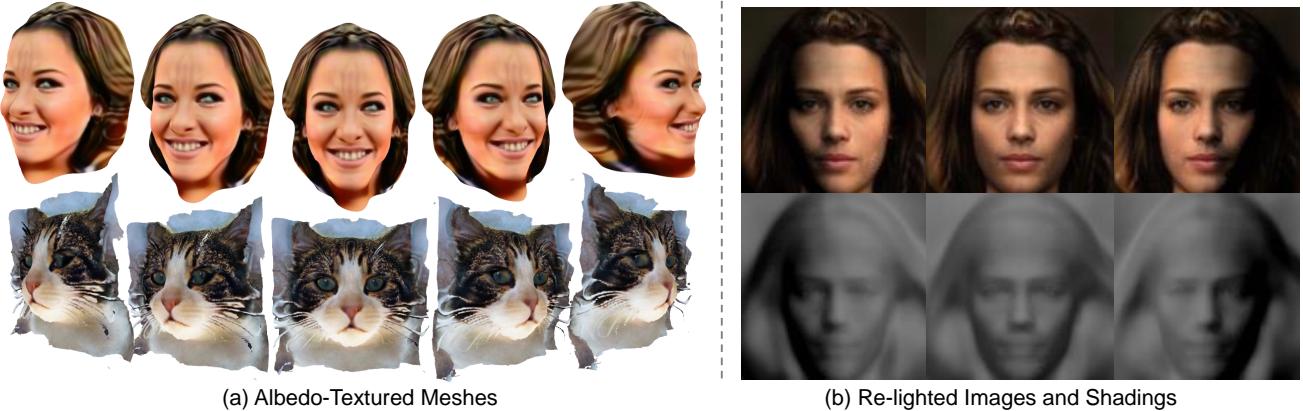


Figure 1. Generation results of 3D-GIF. (a) indicates generated albedo-textured meshes and (b) presents re-lighted images and shadings.

## Abstract

While NeRF-based 3D-aware image generation methods enable viewpoint control, limitations still remain to be adopted to various 3D applications. Due to their view-dependent and light-entangled volume representation, the 3D geometry presents unrealistic quality, and the color should be re-rendered for every desired viewpoint. To broaden the 3D applicability from 3D-aware image generation to 3D-controllable object generation, we propose the factorized representations which are view-independent and light-disentangled, and training schemes with randomly sampled light conditions. We demonstrate the superiority of our method by visualizing factorized representations, re-lighted images, and albedo-textured meshes. Furthermore, we show that our approach improves the quality of the generated geometry via visualization and quantitative comparison. To the best of our knowledge, this is the first work that extracts albedo-textured meshes utilizing unposed 2D images without any additional labels or assumptions.

## 1. Introduction

Generating photo-realistic images has been a long-standing problem in computer vision. Addressing this concern, generative adversarial networks (GANs) [13] have shown their remarkable capability of generating high-quality images [3, 19–21, 41]. However, due to the lack of

learning 3D features of an object, the previous GAN-based approaches have limitations on editing the content of a generated image in a 3D space (*e.g.* viewpoint changing, re-lighting, etc.). To enable controlling contents of the images in a 3D space, 3D-aware image synthesis [17, 30, 31] has been proposed, generating a 3D representation of an object and then producing an image with a neural projection. Among them, recent approaches [4, 32, 43] using neural radiance fields (NeRF) [29] presented a successful viewpoint control while preserving the multi-view consistency. Particularly, Chan *et al.* [4] showed a potential that the idea can be extended to 3D applications (*e.g.* games and movies) by extracting the 3D mesh from the generated representation.

However, as also pointed out in previous work [4], such approaches bear limitations as follows: 1) their texture representation is difficult to be rendered under various external environments (*i.e.* view directions and light conditions) via a conventional computer graphics (CG) renderer, and 2) their 3D representations contain unrealistic geometry.

First, since they generate *view-dependent* and *light-entangled* 3D representations, the outputs are hardly applicable to the CG rendering pipeline that is necessary for various 3D applications. The CG rendering process employs physical reflectance properties (*i.e.* ambient, diffuse, and specular reflection) to render an object with diverse external environments. Hence, we need an object with inherent

feature representations (*e.g.* albedo, normal, and specular coefficients), which are independent of view directions and light conditions; we call such an object a *3D-controllable* object. However, the previous approaches predict the color conditioned on the viewing direction to represent the plausible light reflection, which causes redundant rendering problem when utilized in CG rendering pipelines. Fig. 2 illustrates textured meshes obtained from pi-GAN [4], where the predicted colors are conditioned on two different view directions (V1 and V2). As shown in Fig. 2(a), when we observe a mesh with V1-conditioned texture from the same direction (V1), it looks realistic. However, (c) shows that the same mesh observed from a different view direction (V2) seems unrealistic. To alleviate this issue, the colors have to be regenerated from the given view direction (V2), as presented in Fig. 2(d). That is, we need to re-render an object for every desired view direction, which limits the applicability for 3D applications. Also, since the light effects are already included in the generated colors, which is not albedo, adding different light conditions to an object produces unnatural visual quality. To resolve these issues, we generate inherent 3D feature representations (*i.e.* albedo, normal, and specular coefficients) without view dependency and explicitly compute the light reflection. This allows our model to generate 3D-controllable objects beyond 3D-aware image generation.

Moreover, the existing methods have unrealistic surface geometry of their generated 3D representation due to the ambiguity in 3D-to-2D projection. Light reflection (*e.g.* diffuse and specular reflection) on the object depends on its surface geometry. However, the previous methods predict colors only to compose a realistic 2D image at the given view direction without considering its geometry and light reflection. This allows the models to generate an object with unrealistic geometry as presented in Fig. 2(e), while the projected 2D images look realistic (*i.e.* ambiguity in 3D-to-2D projection). Such defects can be manifested by simple light controls as described in Fig. 2(a) and (b), where (b) is the result of adding light to (a). In other words, we can resolve the ambiguity by adding light reflection to the generated image. Based on these insights, we force the geometry of generated 3D representation to be realistic by modeling the color from physical reflectance properties and explicitly controlling the light conditions while training.

In this paper, we propose 3D-GIF that removes view and light dependency of the color representation and addresses 3D-to-2D projection ambiguity in a 3D-aware image generation task. Addressing the limitations of the color representation of the previous work, our model generates *view-independent* and *light-disentangled* factorized representations (*i.e.* albedo, normal, and specular coefficients) and renders the final image via photometric rendering (*i.e.* shading model with light conditions and the factorized repre-

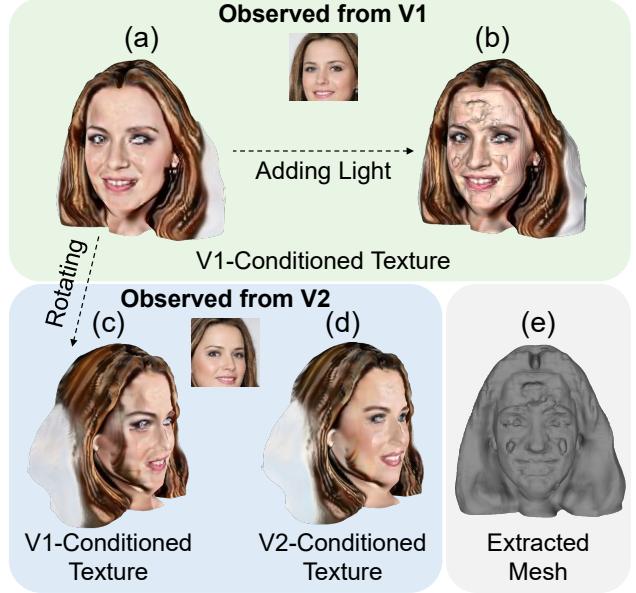


Figure 2. Generated textured mesh from pi-GAN [4]. (a) and (d) are the meshes observed from different view directions and textured with colors conditioned on their view directions (V1, V2). (c) shows the rotated mesh of (a), and (b) is a result of adding light on (a) via the conventional rendering pipeline. (e) describes the extracted geometry.

sentsations using physical reflectance properties). Randomly sampling light conditions and deceiving a discriminator with rendered images, our model can generate multi-view realistic images with accurate geometry of their 3D representations. Also, we introduce the staged training and the normal prediction layer to improve the training stability. We demonstrate that 3D-GIF successfully generates not only realistic multi-view 2D images but also factorized 3D representations with accurate geometry on various datasets. Besides, since our model explicitly simulates the light reflection, it improves the 3D applicability of the generative model by producing re-lighted images and high-quality albedo-textured meshes. To the best of our knowledge, this is the first work that generates 3D-controllable objects using unposed 2D images without any additional supervision. We summarize our contributions as follows:

- We present 3D-GIF that generates view-independent and light-disentangled volume representations (*i.e.* albedo, normal, and specular coefficients), broadening its 3D applicability.
- We improve the geometry quality of generated 3D representations, resolving the 3D-to-2D projection ambiguity by modeling the color from physical reflectance properties.
- We validate the effectiveness of our model by presenting re-lighted images and albedo-textured meshes and providing qualitative and quantitative comparisons on the geometry quality of generated representations.

