



# VoLux-GAN: A Generative Model for 3D Face Synthesis with HDRI Relighting

Feitong Tan  
Google  
USA  
feitongtan@google.com

Sean Fanello  
Google  
USA  
seanfanello@google.com

Abhimitra Meka  
Google  
USA  
abhim@google.com

Sergio Orts-Escolano  
Google  
USA  
sorts@google.com

Danhang Tang  
Google  
USA  
danhangtang@google.com

Rohit Pandey  
Google  
USA  
rohitpandey@google.com

Jonathan Taylor  
Google  
USA  
jontaylor@google.com

Ping Tan  
Simon Fraser University  
Canada  
pingtan@sfu.ca

Yinda Zhang  
Google  
USA  
yindaz@google.com

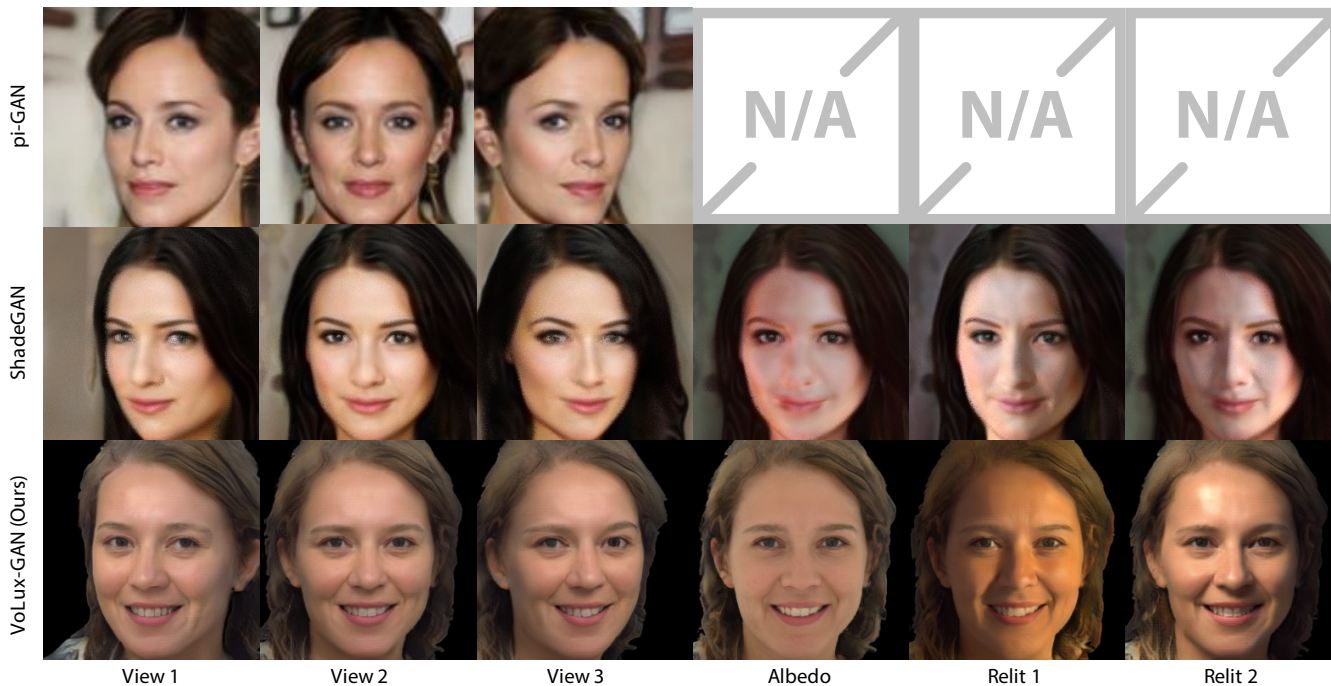


Figure 1: We propose VoLux-GAN, a 3D-aware generator that produces faces with full HDRI relighting capability. Here we show a comparison of images generated by VoLux-GAN and related work pi-GAN [Chan et al. 2021] (which does not support relighting) and ShadeGAN [Pan et al. 2021].



This work is licensed under a Creative Commons Attribution International 4.0 License.

SIGGRAPH '22 Conference Proceedings, August 7–11, 2022, Vancouver, BC, Canada  
© 2022 Copyright held by the owner/author(s).  
ACM ISBN 978-1-4503-9337-9/22/08.  
<https://doi.org/10.1145/3528233.3530751>

## ABSTRACT

We propose VoLux-GAN, a generative framework to synthesize 3D-aware faces with convincing relighting. Our main contribution is a volumetric HDRI relighting method that can efficiently accumulate albedo, diffuse and specular lighting contributions along each 3D ray for any desired HDR environmental map. Additionally, we show the importance of supervising the image decomposition

process using multiple discriminators. In particular, we propose a data augmentation technique that leverages recent advances in single image portrait relighting to enforce consistent geometry, albedo, diffuse and specular components. Multiple experiments and comparisons with other generative frameworks show how our model is a step forward towards photorealistic relightable 3D generative models. Code and pre-trained models are available at: <https://github.com/google/volux-gan>.

## CCS CONCEPTS

• **Computing methodologies** → **Machine learning; Computer vision; Rendering.**

## KEYWORDS

Neural Rendering, Relighting, Generative Model

### ACM Reference Format:

Feitong Tan, Sean Fanello, Abhimitra Meka, Sergio Orts-Escolano, Danhang Tang, Rohit Pandey, Jonathan Taylor, Ping Tan, and Yinda Zhang. 2022. Volux-GAN: A Generative Model for 3D Face Synthesis with HDRI Relighting. In *Special Interest Group on Computer Graphics and Interactive Techniques Conference Proceedings (SIGGRAPH '22 Conference Proceedings)*, August 7–11, 2022, Vancouver, BC, Canada. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3528233.3530751>

## 1 INTRODUCTION

Generating synthetic novel human subjects with convincing photorealism is one of the most desired capabilities for automatic content generation and pseudo ground truth synthesis for machine learning. Such data generation engines can thus benefit many areas including the gaming and movie industries, telepresence in mixed reality, and computational photography. In order to achieve realism and flexibility when delivered in specific applications, the generated images should: 1) be enriched in details, e.g. with high resolution; 2) support free viewpoint rendering to deliver immersive 3D experiences; 3) adapt to novel environmental illumination for realism; 4) synthesize novel identities for scalable data diversity.

Motivated by these principles, in this paper we propose a neural human portrait generator, which delivers compelling rendering quality on arbitrary camera viewpoints and under any desired illumination. With the success of Neural Radiance Field (NeRF) on volumetric rendering [Mildenhall et al. 2020] and Generative Adversarial Networks (GAN) on image generation [Karras et al. 2019], 3D-aware generators [Chan et al. 2021; Gu et al. 2021; Pan et al. 2021] have been proposed as a promising solution, which combine the merits of both. By learning from a collection of portrait images, these methods are able to generate NeRF models from randomly sampled latent codes, which result in impressive free viewpoint rendering capabilities despite arguable underlying geometry quality and multi-view consistency. Concurrent work proposed by [Pan et al. 2021] adds a shading model to enforce multi-lighting constraints during training, however the method shows substantial limitations in terms of photorealism and does not allow for full HDRI relighting.

