# Generative Multiplane Images: Making a 2D GAN 3D-Aware

Xiaoming Zhao[1,2*], Fangchang Ma[1], David Güera[1], Zhile Ren[1],
Alexander G. Schwing[2], Alex Colburn[1]

[1] Apple
[2] University of Illinois, Urbana-Champaign
https://github.com/apple/ml-gmpi

**Abstract.** What is really needed to make an existing 2D GAN 3D-aware? To answer this question, we modify a classical GAN, *i.e.*, Style-GANv2, as little as possible. We find that only two modifications are absolutely necessary: 1) a multiplane image style generator branch which produces a set of alpha maps conditioned on their depth; 2) a pose-conditioned discriminator. We refer to the generated output as a 'generative multiplane image' (GMPI) and emphasize that its renderings are not only high-quality but also guaranteed to be view-consistent, which makes GMPIs different from many prior works. Importantly, the number of alpha maps can be dynamically adjusted and can differ between training and inference, alleviating memory concerns and enabling fast training of GMPIs in less than half a day at a resolution of $1024^2$. Our findings are consistent across three challenging and common high-resolution datasets, including FFHQ, AFHQv2 and MetFaces.

**Keywords:** GANs, 3D-aware generation, multiplane images

## 1 Introduction

Generative adversarial networks (GANs) [22] have been remarkably successful at sampling novel images which look 'similar' to those from a given training dataset. Notably, impressive advances have been reported in recent years which improve quality and resolution of the generated images. Most of these advances focus on the setting where the output space of the generator and the given dataset are identical and, often, these outputs are either images or occasionally 3D volumes.

The latest literature, however, has focused on generating novel outputs that differ from the available training data. This includes methods that generate 3D geometry and the corresponding texture for one class of objects, *e.g.*, faces, while the given dataset only contains widely available single-view images [58, 53, 10, 55, 16, 9, 67, 66]. No multi-view images or 3D geometry are used to supervise the training of these 3D-aware GANs. To learn the 3D geometry from such a limited form of supervision, prior work typically combines 3D-aware inductive

---

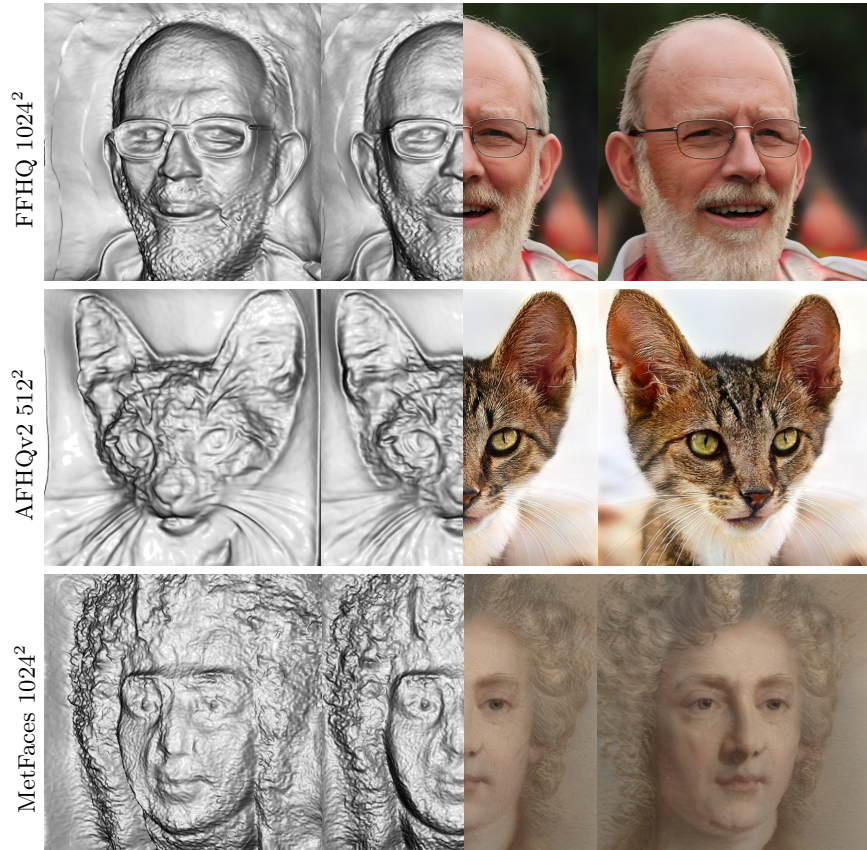[*] Work done as part of an internship at Apple.

Fig. 1: Making a 2D GAN 3D-aware with a minimal set of changes: on three datasets (FFHQ, AFHQv2 and MetFaces) we observe that it suffices to augment StyleGANv2 with an additional branch that generates alpha maps conditioned on depths, and to modify training by using a discriminator conditioned on pose.

biases such as a 3D voxel grid or an implicit representation with a rendering engine [58, 53, 10, 55, 16, 9, 67, 66].

Nevertheless, improving the quality of the results of these methods remains challenging: 3D-aware inductive biases are often memory-intensive explicit or implicit 3D volumes, and/or rendering is often computationally demanding, *e.g.*, involving a two-pass importance sampling in a 3D volume and a subsequent decoding of the obtained features. Moreover, lessons learned from 2D GANs are often not directly transferable because the generator output or even its entire structure has to be adjusted. This poses the question:

*"What is really needed to make an existing 2D GAN 3D-aware?"*

To answer this question, we aim to modify an existing 2D GAN as little as possible. Further, we aim for an efficient inference and training procedure. As a

starting point, we chose the widely used StyleGANv2 [35], which has the added benefit that many training checkpoints are publicly available.

More specifically, we develop a new generator branch for StyleGANv2 which yields a set of fronto-parallel alpha maps, in spirit similar to multiplane images (MPIs) [68]. As far as we know, we are the first to demonstrate that MPIs can be used as a scene representation for *unconditional* 3D-aware generative models. This new alpha branch is trained from scratch while the regular StyleGANv2 generator and discriminator are simultaneously fine-tuned. Combining the generated alpha maps with the single standard image output of StyleGANv2 in an end-to-end differentiable multiplane style rendering, we obtain a 3D-aware generation from different views while *guaranteeing view-consistency*. Although alpha maps have limited ability to handle occlusions, rendering is very efficient. Moreover, the number of alpha maps can be dynamically adjusted and can even differ between training and inference, alleviating memory concerns. We refer to the generated output of this method as a 'generative multiplane image' (GMPI).

To obtain alpha maps that exhibit an expected 3D structure, we find that only two adjustments of StyleGANv2 are really necessary: (a) the alpha map prediction of any plane in the MPI has to be conditioned on the plane's depth or a learnable token; (b) the discriminator has to be conditioned on camera poses. While these two adjustments seem intuitive in hindsight, it is still surprising that an alpha map with planes conditioned on their depth and use of camera pose information in the discriminator are sufficient inductive biases for 3D-awareness.

An additional inductive bias that improves the alpha maps is a 3D rendering that incorporates shading. Albeit useful, we didn't find this inductive bias to be necessary to obtain 3D-awareness. Moreover, we find that metrics designed for classic 2D GAN evaluation, *e.g.*, the Fréchet Inception Distance (FID) [28] and the Kernel Inception Distance (KID) [6] may lead to misleading results since they do not consider geometry.

In summary, our contribution is two-fold: 1) We are the first to study an MPI-like 3D-aware generative model trained with standard single-view 2D image datasets; 2) We find that conditioning the alpha planes on depth or a learnable token and the discriminator on camera pose are sufficient to make a 2D GAN 3D-aware. Other information provides improvements but is not strictly necessary. We study the aforementioned ways to encode 3D-aware inductive biases on three high-resolution datasets: FFHQ [34], AFHQv2 [12], and MetFaces [32]. Across all three datasets, as illustrated in Fig. 1, our findings are consistent.

## 2   Related Work and Background

In the following, we briefly review recent advances in classical and neural scene rendering, as well as the generation of 2D and 3D data with generative models. We then discuss the generation of 3D data using only 2D supervision. We also provide a brief review of single image reconstruction techniques before we review StyleGANv2 in greater detail.

**Scene Rendering.** Image-Based Rendering (IBR) is well-studied [11]. IBR 1) models a scene given a set of images; and 2) uses the scene model to render novel views. Methods can be grouped based on their use of scene geometry: explicit, implicit, or not using geometry. Texture-mapping methods use explicit geometry, whereas layered depth images (LDIs) [60], flow-based [11], lumigraphs [8], and tensor-based methods [4] use geometry implicitly. In contrast, light-field methods [40] do not rely on geometry. Hybrid methods [14] have also been studied. More recently, neural representations have been used in IBR, for instance, neural radiance fields (NeRFs) [47] and multiplane images (MPIs) [68, 62, 64, 21, 42]. In common to both is the goal to extract from a given set of images a volumetric representation of the observed scene. The volumetric representation in MPIs is discrete and permits fast rendering, whereas NeRFs use a continuous spatial representation.

These works differ from the proposed method in two main aspects. 1) The proposed method is unconditionally generative, *i.e.*, novel, never-before-seen scenes can be synthesized without requiring any color [64], depth, or semantic images [25]. In contrast, IBR methods and related recent advances focus on reconstructing a scene representation from a set of images. 2) The proposed method uses a collection of 'single-view images' from different scenes during training. In contrast, IBR methods use multiple viewpoints of a single scene for highly accurate reconstruction.

**2D Generative Models.** Generative adversarial networks (GANs) [22] and variational auto-encoders (VAEs) [37] significantly advanced modeling of probability distributions. In the early days, GANs were notably difficult to train whereas VAEs often produced blurry images. However, in the last decade, theoretical and practical understanding of these methods has significantly improved. New loss functions and other techniques have been proposed [41, 24, 38, 18, 13, 50, 5, 49, 44, 29, 57, 46, 3, 19, 2, 43, 63] to improve the stability of GAN optimization and to address mode-collapse, some theoretically founded and others empirically motivated. We follow the architectural design and techniques in StyleGANv2 [35], including exponential moving average (EMA), gradient penalties [46], and minibatch standard deviation [31].

**3D Generative Models from 2D Data.** 3D-aware image synthesis has gained attention recently. Sparse volume representations are used to generate photorealistic images based on given geometry input [26]. Many early approaches use voxel-based representations [69, 65, 52, 51, 27, 20], where scaling to higher resolutions is prohibitive due to the high memory footprint. Rendering at low-res followed by 2D CNN-based upsampling [54] has been proposed as a workaround, but it leads to view inconsistency. As an alternative, methods built on implicit functions, *e.g.*, NeRF, have been proposed [10, 58, 56]. However, their costly querying and sampling operations limit training efficiency and image resolutions. To generate high-resolution images, concurrently, EG3D [9], StyleNeRF [23], CIPS-3D [67], VolumeGAN [66], and StyleSDF [55] have been developed. Our work differs primarily in the choice of scene representation: EG3D uses a hybrid tri-plane representation while the others follow a NeRF-style implicit representa-

tion. In contrast, we study an MPI-like representation. In our experience, MPIs provide extremely fast rendering speed without incurring quality degradation. Most related to our work are GRAM [16] and LiftedGAN [61]. GRAM uses a NeRF-style scene representation and learns scene-specific isosurfaces. Queries of RGB and density happen on those isosurfaces. Although isosurfaces are conceptually similar to MPIs, the queries of this NeRF-style representation are expensive, limiting the image synthesis to low resolution at $256 \times 256$. LiftedGAN reconstructs the geometry of an image by distilling intermediate representations from a fixed StyleGANv2 to a separate 3D generator which produces a depth map and a transformation map in addition to an image. Different from our proposed approach, because of the transformation map, LiftedGAN is not strictly view-consistent. Moreover, the use of a single depth-map is less flexible.