## 2. Related Work

**Neural 3D representations.** Implicit representation has been employed by recent work [5, 11, 12, 27, 28, 33, 34, 36, 42] for their memory efficiency and continuity, training the network to map 3D coordinates to occupancy field or signed distance functions. Several methods [15, 39, 42, 55] extended this to texture representation and successfully extracted textured mesh from 2D images. However, they still require 3D ground truth for training. To mitigate this undesirable necessity, differentiable rendering techniques are adopted to enable 2D image-based training [23, 46, 51]. Mildenhall *et al.* [29] proposed neural radiance fields (NeRF) with volume rendering, showing impressive performance for novel-view synthesis. To further improve this seminal work, some approaches [2, 47, 54] decomposed the view-dependent color of NeRF into reflectance properties and light conditions. NeRD [2] extracts albedo-textured meshes from the implicit 3D representations for real-time rendering and re-lighting. However, the aforementioned approaches require multi-view and posed images for their training. Therefore, we propose a novel generative architecture that only requires unposed single-view images for training.

**3D-aware image generation.** Recent generative models utilized 3D representation to generate view-controllable images. Several approaches employed voxel [17, 30, 31] as their 3D representation. Meanwhile, they generated coarse results without fine details due to the memory inefficiency of voxel representation. Mitigating this issue, NeRF-based methods [4, 43] emerged, generating view-consistent 2D images in high quality. Niemeyer *et al.* [32] extended this radiance field to feature field and adopted a 2D convolution neural renderer to further enhance the quality of generated images and rendering efficiency. However, convolution operation on adjacent pixels hurts the internal 3D information, which leads to multi-view inconsistency and noisy geometry. In addition, these NeRF-based methods output view-dependent and light-entangled colors, impeding them from being further extended to 3D object generation.

**3D object generation.** Early 3D object generation methods required 3D supervision from voxel [50] and point cloud [5, 18, 45]. However, constructing 3D ground truth datasets is time-consuming and labor-intensive. To avoid 3D supervision, several methods utilized 3D-morphable models [8, 48] or 2D supervision [10, 49]. Some of them [16, 37] successfully generated textured mesh from 2D images with additional supervisions such as silhouette images and camera poses via differentiable rendering process. Despite their success, they utilized ellipsoid template mesh as a starting point and still required additional supervisions. These 3D objects can also be generated from the aforementioned NeRF-based models [4, 43] without requiring additional supervision other than 2D images. Chan *et al.* [4] presented generated meshed in its paper, showing a poten-

Paper	Supervision	Goals	Assumption
[50]	3DV	3DV (w/o T)	
[18, 45]	3DP	3DP (w/o T)	
[5]	3DP	3DM* (w/o T)	
[8, 48]	I, 3DMM	MI, LI, 3DMM	
[10]	SI	3DV (w/o T)	
[49]	I, KP	3DM (w/o T)	Ellipsoid
[16]	I, KP, BG, SI	3DM	Ellipsoid
[37]	I, SI, KP	3DM	Ellipsoid
[44]	I	MI, LI, A, D, N	Symmetry
[4, 43]	I	MI, 3DM* (w/o T)	
Ours	I	MI, LI, 3DM*, A, N, S	

Table 1. Comparison with the previous 3D object generation work. I: image, 3DV: 3D voxel, 3DP: 3D point cloud, 3DM: 3D mesh, 3DMM: 3D morphable model, KP: 2D keypoints or camera matrix, SI: silhouette image, BG: background image, MI: multi-view images, LI: re-lighting images, T: texture, A: albedo, N: normal, D: depth, S: specular coefficients. 3DM\* denotes extracted 3D mesh from implicit representation via marching cube algorithm. tial of applicability to 3D applications. However, the generated meshes have a low quality of surface geometry, and even worse, the texture of mesh was difficult to be rendered with various external environments since it predicts view-dependent colors. In this paper, we successfully address these concerns by factorizing view-dependent colors into inherent representations and explicitly controlling light conditions. Table 1 shows our contributions by comparing the types of supervisions, assumptions, and goals of ours and previous 3D object generation methods.

## 3. Proposed Method

This section presents how we generate view-independent and light-disentangled volume representations for enhancing the 3D controllability of generated objects.

### 3.1. Preliminaries

**Neural radiance field.** Radiance field is a continuous function that maps a 3D location  $\mathbf{x} \in \mathbb{R}^3$  and a viewing direction  $\mathbf{d} \in \mathbb{S}^2$  to volume density  $\sigma \in \mathbb{R}^+$  and view-dependent color  $\mathbf{c} \in \mathbb{R}^3$ . Inspired by NeRF [29], recent approaches [32, 43] use multi-layer perceptrons to represent this mapping function as

$$f_\theta : \mathbb{R}^3 \times \mathbb{S}^2 \rightarrow \mathbb{R}^3 \times \mathbb{R}^+, \quad f_\theta(\mathbf{x}, \mathbf{d}) = (\mathbf{c}, \sigma). \quad (1)$$

To represent non-Lambertian effects of color,  $\mathbf{d}$  is concatenated with intermediate feature vector of  $f_\theta$  and passed to additional layers that output view-dependent RGB color  $\mathbf{c}$ .

**Volume rendering.** With sampled camera pose  $\xi$ ,  $\{\mathbf{x}_j\}_{j=1}^{N_x}$  denotes sampled points along the camera ray for a pixel location  $(u, v) \in \Omega$ , where  $\Omega$  indicates the set of pixel locations in the image plane, and  $N_x$  denotes the number of sampled points in a camera ray. To calculate the expected color  $\mathbf{c}_{uv}$  for each pixel location  $(u, v)$ , the non-parametric

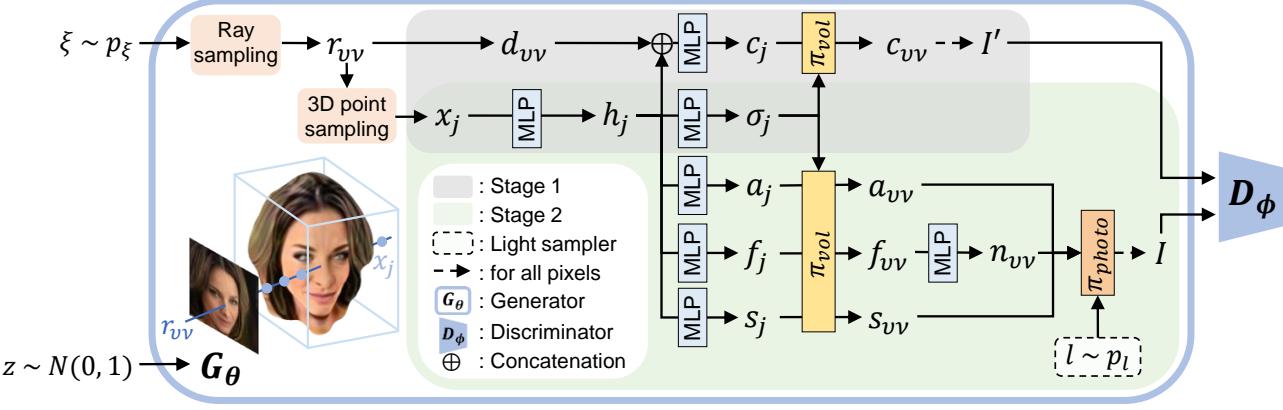


Figure 3. Overview of 3D-GIF. To generate view-independent and light-disentangled colors, we disentangle the light and factorize the volume representations into albedo, normal, and specular coefficients.  $\pi_{vol}$  and  $\pi_{photo}$  denote volume rendering in Sec. 3.1 and photometric image rendering in Sec. 3.2, respectively.

volume rendering function  $\pi_{vol}(\cdot)$  that maps the color  $\mathbf{c}_j$  and the volume density  $\sigma_j$  of the sampled point  $\mathbf{x}_j$  to  $\mathbf{c}_{uv}$  can be defined as

$$\pi_{vol} : (\mathbb{R}^+, \mathbb{R}^3)^{N_x} \rightarrow \mathbb{R}^3, \quad \pi_{vol}(\{\sigma_j, \mathbf{c}_j\}_{j=1}^{N_x}) = \mathbf{c}_{uv}. \quad (2)$$

To be specific, we obtain the estimated color  $\mathbf{c}_{uv}$  via numerical integration [29] as

$$\begin{aligned} \mathbf{c}_{uv} &= \sum_{j=1}^{N_x} w_j \mathbf{c}_j, \quad w_j = \tau_j \alpha_j, \\ \tau_j &= \prod_{k=1}^{j-1} (1 - \alpha_k), \quad \alpha_j = 1 - e^{-\sigma_j \delta_j}, \end{aligned} \quad (3)$$

where  $\tau_j, \alpha_j$  denote the transmittance and the alpha value for  $\mathbf{x}_j$ , and  $\delta_j = \|\mathbf{x}_{j+1} - \mathbf{x}_j\|_2$  is the L2-distance between the adjacent sampled points.

To this end, we extend this definition to the volume parameter  $\mathbf{p}$ , similar to Boss *et al.* [2], which can be one of our factorized volume representations such as albedo  $\mathbf{a}$ , intermediate feature vector  $\mathbf{f}$ , and specular coefficients  $\mathbf{s}$ . These factorized representations obtained from our network are volume rendered as

$$\mathbf{p}_{uv} = \sum_{j=1}^{N_x} w_j \mathbf{p}_j, \quad w_j = \tau_j \alpha_j. \quad (4)$$

**Generative NeRF.** Schwarz *et al.* [43] proposed a NeRF-based generative method that conditions the model on the latent code  $\mathbf{z} \in \mathbb{R}^M$ , formulated as

$$g_\theta : \mathbb{R}^3 \times \mathbb{S}^2 \times \mathbb{R}^M \rightarrow \mathbb{R}^3 \times \mathbb{R}^+, \quad g_\theta(\mathbf{x}, \mathbf{d}, \mathbf{z}) = (\mathbf{c}, \sigma). \quad (5)$$

The model generates implicit radiance field that maps a 3D location and a view direction to density and color conditioned on randomly-sampled latent vector  $\mathbf{z}$ . The generator is trained to deceive the discriminator with volume-rendered 2D images under a sampled viewpoint.

