

VQ3D: Learning a 3D-Aware Generative Model on ImageNet

Kyle Sargent
Stanford University

Jing Yu Koh
Carnegie Mellon University

Han Zhang
Google Research

Huiwen Chang
Google Research

Charles Herrmann
Google Research

Pratul Srinivasan
Google Research

Jiajun Wu
Stanford University

Deqing Sun
Google Research

Abstract

Recent work has shown the possibility of training generative models of 3D content from 2D image collections on small datasets corresponding to a single object class, such as human faces, animal faces, or cars. However, these models struggle on larger, more complex datasets. To model diverse and unconstrained image collections such as ImageNet, we present VQ3D, which introduces a NeRF-based decoder into a two-stage vector-quantized autoencoder. Our Stage 1 allows for the reconstruction of an input image and the ability to change the camera position around the image, and our Stage 2 allows for the generation of new 3D scenes. VQ3D is capable of generating and reconstructing 3D-aware images from the 1000-class ImageNet dataset of 1.2 million training images. We achieve an ImageNet generation FID score of 16.8, compared to 69.8 for the next best baseline method. For video results, please see the project [webpage](#).

1. Introduction

3D assets are an important part of popular media formats such as video games, movies, and computer graphics. Given that 3D content can be time-consuming to create by hand, leveraging machine learning techniques to automatically generate 3D content is an active area of research. While machine learning techniques benefit from training on large amounts of data, existing 3D datasets have noisy labels and are orders of magnitude smaller than those of 2D images. To get around the limitations of 3D datasets, recent work has shown the possibility of learning generative models of 3D scenes from images with limited or no 3D labels [5, 15, 22, 23].

These GAN-based approaches demonstrate the promise of learning 3D representations from 2D data. However, these methods require fine-tuning of prior pose distributions for individual models and datasets [4, 15, 22, 23], or the usage

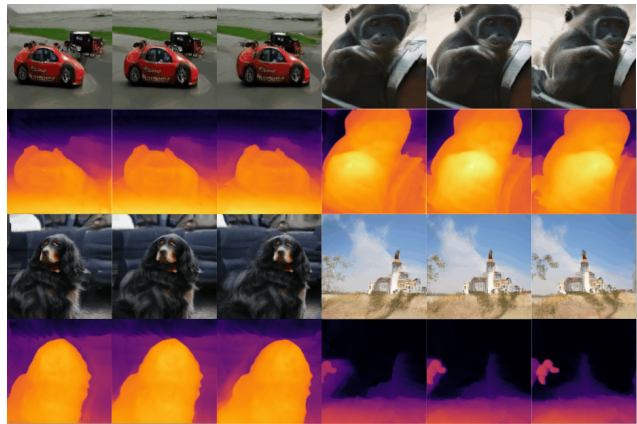


Figure 1. Fully generated 3D-aware images from our Stage 2 model on ImageNet. Please see supplemental materials for video results.

of ground truth pose data [5], and thereby typically operate on single-class datasets, e.g., human faces [16], animal faces [6], or cars [41]. In contrast, many 2D generative models, such as text-to-image generation models [24, 29, 44] and two-stage image models [12, 43] show impressive performance on very large and diverse image collections. The most recent state-of-the-art 2D models leverage diffusion or vector quantization rather than GANs to scale well to large datasets. This motivates us to pursue vector quantization as an alternative to GANs for learning 3D generative models.

In this paper, we propose VQ3D, a strong 3D-aware generative model that can be learnt from large and diverse 2D image collections, such as ImageNet [7]. To encourage stability and higher reconstruction quality, we forgo GAN-based [14] approaches [4, 5, 15, 22, 23], in favor of the 2-stage autoencoder formulation of VQGAN [12] and ViT-VQGAN [43]. But, different from these 2D autoencoder models, we learn 3D geometry by introducing a conditional NeRF decoder and modified triplane representation which can handle unbounded scenes, and training with a novel loss formulation

which encourages high-quality geometry and novel views.

Our formulation has three advantages, ensuring it to scale well to ImageNet. First, separating the training into two stages (reconstruction and generation) enables us to directly supervise the first stage training via a novel depth loss, using pseudo-GT depth. This is possible because in the first stage, as our conditional NeRF decoder learns to reconstruct the input, it also predicts the depth of each image.

Second, we do not require hand-tuning of pose sampling distributions or ground-truth pose data, which are required by previous GAN-based approaches [4, 5, 15, 22, 23]. Our training objective simply enforces reconstruction from a canonical camera pose, and plausible novel views within a neighborhood of the canonical pose. While this objective regrettably rules out very large camera motion, it also eliminates the need for excessive tuning of the pose distribution for each dataset, and allows our model to work out-of-the-box for multiple object categories. Thus, our model uses identical pose sampling hyperparameters for each dataset.

Finally, our two-stage formulation is simpler and more reliable than existing techniques for training 3D-aware generative models. Previous work [2, 30] has identified difficulties in scaling up GANs to large datasets (such as ImageNet). We verify that baseline 3D-aware GAN methods [4, 5, 15, 23], while working well on single-object datasets, fail to learn good generative models for ImageNet. Our formulation does not use progressive growing [4, 5], a neural upsampler [4, 5, 23], pose conditioning [5, 36], or patch-wise discriminators [31, 36], but still learns meaningful 3D representations. Compared to the best existing 3D-aware baseline, VQ3D attains a 75.9% relative improvement on FID scores for 3D-aware ImageNet images (69.8 for StyleNeRF [15] to 16.8 for VQ3D).

In summary, we make the following three contributions:

- We present a novel 3D-aware generative model that can be trained on large and diverse 2D image collections. Our model does not require tuning pose hyperparameters for each dataset or ground truth poses, and can leverage a pseudo-depth estimator during training.
- We obtain state-of-the-art generation results on ImageNet and competitive results on CompCars, demonstrating that our 3D-aware generative model is capable of fitting a dataset at the scale and diversity of ImageNet. Our model significantly outperforms the next best baseline.
- The Stage 1 of our model enables 3D-aware image editing and manipulation. One forward pass through our network converts a single RGB image into a manipulable NeRF, without relying on an expensive inversion optimization used in prior work [4, 5].

2. Related Work

3D-aware generative models. Several recent papers tackle the task of modeling 3D-aware generation, primarily through the GAN framework [14]. HoloGAN [22] learns perspective projection and rendering of 3D features, and applies 3D rigid-body transforms to generate new images from different poses. More recently, several papers use NeRF [21] as the 3D backbone [4, 15, 23, 40], which allows the 3D scene to be defined as a 3D volume parameterized by an MLP. EG3D [5] proposes a hybrid triplane representation which scales well with resolution, and enables greater generation detail. Disentangled3D [37] learns a 3D-aware GAN from monocular images with disentangled geometry, appearance, and pose. Pix2NeRF [3] proposes a method for unsupervised learning of neural representations with a shared pose prior, which enables rendering of novel views from a single input image. GRAF [31] and EpiGRAF [36] train 3D GANs via patch-wise representations to save on the expense of volume rendering. GRAM [8] proposes learning a set of implicit surfaces, shared for the training object category. At inference time, images are generated by accumulating the radiance along each ray using ray-surface intersections as samples.