**StyleGANv2 Revisited.** Since our method is built on StyleGANv2, we provide some background on its architecture. StyleGANv2 generates a square 2D image $C \in \mathbb{R}^{H \times H \times 3}$ by upsampling and accumulating intermediate GAN generation results from various resolutions. Formally, the image $C \triangleq C^H$ is obtained via

$$C^h = \begin{cases} \tilde{C}^h + \texttt{UpSample}_{\frac{h}{2} \to h}(C^{\frac{h}{2}}), & \text{if } h \in \mathcal{R} \setminus \{4\}, \\ \tilde{C}^4, & \text{if } h = 4, \end{cases} \qquad (1)$$

where $C^h \in \mathbb{R}^{h \times h \times 3}$ is the GAN image generation at resolution $h$ and $\tilde{C}^h$ is the generated residual at the same resolution. $\mathcal{R} = \{4, 8, \ldots, H\}$ is the set of available resolutions whose values are powers of 2. $\texttt{UpSample}_{\frac{h}{2} \to h}$ refers to the operation that upsamples from $\mathbb{R}^{\frac{h}{2} \times \frac{h}{2} \times 3}$ to $\mathbb{R}^{h \times h \times 3}$. The residual generation $\tilde{C}^h$ at resolution $h$ is generated with a single convolutional layer $f_{\texttt{ToRGB}}^h$, *i.e.*,

$$\tilde{C}^h = f_{\texttt{ToRGB}}^h(\mathcal{F}^h), \qquad (2)$$

where $\mathcal{F}^h \in \mathbb{R}^{h \times h \times \texttt{dim}_h}$ is the intermediate GAN feature representation at resolution $h$ with $\texttt{dim}_h$ channels. These intermediate GAN feature representations at all resolutions $\mathcal{R}$ are computed with a synthesis network $f_{\texttt{Syn}}$, *i.e.*,

$$\{\mathcal{F}^h : h \in \mathcal{R}\} = f_{\texttt{Syn}}(\boldsymbol{\omega}), \qquad (3)$$

where $\boldsymbol{\omega}$ is the style embedding. $\boldsymbol{\omega}$ is computed via the mapping network $f_{\texttt{Mapping}}$ which operates on the latent variable $\boldsymbol{z}$, *i.e.*, $\boldsymbol{\omega} = f_{\texttt{Mapping}}(\boldsymbol{z})$.

## 3    Generative Multiplane Images (GMPI)

Our goal is to adjust an existing GAN such that it generates images which are 3D-aware and view-consistent, *i.e.*, the image $I_{v_{\texttt{tgt}}}$ can illustrate the *exactly identical* generated object from different camera poses $v_{\texttt{tgt}}$. In order to achieve 3D-awareness and guaranteed view-consistency, different from existing prior work, we aim to augment an existing generative adversarial network, in our case StyleGANv2, as little as possible. For this, we modify the classical
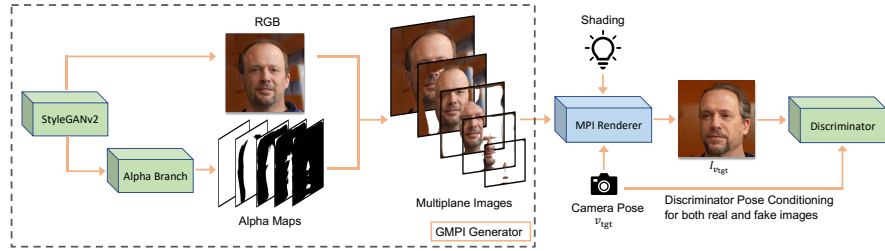
Fig. 2: Overview of the developed 3D-aware generative multiplane images (GMPI). We find that two components are necessary to make a 2D StyleGANv2 3D-aware: accompanying StyleGANv2 with an MPI branch (Sec. 3.2) and conditioning the discriminator on pose during training (Sec. 3.4). We find adding shading during training (Sec. 3.5) improves the generated geometry. Please also refer to Sec. 3.3 for details about the employed MPI rendering. Green blocks denote trainable components.

generator by adding an 'alpha branch' and incorporate a simple and efficient alpha composition rendering. Specifically, the 'alpha branch' produces alpha maps of a multiplane image while the alpha composition rendering step transforms generated alpha maps and generated image into the target view $I_{v_{\text{tgt}}}$ given a user-specified pose $v_{\text{tgt}}$. We refer to the output of the generator and the renderer as a 'generative multiplane image' (GMPI). To achieve 3D-awareness, we also find the pose conditioning of the discriminator to be absolutely necessary. Moreover, additional miscellaneous adjustments like the use of shading help to improve results. We first discuss the generator, specifically our alpha branch (Sec. 3.2) and our rendering (Sec. 3.3). Subsequently we detail the discriminator pose conditioning (Sec. 3.4) and the miscellaneous adjustments (Sec. 3.5).

### 3.1    Overview

An overview of our method to generate GMPIs is shown in Fig. 2. Its generator and the subsequent alpha-composition renderer produce an image $I_{v_{\text{tgt}}}$ illustrating the generated object from a user-specified pose $v_{\text{tgt}}$. Images produced for different poses are guaranteed to be view-consistent. The generator and rendering achieve 3D-awareness and guaranteed view consistency in two steps. First, a novel 'alpha branch' uses intermediate representations to produce a multiplane image representation $\mathcal{M}$ which contains alpha maps at various depths in addition to a single image. Importantly, to obtain proper 3D-awareness we find that it is necessary to condition alpha maps on their depth. We discuss the architecture of this alpha branch and its use of the plane depth in Sec. 3.2. Second, to guarantee view consistency, we employ a rendering step (Sec. 3.3). It converts the representation $\mathcal{M}$ obtained from the alpha branch into the view $I_{v_{\text{tgt}}}$, which shows the generated object from a user-specified pose $v_{\text{tgt}}$. During training, the generated image $I_{v_{\text{tgt}}}$ is compared to real images from a single-view dataset via a pose conditioned discriminator. We discuss the discriminator details in Sec. 3.4.
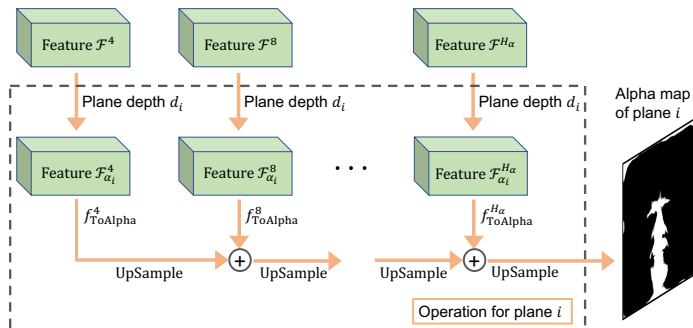
Fig. 3: The alpha branch proposed in Sec. 3.2. Here we show the generation of the alpha map of plane $i$. At each intermediate resolution $h \in \mathcal{R}_\alpha = \{4, 8, \ldots, H_\alpha\}$, we utilize the plane's depth $d_i$ to transform the feature from $\mathcal{F}^h$ (Eq. (3)) to $\mathcal{F}^h_{\alpha_i}$ (Eq. (8)). The final alpha map $\alpha_i$ is obtained by accumulating all intermediate results, which are generated by the single convolutional layer $f^h_{\texttt{ToAlpha}}$.

Finally, in Sec. 3.5, we highlight some miscellaneous adjustments which we found to help improve 3D-awareness while not being strictly necessary.

### 3.2   GMPI: StyleGANv2 with an Alpha Branch

Given an input latent code $\boldsymbol{z}$, our goal is to synthesize a multiplane image inspired representation $\mathcal{M}$ from which realistic and view-consistent 2D images can be rendered at different viewing angles. A classical multiplane image refers to a set of tuples $(C_i, \alpha_i, d_i)$ for $L$ fronto-parallel planes $i \in \{1, \ldots, L\}$. Within each tuple, $C_i \in \mathbb{R}^{H \times H \times 3}$ denotes the color texture for the $i^{\text{th}}$ plane and is assumed to be a square image of size $H \times H$, where $H$ is independent of the plane index $i$. Similarly, $\alpha_i \in [0,1]^{H \times H \times 1}$ and $d_i \in \mathbb{R}$ denote the alpha map and the depth of the corresponding plane, $i.e.$, its distance from a camera. $i = 1$ and $i = L$ denote the planes closest and farthest from the camera.

   As a significant simplification, we choose to reuse the same color texture image $C_i \triangleq C \forall i$ across all planes. $C$ is synthesized by the original StyleGANv2 structure as specified in Eq. (1). Consequently, the task of the generator $f_{\mathsf{G}}$ has been reduced from predicting an entire multiplane image representation to predicting a single color image and a set of alpha planes, $i.e.$,

$$\mathcal{M} \triangleq \{C, \{\alpha_1, \alpha_2, \ldots, \alpha_L\}\} = f_{\mathsf{G}}(\boldsymbol{z}, \{d_1, d_2, \ldots, d_L\}). \tag{4}$$

   For this, we propose a simple modification of the original StyleGANv2 network by adding an additional alpha branch, as illustrated in Fig. 3. For consistency we follow StyleGANv2's design: the alpha map $\alpha_i \triangleq \alpha_i^H$ of plane $i$ is obtained by upsampling and accumulating alpha maps at different resolutions. Notably, we do not generate all the way up to the highest resolution $H$, but instead use the final upsampling step

$$\alpha_i \triangleq \alpha_i^H = \texttt{UpSample}_{H_\alpha \to H}(\alpha_i^{H_\alpha}), \tag{5}$$

if $H_\alpha < H$. Here $H_\alpha \leq H$ refers to a possibly lower resolution. We will explain the reason for this design below. Following Eq. (1), we have

$$\alpha_i^h = \begin{cases} \tilde{\alpha}_i^h + \texttt{UpSample}_{\frac{h}{2} \to h}(\alpha_i^{\frac{h}{2}}), & \text{if } h \in \mathcal{R}_\alpha \setminus \{4\}, \\ \tilde{\alpha}_i^4, & \text{if } h = 4. \end{cases} \tag{6}$$

Here, $\alpha_i^h \in [0,1]^{h \times h \times 1}$ and $\tilde{\alpha}_i^h$ are the alpha map and the residual at resolution $h$ respectively. $\mathcal{R}_\alpha \triangleq \{4, 8, \dots, H_\alpha\} \subseteq \mathcal{R}$ denotes the set of alpha maps' intermediate resolutions. Inspired by Eq. (2), $\tilde{\alpha}_i^h$ is generated from an intermediate feature representation $\mathcal{F}_{\alpha_i}^h$ through a single convolutional layer $f_{\texttt{ToAlpha}}^h$:

$$\tilde{\alpha}_i^h = f_{\texttt{ToAlpha}}^h(\mathcal{F}_{\alpha_i}^h). \tag{7}$$

Note, $f_{\texttt{ToAlpha}}^h$ is shared across all planes $i$, while the input feature $\mathcal{F}_{\alpha_i}^h$ is plane-specific. Inspired by AdaIn [30], we construct this plane-aware feature as follows:

$$\mathcal{F}_{\alpha_i}^h = \frac{\mathcal{F}^h - \mu(\mathcal{F}^h)}{\sigma(\mathcal{F}^h)} + f_{\texttt{Embed}}(d_i, \boldsymbol{\omega}), \tag{8}$$

where $\mu(\mathcal{F}^h) \in \mathbb{R}^{\texttt{dim}_h}$ and $\sigma(\mathcal{F}^h) \in \mathbb{R}^{\texttt{dim}_h}$ are mean and standard deviation of the feature $\mathcal{F}^h \in \mathbb{R}^{h \times h \times \texttt{dim}_h}$ from Eq. (3). Meanwhile, $f_{\texttt{Embed}}$ uses the depth $d_i$ of plane $i$ and the style embedding $\boldsymbol{\omega}$ from Eq. (3) to compute a plane-specific embedding in the space of $\mathbb{R}^{\texttt{dim}_h}$. Note, our design of $\mathcal{F}_{\alpha_i}^h$ disentangles alpha map generation from a pre-determined number of planes as $f_{\texttt{ToAlpha}}^h$ operates on each plane individually. This provides the ability to use an arbitrary number of planes which helps to reduce artifacts during inference as we will show later.

Note, the plane specific feature $\mathcal{F}_{\alpha_i}^h \forall i$ will in total occupy $L$ times the memory used by the feature $\mathcal{F}^h$. This might be prohibitive if we were to generate these intermediate features up to a resolution of $H$. Therefore, we only use this feature until some lower resolution $H_\alpha$ in Eq. (5). We will show later that this design choice works well on real-world data.

