

---

# Rethinking training of 3D GANs

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We are witnessing a surge of works on building and improving 3D-aware generators.  
2 To induce a 3D-aware bias, such models rely on volumetric rendering, which  
3 is expensive to employ at high resolutions. The dominant strategy to address  
4 the scaling issue is to train a separate 2D decoder to upsample a low-resolution  
5 volumetrically rendered representation. But this solution comes at a cost. Not only  
6 does it break multi-view consistency, e.g. shape and texture change when a camera  
7 moves, but it also learns the geometry in a low fidelity. In this work, we take a  
8 different route to 3D synthesis and develop a non-upsampler-based generator with  
9 state-of-the-art image quality, high-resolution geometry and which trains  $2.5\times$   
10 *faster*. For this, we revisit and improve patch-based optimization in two ways.  
11 First, we design a location- and scale-aware discriminator by modulating its filters  
12 with a hypernetwork. Second, we modify the patch sampling strategy based on an  
13 annealed beta distribution to stabilize training and accelerate the convergence. We  
14 train on four datasets (two introduced in this work) at  $256^2$  and  $512^2$  resolutions,  
15 directly, without the need of a 2D upsampler, and our model attains better or  
16 comparable FID and has higher fidelity geometry than the current SotA.

17 Code/data/visualizations: <https://rethinking-3d-gans.github.io>

## 18 1 Introduction

19 Generative models for image synthesis achieved remarkable success in recent years and enjoy a lot of  
20 practical applications [55, 24]. While initially they mainly focused on 2D images [21, 65, 25, 4, 28],  
21 recent research explored generative frameworks with partial 3D control over the underlying object  
22 in terms of texture/structure decomposition, novel view synthesis or lighting manipulation (e.g.,  
23 [58, 56, 7, 67, 6, 12, 49]). These techniques are typically built on top of the recently emerged neural  
24 radiance fields (NeRF) [38] to explicitly represent the object (or its latent features) in 3D space.

25 NeRF is a powerful framework, which made it possible to built expressive 3D-aware generators  
26 from challenging RGB datasets [7, 12, 6]. Under the hood, it trains a multi-layer perceptron (MLP)  
27  $F(\mathbf{x}; \mathbf{d}) = (\mathbf{c}, \sigma)$  to represent a scene by encoding a density  $\sigma \in \mathbb{R}_+$  for each coordinate position  
28  $\mathbf{x} \in \mathbb{R}^3$  and a color value  $\mathbf{c} \in \mathbb{R}^3$  from  $\mathbf{x}$  and view direction  $\mathbf{d} \in \mathbb{S}^2$  [38]. To synthesize an image,  
29 one renders each pixel independently by casting a ray  $\mathbf{r}(q) = \mathbf{o} + q\mathbf{d}$  (for  $q \in \mathbb{R}_+$ ) from origin  
30  $\mathbf{o} \in \mathbb{R}^3$  into the direction  $\mathbf{d} \in \mathbb{S}^2$  and aggregating many color values along it with their corresponding  
31 densities. Such a representation is very expressive, but comes at a cost: rendering a single pixel is  
32 computationally expensive and makes it intractable to produce a lot of pixels in one forward pass. It  
33 is *not* fatal for reconstruction tasks where the loss can be robustly computed on a subset of pixels,  
34 but it creates significant scaling problems for generative NeRFs: they are typically formulated in a  
35 GAN-based framework [14] with 2D convolutional discriminators requiring a full image as input.

36 People address these scaling issues of NeRF-based GANs in different ways, but the dominating  
37 approach is to train a separate 2D decoder to produce a high-resolution image from a low-resolution

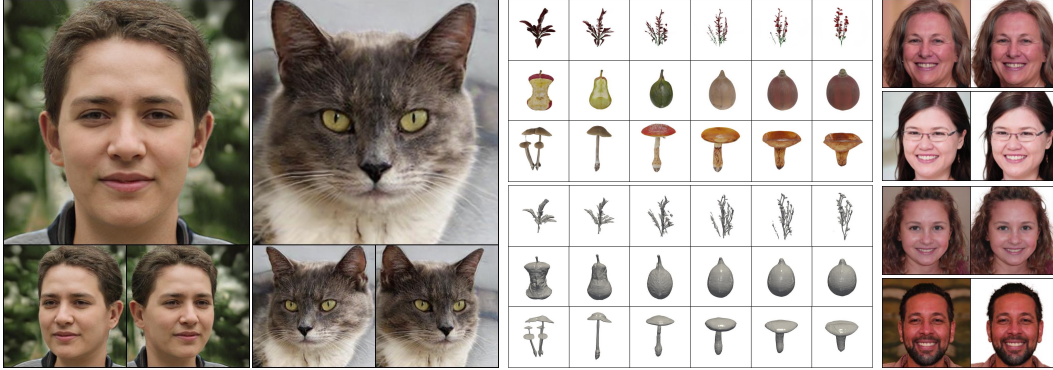


Figure 1: We build a pure NeRF-based generator trained in a patch-wise fashion. Left two grids: samples on FFHQ  $512^2$  [25] and Cats  $256^2$  [76]. Middle grids: interpolations between samples on M-Plants and M-Food (upper) and corresponding geometry interpolations (lower). Right grid: background separation examples. In contrast to the upsampler-based methods, one can naturally incorporate the techniques from the traditional NeRF literature into our generator: for background separation, we simply copy-pasted the corresponding code from NeRF++ [75].

38 image or feature grid rendered from a NeRF backbone [43]. During the past six months, there  
 39 appeared *more than a dozen* of methods which follow this paradigm (e.g., [6, 15, 70, 47, 78, 35,  
 40 74, 23, 71, 77, 63]). While using the upsampler allows to scale the model to high resolution, it  
 41 comes with two severe limitations: 1) it breaks multi-view consistency of a generated object, i.e. its  
 42 texture and shape change when the camera moves; and 2) the geometry gets only represented in a  
 43 low resolution ( $\approx 64^3$ ). In our work, we show that by dropping the upsampler and using a simple  
 44 patch-wise optimization scheme, one can build a 3D generator with better image quality, faster  
 45 training speed and without the above limitations.

46 Patch-wise training of NeRF-based GANs was originally proposed by GRAF [56] and got largely  
 47 neglected by the community since then. The idea is simple: instead of training the generative model on  
 48 full-size images, one does this on small random crops. Since the model is coordinate-based [59, 64],  
 49 it does not face any issues to synthesize only a subset of pixels. This serves as a good way to save  
 50 computation for both the generator and the discriminator, since it makes both of them operate on  
 51 patches of small spatial resolution. To make the generator learn both the texture and the structure,  
 52 crops are sampled to be of variable scales (but having the *same* number of pixels): in some sense, this  
 53 can be seen as optimizing the model on low-resolution images + high-resolution patches.

54 In our work, we improve patch-wise training in two crucial ways. First, we redesign the discriminator  
 55 by making it better suited to operating on image patches of variable scales and locations. Convolutional  
 56 filters of a neural network learn to capture different patterns in their inputs depending on their  
 57 semantic receptive fields [30, 46]. That’s why it is detrimental to reuse the same discriminator to  
 58 judge both high-resolution local and low-resolution global patches, inducing additional burden on it to  
 59 mix filters’ responses of different scales. To mitigate this, we propose to modulate the discriminator’s  
 60 filters with a hypernetwork [16], that predicts which filters to suppress or reinforce from a given patch  
 61 scale and location.

62 Second, we change the random scale sampling strategy from an annealed uniform to an annealed beta  
 63 distribution. Typically, patch scales are sampled from a uniform distribution  $s \sim \mathcal{U}[s(t), 1]$  [56, 36, 5],  
 64 where the minimum scale  $s(t)$  is gradually decreased (i.e. annealed) till some iteration  $T$  from  
 65  $s(0) = 0.9$  to a smaller value  $s(T)$  (in the interval  $[0.125 - 0.5]$ ) during training. This sampling  
 66 strategy prevents learning high-frequency details early on in training and puts too little attention on the  
 67 structure after  $s(t)$  reached its final value  $s(T)$ . This makes the overall convergence of the generator  
 68 slower and less stable that’s why we propose to sample patch scales using the beta distribution  
 69  $\text{Beta}(1, \beta(t))$  instead, where  $\beta(t)$  is gradually annealed from  $\beta(0) \approx 0$  to some maximum value  
 70  $\beta(T)$ . In this way, the model starts learning high-frequency details immediately after the training  
 71 starts and has more focus on the structure after the growth is done. This simple change stabilizes the  
 72 training and allows to converge faster than the typically used uniform distribution [56, 5, 36].

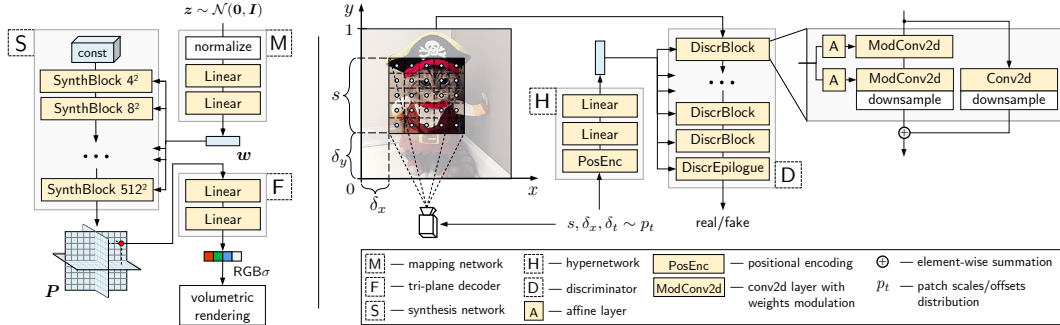


Figure 2: Illustration of the architecture and the patch sampling process. See §3 for details.

73 We use those two ideas to develop a novel state-of-the-art 3D GAN: Patch-wise Generative Radiance Fields (PGRF). We employ it for high-resolution 3D-aware image synthesis on four datasets: 74 FFHQ [25], Cats [76], Megascans Plants and Megascans Food. The last two benchmarks are introduced in our work and contain  $360^\circ$  renderings of photo-realistic scans of different plants and 75 objects (described in §4). They are much more difficult in terms of geometry and are well suited for 76 assessing the structural limitations of modern 3D-aware generators. 77

78 Our model uses a pure NeRF-based backbone, that’s why it represents geometry in high resolution 79 and does not suffer from multi-view synthesis artifacts, as opposed to upsampler-based generators. 80 Moreover, it has higher or comparable image quality (as measured by FID [20]) and  $2.5\times$  lower 81 training cost. Also, in contrast to upsampler-based 3D GANs, our generator can naturally incorporate 82 the techniques from the traditional NeRF literature. To demonstrate this, we incorporate background 83 separation into our framework by simply copy-pasting the corresponding code from NeRF++ [75]. 84

## 85 2 Related work

86 **Neural Radiance Fields.** Neural Radiance Fields (NeRF) is an emerging area [38], which combines 87 neural networks with volumetric rendering techniques to perform novel-view synthesis [38, 75, 2], 88 image-to-scene generation [73], surface reconstruction [45, 68, 44] and other tasks [9, 17, 50]. In our 89 work, we employ them in the context of 3D-aware generation from a dataset of RGB images [56, 7].

