# Mimic3D: Thriving 3D-Aware GANs via 3D-to-2D Imitation

Xingyu Chen*     Yu Deng*     Baoyuan Wang

Xiaobing.AI

## Abstract

*Generating images with both photorealism and multi-view 3D consistency is crucial for 3D-aware GANs, yet existing methods struggle to achieve them simultaneously. Improving the photorealism via CNN-based 2D super-resolution can break the strict 3D consistency, while keeping the 3D consistency by learning high-resolution 3D representations for direct rendering often compromises image quality. In this paper, we propose a novel learning strategy, namely 3D-to-2D imitation, which enables a 3D-aware GAN to generate high-quality images while maintaining their strict 3D consistency, by letting the images synthesized by the generator's 3D rendering branch mimic those generated by its 2D super-resolution branch. We also introduce 3D-aware convolutions into the generator for better 3D representation learning, which further improves the image generation quality. With the above strategies, our method reaches FID scores of 5.4 and 4.3 on FFHQ and AFHQ-v2 Cats, respectively, at 512×512 resolution, largely outperforming existing 3D-aware GANs using direct 3D rendering and coming very close to the previous state-of-the-art method that leverages 2D super-resolution. Project website: https://seanchenxy.github.io/Mimic3DWeb.*

## 1. Introduction

3D-aware GANs [37, 9, 6, 3] have experienced rapid development in recent years and shown great potential for large-scale realistic 3D content creation. The core of 3D-aware GANs is to incorporate 3D representation learning and differentiable rendering into image-level adversarial learning [8]. In this way, the generated 3D representations are forced to mimic real image distribution from arbitrary viewing angles, resulting in their faithful reconstruction of the underlying 3D structures of the subjects for free-view image synthesis. Among different 3D representations, neural radiance field (NeRF) [24] has been proven to be effective in the 3D-aware GAN scenario [37, 4], which guarantees strong 3D consistency when synthesizing multiview
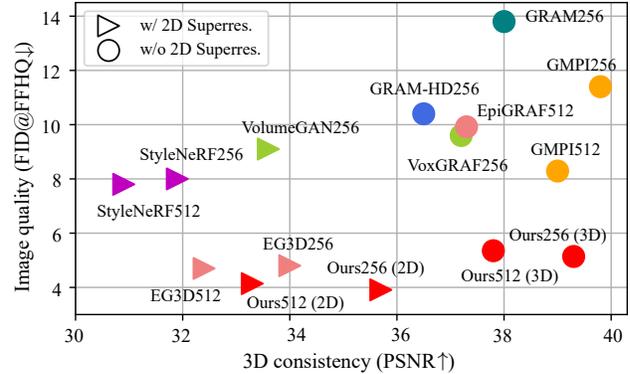


Figure 1. Comparison between different 3D-aware GANs on image generation quality and multiview 3D consistency. The image generation quality is evaluated via FID between generated and real images. The 3D consistency is measured by conducting 3D reconstruction [45] on generated multiview images and calculating PSNR between them and the re-rendered reconstruction results. Our method inherits the high image quality of approaches leveraging 2D super-resolution meanwhile maintains strict 3D consistency by taking the advantage of direct 3D rendering.

images via volume rendering [15].

However, NeRF's volumetric representation also brings high computation costs to GAN training. This hinders the generative models from synthesizing high-resolution images with fine details. Several attempts have been made to facilitate NeRF-based GAN training at high resolution, via sparse representations [38, 6, 47, 54] or patch-wise adversarial learning [43], yet the performance is still unsatisfactory and lags far behind state-of-the-art 2D GANs [19, 17].

Along another line, instead of using direct NeRF rendering, plenty of works [26, 9, 28, 3, 50] introduce 2D super-resolution module to deal with 3D-aware GAN training at high resolution. A typical procedure is to first render a NeRF-like feature field into low-resolution feature maps, then apply a 2D CNN to generate high-resolution images from them. The representative work among this line, namely EG3D [3], utilizes tri-plane representation to effectively model the low-resolution feature field and leverages StyleGAN2-like [19] super-resolution block to achieve image synthesis at high-quality. It sets a record for image

arXiv:2303.09036v2 [cs.CV] 7 Aug 2023

Figure 2. Our method enables high-quality image generation at $512 \times 512$ resolution without using a 2D super-resolution module.

quality among 3D-aware GANs and gets very close to that of state-of-the-art 2D GANs. However, a fatal drawback of this line of works is a sacrifice of strict 3D consistency, due to leveraging a black-box 2D CNN for image synthesis.

A question naturally arises —— *Is there any way to combine the above two lines to achieve strict 3D consistency and high-quality image generation simultaneously?* The answer, as we will show in this paper, is arguably yes. The key intuition is to let the images synthesized by direct NeRF rendering to mimic those generated by a 2D super-resolution module, which we name *3D-to-2D imitation*.

Specifically, we start from an EG3D backbone that adopts 2D super-resolution to generate high-resolution images from a low-resolution feature field. Based on this architecture, we add another 3D super-resolution module to generate high-resolution NeRF from the low-resolution feature field and force the images rendered by the former to imitate those generated by the 2D super-resolution branch. This process can be seen as a multiview reconstruction process —— images sharing the same latent code from different views produced by the 2D branch are pseudo multiview data, and the high-resolution NeRF branch represents the 3D scene to be reconstructed. Previous methods [29, 53, 33] have shown that this procedure can obtain reasonable 3D reconstruction, even if the multiview data are not strictly 3D consistent. We believe this is partially due to the inductive bias (*e.g.*, continuity and sparsity) of the underlying 3D representation. With the above process, the high-resolution NeRF learns to depict fine details of the 2D-branch images, thus enabling high-quality image rendering. The 3D consistency across different views can also be preserved thanks to the intrinsic property of NeRF. Note that if the rendered images try to faithfully reconstruct every detail of the 2D-branch images across different views, it is likely to obtain blurry results due to detail-level 3D inconsistency of the latter. To avoid this problem, we only let the images produced by the two branches be perceptually similar (*i.e.* by LPIPS loss [51]), and further enforce adversarial loss between the rendered images from the high-resolution NeRF and real images to maintain high-frequency details. In addition, we only render small image patches to conduct the imitative

learning to reduce memory costs.

Apart from the above learning strategy, we introduce 3D-aware convolutions to the EG3D backbone to improve tri-plane learning, motivated by a recent 3D diffusion model [46]. The original EG3D generates tri-plane features to model the low-resolution feature field via a StyleGAN2-like generator. The generator is forced to learn 2D-unaligned features on the three orthogonal planes via 2D convolutions, which is inefficient. The 3D-aware convolution considers associated features in 3D space when performing 2D convolution, which improves feature communications and helps to produce more reasonable tri-planes. Nevertheless, directly applying 3D-aware convolution in all layers in the generator is unaffordable. As a result, we only apply them after the output layers at each resolution in the tri-plane generator. This helps us to further improve the image generation quality with only a minor increase in the total memory consumption.

With the above strategies, our generator is able to synthesize 3D-consistent images of virtual subjects with high image quality (Fig. 2). It reaches FID scores [13] of 5.4 and 4.3 on FFHQ [18] and AFHQ-v2 Cats [5], respectively, at $512 \times 512$ resolution, largely outperforming previous 3D-aware GANs with direct 3D rendering and even surpassing many leveraging 2D super-resolution (Fig. 1). A by-product of our method is a more powerful 2D-branch generator, which reaches an FID of 4.1 on FFHQ, exceeding previous state-of-the-art EG3D. Though our method presented in this paper is mostly based on EG3D backbone, its 3D-to-2D imitation strategy can be extended to learning other 3D-aware GANs as well. We believe this would largely close the quality gap between 3D-aware GANs and traditional 2D GANs, and pave a new way for realistic 3D generation.

## 2. Related Works

**3D-aware GAN.** 3D-aware GANs [12, 25, 37, 4, 26, 9, 52, 6, 3, 43, 54] aim to generate multiview images of an object category, given only in-the-wild 2D images as training data. The key is to represent the generated scenes via a 3D representation and leverage corresponding rendering

techniques to synthesize images at different viewpoints for image-level adversarial learning [8]. Initially, explicit representations such as voxels [25, 12] and meshes [44] are used to describe scenes. With the development of neural implicit fields [32, 23, 42, 41, 24, 45, 27], implicit scene representations, especially NeRF [24], gradually overtake explicit ones in 3D-aware GANs [4, 28, 3]. Nevertheless, one great hurdle of NeRF-based GANs is the high computation cost, which restricts earlier works [37, 4, 7, 48, 31] from synthesizing high-quality images. Consequently, a large number of follow-up works [26, 9, 55, 50, 28, 3, 49] avoid rendering NeRF at high resolution by conducting 2D super-resolution from a low-resolution image or feature map rendered by NeRF-like fields. This is only a stopgap as the black-box 2D super-resolution module sacrifices the important 3D consistency brought by NeRF. To keep the strict 3D consistency, several works [6, 47, 38, 43, 54] turn to more sparse 3D representations such as sparse voxel [38], radiance manifolds [6], and multi-plane images [54] to allow direct rendering at high resolution. Carefully designed training strategies such as two-stage training [47] or patch-wise optimization [43] are also introduced to facilitate the learning process. However, their image generation quality still lags behind those with 2D super-resolution. Our method combines the advantages of both lines of works to achieve high-quality image generation and strict 3D consistency at once, by leveraging the proposed 3D-to-2D imitation.