Conditional NeRF and other 3D representations. Recent work has focused on the appropriate way to condition NeRF to achieve maximum expressiveness. GIRAFFE [23] demonstrated success with the “conditioning-by-concatenation” approach [35], in which the scene’s latent codes are fed into the first layer of the NeRF MLP and not thereafter. Other work such as pi-GAN [4] transforms the latent code into a vector of frequencies and phase shifts for each layer of a SIREN [34]. Other work has used hypernetworks [33, 35] to parameterize 3D representations, and MetaSDF [33] showed that many forms of conditioning are special cases of the hypernetwork approach. Our model can be seen as a conditional NeRF. We show that our novel decoder architecture, consisting of a ViT-L [10] and contracted triplane representation, is powerful enough to encode and reconstruct all of ImageNet. Given a single image, we show that in a single forward pass and without any optimization, our model can create a NeRF of an input RGB image with reasonable reconstruction at the main view and plausible novel views.

Quantization models. Image quantization is a powerful paradigm used in recent state-of-the-art generative models. In this setup, an image is encoded into a discrete latent representation [38], which improves generation quality when paired with an autoregressive generative prior (most often a transformer [39]). This has led to impressive results in image generation [13, 20, 43], text-to-image generation [9, 25, 44], and other tasks. Recent image quantization models improve reconstruction quality by introducing adversarial losses [13], using vision transformer encoders and

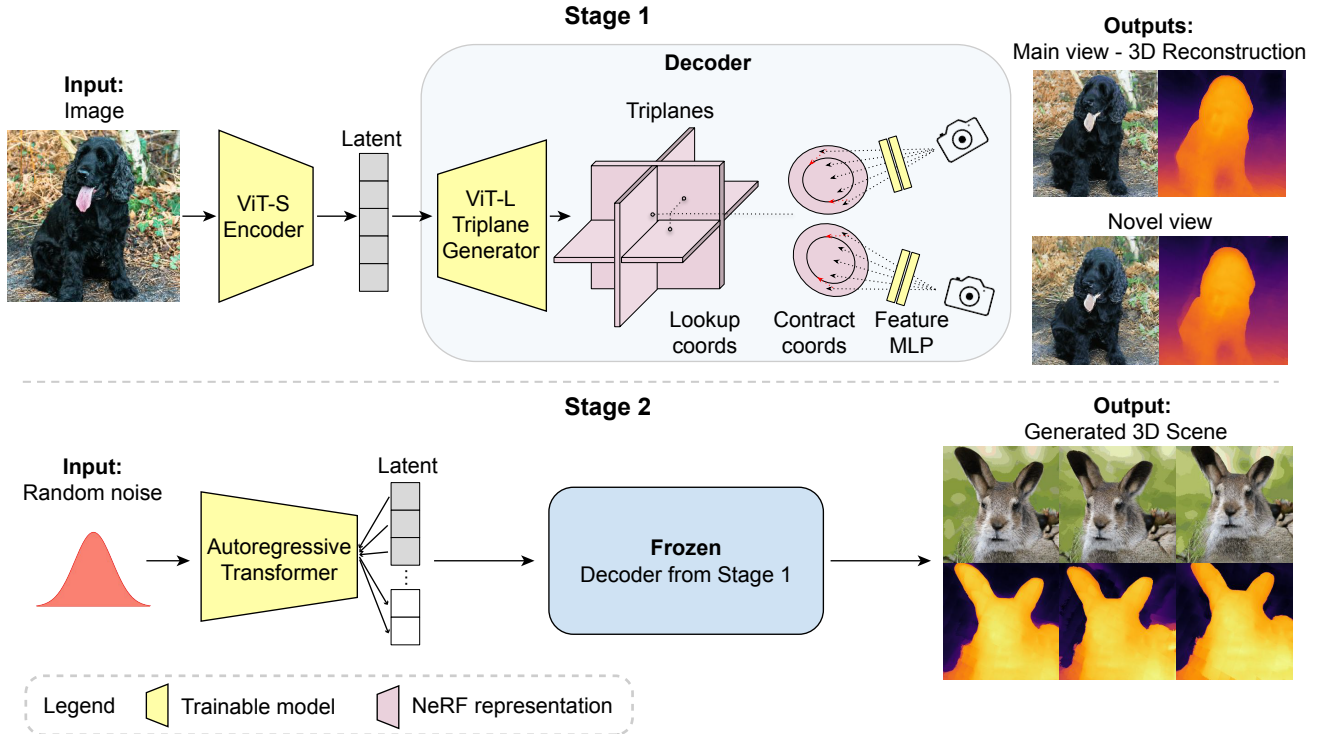


Figure 2. Diagram of our model architecture.

decoders (ViT) [10, 43] as both encoder and decoder, representing discrete codes as a stacked map [20], and more. Such quantization architectures typically use powerful CNNs [12] or ViT [43] encoders and decoders; ViT and CNN-based architectures show good performance reconstructing large image datasets; in this paper, we show that our NeRF-based decoder can also work well in the quantization framework. It has the capacity to encode and reconstruct a large and diverse dataset such as ImageNet, and also learns a discrete latent codebook that can be used to train a powerful fully generative Stage 2 model.

Single-view 3D reconstruction and novel view synthesis.

Various approaches for 3D reconstruction or novel view synthesis in the context of generative or auto-encoder models have been proposed. Kato et. al [17] propose an adversarial training scheme using two discriminators for single-view 3D reconstruction. Their scheme, in which the main discriminator critiques real and reconstructed views, while an auxiliary discriminator distinguishes between the reconstructed input view and predicted novel views, inspires our use of two discriminators for similar reasons. However, their model cannot sample totally new scenes. More recently, uORF [42] uses NeRFs as 3D object representations to enable 3D scene decomposition. uORF represents a 3D scene as a composition of an object radiance field for each object, and a background

radiance field for the remainder of the scene. This enables re-rendering and editing of 3D scenes from an input image. However, uORF also cannot sample new scenes, and moreover requires multi-view training datasets.

In the domain of novel scene generation, Generative Query Networks (GQN) [11] use CNNs to represent and generate scenes. GQNs can imagine and re-render scenes from novel viewpoints, but due to the usage of CNNs, do not explicitly embed 3D geometry or have any guarantees of scene consistency. NeRF-VAE [19] proposes an improved representation using a VAE which models multiple scenes. This enables efficient inference-time sampling of novel scenes, as well as re-rendering from multiple viewpoints. Unlike GQNs, which have no 3D prior, NeRF-VAE uses NeRF to achieve 3D consistency. However, it relies on multi-view training data. LOLNeRF [28] learns a generative model of 3D face images but requires a pretrained keypoint estimator and the auto-decoder formulation requires an optimization to be applied to examples outside its training set. By contrast, our method can be applied to single RGB images and requires only 2D training data and an off-the-shelf depth estimator for training.