90 **3D generative models.** A popular way to learn a 3D generative model is to train it on 3D data or in 91 an autoencoder’s latent space (e.g., [10, 69, 1, 34, 31, 39, 29]). This requires explicit 3D supervision 92 and there appeared methods which train from RGB datasets with segmentation masks, keypoints 93 or multiple object views [13, 32, 54]. Recently, there appeared works which train from single-view 94 RGB only, including mesh-generation methods [19, 72, 53] and methods that extract 3D structure 95 from pretrained 2D GANs [58, 48]. And recent neural rendering advancements allowed to train 96 NeRF-based generators [56, 7, 42] from purely RGB data from scratch, which became the dominating 97 direction since then and which are typically formulated in the GAN-based framework [14].

98 **NeRF-based GANs.** HoloGAN [41] generates a 3D feature voxel grid which is projected on a 99 plane and then upsampled. GRAF [56] trains a noise-conditioned NeRF in an adversarial manner. 100  $\pi$ -GAN [7] builds upon it and uses progressive growing and hypernetwork-based [16] conditioning 101 in the generator. GRAM [12] builds on top of  $\pi$ -GAN and samples ray points on a set of learnable 102 iso-surfaces. GNeRF [36] adapts GRAF for learning a scene representation from RGB images 103 without known camera parameters. GIRAFFE [43] uses a composite scene representation for better 104 controllability. CAMPARI [42] learns a camera distribution and a background separation network 105 with inverse sphere parametrization [75]. To mitigate the scaling issue of volumetric rendering, many 106 recent works train a 2D decoder under different multi-view consistency regularizations to upsample 107 a low-resolution volumetrically rendered feature grid [6, 15, 70, 47, 78, 71, 77]. However, none of such 108 regularizations can currently provide the multi-view consistency of pure-NeRF-based generators.

109 **Patch-wise generative models.** Patch-wise training had been routinely utilized to learn the textural 110 component of image distribution when the global structure is provided from segmentation masks, 111 sketches, latents or other sources (e.g., [22, 57, 11, 66, 52, 51, 33, 61]). Recently, there appeared 112 works which sample patches at variable scales, in which way a patch can carry global information

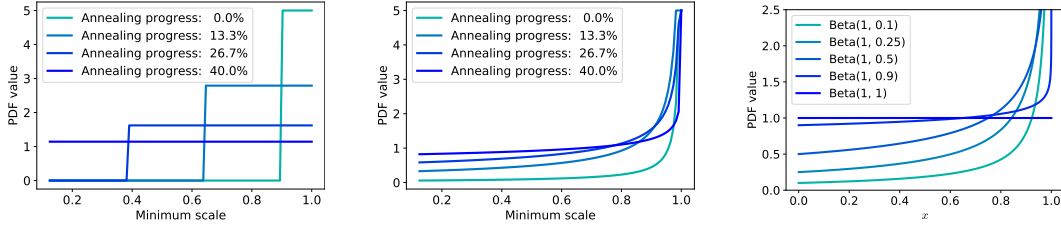


Figure 3: Comparing uniform (left) and beta (middle) annealed patch scale sampling in terms of their probability density function (PDF) (for visualization purposes, we clamp the maximum density value to 5); (right) PDF of Beta(1,  $\beta$ ), provided for completeness. Uniform distribution with annealed  $s_{\min}(0) = 0.9$  from 0.9 to  $s_{\min}(T) = 0.125$  does not put any attention to high-frequency details in the beginning and treats small-scale and large-scale patches equally at the end of the annealing. Beta distribution with annealed  $\beta(0) \approx 0$  to  $\beta(T) \approx 1$ , in contrast, lets the model learn high-resolution texture immediately after the training starts, and puts more focus on the structure at the end.

113 about the whole image. Recent works use it to train a generative NeRF [56], fit a neural representation  
 114 in an adversarial manner [36] or to train a 2D GAN on a dataset of variable resolution [5].

### 115 3 Model

116 We build upon StyleGAN2 [26], replacing its generator with the tri-plane-based NeRF model [6]  
 117 and using its discriminator as the backbone. We train the model on  $r \times r$  patches (we use  $r = 64$   
 118 everywhere) of random scales instead of the full images of resolution  $R \times R$ . Scales  $s \in [\frac{r}{R}, 1]$  are  
 119 randomly sampled from a time-varying distribution  $s \sim p_t(s)$ .

#### 120 3.1 NeRF-based generator

121 Compared to upsampler-based 3D GANs [15, 43, 71, 78, 6, 77], we use a pure NeRF [38] as our  
 122 generator G and utilize the tri-plane representation [6, 8] as the backbone. It consists of three  
 123 components: 1) mapping network  $M : z \mapsto w$  which transforms a noise vector  $z \sim \mathbb{R}^{512}$  into  
 124 the latent vector  $w \sim \mathbb{R}^{512}$ ; 2) synthesis network  $S : w \mapsto P$  which takes the latent vector  $w$   
 125 and synthesizes three 32-dimensional feature planes  $P = (P_{xy}, P_{yz}, P_{xz})$  of resolution  $R_p \times R_p$   
 126 (i.e.  $P_{(*)} \in \mathbb{R}^{R_p \times R_p \times 32}$ ); 3) tri-plane decoder network  $F : (x, P) \mapsto (c, \sigma) \in \mathbb{R}^4$ , which takes  
 127 the space coordinate  $x \in \mathbb{R}^3$  and tri-planes  $P$  as input and produces the RGB color  $c \in \mathbb{R}^3$  and  
 128 density value  $\sigma \in \mathbb{R}_+$  at that point by interpolating the tri-plane features in the given coordinate and  
 129 processing them with a tiny MLP. In contrast to classical NeRF [38], we do not utilize view direction  
 130 conditioning since it worsens multi-view consistency [7] in GANs which are trained on RGB datasets  
 131 with a single view per instance. To render a single pixel, we follow the classical volumetric rendering  
 132 pipeline with hierarchical sampling [38, 7], using 48 ray steps in coarse and 48 ones in fine sampling  
 133 stages. See the accompanying source code for more details.

#### 134 3.2 Scale/location-aware discriminator

135 Our discriminator D is built on top of StyleGAN2 [26]. Since we train the model in a patch-wise  
 136 fashion, the original backbone is not well suited for this: convolutional filters are forced to adapt  
 137 to signals of very different scales and extracted from different locations. A natural way to resolve  
 138 this problem is to use separate discriminators depending on the scale, but that strategy has three  
 139 limitations: 1) each particular discriminator receives less overall training signal (since the batch size  
 140 is limited); 2) from an engineering perspective, it is more expensive to evaluate a convolutional kernel  
 141 with different parameters on different inputs; 3) one can use only a small fixed amount of possible  
 142 patch scales. This is why we develop a novel hypernetwork-modulated [16] discriminator architecture  
 143 to operate on patches with continuously varying scale.

144 To modulate the convolutional kernels of D, we define a hypernetwork  $H : (s, \delta_x, \delta_y) \mapsto (\sigma_1, \dots, \sigma_L)$   
 145 as a 2-layer MLP with tanh non-linearity at the end which takes patch scale  $s$  and its cropping offsets  
 146  $\delta_x, \delta_y$  as input and produces modulations  $\sigma_\ell \in (0, 2)^{c_{\text{out}}}$  (we shift the tanh output by 1 to map into



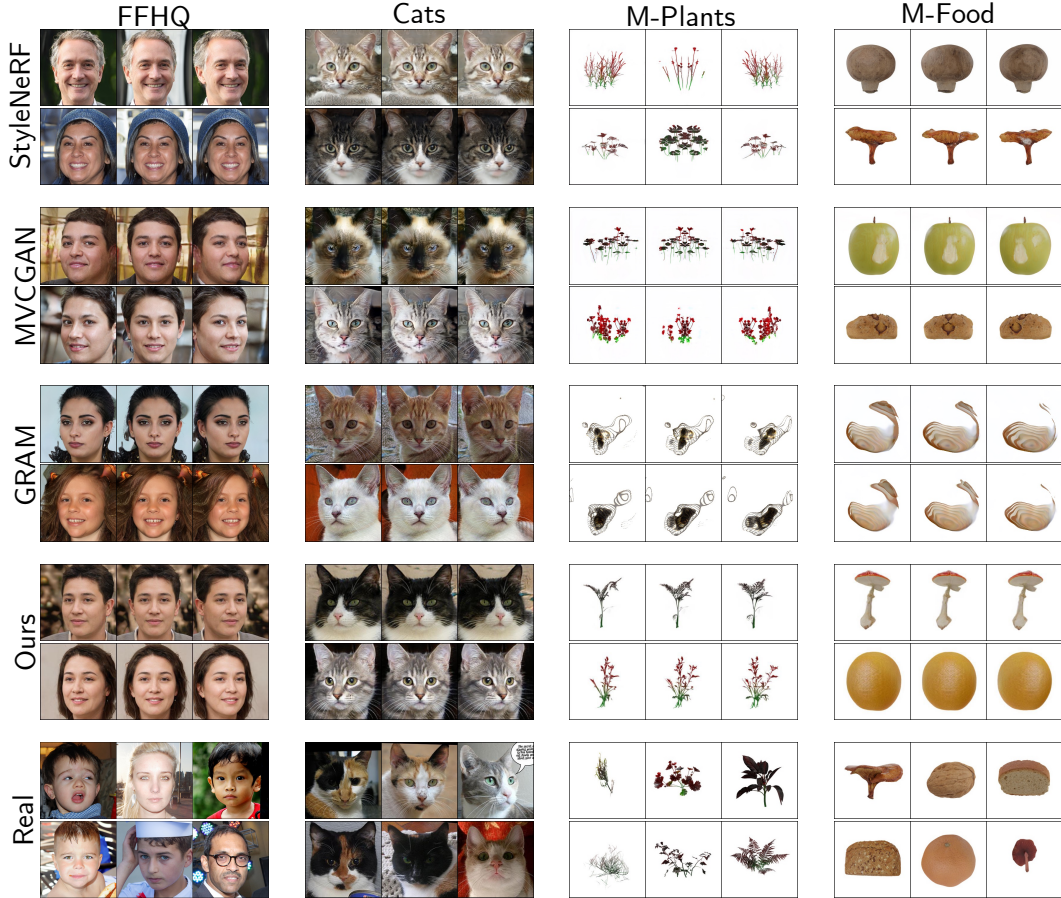


Figure 4: Comparing samples of PGRF and current 3D-aware generators. Our method attains state-of-the-art image quality, recovers high-fidelity geometry and preserves multi-view consistency for both simple-shape (FFHQ and Cats) and variable-shape (M-Plants and M-Food) datasets. We refer the reader to the supplementary for the video comparisons to evaluate multi-view consistency.