### 3.2. Generative factorized volume representations

To overcome the aforementioned limitations of the previous methods in Sec. 1, we propose 3D-GIF that aims to generate view-independent and light-disentangled volume representations. We factorize the representations into albedo, normal, and specular coefficients to obtain the final texture by leveraging physical reflectance properties.

**Factorized volume representations.** As shown in Fig. 3, we re-parameterize the view-dependent color of the previous methods [4, 29, 43] with normal  $\mathbf{n} \in \mathbb{R}^3$ , specular coefficients  $\mathbf{s} \in \mathbb{R}^2$ , and view-independent color  $\mathbf{a} \in \mathbb{R}^3$  that is interpreted as albedo. The factorized volume representations of our model are predicted as

$$\begin{aligned} h_\theta : \mathbb{R}^3 \times \mathbb{R}^M &\rightarrow \mathbb{R}^H, & h_\theta(\mathbf{x}_j, \mathbf{z}) &= \mathbf{h}_j \\ \sigma_\theta : \mathbb{R}^H &\rightarrow \mathbb{R}^+, & \sigma_\theta(\mathbf{h}_j) &= \sigma_j \\ a_\theta : \mathbb{R}^H &\rightarrow \mathbb{R}^3, & a_\theta(\mathbf{h}_j) &= \mathbf{a}_j \\ f_\theta : \mathbb{R}^H &\rightarrow \mathbb{R}^4, & f_\theta(\mathbf{h}_j) &= \mathbf{f}_j \\ s_\theta : \mathbb{R}^H &\rightarrow \mathbb{R}^2, & s_\theta(\mathbf{h}_j) &= \mathbf{s}_j, \end{aligned} \quad (6)$$

where  $\mathbf{h}_j$  is an intermediate feature vector conditioned on  $\mathbf{z}$ . The prediction layers  $\sigma_\theta, a_\theta, f_\theta$ , and  $s_\theta$  are implemented with a linear layer. We obtain the normal after volume rendering the predicted normal features  $\mathbf{f}$  and passing it through additional linear layers as

$$n_\theta : \mathbb{R}^4 \rightarrow \mathbb{R}^3, \quad n_\theta(\mathbf{f}_{uv}) = \mathbf{n}_{uv}, \quad (7)$$

where  $n_\theta$  consists of two linear layers and normalization to a unit vector.

Although the normal of the implicit function at a point  $j$  can be obtained by  $\frac{-\nabla_{\mathbf{x}} \sigma_j}{\|\nabla_{\mathbf{x}} \sigma_j\|}$  [2, 52, 54], utilizing this directly for the photometric rendering can lead to sub-optimal results for the following reasons. First, as also pointed out in Zhang *et al.* [54], the gradient of volume density tends to be noisy. This makes it difficult to render clear images when

exploiting the gradient of  $\sigma$  for rendering. Moreover, the second-order gradient has to be computed for optimization, which is computationally demanding. To address such concerns for stabilizing the training procedure, we introduce an additional normal prediction layer and use the output (*i.e.* predicted normal) as a proxy of the gradient of  $\sigma$ . Afterward, we devise a regularization term that utilizes the predicted normal as weak supervision for the gradient of  $\sigma$ , stabilizing the optimization process. Our final normal loss is defined as

$$\mathcal{L}_{\text{normal}} = \frac{1}{|\Omega|} \sum_{(u,v) \in \Omega} \left\{ \left\| \mathbf{n}_{uv} - \sum_{j=1}^{N_x} w_j \frac{-\nabla_{\mathbf{x}} \sigma_j}{\|\nabla_{\mathbf{x}} \sigma_j\|_2} \right\|_2 + \lambda \mathcal{L}_{\text{TV}}(\mathbf{n}_{uv}) \right\}, \quad (8)$$

where  $\lambda$  denotes the hyper-parameter indicating the relative importance between two losses, and  $\mathcal{L}_{\text{TV}}$  represents the difference between the center pixel value and its surrounding ones for smoothing the normal in the local region.

**Photometric image rendering.** The final textured image is rendered from a combination of albedo, normal, specular coefficients, and external environment via photometric rendering function  $\pi_{photo}(\cdot)$ . We sample light condition  $l = (k_a, k_d, l_d)$  from the prior distribution, where  $l$  consists of the intensity of ambient light  $k_a \in \mathbb{R}^+$ , the intensity of directional light  $k_d \in \mathbb{R}^+$ , and the direction of light  $l_d \in \mathbb{R}^3$ . We calculate the dot product of the normal  $\mathbf{n}_{uv}$  and light direction  $l_d$  to obtain a value for the directional illumination. Afterward, the shading map  $\mathbf{H}$  is obtained by multiplying the intensity of directional light and adding the intensity of ambient light to the obtained value:  $\mathbf{H}_{uv} = (k_a + k_d \cdot \max\{0, \langle l_d, \mathbf{n}_{uv} \rangle\})$ . The specular map  $\mathbf{S}$  is obtained from the predicted specular coefficients, normal, directional light and view direction:  $\mathbf{S}_{uv} = k_d \cdot \mathbf{s}_{uv}^{(0)} \cdot (\max\{0, \cos \theta\}) \mathbf{s}_{uv}^{(1)}$ , where  $\theta$  denotes the angle between the view direction and the light direction after reflected by normal, and  $\mathbf{s}^{(0)}, \mathbf{s}^{(1)}$  denote the coefficient for the specular reflection and shininess, respectively. Note that this process is calculated for each pixel  $(u, v)$  in the image plane. Finally, we render the final textured image from a combination of albedo, shading map, and specular map, formulated as

$$\mathbf{I} = \mathbf{a} \cdot \mathbf{H} + \mathbf{S}. \quad (9)$$

### 3.3. Training

**Training objective function.** Fig. 3 depicts the overall training procedure. Our model first samples a latent variable  $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I})$  to condition the model. We sample points  $\{\mathbf{x}_{j=1}\}_j^{N_x}$  on the ray direction from a random camera pose  $\xi \sim p_\xi$  and sample the random light condition  $l \sim p_l$ . Then, we train our model  $G_\theta$  using the non-saturating GAN [13]

loss with R1 regularization [26]:

$$\begin{aligned} \mathcal{L}_{\text{GAN}}(\theta, \phi) = & \mathbf{E}_{\mathbf{z} \sim p_z, \xi \sim p_\xi, l \sim p_l} [f(D_\phi(G_\theta(\mathbf{z}, \xi, l)))] \\ & + \mathbf{E}_{\mathbf{I} \sim p_D} [f(-D_\phi(\mathbf{I})) + \lambda \|\nabla D_\phi(\mathbf{I})\|^2], \end{aligned} \quad (10)$$

where  $f(x) = -\log(1 + \exp(-x))$ ,  $p_D$  denotes the real data distribution, and  $D_\phi$  indicates the discriminator.  $D_\phi$  is trained by maximizing Eq. (10), and the generator is trained by minimizing both Eq. (8) and Eq. (10):

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{GAN}} + \lambda_n \mathcal{L}_{\text{normal}}, \quad (11)$$

where  $\lambda_n$  denotes the hyper-parameter indicating the relative importance between two losses.

**Staged training.** To initialize the geometry of volume representation, we first train our network to generate the view-dependent color image  $\mathbf{I}'$  similar to the previous work [4, 43]. Next, for the second stage, we attach the newly proposed layers for  $(\mathbf{a}, \mathbf{n}, \mathbf{s})$  to the generator, and fine-tune the whole model to generate the images  $\mathbf{I}$  from photometric rendering. Utilizing the volume representation that has roughly predicted geometry as the initial point, we can stabilize the training of the factorized representations.

## 3.4. Implementation details

For our first stage training, we adopt the FiLMed SIREN architecture of Chan *et al.* [4], which leverages a mapping network for conditioning  $\mathbf{z}$  and sinusoidal activation for improving image quality. Regarding newly attached layers in the second stage, we use the different activation functions for the linear layers that map the intermediate feature vectors to our factorized representations:  $\{\sigma_\theta, s_\theta : \text{ReLU}\}$ ,  $\{a_\theta, f_\theta : \text{Sigmoid}\}$ , and  $\{n_\theta : \text{Tanh}\}$ . We reduce the batch size to half for the second stage training. Further implementation details are presented in supplementary materials.

## 4. Experiments

### 4.1. Experiment setup

**Datasets.** We conduct experiments on two real-world datasets, CelebA [24] and CATs [53], and a synthetic dataset, BFM (Basel Face Model) [38]. CelebA dataset contains more than 200K face images of 10,177 unique identities. We utilize the aligned and cropped images with the resolution of  $128 \times 128$ . CATs dataset consists of over 6K in-the-wild cat images. We align and crop the cat faces based on the given annotations indicating the head of cats. Additionally, we utilize a synthetic face model, BFM dataset, to evaluate the quality of geometry with the ground truth depth. For all datasets, the images are resized to  $64 \times 64$  during the training, and we generate images at the resolution of  $128 \times 128$  and volume representations at  $128 \times 128 \times 128$  for mesh extraction. In the training, we only utilize the unposed 2D images without any additional labels or assumptions.