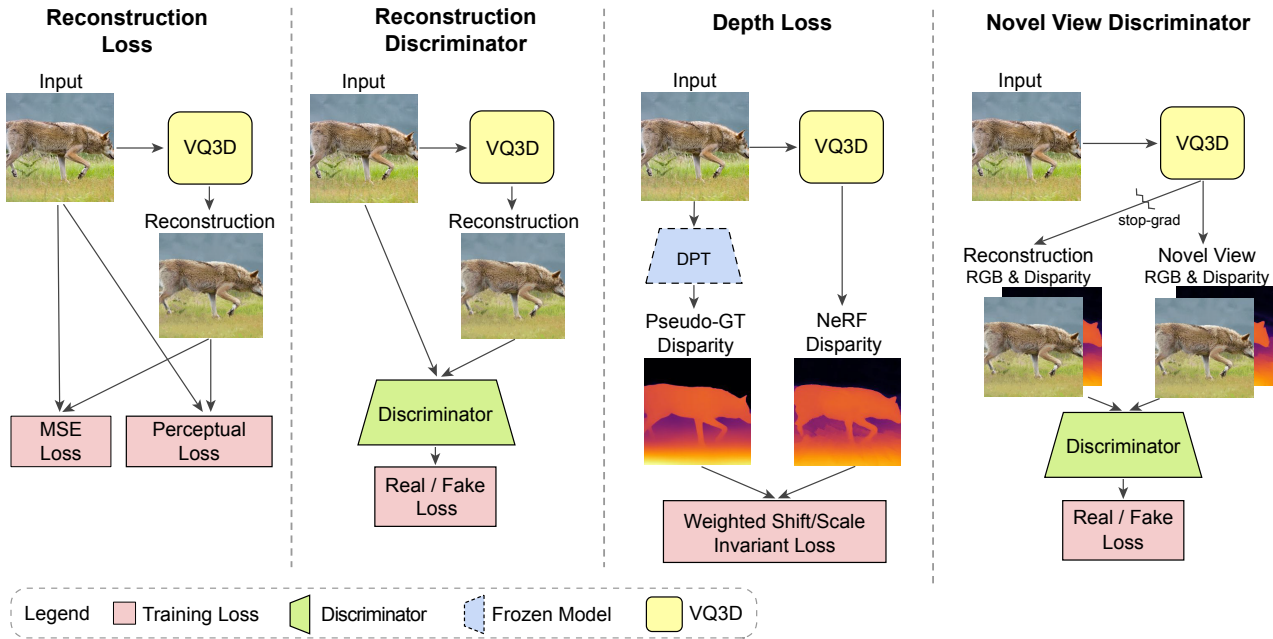


Figure 3. Diagram of the key losses in Stage 1 optimization.

3. Model

3.1. Overview of VQ3D

Our model is a vector-quantized autoencoder [12, 43], which is trained in two stages. Stage 1 of our model consists of an encoder and decoder. The encoder encodes RGB images into a learned latent codebook, and the decoder reconstructs them. A diagram of the inputs, outputs, and architecture of the first stage is given in the top of Figure 2. The encoder of our first stage is a ViT similar to VIM [43], but the decoder is a conditional NeRF. The first stage is trained end-to-end by encoding and reconstructing RGB training images while minimizing reconstruction and adversarial losses. Because the decoder is a NeRF, we are able to supervise the NeRF geometry with an additional training loss using pseudo-GT disparity. We also render novel views of decoded images and critique them with an additional adversarial loss. A diagram of the key losses used in Stage 1 training is shown in Figure 3. After training, the first stage can be used to encode unseen single RGB images and then reconstruct them in 3D, which enables novel view synthesis, image editing and manipulations.

Stage 2 is a generative autoregressive transformer which predicts sequences of latent tokens. A diagram of the inputs, outputs, and architecture is shown in the bottom of Figure 2. The architectural and training details are generally the same as [43]. We train it on the sequences of latent codes produced by our Stage 1 encoder. After training, the autoregressive transformer can be used to generate totally new 3D images by first sampling a sequence of latent tokens and then applying

our NeRF-based decoder. Importantly, our Stage 2 model inherits the properties optimized in Stage 1, so the fully generated images have high quality geometry and plausible novel views.

3.2. Training

We now provide additional training details for the two stages of our model.

Stage 1. The goal of the first stage is to learn a model which can compress image pixels into a sequence of discrete indices corresponding to a learnt latent codebook [12, 43]. Since we desire our model to be 3D-aware, we impose several additional criteria:

- 1. Good reconstruction from a canonical view.** On ImageNet, ground truth camera extrinsics are unknown and probably not even well-defined due to the presence of deformable and ambiguous object categories and scenes without salient objects. Therefore, we simply fix a single ‘canonical pose’ for reconstruction, and our criterion is that our conditional NeRF-based autoencoder should successfully reconstruct the dataset from this view.
- 2. Reasonable novel views.** We expect that images decoded at novel views within a specified range of the canonical view will have similar quality to images decoded at the canonical view.

3. **Correct geometry.** The geometry of the scene as represented by the NeRF should correspond to the unknown ground truth geometry of the RGB image up to scale and shift.

We enforce these criteria by introducing several auxiliary models and losses, summarized in Figure 3. To enforce (1) good reconstruction at the canonical view, we train with a combination of the MSE, perceptual, and logit-laplace loss following [43], the combination of which we term \mathcal{L}_{rec} .

To enforce (2) reasonable novel views, we leverage a main and auxiliary discriminator similar to [17]. The first discriminator distinguishes between real and reconstructed images at the canonical viewpoint, while the second distinguishes between reconstructed images at the canonical viewpoint and novel views. In this way, the model cannot allocate all its capacity to reconstructing images at the canonical viewpoint without also having high-quality novel views. As noted by [17], the generator may slightly corrupt the main view in order to collaborate with the novel view branch to fool the discriminator; thus, we add a stop-grad between the main view and the novel view discriminator. Unlike [4, 5, 15, 23], we find it unnecessary to tune a separate distribution of novel views for each dataset, and instead sample novel views uniformly in a disc tangent to a sphere at the canonical camera pose. We use the non-saturating GAN objective \mathcal{L}_{gan} [14] for both discriminators. We additionally concatenate the predicted depth as input to the auxiliary discriminator to ensure the distribution of depths does not change depending on the camera viewpoint.