In this work, we propose a 3D aware generative model with HDRI relighting supervised by adversarial losses. To overcome the limitations of prior arts, we contributed to two main aspects:

*Volumetric HDRI Relighting.* We propose a novel approach of the volumetric rendering function that naturally supports efficient HDRI relighting. The core idea relies on the intuition that diffuse and specular components can be efficiently accumulated per-pixel when pre-filtered HDR lighting environments are used [Greene 1986; Ramamoorthi and Hanrahan 2001]. This was successfully applied to single image portrait relighting [Pandey et al. 2021], and here we introduce an alternative formulation to allow for volumetric HDRI relighting. Differently from [Nestmeyer et al. 2020; Pandey et al. 2021; Wang et al. 2020] that predict surface normals and calculate the shading with respect to the light sources (for a given HDR environment map), we propose to directly integrate the diffuse and specular components at each 3D location along the ray according to their local surface normal and viewpoint direction. Simultaneously, an albedo image and neural features are accumulated along the 3D ray. Finally, a neural renderer combines the generated outputs to infer the final image.

*Supervised Image Decomposition.* Though producing impressive rendering quality, the geometry from 3D-aware generators is often incomplete or inaccurate [Chan et al. 2021; Gu et al. 2021]. As a result, the model tends to bias the image quality for highly sampled camera views (e.g. front facing), but starts to show unsatisfactory multi-view consistency and 3D perception, breaking the photorealism when rendered from free-viewpoint camera trajectories. Additionally, high quality geometry is particularly important for relighting since any underlying reflectance models rely on accurate surface normal directions in order to correctly accumulate the light contributions from the HDR environment map.

Similarly, decomposing an image into albedo, diffuse and specular components without explicit supervision could lead to artifacts and inconsistencies, since, without any explicit constraints, the network could encode details in any channel even though it does not follow light transport principles. For instance in Fig. 1, the albedo image generated by previous methods [Pan et al. 2021] contains clear shading information, whereas the expected albedo (i.e. flat lit image) should be closer to ours. At the same time, such supervision is not available for in-the-wild datasets like FFHQ [Karras et al. 2019].

Motivated by this, and inspired by other works that apply pseudo-groundtruth labels [Chogovadze et al. 2021] or synthetic renderings [Saito et al. 2020; Tan et al. 2021; Wood et al. 2021; Zhu et al. 2020] for in-the-wild tasks, we propose a data augmentation technique to explicitly supervise the image decomposition in geometry, albedo, diffuse and specular components. In particular, we employ the work of [Pandey et al. 2021] to generate albedo, geometry, diffuse, specular and relit images for each image of the dataset, and have additional discriminators guide the intrinsic decomposition during the training. This technique alone, however, would guide the generative model to synthesize images that are less photorealistic since their quality upper bound would depend on the specific image decomposition and relighting algorithm used as supervision (e.g. [Pandey et al. 2021]). In order to address this, we also add a final discriminator on the original images, which will guide the network towards real photorealism and higher order light transport effects such as specular highlights and subsurface scattering.

We summarize the contributions of this paper: 1) We propose a novel approach to generate HDRI relightable 3D faces with a volumetric rendering framework. 2) Supervised adversary losses are leveraged to increase the geometry and relighting quality, which also improves multi-view consistency. 3) Exhaustive experiments demonstrated the effectiveness of the framework for image synthesis and relighting.

## 2 RELATED WORK

*2D Image Generation.* Generating convincing renderings of humans is a very active trend in the field of neural rendering [Tewari et al. 2020]. Here, we consider works that rely on a generative adversarial framework [Goodfellow et al. 2014] to synthesize photorealistic portraits. High quality results have been demonstrated by multiple early works [Durugkar et al. 2017; Mordido et al. 2018; Zhang et al. 2019] and since the groundbreaking work of StyleGAN [Karras et al. 2019], the community has made tremendous progress in synthesizing photorealistic and high resolution images [Brock et al. 2019; Choi et al. 2020; Karras et al. 2021, 2020] with methods focusing on addressing most of the common issues with GANs including stability [Karras et al. 2020], resolution [Brock et al. 2019] and aliasing [Karras et al. 2021]. These approaches generate impressive photorealistic images, but results typically lack free-viewpoint rendering and/or multi-view consistency.

*3D Aware Generation.* Many recent approaches incorporated the use of geometry and its multi-view consistency to allow for 3D aware synthesis. [Alhaja et al. 2018; Chan et al. 2021; Chen et al. 2021; Gu et al. 2021; Liao et al. 2020; Nguyen-Phuoc et al. 2019; Niemeyer and Geiger 2021; Phuoc et al. 2020; Zhou et al. 2021; Zhu et al. 2018]. Past works rely on voxels [Gadelha et al. 2017; Nguyen-Phuoc et al. 2019; Phuoc et al. 2020; Zhu et al. 2018], meshes [Szabó et al. 2019], face models [Buehler et al. 2021] or shape primitives [Liao et al. 2020] as the 3D representation for image generation, but the majority have been limited to low resolution image generation. Inspired by the success of NeRF [Mildenhall et al. 2020], state-of-art approaches adopt implicit volumetric rendering framework, and require only unconstrained images for training [Chan et al. 2021; Gu et al. 2021; Niemeyer and Geiger 2021; Schwarz et al. 2020]. Though the proposed representations are essentially 3D, the underlying geometry could still be inconsistent across camera views. Additionally, these methods lack relighting capabilities.

*Relightable Generative Models.* Relightable NeRF models [Boss et al. 2020, 2021; Zhang et al. 2021b,a] have shown that full image decomposition is possible when explicit multi-view imagery is provided as supervision. As for generative networks, the concurrent work of [Pan et al. 2021] is, to the best of our knowledge, the first at enabling relightability into generative model in a volumetric 3D framework. The method enforces both multi-view and multi-lighting consistency to allow controllable viewpoint and illumination. This approach, however, adopts a simplified Lambertian model and only supports one specific light direction at the time and extending it to full HDR relighting is computationally prohibitive.

*Intrinsic Image Decomposition.* Decomposing an image into albedo, geometry and reflectance components has achieved using model-fitting techniques [Barron and Malik 2015; Meka et al. 2017] and

deep learning based approaches [Kanamori and Endo 2018; Meka et al. 2018; Ren et al. 2015; Xu et al. 2018] that attempt at inferring image properties from one or multiple images. Very recently, state-of-art image based portrait relighting methods [Nestmeyer et al. 2020; Pandey et al. 2021; Tajima et al. 2021; Wang et al. 2020] have shown impressive results by predicting explicit surface normals, albedo and shading information to formulate the interaction between light sources and geometry. These approaches usually rely on a specific shading model (e.g. Phong) and a neural renderer to synthesize the final image.