### 3.3   Differentiable Rendering in GMPI

We obtain the desired image $I_{v_{\texttt{tgt}}}$ which illustrates the generated MPI representation $\mathcal{M} = \{C, \{\alpha_1, \dots, \alpha_L\}\}$ from the user-specified target view $v_{\texttt{tgt}}$ with an MPI renderer in two steps: 1) a warping step transforms the representation $\mathcal{M}$ from its canonical pose $v_{\texttt{cano}}$ to the target pose $v_{\texttt{tgt}}$; 2) a compositing step combines the planes into the desired image $I_{v_{\texttt{tgt}}}$. Importantly, both steps entail easy computations which are end-to-end differentiable such that they can be included into any generator. Please see Appendix A for details.

### 3.4   Discriminator Pose Conditioning

For the discriminator $f_{\texttt{D}}$ to encourage 3D-awareness of the generator, we find the conditioning of the discriminator on camera poses to be essential. Formally,

inspired by Miyato and Koyama [48], the final prediction of the discriminator is

$$\log P(y = \mathsf{real}|I_{v_{\mathsf{tgt}}}, v_{\mathsf{tgt}}) \propto \mathsf{Normalize}(f^{\mathsf{D}}_{\mathsf{Embed}}(v_{\mathsf{tgt}})) \cdot f_{\mathsf{D}}(I_{v_{\mathsf{tgt}}})^{\top}, \qquad (9)$$

where $P(y = \mathsf{real}|I_{v_{\mathsf{tgt}}}, v_{\mathsf{tgt}})$ denotes the probability that image $I_{v_{\mathsf{tgt}}}$ from the camera pose $v_{\mathsf{tgt}}$ is real. $v_{\mathsf{tgt}} \in \mathbb{R}^{16}$ denotes the flattened extrinsic matrix of the camera pose. $f^{\mathsf{D}}_{\mathsf{Embed}} : \mathbb{R}^{16} \mapsto \mathbb{R}^{16}$ denotes an embedding function, while $\mathsf{Normalize}(\cdot)$ results in a zero-mean, unit-variance embedding. $f_{\mathsf{D}}(\cdot) \in \mathbb{R}^{16}$ denotes the feature computed by the discriminator $\mathsf{D}$. For real images of faces, $e.g.$, humans and cats, $v_{\mathsf{tgt}}$ can be estimated via off-the-shelf tools [17, 1].[3]

### 3.5   Miscellaneous Adjustment: Shading-guided Training

Inspired by [56], we incorporate shading into the rendering process introduced in Sec. 3.3. Conceptually, shading will amplify artifacts of the generated geometry that might be hidden by texture, encouraging the alpha branch to produce better results. To achieve this, we adjust the RGB component $C \in \mathbb{R}^{H \times W \times 3}$ via

$$\widehat{C} = C \cdot (k_a + k_d\, \boldsymbol{l} \cdot \mathsf{N}\,(D_{v_{\mathsf{cano}}}))\,, \qquad (10)$$

where $k_a$ and $k_d$ are coefficients for ambient and diffuse shading, $\boldsymbol{l} \in \mathbb{R}^3$ indicates lighting direction, and $\mathsf{N}\,(D_{v_{\mathsf{cano}}}) \in \mathbb{R}^{H \times W \times 3}$ denotes the normal map computed from the depth map $D_{v_{\mathsf{cano}}}$ (obtained by using the canonical alpha maps $\alpha_i$, see Eq. (S4) in Appendix A). We find shading to slightly improve results while not being required for 3D-awareness. Implementation details are in Appendix C.

### 3.6   Training

**Model structure.** The trainable components of GMPI are $f^h_{\mathsf{ToRGB}}$ (Eq. (2)), $f_{\mathsf{Syn}}$ and $f_{\mathsf{Mapping}}$ (Eq. (3)), $f^h_{\mathsf{ToAlpha}}$ (Eq. (7)), $f_{\mathsf{Embed}}$ (Eq. (8)), and $f_{\mathsf{D}}$ and $f^{\mathsf{D}}_{\mathsf{Embed}}$ (Eq. (9)). Please check the appendix for implementation details. Our alpha maps (Eq. (4)) are equally placed in the disparity (inverse depth) space during training and we set $H_\alpha = 256$ (Eq. (5)).

**Initialization.** For any training, we initialize weights from the officially-released StyleGANv2 checkpoints.[4] This enables a much faster training as we will show.

**Loss.** We use $\theta$ to subsume all trainable parameters. Our training loss consists of a non-saturating GAN loss with R1 penalties [46], $i.e.$,

$$\mathcal{L}_\theta = \mathbb{E}_{I_{v_{\mathsf{tgt}}}, v_{\mathsf{tgt}}} \left[ f(\log P(y = \mathsf{real}|I_{v_{\mathsf{tgt}}}, v_{\mathsf{tgt}})) \right]$$
$$+ \mathbb{E}_{I, v_{\mathsf{tgt}}} \left[ f(\log P(y = \mathsf{real}|I, v_{\mathsf{tgt}})) + \lambda |\nabla_I \log P(y = \mathsf{real}|I, v_{\mathsf{tgt}})|^2 \right]. \quad (11)$$

Here $f(x) = -\log(1 + \exp(-x))$, $\lambda = 10.0$ in all studies, $I$ refers to real images and $v_{\mathsf{tgt}}$ denotes the corresponding observer's pose information.

---

[3] Concurrently, EG3D [9] also finds that pose conditioning of the discriminator is required for their tri-plane representation to produce 3D-aware results, corroborating that this form of inductive bias is indeed necessary.

[4] https://github.com/NVlabs/stylegan2-ada-pytorch

## 4    Experiments

We analyze GMPI on three datasets (FFHQ, AFHQv2 and MetFaces) and across a variety of resolutions. We first provide details regarding the three datasets before discussing evaluation metrics and quantitative as well as qualitative results.

### 4.1    Datasets

**FFHQ.** The FFHQ dataset [34] consists of 70,000 high-quality images showing real-world human faces from different angles at a resolution of $1024 \times 1024$. To obtain the pose of a face we use an off-the-shelf pose estimator [17].
**AFHQv2.** The AFHQv2-Cats dataset [12, 33] consists of 5,065 images showing faces of cats from different views at a resolution of $512 \times 512$. We augment the dataset by horizontal flips and obtain the pose of a cat's face via an off-the-shelf cat face landmark predictor [1] and OpenCV's perspective-n-point algorithm.
**MetFaces.** The MetFaces dataset [32] consists of 1,336 images showing high-quality faces extracted from the collection of the Metropolitan Museum of Art.[5] To augment the dataset we again use horizontal flips. To obtain the pose of a face we use an off-the-shelf pose estimator [17].

### 4.2    Evaluation Metrics

We follow prior work and assess the obtained results using common metrics:
**2D GAN metrics.** We report Fréchet Inception Distance (FID) [28] and Kernel Inception Distance (KID) [6], computed by using 50k artificially generated images that were rendered from random poses and 1) 50k real images for FFHQ; 2) all 10,130 real images in the $x$-flip-augmented dataset for AFHQv2-Cat [12].
**Identity similarity (ID).** Following EG3D [9], we also evaluate the level of facial identity preservation. Concretely, we first generate 1024 MPI-like representations $\mathcal{M}$. For each representation we then compute the identity cosine similarity between two views rendered from random poses using ArcFace [15, 59].
**Depth accuracy (Depth).** Similar to [9, 61], we also assess geometry and depth accuracy. For this we utilize a pre-trained face reconstruction model [17] to provide facial area mask and pseudo ground-truth depth map $\widehat{D}_{v_{\mathrm{tgt}}}$. We report the MSE error between our rendered depth $D_{v_{\mathrm{tgt}}}$ (Eq. (S4)) and $\widehat{D}_{v_{\mathrm{tgt}}}$ on areas constrained by the face mask. Note, following prior work, we normalize both depth maps to zero-mean and unit-variance. The result is obtained by averaging over 1024 representations $\mathcal{M}$.
**Pose accuracy (Pose).** Following [9, 61], we also study the 3D geometry's pose accuracy. Specifically, for each MPI, we utilize a pose predictor [17] to estimate the yaw, pitch, and roll of a rendered image. The predicted pose is then compared to the pose used for rendering via the MSE. The reported result is averaged over 1024 representations $\mathcal{M}$. Notably but in hindsight expected, 2D metrics, as well as ID, lack the ability to capture 3D errors which we will show next.

---

[5] https://metmuseum.github.io/

Table 1: **Speed comparison.** '−' indicates that corresponding papers do not report this result. **1) Training:** EG3D, GRAM, StyleNeRF and GMPI report training time when using 8 Tesla V100 GPUs. pi-GAN uses two RTX 6000 GPUs or a single RTX 8000 GPU. For GMPI, the results reported in this paper come from 3/5/11-hour training for a resolution of $256^2/512^2/1024^2$ with initialization from official pretrained-checkpoints. **2) Inference:** we measure frames-per-second (FPS) for each model. GMPI uses 96 planes. GRAM reports speed by utilizing a specified mesh rasterizer [39] while the others use pure forward passes. EG3D uses an RTX 3090 GPU. GRAM and StyleNeRF use a Tesla V100 GPU, which we also utilize to run GIRAFFE, pi-GAN, LiftedGAN, and our GMPI. We observe: GMPI is quick to train and renders the fastest among all approaches.

| | Res. | Unit | GIRAFFE [54] | pi-GAN [10] | LiftedGAN [61] | EG3D† [9] | GRAM† [16] | StyleNeRF† [23] | GMPI |
|---|---|---|---|---|---|---|---|---|---|
| Train - | | Time↓ | − | 56h | − | 8.5d | 3-7d | 3d | **3/5/11h** |
| Inference | $256^2$ | FPS↑ | 250 | 1.63 | 25 | 36 | 180 | 16 | **328**$^*$ |
| | $512^2$ | FPS↑ | − | 0.41 | − | 35 | − | 14 | **83.5**$^*$ |
| | $1024^2$ | FPS↑ | − | 0.10 | − | − | − | 11 | **19.4**$^*$ |

† We quote results from their papers.

$^*$ We report the rendering speed. GMPI, different from StyleNeRF, pi-GAN, and LiftedGAN, only needs a single forward pass to generate the scene representation. Further rendering doesn't involve the generator. We hence follow EG3D to report inference FPS without forward pass time that is 82.34 ms, 99.97 ms, and 115.10 ms for 96 planes at $256^2$, $512^2$, and $1024^2$ on a V100 GPU.

### 4.3  Results

We provide speed comparison in Tab. 1 and a quantitative evaluation in Tab. 2. With faster training, GMPI achieves on-par or better performance than start-of-the-art when evaluating on $256^2$ images and can generate high-resolution results up to $1024^2$ which most baselines fail to produce. Specifically, GMPI results on resolutions $256^2$, $512^2$, and $1024^2$ are reported after 3/5/11-hours of training. Note, the pretrained StyleGANv2 initialization for FFHQ (see Sec. 3.6) requires training for 1d 11h ($256^2$), 2d 22h ($512^2$), and 6d 03h ($1024^2$) with 8 Tesla V100 GPUs respectively, as reported in the official repo.[4] In contrast, EG3D, GRAM, and StyleNeRF require training of at least three days. At a resolution of $256^2$, 1) GMPI outperforms GIRAFFE, pi-GAN, LiftedGAN, and GRAM on FID/KID while outperforming StyleSDF on FID; 2) GMPI demonstrates better identity similarity (ID) than GIRAFFE, pi-GAN, and LiftedGAN; 3) GMPI outperforms GIRAFFE regarding depth; 4) GMPI performs best among all baselines on pose accuracy. Overall, GMPI demonstrates that it is a flexible architecture which achieves 3D-awareness with an affordable training time.