To enforce (3) correct geometry, we supervise the NeRF depth with pseudo-GT geometry at the main viewpoint. We employ the pretrained depth prediction transformer model DPT [26] which produces pseudo-GT disparity estimates for the images in our training datasets. Thus, our model is limited to some extent by the quality of the depth estimator chosen. [27] proposed a shift- and scale- invariant l_2 loss for training monocular depth estimation in which the shift and scale are determined by solving a closed-form least squares alignment with the GT depth. We propose a novel formulation of this shift- and scale- invariant loss adapted to the NeRF setting, in which we supervise the weight of every sample along each ray rather than the accumulated depth. For a given image, let $i \in \{1 \dots N\}$ and $k \in \{1 \dots L\}$ be indices which range over the image plane and ray samples respectively, let D_{ik} be the pointwise disparities of the NeRF sample locations, let W_{ik} be corresponding NeRF weights from volumetric rendering [21], and let d_i be the pseudo-GT depth from DPT. Then we define s^*, t^* to be the closed-form solution of the weighted least squares problem

$$s^*, t^* = \arg \min_{s, t} \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^L W_{ik} (s D_{ik} + t - d_i)^2 \quad (1)$$

And set our depth loss to be the weighted scale- and shift-

invariant loss

$$\mathcal{L}_{\text{depth}} = \frac{1}{N} \sum_{i=1}^N \sum_{d=1}^L W_{ik} (s^* D_{ik} + t^* - d_i)^2 \quad (2)$$

Assuming the weight sum to 1 along each ray, this loss is minimized when the NeRF allocates 0 weight to all but one sample location along each ray, and the expectation with respect to the weights of the disparity is equal to the GT disparity map up to a scale and shift. In this way it functions similarly to the distortion loss proposed in [1] by penalizing weight distributions which are too spread out, but also encourages the weights to be concentrated near the correct geometry. Importantly, this formulation still allows for more than one surface along each ray and thus for occlusion and disocclusion, because the penalty is applied to the volumetric rendering weights and not the predicted density. We find this depth loss formulation to be critical for good performance. In particular, supervising the accumulated disparity rather than the pointwise disparities leads to poor performance, and we provide an ablation of this and other design choices in the supplementary material. We additionally introduce two penalties on the scale determined by this alignment:

$$\mathcal{L}_{\text{scale}} = \lambda_{s1} \max(0, -s_{\text{scale}}^*) + \lambda_{s2} \max(s_{\text{scale}}^* - 1, 0) \quad (3)$$

λ_{s1} is the weight of a small penalty to prevent the sign of the disparity scale from flipping negative, which we found necessary unlike in [27]. λ_{s2} weights a penalty preventing the disparity maps from becoming too flat, which encourages perceptually pleasing novel views. We additionally include the same vector-quantization loss \mathcal{L}_{vq} as [43], and the distortion and interlevel losses of MipNeRF360 [1], given by $\mathcal{L}_{\text{nerf}}$. The loss for our autoencoder is thus:

$$\mathcal{L} = \mathcal{L}_{\text{rec}} + \mathcal{L}_{\text{gan}} + \mathcal{L}_{\text{depth}} + \mathcal{L}_{\text{scale}} + \mathcal{L}_{\text{vq}} + \mathcal{L}_{\text{nerf}} \quad (4)$$

Stage 2. The goal of Stage 2 is to learn an autoregressive model over the discrete encodings produced by the Stage 1 encoder, so that completely new 3D scenes can be generated. Our Stage 2 transformer and training details follow [43]. We verify experimentally that our fully generative Stage 2 model inherits the properties optimized in Stage 1; namely, 3D-consistent novel views and high quality geometry. We also apply top- k and top- p filtering similar to [13].

3.3. Architecture

A full architecture diagram is shown in Figure 2. Similar to [43], we leverage the powerful vision transformer [10] architecture in both the encoder and decoder. Different from [43], which is trained on 2D images, we utilize a novel decoder with 3D inductive bias to facilitate the learning of 3D representations. We now give an overview of the individual components of our architecture.

Encoder and triplane generator. For the encoder, we use a ViT-S model. For the decoder, we use a ViT-L model to decode the latent codes into 3 triplanes of size 512x512 with feature dimension 32. We find that the triplane construction stage of the decoder benefits from the increased capacity of the ViT-L model.

Contracted triplane representation & NeRF MLP. We must reconstruct and generate potentially unbounded ImageNet scenes, but we are motivated to leverage the powerful triplane representation [5], Therefore, we propose an adapted triplane representation borrowing from both [5] and [1]. We apply the contraction function of MipNeRF360 to bound coordinates within the triplanes before looking up their values, and use the linear-in-disparity sampling scheme with separate proposal and NeRF MLP. The MLPs convert interpolated triplane features to density and, in the case of the NeRF MLP, RGB color. Similar to [5], our MLPs are lightweight, with 2 layers and 32 hidden units each; unlike [5], we directly render RGB color rather than using a neural upsampler, as we found neural upsampling to be a source of myriad and confusing artifacts not fixable via dual discriminators [5] or consistency losses [15].

Autoregressive transformer. We train transformer [39] to autoregressively predict the next image token. We follow the hyperparameters in the base model of VIM [43]. For ImageNet, we train a conditional model, and for other datasets we train unconditional generative models.

4. Experiments

4.1. Main results

We study the performance of our method and the baseline methods on ImageNet. The ImageNet dataset [7] is a well-known classification benchmark which consists of 1.28M images of 1000 object classes. It is a standard benchmark for 2D image generation, for both conditional and unconditional generation. We compare against pi-GAN [4], GIRAFFE [23], EG3D [5], and StyleNeRF [15]. We re-implemented pi-GAN and GIRAFFE using our internal framework, and ran the provided code for EG3D and StyleNeRF. Since ImageNet does not have GT poses and pseudo-GT poses are not possible to compute, we disable generator and discriminator pose conditioning for EG3D and sample from a pre-defined pose distribution. We note that EG3D exhibits significant inter-run variance in ImageNet FID even for the same config, and provide more details in the supplementary material.

Our main results for generation on ImageNet compared against the benchmarks are given in Table 1. Notably, our FID score on ImageNet is the best by a wide margin. We show generated examples from our method and the benchmarks in Figure 4 and note our method generates superior

Generation	FID ↓
pi-GAN [4]	97.8
GIRAFFE [23]	132.0
StyleNeRF [15]	69.8
EG3D [5]	82.2
VQ3D (Ours)	16.8

Table 1. FID scores of 3D generative models on ImageNet. We set a new state of the art on ImageNet with a more than fourfold improvement over the next best baseline.

samples.

In addition to generating high quality scenes, Stage 1 of our method can also be used for single-view 3D reconstruction and manipulation. Figure 5 shows single RGB images reconstructed by our Stage 1 with estimated geometry. Our network performs well at reconstruction and needs only a single forward pass to compute a NeRF for an input image, unlike prior work [4, 5] which requires an inversion optimization. Moreover, the reconstructed NeRFs can be manipulated, for instance to render novel views. We show examples of novel views in Figure 6.