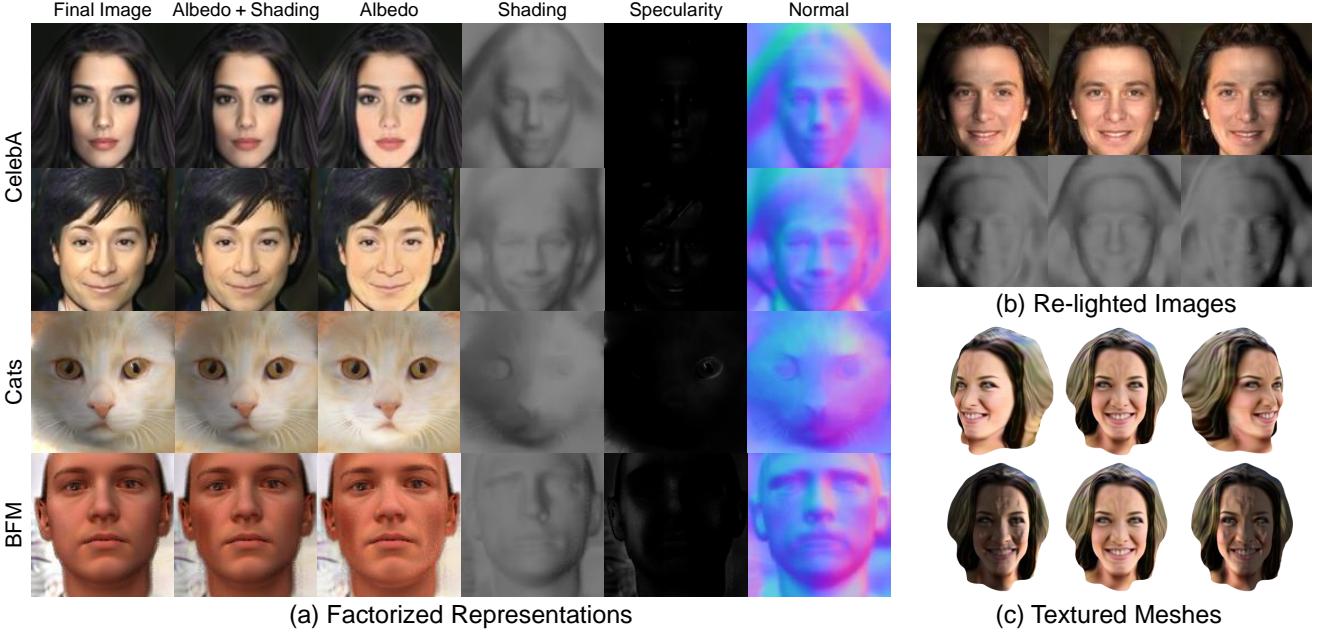


Figure 4. Visualization of the (a) factorized representations, (b) re-lighted images, and (c) textured meshes. Our method generates high-quality factorized representations and conducts re-lighting. The albedo-textured mesh generated from our method can be adopted to the conventional rendering pipelines. Note that we train the model without disentangling specularity for BFM dataset.

**Baselines.** As our baselines, we adopt the previous NeRF based 3D-aware image synthesis models: GRAF [43], pi-GAN [4], and GIRAFFE [32]. Note that GIRAFFE involves convolution operations for their neural rendering process, thus, we do not extract meshes from this model. We implement and train each model following the official codes.

## 4.2. Interpreting the factorized representations

**Disentanglement.** We present the qualitative results of factorized representations and light disentanglement in Fig. 4. As shown in Fig. 4(a), ours generates reasonable albedo, shading, specularity, and normal images. Surprisingly, even though we do not provide the model with additional supervision, ours successfully achieves factorization. Combining those factorized representations, our model outputs realistic final images. Note that predicted specularity makes the final images more natural.

Furthermore, from these factorized representations, we successfully re-light 2D images under various light conditions, as illustrated in Fig. 4(b). One can observe the shading image effectively reflects the given light conditions.

**Textured mesh.** Thanks to our view-independent and light-disentangled representations, we successfully generate high-quality albedo-textured meshes. Fig. 4(c) presents the meshes controlled via a conventional rendering pipeline. The rendered meshes look realistic in various viewpoints and under diverse light conditions. This is because the meshes are textured with albedo where all the effects of external environments are removed.

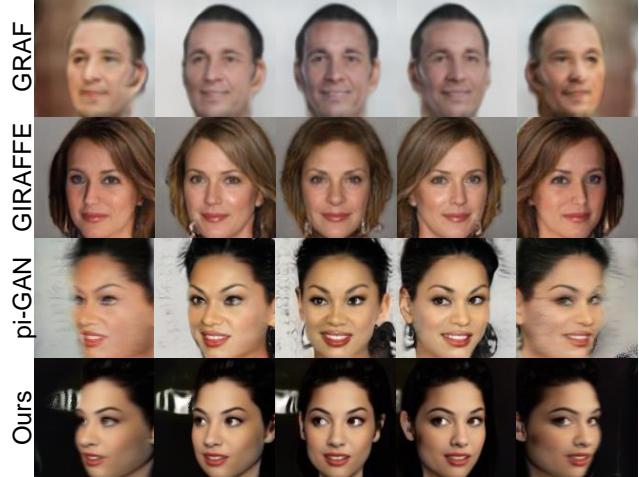


Figure 5. Comparison of multi-view images with baselines on CelebA and BFM dataset. Our method preserves the multi-view consistency while other baselines fail in extreme cases.

## 4.3. Comparison with baselines

**Metrics.** For the quantitative comparison, we show 1) the Frechet inception distance (FID) score for the image quality, 2) scale-invariant depth error (SIDE) for the geometry quality, and 3) projection FID for the quality of the textured mesh. To measure SIDE [9], we obtain depth maps of the generated images by volume rendering the distance between the given camera position and sampled points, as

View dependency	Model	CelebA	CATs	BFM	
		FID <sub>↓</sub>	FID <sub>↓</sub>	FID <sub>↓</sub>	SIDE <sub>↓</sub>
✓	GRAF	38.86	26.12	56.11	3.19
	pi-GAN	<b>14.25</b>	<b>10.77</b>	<b>14.83</b>	1.78
	GIRAFFE	22.02	21.01	18.46	-
✗	pi-GAN*	<i>18.10</i>	15.77	21.54	-
	Ours	20.17	<i>13.56</i>	<i>19.94</i>	<b>1.41</b>

Table 2. Quantitative evaluations with the FID and SIDE scores using CelebA, CATs, and BFM dataset. pi-GAN\* denotes a variant where the view-dependency of color representation is removed. Note that SIDE score is scaled by  $10^2$  for better visualization.

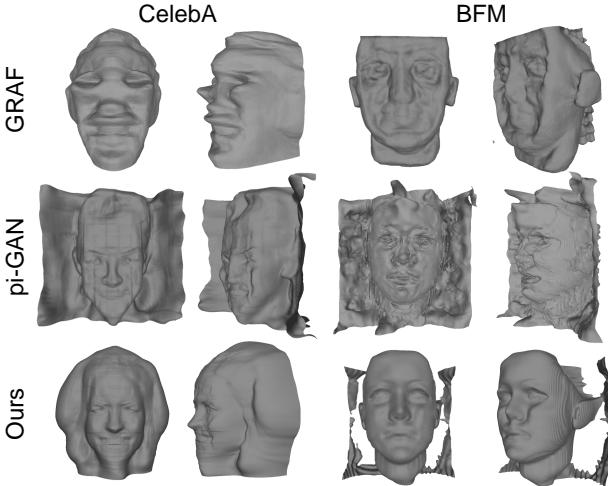


Figure 6. Comparison of the extracted meshes with baselines on CelebA and BFM dataset. Unlike other baselines, we can observe that our model generates more convex meshes.

done in the previous approaches. Moreover, for the ground truth depth of the generated images, we use the depth estimation model [1] pre-trained with BFM dataset. Projection FID score is the FID score of images projected from the generated textured meshes via the conventional rendering pipeline. To obtain textured meshes, we apply a marching cube algorithm on volume density to extract meshes and use albedo at surface points as texture. For the baselines, which require a view direction to predict the color, we use the color from the frontal view direction to obtain the texture. Afterward, we project the textured meshes into images with randomly-sampled viewpoints using Blender [6] and then measure the projection FID score. Further details are presented in supplementary materials.

**Quality of generated images** To evaluate the quality of generated images, we measure the FID scores on CelebA, CATs, and BFM dataset as illustrated in Table 2. Our method achieves reasonable FID scores, comparable to other baselines, involving marginal performance drop. We conjecture that this is mainly due to the view independency of our representations. The view dependent color rep-

Thresholds	pi-GAN							Ours						
	60 Proj. FID	56	70	55	80	90	100	110	20	30	40	50	60	70
Best	53.84							<b>50.17</b>						

Table 3. Projection FID scores of ours and baselines on CelebA dataset along with the sigma threshold of the marching cube algorithm.

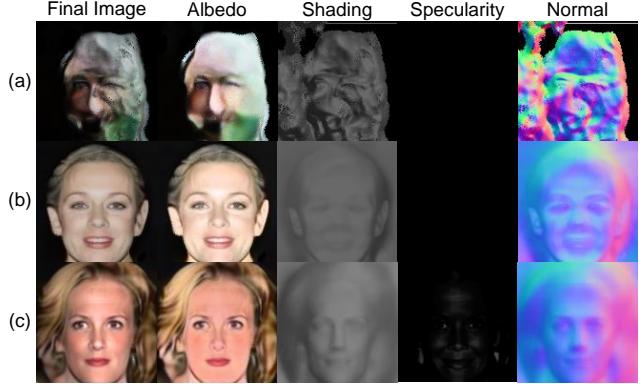


Figure 7. Visualization of factorized representations for ablation study. Images in (a) are from the model trained without normal prediction layer. (b) indicates generated images from our model without specularity disentanglement. The last row (c) is the results of our final model.

resentation of previous work effectively accounts for dataset bias (*e.g.* eye movement toward the front) and lowers the FID scores. For this reason, the variant of pi-GAN (*i.e.* pi-GAN\*) where the view-dependency of color representation is removed shows performance drop and the similar FID scores with the ones of ours. That is, our proposed approach has a comparable capability to pi-GAN in generating high-quality images.