**3D generation by 3D-to-2D imitation.** Recent studies [14, 39, 10] reveal that 2D generative models [2, 18] have the ability to generate pseudo multiview images of a subject. Based on this observation, several methods [29, 53, 40, 30] propose to distill the knowledge from a pre-trained 2D generative model for 3D generation by performing 3D reconstruction on the generated "multiview" images. A standard procedure is to render the 3D representation of an object from multiple views, and compare them with the closest samples falling in the latent space of the pre-trained 2D generator for iterative optimization. The 2D generator ensures that the rendered results are photorealistic from different views, meanwhile the intrinsic property of the 3D representation guarantees reasonable 3D structure, thus leading to high-quality 3D generation. Some recent methods [33, 20] also combine this idea with text-to-image diffusion models [35, 36] to achieve text-driven 3D creation. Our method shares a similar spirit, which distills the knowledge from the generator's 2D super-resolution branch to its 3D rendering branch, thus achieving image generation with both photorealism and strict 3D consistency.

## 3. Approach

Given a collection of 2D images, we aim to learn a 3D-aware generator $G$ for free-view image synthesis. The generator takes a random code $\boldsymbol{z} \in \mathbb{R}^{d_z}$ and an explicit camera pose $\boldsymbol{\theta} \in \mathbb{R}^{d_\theta}$ as input, and generates a 2D image $I$:

$$G : (\boldsymbol{z}, \boldsymbol{\theta}) \in \mathbb{R}^{d_z} \times \mathbb{R}^{d_\theta} \to I \in \mathbb{R}^{H \times W \times 3}. \quad (1)$$

To enable high-quality image synthesis, we adopt EG3D [3] as the backbone of the generator, which synthesizes low-resolution feature fields via the tri-plane representation [3], and leverages 2D super-resolution for high-resolution image generation (Sec. 3.1). Based on EG3D, we propose a 3D-to-2D imitation strategy to synthesize high-resolution NeRF for 3D-consistent image rendering. We leverage a 3D super-resolution branch to predict high-resolution tri-planes from the low-resolution ones, and force the rendered images from the former to mimic the images generated by the 2D super-resolution branch (Sec. 3.2). In addition, we introduce 3D-aware convolution [46] to the generator for better tri-plane learning via cross-plane communications, which helps to further improve the image generation quality (Sec. 3.3). The overview of our method is illustrated in Fig. 3. We describe each part in detail below.

### 3.1. Preliminaries: EG3D

EG3D adopts a StyleGAN2-based [19] generator $\mathcal{E}$ to efficiently synthesize the low-resolution feature field of a subject. The feature field is represented by the tri-plane representation which consists of three orthogonal 2D planes produced by reshaping the output feature map of $\mathcal{E}$, given the latent code $\boldsymbol{z}$ as input. For a point $\boldsymbol{x} \in \mathbb{R}^3$ in the 3D space, its corresponding feature $\mathbf{f}$ can be obtained by projecting itself onto the three planes $\mathbf{P}_{xy}, \mathbf{P}_{yz}, \mathbf{P}_{zx}$, and summing the retrieved features $\mathbf{f}_{xy}, \mathbf{f}_{yz}, \mathbf{f}_{zx}$. A small MLP $\mathcal{M}$ then maps this intermediate feature to volume density $\sigma \in \mathbb{R}$ and color feature $\boldsymbol{c} \in \mathbb{R}^{d_c}$ (the first three dimensions represent $RGB$ color), forming the low-resolution feature field:

$$\mathcal{M} : \mathbf{f} \in \mathbb{R}^{d_{\mathbf{f}}} \to (\boldsymbol{c}, \sigma) \in \mathbb{R}^{d_c} \times \mathbb{R}. \quad (2)$$

To generate high-resolution images, EG3D enforces volume rendering [15, 24] to render the above feature field to a low-resolution feature map $C$, where each pixel value $C(\boldsymbol{r})$ corresponding to a viewing ray $\boldsymbol{r}$ can be obtained via

$$C(\boldsymbol{r}) = \sum_{i=1}^{N} T_i(1 - \exp(-\sigma_i \delta_i))\boldsymbol{c}_i, \ T_i = \exp(-\sum_{j=1}^{i-1} \sigma_j \delta_j). \quad (3)$$

Here, $i$ is the index of points along ray $\boldsymbol{r}$ sorted from near to far, and $\delta$ is the distance between adjacent points. Then, the rendered feature map $C$ is sent to a 2D super-resolution module $\mathcal{S}^{2D}$ consisting of several StyleGAN2-modulated convolutional layers to generate the final image $I^{2D}$.

Although EG3D can generate free-view images of high quality, it cannot well maintain their 3D consistency across different views. This is inevitable due to incorporating the black-box CNN-based 2D super-resolution mod-
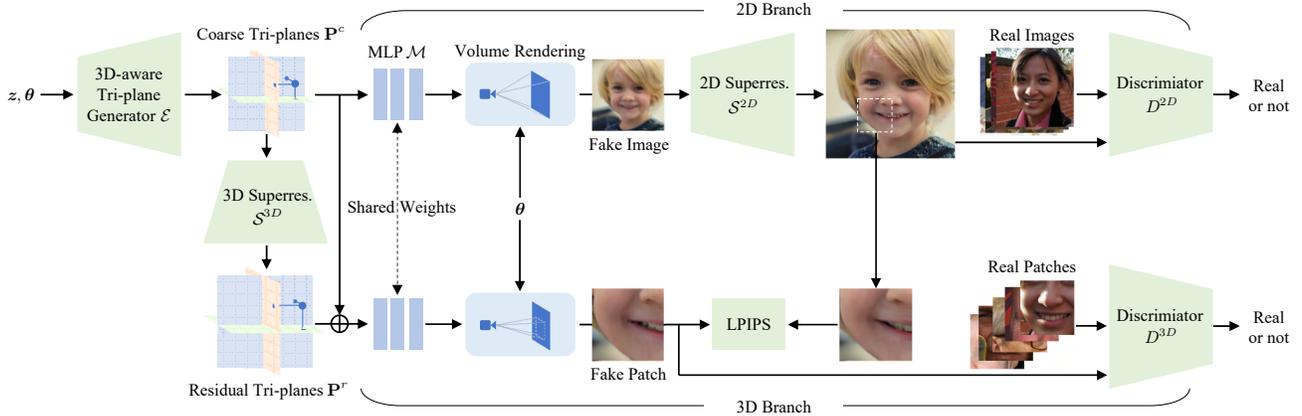
Figure 3. Overview of our framework. 3D-to-2D imitation strategy is enforced to let the generator's 3D branch to mimic the results of its 2D branch, thus leading to image generation of high quality and strict 3D consistency. 3D-aware convolutions are also introduced to the tri-plane generator to enhance 3D representation learning, which further improves the image generation quality.

ule, which breaks the physical rules of the volume rendering process. Despite that EG3D further proposes a dual-discrimination [3] strategy to force the high-resolution images to be consistent with their low-resolution counterparts, detail-level 3D inconsistency (*i.e.* texture flickering) still cannot be eliminated. During continuous camera variation, these artifacts can be easily captured by human eyes, differing the synthesized results from a real video sequence. To maintain the 3D consistency meanwhile keep the high-quality image generation to the maximum extent, we propose a 3D-to-2D imitation strategy described below.

### 3.2. 3D-to-2D Imitation

To keep the strict 3D consistency, a better way is to directly render the 3D representation instead of resorting to a 2D CNN for image synthesis. Noticing that the images generated by EG3D contain rich details, it is natural to use them as guidance for images synthesized by direct 3D rendering. If the directly-rendered images well mimic those fine details, their quality should get very close to that of EG3D. Meanwhile, since they are rendered from a continuous 3D representation, their 3D consistency across different views should be trivially maintained. This motivates us to design the 3D-to-2D imitation strategy, as depicted in Fig. 3.