### 4.4  Ablation Studies

In order to show the effects of various design choices, we run ablation studies and selectively drop out the discriminator conditioned on pose (DPC), the plane-specific feature $\mathcal{F}_{\alpha_i}^h$, and shading. The following quantitative (Tab. 3) and qualitative (Fig. 4 and Fig. 5) studies were run using a resolution of $512^2$ to set baseline metrics, and answer the following questions:

Table 2: **Quantitative results.** GMPI uses 96 planes and no truncation trick [34, 7] is applied. '–' indicates that the corresponding paper does not report this result. KID is reported in KID×100. For GIRAFFE, pi-GAN, and LiftedGAN, we quote results from [9]. Compared to baselines, GMPI achieves on-par or even better performance despite much faster training as reported in Tab. 1. Moreover, GMPI can produce high-resolution images which most baselines fail to achieve, demonstrating the flexibility.

| | | | FFHQ | | | | | AFHQv2-Cat | |
|---|---|---|---|---|---|---|---|---|---|
| | | | FID↓ | KID↓ | ID↑ | Depth↓ | Pose↓ | FID↓ | KID↓ |
| $256^2$ | 1 | GIRAFFE [54] | 31.5 | 1.992 | 0.64 | 0.94 | 0.089 | 16.1 | 2.723 |
| | 2 | pi-GAN $128^2$ [10] | 29.9 | 3.573 | 0.67 | 0.44 | 0.021 | 16.0 | 1.492 |
| | 3 | LiftedGAN [61] | 29.8 | – | 0.58 | 0.40 | 0.023 | – | – |
| | 4 | GRAM[†] [16] | 29.8 | 1.160 | – | – | – | – | – |
| | 5 | StyleSDF[†] [55] | 11.5 | 0.265 | – | – | – | 12.8[*] | 0.447[*] |
| | 6 | StyleNeRF[†] [23] | 8.00 | 0.370 | – | – | – | 14.0[*] | 0.350[*] |
| | 7 | CIPS-3D[†] [67] | 6.97 | 0.287 | – | – | – | – | – |
| | 8 | EG3D[†] [9] | 4.80 | 0.149 | 0.76 | 0.31 | 0.005 | 3.88 | 0.091 |
| | 9 | **GMPI** | 11.4 | 0.738 | 0.70 | 0.53 | 0.004 | n/a[§] | n/a[§] |
| $512^2$ | 10 | EG3D[†] [9] | 4.70 | 0.132 | 0.77 | 0.39 | 0.005 | 2.77 | 0.041 |
| | 11 | StyleNeRF[†] [23] | 7.80 | 0.220 | – | – | – | 13.2[*] | 0.360[*] |
| | 12 | **GMPI** | 8.29 | 0.454 | 0.74 | 0.46 | 0.006 | 7.79 | 0.474 |
| $1024^2$ | 13 | CIPS-3D[†] [67] | 12.3 | 0.774 | – | – | – | – | – |
| | 14 | StyleNeRF[†] [23] | 8.10 | 0.240 | – | – | – | – | – |
| | 15 | **GMPI** | 7.50 | 0.407 | 0.75 | 0.54 | 0.007 | n/a[§] | n/a[§] |

[†] We quote results from their papers.
[*] Performance is reported on the whole AFHQv2 dataset instead of only cats.
[§] No GMPI results on AFHQv2 for Row 9 and 13 since there are no available pre-trained StyleGANv2 checkpoints for AFHQv2 with a resolution of $256^2$ or $1024^2$.

**Baseline condition** (Tab. 3 row a) – a pre-trained StyleGANv2 texture without transparency ($\alpha_i \forall i$ is set to 1). This representation is not truly 3D, as it is rendered as a textured block in space. It is important to note that the commonly used metrics for 2D GAN evaluation, *i.e.*, FID, KID, as well as ID yield good results. They are hence not sensitive to 3D structure. In contrast, the depth and pose metrics do capture the lack of 3D structure in the scene.

**Can a naïve MPI generator learn 3D?** (Tab. 3 row b) – a generator which only uses $\mathcal{F}^h$ (Eq. (2)) without being trained with pose-conditioned discriminator, plane-specific features, or a shading loss. The Depth (2.190 *vs.* 2.043) and Pose (0.062 *vs.* 0.060) metrics show that this design fails, yielding results similar to the baseline (row a).

**Does `DPC` alone enable 3D-awareness?** (Tab. 3 row c) – we now use discriminator pose conditioning (DPC) while ignoring plane-specific features $\mathcal{F}^h_{\alpha_i}$. This design improves the Depth (0.501 *vs.* 2.190) and Pose (0.006 *vs.* 0.062) metrics. However upon closer inspection of the alpha maps shown in Fig. 5a we observe that the generator fails to produce realistic structures. The reason that this design performs well on Depth and Pose metrics is primarily due to these two metrics only evaluating the surface while not considering the whole volume.

Table 3: **Ablation studies.** We evaluate at a resolution of $512^2$. DPC refers to discriminator pose conditioning (Sec. 3.4), $\mathcal{F}^h_{\alpha_i}$ refers to the plane-specific feature introduced in Eq. (8), and Shading indicates the shading-guided training discussed in Sec. 3.5. #planes denotes the number of planes we used during evaluation. Note, rows (b) and (c) can only use 32 planes during training and inference. Therefore, to make the comparison fair, all ablations use 32 planes. We provide an additional 96-plane result for our full model, which is used in Tab. 2. Please check Sec. 4.4 and Fig. 4 for a detailed discussion regarding this table as we find that 2D metrics, *e.g.*, FID/KID as well as ID can be misleading when evaluating 3D generative models.

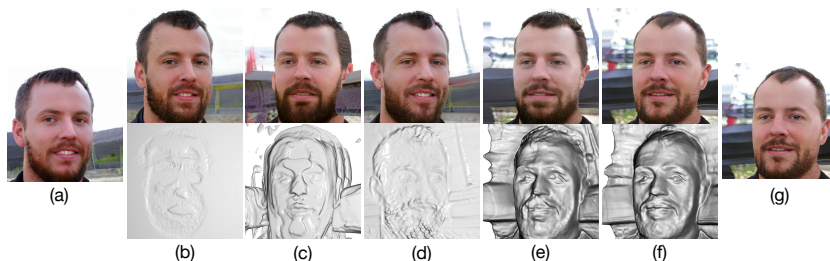| | #planes | DPC | $\mathcal{F}^h_{\alpha_i}$ | Shading | FFHQ | | | | | AFHQv2-Cat | |
| | | | | | FID↓ | KID↓ | ID↑ | Depth↓ | Pose↓ | FID↓ | KID↓ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| (a) | 32 | | | | 6.64 | 0.368 | **0.91** | 2.043 | 0.060 | 4.29 | 0.199 |
| (b) | 32 | | | | 4.61 | 0.167 | 0.89 | 2.190 | 0.062 | 4.00 | 0.166 |
| (c) | 32 | ✓ | | | 7.98 | 0.347 | 0.75 | 0.501 | **0.006** | 7.13 | 0.385 |
| (d) | 32 | | ✓ | | **4.35** | **0.150** | 0.89 | 2.140 | 0.061 | 3.70 | 0.132 |
| (e) | 32 | ✓ | ✓ | | 7.19 | 0.313 | 0.73 | 0.462 | **0.006** | 7.54 | 0.433 |
| (f) | 32 | ✓ | ✓ | ✓ | 7.40 | 0.337 | 0.74 | **0.457** | **0.006** | 7.93 | 0.489 |
| (g) | 96 | ✓ | ✓ | ✓ | 8.29 | 0.454 | 0.74 | **0.457** | **0.006** | 7.79 | 0.474 |



Fig. 4: The images in this figure correspond the the ablation studies in Tab. 3. Panels (a)-(g) correspond to Tab. 3's rows (a)-(g). All results are generated from the same latent code $z$. The face is rendered with a camera positioned to the right of the subject, *i.e.*, the rendered face should look to the left of the viewer. (a) is the 2D image produced by the pre-trained StyleGANv2. Note how GMPI becomes 3D-aware in (e), and generates geometry and texture that is occluded in the pre-trained StyleGANv2 image. (e) *vs*. (f): shading-guided training (Sec. 3.5) alleviates geometric artifacts such as concavities in the forehead. (f) *vs*. (g): the ability to use more planes during inference (Sec. 3.2) reduces "stair step" artifacts, visible on the cheek and the ear.

**Does $\mathcal{F}^h_{\alpha_i}$ alone enable 3D-awareness?** (Tab. 3 row d) – we use only the plane-specific feature $\mathcal{F}^h_{\alpha_i}$ without DPC. Geometry related metrics (Depth and Pose) perform poorly, which is corroborated by Fig. 4 (d), indicating an inability to model 3D information.

**Do DPC and $\mathcal{F}^h_{\alpha_i}$ enable 3D-awareness?** (Tab. 3 row e) – we combine both plane-specific features $\mathcal{F}^h_{\alpha_i}$ and DPC. This produce good values for Depth and Pose metrics. Visualizations in Fig. 4e, and Fig. 5b verify that GMPI successfully generates 3D-aware content. However, the FID/KID, as well as ID values are generally worse than non-3D-aware generators (row a-e).
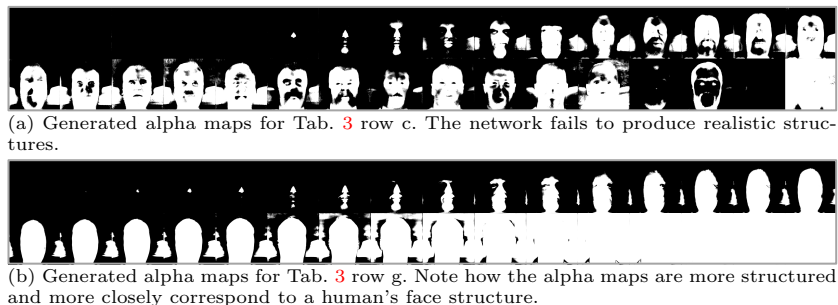
(a) Generated alpha maps for Tab. 3 row c. The network fails to produce realistic structures.



(b) Generated alpha maps for Tab. 3 row g. Note how the alpha maps are more structured and more closely correspond to a human's face structure.

Fig. 5: **Qualitative results accompanying Tab. 3.** For alpha maps, from top to bottom, left to right, we show $\alpha_1$ to $\alpha_L$ respectively (Eq. (4)): the whiter, the denser the occupancy. Grey boundaries are added for illustration.

**Does shading improve 3D-awareness?** (Tab. 3 row f) – we use plane-specific features $\mathcal{F}^h_{\alpha_i}$, DPC and shading loss. As discussed in Sec. 3.1, DPC and $\mathcal{F}^h_{\alpha_i}$ are sufficient to make a 2D GAN 3D-aware. However, inspecting the geometry reveals artifacts such as the concave forehead in Fig. 4e. Shading-guided rendering tends to alleviate these issues (Fig. 4e *vs.* f) while not harming the quantitative results (Tab. 3's row e *vs.* f).

**Can we reduce aliasing artifacts?** (Tab. 3 row g) – we use plane-specific features $\mathcal{F}^h_{\alpha_i}$, DPC, shading loss and 96 planes. Due to the formulation of $\mathcal{F}^h_{\alpha_i}$, we can generate an arbitrary number of planes during inference, which helps avoid "stair step" artifacts that can be observed in Fig. 4f *vs.* g.

**Qualitative results** are shown in Fig. 1. Please see the appendix for more.

## 5   Conclusion

To identify what is really needed to make a 2D GAN 3D-aware we develop generative multiplane images (GMPIs) which *guarantee view-consistency*. GMPIs show that a StyleGANv2 can be made 3D-aware by 1) adding a multiplane image style branch which generates a set of alpha maps conditioned on their depth in addition to a single image, both of which are used for rendering via an end-to-end differentiable warping and alpha compositing; and by 2) ensuring that the discriminator is conditioned on the pose. We also identify shortcomings of classical evaluation metrics used for 2D image generation. We hope that the simplicity of GMPIs inspires future work to fix limitations such as occlusion reasoning.