For our main results on ImageNet, we training for the longest possible time and use the most optimal top-p and top-k sampling parameters. We conduct additional analysis experiments on the learning of geometry and model ablation, for which we use a consistent Stage 1 step, Stage 2 step across each study and do not using top-p or top-k sampling unless ablating it directly as in Table 4.

First, we study the learning of good geometry, both for our model and the baseline methods. One potential concern may be that the use of pseudo-GT depth limits the comparability of our technique with the baseline GAN methods. We address this concern by analyzing both the FID score and the depth accuracy metric used in [5, 32]. This metric is defined as the mean- and variance-normalized MSE between the NeRF depth and the predicted depth of the generated image. Table 2 gives the result for generative models with and without depth losses. Note that EG3D’s FID without depth loss is different from Table 1 due to the significant inter-run variance for EG3D’s performance. For the GAN methods, we find that our pointwise disparity loss works poorly but the original scale- and shift- invariant MSE loss from [27] improves geometry. For our method, we show the Stage 2 performance with and without our novel pointwise weighted depth loss. While performance on the depth accuracy metric can improve when various depth losses are incorporated training, the effect on FID is negligible. In this way we see that incorporating pseudo-GT depth is unlikely to meaningfully improve the FID for the baseline methods without substantial changes. We were unable to design a depth loss which prevented flat depths for StyleNeRF.

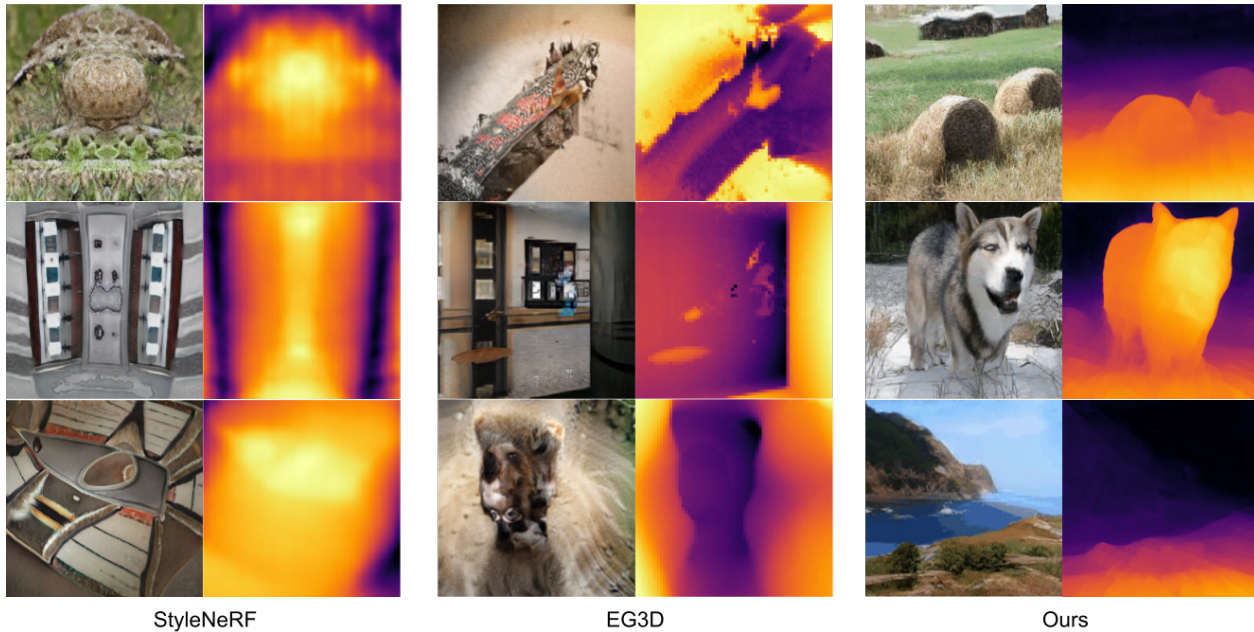


Figure 4. Generated samples and disparity from models trained on ImageNet. Ours model generates high-quality images and geometry.

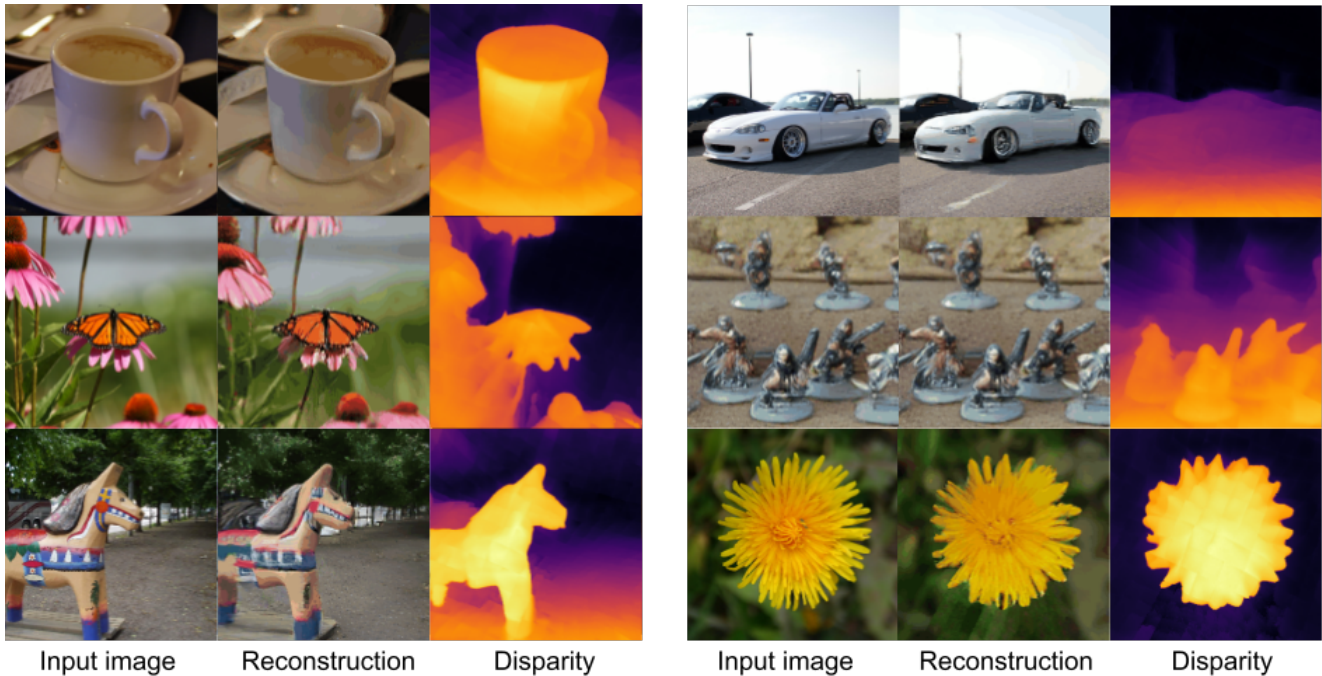


Figure 5. Reconstructions and estimated disparity on single images by our conditional NeRF-based autoencoder. Though our model is trained on ImageNet and achieves comparable performance on unseen ImageNet images, we show OpenImages results for licensing reasons.