Specifically, we introduce a 3D super-resolution module $\mathcal{S}^{3D}$ to generate residual tri-planes $\mathbf{P}^r$ from the coarse tri-planes $\mathbf{P}^c$ produced by the tri-plane generator $\mathcal{E}$:

$$\mathcal{S}^{3D} : \mathbf{P}^c \in \mathbb{R}^{3 \times H^c \times W^c \times d^{\mathbf{f}}} \to \mathbf{P}^r \in \mathbb{R}^{3 \times H^r \times W^r \times d^{\mathbf{f}}}. \quad (4)$$

The $\mathcal{S}^{3D}$ adopts several StyleGAN2-modulated convolutional layers conditioned on a latent code $\boldsymbol{w}$ mapped from the random code $\boldsymbol{z}$, similar to the 2D super-resolution module $\mathcal{S}^{2D}$ in EG3D. The difference is that $\mathcal{S}^{3D}$ conducts super-resolution on the triplane-based 3D representation instead of the rendered 2D feature map. In this way, we can

generate a high-resolution 3D field for direct 3D rendering. Given the coarse and residual tri-planes (*i.e.* $\mathbf{P}^c$ and $\mathbf{P}^r$), we obtain a more detailed intermediate feature $\mathbf{f} = \mathbf{f}^c + \mathbf{f}^r$ for a 3D point $\boldsymbol{x}$, and further obtain the high-resolution feature field by sending the intermediate feature into the MLP-based decoder $\mathcal{M}$ following Eq. (2). The first three feature dimensions of the field derive the high-resolution NeRF for rendering 3D-consistent fine image $I^{3D}$ via Eq. (3).

To ensure that $I^{3D}$ contains reasonable geometry structure with rich texture details, we let it to mimic the contents of $I^{2D}$ generated by the 2D branch $\mathcal{S}^{2D}$. For a pair of $I^{3D}$ and $I^{2D}$ synthesized with the same latent code $\boldsymbol{z}$ and camera pose $\boldsymbol{\theta}$, we enforce imitation loss between them to guarantee their perceptual similarity:

$$\mathcal{L}_{imitation} = \text{LPIPS}(I^{3D}, \text{sg}(I^{2D})), \quad (5)$$

where $\text{LPIPS}(\cdot, \cdot)$ is the perceptual loss defined in [51], and sg denotes stopping gradient to avoid undesired influence of $I^{3D}$ on the 2D branch. This process is very similar to a standard multiview reconstruction process. During training, $I^{2D}$ sharing the same code $\boldsymbol{z}$ are generated under different camera views from a statistical aspect, forming the multi-view supervision. The high-resolution NeRF from the 3D branch renders $I^{3D}$ under the same camera views to compare with the multiview data for 3D reconstruction. Considering that $I^{2D}$ are nearly 3D-consistent, they should help to learn a reasonable NeRF for 3D-consistent image rendering.

Nevertheless, since $I^{2D}$ are not strictly 3D-consistent, faithfully reconstructing their image contents leads to blurry results where the texture details across different views are averaged out. Therefore, we further introduce the non-saturating GAN loss with R1 regularization [22] between $I^{3D}$ and real images $\hat{I}$ to maintain the high-frequency de-

tails:

$$\mathcal{L}_{adv}^{3D} = \mathbb{E}_{\boldsymbol{z} \sim p_z, \boldsymbol{\theta} \sim p_\theta}[f(D^{3D}(G^{3D}(\boldsymbol{z}, \boldsymbol{\theta})))]$$
$$+ \mathbb{E}_{\hat{I} \sim p_{real}}[f(-D^{3D}(\hat{I})) + \lambda \|\nabla D^{3D}(\hat{I})\|^2], \quad (6)$$

where $f(u) = \log(1 + \exp(u))$ is the Softplus function, $G^{3D}$ including $\{\mathcal{E}, \mathcal{M}, \mathcal{S}^{3D}\}$ is the 3D rendering branch of the generator, and $D^{3D}$ is the corresponding discriminator.

An advantage of the above imitation learning is that we can render small patches (*i.e.* $64 \times 64$) to compute Eq. (5) and Eq. (6), as shown in Fig. 3, with only minor influence to the final image quality. This largely reduces the memory cost during training and enables learning the 3D branch at high resolution (*e.g.* $512 \times 512$). By contrast, solely applying adversarial loss at patch-level often leads to large quality drops as shown in previous methods [37, 43] and Tab. 2.

Finally, we apply image-level adversarial loss to the 2D branch following EG3D to ensure that $I^{2D}$, as the supervision for the 3D branch, are of high quality:

$$\mathcal{L}_{adv}^{2D} = \mathbb{E}_{\boldsymbol{z} \sim p_z, \boldsymbol{\theta} \sim p_\theta}[f(D^{2D}(G^{2D}(\boldsymbol{z}, \boldsymbol{\theta})))]$$
$$+ \mathbb{E}_{I \sim p_{real}}[f(-D^{2D}(\hat{I})) + \lambda \|\nabla D^{2D}(\hat{I})\|^2], \quad (7)$$

where $G^{2D}$ is the 2D branch generator consisting of $\{\mathcal{E}, \mathcal{M}, \mathcal{S}^{2D}\}$, and $D^{2D}$ is the corresponding discriminator. The same dual discrimination is adopted as done in EG3D.

Overall, the training objective is

$$\mathcal{L}_{total} = \mathcal{L}_{imitation} + \mathcal{L}_{adv}^{3D} + \mathcal{L}_{adv}^{2D}. \quad (8)$$

In practice, we first learn the 2D branch via $\mathcal{L}_{adv}^{2D}$ to obtain reasonable synthesized images $I^{2D}$, then leverage $\mathcal{L}_{total}$ to simultaneously learn the 2D and 3D branches for high-quality and 3D-consistent image synthesis.

### 3.3. 3D-Aware Tri-plane Generator

As depicted in Sec. 3.2, the tri-plane generator $\mathcal{E}$ is responsible for synthesizing the coarse tri-planes $\mathbf{P}^c$ shared by both the 2D and 3D branches, which is an important component that would affect the final image generation quality. However, in EG3D, $\mathcal{E}$ takes a StyleGAN2 architecture originally designed for 2D generative tasks. As shown in Fig. 4(a), the original tri-plane generator only contains the main stream and the output stream. The tri-planes are obtained from latent feature maps in the main stream via 2D convolutions (*i.e.* $toRGB$ layers), and the latent feature maps are also produced by a serials of 2D synthesis blocks. Consequently, the latent feature maps are forced to learn 3D unaligned features of the three orthogonal planes and the latters also lack feature communications with each other. Inspired by a recent 3D diffusion model [46], we introduce 3D-aware convolutions into our tri-plane generator $\mathcal{E}$ to enhance feature communications between 3D-associated positions across different planes, for better tri-plane generation.
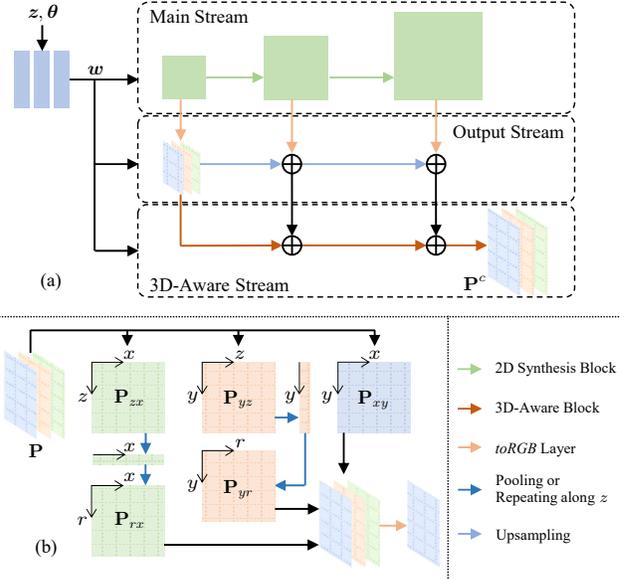


Figure 4. (a) Structure of our 3D-aware tri-plane generator. (b) Operations of the 3D-aware block on $xy$ plane.