# References

1. Cat hipsterizer. https://github.com/kairess/cat_hipsterizer (2022), accessed: 2022-03-06 9, 10
2. Aneja, J., Schwing, A.G., Kautz, J., Vahdat, A.: A contrastive learning approach for training variational autoencoder priors. ICLR (2020) 4
3. Arjovsky, M., Chintala, S., Bottou, L.: Wasserstein GAN. ArXiv (2017) 4
4. Avidan, S., Shashua, A.: Novel view synthesis in tensor space. In: CVPR (1997) 4
5. Berthelot, D., Schumm, T., Metz, L.: Began: Boundary equilibrium generative adversarial networks. In: arXiv preprint arXiv:1703.10717 (2017) 4
6. Binkowski, M., Sutherland, D.J., Arbel, M., Gretton, A.: Demystifying mmd GANs. ICLR (2018) 3, 10
7. Brock, A., Donahue, J., Simonyan, K.: Large scale gan training for high fidelity natural image synthesis. ICLR (2019) 12, 21
8. Buehler, C., Bosse, M., McMillan, L., Gortler, S., Cohen, M.: Unstructured Lumigraph Rendering. SIGGRAPH (2001) 4
9. Chan, E.R., Lin, C.Z., Chan, M.A., Nagano, K., Pan, B., Mello, S.D., Gallo, O., Guibas, L., Tremblay, J., Khamis, S., Karras, T., Wetzstein, G.: Efficient Geometry-aware 3D Generative Adversarial Networks. In: CVPR (2022) 1, 2, 4, 9, 10, 11, 12, 24
10. Chan, E.R., Monteiro, M., Kellnhofer, P., Wu, J., Wetzstein, G.: pi-GAN: Periodic implicit generative adversarial networks for 3D-aware image synthesis. In: CVPR (2021) 1, 2, 4, 11, 12
11. Chen, S.E., Williams, L.: View Interpolation for Image Synthesis. SIGGRAPH (1993) 4
12. Choi, Y., Uh, Y., Yoo, J., Ha, J.W.: Stargan v2: Diverse image synthesis for multiple domains. In: CVPR (2020) 3, 10
13. Cully, R.W.A., Chang, H.J., Demiris, Y.: MAGAN: Margin adaptation for generative adversarial networks. ArXiv (2017) 4
14. Debevec, P.E., Taylor, C.J., Malik, J.: Modeling and Rendering Architecture from Photographs: A hybrid geometry- and image-based approach. SIGGRAPH (1996) 4
15. Deng, J., Guo, J., Zafeiriou, S.: ArcFace: Additive Angular Margin Loss for Deep Face Recognition. CVPR (2019) 10
16. Deng, Y., Yang, J., Xiang, J., Tong, X.: GRAM: Generative Radiance Manifolds for 3D-Aware Image Generation. CVPR (2022) 1, 2, 5, 11, 12
17. Deng, Y., Yang, J., Xu, S., Chen, D., Jia, Y., Tong, X.: Accurate 3D Face Reconstruction With Weakly-Supervised Learning: From Single Image to Image Set. CVPRW (2019) 9, 10, 23
18. Deshpande, I., Zhang, Z., Schwing, A.G.: Generative Modeling using the Sliced Wasserstein Distance. In: CVPR (2018) 4
19. Deshpande, I., Hu, Y.T., Sun, R., Pyrros, A., Siddiqui, N., Koyejo, O., Zhao, Z., Forsyth, D.A., Schwing, A.G.: Max-Sliced Wasserstein Distance and Its Use for GANs. CVPR (2019) 4
20. Gadelha, M., Maji, S., Wang, R.: 3D shape induction from 2D views of multiple objects. In: 3DV (2017) 4
21. Ghosh, S., Lv, Z., Matsuda, N., Xiao, L., Berkovich, A., Cossairt, O.: Liveview: Dynamic target-centered mpi for view synthesis. arXiv preprint arXiv:2107.05113 (2021) 4

22. Goodfellow, I.J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A.C., Bengio, Y.: Generative adversarial nets. In: NeurIPS (2014) 1, 4
23. Gu, J., Liu, L., Wang, P., Theobalt, C.: StyleNeRF: A style-based 3D-aware generator for high-resolution image synthesis. ICLR (2022) 4, 11, 12
24. Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., Courville, A.: Improved training of wasserstein GANs. In: NeurIPS (2017) 4
25. Habtegebrial, T., Jampani, V., Gallo, O., Stricker, D.: Generative view synthesis: From single-view semantics to novel-view images. NeurIPS (2020) 4
26. Hao, Z., Mallya, A., Belongie, S., Liu, M.Y.: GANcraft: Unsupervised 3D neural rendering of minecraft worlds. In: ICCV (2021) 4
27. Henzler, P., Mitra, N.J., Ritschel, T.: Escaping Plato's cave: 3D shape from adversarial rendering. In: ICCV (2019) 4
28. Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S.: GANs trained by a two time-scale update rule converge to a local nash equilibrium. In: NeurIPS (2017) 3, 10
29. Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S.: GANs trained by a two time-scale update rule converge to a local nash equilibrium. In: NeurIPS (2017) 4
30. Huang, X., Belongie, S.J.: Arbitrary style transfer in real-time with adaptive instance normalization. ICCV (2017) 8
31. Karras, T., Aila, T., Laine, S., Lehtinen, J.: Progressive growing of GANs for improved quality, stability, and variation. ICLR (2018) 4, 26, 27
32. Karras, T., Aittala, M., Hellsten, J., Laine, S., Lehtinen, J., Aila, T.: Training generative adversarial networks with limited data. NeurIPS (2020) 3, 10, 26
33. Karras, T., Aittala, M., Laine, S., Härkönen, E., Hellsten, J., Lehtinen, J., Aila, T.: Alias-free generative adversarial networks. In: NeurIPS (2021) 10
34. Karras, T., Laine, S., Aila, T.: A style-based generator architecture for generative adversarial networks. CVPR (2019) 3, 10, 12, 21, 28, 29, 30
35. Karras, T., Laine, S., Aittala, M., Hellsten, J., Lehtinen, J., Aila, T.: Analyzing and improving the image quality of stylegan. CVPR (2020) 3, 4, 27
36. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. Arxiv (2015) 26
37. Kingma, D.P., Welling, M.: Auto-encoding variational bayes. ICLR (2014) 4
38. Kolouri, S., Rohde, G.K., Hoffman, H.: Sliced wasserstein distance for learning gaussian mixture models. In: CVPR (2018) 4
39. Laine, S., Hellsten, J., Karras, T., Seol, Y., Lehtinen, J., Aila, T.: Modular primitives for high-performance differentiable rendering. TOG (2020) 11
40. Levoy, M., Hanrahan, P.: Light field rendering. SIGGRAPH (1996) 4
41. Li, C.L., Chang, W.C., Cheng, Y., Yang, Y., Póczos, B.: MMD GAN: Towards deeper understanding of moment matching network. In: NeurIPS (2017) 4
42. Li, J., Feng, Z., She, Q., Ding, H., Wang, C., Lee, G.H.: Mine: Towards continuous depth mpi with nerf for novel view synthesis. In: ICCV (2021) 4
43. Li, Y., Schwing, A.G., Wang, K.C., Zemel, R.: Dualing GANs. In: NeurIPS (2017) 4
44. Lin, Z., Khetan, A., Fanti, G., Oh, S.: Pacgan: The power of two samples in generative adversarial networks. In: NeurIPS (2018) 4
45. Lorensen, W.E., Cline, H.E.: Marching cubes: A high resolution 3d surface construction algorithm. ACM siggraph computer graphics (1987) 27
46. Mescheder, L.M., Geiger, A., Nowozin, S.: Which training methods for GANs do actually converge? In: ICML (2018) 4, 9, 26, 27

47. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis. In: ECCV (2020) 4

48. Miyato, T., Koyama, M.: cgans with projection discriminator. ICLR (2018) 9

49. Mroueh, Y., Sercu, T.: Fisher gan. In: NeurIPS (2017) 4

50. Mroueh, Y., Sercu, T., Goel, V.: McGan: Mean and covariance feature matching GAN. ArXiv (2017) 4

51. Nguyen-Phuoc, T., Li, C., Theis, L., Richardt, C., Yang, Y.L.: HoloGAN: Unsupervised learning of 3D representations from natural images. In: ICCV (2019) 4

52. Nguyen-Phuoc, T., Richardt, C., Mai, L., Yang, Y.L., Mitra, N.J.: BlockGAN: Learning 3D object-aware scene representations from unlabelled images. In: NeurIPS (2020) 4

53. Niemeyer, M., Geiger, A.: Giraffe: Representing scenes as compositional generative neural feature fields. CVPR (2021) 1, 2

54. Niemeyer, M., Geiger, A.: GIRAFFE: Representing scenes as compositional generative neural feature fields. In: CVPR (2021) 4, 11, 12

55. Or-El, R., Luo, X., Shan, M., Shechtman, E., Park, J.J., Kemelmacher-Shlizerman, I.: StyleSDF: High-resolution 3d-consistent image and geometry generation. CVPR (2022) 1, 2, 4, 12, 22

56. Pan, X., Xu, X., Loy, C.C., Theobalt, C., Dai, B.: A shading-guided generative implicit model for shape-accurate 3d-aware image synthesis. NeurIPS (2021) 4, 9, 27

57. Salimans, T., Zhang, H., Radford, A., Metaxas, D.: Improving GANs using optimal transport. In: ICLR (2018) 4

58. Schwarz, K., Liao, Y., Niemeyer, M., Geiger, A.: GRAF: Generative radiance fields for 3D-aware image synthesis. In: NeurIPS (2020) 1, 2, 4

59. Serengil, S.I., Ozpinar, A.: Lightface: A hybrid deep face recognition framework. In: 2020 Innovations in Intelligent Systems and Applications Conference (ASYU). IEEE (2020) 10

60. Shade, J., Gortler, S., Hey, L.W., Szeliski, R.: Layered Depth Images. SIGGRAPH (1998) 4

61. Shi, Y., Aggarwal, D., Jain, A.K.: Lifting 2d StyleGAN for 3D-aware face generation. CVPR (2021) 5, 10, 11, 12, 20

62. Srinivasan, P.P., Tucker, R., Barron, J.T., Ramamoorthi, R., Ng, R., Snavely, N.: Pushing the boundaries of view extrapolation with multiplane images. In: CVPR (2019) 4

63. Sun, R., Fang, T., Schwing, A.G.: Towards a Better Global Loss Landscape of GANs. In: Proc. NeurIPS (2020) 4

64. Tucker, R., Snavely, N.: Single-view view synthesis with multiplane images. In: CVPR (2020) 4

65. Wu, J., Zhang, C., Xue, T., Freeman, W.T., Tenenbaum, J.B.: Learning a probabilistic latent space of object shapes via 3D generative-adversarial modeling. In: NeurIPS (2016) 4

66. Xu, Y., Peng, S., Yang, C., Shen, Y., Zhou, B.: 3D-aware Image Synthesis via Learning Structural and Textural Representations. CVPR (2022) 1, 2, 4

67. Zhou, P., Xie, L., Ni, B., Tian, Q.: CIPS-3D: A 3D-Aware Generator of GANs Based on Conditionally-Independent Pixel Synthesis. ArXiv (2021) 1, 2, 4, 12

68. Zhou, T., Tucker, R., Flynn, J., Fyffe, G., Snavely, N.: Stereo magnification: Learning view synthesis using multiplane images. TOG (2018) 3, 4

69. Zhu, J.Y., Zhang, Z., Zhang, C., Wu, J., Torralba, A., Tenenbaum, J.B., Freeman, W.T.: Visual object networks: Image generation with disentangled 3D representations. In: NeurIPS (2018) 4

# Supplementary Material:
# Generative Multiplane Images:
# Making a 2D GAN 3D-Aware

This supplementary is structured as follows:

## A    Differentiable Rendering in GMPI

In Sec. 3.3, we obtain the desired image $I_{v_{\mathsf{tgt}}}$ which illustrates the generated MPI representation $\mathcal{M} = \{C, \{\alpha_1, \ldots, \alpha_L\}\}$ from the user-specified target view $v_{\mathsf{tgt}}$ in two steps: 1) a warping step transforms the representation $\mathcal{M}$ from its canonical pose $v_{\mathsf{cano}}$ to the target pose $v_{\mathsf{tgt}}$; 2) a compositing step combines the planes into the desired image $I_{v_{\mathsf{tgt}}}$. Importantly, both steps entail easy computations which are end-to-end differentiable such that they can be included into any generator. Here we provide details.