147 the 1-centered interval), where  $c_{\text{out}}^{\ell}$  is the number of output channels in the  $\ell$ -th convolutional layer.  
 148 Given a convolutional kernel  $\mathbf{W}^{\ell} \in \mathbb{R}^{c_{\text{out}}^{\ell} \times c_{\text{in}}^{\ell} \times k \times k}$  and input  $\mathbf{x} \in \mathbb{R}^{c_{\text{in}}}$ , a straightforward strategy  
 149 to apply the modulation is to multiply  $\sigma$  on the weights (depicting the convolution operation by  
 150  $\text{conv2d}(\cdot)$  and omitting its other parameters for simplicity):

$$\mathbf{y} = \text{conv2d}(\mathbf{W}^{\ell} \odot \sigma, \mathbf{x}), \quad (1)$$

151 where we broadcast the remaining axes and  $\mathbf{y} \in \mathbb{R}^{c_{\text{out}}}$  is the layer output (before the non-linearity).  
 152 However, using different kernel weights on top of different inputs is not too efficient in modern  
 153 deep learning frameworks (even with the group-wise convolution trick [26]), that’s why we use an  
 154 equivalent strategy by multiplying the weights on  $\mathbf{y}$  instead:

$$\mathbf{y} = \sigma \odot \text{conv2d}(\mathbf{W}^{\ell}, \mathbf{x}). \quad (2)$$

155 This suppresses and reinforces different convolutional filters of the layer depending on the patch  
 156 scale and location. And to incorporate an even stronger conditioning, we also use the projection  
 157 strategy [40] from StyleGAN2-ADA [24]. We depict our discriminator architecture in Fig 2. As we  
 158 show in Tab 2, it allows to obtain  $\approx 15\%$  lower FID compared to the standard discriminator.

### 159 3.3 Patch-wise optimization with Beta-distributed scales

160 Training NeRF-based GANs is computationally expensive, because rendering each pixel via vol-  
 161 umetric rendering requires many evaluations (e.g., in our case, 96) of the underlying MLP. For

162 scene reconstruction tasks, it does not create issues since the typically used  $\mathcal{L}_2$  loss [38, 75, 68]  
 163 can be robustly computed on a sparse subset of the pixels. But for NeRF-based GANs, it becomes  
 164 prohibitively expensive for high resolutions since convolutional discriminators operate on dense  
 165 full-size images. The currently dominating approach to mitigate this is to train a separate 2D decoder  
 166 to upsample a low-resolution image representation, rendered from a NeRF-based MLP. But this  
 167 breaks multi-view consistency (i.e. object’s shape and texture change when the camera is moving)  
 168 and learns the 3D geometry in a low resolution (from  $\approx 16^2$  [71] to  $\approx 128^2$  [6]). This is why we build  
 169 upon the multi-scale patch-wise training scheme [56] and demonstrate that it can give state-of-the-art  
 170 image quality and training speed without the above limitations.

171 Patch-wise optimization works the following way. On each iteration, instead of passing the full-size  
 172  $R \times R$  image to D, we instead input only a small patch with resolution  $r \times r$  of random scale  
 173  $s \in [r/R, 1]$  and extracted with a random offset  $(\delta_x, \delta_y) \in [0, 1 - s]^2$ . We illustrate this procedure in  
 174 Fig 2. Patch parameters are sampled from distribution:

$$s, \delta_x, \delta_y \sim p_t(s, \delta_x, \delta_y) \triangleq p_t(s)p(\delta_x|s)p(\delta_y|s) \quad (3)$$

175 where  $t$  is the current training iteration. In this way, patch scales depend on the current training  
 176 iteration  $t$  and offsets are sampled independently after we know  $s$ . As we show next, the choice of  
 177 distribution  $p_t(s)$  has crucial influence on the learning speed and stability.

178 Typically, patch scales are sampled from the annealed uniform distribution [56, 36, 5]  $s$ :

$$p_t(s) = U[s_{\min}(t), 1], \quad s_{\min}(t) = \text{lerp}[1, r/R, \min(t/T, 1)], \quad (4)$$

179 where  $\text{lerp}$  is the linear interpolation function<sup>1</sup>, and the left interval bound  $s_{\min}(t)$  is gradually  
 180 annealed during the first  $T$  iterations until it reaches the minimum possible value of  $r/R$ .<sup>2</sup> But this  
 181 strategy does not let the model learn high-frequency details early on in training and puts little focus on  
 182 the structure when  $s_{\min}(t)$  is fully annealed to  $r/R$  (which is usually very small, e.g.  $r/R = 0.125$   
 183 for a typical  $64^2$  patch-wise training on  $512^2$  resolution). As we show, the first issue makes the  
 184 generator converge slower, and the second one makes the overall optimization less stable.

185 To mitigate this, we propose a small change in the pipeline by simply replacing the uniform scale  
 186 sampling distribution with:

$$s \sim \text{Beta}(1, \beta(t)) \cdot (1 - r/R) + r/R, \quad (5)$$

187 where  $\beta(t)$  is gradually annealed from  $\beta(0)$  to some final value  $\beta(T)$ . Using beta distribution instead  
 188 of the uniform one gives a very convenient knob to shift the training focus between large patch scales  
 189  $s \rightarrow 1$  (carrying the global information about the whole image) and small patch scales  $r \rightarrow r/R$   
 190 (representing high-resolution local crops).

191 A natural way to do the annealing is to anneal from 0 to 1: at the start, the model focuses entirely on  
 192 the structure, while at the end it transforms into the uniform distribution (See Fig 3). We follow this  
 193 strategy, but from the design perspective instead set  $\beta(T)$  to a value, which is slightly smaller than  
 194 1 (we use  $\beta(T) = 0.8$  everywhere) to keep more focus on the structure at the end of the annealing  
 195 as well. In our initial experiments, we observed that  $\beta(T) \in [0.7, 1]$  perform similarly. The scales  
 196 distributions comparison between beta and uniform sampling is provided in Fig 3 and the convergence  
 197 comparison in Fig 6.

### 198 3.4 Training details

199 We inherit the training procedure from StyleGAN2-ADA [24] with minimal changes. The optimiza-  
 200 tion is performed by Adam [27] with learning rate of 0.002 and betas of 0 and 0.99 for both G and  
 201 D. We use  $\beta(T) = 0.8$  for  $T = 10000$ ,  $z \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and set  $R_p = 512$ . D is trained with R1  
 202 regularization [37] with  $\gamma = 0.05$ . We train with the overall batch size of 64 for  $\approx 15\text{M}$  images  
 203 seen by D for  $256^2$  resolution and  $\approx 20\text{M}$  for  $512^2$ . Similar to previous works [6, 12], we use pose  
 204 supervision for D for the FFHQ and Cats dataset to avoid geometry ambiguity. We train G in full  
 205 precision and use mixed precision for D. Further details can be found in the source code.

Table 1: FID scores for recent 3D GANs. “†” — evaluated on a re-aligned version of FFHQ (different from original FFHQ [25]). Training cost is measured as the amount of NVidia V100 GPU days. In contrast to the existing upsampler-based generators or GRAM [12], our model does not have any geometry constraints, trains in 2 – 3× less time and attains state-of-the-art image quality. Note that all those methods (except for  $\pi$ -GAN), appeared (on arxiv, some are not in the proceedings yet) during the past  $\approx 6$  months. “OOM” denotes out-of-memory error.

Method	FFHQ		Cats	M-Plants	M-Food	Training cost		Geometry constraints
	256 <sup>2</sup>	512 <sup>2</sup>	256 <sup>2</sup>	256 <sup>2</sup>	256 <sup>2</sup>	256 <sup>2</sup>	512 <sup>2</sup>	
StyleNeRF [15]	8.00	7.8	27.91	19.32	16.75	40	56	32 <sup>2</sup> -res + 2D upsampler
StyleSDF [47]	11.5	11.19	–	–	–	42	56	64 <sup>2</sup> -res + 2D upsampler
EG3D [6]	4.8 <sup>†</sup>	4.7 <sup>†</sup>	–	–	–	N/A	76	128 <sup>2</sup> -res + 2D upsampler
VolumeGAN [70]	9.1	–	–	–	–	N/A	N/A	64 <sup>2</sup> -res + 2D upsampler
MVCGAN [77]	13.7	13.4	39.16	31.70	29.29	42	64	64 <sup>2</sup> -res + 2D upsampler
GIRAFFE-HD [71]	11.93	–	12.36	–	–	N/A	N/A	16 <sup>2</sup> -res + 2D upsampler
$\pi$ -GAN [7]	53.2	OOM	68.28	46.72	51.99	56	$\infty$	none
GRAM [12]	13.78	OOM	13.40	211.7	248.3	56	$\infty$	iso-surfaces
PGRF (ours)	9.71	9.92	6.93	19.42	18.15	16	24	none

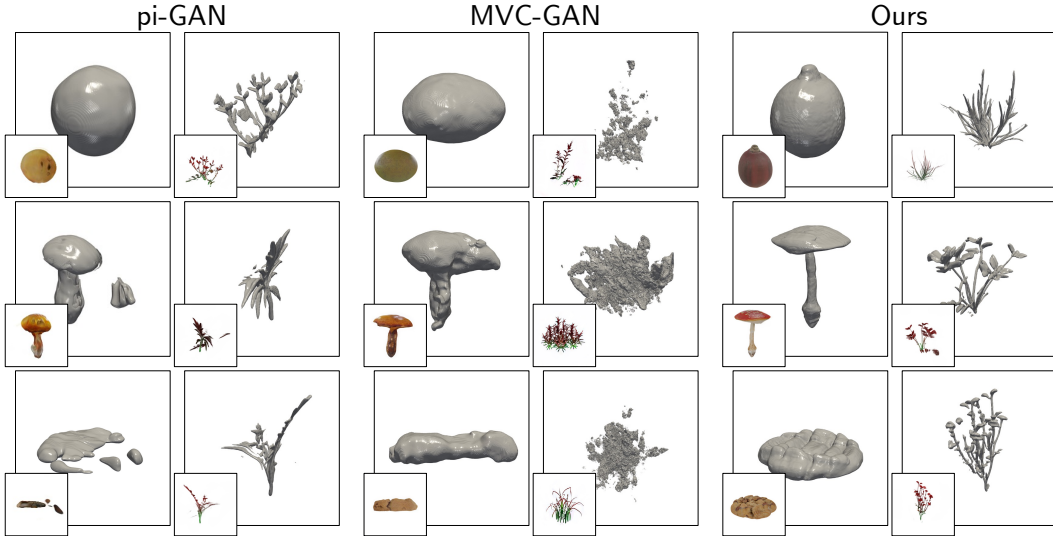


Figure 5: Visualizing the learned geometry for different methods.  $\pi$ -GAN [7] recovers high-fidelity shapes, but has worse image quality (see Table 1) and is much more expensive to train than our model. MVC-GAN [77] fails to capture good geometry because of the 2D upsampler. Our method learns proper geometry and achieves state-of-the-art image quality. We extracted the surfaces using marching cubes from the density fields sampled on 256<sup>3</sup> grid and visualized them in PyVista [62]. We manually optimized the marching cubes contouring threshold for each checkpoint of each method.