Better geometry does not imply better FID. Additionally, learning geometry without a depth loss may be unreliable. For example, StyleNeRF [15] found learning of geometry was unreliable without training tricks such as progressive growing. During our ImageNet experiments, we also observed that StyleNeRF is sensitive to hyperparameters, and

can learn to produce flat depths. EG3D [5] showed that removing GT poses as input to the discriminator is enough to cause the geometry to degenerate to a flat plane.

We conduct ablations on our Stage 1 in Table 3 starting from our baseline architecture (row 1). Using a CNN encoder and decoder rather than ViT (row 2) is unstable and leads to

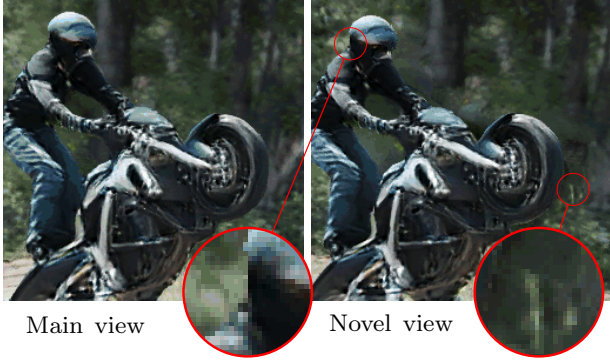


Figure 6. Example camera manipulations of a reconstructed scene. Our approach naturally handles sharp occlusions (left spyglass) and inpainting of disoccluded pixels (right spyglass) without supervision of novel views.

Generation	FID ↓	Depth accuracy ↓
pi-GAN	101.4	1.41
GIRAFFE	132.1	1.78
EG3D	109.3	1.66
VQ3D (Ours)	36.1	1.90
pi-GAN + depth loss	97.8	0.88
GIRAFFE + depth loss	132.0	1.16
EG3D + depth loss	91.8	0.88
VQ3D (Ours) + depth loss	31.7	0.16

Table 2. Evaluation of depth losses on ImageNet. While adding depth losses can improve the quality of geometry, it will not lead to FID improvements significant enough to close the gap between our method and the baselines.

Reconstruction	ImageNet FID ↓	Depth Acc. ↓	Disparity scale ↓
(1) VQ3D (Ours)	11.2	0.18	1.00
(2) CNN enc., dec.	(diverges)	-	-
(3) W/o \mathcal{L}_{gan}	10.6	0.22	1.27
(4) W/o \mathcal{L}_{scale}	9.4	0.23	1.21
(5) W/o \mathcal{L}_{nerf}	9.2	0.28	4.88
(6) W/o \mathcal{L}_{depth}	4.0	1.91	0.61
(7) W/o Triplanes	273.5	1.00	2.15

Table 3. VQ3D ablation study. Removing components compromises the model capacity, 3D awareness, or novel view quality.

divergence. Eliminating the GAN loss (row 3) or depth scale loss (row 4) leads to a higher learned disparity scale causing perceptually flat novel views. Removing the GAN loss (row 3) also leads to artifacts in inpainting disoccluded pixels. Eliminating the NeRF loss (row 5) leads to worse depth accuracy and a very high disparity scale. Eliminating the depth loss (row 6) improves reconstruction FID, but causes the depths to collapse to a flat plane and leads to worse depth accuracy. A fully implicit representation instead of triplanes (row 7) gives very poor FID since we are forced to use a very small MLP due to the expense of volume rendering at 256×256 .

We analyze the performance of VQ3D with top- p and top-

top- k	top- p	FID ↓
1000	1.00	31.5
2000	1.00	33.2
3000	1.00	34.1
4000	1.00	34.7
8192	1.00	36.1
8192	0.98	32.2
8192	0.95	35.7

Table 4. FID scores on ImageNet from sampling over top- k and top- p values. 8192 is the size of our full codebook. In general our model benefits from some restriction of the sampling process, though too restrictive top- p hurts performance.

Generation	CompCars
pi-GAN [4]	16.9
GIRAFFE [23]	26 [†]
StyleNeRF [15]	8 [†]
GIRAFFE HD [40]	7.2 [†]
EG3D [5]	32.2
VQ3D (Ours)	7.3

Table 5. FID scores of 3D generative models. † indicates numbers taken from the respective papers, we trained other models ourselves. Although baseline models use separate, tuned pose hyperparameters for each dataset, our identical, simple pose sampling scheme works well on both ImageNet and CompCars.

k sampling in Table 4, as [12] noted these sampling changes can give significant performance improvements analogous to truncation sampling for GANs [2]. For VQ3D, a top- k of 1000 and top- p of 1.0 gives the best FID results.

4.2. Other 3D benchmark datasets

Two other prominent 3D-aware benchmark datasets are FFHQ [16] and CompCars [41]. Due to the ethical and legal issues associated with manipulation and generative modeling of faces, we do not study FFHQ. On CompCars, our model is competitive with the state of the art (Table 5).

5. Discussion and conclusion

Limitations and ethical considerations. Our work has several limitations. First, while some benchmark methods, e.g. [15, 23] have shown the ability to model 360-degree rotation of the generated scene when trained on specific single-class datasets like CompCars [41], our auto-encoder based formulation makes large viewpoint manipulation difficult. While some baseline approaches such as [5] and [4] have demonstrated the ability to learn high-quality geometry in an unsupervised manner, our approach requires a pretrained depth network for the depth loss. However note that our model works on general object classes, while those focus only on small and single-class image collections and

require tuning hyperparameters of the pose distribution.

We are committed to understanding and promoting positive societal impacts. Although we do not train a generative model on FFHQ and thus avoid many serious ethical considerations, ImageNet does contain some images of humans and human faces, and our model will likely inherit biases which are present in the dataset.

Conclusion. We have presented VQ3D, a framework for 3D-aware representation learning and generation. VQ3D sets a state-of-the-art by a wide margin on the large and diverse ImageNet dataset, relative to existing strong geometry-aware generative model baselines. We conduct extensive analysis and ablation and also show our model performs competitively on the more standard single class benchmark CompCars. Our work shows that it could be a fruitful path to use large and diverse 2D image datasets to train 3D-aware generative models, thereby facilitating 3D content creation.