**Warping.** We warp the RGB image and the alpha map of the $i^{\mathrm{th}}$ plane from the canonical view to the target view via

$$(C_i', \alpha_i') = \mathcal{H}_{i, v_{\mathsf{cano}} \to v_{\mathsf{tgt}}}(C, \alpha_i, d_i). \tag{S1}$$

Here, $\mathcal{H}_{i, v_{\mathsf{cano}} \to v_{\mathsf{tgt}}}$ represents the homography operation. Essentially, the homography $\mathcal{H}_{i, v_{\mathsf{cano}} \to v_{\mathsf{tgt}}}$ specifies a mapping: for each pixel coordinate $(p_x', p_y')$ in the image $C_i'$ and in the alpha map $\alpha_i'$ of the target view $v_{\mathsf{tgt}}$, we obtain corresponding coordinates $(p_x, p_y)$ in the image $C$ and in the alpha map $\alpha_i$ of the canonical view $v_{\mathsf{cano}}$. Bilinear sampling is applied on $(p_x, p_y)$ to obtain values for the pixel locations $(p_x', p_y')$. Concretely,

$$\begin{bmatrix} p_x\ p_y\ 1 \end{bmatrix}^\top = K_{v_{\mathsf{cano}}} \left( R_{v_{\mathsf{tgt}} \to v_{\mathsf{cano}}} - \frac{\boldsymbol{t}_{v_{\mathsf{tgt}} \to v_{\mathsf{cano}}} \boldsymbol{n}^\top}{b_i} \right) K_{v_{\mathsf{tgt}}}^{-1} \begin{bmatrix} p_x'\ p_y'\ 1 \end{bmatrix}^\top, \tag{S2}$$

where $\boldsymbol{n} \in \mathbb{R}^3$ is the normal of the plane defined in target camera coordinate system $v_{\mathsf{tgt}}$, which is identical for all planes. $b_i$ is the depth of the plane from the target camera $v_{\mathsf{tgt}}$. We let $K_{v_{\mathsf{cano}}} \in \mathbb{R}^{3 \times 3}$ and $K_{v_{\mathsf{tgt}}} \in \mathbb{R}^{3 \times 3}$ refer to the intrinsic matrices of the canonical view $v_{\mathsf{cano}}$ and the target view $v_{\mathsf{tgt}}$. Further, $R_{v_{\mathsf{tgt}} \to v_{\mathsf{cano}}} \in \mathbb{R}^{3 \times 3}$ and $\boldsymbol{t}_{v_{\mathsf{tgt}} \to v_{\mathsf{cano}}} \in \mathbb{R}^{3 \times 1}$ are the rotation and the translation from $v_{\mathsf{tgt}}$ to $v_{\mathsf{cano}}$.

**Alpha Compositing.** Given the warped image $C_i'$ and the warped alpha map $\alpha_i'$ for each plane $i$, we compute the final rendered 2D image $I_{v_{\mathsf{tgt}}}$ via

$$I_{v_{\mathsf{tgt}}} = \sum_{i=1}^{L} \left( C_i' \cdot \alpha_i' \cdot \prod_{j=1}^{i-1} (1 - \alpha_j') \right). \tag{S3}$$

Table S1: **Comparing representations on FFHQ.** 1[st] and 3[rd] rows are copied from Tab. 2's 3[rd] and 9[th] rows in the main text. Row 2 and 3 infer with 96 planes. Row 1: depth to renderer; Row 2: depth to MPI to renderer; Row 3: MPI to renderer.

|  | Row | | FID↓ | KID×100↓ | ID↑ | Depth↓ | Pose↓ |
|---|---|---|---|---|---|---|---|
|  | 1 | LiftedGAN [58] | 29.8 | – | 0.58 | **0.40** | 0.023 |
| $256^2$ | 2-1 | D2A ($\epsilon = 1/64$) | 13.4 | 0.920 | 0.69 | 0.60 | **0.004** |
|  | 2-2 | D2A ($\epsilon = 1/128$) | 13.5 | 0.867 | **0.70** | 0.60 | **0.004** |
|  | 2-3 | D2A ($\epsilon = 1/256$) | 11.7 | **0.644** | **0.70** | 0.63 | 0.005 |
|  | 2-4 | D2A ($\epsilon = 1/512$) | 12.6 | 0.684 | 0.69 | 0.62 | 0.005 |
|  | 3 | **GMPI** | **11.4** | 0.738 | **0.70** | 0.53 | **0.004** |

Similarly, we approximate depth $D_{v_{\text{tgt}}}$ via

$$D_{v_{\text{tgt}}} = \sum_{i=1}^{L} \left( b_i \cdot \alpha_i' \cdot \prod_{j=1}^{i-1} (1 - \alpha_j') \right), \tag{S4}$$

where $b_i$ is the distance mentioned in Eq. (S2). Notably, the combination of Eq. (4), Eq. (S1) and Eq. (S3) is end-to-end differentiable and hence straightforward to integrate into a generator. Importantly, the computations are also extremely efficient as only simple matrix multiplications are involved. It is hence easy to augment an existing generator like the one in StyleGANv2.

# B   Additional Quantitative Results

## B.1   Comparisons of Representations

We provide additional ablations of the plane representation. Specifically, we consider the following three ablations:

**Depth map to renderer:** We let the generator predict a single channel depth map (instead of multiple-channel alpha maps) and use a differentiable renderer to supervise the transformed image. LiftedGAN [61] serves as this ablation.

**Depth maps to alpha maps to renderer:** We compute multiple alpha maps from a predicted single-channel depth map in two steps: 1) we predict a depth map $\text{Depth} \in \mathbb{R}^{H \times W}$ in normalized coordinates; 2) we generate $L$ alpha maps from $\text{Depth}$ by computing the alpha value for a pixel $[x, y]$ on the $i$-th alpha map $\alpha_i$ via:

$$\alpha_i[x, y] = \min \left( 1, \max \left( 0, \frac{d_i - (\text{Depth}[x, y] - \epsilon)}{2\epsilon} \right) \right),$$

where $d_i$ is $\alpha_i$'s depth. Essentially, alpha values linearly increase from 0 to 1 within the range $[\text{Depth}[x, y] - \epsilon, \text{Depth}[x, y] + \epsilon]$. Any depth closer than $\text{Depth}[x, y] - \epsilon$ is set to 0. Any depth further than $\text{Depth}[x, y] + \epsilon$ is set to 1. We call this representation D2A.

Table S2: **Ablation studies on truncation level $\psi$ during inference on FFHQ.** All results use the same checkpoint trained with DPC, plane-specific features $\mathcal{F}^h_{\alpha_i}$, and shading-guided training. We evaluate with 96 planes. We show more digits for 'Depth' and 'Pose' than we use in Tab. 2 to emphasize differences. With lower $\psi$, GMPI produces geometries with fewer artifacts (Depth) at the cost of less variety (FID/KID).

|  | $\psi$ | $256^2$ | | | | | $512^2$ | | | | | $1024^2$ | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  |  | FID↓ | KID↓ | ID↑ | Depth↓ | Pose↓ | FID↓ | KID↓ | ID↑ | Depth↓ | Pose↓ | FID↓ | KID↓ | ID↑ | Depth↓ | Pose↓ |
| (a) | 1.0 | **11.4** | **0.738** | **0.70** | 0.533 | 0.0042 | **8.29** | **0.454** | **0.74** | 0.455 | 0.0056 | **7.50** | **0.407** | **0.75** | 0.535 | 0.0068 |
| (b) | 0.9 | 12.6 | 0.837 | 0.70 | 0.515 | 0.0037 | 10.1 | 0.588 | 0.74 | 0.421 | 0.0047 | 8.83 | 0.496 | 0.75 | 0.490 | 0.0058 |
| (c) | 0.8 | 15.4 | 1.079 | 0.70 | 0.497 | 0.0033 | 13.7 | 0.885 | 0.73 | 0.386 | 0.0039 | 11.5 | 0.691 | 0.74 | 0.449 | 0.0049 |
| (d) | 0.7 | 20.3 | 1.510 | 0.70 | 0.477 | 0.0029 | 19.8 | 1.391 | 0.73 | 0.354 | 0.0032 | 16.0 | 1.024 | 0.73 | 0.408 | 0.0040 |
| (e) | 0.6 | 27.8 | 2.172 | 0.70 | 0.459 | 0.0026 | 28.8 | 2.160 | 0.73 | 0.326 | 0.0025 | 22.8 | 1.510 | 0.73 | 0.368 | 0.0031 |
| (f) | 0.5 | 38.5 | 3.112 | 0.70 | **0.445** | **0.0023** | 41.1 | 3.249 | 0.72 | **0.304** | **0.0020** | 32.0 | 2.145 | 0.72 | **0.333** | **0.0024** |

**Alpha maps to renderer:** We directly predict multiple-channel alpha maps. Our GMPI serves as this ablation.

In Tab. S1 we report the results. For D2A, we ablate several values of $\epsilon$ and find that GMPI always outperforms on FID and depth metrics, verifying the fidelity of th texture and the high-quality of the geometry.

## B.2    Depth Score Analysis

We inspect GMPI-synthesized images and find two main reasons for our suboptimal 'Depth' scores in Tab. 2: 1) artifacts produced by StyleGANv2; 2) specular reflections on images deteriorate geometry generation. We discuss both in detail below.

**1) Artifacts in StyleGANv2.** Truncation was introduced in [34, 7] to balance between variety and fidelity. Specifically, using a truncation level $\psi \in [0, 1]$, we replace the style embedding $\boldsymbol{\omega}$ in Eq. (3) with

$$\boldsymbol{\omega}' = \bar{\boldsymbol{\omega}} + \psi \cdot (\boldsymbol{\omega} - \bar{\boldsymbol{\omega}}), \tag{S5}$$

where $\bar{\boldsymbol{\omega}} = \mathbb{E}_{\boldsymbol{z}}[f_{\texttt{Mapping}}(\boldsymbol{z})]$ represents the style embedding space's center of mass. In practice, $\bar{\boldsymbol{\omega}}$ is approximated by computing the moving average of all $\boldsymbol{\omega}$ encountered during training. Without truncation, *i.e.*, for $\psi = 1.0$, we find Style-GANv2 can produce results with significant artifacts as shown in Fig. S1a. This demonstrates that the generator fails to convert the corresponding $\boldsymbol{\omega}$ properly. Although these artifacts are not being reflected in our 2D GAN metrics, *i.e.*, FID and KID, they do affect the alpha map generation adversely. Specifically, the alpha maps are produced based on feature $\mathcal{F}^h_{\alpha_i}$, which is largely determined by the style embedding $\boldsymbol{\omega}$ through Eq. (3) and Eq. (8). Consequently, artifacts cause an inferior 'Depth' score. To understand the effect of the truncation level $\psi$ during inference, we evaluate GMPI using various truncation levels $\psi$ in Tab. S2. As can be seen clearly, with smaller $\psi$, we consistently perform better on geometry metrics, *i.e.*, lower 'Depth' and 'Pose' error, while trading in variety, *i.e.*, higher FID/KID. We provide a qualitative example in Fig. S2. Note, we do not apply truncation for any other quantitative results reported in this paper. *I.e.*, we always use $\psi = 1.0$ for quantitative evaluation results except in Tab. S2.

**2) Specular reflections.** Our rendering model is not designed to handle specular reflections. Strong specular reflections in the training data tend to degrade geometry generation, producing artifacts such as concave foreheads. We provide a qualitative example in Fig. S3. Specular reflections are a common failure mode of geometry generation. For instance, see Sec. 5 and Fig. 8 in StyleSDF [55]. We leave addressing of this issue to future work.



(a) $\psi = 1.0$.

(b) $\psi = 0.7$.

(c) $\psi = 0.5$.

Fig. S1: **Effects of truncation level $\psi$ for StyleGANv2.** Images are generated with the officially-released code and checkpoints.[6]    Without truncation, $i.e.$, for $\psi = 1.0$, StyleGANv2 generates images with significant artifacts. This degrades alpha map generation. The images can be generated by running the command `python generate.py --outdir=out --trunc=1.0 --seeds=10,56,88 --network=https://nvlabs-fi-cdn.nvidia.com/stylegan2-ada-pytorch/ pretrained/ffhq.pkl` while setting `--trunc` to 1.0, 0.7, and 0.5 respectively.
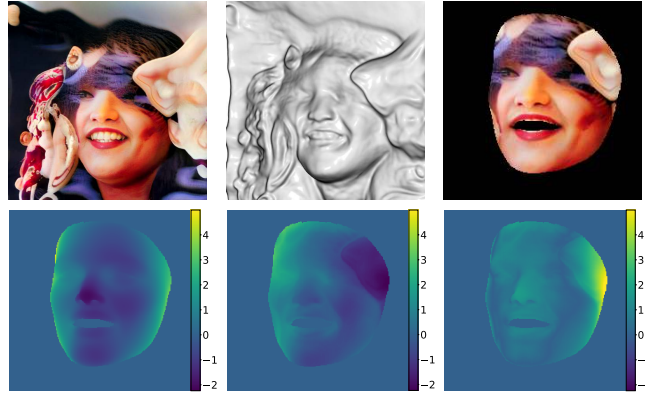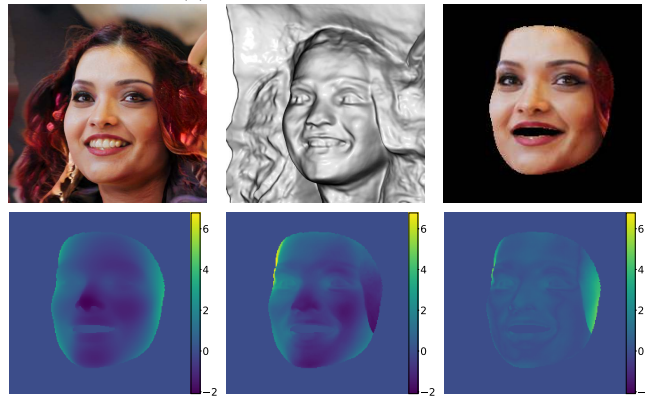
(a) $\psi = 1.0$, Depth score = 1.44.