Specifically, as illustrated in Fig. 4(a), we add an extra 3D-aware stream upon the original output stream after each $toRGB$ layer at different resolutions. At each resolution level $k$, the corresponding tri-planes $\mathbf{P}_k = [\mathbf{P}_{k,xy}, \mathbf{P}_{k,yz}, \mathbf{P}_{k,zx}]$, are summed with the tri-planes produced by the original output stream, and further sent into a 3D-aware block to produce tri-plane features for the next level. The 3D-aware block conducts similar operations on each of the three planes. For brevity, we omit the subscript $k$ here and take $\mathbf{P}_{xy}$ as an example to illustrate the operation process. As shown in Fig. 4(b), to align $\mathbf{P}_{yz}$ and $\mathbf{P}_{zx}$ towards $\mathbf{P}_{xy}$, we first perform global pooling along $z$ axis of the former two to obtain $z$-squeezed feature vectors. These vectors are then repeated along the $z$ dimension to restore the original spatial size, denoted as $\mathbf{P}_{yr}$ and $\mathbf{P}_{rx}$. In this manner, the obtained $\mathbf{P}_{yr}$ and $\mathbf{P}_{rx}$ are aligned with $\mathbf{P}_{xy}$ from a 3D perspective, *i.e.*, a 2D position $uv$ on $\mathbf{P}_{xy}$ is responsible for features in region $uvz, z \in [z_{min}, z_{max}]$ in the 3D space, meanwhile the same $uv$ position on $\mathbf{P}_{yr}$ and $\mathbf{P}_{rx}$ also associate with the features in this 3D region. As a result, we can simply concatenate them along the channel dimension as $[\mathbf{P}_{xy}, \mathbf{P}_{yr}, \mathbf{P}_{rx}]$, and perform modulated 2D convolution [19] on it. The 2D convolution aggregates the 3D-associated features to produce next-level $\mathbf{P}_{xy}$, leading to better feature communications across the planes. $\mathbf{P}_{yz}$ and $\mathbf{P}_{zx}$ can be processed similarly.

Note that in [46], the 3D-aware convolution is applied in all layers in a U-Net structure. However, in our scenario, leveraging 3D-aware convolution for all layers, especially the main stream, introduces unaffordable memory cost during training, as it would produce multiple auxiliary tensors
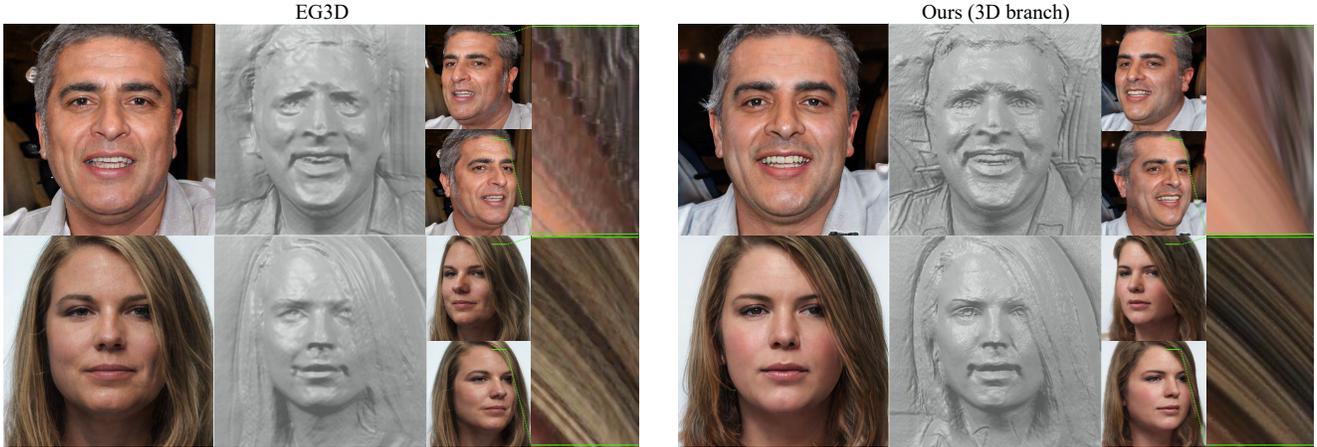
EG3D

Ours (3D branch)

Figure 5. Comparison between our method and EG3D on FFHQ at $512 \times 512$ resolution. Our method generates images with comparable quality to those of EG3D, while producing 3D geometries with finer details and multiview sequences with better 3D-consistency.



GMPI        EpiGRAF        Ours (3D branch)

Figure 6. Comparison between our method and other 3D rendering baselines on FFHQ at $512 \times 512$ resolution. **Best viewed with zoom-in.**

and triples the channel dimension for each processed latent feature map, as shown in Fig. 4(b). Comparing to the latent feature maps in the main stream, the tri-planes after each output layer contains much fewer channels thus more memory-friendly to adopt the 3D-aware convolution. Empirically, our proposed 3D-aware stream helps to learn more reasonable tri-planes and improves the final image generation quality, with only a minor increase in the total memory consumption (see Sec. 4.3).

## 4. Experiments

**Implementation details.** We train our method on two real-world datasets: FFHQ [18] and AFHQ-v2 Cats [5], which consists of 70K human face images of $1024^2$ resolution and 5.5K cat face images of $512^2$ resolution, respectively. We follow the data pre-processing of EG3D [3] to crop and resize the images to $256^2$ or $512^2$ resolution. Experiments are conducted on 8 NVIDIA Tesla A100 GPUs with 40GB memory, following the training configuration of EG3D. For FFHQ, the training process takes around 8 days, where learning the 2D branch takes 5 days and jointly training the whole framework takes additional 3 days. For AFHQ-v2, we finetune the 2D branch initially trained on FFHQ for 1 day, then jointly train the whole framework for extra 3 days. Adaptive data augmentation [16] is applied to AFHQ-v2 to facilitate training with limited data. See the *suppl. material* for more details.

### 4.1. Visual Results

Figure 2 shows the multiview images generated by our 3D branch generator. It can produce high-quality images with fine details at a resolution of $512^2$. Moreover, the images are of strict 3D consistency across different views via directly rendering the generated high-resolution NeRF. More results are in Fig. 5, 6, and the *suppl. material*.

### 4.2. Comparison with Prior Arts

**Baselines.** We compare our method with existing 3D-aware GANs, including methods leveraging 2D super-resolution: StyleSDF [28], VolumeGAN [49], StyleN-eRF [9], and EG3D [3]; and methods with direct 3D rendering: GRAM [6], GRAM-HD [47], GMPI [54], Epi-GRAF [43], and VoxGRAF [38].

**Qualitative comparison.** Figure 5 shows the visual comparison between our method and EG3D. Our generated images via direct rendering have comparable quality with those generated by EG3D via 2D super-resolution. We further visualize the 3D geometry and the spatiotemporal texture images [47] of the two methods. The geometry is extracted via Marching Cubes [21] on the density field at $512^3$ resolution. The spatiotemporal textures are obtained by stacking the pixels of a fixed line segment under continuous camera change, very similar to the Epipolar Line Images [1], where smoothly tilted strips indicate better 3D

| | Method | FFHQ256 | | | FFHQ512 | | | CATS256 | CATS512 |
|---|---|---|---|---|---|---|---|---|---|
| | | FID $\downarrow$ | PSNR$_{mv}$ $\uparrow$ | SSIM$_{mv}$ $\uparrow$ | FID $\downarrow$ | PSNR$_{mv}$ $\uparrow$ | SSIM$_{mv}$ $\uparrow$ | FID $\downarrow$ | FID $\downarrow$ |
| *w/ 2D SR* | StyleSDF [28] | 11.5 | - | - | 11.2 | - | - | - | 7.91 |
| | VolumeGAN [49] | 9.10 | 33.6 | 0.926 | - | - | - | - | - |
| | StyleNeRF [9] | 8.00 | 31.9 | 0.915 | 7.80 | 30.9 | 0.843 | - | - |
| | EG3D [3] | 4.80 | 34.0 | 0.928 | 4.70 | 32.4 | 0.861 | 3.88 | 2.77 |
| | Ours (2D branch) | **3.91** | **35.7** | **0.938** | **4.14** | **33.3** | **0.891** | **3.41** | **2.72** |
| *w/o 2D SR* | GRAM [6] | 13.8 | 38.0 | 0.966 | - | - | - | 13.4 | - |
| | GRAM-HD [47] | 10.4 | 36.5 | 0.955 | - | - | - | - | 7.67 |
| | GMPI [54] | 11.4 | **39.8** | **0.977** | 8.29 | **39.0** | **0.961** | - | 7.79 |
| | EpiGRAF [43] | 9.71 | - | - | 9.92 | 37.3 | 0.949 | 6.93 | - |
| | VoxGRAF [38] | 9.60 | 37.2 | 0.960 | - | - | - | 9.60 | - |
| | Ours (3D branch) | **5.14** | 39.3 | 0.974 | **5.37** | 37.8 | 0.955 | **4.14** | **4.29** |

Table 1. Comparison on image generation quality and 3D consistency among different 3D-aware GANs.

consistency. As shown, our geometries contain finer details in that we directly learn the NeRF of a subject at high resolution. Our spatiotemporal textures are also more reasonable with fewer twisted patterns, thanks to the direct 3D rendering for image synthesis instead of using a black-box 2D super-resolution module.

Figure 6 compares our method with other 3D baselines on FFHQ at $512^2$ resolution. Visually inspected, our 3D branch produces images of higher fidelity compared to existing methods leveraging direct 3D rendering. More analysis and video results can be found in the *suppl. material*.