*Our Approach.* In contrast, we propose a volumetric generative model that supports full HDR relighting. We show how we can efficiently aggregate albedo, diffuse and specular components within the 3D volume. Thanks to the explicit supervision in our adversarial losses, we demonstrate that the method can perform such a full image component decomposition for novel face identities, starting from a randomly sampled latent code.

## 3 VOLUX-GAN FRAMEWORK

In this section, we introduce our neural generator that produces novel faces that can be rendered at free camera viewpoints and relit under an arbitrary HDR environment light map. Our method starts from a neural implicit field that takes a randomly sampled latent vector as input and produces an albedo, volume density, and reflectance properties for any queried 3D location. These outputs are then aggregated via volumetric rendering to produce low resolution albedo, diffuse shading, specular shading, and neural feature maps. These intermediate outputs are then upsampled to high resolution and fed into a neural renderer to produce relit images. The overall framework is depicted in Figure 2.

### 3.1 Preliminaries: Neural Volumetric Rendering.

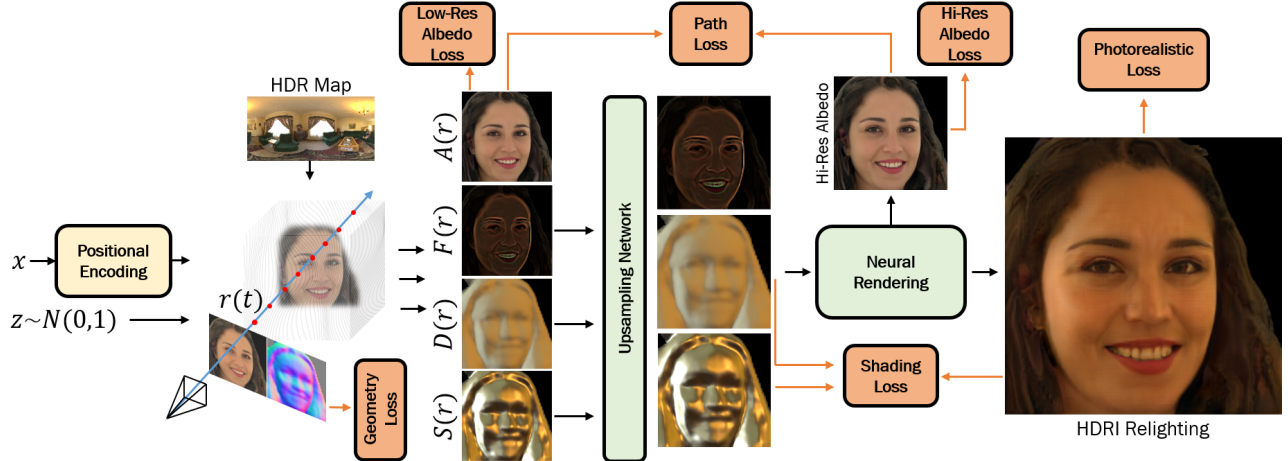
To aid the reader, we first briefly introduce the neural volumetric rendering framework originally presented in [Mildenhall et al. 2020]. There, the 3D appearance of an object of interest is encoded into a neural implicit field implemented using a multilayer perceptron (MLP), which takes a 3D coordinate  $x \in R^3$  and viewing direction  $\mathbf{d} \in S^2$  as inputs and outputs a volume density  $\sigma \in R^+$  and view-dependent color  $\mathbf{c} \in R^3$ . To render an image, the pixel color  $\mathbf{C}$  is accumulated along each camera ray  $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$  as

$$\mathbf{C}(\mathbf{r}, \mathbf{d}) = \int_{t_n}^{t_f} T(t)\sigma(\mathbf{r}(t))\mathbf{c}(\mathbf{r}(t), \mathbf{d})dt, \quad (1)$$

where  $T(t) = \exp(-\int_{t_n}^t \sigma(\mathbf{r}(s))ds)$  and bounds  $t_n$  and  $t_f$ . Compared to surface based rendering, volumetric rendering more naturally handles translucent materials and regions with complex geometry such as thin structures.

### 3.2 Generative Neural Implicit Intrinsic Field.

Similar to other state-of-the-art 3D-aware generators [Chan et al. 2021; Gu et al. 2021; Pan et al. 2021], we train a MLP-based neural implicit field conditioned on a latent code  $\mathbf{z}$  sampled from a Gaussian distribution  $N(0, I)^d$  and extend it to support HDRI relighting.



**Figure 2: VoLux-GAN Framework.** Starting from a latent code we can efficiently accumulate albedo  $A(\mathbf{r})$ , surface normals  $N(\mathbf{r})$ , diffuse  $D(\mathbf{r})$ , specular components  $S(\mathbf{r})$ , and a feature map  $F(\mathbf{r})$  along the 3D ray  $\mathbf{r}(t)$  for any given HDR map. An upsampling strategy and a neural renderer synthesize the final relit image.

We adopt a Phong shading model [Phong 1975], where the illumination of each point is determined by albedo, diffuse, and specular component. Therefore, instead of having the network predict per-point radiance and directly obtaining a color image (via Eq. 1), our network produces per-point albedo ( $\alpha$ ), density ( $\sigma$ ) and reflectance properties from separate MLP heads. The normal directions are obtained via the spatial derivative of the density field, which are used together with HDR illumination to compute diffuse and specular shading. Similar to [Pandey et al. 2021], rather than explicitly using the Phong model for the final rendering, we feed the albedo, diffuse and specular components to a lightweight neural renderer, which can also model higher order light transport effects.

**3.2.1 Efficient Shading Computation.** Concurrent work [Pan et al. 2021] assumes Lambertian shading from a single light source. Extending this to support full HDR illumination would require the integration of the shading contribution from multiple positional lights, making the approach computationally prohibitive, especially when performed at training time for millions of images. Inspired by the success of recent image based portrait relighting work [Pandey et al. 2021], we adopt a method designed for real-time shading rendering under HDR illumination [Greene 1986; Ramamoorthi and Hanrahan 2001]. The core idea is to approximate the diffuse and specular components using a preconvolved HDRI map. Specifically, we first preconvolve the given HDRI map ( $\mathbf{H}$ ) into light maps ( $L_{n_i}, i = 1, 2, \dots, N$ ) with cosine lobe functions corresponding to a set of pre-selected Phong specular exponents ( $n_i, i = 1, 2, \dots, N$ ) [Miller and Hoffman 1984]. The diffuse shading  $D$  is the first light map (i.e.  $n = 1$  above) following the surface normal direction, and the specular shading is defined as a linear combination of all light maps indexed by the reflection direction. In order to capture possible diverse material properties of the face, we let the network estimate the blending weights ( $\omega_i$ ) with another MLP branch, which are then used for the specular component  $S$ .