**Multi-view consistency.** We compare the quality of the generated multi-view images from ours and baselines. As Fig. 5 describes, our method preserves reasonable multi-view consistency even in extreme viewpoints while other baselines generally fail. Since the baselines require view direction to predict the color, they fail in generating color from the extreme viewpoints, which is scarce in training datasets. Meanwhile, our model predicts view-independent color, resulting in more reasonable image generation quality from extreme viewpoints. Also, in the case of GIRAFFE [32], multi-view inconsistency becomes even worse (*e.g.* hair direction changes) since GIRAFFE [32] employs a 2D neural renderer composed of convolution operations, which hurts the internal 3D-aware feature.

**Quality of 3D representation.** To evaluate the quality of the generated 3D representations, we measure the depth error (SIDE) and the projection FID score of textured meshes. As presented in Table 2, our method outperforms baselines with a large gap in the depth error. That is, the quality of geometry are improved by applying our method, and this

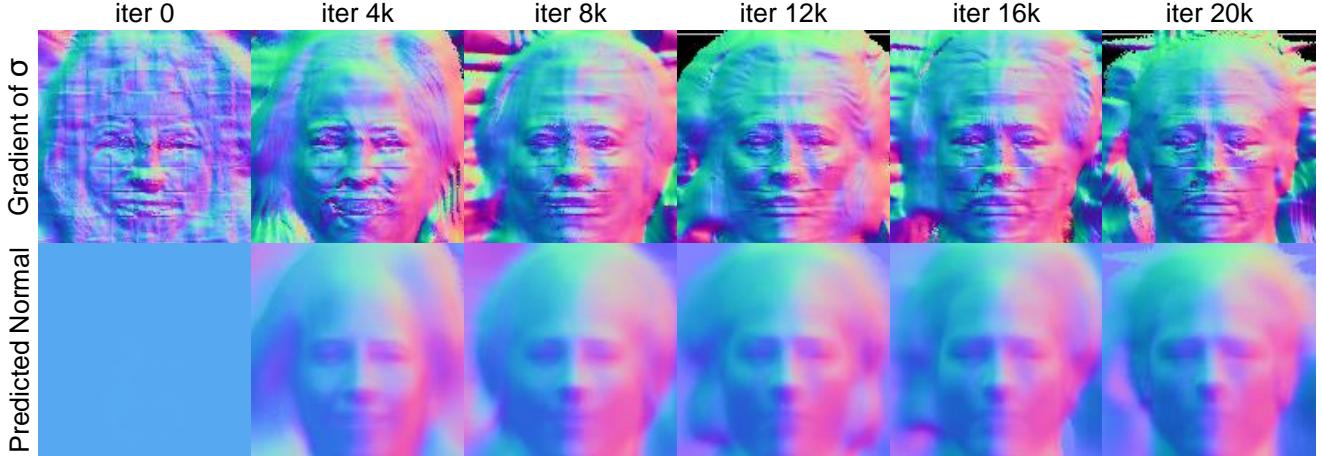


Figure 8. Normal refinements over iteration changes. As shown, gradient of  $\sigma$  gets gradually refined over training. In addition, predicted normal converges fast and works as a guidance to form better geometry.

can also be shown in Fig. 6. Regarding the projection FID scores, while our model shows the lower FID scores than pi-GAN, the improvements look marginal in Table 3. However, this is because the generated textured meshes from pi-GAN can seem realistic in frontal view, which is prevalent in the projected images since we sample viewpoints from Gaussian distribution. As shown in Fig. 6 and the Fig. 9 of supplementary, the dents of the geometry is not easily identified in the frontal view, unlike in other view.

#### 4.4. Ablation studies

Fig. 7 demonstrates the importance of each component of our model. Fig. 7(a) indicate the results of our model trained with the gradient of  $\sigma$  as normal, and (b) illustrates the results of our model trained with predicted normal (*i.e.* the output of normal prediction layer) and without specularity disentanglement. Fig. 7(c) are the results of our final model. As shown in Fig. 7(a) and described in Sec. 3, the training is unstable if we use the gradient of  $\sigma$  as normal. On the other hand, as presented in Fig. 7(b), the training becomes stable if we utilize the predicted normal instead of the gradient of  $\sigma$ . However, specularity still remains in albedo in (b), since we do not separate the specularity here. Finally, our full model generates valid factorized representations as illustrated in Fig. 7(c).

As described in Sec 3.2, the predicted normal from the normal prediction layer acts as weak supervision on the gradient of  $\sigma$ . This results in the improvement of the quality of geometry of the generated volume representations. We demonstrate the effectiveness of the normal prediction layer by visualizing the refinement progress of the gradient of  $\sigma$  and the predicted normal during Stage 2 training, as illustrated in Fig. 8. As the training proceeds, the predicted normal converges fast and presents more reasonable shape, and the gradient of  $\sigma$  gets refined accordingly.

## 5. Discussion

**Limitations.** First, we assume a single directional light for our method, which may not be appropriate for datasets containing complex light conditions. Next, while normal is defined on the surface of an object, we utilize volume rendering process to calculate normal. Although our model achieves successful results, focusing on surface to obtain normal could be one of the future work. As another trial, we attempted to calculate the reflection at each sampled 3D point and apply volume rendering on this reflected light, instead of volume rendering normal and calculating the reflection in images. Though this seems more technically sound, the quality of the generated geometry and images degraded.

**Ethical considerations.** Since we train ours on publicly available datasets, there may exist underlying bias in those datasets. This can result in lacking diversity of our generated images and 3D objects. Moreover, 3D representations can be obtained with inverse rendering techniques, the method can be abused for generating DeepFakes. We are aware of the fact that DeepFake is a substantial threat in our society, and we never approve of employing our model for DeepFake generation with inappropriate intentions.

## 6. Conclusion

We propose 3D-GIF that generates 3D-controllable objects beyond 3D-aware image generation. By combining physical reflectance prior with the model, we generate an object with factorized representations that are independent of external environments. With our method, we can generate 3D-controllable objects that can be adopted in 3D conventional rendering pipelines and have realistic geometry. We believe our work bridges a gap between content generation techniques and real-world 3D applications.

## References

- [1] Ibraheem Alhashim and Peter Wonka. High quality monocular depth estimation via transfer learning. *arXiv preprint arXiv:1812.11941*, 2018. 7, 13
- [2] Mark Boss, Raphael Braun, Varun Jampani, Jonathan T Barron, Ce Liu, and Hendrik Lensch. Nerd: Neural reflectance decomposition from image collections. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12684–12694, 2021. 3, 4
- [3] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018. 1
- [4] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5799–5809, 2021. 1, 2, 3, 4, 5, 6, 12, 13
- [5] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5939–5948, 2019. 3
- [6] Blender Online Community. *Blender - a 3D modelling and rendering package*. Blender Foundation, Stichting Blender Foundation, Amsterdam, 2018. 7, 13
- [7] 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. 12
- [8] Yu Deng, Jiaolong Yang, Dong Chen, Fang Wen, and Xin Tong. Disentangled and controllable face image generation via 3d imitative-contrastive learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5154–5163, 2020. 3
- [9] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. *arXiv preprint arXiv:1406.2283*, 2014. 6, 12
- [10] Matheus Gadelha, Subhransu Maji, and Rui Wang. 3d shape induction from 2d views of multiple objects. In *2017 International Conference on 3D Vision (3DV)*, pages 402–411. IEEE, 2017. 3
- [11] Kyle Genova, Forrester Cole, Avneesh Sud, Aaron Sarna, and Thomas Funkhouser. Local deep implicit functions for 3d shape. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4857–4866, 2020. 3
- [12] Kyle Genova, Forrester Cole, Daniel Vlasic, Aaron Sarna, William T Freeman, and Thomas Funkhouser. Learning shape templates with structured implicit functions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7154–7164, 2019. 3
- [13] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014. 1, 5
- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 12
- [15] Tong He, John Collomosse, Hailin Jin, and Stefano Soatto. Geo-pifu: Geometry and pixel aligned implicit functions for single-view human reconstruction. *arXiv preprint arXiv:2006.08072*, 2020. 3
- [16] Paul Henderson, Vagia Tsiminaki, and Christoph H Lampert. Leveraging 2d data to learn textured 3d mesh generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7498–7507, 2020. 3
- [17] Philipp Henzler, Niloy J Mitra, and Tobias Ritschel. Escaping plato’s cave: 3d shape from adversarial rendering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9984–9993, 2019. 1, 3
- [18] Le Hui, Rui Xu, Jin Xie, Jianjun Qian, and Jian Yang. Progressive point cloud deconvolution generation network. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XV 16*, pages 397–413. Springer, 2020. 3
- [19] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*, 2017. 1
- [20] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4401–4410, 2019. 1
- [21] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8110–8119, 2020. 1
- [22] Rosanne Liu, Joel Lehman, Piero Molino, Felipe Petroski Such, Eric Frank, Alex Sergeev, and Jason Yosinski. An intriguing failing of convolutional neural networks and the coordconv solution. *arXiv preprint arXiv:1807.03247*, 2018. 12
- [23] Shichen Liu, Shunsuke Saito, Weikai Chen, and Hao Li. Learning to infer implicit surfaces without 3d supervision. *arXiv preprint arXiv:1911.00767*, 2019. 3
- [24] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of the IEEE international conference on computer vision*, pages 3730–3738, 2015. 5
- [25] William E Lorensen and Harvey E Cline. Marching cubes: A high resolution 3d surface construction algorithm. *ACM siggraph computer graphics*, 21(4):163–169, 1987. 12, 13
- [26] Lars Mescheder, Andreas Geiger, and Sebastian Nowozin. Which training methods for gans do actually converge? In *International conference on machine learning*, pages 3481–3490. PMLR, 2018. 5
- [27] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4460–4470, 2019. 3
- [28] Mateusz Michalkiewicz, Jhony K Pontes, Dominic Jack, Mahsa Baktashmotlagh, and Anders Eriksson. Implicit sur-