## 206 4 Experiments

### 207 4.1 Experimental setup

208 **Benchmarks.** In our study, we consider four benchmarks: 1) FFHQ [25] in 256<sup>2</sup> and 512<sup>2</sup> resolutions,  
 209 consisting of 70,000 (mostly front-view) human face images; 2) Cats 256<sup>2</sup> [76], consisting of 9,998  
 210 (mostly front-view) cat face images; 3) Megascans Food (M-Food) 256<sup>2</sup> consisting of 199 models  
 211 of different food items with 128 views per model (25472 images in total); and 4) Megascans Plants  
 212 (M-Plants) 256<sup>2</sup> consisting of 1108 different plant models with 128 views per model (141824 images  
 213 in total). The last two datasets are introduced in our work to fix two issues with the modern 3D  
 214 generation benchmarks. First, existing benchmarks have low variability of global object geometry,

<sup>1</sup> $\text{lerp}(x, y, \alpha) = (1 - \alpha) \cdot x + \alpha \cdot y$  for  $x, y \in \mathbb{R}$  and  $\alpha \in [0, 1]$ .

<sup>2</sup>In practice, those methods use a *very* slightly different distribution (see Appx B)

215 focusing entirely on a single class of objects, like human/cat faces or cars, that do not vary much from  
216 instance to instance. Second, they all have limited camera pose distribution: for example, FFHQ [25]  
217 and Cats [76] are completely dominated by the frontal and near-frontal views (see Appx D). That’s  
218 why we obtain and render 1307 Megascans models from Quixel, which are photo-realistic (barely  
219 distinguishable from real) scans of real-life objects with extremely difficult geometry (especially  
220 plants). To have a rich pose coverage we produced 128 views for each model, rendered from a fixed  
221 distance to the object center from uniformly sampled points on the entire sphere (even from below).  
222 Those benchmarks, together with the rendering code, will be made publicly available.

223 **Metrics.** We use FID [20] to measure image quality and also estimate the training cost for each  
224 method in terms of NVidia V100 GPU days needed for it to complete the training process.

225 **Baselines.** For upsampler-based baselines, we compare to the following generators: StyleNeRF [15],  
226 StyleSDF [47], EG3D [6], VolumeGAN [70], MVCGAN [77] and GIRAFFE-HD [71]. Apart from  
227 that, we also compare to pi-GAN [7] and GRAM [12], which are non-upsampler-based GANs. To  
228 compare on Megascans, we train StyleNeRF, MVCGAN, pi-GAN, and GRAM from scratch using  
229 their official code repositories (obtained online or requested from the authors), using their FFHQ or  
230 CARLA hyperparameters, except for the camera distribution and rendering settings. We also train  
231 StyleNeRF, MVCGAN and  $\pi$ -GAN on Cats 256<sup>2</sup>. GRAM [12] restricts the sampling space to a  
232 set of learnable iso-surfaces which makes it not well-suited for datasets with varying geometry. To  
233 optimize its hyperparameters for Megascans, we varied the number of surfaces (24 and 48), field of  
234 view (30 and 12) and batch size (16 and 32), but it was diverging for us in each case.

## 235 4.2 Results

236 *PGRF achieves state-of-the-art image quality.* For Cats 256<sup>2</sup>, M-Plants 256<sup>2</sup> and M-Food 256<sup>2</sup>,  
237 PGRF outperforms all the baselines in terms of FID except for StyleNeRF, performing very similar  
238 to it on M-Plants and M-Food, but greatly surpassing it on Cats. For FFHQ, our model attains  
239 very similar FID scores as the other methods, ranking 4/9 (including older  $\pi$ -GAN [7]), noticeably  
240 losing only to EG3D [6], which trains and evaluates on a different version of FFHQ and uses pose  
241 conditioning in the generator (which potentially improves FID at the cost of multi-view consistency).  
242 We provide a visual comparison for different methods in Fig 4.

243 *PGRF is much faster to train.* As reported in Tab 1, existing methods typically train for  $\approx 1$  week on  
244 8 V100s, PGRF finishes training in just 2 days for 256<sup>2</sup> and 3 days for 512<sup>2</sup> resolutions, which is  
245 2 – 3 $\times$  faster. Note that this high training efficiency is achieved without the use of an upsampler,  
246 which initially enabled high-resolution synthesis of 3D-aware GANs. As to the non-upsampler  
247 methods, we couldn’t train GRAM or  $\pi$ -GAN on 512<sup>2</sup> resolution due to the memory limitations of  
248 the setup with 8 NVidia V100 32GB GPUs (i.e., 256GB of GPU memory in total).

249 *PGRF learns high-fidelity geometry.* Using a pure NeRF-based backbone carries two crucial benefits:  
250 it provides multi-view consistency and allows to learn the geometry in the full dataset resolution. In  
251 Fig 5, we visualize the learned shapes on M-Food and M-Plants for 1)  $\pi$ -GAN: a pure NeRF-based  
252 generator without the geometry constraints; 2) MVC-GAN [77]: an upsampler-based generator with  
253 strong multi-view consistency regularization; 3) our model. We provide the details and analysis in  
254 the caption of Fig 5.

255 *PGRF easily capitalizes on techniques from the NeRF literature.* Since our generator is purely NeRF  
256 based and renders images without a 2D upsampler, it is well coupled with the existing techniques from  
257 the NeRF scene reconstruction field. To demonstrate this, we adopted background separation from  
258 NeRF++ [75] using the inverse sphere parametrization by simply copy-pasting the corresponding  
259 code from their repo. We depict the results in Fig 1 and provide the details in Appx B.

## 260 4.3 Ablations

261 We report the ablations for different discriminator architectures and patch sizes on FFHQ 512<sup>2</sup>  
262 and M-Plants 256<sup>2</sup> in Tab 2. Using a traditional discriminator architecture results in  $\approx 15\%$  worse  
263 performance. Using several ones (via the group-wise convolution trick [26]) results in noticeably  
264 slower training time and degrades the image quality a lot. We hypothesize that the reason of it was  
265 the reduced overall training signal which each discriminator receives, which we tried to alleviate by  
266 increasing the learning rate for them, but that did not improve the results. A too small patch size

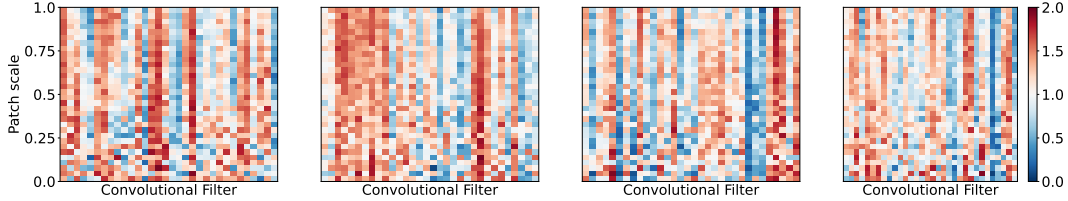
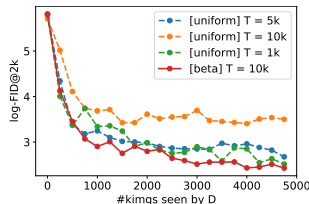


Figure 7: Visualizing modulation weights  $\sigma$ , predicted by H for 2-nd, 6-th, 10-th and 14-th convolutional layers. Each subplot denotes a separate layer and we visualize random 32 filters for it.

267 hampers the learning process and results in a  $\approx 80\%$  worse FID. A too large one provides decent  
 268 image quality, but greatly reduces the training speed.

269 To assess the convergence of our proposed patch sampling scheme, we compared against uniform  
 270 sampling on Cats  $256^2$  for  $T \in \{1000, 5000, 10000\}$ , representing different annealing speeds. We  
 271 show the results for it in Fig 6: our proposed beta scale sampling strategy with  $T = 10k$  schedule  
 272 robustly converges to lower values than the uniform one with  $T = 5k$  or  $T = 10k$  and does not  
 273 fluctuate much compared to the  $T = 1k$  uniform one (where the model reached its final annealing  
 274 stage in just 1k kilo-images seen by D).

275 To analyze how hyper-modulation manipulates the convolutional filters of the discriminator, we  
 276 visualize the modulation weights  $\sigma$ , predicted by H, in Fig 7 (see the caption for the details). These  
 277 visualizations show that some of the filters are always switched on, regardless of the patch scale;  
 278 while others are always switched off providing potential room for pruning [18]. And  $\approx 40\%$  of the  
 279 filters are getting switched on and off depending on the patch scale, which shows that H indeed learns  
 280 to perform meaningful modulation.



281 Figure 6: Convergence comparison on Cats  $256^2$  for PGRF between uniform and beta scale sampling strategies in terms of log-FID measured on 2048 fake images.

Experiment	FFHQ $512^2$	M-Plants $256^2$	Training cost
Discriminator			
- standard	11.57	21.77	24
- 2 scale-specific D-s	10.87	21.02	28
- 4 scale-specific D-s	21.56	43.11	28
- scale/position-aware D	9.92	19.42	24
Patch size			
- $32^2$	17.44	34.32	19
- $64^2$ (default)	9.92	19.42	24
- $128^2$	11.36	18.90	34

Table 2: Ablating the patch size and the discriminator architecture for our model in terms of FID scores and training cost (V100 GPU days) on  $512^2$  resolution.

## 282 5 Conclusion

283 In this work, we showed that it is possible to build a state-of-the-art 3D GAN framework without a  
 284 2D upsampler, but using a pure NeRF-based generator trained in a multi-scale patch-wise fashion.  
 285 For this, we improved the traditional patch-wise training scheme in two important ways. First,  
 286 we proposed to use a scale/location-aware discriminator with convolutional filters modulated by a  
 287 hypernetwork depending on the patch parameters. Second, we developed a schedule for patch scale  
 288 sampling based on the beta distribution, that leads to faster and more robust convergence. We believe  
 289 that the future of 3D GANs is a combination of efficient volumetric representations, regularized 2D  
 290 upsamplers, and patch-wise training. We propose this avenue of research for future work.