**3.2.2 Volumetric Shading Rendering.** Typically, a reflectance model is defined on a surface [Phong 1975] and relighting methods [Nestmeyer et al. 2020; Pandey et al. 2021; Wang et al. 2020] explicitly estimate surface normals from a single image. Here, we propose a volumetric formulation to compute albedo, diffuse and view-dependent specular shading maps as:

$$\begin{aligned}
 A(\mathbf{r}) &= \int_{t_n}^{t_f} T(t)\sigma(\mathbf{r}(t))\alpha(\mathbf{r}(t))dt \\
 D(\mathbf{r}) &= \int_{t_n}^{t_f} T(t)\sigma(\mathbf{r}(t))L_{n=1}(\mathbf{n}(t))dt \\
 S(\mathbf{r}) &= \int_{t_n}^{t_f} T(t)\sigma(\mathbf{r}(t))\sum_i^N \omega_i L_{n_i}(\mathbf{n}(t), \mathbf{d})dt \\
 F(\mathbf{r}) &= \int_{t_n}^{t_f} T(t)\sigma(\mathbf{r}(t))f(\mathbf{r}(t))dt
 \end{aligned} \tag{2}$$

where  $\mathbf{n}(t)$  is the normal direction estimated via  $\nabla\sigma(\mathbf{r}(t))$ ,  $L_{n=1}(\mathbf{n}(t))$  is the diffuse light map indexed by the normal direction  $\mathbf{n}(t)$ , and  $L_{n_i}(\mathbf{n}(t), \mathbf{d})$  is the specular component  $n_i$  indexed by the inbound reflection direction depending on the local normal and viewing direction  $\mathbf{d}$ . Finally,  $\alpha, \sigma, \omega$ , and a per-location feature  $f$  are the network outputs conditioned on the sampled latent code  $z$ . We restrict the albedo to be view and lighting independent and encourage multi-view consistency. Note that in addition to rendering components such as albedo, diffuse and specular components, we let our network accumulate additional features  $F(\mathbf{r})$ , so that it can capture high frequency details and material properties in an unsupervised fashion.

**3.2.3 Volumetric Model Network Architecture.** We extend the architecture from concurrent work proposed by [Gu et al. 2021] for our neural implicit field. Rather than explicitly use the low resolution albedo  $A(\mathbf{r})$  following Eq. 2, our network produces a feature vector  $f(\mathbf{r}(t)) \in R^{256}$  via 6 fully-connected layers from the positional encoding on the 3D coordinates. A linear-layer is attached

to the output of the 4-th layer to produce the volume density, and an additional two-layer MLP is attached to 6-th layer to produce the albedo and reflectance properties. Diffuse component  $D$  and Specular Component  $S$  are estimated following Eq. 2, where the blending weights  $\omega_i$  are estimated by the network.

**3.2.4 Neural Rendering Network.** To reduce the memory consumption and computation cost, we render albedo, diffuse, and specular shading in low resolution and upsample them to high resolution for relighting. The specific low and high resolutions depend on the dataset used and details can be found in the Section 4. To generate the high resolution albedo, we upsample the feature map  $F(\mathbf{r})$  and enforce its first 3 channels to correspond to the albedo image, similar to some other works in the literature [Meka et al. 2020; Thies et al. 2019]. Each upsampling unit consists of two  $1 \times 1$  convolutions modulated by the latent code  $z$ , a pixelshuffle upsampler [Shi et al. 2016] and a BlurPool [Zhang 2019] with stride 1. The low resolution albedo  $A(\mathbf{r})$  is still used to enforce consistency with the upsampled high resolution albedo (see Section 3.3). For shading maps, we directly apply bilinear upsampling.

Finally, a relighting network takes as input the albedo map  $A$ , the diffuse map  $D$ , the specular component map  $S$  and the features  $F$  and generates the final  $I_{relit}$  image. The architecture of Relighting Network is a shallow U-Net [Ronneberger et al. 2015].

### 3.3 Supervised Adversarial Training

Here we introduce the scheme to train our pipeline from a collection of unconstrained in-the-wild images. While it is possible to train the full pipeline with a single adversarial loss on the relit image, we found empirically that adding additional supervision on intermediate outputs significantly improves the training convergence and rendering quality.

*Pseudo Ground Truth Generation.* Large scale in-the-wild images provides great data diversity, which is critical for training a generator. However, the groundtruth labels for geometry and shading are usually missing. In our case, we are particularly interested in having “real examples” of the albedo and geometry to supervise our method. To this end, we resort to the state-of-the-art image based relighting algorithm, Total Relighting [Pandey et al. 2021], to produce pseudo ground truth albedo and normals and to also further increase the data diversity. Specifically, for each image in our training set, we randomly select an HDRI map from a collection of 400 maps sourced from public repository [Zaal et al. 2020], apply a random rotation, and run Total Relighting to generate the albedo, surface normal and a relit image with the associated light maps (diffuse and specular components). Example augmented data can be found in supplementary materials.

*Albedo Adversarial Loss  $\mathcal{L}_A$ :*  $D_A(A(\mathbf{r})) + D_A(A_{hi-res})$ . We supervise the output albedo images in both low and high resolution with adversarial loss using the pseudo ground truth generated with [Pandey et al. 2021]. A standard non-saturating logistic GAN loss with R1 penalty is applied to train the generator and discriminator. The discriminator architecture  $D_*$  for all the losses follows the one proposed in [Karras et al. 2020].

*Geometry Adversarial Loss  $\mathcal{L}_G$ :*  $D_G(\nabla\sigma(\mathbf{r}(t)))$ . We also supervise the geometry as it is crucial for multi-view consistent rendering and relighting realism. While the density  $\sigma$  is the immediate output from the network that measures the geometry, we find it is more convenient to supervise the surface normals computed via  $\nabla\sigma(\mathbf{r}(t))$ . Therefore, we add an adversarial loss between the volumetric rendered normal from the derivative of the density and the pseudo ground truth normal obtained from [Pandey et al. 2021].

*Shading Adversarial Loss  $\mathcal{L}_S$ :*  $D_S(D(\mathbf{r}), S(\mathbf{r}), I_{relit})$ . Since we are resorting to a generative model, directly supervising the albedo and relit pair with a reconstruction loss is not possible. Indeed the network produces new identity from a randomly sampled latent code where direct supervision is not available. Therefore, to enforce the Relight Network faithfully integrating shading with albedo, we apply a conditional adversarial loss on the relit image. This is achieved by adding a discriminator  $D_S$  that takes the concatenation of the relit image  $I_{relit}$ , diffuse map  $D(\mathbf{r})$  and specular map  $S(\mathbf{r})$  as the inputs and discriminate if the group is fake, i.e. from our model, or true, i.e. from [Pandey et al. 2021]. The training gradients are only allowed to back-propagate to the relit image but not the other inputs (i.e. set to zero) as they are the data to be conditioned on.

*Photorealistic Adversarial Loss  $\mathcal{L}_P$ :*  $D_P(I_{relit})$ . A downside of the Shading Adversarial Loss is that the model performance is upper-bounded by the specific algorithm used to generate pseudo-groundtruth labels, in our case [Pandey et al. 2021]. As a result, inaccuracies in the relighting examples, e.g. overly smoothed shading and lack of specular highlights, may affect our rendering quality. To enhance the photorealism, we add an additional adversarial loss directly on the generated relit images with the original images from the dataset.