- face representations as layers in neural networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4743–4752, 2019. 3
- [29] 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 *European conference on computer vision*, pages 405–421. Springer, 2020. 1, 3, 4, 12
- [30] Thu Nguyen-Phuoc, Chuan Li, Lucas Theis, Christian Richardt, and Yong-Liang Yang. Hologan: Unsupervised learning of 3d representations from natural images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7588–7597, 2019. 1, 3
- [31] Thu Nguyen-Phuoc, Christian Richardt, Long Mai, Yong-Liang Yang, and Niloy Mitra. Blockgan: Learning 3d object-aware scene representations from unlabelled images. *arXiv preprint arXiv:2002.08988*, 2020. 1, 3
- [32] Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11453–11464, 2021. 1, 3, 6, 7
- [33] 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. 3
- [34] 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. 3
- [35] Xingang Pan, XU Xudong, Chen Change Loy, Christian Theobalt, and Bo Dai. A shading-guided generative implicit model for shape-accurate 3d-aware image synthesis. In *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021. 13
- [36] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 165–174, 2019. 3
- [37] Dario Pavllo, Graham Spinks, Thomas Hofmann, Marie-Francine Moens, and Aurelien Lucchi. Convolutional generation of textured 3d meshes. *arXiv preprint arXiv:2006.07660*, 2020. 3
- [38] Pascal Paysan, Reinhard Knothe, Brian Amberg, Sami Romdhani, and Thomas Vetter. A 3d face model for pose and illumination invariant face recognition. In *2009 sixth IEEE international conference on advanced video and signal based surveillance*, pages 296–301. Ieee, 2009. 5
- [39] Sida Peng, Yuanqing Zhang, Yinghao Xu, Qianqian Wang, Qing Shuai, Hujun Bao, and Xiaowei Zhou. Neural body: Implicit neural representations with structured latent codes for novel view synthesis of dynamic humans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9054–9063, 2021. 3
- [40] Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual reasoning with a general conditioning layer. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018. 12
- [41] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015. 1
- [42] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Moriguchi, Angjoo Kanazawa, and Hao Li. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2304–2314, 2019. 3
- [43] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. Graf: Generative radiance fields for 3d-aware image synthesis. *arXiv preprint arXiv:2007.02442*, 2020. 1, 3, 4, 5, 6, 12
- [44] Yichun Shi, Divyansh Aggarwal, and Anil K Jain. Lifting 2d stylegan for 3d-aware face generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6258–6266, 2021. 3
- [45] Dong Wook Shu, Sung Woo Park, and Junseok Kwon. 3d point cloud generative adversarial network based on tree structured graph convolutions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3859–3868, 2019. 3
- [46] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. *arXiv preprint arXiv:1906.01618*, 2019. 3
- [47] Pratul P Srinivasan, Boyang Deng, Xiuming Zhang, Matthew Tancik, Ben Mildenhall, and Jonathan T Barron. Nerv: Neural reflectance and visibility fields for relighting and view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7495–7504, 2021. 3
- [48] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6142–6151, 2020. 3
- [49] Nanyang Wang, Yinda Zhang, Zhuwen Li, Yanwei Fu, Wei Liu, and Yu-Gang Jiang. Pixel2mesh: Generating 3d mesh models from single rgb images. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 52–67, 2018. 3
- [50] 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. 3
- [51] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Ronen Basri, and Yaron Lipman. Multiview neural surface reconstruction by disentangling geometry and appearance. *arXiv preprint arXiv:2003.09852*, 2020. 3

- [52] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Ronen Basri, and Yaron Lipman. Multiview neural surface reconstruction by disentangling geometry and appearance. *arXiv preprint arXiv:2003.09852*, 2020. [4](#)
- [53] Weiwei Zhang, Jian Sun, and Xiaoou Tang. Cat head detection-how to effectively exploit shape and texture features. In *European Conference on Computer Vision*, pages 802–816. Springer, 2008. [5](#)
- [54] Xiuming Zhang, Pratul P Srinivasan, Boyang Deng, Paul Debevec, William T Freeman, and Jonathan T Barron. Nerfactor: Neural factorization of shape and reflectance under an unknown illumination. *arXiv preprint arXiv:2106.01970*, 2021. [3, 4](#)
- [55] Zerong Zheng, Tao Yu, Yebin Liu, and Qionghai Dai. Pamir: Parametric model-conditioned implicit representation for image-based human reconstruction. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021. [3](#)

## A. Overview

In this supplementary material, we provide further implementation and experiment details. Afterward, we illustrate additional experiment results and analyses. Finally, we introduce a concurrent work, ShadeGAN, and present comparisons with them. Furthermore, we present qualitative results on viewpoint and light condition control with videos, attached to supplementary materials.

## B. Implementation Details

### B.1. Architecture Details

**FiLMed SIREN network.** We adopt FiLMed SIREN network proposed by Chan *et al.* [4] as our backbone, which leverages a mapping network conditioning and sinusoidal activations. The FiLMed SIREN backbone is formulated as

$$\begin{aligned} h_\theta(\mathbf{x}_j) &= \phi_{n-1} \circ \phi_{n-2} \circ \dots \circ \phi_0(\mathbf{x}_j), \\ \phi_i(\mathbf{x}_j^{(i)}) &= \sin(\gamma_i \cdot (\mathbf{W}_i \mathbf{x}_j^{(i)} + \mathbf{b}_i) + \beta_i), \end{aligned} \quad (12)$$

where  $\mathbf{x}_j$  indicates the  $j^{\text{th}}$  sampled point along the given camera ray,  $\phi_i$  denotes the  $i^{\text{th}}$  FiLMed linear layer of the MLP followed by sinusoidal activation function (*i.e.* Sine nonlinearity), and  $\mathbf{x}_j^{(i)}$  is the input of  $\phi_i$ .  $\gamma^i$  and  $\beta^i$  are the affine transformation parameters of the FiLM layer [40], which are the output of the mapping network conditioned on the latent vector  $\mathbf{z} \in \mathbb{R}^{256}$ . Here, we implement the mapping network using an MLP of three hidden layers with 256 units each followed by leaky-ReLU activation with a slope of 0.2. The FiLMed SIREN network is composed of eight FiLMed SIREN hidden layers of 256 units each.

**Normal prediction layer.** We adopt the normal prediction layer that consists of MLP with two linear layers. The first linear layer has 16 units followed by leaky-ReLU activation with a slope of 0.2. The second linear layer followed by Tanh activation and normalization to a unit vector maps the intermediate feature to the predicted normal  $\mathbf{n} \in \mathbb{R}^3$ .

**Discriminator.** We leveraged CoordConv [22] layers and residual connections [14] for the discriminator, referring to pi-GAN [4]. Unlike pi-GAN, we did not apply progressive discriminator.

### B.2. Training Details

For the training, we utilized Adam optimizer with  $\beta_1 = 0$  and  $\beta_2 = 0.9$ . We set learning rates to  $6 \times 10^{-5}$  and  $2 \times 10^{-4}$  for the generator and the discriminator, respectively. Also, we set  $\lambda$  as 0.5 in Eq. (8) and  $\lambda_n$  as 20 in Eq. (11) in our main paper. The number of iterations is set to 200,000 for the first stage and 20,000 for the second stage. We conduct staged training where we first train Stage 1 and then jointly optimize Stage 1 and Stage 2. For stable training, hierarchical sampling [29] is adopted to our training procedure. Also, we used Automatic Mixed Precision

(AMP) in PyTorch for efficiency in computation and memory usage. We utilize a pinhole perspective camera with a field of view (FOV) of  $12^\circ$  for CelebA, CATs, and BFM. Batch size and the number of samples along each ray are set to 56 and 12 for CelebA and BFM, and 28 and 24 for CATs. For all datasets, the images are resized into  $64 \times 64$  for the training. Camera parameters and light condition parameters were adjusted according to each dataset, as described in Table 4.

Dataset	$dist_{cam}$	Camera param.		Light condition param.			
		$\sigma_v / \text{range}_v$	$\sigma_h / \text{range}_h$	$\mu_x$	$\sigma_x$	$\mu_y$	$\sigma_y$
CelebA	Gauss.	0.15	0.3	0	0.27	0.39	0.07
CATs	Uni.	(-0.5, 0.5)	(-0.4, 0.4)	0	0.27	0.39	0.07
BFM	Gauss.	0.15	0.3	0	0.35	0	0.15

Table 4. Camera parameters and light condition parameters for each dataset.  $dist_{cam}$ ,  $v$ , and  $h$  denote camera pose distribution, pitch, and yaw.  $x$  and  $y$  denote the direction of light. Note that  $\sigma_{v,h}$  in the camera parameters are for Gaussian distributions, and  $\text{range}_{v,h}$  are for uniform distribution.