**Quantitative comparison.** Table 1 and Fig. 1 show the quantitative results of different methods in terms of image generation quality and 3D consistency. For image generation quality, We calculate the Fréchet Inception Distance (FID) [13] between 50K generated images and all available real images in the training set. For 3D consistency, we follow GRAM-HD [47] to generate multiview images of 50 random subjects and train the multiview reconstruction method NeuS [45] on each of them. We report the average PSNR and SSIM scores between our generated multiview images and the re-rendered images of NeuS (denoted as PSNR$_{mv}$ and SSIM$_{mv}$). Theoretically, better 3D consistency facilitates the 3D reconstruction process of NeuS, thus leading to higher PSNR and SSIM.

As shown, our 2D branch generator demonstrates better results compared to EG3D in all metrics across different datasets, thanks to our 3D-aware stream in the tri-plane generator. Moreover, with the 3D-to-2D imitation strategy, our 3D branch generator largely improves the image generation quality among methods using direct 3D rendering, while maintaining competitive 3D consistency. Its image quality even surpasses most of the methods with 2D super-resolution and comes very close to that of EG3D.

## 4.3. Ablation Study

We conduct ablation studies to validate the efficacy of our proposed 3D-to-2D imitation and the 3D-aware tri-

| Label | $\mathcal{L}_{imitation}$ | $\mathcal{S}^{3D}$ | $\mathcal{L}_{adv}^{3D}$ | FID (3D branch) |
|---|---|---|---|---|
| (A) | | | | 30.6 |
| (B) | ✓ | | | 29.9 |
| (C) | ✓ | ✓ | | 9.29 |
| (D) | | ✓ | ✓ | 22.8 |
| (E) | ✓ | ✓ | ✓ | **5.14** |

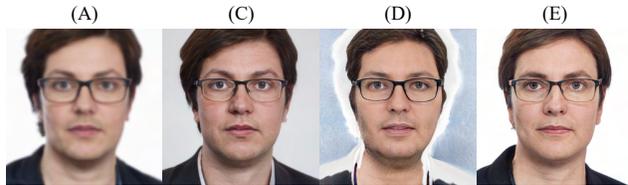Table 2. Ablation study on 3D-to-2D imitation strategy.



Figure 7. Generated images under different learning strategies. The labels are consistent with Tab. 2.

plane generator. For efficiency, all experiments are conducted on FFHQ dataset at $256^2$ resolution.

**3D-to-2D imitation strategy.** As shown in Tab. 2 and Fig. 7, We start from a generator without using the 3D-to-2D imitation and the 3D super-resolution module $\mathcal{S}^{3D}$ (setting A), by directly rendering the coarse tri-planes $\mathbf{P}^c$ for image synthesis. The rendered images in this way are blurry and lack fine details, leading to a high FID score of 30.6. Naively introducing the imitation loss (setting B) to improve the rendered images of $\mathbf{P}^c$ has minor influence, as the capacity of the coarse tri-planes are limited. Further incorporating the 3D super-resolution module (setting C) effectively releases the potential of the imitation loss and largely improves the image generation quality in terms of FID. However, the rendered images still lack rich details limited by the 3D-inconsistent 2D branch supervisions. Then, if the imitation loss is replaced with the adversarial loss (setting D), the image quality decreases significantly. This is due to that we only render small image patches to compute the corresponding losses for memory consideration. Under this circumstance, the adversarial loss is less

| Method | FID (2D) | FID (3D) | #Param | Mem. |
|---|---|---|---|---|
| w/o 3D-aware | 4.80 | 6.71 | 29.0M | 2.3G |
| 3D-aware latent | OOM | OOM | 111.7M | 11.6G |
| 3D-aware tri-plane | 4.14 | - | 32.6M | 2.4G |
| 3D-aware stream (Ours) | **3.91** | **5.14** | 32.6M | 2.4G |

Table 3. Ablation study on designs of 3D-aware tri-plane generator. The FID scores are from 2D or 3D branch; #Param only considers the tri-plane generator $\mathcal{E}$ and Mem. indicates the GPU memory cost for generating the coarse tri-planes.
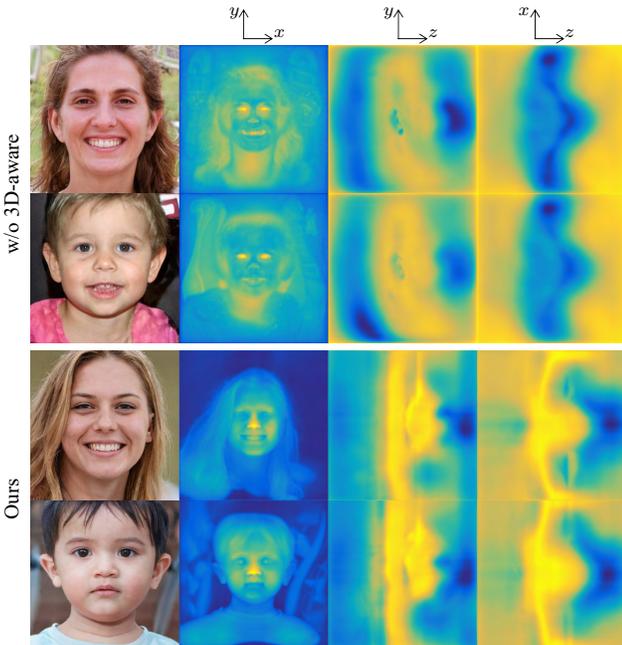


Figure 8. Generated tri-planes with or w/o 3D-aware convolutions.

stable compared to the imitation loss which is a perceptual-level reconstruction loss. This reveals the advantage of our imitation strategy, which could be extended to higher resolution via patch-wise optimization while maintaining a good image generation quality. Finally, leveraging all the three components (setting E) yields the best result, where the imitation loss keeps the overall structure reasonable and the adversarial loss helps with fine details learning.

**3D-aware tri-plane generator.** Table 3 shows the ablation study on the 3D-aware tri-plane generator. We compare our design with two alternatives and one without 3D-aware convolutions originally adopted by EG3D. We report the parameter size of the tri-plane generators, the inference memory cost to generate the coarse tri-planes, as well as the final image generation quality in terms of FID. In the first alternative, we remove our 3D-aware stream, and leverage 3D-aware convolutions for the latent feature maps in the main stream, namely *3D-aware latent*. Since the main stream feature maps have relatively larger feature channels, and the 3D-aware convolution requires to concatenate two additional tensors with the same size as the input tensor, this

design increases the parameter size and memory consumption significantly, and raises the out-of-memory issue during training. In the second alternative, namely *3D-aware tri-plane*, we directly apply 3D-aware convolutions in the output stream, by inserting them after the upsampling operations at each resolution, instead of using the additional 3D-aware stream. This strategy leads to an improvement of the image generation quality of the 2D branch, and largely reduces the parameter size and memory cost compare to the first design. Finally, our 3D-aware stream design further improves the image generation quality without introducing extra parameters and memory costs. Therefore, we adopt it as our final 3D-aware tri-plane generator for 3D-to-2D imitation. It effectively lowers the FID score of both the 2D and 3D branches compared to the original structure without 3D-aware convolutions, with only a minor increase of the parameter size and memory cost.

Figure 8 further shows the synthesized tri-planes, where we visualize the L2 norm of each spatial location on the three orthogonal planes. Our method leveraging the 3D-aware stream produces more informative tri-planes. The generated planes of the side-views better depict the characters of different instances (*e.g.*, see the difference of the profiles on the $yz$ planes). Our frontal planes (*i.e.* $xy$ planes) also demonstrate more clear head silhouettes compared to those without using the 3D-aware convolutions.

## 5. Conclusions

We presented a novel learning strategy for 3D-aware GANs to achieve image synthesis of high-quality and strict 3D consistency. The core idea is to enforce the images synthesized by the generator's 3D rendering branch to mimic those generated by its 2D super-resolution branch. We also introduced 3D-aware convolutions to the generator to further improve the image generation quality. With the above strategies, our method largely improves the image quality among methods using direct 3D rendering, which we believe enables a new way for more realistic 3D generation.

**Limitation and future works.** Our method has several limitations. The image generation quality of its 3D branch still lags behind that of the 2D branch. Certain generated 3D structures such as hairs and cat whiskers are stuck to the geometry surfaces instead of correctly floating in the volumetric space. The 3D-to-2D imitation strategy also introduces extra training time and memory costs compared to only learning the 2D branch. We expect more effective learning strategies and more advanced 3D representations to alleviate these problems.

**Ethics consideration.** The goal of this paper is to generate images of virtual subjects. It is not intended for creating misleading or deceptive contents of real people and we do not condone any such harmful behavior.