*Path Loss  $\mathcal{L}_{path}$ :*  $\ell_1(A(\mathbf{r}), A_{hi-res})$ . Following StyleNeRF [Gu et al. 2021], we add a loss to ensure the consistency between the albedo maps in low and high resolutions. Specifically, we downsample the high resolution to the low resolution, and add a per-pixel  $\ell_1$  loss.

The final loss function is a weighted sum of all above mentioned terms:  $\mathcal{L} = \lambda_1\mathcal{L}_A + \lambda_2\mathcal{L}_G + \lambda_3\mathcal{L}_S + \lambda_4\mathcal{L}_P + \lambda_5\mathcal{L}_{path}$ , where for our experiments we empirically determined these weights to be 1.0, 0.5, 0.25, 0.75, 0.5. Training details can be found in supplementary materials.

## 4 EXPERIMENTS

In this section, we compare our rendering quality and relighting performance with state-of-the-art methods. We also provide ablation study showing the contribution of major system design choice to the final performance.

*Datasets.* We train our model on CelebA dataset [Liu et al. 2015] which is widely used for such comparisons, and on the FFHQ [Karras et al. 2019] where a comparison of high resolution results can be made. On CelebA, our model produces volumetric renderings at  $64 \times 64$  and final outputs at  $128 \times 128$ . On FFHQ, the volumetric renderings and final resolution are  $64 \times 64$  and  $256 \times 256$  respectively.

*Baseline Methods.* We show qualitative and quantitative comparison with ShadeGAN [Pan et al. 2021] since, to the best of our



Figure 3: Our synthesized images under rotating camera or rotating lighting. Note the relighting consistency and view-dependent effects.

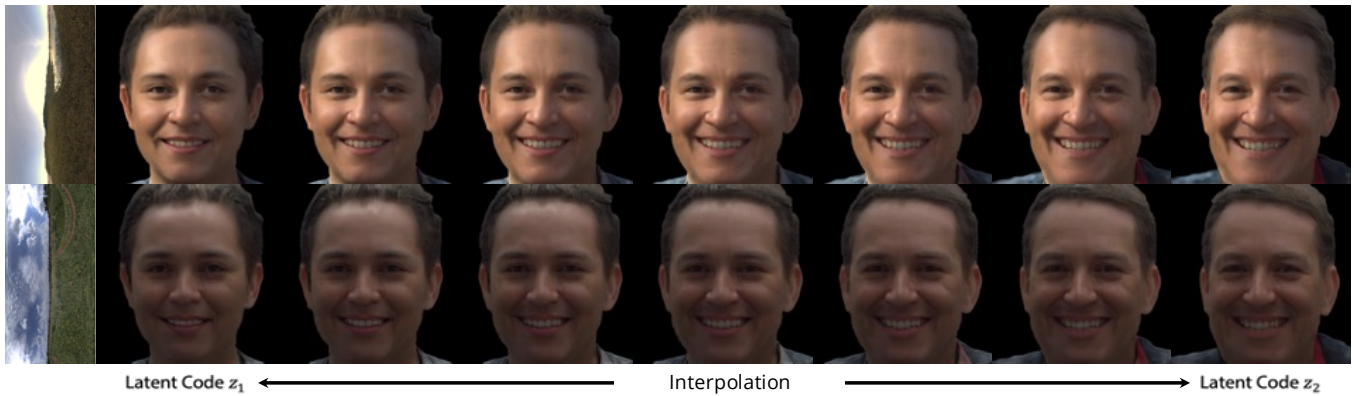


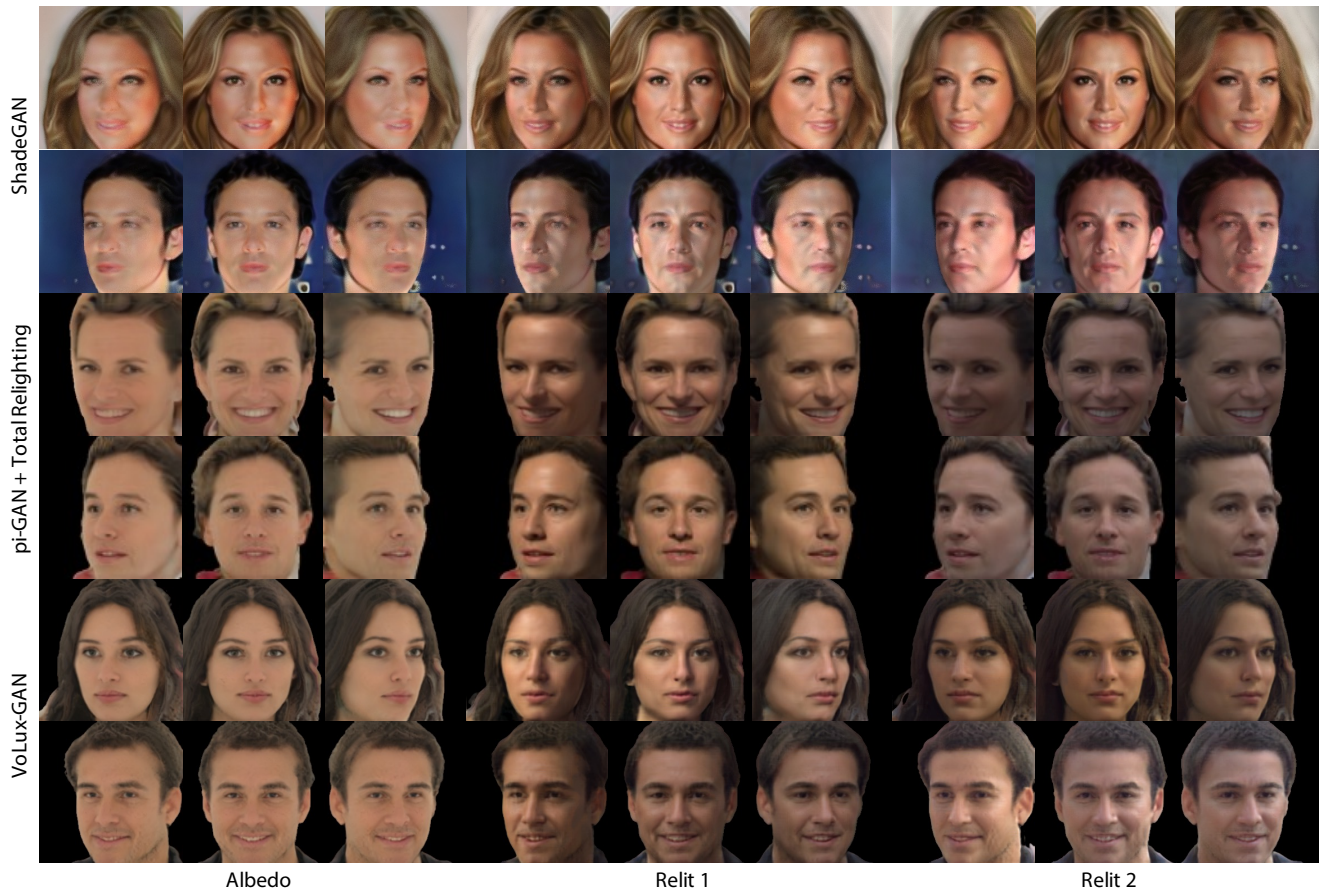
Figure 4: Our result with interpolated latent code under different HDR environment maps. Note the relighting effects are correctly transferred across identities.