291 Our method also has several limitations. Before switching to training 3D-aware generators, we  
 292 spent a considerable amount of time exploring our ideas on top of StyleGAN2 for traditional 2D  
 293 generation. This always resulted in increased FID scores (see Appx A). Further, the discriminator  
 294 loses information about global context. We tried multiple ideas to incorporate global context, but it  
 295 did not lead to an improvement. Finally, 3D GANs generating faces and humans may have negative  
 296 societal impact as discussed in Appx G.



297 **References**

- 298 [1] P. Achlioptas, O. Diamanti, I. Mitliagkas, and L. Guibas. Learning representations and generative models  
299 for 3d point clouds. In *International conference on machine learning*, pages 40–49. PMLR, 2018.
- 300 [2] J. T. Barron, B. Mildenhall, M. Tancik, P. Hedman, R. Martin-Brualla, and P. P. Srinivasan. Mip-nerf:  
301 A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF*  
302 *International Conference on Computer Vision*, pages 5855–5864, 2021.
- 303 [3] Blender Online Community. *Blender - a 3D modelling and rendering package*. Blender Foundation,  
304 Blender Institute, Amsterdam, 2022.
- 305 [4] A. Brock, J. Donahue, and K. Simonyan. Large scale gan training for high fidelity natural image synthesis.  
306 *arXiv preprint arXiv:1809.11096*, 2018.
- 307 [5] L. Chai, M. Gharbi, E. Shechtman, P. Isola, and R. Zhang. Any-resolution training for high-resolution  
308 image synthesis. *arXiv preprint arXiv:2204.07156*, 2022.
- 309 [6] E. R. Chan, C. Z. Lin, M. A. Chan, K. Nagano, B. Pan, S. D. Mello, O. Gallo, L. Guibas, J. Tremblay,  
310 S. Khamis, T. Karras, and G. Wetzstein. Efficient geometry-aware 3D generative adversarial networks. In  
311 *arXiv*, 2021.
- 312 [7] E. R. Chan, M. Monteiro, P. Kellnhofer, J. Wu, and G. Wetzstein. pi-gan: Periodic implicit generative  
313 adversarial networks for 3d-aware image synthesis. In *Proceedings of the IEEE/CVF Conference on*  
314 *Computer Vision and Pattern Recognition*, pages 5799–5809, 2021.
- 315 [8] A. Chen, Z. Xu, A. Geiger, J. Yu, and H. Su. Tensorf: Tensorial radiance fields. *arXiv preprint*  
316 *arXiv:2203.09517*, 2022.
- 317 [9] H. Chen, B. He, H. Wang, Y. Ren, S.-N. Lim, and A. Shrivastava. Nerv: Neural representations for videos.  
318 *arXiv preprint arXiv:2110.13903*, 2021.
- 319 [10] Z. Chen and H. Zhang. Learning implicit fields for generative shape modeling. In *Proceedings of the*  
320 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5939–5948, 2019.
- 321 [11] Y. Choi, M. Choi, M. Kim, J.-W. Ha, S. Kim, and J. Choo. Stargan: Unified generative adversarial networks  
322 for multi-domain image-to-image translation. In *Proceedings of the IEEE conference on computer vision*  
323 *and pattern recognition*, pages 8789–8797, 2018.
- 324 [12] Y. Deng, J. Yang, J. Xiang, and X. Tong. Gram: Generative radiance manifolds for 3d-aware image  
325 generation. In *IEEE Computer Vision and Pattern Recognition*, 2022.
- 326 [13] M. Gadelha, S. Maji, and R. Wang. 3d shape induction from 2d views of multiple objects. In *2017*  
327 *International Conference on 3D Vision (3DV)*, pages 402–411. IEEE, 2017.
- 328 [14] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio.  
329 Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014.
- 330 [15] J. Gu, L. Liu, P. Wang, and C. Theobalt. Stylenerf: A style-based 3d aware generator for high-resolution  
331 image synthesis. In *International Conference on Learning Representations*, 2022.
- 332 [16] D. Ha, A. Dai, and Q. V. Le. Hypernetworks. *arXiv preprint arXiv:1609.09106*, 2016.
- 333 [17] Z. Hao, A. Mallya, S. Belongie, and M.-Y. Liu. GANcraft: Unsupervised 3D Neural Rendering of  
334 Minecraft Worlds. In *ICCV*, 2021.
- 335 [18] Y. He, X. Zhang, and J. Sun. Channel pruning for accelerating very deep neural networks. In *Proceedings*  
336 *of the IEEE international conference on computer vision*, pages 1389–1397, 2017.
- 337 [19] P. Henderson, V. Tsiminaki, and C. H. Lampert. Leveraging 2d data to learn textured 3d mesh generation. In  
338 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7498–7507,  
339 2020.
- 340 [20] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter. Gans trained by a two time-scale  
341 update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30,  
342 2017.
- 343 [21] J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information*  
344 *Processing Systems*, 33:6840–6851, 2020.

- 345 [22] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial  
346 networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages  
347 1125–1134, 2017.
- 348 [23] K. Jo, G. Shim, S. Jung, S. Yang, and J. Choo. Cg-nerf: Conditional generative neural radiance fields.  
349 *arXiv preprint arXiv:2112.03517*, 2021.
- 350 [24] T. Karras, M. Aittala, J. Hellsten, S. Laine, J. Lehtinen, and T. Aila. Training generative adversarial  
351 networks with limited data. *arXiv preprint arXiv:2006.06676*, 2020.
- 352 [25] T. Karras, S. Laine, and T. Aila. A style-based generator architecture for generative adversarial networks. In  
353 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4401–4410,  
354 2019.
- 355 [26] T. Karras, S. Laine, M. Aittala, J. Hellsten, J. Lehtinen, and T. Aila. Analyzing and improving the  
356 image quality of stylegan. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
357 Recognition*, pages 8110–8119, 2020.
- 358 [27] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*,  
359 2014.
- 360 [28] D. P. Kingma and P. Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. *Advances in  
361 neural information processing systems*, 31, 2018.
- 362 [29] A. R. Kosiorok, H. Strathmann, D. Zoran, P. Moreno, R. Schneider, S. Mokrá, and D. J. Rezende. Nerf-vae:  
363 A geometry aware 3d scene generative model. *arXiv preprint arXiv:2104.00587*, 2021.
- 364 [30] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural  
365 networks. *Advances in neural information processing systems*, 25, 2012.
- 366 [31] R. Li, X. Li, K.-H. Hui, and C.-W. Fu. Sp-gan: Sphere-guided 3d shape generation and manipulation.  
367 *ACM Transactions on Graphics (TOG)*, 40(4):1–12, 2021.
- 368 [32] X. Li, Y. Dong, P. Peers, and X. Tong. Synthesizing 3d shapes from silhouette image collections using multi-  
369 projection generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer  
370 Vision and Pattern Recognition*, pages 5535–5544, 2019.
- 371 [33] C. H. Lin, H.-Y. Lee, Y.-C. Cheng, S. Tulyakov, and M.-H. Yang. Infinitygan: Towards infinite-resolution  
372 image synthesis. *arXiv preprint arXiv:2104.03963*, 2021.
- 373 [34] S. Luo and W. Hu. Diffusion probabilistic models for 3d point cloud generation. In *Proceedings of the  
374 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021.
- 375 [35] Y. A. Mejjati, I. Milefchik, A. Gokaslan, O. Wang, K. I. Kim, and J. Tompkin. Gaussigan: Controllable  
376 image synthesis with 3d gaussians from unposed silhouettes. *arXiv preprint arXiv:2106.13215*, 2021.
- 377 [36] Q. Meng, A. Chen, H. Luo, M. Wu, H. Su, L. Xu, X. He, and J. Yu. Gnerf: Gan-based neural radiance field  
378 without posed camera. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,  
379 pages 6351–6361, 2021.
- 380 [37] L. Mescheder, A. Geiger, and S. Nowozin. Which training methods for gans do actually converge? In  
381 *International conference on machine learning*, pages 3481–3490. PMLR, 2018.
- 382 [38] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng. Nerf: Representing  
383 scenes as neural radiance fields for view synthesis. In *European conference on computer vision*, pages  
384 405–421. Springer, 2020.
- 385 [39] P. Mittal, Y.-C. Cheng, M. Singh, and S. Tulsiani. AutoSDF: Shape priors for 3d completion, reconstruction  
386 and generation. In *CVPR*, 2022.
- 387 [40] T. Miyato and M. Koyama. cgans with projection discriminator. *arXiv preprint arXiv:1802.05637*, 2018.
- 388 [41] T. Nguyen-Phuoc, C. Li, L. Theis, C. Richardt, and Y.-L. Yang. Hologan: Unsupervised learning of 3d  
389 representations from natural images. In *The IEEE International Conference on Computer Vision (ICCV)*,  
390 Nov 2019.
- 391 [42] M. Niemeyer and A. Geiger. Campari: Camera-aware decomposed generative neural radiance fields. In  
392 *2021 International Conference on 3D Vision (3DV)*, pages 951–961. IEEE, 2021.