# References

[1] Robert C Bolles, H Harlyn Baker, and David H Marimont. Epipolar-plane image analysis: An approach to determining structure from motion. *International journal of computer vision*, 1(1):7–55, 1987. 6

[2] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis. In *Proceedings of the International Conference on Learning Representations*, 2019. 3

[3] Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas Guibas, Jonathan Tremblay, Sameh Khamis, et al. Efficient geometry-aware 3D generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. 1, 2, 3, 4, 6, 7, 12, 13, 14

[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*, 2021. 1, 2, 3

[5] Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. StarGAN v2: Diverse image synthesis for multiple domains. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020. 2, 6

[6] Yu Deng, Jiaolong Yang, Jianfeng Xiang, and Xin Tong. GRAM: Generative radiance manifolds for 3D-aware image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. 1, 2, 3, 6, 7, 12, 13

[7] Terrance DeVries, Miguel Angel Bautista, Nitish Srivastava, Graham W Taylor, and Joshua M Susskind. Unconstrained scene generation with locally conditioned radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021. 3

[8] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Proceedings of the Advances in Neural Information Processing Systems*, 2014. 1, 3

[9] Jiatao Gu, Lingjie Liu, Peng Wang, and Christian Theobalt. StyleNeRF: A style-based 3D aware generator for high-resolution image synthesis. In *Proceedings of the International Conference on Learning Representations*, 2022. 1, 2, 3, 6, 7

[10] Erik Härkönen, Aaron Hertzmann, Jaakko Lehtinen, and Sylvain Paris. GANSpace: Discovering interpretable GAN controls. In *Proceedings of the Advances in Neural Information Processing Systems*, 2020. 3

[11] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020. 13

[12] 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*, 2019. 2, 3

[13] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. GANs trained by a two time-scale update rule converge to a local nash equilibrium. In *Proceedings of the Advances in Neural Information Processing Systems*, 2017. 2, 7

[14] Ali Jahanian, Lucy Chai, and Phillip Isola. On the "steerability" of generative adversarial networks. In *Proceedings of the International Conference on Learning Representations*, 2020. 3

[15] James T Kajiya and Brian P Von Herzen. Ray tracing volume densities. *ACM SIGGRAPH*, 18(3):165–174, 1984. 1, 3

[16] Tero Karras, Miika Aittala, Janne Hellsten, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Training generative adversarial networks with limited data. In *Proceedings of the Advances in Neural Information Processing Systems*, 2020. 6

[17] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Alias-free generative adversarial networks. In *Proceedings of the Advances in Neural Information Processing Systems*, 2021. 1

[18] 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. 2, 3, 6

[19] 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*, 2020. 1, 3, 5, 12

[20] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3D: High-resolution text-to-3D content creation. *arXiv:2211.10440*, 2022. 3

[21] William E Lorensen and Harvey E Cline. Marching cubes: A high resolution 3D surface construction algorithm. *ACM SIGGRAPH*, 21(4):163–169, 1987. 6

[22] Lars Mescheder, Andreas Geiger, and Sebastian Nowozin. Which training methods for GANs do actually converge? In *Proceedings of the International Conference on Machine Learning*, 2018. 4

[23] 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*, 2019. 3

[24] 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 *Proceedings of the European Conference on Computer Vision*, 2020. 1, 3

[25] 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. 2, 3

[26] 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, 2, 3

[27] Michael Oechsle, Songyou Peng, and Andreas Geiger. Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021. 3

[28] Roy Or-El, Xuan Luo, Mengyi Shan, Eli Shechtman, Jeong Joon Park, and Ira Kemelmacher-Shlizerman. StyleSDF: High-resolution 3D-consistent image and geometry generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2022. 1, 3, 6, 7

[29] Xingang Pan, Bo Dai, Ziwei Liu, Chen Change Loy, and Ping Luo. Do 2D GANs know 3D shape? unsupervised 3D shape reconstruction from 2d image gans. In *Proceedings of the International Conference on Learning Representations*, 2021. 2, 3

[30] Xingang Pan, Ayush Tewari, Lingjie Liu, and Christian Theobalt. Gan2x: Non-lambertian inverse rendering of image gans. In *Proceedings of the International Conference on 3D Vision*, 2022. 3

[31] Xingang Pan, Xudong Xu, Chen Change Loy, Christian Theobalt, and Bo Dai. A shading-guided generative implicit model for shape-accurate 3D-aware image synthesis. In *Proceedings of the Advances in Neural Information Processing Systems*, 2021. 3

[32] 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*, 2019. 3

[33] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3D using 2D diffusion. *arXiv:2209.14988*, 2022. 2, 3

[34] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *Proceedings of the International Conference on Machine Learning*, 2021. 13

[35] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. 3

[36] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Raphael Gontijo-Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. In *Proceedings of the Advances in Neural Information Processing Systems*, 2022. 3

[37] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. GRAF: Generative radiance fields for 3D-aware image synthesis. In *Proceedings of the Advances in Neural Information Processing Systems*, 2020. 1, 2, 3, 5

[38] Katja Schwarz, Axel Sauer, Michael Niemeyer, Yiyi Liao, and Andreas Geiger. VoxGRAF: Fast 3D-aware image synthesis with sparse voxel grids. In *Proceedings of the Advances in Neural Information Processing Systems*, 2022. 1, 3, 6, 7

[39] Yujun Shen, Jinjin Gu, Xiaoou Tang, and Bolei Zhou. Interpreting the latent space of GANs for semantic face editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020. 3

[40] 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*, 2021. 3

[41] Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. In *Proceedings of the Advances in Neural Information Processing Systems*, 2020. 3

[42] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: continuous 3D-structure-aware neural scene representations. In *Proceedings of the Advances in Neural Information Processing Systems*, 2019. 3

[43] Ivan Skorokhodov, Sergey Tulyakov, Yiqun Wang, and Peter Wonka. EpiGRAF: Rethinking training of 3D GANs. In *Proceedings of the Advances in Neural Information Processing Systems*, 2022. 1, 2, 3, 5, 6, 7, 12, 15

[44] Attila Szabó, Givi Meishvili, and Paolo Favaro. Unsupervised generative 3D shape learning from natural images. *arXiv:1910.00287*, 2019. 3

[45] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. NeuS: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. In *Proceedings of the Advances in Neural Information Processing Systems*, 2021. 1, 3, 7

[46] Tengfei Wang, Bo Zhang, Ting Zhang, Shuyang Gu, Jianmin Bao, Tadas Baltrusaitis, Jingjing Shen, Dong Chen, Fang Wen, Qifeng Chen, et al. Rodin: A generative model for sculpting 3D digital avatars using diffusion. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2023. 2, 3, 5

[47] Jianfeng Xiang, Jiaolong Yang, Yu Deng, and Xin Tong. GRAM-HD: 3D-consistent image generation at high resolution with generative radiance manifolds. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023. 1, 3, 6, 7

[48] Xudong Xu, Xingang Pan, Dahua Lin, and Bo Dai. Generative occupancy fields for 3D surface-aware image synthesis. In *Proceedings of the Advances in Neural Information Processing Systems*, 2021. 3

[49] Yinghao Xu, Sida Peng, Ceyuan Yang, Yujun Shen, and Bolei Zhou. 3D-aware image synthesis via learning structural and textural representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18430–18439, 2022. 3, 6, 7

[50] Yang Xue, Yuheng Li, Krishna Kumar Singh, and Yong Jae Lee. Giraffe HD: A high-resolution 3D-aware generative

model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18440–18449, 2022. 1, 3

[51] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018. 2, 4

[52] Xuanmeng Zhang, Zhedong Zheng, Daiheng Gao, Bang Zhang, Pan Pan, and Yi Yang. Multi-view consistent generative adversarial networks for 3D-aware image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. 2

[53] Yuxuan Zhang, Wenzheng Chen, Huan Ling, Jun Gao, Yinan Zhang, Antonio Torralba, and Sanja Fidler. Image gans meet differentiable rendering for inverse graphics and interpretable 3D neural rendering. In *Proceedings of the International Conference on Learning Representations*, 2021. 2, 3

[54] Xiaoming Zhao, Fangchang Ma, David Güera, Zhile Ren, Alexander G Schwing, and Alex Colburn. Generative multiplane images: Making a 2D GAN 3D-aware. In *Proceedings of the European Conference on Computer Vision*, 2022. 1, 2, 3, 6, 7, 12, 15

[55] Peng Zhou, Lingxi Xie, Bingbing Ni, and Qi Tian. Cips-3d: A 3D-aware generator of gans based on conditionally-independent pixel synthesis. *arXiv:2110.09788*, 2021. 3

## A. More Implementation Details

### A.1. Network Structure

Figure 12 illustrates our network designs, including the 3D super-resolution module $\mathcal{S}^{3D}$ and the 3D-aware block in the tri-plane generator $\mathcal{E}$.