Table 1: Identity consistency across camera poses around the yaw axis. The scores indicate the similarity calculated as the dot product between normalized embeddings from a state-of-the-art face recognition network [Deng et al. 2019] (higher is better). While our method performs comparably or marginally better at small view changes, we significantly outperform the state-of-the-art at more extreme viewpoints.

Method	Relit image identity similarity $\uparrow$					Albedo image identity similarity $\uparrow$				
	-0.5 rad	-0.25 rad	0 rad	0.25 rad	0.5 rad	-0.5 rad	-0.25 rad	0 rad	0.25 rad	0.5 rad
ShadeGAN[Pan et al. 2021]	0.4814	0.7513	-	0.7628	0.4997	0.4818	0.7582	-	0.7702	0.5091
pi-GAN [Chan et al. 2021] + TR[Pandey et al. 2021]	0.5215	0.7472	-	0.7419	0.4981	0.5135	0.7378	-	0.7438	0.4898
VoLux-GAN+Surface Relighting	0.4471	0.7065	-	0.7388	0.4959	0.5531	0.7611	-	0.7796	0.5585
VoLux-GAN - $\mathcal{L}_P$	0.4827	0.6487	-	0.6721	0.5156	0.5467	0.7809	-	0.8151	0.5901
VoLux-GAN - $\mathcal{L}_S$	0.4071	0.6996	-	0.7788	0.4718	0.4886	0.7292	-	0.8095	0.5581
VoLux-GAN - $\mathcal{L}_G$	0.4776	0.7182	-	0.7564	0.5015	0.4652	0.7278	-	<b>0.8099</b>	0.5398
<b>VoLux-GAN</b>	<b>0.6064</b>	<b>0.7736</b>	-	<b>0.7997</b>	<b>0.5985</b>	<b>0.6389</b>	<b>0.7919</b>	-	0.7863	<b>0.6162</b>

knowledge, it is the only 3D-aware generator that supports relighting. In addition, we consider an alternative strong baseline where we use pi-GAN [Chan et al. 2021] to render multi-view images and then run a single image based portrait relighting method [Pandey et al. 2021] for HDR relighting.

*Metrics.* Many evaluation metrics relying on perception features have been proposed to measure the rendering quality [Bińkowski et al. 2018; Heusel et al. 2017]. While these metrics indeed measure the similarity between two collections of images, they are very sensitive to implementation details such as training image resolution, image post-processing e.g. cropping, or the choice of training



**Figure 5: Qualitative comparisons on CelebA with ShadeGAN [Pan et al. 2021] and pi-GAN [Chan et al. 2021] + portrait relighting [Pandey et al. 2021]: note how our method produces more consistent albedo and relighting results across multiple views.**

dataset, as has been shown in literature [Parmar et al. 2021]. As a result, these metrics are not suitable to evaluate our model since our pipeline is not trained directly on publicly available datasets but on a specifically tailored augmented dataset for good rendering and relighting performance.

Instead, we evaluate our pipeline and other methods by measuring the perceptual impact of view-synthesis and relighting on a subject’s identity. Specifically, we use a similarity metric based on the embedding space of a state-of-the-art face recognition network [Deng et al. 2019]. This stability metric indicates how well the subject identity is preserved when we synthesize novel views and novel light renderings for a synthesized face.

#### 4.1 Relightable Face Generation

In this section we demonstrate the capabilities of our framework. In Fig. 3, we show one subject randomly sampled from latent space  $z \sim N(0, I)$  trained on the FFHQ dataset. The first row shows the faces rendered at different camera poses. Our network successfully renders consistent faces even under a large yaw angle (e.g. 45°) thanks to better geometry supervised by the geometry adversarial

loss. The second row shows the same subject rendered under a rotating HDR map. Note how the specularities and shading on the face respond correctly to the HDR environment maps.

Our latent space also supports interpolation. In Fig. 4, we linearly interpolate between two randomly sampled latent codes, and show relit images of each subject under two HDR lighting conditions. As seen, the appearance of the subject transitions smoothly and the intermediate identities are successfully relit. Note also the consistent relighting, where view dependent effects and specularities are successfully transferred between different latent codes.

#### 4.2 Comparisons with State-of-the-Art

We compare to ShadeGAN [Pan et al. 2021], pi-GAN [Chan et al. 2021] coupled to an image based relighting [Pandey et al. 2021], and show the qualitative results in Fig. 5 (Due to the page limit, the target HDRI images are visualized in the supplementary materials). For a fair comparison in terms of image resolution, we trained our model on CelebA, like ShadeGAN and pi-GAN. In each row, we show the albedo map and color images rendered under two different lightings from three camera viewpoints.

**Table 2: Relighting consistency. Our method better preserve the original identity (albedo) under different illuminations. Note that ShadeGAN is evaluated over three different positional light sources instead of HDR maps.**

Method	HDR map 1	HDR map 2	HDR map 3
ShadeGAN[Pan et al. 2021]	0.5806	0.6486	0.6559
pi-GAN[Chan et al. 2021] + TR[Pandey et al. 2021]	0.7014	0.8796	0.7677
VoLux-GAN+Surface Relighting	0.5510	0.5548	0.5349
VoLux-GAN - $\mathcal{L}_P$	0.4132	0.4471	0.4712
VoLux-GAN - $\mathcal{L}_S$	<b>0.8917</b>	0.8846	0.7387
VoLux-GAN - $\mathcal{L}_G$	0.5331	0.5864	0.6158
<b>VoLux-GAN</b>	0.7600	<b>0.8900</b>	<b>0.8082</b>

Our method produces significantly better albedo than ShadeGAN thanks to our supervised albedo adversarial loss. Moreover, our results contain more high frequency specular components and can respond to more diverse global illumination. We also coupled pi-GAN with Pandey et al. [2021] as a baseline. Compared to ShadeGAN, pi-GAN generates faces that exhibit overall diffuse illumination, which works better in conjunction with Pandey et al. [2021]. Although the coupled model can produce plausible relighting results, they are often inconsistent across views (e.g. in the 3rd row middle, the cheek is dark in one view but bright in another). In contrast, our results show more consistency across views and lights thanks to the proposed volumetric relighting formulation, which is also reflected by the quantitative metric below.

**4.2.1 Quantitative Results.** We evaluate our method with quantitative metrics. The goal of the following experiments is to demonstrate that our method is able to synthesize images that are consistent across views in terms of geometry and relighting.

**Geometry Consistency.** To demonstrate the geometry consistency, we render a fixed random latent code to multiple views. We then compute the similarity score [Deng et al. 2019] of yaw-posed renderings with the frontal facing rendering, and average it over 100 such randomly sampled latent codes. We compute the score for both relit images and the high-res albedo images. We also show the score computed with the same scheme for ShadeGAN [Pan et al. 2021] and the baseline of pi-GAN [Chan et al. 2021] + portrait relighting [Pandey et al. 2021]. The results are showed in Table 1. Note how our method consistently outperforms the other state-of-art approaches, demonstrating better multi-view consistency for each identity on generated albedo and relit images.