References

- [1] J. T. Barron, B. Mildenhall, D. Verbin, P. P. Srinivasan, and P. Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. *CVPR*, 2022. 5, 6, 12
- [2] A. Brock, J. Donahue, and K. Simonyan. Large scale gan training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018. 2, 8
- [3] S. Cai, A. Obukhov, D. Dai, and L. Van Gool. Pix2nerf: Unsupervised conditional p-gan for single image to neural radiance fields translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3981–3990, 2022. 2
- [4] E. Chan, M. Monteiro, P. Kellnhofer, J. Wu, and G. Wetzstein. pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *arXiv*, 2020. 1, 2, 5, 6, 8, 12
- [5] E. R. Chan, C. Z. Lin, M. A. Chan, K. Nagano, B. Pan, S. D. Mello, O. Gallo, L. Guibas, J. Tremblay, S. Khamis, T. Karras, and G. Wetzstein. Efficient geometry-aware 3D generative adversarial networks. In *arXiv*, 2021. 1, 2, 5, 6, 7, 8, 12, 13
- [6] Y. Choi, Y. Uh, J. Yoo, and J.-W. Ha. Stargan v2: Diverse image synthesis for multiple domains. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 1
- [7] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 1, 6
- [8] Y. Deng, J. Yang, J. Xiang, and X. Tong. Gram: Generative radiance manifolds for 3d-aware image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10673–10683, 2022. 2
- [9] M. Ding, Z. Yang, W. Hong, W. Zheng, C. Zhou, D. Yin, J. Lin, X. Zou, Z. Shao, H. Yang, et al. Cogview: Mastering text-to-image generation via transformers. *Advances in Neural Information Processing Systems*, 34:19822–19835, 2021. 2
- [10] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 2, 3, 5
- [11] S. A. Eslami, D. Jimenez Rezende, F. Besse, F. Viola, A. S. Morcos, M. Garnelo, A. Ruderman, A. A. Rusu, I. Danihelka, K. Gregor, et al. Neural scene representation and rendering. *Science*, 360(6394):1204–1210, 2018. 3
- [12] P. Esser, R. Rombach, and B. Ommer. Taming transformers for high-resolution image synthesis, 2020. 1, 3, 4, 8
- [13] P. Esser, R. Rombach, and B. Ommer. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12873–12883, 2021. 2, 5
- [14] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014. 1, 2, 5
- [15] J. Gu, L. Liu, P. Wang, and C. Theobalt. Stylenerf: A style-based 3d-aware generator for high-resolution image synthesis. *arXiv preprint arXiv:2110.08985*, 2021. 1, 2, 5, 6, 7, 8, 13
- [16] T. Karras, S. Laine, and T. Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019. 1, 8, 12
- [17] H. Kato and T. Harada. Learning view priors for single-view 3d reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9778–9787, 2019. 3, 5
- [18] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 12

- [19] A. R. Kosiorok, H. Strathmann, D. Zoran, P. Moreno, R. Schneider, S. Mokra, and D. J. Rezende. Nerf-vae: A geometry aware 3d scene generative model. In M. Meila and T. Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 5742–5752. PMLR, 18–24 Jul 2021. 3
- [20] D. Lee, C. Kim, S. Kim, M. Cho, and W.-S. Han. Autoregressive image generation using residual quantization. *arXiv preprint arXiv:2203.01941*, 2022. 2, 3
- [21] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European conference on computer vision*, pages 405–421. Springer, 2020. 2, 5
- [22] T. Nguyen-Phuoc, C. Li, L. Theis, C. Richardt, and Y.-L. Yang. Hologan: Unsupervised learning of 3d representations from natural images. *arXiv*, 2019. 1, 2
- [23] M. Niemeyer and A. Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021. 1, 2, 5, 6, 8, 12
- [24] A. Ramesh, P. Dhariwal, A. Nichol, C. Chu, and M. Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022. 1
- [25] A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021. 2
- [26] R. Ranftl, A. Bochkovskiy, and V. Koltun. Vision transformers for dense prediction. *ArXiv preprint*, 2021. 5, 13
- [27] R. Ranftl, K. Lasinger, D. Hafner, K. Schindler, and V. Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(3), 2022. 5, 6
- [28] D. Rebain, M. Matthews, K. M. Yi, D. Lagun, and A. Tagliasacchi. Lolnerf: Learn from one look, 2022. 3
- [29] C. Saharia, W. Chan, S. Saxena, L. Li, J. Whang, E. Denton, S. K. S. Ghasemipour, B. K. Ayan, S. S. Mahdavi, R. G. Lopes, T. Salimans, J. Ho, D. J. Fleet, and M. Norouzi. Photorealistic text-to-image diffusion models with deep language understanding, 2022. 1
- [30] A. Sauer, K. Schwarz, and A. Geiger. Stylegan-xl: Scaling stylegan to large diverse datasets. *CoRR*, abs/2202.00273, 2022. 2
- [31] K. Schwarz, Y. Liao, M. Niemeyer, and A. Geiger. Graf: Generative radiance fields for 3d-aware image synthesis. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 2
- [32] Y. Shi, D. Aggarwal, and A. K. Jain. Lifting 2d stylegan for 3d-aware face generation. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. arXiv, 2020. 6, 13
- [33] V. Sitzmann, E. Chan, R. Tucker, N. Snavely, and G. Wetzstein. Metasdf: Meta-learning signed distance functions. *Advances in Neural Information Processing Systems*, 33:10136–10147, 2020. 2
- [34] V. Sitzmann, J. Martel, A. Bergman, D. Lindell, and G. Wetzstein. Implicit neural representations with periodic activation functions. *Advances in Neural Information Processing Systems*, 33:7462–7473, 2020. 2
- [35] V. Sitzmann, M. Zollhöfer, and G. Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. *Advances in Neural Information Processing Systems*, 32, 2019. 2
- [36] I. Skorokhodov, S. Tulyakov, Y. Wang, and P. Wonka. Epigraf: Rethinking training of 3d gans. *arXiv preprint arXiv:2206.10535*, 2022. 2
- [37] A. Tewari, X. Pan, O. Fried, M. Agrawala, C. Theobalt, et al. Disentangled3d: Learning a 3d generative model with disentangled geometry and appearance from monocular images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1516–1525, 2022. 2
- [38] A. Van Den Oord, O. Vinyals, et al. Neural discrete representation learning. *Advances in neural information processing systems*, 30, 2017. 2
- [39] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 2, 6
- [40] Y. Xue, Y. Li, K. K. Singh, and Y. J. Lee. Giraffe hd: A high-resolution 3d-aware generative model. In *CVPR*, 2022. 2, 8
- [41] L. Yang, P. Luo, C. C. Loy, and X. Tang. A large-scale car dataset for fine-grained categorization and verification, 2015. 1, 8

- [42] H.-X. Yu, L. J. Guibas, and J. Wu. Unsupervised discovery of object radiance fields. *arXiv preprint arXiv:2107.07905*, 2021. [3](#)
- [43] J. Yu, X. Li, J. Y. Koh, H. Zhang, R. Pang, J. Qin, A. Ku, Y. Xu, J. Baldrige, and Y. Wu. Vector-quantized image modeling with improved vqgan. *arXiv preprint arXiv:2110.04627*, 2021. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [12](#)
- [44] J. Yu, Y. Xu, J. Y. Koh, T. Luong, G. Baid, Z. Wang, V. Vasudevan, A. Ku, Y. Yang, B. K. Ayan, B. Hutchinson, W. Han, Z. Parekh, X. Li, H. Zhang, J. Baldrige, and Y. Wu. Scaling autoregressive models for content-rich text-to-image generation, 2022. [1](#), [2](#)

Thanks for checking the supplementary materials. in which we provide additional details for the ease of replicating the results of our method. For video results, **we encourage the reader to consult the project webpage.**

6. Implementation details

6.1. NeRF Model

We train and evaluate all models at 256x256 resolution, except pi-GAN [4] which we train and evaluate at 128x128 following [5].