(b) $\psi = 0.5$, Depth score = 0.83. The large difference in depth prediction appearing on the cheek is primarily due to hair. The parametric facial model does not reconstruct hair and accessories such as glasses (see Fig. 2 in [17]).

Fig. S2: **Truncation level $\psi$ affects geometry generation.** GMPI generated scenes with the same latent variable $z$ using two different truncation levels. Each scene is rendered with the same pose where the 'Depth' score is computed. For each subplot, from top to bottom, left to right, we show 1) rendered image; 2) corresponding geometry; 3) predicted face model [17]; 4) normalized pseudo ground-truth depth map; 5) normalized GMPI generated depth map (the smaller value indicates closer distance to the camera); and 6) difference between normalized depth maps 4) and 5). Smaller truncation levels $\psi$ benefit 'Depth' scores at the cost of increased FID/KID values.
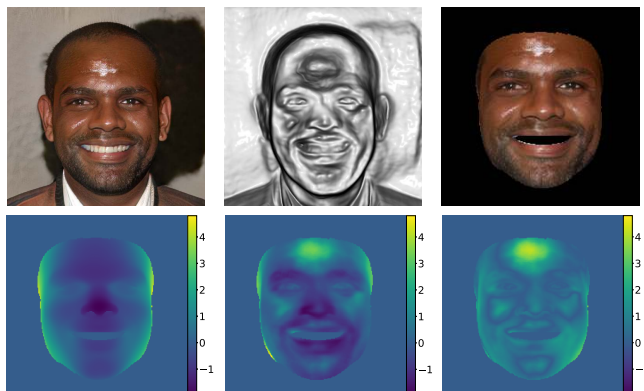
Fig. S3: **Specular reflection deteriorates geometry generation.** For this example, the depth score is 1.26. The strong lighting effect on the forehead breaks the Lambertian assumption and degrades alpha map generation.

Table S3: **Ablation studies on #planes during training.** No truncation is applied. We evaluate at a resolution of $512^2$. DPC refers to discriminator pose conditioning (Sec. 3.4), $\mathcal{F}^h_{\alpha_i}$ refers to the plane-specific feature introduced in Eq. (8), and Shading indicates the shading-guided training discussed in Sec. 3.5. #planes denotes the number of planes we used during training. During inference, we use 32 planes. We observe, the more planes we provide during training, the better the results of GMPI. Please see Fig. S5 for qualitative examples.

|     | #planes | DPC | $\mathcal{F}^h_{\alpha_i}$ | Shading | FFHQ | | | | | AFHQv2-Cat | |
|-----|---------|-----|------|---------|------|------|------|------|------|------|------|
|     |         |     |      |         | FID↓ | KID↓ | ID↑ | Depth↓ | Pose↓ | FID↓ | KID↓ |
| (a) | 32 | ✓ | ✓ | ✓ | **7.40** | **0.337** | **0.74** | **0.457** | **0.006** | 7.93 | 0.489 |
| (b) | 16 | ✓ | ✓ | ✓ | 9.83 | 0.575 | 0.74 | 0.574 | 0.007 | **6.82** | **0.358** |
| (c) | 8 | ✓ | ✓ | ✓ | 11.4 | 0.657 | 0.72 | 0.778 | 0.008 | 7.26 | 0.384 |
| (d) | 4 | ✓ | ✓ | ✓ | 16.4 | 1.043 | 0.65 | 0.992 | 0.007 | 8.53 | 0.467 |

### B.3   More Ablations

**#Planes During Training.** We ablate the #planes during training in Tab. S3 and Fig. S5. We observe: the more planes we can provide during training, the better the results of GMPI.

**Robustness to inaccurate camera poses.** To understand the robustness to inaccurate camera poses, we add noise to the observer's camera pose. Specifically, we follow [9] to first compute per-element standard deviation $\sigma$ of the estimated camera pose matrices for real images. During training, we add noise of $\{1\sigma, 2\sigma, 3\sigma, 4\sigma\}$ to each element of the camera pose matrix. We report results in Tab. S4 and Fig. S4. We want to emphasize, the estimated camera poses for real images are not perfect. Therefore, $0\sigma$ does not indicate fully-accurate camera pose information. Fig. S4 verifies that the more noise we add, the less photo-realistic the geometry we obtain. This is aligned with the trend of the Depth

Table S4: **Ablation study on camera pose accuracy.** No truncation is applied. We evaluate at a resolution of $512^2$. `Noise` indicates the level of noise we added to the camera poses during training. Specifically, $1\sigma$ denotes camera poses are corrupted with one standard deviation of the estimated camera pose matrices for images from the corresponding dataset (see Sec. B.3). All results are trained with `DPC`, plane-specific features $\mathcal{F}^h_{\alpha_i}$, and shading-guided training with 32 planes during training.

|     | Noise | FFHQ | | | | | AFHQv2-Cat | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
|     |       | FID↓ | KID↓ | ID↑ | Depth↓ | Pose↓ | FID↓ | KID↓ |
| (a) | $0\sigma$ | 8.29 | 0.454 | 0.74 | 0.46 | 0.006 | 7.79 | 0.474 |
| (b) | $1\sigma$ | 6.02 | 0.269 | 0.81 | 0.81 | 0.017 | 5.95 | 0.296 |
| (c) | $2\sigma$ | 5.38 | 0.221 | 0.86 | 1.00 | 0.027 | 5.51 | 0.266 |
| (d) | $3\sigma$ | 5.26 | 0.191 | 0.89 | 1.37 | 0.038 | 7.05 | 0.397 |
| (e) | $4\sigma$ | 4.72 | 0.176 | 0.90 | 1.68 | 0.046 | 6.92 | 0.393 |



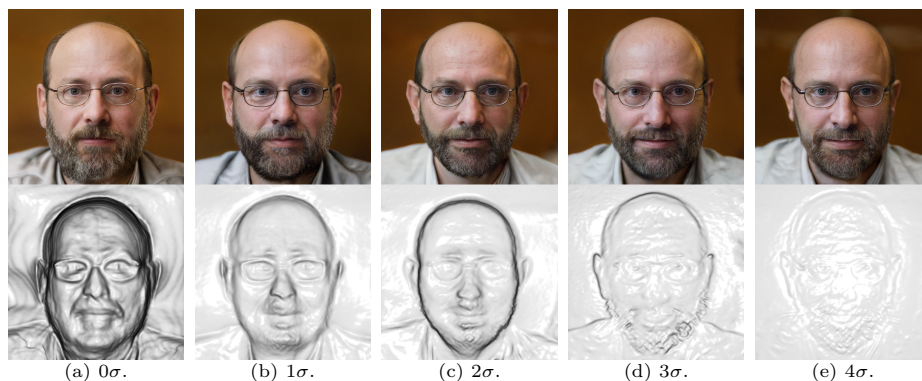(a) $0\sigma$.　　(b) $1\sigma$.　　(c) $2\sigma$.　　(d) $3\sigma$.　　(e) $4\sigma$.

Fig. S4: **Ablate robustness to inaccurate camera poses.** Panels (a)-(e) correspond to Tab. S4's rows (a)-(e). All results are generated from the same latent code $\boldsymbol{z}$. Within expectation, the more inaccurate camera poses the model is trained with, the less 3D geometry we can obtain.

and Pose metrics in Tab. S4. Note, FID, KID, and ID metrics are misleading as we do not observe much difference. This verifies again that 2D metrics lack the ability to capture 3D errors.

## C   Implementation Details

**Hyperparameters.** Tab. S5 summarizes the hyperparameters we used for training GMPI. We perform all training on 8 Tesla V100 GPUs using 5000 iterations. For all experiments except FFHQ256, we use a batch size of 32. This exposes the discriminator to 0.16 million real images. FFHQ256 is trained on 0.32 million real images since it uses a larger batch size of 64. We apply the minibatch

---

[6] https://github.com/NVlabs/stylegan2-ada-pytorch

(a) 32 planes.



(b) 16 planes.



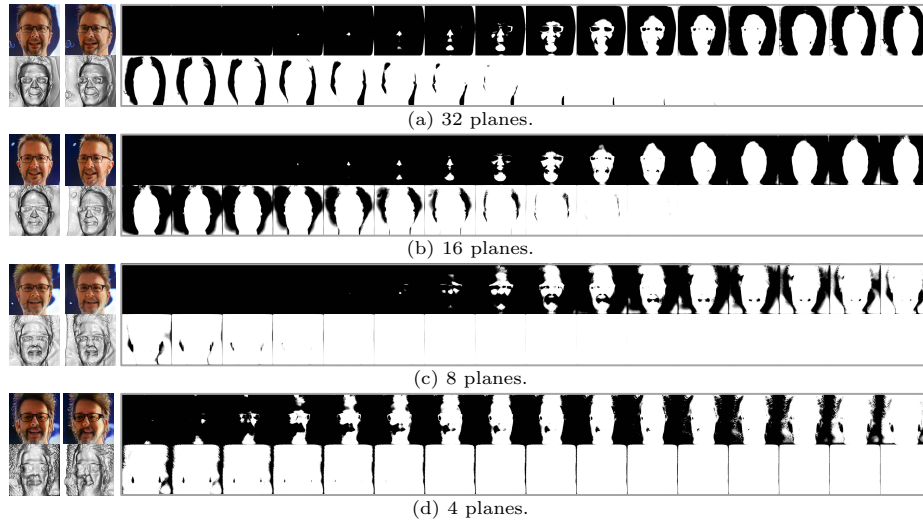(c) 8 planes.



(d) 4 planes.

Fig. S5: **Ablate #planes during training.** For each subplot, from left to right, we show two views and the 32 alpha maps that GMPI produces during inference. Each subplot's caption indicates the number of planes used during training. We observe: the more planes we can provide during training, the better the results of GMPI.

standard deviation layer [31] independently on each image. We apply $x$-flip data augmentation on AFHQv2 and MetFaces. To determine the channel number $\mathtt{dim}_h$ of $\mathcal{F}^h$ (Eq. (2)) and $\mathcal{F}^h_{\alpha_i}$ (Eq. (8)), we follow [32]'s official implementation[7] using:

$$\mathtt{dim}_h = \min(\frac{2^{15}}{h}, 512),  \tag{S6}$$

where $2^{15}$ is the channel base number. The only exception is FFHQ256, where we use a channel base number of $2^{14}$ following the official implementation. All experiments utilize R1 penalties [46] with weight 10.0 and learning rate $2 \times 10^{-3}$. We use Adam [36] as the optimizer. In all experiments, we use 32 planes during training.

**Mixed precision.** We follow StyleGANv2's official code and use half precision for both generator and discriminator layers corresponding to the four highest resolutions.

**Near/far depths.** We set the near and far depths of MPI for each dataset as 0.95/1.12 (FFHQ and MetFaces), 2.55/2.8 (AFHQv2-Cat).

**Depth normalization.** The depth $d_i$ $\forall i$ is normalized to the range $[0, 1]$ before being used in the embedding function $f_{\mathsf{Embed}}$ in Eq. (8). Specifically, we use $d'_i = \frac{d_i - d_1}{d_L - d_1}$.