**Relighting Consistency.** Similarly, we also evaluate the stability of our relighting. Following the geometry consistency experiment, we use the embedding space of a face recognition network [Deng et al. 2019] to generate the identity similarity score between the albedo and 3 renderings under different environment maps (as shown in Fig. 4). We report an average score over 100 randomly sampled latent codes. A stable relighting method should give a high similarity score, since relighting does not change the identity.

The results are reported in Table 2, showing that our approach is also able to generate relit images that are more consistent with the original albedo identity. At the same time our relit images look more photorealistic as shown in Figure 5, where we better capture higher order light transport effects.

**Ablation Study.** We report an ablation study of the loss functions used for training our pipeline and Table 2. As demonstrated by these quantitative results, the full framework and the proposed supervised adversarial losses are all contributing to the final rendering quality. In particular, we show that removing the shading loss  $\mathcal{L}_S$ , or geometry loss  $\mathcal{L}_G$  or photorealistic  $\mathcal{L}_P$  all lead to lower metrics. Similarly, when we first accumulate the surface geometry and then perform image based relighting, the results are unsatisfactory (see comparison with VoLux-GAN+Surface Relighting in Table 1 and Table 2).

## 5 DISCUSSION

We proposed a generative model of face images, that internally leverages a volumetric representation to facilitate multi-view generation and full HDRI relighting. Of particular note is that we have shown how to efficiently perform the aggregation of albedo, specular and diffuse components that helps to preserve the identity. Furthermore, a proposed supervised adversarial framework guides the network to generate the right intrinsic properties of faces. Our results prove the effectiveness of the approach for synthesizing novel identities. Future work could explore the use of semantic information to allow for expression control similar to StyleGAN [Karras et al. 2019].

## REFERENCES

- Hassan Abu Alhaja, Siva Karthik Mustikovela, Andreas Geiger, and Carsten Rother. 2018. Geometric Image Synthesis. *ACCV* (2018).
- Jonathan T. Barron and Jitendra Malik. 2015. Shape, Illumination, and Reflectance From Shading. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2015).
- Mikołaj Bińkowski, Danica J Sutherland, Michael Arbel, and Arthur Gretton. 2018. Demystifying mmd gans. *arXiv preprint arXiv:1801.01401* (2018).
- Mark Boss, Raphael Braun, Varun Jampani, Jonathan T. Barron, Ce Liu, and Hendrik P.A. Lensch. 2020. NeRD: Neural Reflectance Decomposition from Image Collections. *ICCV* (2020).
- Mark Boss, Varun Jampani, Raphael Braun, Ce Liu, Jonathan T. Barron, and Hendrik P.A. Lensch. 2021. Neural-PIL: Neural Pre-Integrated Lighting for Reflectance Decomposition. In *NeurIPS*.
- Andrew Brock, Jeff Donahue, and Karen Simonyan. 2019. Large Scale GAN Training for High Fidelity Natural Image Synthesis. In *ICLR*.
- Marcel C. Buehler, Abhimitra Meka, Gengyan Li, Thabo Beeler, and Otmar Hilliges. 2021. VariTex: Variational Neural Face Textures. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*.
- Eric R Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. 2021. 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*. 5799–5809.
- Xuelin Chen, Daniel Cohen-Or, Baoquan Chen, and Niloy J. Mitra. 2021. Towards a Neural Graphics Pipeline for Controllable Image Generation. *Computer Graphics Forum* 40, 2 (2021).
- George Chogovadze, Rémi Pautrat, and Marc Pollefeys. 2021. Controllable Data Augmentation Through Deep Relighting.
- Yunjeong Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. 2020. StarGAN v2: Diverse Image Synthesis for Multiple Domains.
- Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. 2019. Arcface: Additive angular margin loss for deep face recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4690–4699.
- Ishan Durugkar, Ian Gemp, and Sridhar Mahadevan. 2017. Generative Multi-Adversarial Networks.
- Matheus Gadelha, Subhransu Maji, and Rui Wang. 2017. 3d shape induction from 2d views of multiple objects. In *2017 International Conference on 3D Vision (3DV)*. IEEE, 402–411.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative Adversarial Nets. In *Advances in Neural Information Processing Systems*.
- Ned Greene. 1986. Environment mapping and other applications of world projections. *IEEE Computer Graphics and Applications* 6, 11 (1986), 21–29.
- Jiatao Gu, Lingjie Liu, Peng Wang, and Christian Theobalt. 2021. StyleNeRF: A Style-based 3D-Aware Generator for High-resolution Image Synthesis. *arXiv preprint*