- 393 [43] M. Niemeyer and A. Geiger. Giraffe: Representing scenes as compositional generative neural feature  
394 fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages  
395 11453–11464, 2021.
- 396 [44] M. Niemeyer, L. Mescheder, M. Oechsle, and A. Geiger. Differentiable volumetric rendering: Learning  
397 implicit 3d representations without 3d supervision. In *Proceedings of the IEEE/CVF Conference on*  
398 *Computer Vision and Pattern Recognition*, pages 3504–3515, 2020.
- 399 [45] M. Oechsle, S. Peng, and A. Geiger. Unisurf: Unifying neural implicit surfaces and radiance fields for  
400 multi-view reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,  
401 pages 5589–5599, 2021.
- 402 [46] C. Olah, A. Mordvintsev, and L. Schubert. Feature visualization. *Distill*, 2017.  
403 <https://distill.pub/2017/feature-visualization>.
- 404 [47] R. Or-El, X. Luo, M. Shan, E. Shechtman, J. J. Park, and I. Kemelmacher-Shlizerman. StyleSDF:  
405 High-Resolution 3D-Consistent Image and Geometry Generation. *arXiv preprint arXiv:2112.11427*, 2021.
- 406 [48] X. Pan, B. Dai, Z. Liu, C. C. Loy, and P. Luo. Do 2d gans know 3d shape? unsupervised 3d shape  
407 reconstruction from 2d image gans. *arXiv preprint arXiv:2011.00844*, 2020.
- 408 [49] X. Pan, X. Xu, C. C. Loy, C. Theobalt, and B. Dai. A shading-guided generative implicit model for shape-  
409 accurate 3d-aware image synthesis. In *Advances in Neural Information Processing Systems (NeurIPS)*,  
410 2021.
- 411 [50] K. Park, U. Sinha, J. T. Barron, S. Bouaziz, D. B. Goldman, S. M. Seitz, and R. Martin-Brualla. Deformable  
412 neural radiance fields. *arXiv preprint arXiv:2011.12948*, 2020.
- 413 [51] T. Park, M.-Y. Liu, T.-C. Wang, and J.-Y. Zhu. Semantic image synthesis with spatially-adaptive normal-  
414 ization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages  
415 2337–2346, 2019.
- 416 [52] T. Park, J.-Y. Zhu, O. Wang, J. Lu, E. Shechtman, A. Efros, and R. Zhang. Swapping autoencoder for deep  
417 image manipulation. *Advances in Neural Information Processing Systems*, 33:7198–7211, 2020.
- 418 [53] D. Pavllo, J. Kohler, T. Hofmann, and A. Lucchi. Learning generative models of textured 3d meshes from  
419 real-world images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*,  
420 pages 13879–13889, October 2021.
- 421 [54] D. Pavllo, G. Spinks, T. Hofmann, M.-F. Moens, and A. Lucchi. Convolutional generation of textured 3d  
422 meshes. In *Neural Information Processing Systems (NeurIPS)*, 2020.
- 423 [55] A. Ramesh, P. Dhariwal, A. Nichol, C. Chu, and M. Chen. Hierarchical text-conditional image generation  
424 with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- 425 [56] K. Schwarz, Y. Liao, M. Niemeyer, and A. Geiger. Graf: Generative radiance fields for 3d-aware image  
426 synthesis. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- 427 [57] T. R. Shaham, T. Dekel, and T. Michaeli. Singan: Learning a generative model from a single natural image.  
428 In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4570–4580, 2019.
- 429 [58] Y. Shi, D. Aggarwal, and A. K. Jain. Lifting 2d stylegan for 3d-aware face generation. In *Proceedings of*  
430 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6258–6266, 2021.
- 431 [59] V. Sitzmann, J. Martel, A. Bergman, D. Lindell, and G. Wetzstein. Implicit neural representations with  
432 periodic activation functions. *Advances in Neural Information Processing Systems*, 33, 2020.
- 433 [60] I. Skorokhodov, S. Ignatyev, and M. Elhoseiny. Adversarial generation of continuous images. In *Pro-*  
434 *ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10753–10764,  
435 2021.
- 436 [61] I. Skorokhodov, G. Sotnikov, and M. Elhoseiny. Aligning latent and image spaces to connect the uncon-  
437 nectable. *arXiv preprint arXiv:2104.06954*, 2021.
- 438 [62] C. B. Sullivan and A. Kaszynski. PyVista: 3d plotting and mesh analysis through a streamlined interface  
439 for the visualization toolkit (VTK). *Journal of Open Source Software*, 4(37):1450, may 2019.
- 440 [63] F. Tan, S. Fanello, A. Meka, S. Orts-Escolano, D. Tang, R. Pandey, J. Taylor, P. Tan, and Y. Zhang.  
441 Volux-gan: A generative model for 3d face synthesis with hdri relighting. *arXiv preprint arXiv:2201.04873*,  
442 2022.

- 443 [64] M. Tancik, P. P. Srinivasan, B. Mildenhall, S. Fridovich-Keil, N. Raghavan, U. Singhal, R. Ramamoorthi,  
444 J. T. Barron, and R. Ng. Fourier features let networks learn high frequency functions in low dimensional  
445 domains. *arXiv preprint arXiv:2006.10739*, 2020.
- 446 [65] A. Van den Oord, N. Kalchbrenner, L. Espeholt, O. Vinyals, A. Graves, et al. Conditional image generation  
447 with pixelcnn decoders. *Advances in neural information processing systems*, 29, 2016.
- 448 [66] Y. Vinker, E. Horwitz, N. Zabari, and Y. Hoshen. Image shape manipulation from a single augmented  
449 training sample. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*,  
450 pages 13769–13778, October 2021.
- 451 [67] C. Wang, M. Chai, M. He, D. Chen, and J. Liao. Clip-nerf: Text-and-image driven manipulation of neural  
452 radiance fields. *arXiv preprint arXiv:2112.05139*, 2021.
- 453 [68] P. Wang, L. Liu, Y. Liu, C. Theobalt, T. Komura, and W. Wang. Neus: Learning neural implicit surfaces by  
454 volume rendering for multi-view reconstruction. *arXiv preprint arXiv:2106.10689*, 2021.
- 455 [69] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, and J. Xiao. 3d shapenets: A deep representation for  
456 volumetric shapes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
457 pages 1912–1920, 2015.
- 458 [70] Y. Xu, S. Peng, C. Yang, Y. Shen, and B. Zhou. 3d-aware image synthesis via learning structural and  
459 textural representations. *arXiv preprint arXiv:2112.10759*, 2021.
- 460 [71] Y. Xue, Y. Li, K. K. Singh, and Y. J. Lee. Giraffe hd: A high-resolution 3d-aware generative model. *arXiv*  
461 *preprint arXiv:2203.14954*, 2022.
- 462 [72] Y. Ye, S. Tulsiani, and A. Gupta. Shelf-supervised mesh prediction in the wild. In *Proceedings of the*  
463 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8843–8852, June  
464 2021.
- 465 [73] A. Yu, V. Ye, M. Tancik, and A. Kanazawa. pixelnerf: Neural radiance fields from one or few images. In  
466 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages  
467 4578–4587, June 2021.
- 468 [74] J. Zhang, E. Sangineto, H. Tang, A. Siarohin, Z. Zhong, N. Sebe, and W. Wang. 3d-aware semantic-guided  
469 generative model for human synthesis. *arXiv preprint arXiv:2112.01422*, 2021.
- 470 [75] K. Zhang, G. Riegler, N. Snavely, and V. Koltun. Nerf++: Analyzing and improving neural radiance fields.  
471 *arXiv preprint arXiv:2010.07492*, 2020.
- 472 [76] W. Zhang, J. Sun, and X. Tang. Cat head detection-how to effectively exploit shape and texture features.  
473 In *European conference on computer vision*, pages 802–816. Springer, 2008.
- 474 [77] X. Zhang, Z. Zheng, D. Gao, B. Zhang, P. Pan, and Y. Yang. Multi-view consistent generative adversarial  
475 networks for 3d-aware image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision*  
476 *and Pattern Recognition*, 2022.
- 477 [78] P. Zhou, L. Xie, B. Ni, and Q. Tian. Cips-3d: A 3d-aware generator of gans based on conditionally-  
478 independent pixel synthesis. *arXiv preprint arXiv:2110.09788*, 2021.

## 479 Checklist

- 480 1. For all authors...
- 481 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
482 contributions and scope? [Yes]
- 483 (b) Did you describe the limitations of your work? [Yes] See §5 and Appx A.
- 484 (c) Did you discuss any potential negative societal impacts of your work? [Yes] We do  
485 this in Appendix F.
- 486 (d) Have you read the ethics review guidelines and ensured that your paper conforms  
487 to them? [Yes] We discuss the potential ethical concerns of using our model in  
488 Appendix F.
- 489 2. If you are including theoretical results...
- 490 (a) Did you state the full set of assumptions of all theoretical results? [N/A]

- 491 (b) Did you include complete proofs of all theoretical results? [N/A]
- 492 3. If you ran experiments...
- 493 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
- 494 mental results (either in the supplemental material or as a URL)? [Yes] We provide
- 495 the code/data and additional visualizations on <https://rethinking-3d-gans.github.io>(as
- 496 specified in the introduction).
- 497 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
- 498 were chosen)? [Yes] We provide the most important training details in §3.4. The rest
- 499 of the details are provided in Appx B and the provided source code.
- 500 (c) Did you report error bars (e.g., with respect to the random seed after running exper-
- 501 iments multiple times)? [No] . That's too computationally expensive and single-run
- 502 results are typically reliable in the GAN field.
- 503 (d) Did you include the total amount of compute and the type of resources used (e.g., type
- 504 of GPUs, internal cluster, or cloud provider)? [Yes] We report this numbers in Appx B.
- 505 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 506 (a) If your work uses existing assets, did you cite the creators? [Yes] We cite all the
- 507 sources of the datasets which were used or mentioned in our submission.
- 508 (b) Did you mention the license of the assets? [Yes] In this work, we release two new
- 509 datasets: Megascans Plants and Megascans Food. We discuss their licensing in Appx D.
- 510 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 511 We provide our datasets on the project website: <https://rethinking-3d-gans.github.io>.
- 512 (d) Did you discuss whether and how consent was obtained from people whose data you're
- 513 using/curating? [Yes] We specify the information on dataset collection in Appx D.
- 514 (e) Did you discuss whether the data you are using/curating contains personally identifiable
- 515 information or offensive content? [Yes] As discussed in Appx D, the released data
- 516 does not contain personally identifiable information or offensive content.
- 517 5. If you used crowdsourcing or conducted research with human subjects...
- 518 (a) Did you include the full text of instructions given to participants and screenshots, if
- 519 applicable? [N/A]
- 520 (b) Did you describe any potential participant risks, with links to Institutional Review
- 521 Board (IRB) approvals, if applicable? [N/A]
- 522 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 523 spent on participant compensation? [N/A]



524 **A Limitations**

525 Multi-scale patch-wise training, studied in this work, has both practical and theoretical limitations.  
 526 The practical ones include “engineering” difficulties when trying to squeeze the best performance,  
 527 and theoretical ones are related to the issues which could be faced in asymptotic cases.

528 **A.1 Practical limitations**

529 **Performance drop for 2D generation.** Before switching to training 3D-aware generators, we spent  
 530 a considerable amount of time, exploring our ideas on top of StyleGAN2 [24] for traditional 2D  
 531 generation since it is faster, less error-prone and more robust to a hyperparameters choice. What we  
 532 observed is that despite our best efforts (see C) and even with longer training, we couldn’t obtain the  
 533 same image quality as the full resolution StyleGAN2 generator.

Table 3: Trying to train a traditional StyleGAN2 [26] generator in the patch-wise fashion. We tried to train longer to compensate for a smaller learning signal overall (a  $64^2$  patch is  $1/64$  of information compared to a  $512^2$  image), but this didn’t allow to catch up. Note, however, that AnyResGAN [5] reaches SotA when training on  $256^2$  patches compared to  $1024^2$  images.

Method	FFHQ $512^2$		LSUN Bedroom $256^2$	
	FID	Training cost	FID	Training cost
StyleGAN2-ADA [24]	3.83	8	4.12	5
+ multi-scale $64^2$ patch-wise training	7.11	6	6.73	4
+ $\times 2$ longer training	5.71	12	5.42	8
+ $\times 4$ longer training	4.76	24	4.31	16

534 **A range of possible patch sizes is restricted.** Tab 2 shows the performance drop when using the  
 535  $32^2$  patch size instead of the default  $64^2$  one without any dramatic improvement in speed. Trying to  
 536 decrease it further would produce even worse performance (imagine training in the extreme case of  
 537  $2^2$  patches). Increasing the patch size is also not desirable since it decreases the training speed a lot:  
 538 going from  $64^2$  to  $128^2$  resulted in 30% cost increase without clear performance benefits. In this way,  
 539 we are very constrained in what patch size one can use.