### B.3. Albedo-Textured Mesh Extraction

Since our model generates view-independent and light-disentangled representations, we can obtain high-quality albedo-textured meshes, which are easily applicable to the conventional computer graphics rendering pipelines. To obtain the textured mesh, our model first predicts the volume density of each point in a  $128 \times 128 \times 128$  voxel grid with the randomly-sampled latent vector  $\mathbf{z}$  as a condition. Here, we set a length of the voxel grid to include all the sampled points considering hyper-parameters (*e.g.*, ray start, ray end, FOV, etc.) of each dataset. Then, we extract a mesh (*i.e.*, surface points and faces) from the volume density using a marching cube algorithm [25]. Lastly, we texture the obtained mesh with the albedo predicted at each surface point by our model. Unlike the previous NeRF-based methods [4, 43] that generate view-dependent and light-entangled colors, 3D-GIF successfully generates meshes textured with albedo since our approach does not need view direction when predicting the albedo and separately controls the light condition.

## C. Experiment Details

As quantitative evaluations, we compare the FID, depth error, and projection FID of our model and the baselines. For each evaluation, we generated 2,000 images using each model.

**Depth error.** To evaluate the quality of the generated geometry, we measure scale-invariant depth error (SIDE) [9] between the depth maps of the generated images and their ground truth. As done in the previous approaches [7], we extract depth maps of the generated images by volume ren-



Figure 9. 3D-rendered images from pi-GAN and our method. As illustrated, ours shows more realistic results than pi-GAN. Those images are utilized to measure the projection FID.

dering the distance between the given camera position and each sampled point. In addition, to obtain the ground truth depth maps of the generated images, we leverage a depth estimation model [1]. We trained the depth estimation model with the training dataset of BFM. Then, we estimated the ground truth depth map of each generated image using the pre-trained depth estimation model.

**Projection FID.** As mentioned in Sec. 4.3 in our main paper, we measure projection FID of textured meshes extracted from ours and state-of-the-art model [4] as a quantitative comparison. Here, projection FID indicates the FID score of projected images from the textured meshes. We first extract surface meshes using a marching cube algorithm [25]. Since the estimated surface changes according to the sigma threshold of the algorithm, we extract meshes at different thresholds from 10 to 200 with an interval of 10. Then, we texture the meshes with the generated colors. Since the baseline model requires view direction to predict colors, we provide the model with the frontal view direction as a condition. To project images from the textured mesh, we utilize Blender [6] with sampled viewpoints and light conditions following our training setting. Moreover, we added random solid colors as backgrounds to the projected images since they lack valid backgrounds, unlike the real images. Table 3 in our main paper presents projection FID scores according to different sigma thresholds and the best score of each model. Fig. 9 presents the examples of projected images of ours and the baseline at the sigma threshold of the best score.

## D. Additional Experiments

We report additional qualitative results using CelebA and CATs datasets. Fig. 12 and Fig. 13 present our generated albedo-textured meshes and their re-lighted results. As shown, our model generates realistic albedo-textured meshes applicable to the conventional rendering pipeline.

To further illustrate our robustness in generating high-quality factorized representations, we show randomly sampled results in Fig. 14 and Fig. 15. Note that those results are obtained by varying a random seed from 0 to 25.

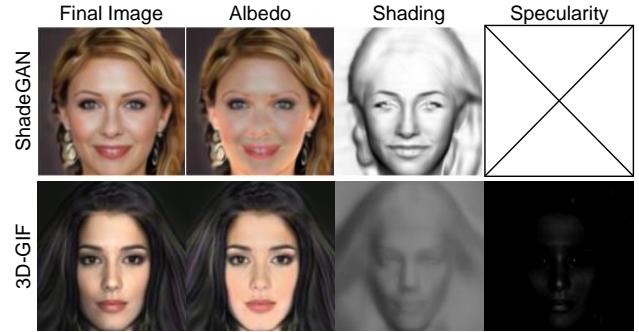


Figure 10. Visualization of factorized representations from ShadeGAN and our method. Generated albedo from ShadeGAN contains specularity while ours does not.

## E. Comparison with ShadeGAN

ShadeGAN [35], which is concurrent work presented in arXiv recently, proposes a framework to improve the quality of geometry from the previous NeRF-based generative methods by generating light-disentangled albedo under the Lambertian assumption. While albedo is independent of view direction by its definition, ShadeGAN predicts view-dependent albedo conditioned on a certain view direction. This view dependency enables the model to express specularity even under the Lambertian shading assumption where the specular reflection is not involved. However, as also described in Sec. 3.2 of our main paper, the view dependency of the generated volume representation hampers its applicability in various 3D applications. This is because they still need redundant rendering for every desired viewpoint, and their albedo contains specularity. Although ShadeGAN improved the geometry quality, as shown in Fig. 10 and Fig. 11, their generated albedo image still contains specularity, and their re-lighted images cannot fully reflect the given light conditions. Note that we brought the result images of ShadeGAN from their paper to compare with ours since their official code is not publicly available yet.

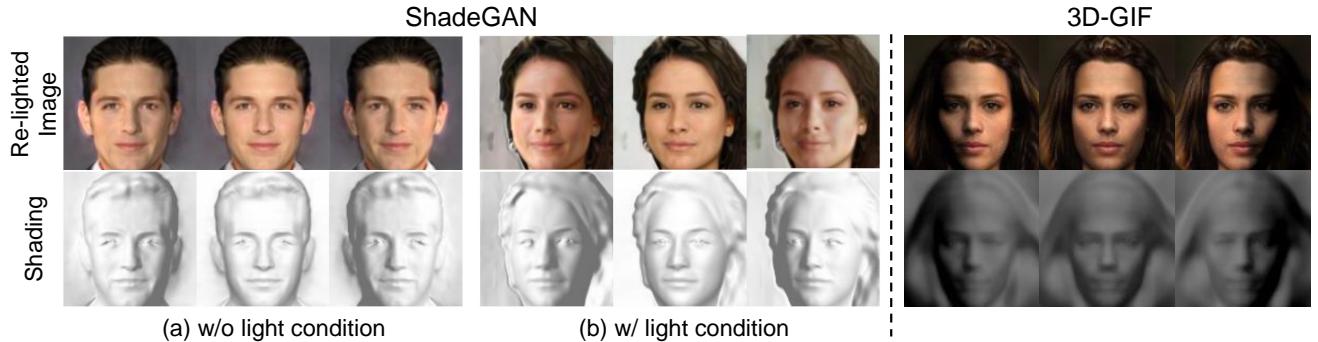


Figure 11. Re-lighted images from ShadeGAN and our method. As shown, our model reflects the light changes more effectively than ShadeGAN. The ShadeGAN presents two settings, whether light condition is concatenated to intermediate feature or not. (a) images are generated without light concatenation, and (b) images are generated with light concatenation.

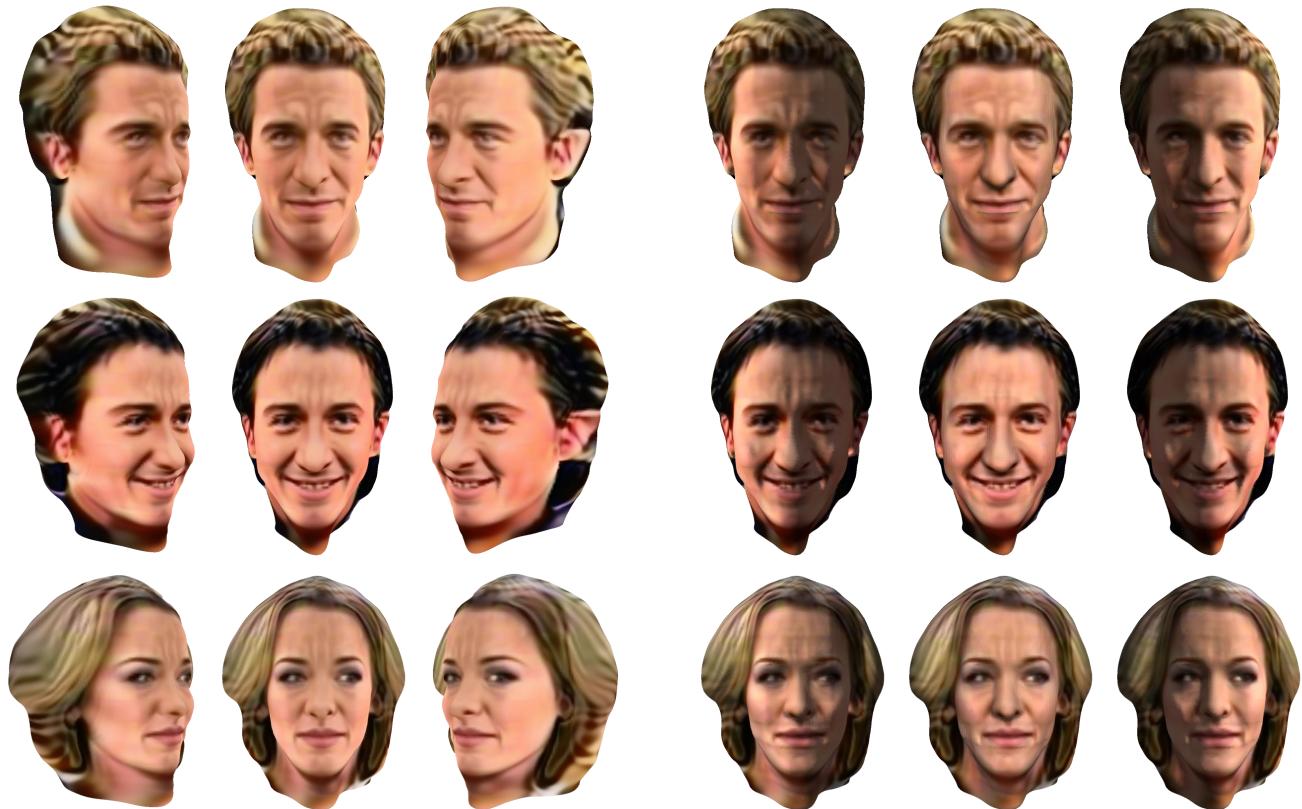
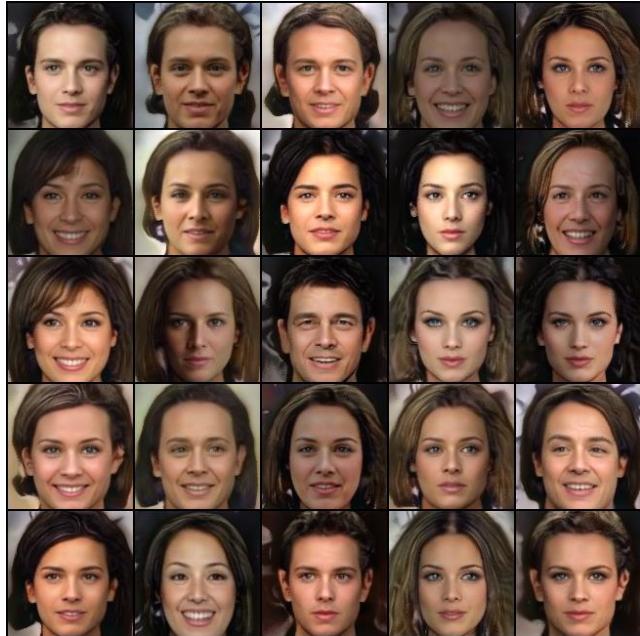