- arXiv:2110.08985* (2021).
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. 2017. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems* 30 (2017).
- Yoshihiro Kanamori and Yuki Endo. 2018. Relighting Humans: Occlusion-Aware Inverse Rendering for Full-Body Human Images. *ACM Transactions Graphics (Proc. SIGGRAPH Asia)* (2018).
- Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. 2021. Alias-Free Generative Adversarial Networks. In *NeurIPS*.
- Tero Karras, Samuli Laine, and Timo Aila. 2019. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4401–4410.
- Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. 2020. Analyzing and Improving the Image Quality of StyleGAN. In *Proc. CVPR*.
- Yiyi Liao, Katja Schwarz, Lars Mescheder, and Andreas Geiger. 2020. Towards Unsupervised Learning of Generative Models for 3D Controllable Image Synthesis. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*.
- Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. 2015. Deep Learning Face Attributes in the Wild. In *Proceedings of International Conference on Computer Vision (ICCV)*.
- Abhimitra Meka, Gereon Fox, Michael Zollhöfer, Christian Richardt, and Christian Theobalt. 2017. Live User-Guided Intrinsic Video for Static Scene. *IEEE Transactions on Visualization and Computer Graphics* (2017).
- Abhimitra Meka, Maxim Maximov, Michael Zollhoefer, Avishek Chatterjee, Hans-Peter Seidel, Christian Richardt, and Christian Theobalt. 2018. LIME: Live Intrinsic Material Estimation. In *Proc. Computer Vision and Pattern Recognition*.
- Abhimitra Meka, Rohit Pandey, Christian Häne, Sergio Orts-Escolano, Peter Barnum, Philip David-Son, Daniel Erickson, Yinda Zhang, Jonathan Taylor, Sofien Bouaziz, Chloe LeGendre, Wan-Chun Ma, Ryan Overbeck, Thabo Beeler, Paul Debevec, Shahram Izadi, Christian Theobalt, Christoph Rhemann, and Sean Fanello. 2020. Deep Relightable Textures: Volumetric Performance Capture with Neural Rendering. *ACM Transactions on Graphics* (2020).
- Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. 2020. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European conference on computer vision*. Springer, 405–421.
- Gene S Miller and CR Hoffman. 1984. Illumination and reflection maps. In *ACM SIGGRAPH*.
- Gonçalo Mordido, Haojin Yang, and Christoph Meinel. 2018. Dropout-GAN: Learning from a Dynamic Ensemble of Discriminators.
- Thomas Nestmeyer, Jean-François Lalonde, Iain Matthews, and Andreas M. Lehrmann. 2020. Learning Physics-guided Face Relighting under Directional Light. In *CVPR*.
- Thu Nguyen-Phuoc, Chuan Li, Lucas Theis, Christian Richardt, and Yong-Liang Yang. 2019. HoloGAN: Unsupervised Learning of 3D Representations From Natural Images. In *The IEEE International Conference on Computer Vision (ICCV)*.
- Michael Niemeyer and Andreas Geiger. 2021. Giraffe: Representing scenes as compositional generative neural feature fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 11453–11464.
- Xingang Pan, Xudong Xu, Chen Change Loy, Christian Theobalt, and Bo Dai. 2021. A Shading-Guided Generative Implicit Model for Shape-Accurate 3D-Aware Image Synthesis. In *NeurIPS*.
- Rohit Pandey, Sergio Orts Escolano, Chloe Legendre, Christian Haene, Sofien Bouaziz, Christoph Rhemann, Paul Debevec, and Sean Fanello. 2021. Total relighting: learning to relight portraits for background replacement. *ACM Transactions on Graphics (TOG)* 40, 4 (2021), 1–21.
- Gaurav Parmar, Richard Zhang, and Jun-Yan Zhu. 2021. On Buggy Resizing Libraries and Surprising Subtleties in FID Calculation. *arXiv preprint arXiv:2104.11222* (2021).
- Bui Tuong Phong. 1975. Illumination for computer generated pictures. *Commun. ACM* 18, 6 (1975), 311–317.
- Thu Nguyen Phuoc, Christian Richardt, Long Mai, Yongliang Yang, and Niloy J Mitra. 2020. BlockGAN: Learning 3D Object-aware Scene Representations from Unlabelled Images. In *NeurIPS 2020: Conference on Neural Information Processing Systems*.
- Ravi Ramamoorthi and Pat Hanrahan. 2001. An efficient representation for irradiance environment maps. In *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*. 497–500.
- Peiran Ren, Yue Dong, Stephen Lin, Xin Tong, and Baining Guo. 2015. Image Based Relighting Using Neural Networks. *ACM Transactions on Graphics* (2015).
- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. 2015. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*. Springer, 234–241.
- Shunsuke Saito, Tomas Simon, Jason Saragih, and Hanbyul Joo. 2020. PIFuHD: Multi-Level Pixel-Aligned Implicit Function for High-Resolution 3D Human Digitization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. 2020. GRAF: Generative Radiance Fields for 3D-Aware Image Synthesis.
- Wenzhe Shi, Jose Caballero, Ferenc Huszár, Johannes Totz, Andrew P Aitken, Rob Bishop, Daniel Rueckert, and Zehan Wang. 2016. Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1874–1883.
- Attila Szabó, Givi Meishvili, and Paolo Favaro. 2019. Unsupervised generative 3d shape learning from natural images. *arXiv preprint arXiv:1910.00287* (2019).
- Daichi Tajima, Yoshihiro Kanamori, and Yuki Endo. 2021. Relighting Humans in the Wild: Monocular Full-Body Human Relighting with Domain Adaptation.
- Feitong Tan, Danhang Tang, Dou Mingsong, Guo Kaiwen, Rohit Pandey, Cem Keskin, Ruofei Du, Deqing Sun, Sofien Bouaziz, Sean Fanello, Ping Tan, and Yinda Zhang. 2021. HumanGPS: Geodesic PreServing Feature for Dense Human Correspondences. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Ayush Tewari, Ohad Fried, Justus Thies, Vincent Sitzmann, Stephen Lombardi, Kalyan Sunkavalli, Ricardo Martin-Brualla, Tomas Simon, Jason Saragih, Matthias Nießner, Rohit Pandey, Sean Fanello, Gordon Wetzstein, Jun-Yan Zhu, Christian Theobalt, Maneesh Agrawala, Eli Shechtman, Dan B Goldman, and Michael Zollhoefer. 2020. State of the Art on Neural Rendering. In *Eurographics*.
- Justus Thies, Michael Zollhöfer, and Matthias Niessner. 2019. Deferred Neural Rendering: Image Synthesis Using Neural Textures. *SIGGRAPH and ACM TOG* (2019).
- Zhibo Wang, Xin Yu, Ming Lu, Quan Wang, Chen Qian, and Feng Xu. 2020. Single Image Portrait Relighting via Explicit Multiple Reflectance Channel Modeling. *ACM SIGGRAPH Asia and Transactions on Graphics* (2020).
- Erroll Wood, Tadas Baltrušaitis, Charlie Hewitt, Sebastian Dziadzio, Matthew Johnson, Virginia Estellers, Thomas J. Cashman, and Jamie Shotton. 2021. Fake It Till You Make It: Face analysis in the wild using synthetic data alone.
- Zexiang Xu, Kalyan Sunkavalli, Sunil Hadap, and Ravi Ramamoorthi. 2018. Deep Image-based Relighting from Optimal Sparse Samples. *ACM Transactions on Graphics* (2018).
- Greg Zaal, Sergej Majboroda, and Andreas Mischok. 2020. HDRI Haven. <https://www.hdr Haven.com/>. Accessed: 2021-11-13.
- Han Zhang, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. 2019. Self-Attention Generative Adversarial Networks. In *International Conference on Machine Learning*.
- Jason Y. Zhang, Gengshan Yang, Shubham Tulsiani, and Deva Ramanan. 2021b. NeRS: Neural Reflectance Surfaces for Sparse-view 3D Reconstruction in the Wild. In *Conference on Neural Information Processing Systems*.
- Richard Zhang. 2019. Making convolutional networks shift-invariant again. In *International conference on machine learning*. PMLR, 7324–7334.
- Xiuming Zhang, Pratul P Srinivasan, Boyang Deng, Paul Debevec, William T Freeman, and Jonathan T Barron. 2021a. NeRFactor: Neural Factorization of Shape and Reflectance Under an Unknown Illumination. *arXiv preprint arXiv:2106.01970* (2021).
- Peng Zhou, Lingxi Xie, Bingbing Ni, and Qi Tian. 2021. CIPS-3D: A 3D-Aware Generator of GANs Based on Conditionally-Independent Pixel Synthesis. *arXiv preprint arXiv:2110.09788* (2021).
- Jun-Yan Zhu, Zhoutong Zhang, Chengkai Zhang, Jiajun Wu, Antonio Torralba, Josh Tenenbaum, and Bill Freeman. 2018. Visual Object Networks: Image Generation with Disentangled 3D Representations. In *Advances in Neural Information Processing Systems*.
- Tyler Zhu, Per Karlsson, and Christoph Bregler. 2020. SimPose: Effectively Learning DensePose and Surface Normals of People from Simulated Data. In *European Conference on Computer Vision*. Springer, 225–242.