540 **Discriminator does not see the global context.** When the discriminator classifies patches of small  
 541 scale, it is forced to do so without relying on the global image information, which could be useful  
 542 for this. Our attempts to incorporate it (see Appx C) did not improve the performance (though we  
 543 believe we under-explored this).

544 **A.2 Theoretical limitations**

545 Multi-scale patch-wise training is not theoretically equivalent to full-resolution training. Imagine  
 546 that we have a distribution  $p(\mathbf{x})$  for  $\mathbf{x} \in \mathbb{R}^{R \times R}$ , where we consider images to be single-channel (for  
 547 simplicity) and  $R$  is the image size. If we train with patch size of  $r$ , then each optimization step uses  
 548 random  $r \times r$  pixels out of  $R \times R$ , i.e. we optimize over the distribution of all possible marginals  
 549  $p(\mathbf{x}_p)$  where  $\mathbf{x}_p = e(\mathbf{x}; \xi) \in \mathbb{R}^{r \times r}$  is an image patch and  $e(\mathbf{x}, \xi)$  is a patch extraction function with  
 550 random seed  $\xi$ .<sup>3</sup> This means that our minimax objective becomes:

$$\min_G \max_D \mathbb{E}_{p(\mathbf{x}_p)} [\log D(\mathbf{x}_p)] + \mathbb{E}_{p(\mathbf{z}), p(\xi)} [\log(1 - D(e(G(\mathbf{z}), \xi)))] \quad (6)$$

551 If we rely on the GAN convergence theorem [14], stating that we recover the training distribution  
 552 as the solution of the minimax problem, then G will learn to approximate all the possible marginal  
 553 distributions  $p(\mathbf{x}_p)$  instead of the full joint distribution  $p(\mathbf{x})$ , that we seek.

<sup>3</sup>For brevity, we “hide” all the randomness of the patch extraction process into  $\xi$ .

## 554 B Training details

### 555 B.1 Hyper-parameters and optimization details

556 We inherit most of the hyperparameters from the StyleGAN2-ADA repo [24] repo which we build  
557 on top<sup>4</sup>. In this way, we use the dimensionalities of 512 for both  $z$  and  $w$ . The mapping network  
558 has 2 layers of dimensionality 512 with LeakyReLU non-linearities with the negative slope of  $-0.2$ .  
559 Synthesis network S produced three  $512^2$  planes of 32 channels each. We use the SoftPlus non-  
560 linearity instead of typically used ReLU [38] as a way to clamp the volumetric density. Similar to  
561  $\pi$ -GAN, we also randomize

562 For FFHQ and Cats, we also use camera conditioning in D. For this, we encode yaw and pitch angles  
563 (roll is always set to 0) with Fourier positional encoding [59, 64], apply dropout with 0.5 probability  
564 (otherwise, D can start judging generations from 3D biases in the dataset, hurting the image quality),  
565 pass through a 2-layer MLP with LeakyReLU activations to obtain a 512-dimensional vector, which  
566 is finally as a projection conditioning [40]. Cameras positions were extracted in the same way as in  
567 GRAM [12].

568 We optimize both G and D with the batch size of 64 until D sees 25,000,000 real images, which is  
569 the default setting from StyleGAN2-ADA. We the default setup of adaptive augmentations, except  
570 for random horizontal flipping, since it would require the corresponding change in the yaw angle  
571 at augmentation time, which was not convenient to incorporate from the engineering perspective.  
572 Instead, random horizontal flipping is used non-adaptively as a dataset mirroring where flipping the  
573 yaw angles it more accessible. We train G in full precision, while D uses mixed precision.

574 Hypernetwork H is structured very similar to the generator’s mapping network. It consists on 2 layers  
575 with LeakyReLU non-linearities with the negative slope of  $-0.2$ . Its input is the positional embedding  
576 of the patch scales and offsets  $s, \delta_x, \delta_y$ , encoded with Fourier features [59, 64] and concatenated into  
577 a single vector of dimensionality 828. It produces a patch representation vector  $p \in \mathbb{R}^{512}$ , which is  
578 then adapted for each convolutional layer via:

$$\sigma = \tanh(W_\ell p + b_\ell) + 1, \tag{7}$$

579 where  $\sigma \in [0, 2]^{c_{\text{out}}^\ell}$  is the modulation vector,  $(W_\ell, b_\ell)$  is the layer-specific affine transformation,  $c_{\text{out}}^\ell$   
580 is the amount of output filters in the  $\ell$ -th layer. In this way, H has layer-specific adapters.

581 For the background separation experiment, we adapt the neural representation MLP from INR-  
582 GAN [60], but passing 4 coordinates (for the inverse sphere parametrization [75]) instead of 2 as an  
583 input. It consists on 2 blocks with 2 linear layers each. We use 16 steps per ray for the background  
584 without hierarchical sampling.

585 Further details could be found in the accompanying source code.

### 586 B.2 Utilized computational resources

587 While developing our model, we had been launching experiments on  $4 \times$  NVidia A100 81GB or  
588 Nvidia V100 32GB GPUs with the AMD EPYC 7713P 64-Core processor. We found that in practice,  
589 running the model on A100s gives a  $2 \times$  speed-up compared to V100s due to the possibility of  
590 increasing the batch size from 32 to 64. In this way, training PGRF on  $4 \times$  A100s gives the same  
591 training speed as training it  $8 \times$  V100s.

592 For the baselines, we were running them on  $4-8 \times$  V100s GPUs as was specified by the original  
593 papers unless the model could fit into 4 V100s without decreasing the batch size (it was only possible  
594 for StyleNeRF [15]).

595 For rendering Megascans, we used  $4 \times$  NVIDIA TITAN RTX with 24GB memory each. But resource  
596 utilization for rendering is negligible compared to training the generators.

597 In total, the project consumed  $\approx 4$  A100s GPU-years,  $\approx 4$  V100s GPU-years, and  $\approx 20$  TITAN RTX  
598 GPU-days. Note, that out of this time, training the baselines consumed  $\approx 1.5$  V100s GPU-years.

---

<sup>4</sup><https://github.com/NVlabs/stylegan2-ada-pytorch>

599 **B.3 Annealing schedule details**

600 As being said in §3.3, the existing multi-scale patch-wise generators [56, 36] use uniform distribution  
601  $U[s_{\min}(t), 1]$  to sample patch scales, where  $s_{\min}(t)$  is gradually annealed during training from 0.9  
602 (or 0.8 [36]) to  $r/R$  with different speeds. We visualize the annealing schedule for both GRAF  
603 and GNeRF on Fig 8, which demonstrates that their schedules are very close to lerp-based one,  
604 described in §3.3.

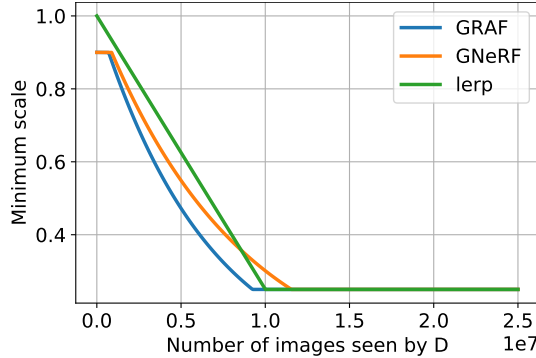


Figure 8: Comparing annealing schedules for GRAF [56], GNeRF [36] and the lerp-based schedule from §3.3. We simplified the exposition by stating that GRAF and GNeRF use the lerp-based schedule, which is *very* close to reality.

605 **C Failed experiments**

606 Modern GANs are a lot of engineering and it often takes a lot of futile experiments to get to a point  
607 where the obtained performance is acceptable. We want to enumerate some experiments which did  
608 not work out (despite looking like they should work) — either because the idea was fundamentally  
609 flawed on its own or because we’ve under-explored it (or both).

610 **Conditioning D on global context worsened the performance.** In Appx A, we argued that when D  
611 processes a small-scale patch, it does not have access to the global image information, which might  
612 be a source of decreased image quality. We tried several strategies to compensate for this. Our first  
613 attempt was to generate a low-resolution image, bilinearly upsample it to the target size, and then  
614 “grid paste” a high-resolution patch into it. The second attempt was to simply always concatenating  
615 a low-resolution version of an image as 3 additional channels. However, in both cases, generator  
616 learned to produce low-resolution version of images well, but the texture was poor. We hypothesize  
617 that it was due to D starting to produce its prediction almost entirely based on the low-resolution  
618 image, ignoring the high-resolution patches since they are harder to discriminate.

619 **Patch importance sampling did not work.** Almost all the datasets used for 3D-aware image  
620 synthesis has regions of difficult content and regions with simpler content — it is especially noticeable  
621 for CARLA [56] and our Megascans datasets, which contain a lot of white background. That’s why,  
622 patch-wise sampling could be improved if we sample patches from the more difficult regions more  
623 frequently. We tried this strategy in the GNeRF [36] problem setup on the NeRF-Synthetic dataset [38]  
624 of fitting a scene without known camera parameters. We sampled patches from regions with high  
625 average gradient norm more frequently. For some scenes, it helped, for other ones, it worsened the  
626 performance.

627 **View direction conditioning breaks multi-view consistency.** Similar to the prior works [56, 7],  
628 our attempt to condition the radiance (but not density) MLP on ray direction (similar to NeRF [38])  
629 led to poor multi-view consistency with radiance changing with camera moving. We tested this  
630 on FFHQ [25], which has only a single view per object instance and suspect that it wouldn’t be  
631 happening on Megascans, where view coverage is very rich.

632 **Tri-planes produced from convolutional layers are much harder to optimize for reconstruction.**  
633 While debugging our tri-plane representation, we found that tri-planes produced with convolutional  
634 layers are extremely difficult to optimize for reconstruction. I.e., if one fits a 3D scene while

635 optimizing tri-planes directly, then everything goes smoothly, but when those tri-plane are being  
636 produced by the synthesis network of StyleGAN2 [26], then PNSR scores (and the loss values) are  
637 plateauing very soon.

## 638 **D Datasets details**

### 639 **D.1 Megascans dataset**

640 Modern 3D-aware image synthesis benchmarks has two issues: 1) they contain objects of very similar  
641 global geometry (like, human or cat faces, cars and chairs), and 2) they have poor camera coverage.  
642 Moreover, some of them (e.g., FFHQ), contain 3D-biases, when an object features (e.g., smiling  
643 probability, gaze direction, posture or haircut) correlate with the camera position [6]. As a result, this  
644 does not allow to evaluate a model’s ability to represent the underlying geometry and makes it harder  
645 understand whether performance come from methodological changes or better data preprocessing.