**Shading-guided training.** As mentioned in Sec. 3.5, we utilize Eq. (10) to apply shading. For the first 1000 iterations, we set $k_a = 1.0$ and $k_d = 0.0$. We

---

[7] https://github.com/NVlabs/stylegan2-ada-pytorch

Table S5: **Hyperparameters used for training GMPI.**

|  | FFHQ256 | FFHQ512 | FFHQ1024 | AFHQv2 | MetFaces |
|---|---|---|---|---|---|
| Resolution | $256^2$ | $512^2$ | $1024^2$ | $512^2$ | $1024^2$ |
| #GPUs | 8 | 8 | 8 | 8 | 8 |
| Training length (iters) | 5k | 5k | 5k | 5k | 5k |
| Training length (#imgs) | 0.32M | 0.16M | 0.16M | 0.16M | 0.16M |
| Batch size | 64 | 32 | 32 | 32 | 32 |
| Minibatch stddev [31] | 1 | 1 | 1 | 1 | 1 |
| Dataset $x$-flips | ✗ | ✗ | ✗ | ✓ | ✓ |
| Channel base | $\frac{1}{2}\times$ | $1\times$ | $1\times$ | $1\times$ | $1\times$ |
| Learning rate ($\times 10^{-3}$) | 2 | 2 | 2 | 2 | 2 |
| R1 penalty weight [46] | 10 | 10 | 10 | 10 | 10 |
| Mixed-precision | ✓ | ✓ | ✓ | ✓ | ✓ |

linearly reduce $k_a$ to 0.9 and linearly increase $k_d$ to 0.1 during iteration 1001 to 2000. Starting from the $2001^{st}$ iteration, we fix $k_a = 0.9$ and $k_d = 0.1$. Meanwhile, lighting direction $l$ is randomly sampled. We follow [56] to use $l_h \sim \mathcal{N}(0.0, 0.2)$ and $l_v \sim \mathcal{N}(0.2, 0.05)$, where $l_h$ and $l_v$ represent horizontal and vertical angles for lighting directions respectively.

**Structure of $f_{\mathsf{Embed}}$ (Eq. (8)).** For each resolution $h \times h$, we utilize three modulated-convolutional layers [35] to construct $f_{\mathsf{Embed}}$, whose channel numbers are $\mathtt{dim}_{h/4}$, $\mathtt{dim}_{h/2}$, $\mathtt{dim}_h$ respectively.

**Alternative representation of $f_{\mathsf{Embed}}$ (Eq. (8)).** Besides conditioning $f_{\mathsf{Embed}}$ on specific depth $d_i$, we also studied use of a learnable token. Concretely, $f_{\mathsf{Embed}}$ was represented using a tensor of shape $h \times h \times \mathtt{dim}_h$. This resulted in reasonable geometry but we find conditioning of alpha maps on depth to provide better results. We leave further exploration of learnable tokens to future work.

**Background plane.** We treat the last plane of the MPI-like representation, *i.e.*, the $L^{\mathrm{th}}$ plane, as the background plane. To color this plane, we treat the left-most and right-most 5% pixels of the synthesized image $C$ as the left and right boundary. RGB values of all remaining pixels on this plane are linearly interpolated between the left and the right boundary. We find this simple procedure to work well in our experiments.

**Mesh generation.** We utilize the marching cube algorithm [45] implemented in PyMCubes for generating meshes.[8] We also utilize its smoothing function for better visualization.

# D   Additional Qualitative Results

## D.1   Uncurated Results

We provide uncurated results on FFHQ (Fig. S6), AFHQv2 (Fig. S7), and Met-Faces (Fig. S8). We observe GMPI to generate high-quality geometry.

---

[8] https://github.com/pmneila/PyMCubes

(a) Renderings for corresponding geometries in Fig. S6b.



(b) Geometries for corresponding renderings in Fig. S6a.

Fig. S6: **Uncurated results on FFHQ.** From top to bottom, left to right, we show generations with seed 1-32. Results are generated with truncation $\psi = 0.5$ [34].

## D.2    Style Mixing

We illustrate style mixing [34] results on FFHQ (Fig. S9), AFHQv2 (Fig. S10), and MetFaces (Fig. S11). GMPI successfully disentangles coarse and fine levels of generations.

## D.3    More Results

This supplementary material also includes an interactive viewer for the generated MPI representations and an HTML page with videos to illustrate generations from GMPI.

(a) Renderings for corresponding geometries in Fig. S7b.



(b) Geometries for corresponding renderings in Fig. S7a.

Fig. S7: **Uncurated results on AFHQv2.** From top to bottom, left to right, we show generations with seed 1-32. Results are generated with truncation $\psi = 0.7$ [34].
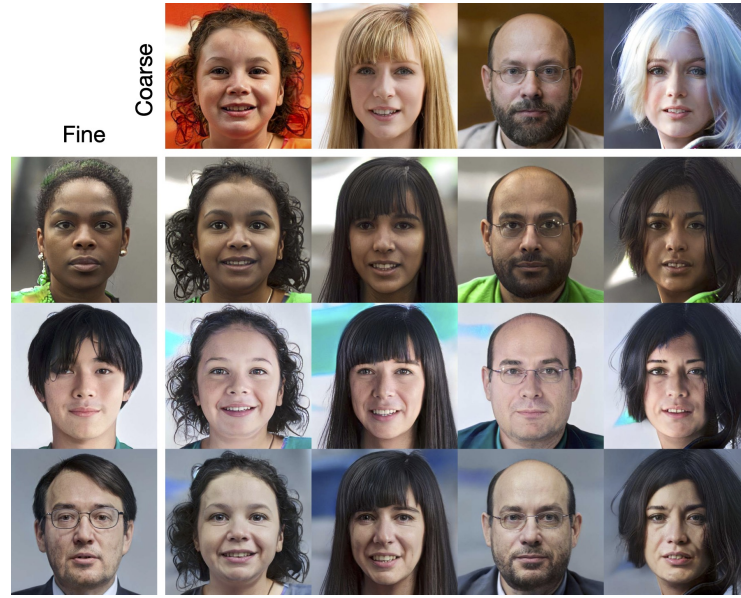
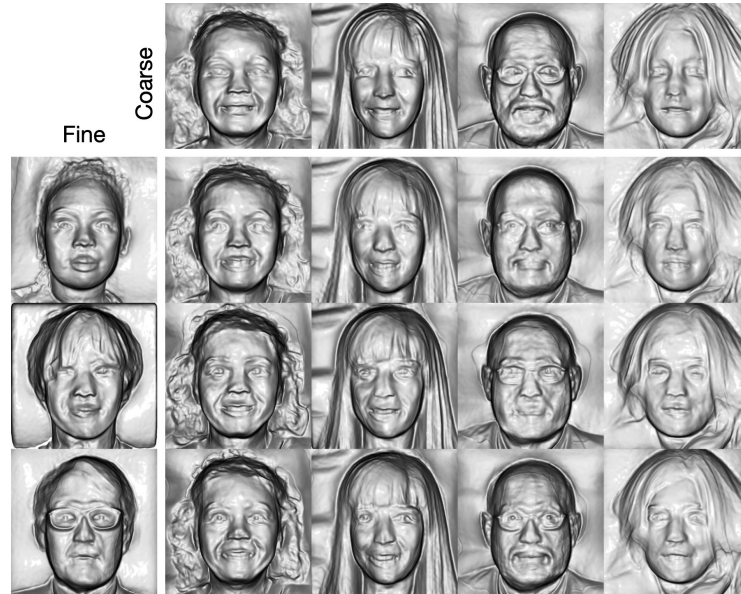(a) Renderings for corresponding geometries in Fig. S8b.



(b) Geometries for corresponding renderings in Fig. S8a.

Fig. S8: **Uncurated results on MetFaces.** From top to bottom, left to right, we show generations with seed 1-32. Results are generated with truncation $\psi = 0.7$ [34].

(a) Renderings for corresponding geometries in Fig. S9b.
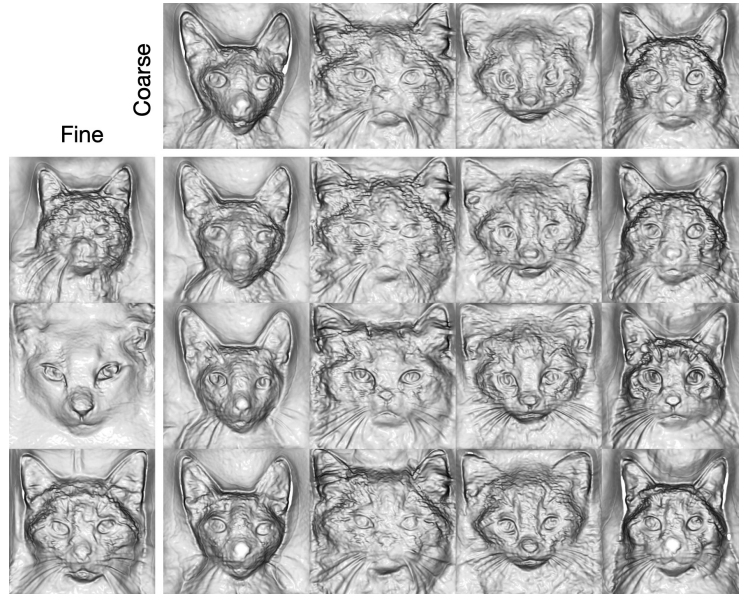


(b) Geometries for corresponding renderings in Fig. S9a.

Fig. S9: **Style mixing on FFHQ.** Results don't use truncation, *i.e.*, $\psi = 1.0$. To obtain each cell in the bottom right grid, we replace lower-level style embeddings $\boldsymbol{\omega}$ (Eq. (3)) in the *Fine* column with the corresponding $\boldsymbol{\omega}$ from the *Coarse* row. We observe, GMPI enables semantic editing: lower-level $\boldsymbol{\omega}$ (layers $0-6$) control the shape while upper-level $\boldsymbol{\omega}$ (layers 7 and higher) determine fine-grained styles.

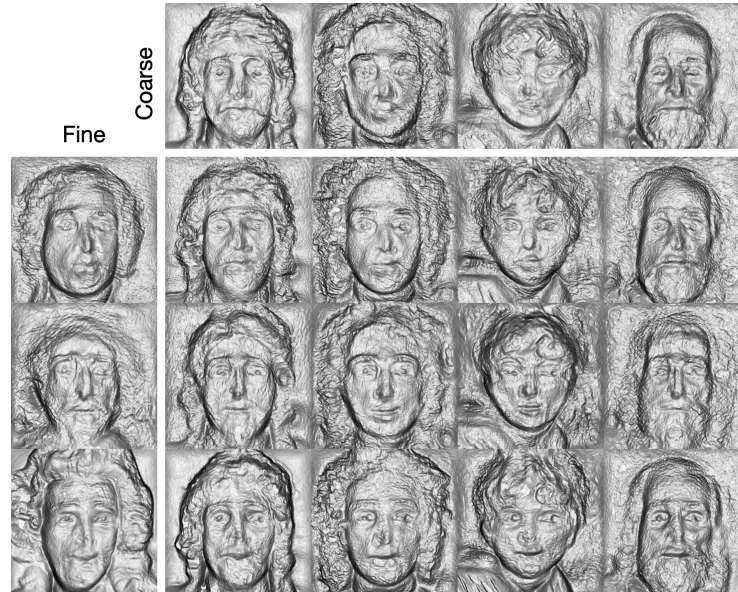(a) Renderings for corresponding geometries in Fig. S10b.



(b) Geometries for corresponding renderings in Fig. S10a.

Fig. S10: **Style mixing on AFHQv2.** Results don't use truncation, *i.e.*, $\psi = 1.0$. To obtain each cell in the bottom right grid, we replace lower-level style embeddings $\boldsymbol{\omega}$ (Eq. (3)) in the *Fine* column with the corresponding $\boldsymbol{\omega}$ from the *Coarse* row. We observe, GMPI enables semantic editing: lower-level $\boldsymbol{\omega}$ (layers $0-6$) control the shape while upper-level $\boldsymbol{\omega}$ (layers 7 and higher) determine fine-grained styles.

(a) Renderings for corresponding geometries in Fig. S11b.



(b) Geometries for corresponding renderings in Fig. S11a.

Fig. S11: **Style mixing on MetFaces.** Results don't use truncation, *i.e.*, $\psi = 1.0$. To obtain each cell in the bottom right grid, we replace lower-level style embeddings $\boldsymbol{\omega}$ (Eq. (3)) in the *Fine* column with the corresponding $\boldsymbol{\omega}$ from the *Coarse* row. We observe, GMPI enables semantic editing: lower-level $\boldsymbol{\omega}$ (layers $0-6$) control the shape while upper-level $\boldsymbol{\omega}$ (layers 7 and higher) determine fine-grained styles.