For $\mathcal{S}^{3D}$ (Fig. 12(a)), we use two modulated 2D convolution blocks [19] to upsample the tri-planes.

For the 3D-aware block (Fig. 12(b)), we re-organize the tri-planes according to Fig. 4 in the main text, and apply modulated 2D convolutions for each of the three planes. We use different affine layers to generate style codes for the three modulated convolutions, respectively.

### A.2. Training Details

We randomly sample latent code $z$ from the normal distribution and camera pose $\boldsymbol{\theta}$ from those of the training datasets to synthesize fake images, following EG3D [3]. For each viewing ray, we sample 96 points to calculate the volume rendering equation, including 48 points with stratified sampling and 48 points with importance sampling. The learning rates of the generator and the two discriminators are set to 0.0025 and 0.002, respectively. We train the 2D branch with 25M images in total, and then jointly train the whole framework with additional 15M images. The batch size during training is set as 32. Other training settings are identical to those of EG3D [3].

### A.3. Patch Scale

To reduce GPU memory costs and enable training at high resolution, we render $64^2$ patches for the 3D-to-2D imitation. Thus, the patch scale is $1/4$ or $1/8$ of the whole image for the $256^2$ or $512^2$ experiments, respectively. The patch center is uniformly sampled from the whole image space.

### A.4. The necessity of 2D super-resolution module

The function of the 2D super-resolution in the 2D branch is to provide stable and high-quality guidance for the 3D branch. Previous studies have attempted to directly learn in 3D space without 2D super-resolution via the adversarial loss. However, due to the restriction of modern GPU memory, they either adopted more efficient 3D representations (*e.g.*, radiance manifolds [6] or MPI [54]) or used patchwise loss (*e.g.*, EpiGRAF [43]), yet these strategies often lead to worse diversity and image quality due to the instability of the GAN loss. By contrast, our imitation with the 2D branch via LPIPS loss provides stable gradients for learning the 3D representation, and thus supports patch-wise training without sacrificing the generation quality, which is the key to our superior results. Furthermore, our strategy also avoids troublesome training tricks (*e.g.*, the annealed strategy in EpiGRAF [43]) thus easier to be adapted to other frameworks.
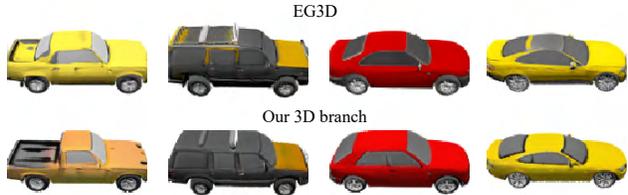


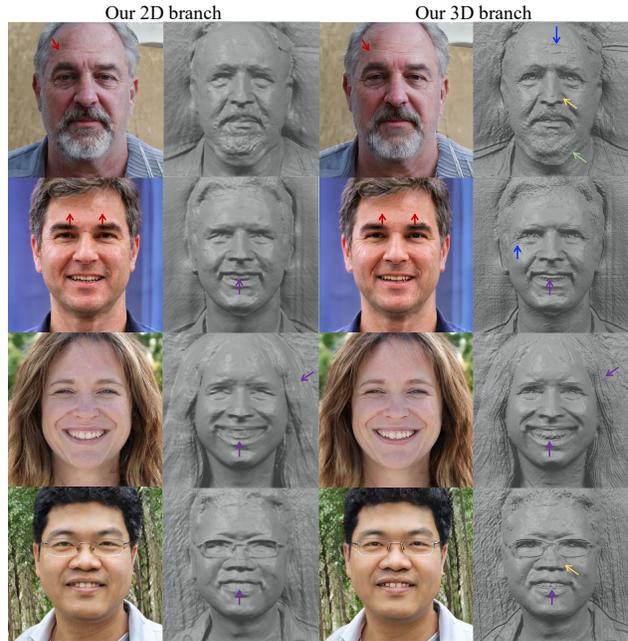Figure 9. Comparison with EG3D on ShapeNet-Cars.



Figure 10. Comparison of 2D and 3D branches. (Zoom in for better visualization.)

### A.5. Training time/memory of 3D-to-2D imitation

Our method requires 31 GB memory at $256^2$ resolution with a batch size of 32 when trained on 8 GPUs, compared to 27 GB memory without the 3D-to-2D imitation. Also, our training time is 1.5 times longer than that of EG3D.

## B. More Results and Comparisons

### B.1. End-to-end 3D-to-2D imitation learning

Our initial motivation for the two-stage training is to leverage the powerful prior of an existing 2D generator (with 2D super-resolution) to guide our 3D branch. In fact, the overall framework (including both 2D and 3D branches) can be trained end-to-end from scratch. We conduct a simple experiment on FFHQ at $256^2$ with identical hyper parameters as described in the main paper and achieve an FID of 5.03 for the 3D branch, which is comparable to the two-stage training result.

Figure 11. Failure case.

## B.2. More results on faces

Figures 13 and 14 illustrate more visual comparisons. Compared to EG3D [3], we have more detailed geometry and smoothly tilted strips in spatiotemporal texture images, indicating better 3D consistency. Similar to ours, EpiGRAF and GMPI also generate high-resolution images via direct rendering. Yet, we have superior image quality as shown in Fig. 14.

Figures 15 and 16 show more of our results on FFHQ and AFHQ -v2 Cats datasets, respectively.

***Referring to the supplemental video for animations.***

## B.3. Results on general objects.

Our method can handle general objects with wider range of camera views. In Fig. 9, we compare our 3D branch with EG3D on ShapeNet-Cars ($128^2$) and achieve comparable image generation quality.

## B.4. Comparison of our 2D and 3D branches

Our 3D branch can generate fine details comparable to the 2D branch. In Fig. 10 (red arrows), we show details produced by the 3D branch that are not visible in the 2D branch.

Our 3D branch clearly produces finer geometry details compared to the alternatives with 2D super-resolution (see Fig. 10). As shown, the finer geometry details are not random noises but features of hair, teeth, wrinkles, etc (purple arrows). Furthermore, we can generate diverse nose shapes (yellow arrows), complex jaws with beards (green arrows), and wrinkles (blue arrows) on the geometries.

## C. Limitations and Future Works

We thoroughly discuss the limitations of our method and possible future improvements.

First, our learned 3D branch still has inferior image quality in terms of FID compared to the 2D branch. This may come from the current design of the 3D super-resolution module and the learning strategy. Specifically, our 3D super-resolution module adopts a similar structure to that of the 2D branch in order for a fair comparison, which may not be the optimal solution. More advanced structures, including leveraging 3D-aware convolutions could be

further explored for better 3D super-resolution. Besides, the LPIPS loss during 3D-to-2D imitation leverages a pre-trained VGG network which is trained on images of $224^2$ resolution. It may not well capture the perceptual information of a small image patch. Leveraging more recent pretrained models [11, 34] or even multiple feature extractors could be a possible choice. Exploring better discriminators for the patch-level adversarial loss in the 3D branch could also benefit the training process.

Second, our method can produce incorrect geometries in certain cases. As shown in Fig. 11, a typical failure case is geometry discontinuity, where the face region is not smoothly connected with the head region, leading to obvious artifacts at side views. These artifacts also occur in the original EG3D. We believe this problem can be alleviated by introducing more profile images for training, as currently the training data are mostly frontal images so that the planes for depicting side-view features may not be welltrained. In addition, certain generated geometry structures such as hairs and cat whiskers are stuck to the surfaces instead of correctly floating in the volumetric space, as shown in Fig. 15 and 16. We conjecture this is due to that the random sampling strategy with limited queries during volume rendering is hard to model thin structures, as also indicated by a previous method [6]. Therefore, a more advanced 3D representation that could efficiently capture these complex structures is worthy of ongoing exploration.