646 To mitigate these issues, we introduce two new datasets: Megascans Plants (M-Plants) and Megascans  
647 Food (M-Food). To build them, we obtain  $\approx 1,500$  models from Quixel Megascans<sup>5</sup> from Plants,  
648 Mushrooms and Food categories. Megascans are very high-quality scans of real objects which are  
649 almost indistinguishable from real. For Mushrooms and Plants, we merge them into the same Food  
650 category since they have too few models on their own.

651 We render all the models in Blender [3] with cameras, distributed uniformly at random over the  
652 sphere of radius 3.5 and field-of-view of  $\pi/4$ . While rendering, we scale each model into  $[-1, 1]^3$   
653 cube and discard those models, which has the dimension produce of less than 2. We render 128  
654 views per object. For M-Plants, we additionally remove those models which has less than 0.03 pixel  
655 intensity on average (computed as the mean alpha value over the pixels and views). This is needed to  
656 remove small grass or leaves which will be occupying a too small amount of pixels. As a result, this  
657 procedure produces 1,108 models for the Plants category and 199 models for the Food category.

658 We include the rendering script as a part of the released source code. We cannot release the source  
659 models or textures due to the copyright restrictions. We release all the images under the CC BY-NC-  
660 SA 4.0 license<sup>6</sup>. Apart from the images, we also release the class categories for both M-Plants and  
661 M-Food.

662 The released datasets does not contain any personally identifiable information or offensive content  
663 since it does not have any human subjects, animals or other creatures with scientifically proved  
664 cognitive abilities. One concern that might arise is the inclusion of Amanita muscaria<sup>7</sup> into the  
665 Megascans Food dataset, which is poisonous (when consumed by ingestion without any specific  
666 preparation). This is why we urge the reader not to treat the included objects as edible items, even  
667 though they are a part of the “food” category. We provide random samples from both of them in Fig 9  
668 and Fig 10. Note that they are almost indistinguishable from real objects.

### 669 **D.2 Datasets statistics**

670 We provide the datasets statistics in Tab 4. For CARLA [56], we provide them for comparison and do  
671 not use this dataset as a benchmark since it is small, has simple geometry and texture.

## 672 **E Additional samples**

673 We provide random non-cherry-picked samples from our model in Fig 12, but we recommend visiting  
674 the website for video illustrations: <https://rethinking-3d-gans.github.io>.

## 675 **F Potential negative societal impacts**

676 Our developed method is in the general family of media synthesis algorithms, that could be used for  
677 automatized creation and manipulation of different types of media content, like images, videos or 3D

---

<sup>5</sup><https://quixel.com/megascans>

<sup>6</sup><https://creativecommons.org/licenses/by-nc-sa/4.0>

<sup>7</sup>[https://en.wikipedia.org/wiki/Amanita\\_muscaria](https://en.wikipedia.org/wiki/Amanita_muscaria)



Figure 9: Real images from the Megascans Plants dataset. This dataset contains very complex geometry and texture, while having good camera coverage.

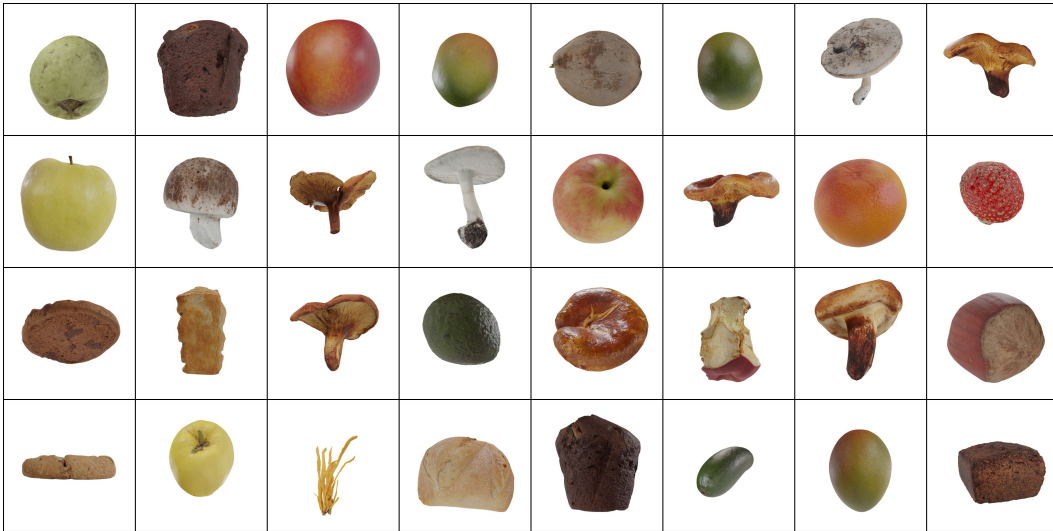


Figure 10: Real images from the Megascans Food dataset. Caution: some objects in this dataset could be poisonous.

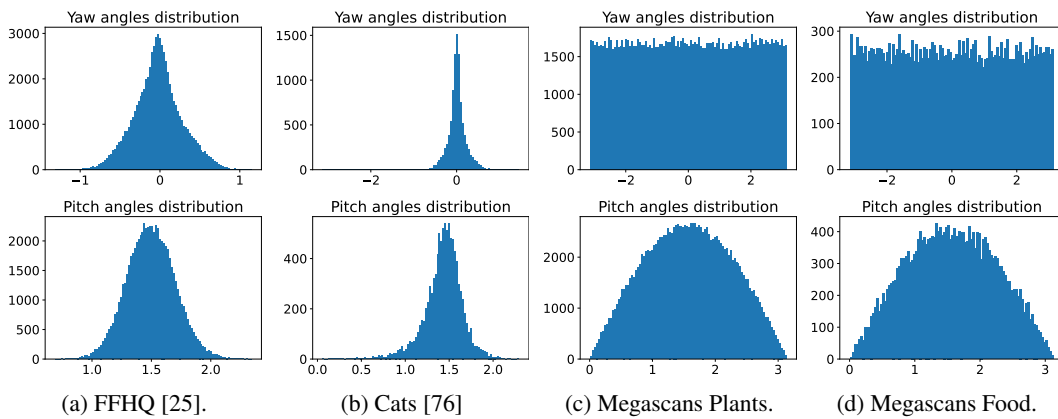
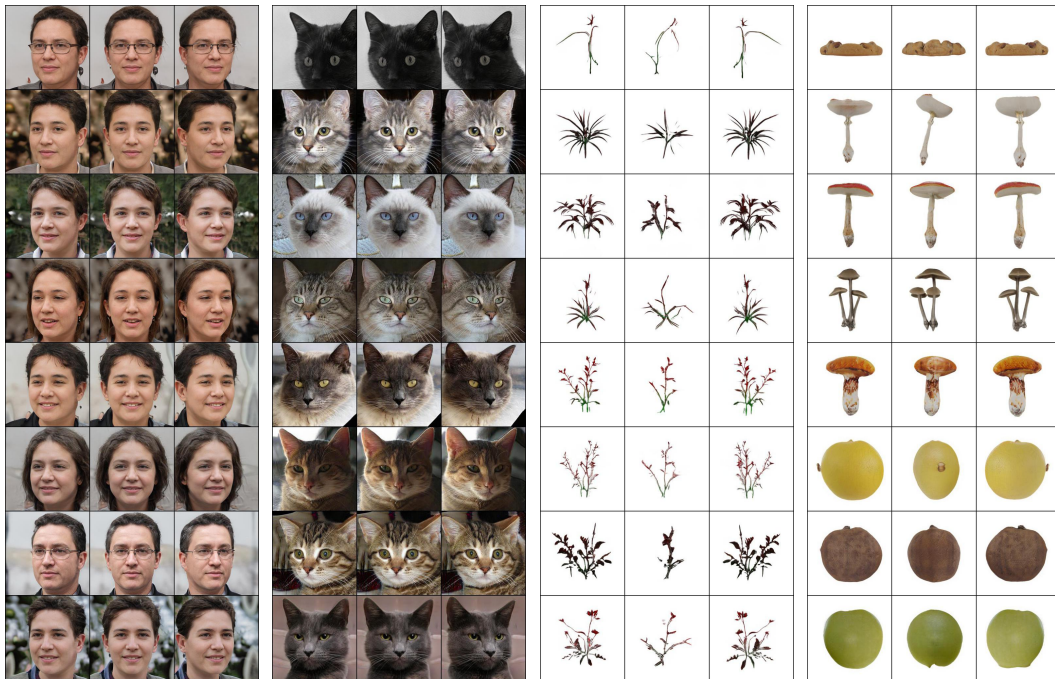


Figure 11: Comparing yaw/pitch angles distribution for different datasets.



Table 4: Comparing 3D datasets. Megascans Plants and Megascans Food are much more complex in terms of geometry and has much better camera coverage than FFHQ [25] or Cats [76]. The abbreviation “USphere( $\mu, \zeta$ )” denotes uniform distribution on a sphere (see  $\pi$ -GAN [7]) with mean  $\mu$  and pitch interval of  $[\mu - \zeta, \mu + \zeta]$ . For Cats, the final resolution depends on the cropping and we report the original dataset resolution.

Dataset	Number of images	Yaw distribution	Pitch distribution	Resolution
FFHQ [25]	70,000	Normal(0, 0.3)	Normal( $\pi/2, 0.2$ )	1024 <sup>2</sup>
Cats	10,000	Normal(0, 0.2)	Normal( $\pi/2, 0.2$ )	$\approx 604 \times 520$
CARLA	10,000	USphere(0, $\pi$ )	USphere( $\pi/4, \pi/4$ )	512 <sup>2</sup>
M-Plants	141,824	USphere(0, $\pi$ )	USphere( $\pi/2, \pi/2$ )	1024 <sup>2</sup>
M-Food	25,472	USphere(0, $\pi$ )	USphere( $\pi/2, \pi/2$ )	1024 <sup>2</sup>



(a) FFHQ 512<sup>2</sup> [25]. (b) Cats 256<sup>2</sup> [76] (c) Megascans Plants 256<sup>2</sup> (d) Megascans Food 256<sup>2</sup>

Figure 12: Random samples (without any cherry-picking) for our model. Zoom-in is recommended.

678 scenes. Of particular concern is creation of deepfakes<sup>8</sup> — photo-realistic replacing of one person’s  
 679 identity with another one in images and videos. While our model does not yet rich good enough  
 680 quality to have perceptually indistinguishable generations from real media, such concerns should be  
 681 kept in mind when developing this technology further.

## 682 G Ethical concerns

683 We have reviewed the ethics guidelines<sup>9</sup> and confirm that our work complies with them. As being  
 684 discussed in Appx D.1, our released datasets are not human-derived and hence do not contain any  
 685 personally identifiable information and are not biased against any groups of people.

<sup>8</sup><https://en.wikipedia.org/wiki/Deepfake>

<sup>9</sup><https://nips.cc/public/EthicsGuidelines>