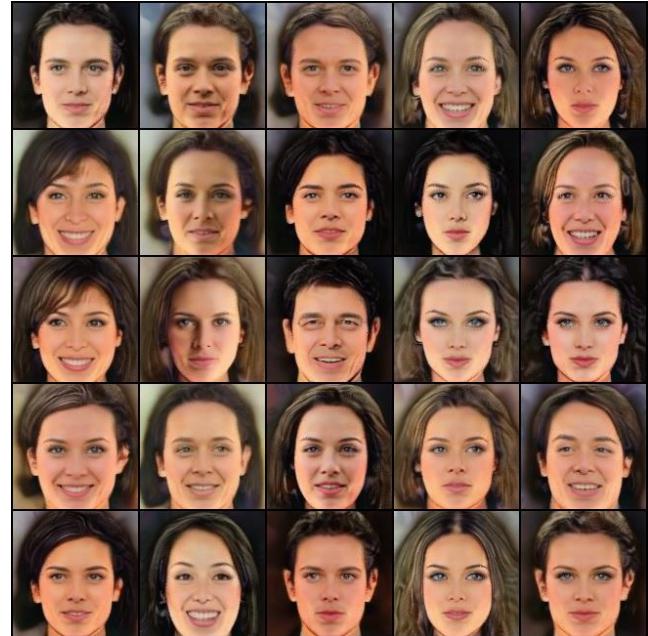
Figure 12. Our generated albedo-textured meshes and their re-lighted results on CelebA dataset.



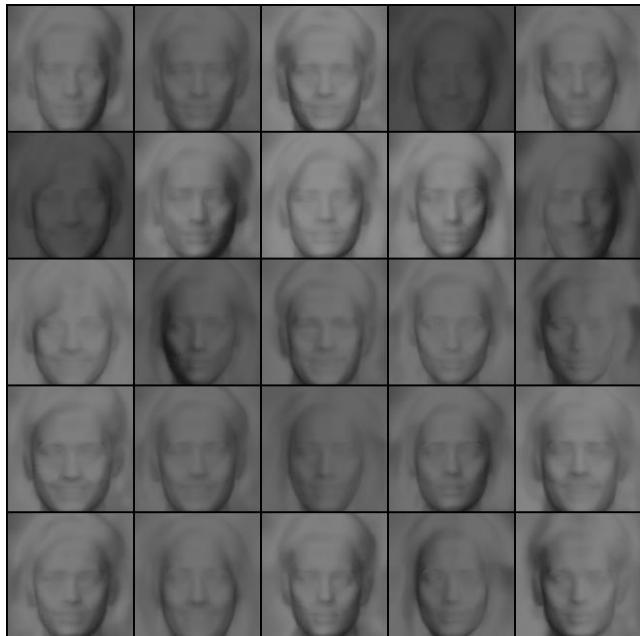
Figure 13. Our generated albedo-textured meshes and their re-lighted results on CATs dataset.



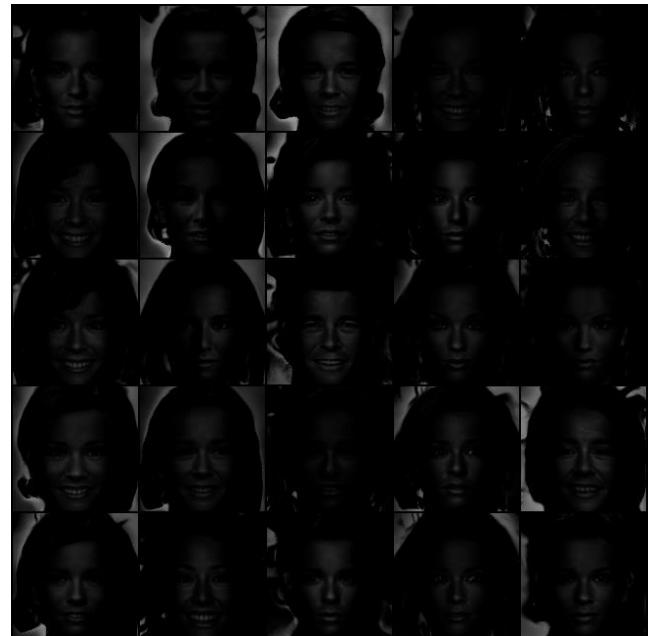
(a) Final Image



(b) Albedo



(c) Shading



(d) Specularity

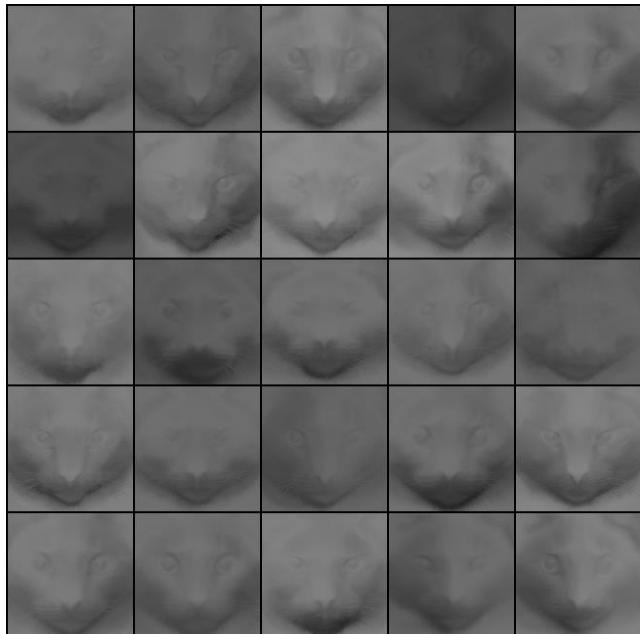
Figure 14. Random samples of CelebA. (a), (b), (c), and (d) indicate final images, albedo, shading, and specularity, respectively.



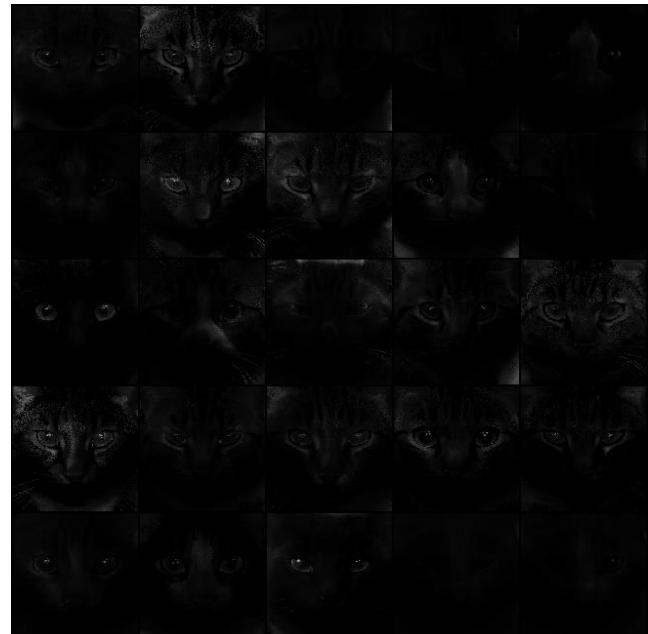
(a) Final Image



(b) Albedo



(c) Shading



(d) Specularity

Figure 15. Random samples of CATs. (a), (b), (c), and (d) indicate final images, albedo, shading, and specularity, respectively.