Finally, our training strategy also requires training the 2D branch in advance, which increases the overall training time compared to the original EG3D. A possible way to reduce the training time is to jointly train the 2D and 3D branches from scratch. We leave it for our future work.
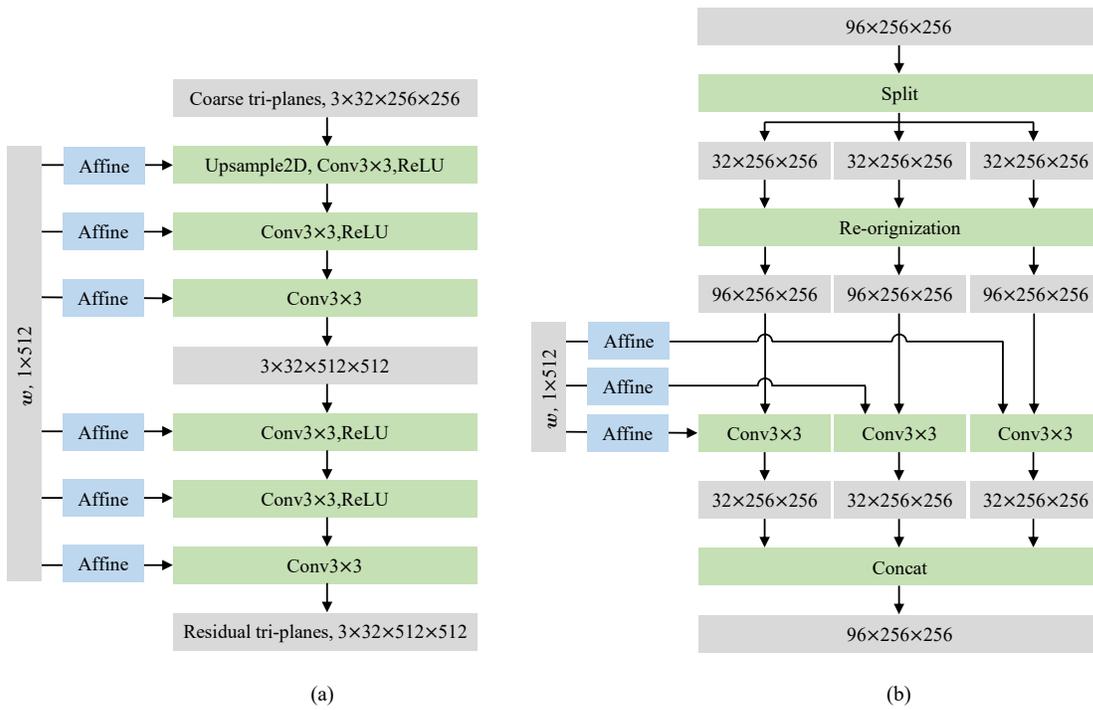
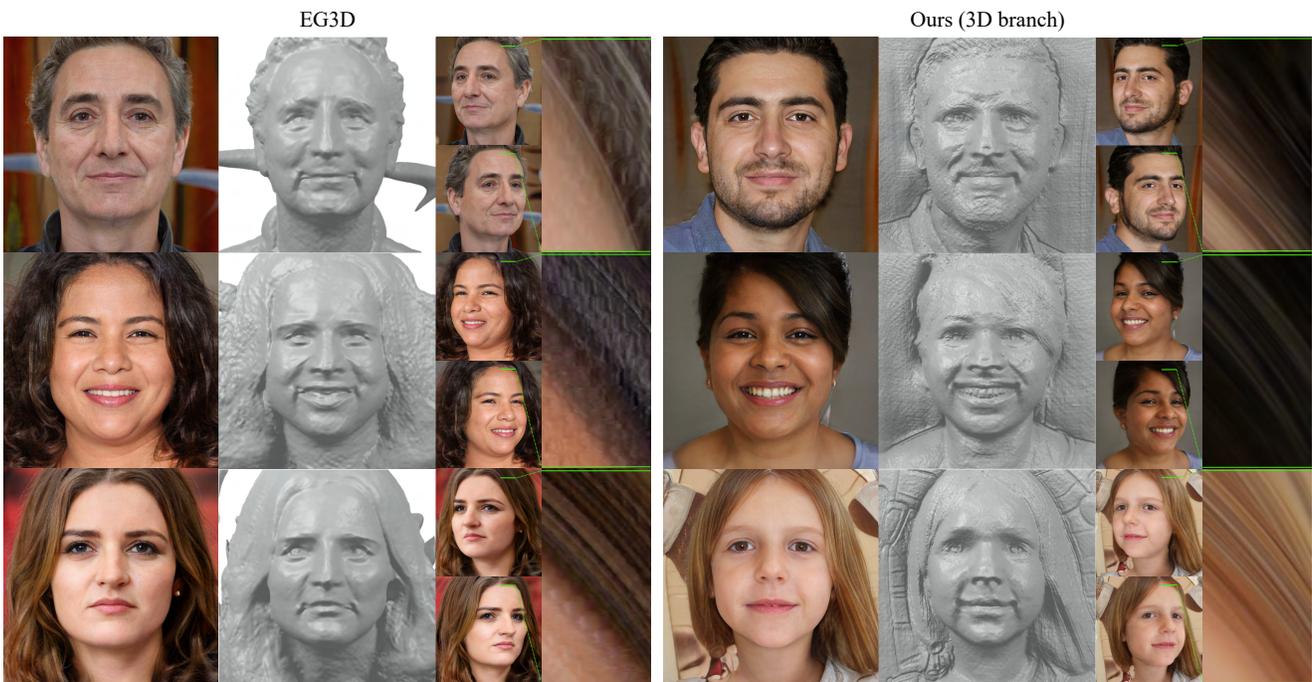Figure 12. Network designs. (a) 3D super-resolution module $\mathcal{S}^{3D}$. (b) 3D-aware block.



Figure 13. Comparison w/ EG3D [3]. Our method generates images with comparable quality to those of EG3D, while producing 3D geometries with finer details and multiview sequences with better 3D-consistency. Referring to the supplemental video for animations.
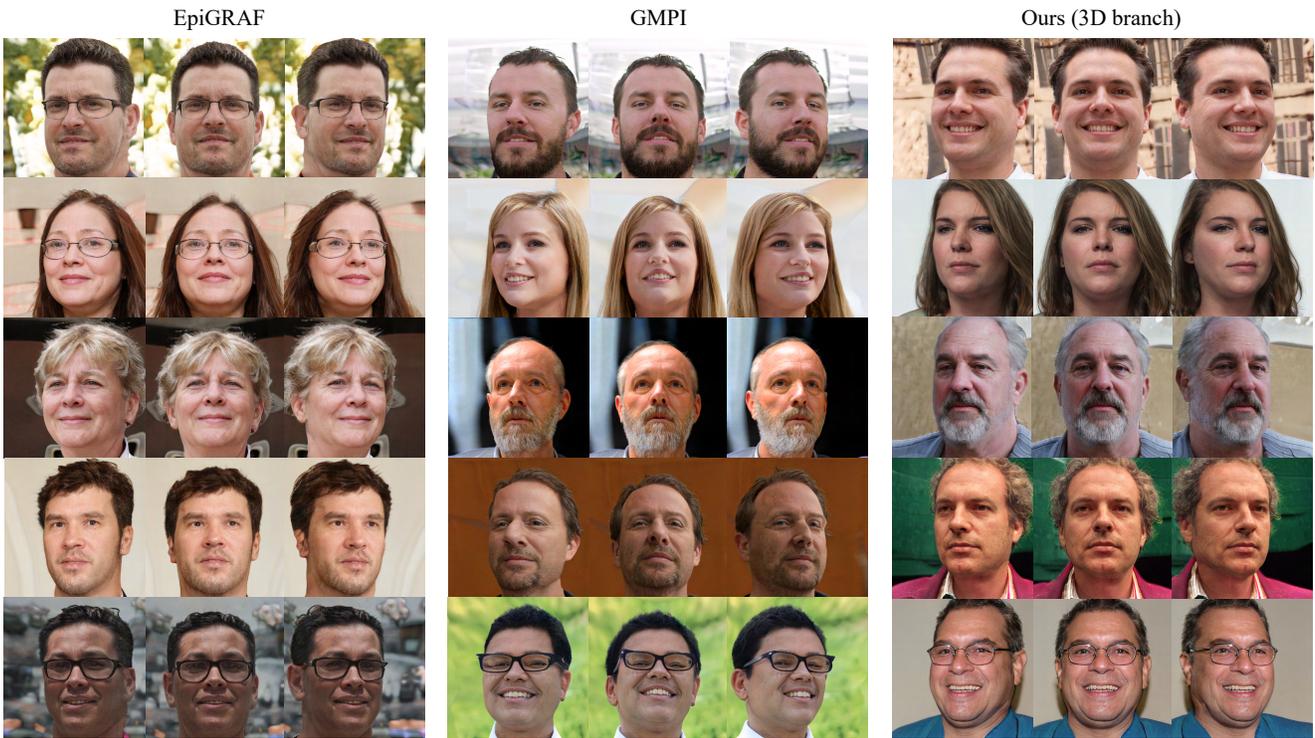
EpiGRAF            GMPI            Ours (3D branch)

Figure 14. Comparison w/ EpiGRAF [43] and GMPI [54] . Referring to the supplemental video for animations.

Figure 15. Our results on FFHQ dataset. Referring to the supplemental video for animations.

Figure 16. Our results on AFHQ-v2 Cats. Referring to the supplemental video for animations.