We use a constant 49.13 degree field of view and pinhole camera model. We use a camera radius of 2.732 following [23] and a canonical pose at $(-2.732, 0, 0)$. All views canonical and novel are looking at $(0, 0, 0)$ and have a constant camera up vector of $(0, 0, 1)$. We sample novel view camera locations uniformly in a disc in the YZ-plane centered at the canonical pose with radius .4. We use a near plane of .7 and far plane of $1e6$. We find that using the slightly large near plane of .7 was necessary in order to avoid a failure mode where all the content was clustered very close to the camera leading to poor novel views; we hope to eliminate this failure mode in future work.

We perform volume rendering at the full 256x256 resolution using the importance sampling scheme of [1]. We have a separate proposal and NeRF MLP and render in two stages, the first stage using the proposal MLP to evaluate a wide range of sample locations, and the second stage using the NeRF MLP queried at locations determined by importance sampling of the weights and locations from the first stage. During training, we add a stop-grad between the proposal and NeRF MLP like [1] and supervise the Proposal MLP with the interlevel loss. Our NeRF MLP is not view dependent and the only input it receives is triplane features which are determined by looking up the contracted 3D points of the sample locations. We apply a fixed orthonormal transformation to all points before triplane lookup because our canonical pose is axis-aligned, so we desire that our triplanes are not axis-aligned to avoid artifacts.

We evaluate 32 samples along each ray for each sampling stage. Thus, rendering a full 256x256 RGB image takes 256x256x64 triplane lookups and MLP evaluations. We use the same number of ray samples, 32, for training, FID evaluation, and rendering videos.

6.2. Setup and hyperparameters

We train with the Adam optimizer [18] with $\beta_1 = .9, \beta_2 = .99$, and cosine learning rate schedule with 50K warmup steps, similar to [43], with an initial autoencoder LR of 0 and max LR of $1e-4$. We use codebook size 8192 and l_2 -normalized, factorized codebook with embedding dimension 8.

Different from [43], we do not use weight decay, and

Loss	Weight
l_2 [43]	1
Perceptual [43]	1e-1
Logit-laplace [43]	1e-1
Discriminator [43]	1e-1
Novel discriminator	1e-1
Quantizer [43]	1
Weighted pointwise depth (λ_{depth})	1e1
Negative depth scale penalty (λ_{s1})	1
Large depth scale penalty (λ_{s2})	1e-3
Interlevel [1]	1
Distortion [1]	2.5e-1

Table 6. Weights of various losses used in Stage 1 training of our autoencoder.

our discriminator LR is scaled down from the autoencoder LR by .5 so that the discriminator does not overpower the autoencoder, which was an issue especially in early training.

Due to the many losses in our Stage 1 training, we outline their weights in Table 6 and reference the original implementation if they are not losses designed by us.

6.3. Discriminators

We use StyleGAN [16] discriminators for both the main and novel view discriminator. They are identical except that the novel view discriminator accepts 4-channel RGBD images, and the main view discriminator accepts 3-channel RGB images.

6.4. Timing and throughput

We train our main model on ImageNet for 180K steps in Stage 1, and 140K steps in Stage 2. On a single V100, our Stage 1 model renders 8.7 img/s. We train with a Stage 1 batch size of 128 and Stage 2 batch size of 512. For each batch in Stage 1, we render 256 images; 128 to reconstruct the full batch at the canonical view, and an additional 128 novel views to be critiqued by the novel view discriminator. Though this is expensive, our volume rendering stage is made cheaper even than [5] by using 32 instead of 64 hidden units for the feature MLPs and using 32 instead of 48 samples per ray. We leverage gradient accumulation in Stage 2 training in order to train with 512 batch size.

6.5. Evaluation

As is standard [43], we compute Stage 1 metrics (reconstruction) over the ImageNet validation set and Stage 2 metrics (generation) over real samples from the train set and generated samples. We use 50K samples to evaluate FID for all methods. We sample views for Stage 2 FID computation uniformly in a disc of radius .2 tangent to the sphere at the canonical pose.

EG3D Tuning	Sweep	ImageNet FID
R1 gamma	{.3, .6}	{ 82 , 99}
Density reg.	{.125, .25, .5}	{91, 82 , 96}
Disc. LR (1e-3)	{.5, 1, 2, 4}	{122, 82 , 116, 113}
Gen. LR (1e-3)	{.625, 1.25, 2.5, 5}	{111, 82 , 106, 136}

Table 7. Hyperparameter tuning of EG3D on ImageNet.

StyleNeRF Tuning	Sweep	ImageNet FID
R1 gamma	{.15, .3, .6}	{75, 73 , 74}
Disc. LR (1e-3)	{.625, 1.25, 2.5, 5}	{96, 87, 73, 69 }
Gen. LR (1e-3)	{.625, 1.25, 2.5, 5}	{78, 74, 73 , 107}

Table 8. Hyperparameter tuning of StyleNeRF on ImageNet.

We use the Depth Accuracy metric used in [5, 32], but differently we don’t mask out any invalid regions because our monocular depth estimator DPT [26] predicts a dense depth map over the input and every pixel is assumed to be valid. We also use disparity instead of depth because we model much larger scenes than either [32] or [5].

We experiment with classifier guidance but find it gives only a small performance boost, and so investigating model improvements was more worthwhile to improve the FID than tweaking classifier guidance settings.

7. Additional experiments

Although the strongest baselines, EG3D [5] and StyleNeRF [15], perform poorly on ImageNet, they may need to be tuned to perform well on this new dataset. To verify that the limitation of the baseline methods is fundamental, we extensively tune both on Imagenet for a range of hyperparameters in Tables 7 and 8. We see the baselines do not achieve good performance for a range of hyperparameter settings. Additionally, we observe that EG3D has significant inter-run variance in terms of FID on ImageNet, even when rerunning the same configuration, which may indicate instability for large datasets such as ImageNet. When running the same config multiple times, we report the best value achieved among all runs.

8. Additional samples

We show additional uncurated generated samples with geometry in Figure 7 and Figure 8.



Figure 7. Uncurated fully generated samples from our Stage 2 model.



Figure 8. More uncurated fully generated samples from our Stage 2 model.