
DT-NVS: Diffusion Transformers for Novel View Synthesis

Wonbong Jang^{*} Jonathan Tremblay[†] Lourdes Agapito^{*}

^{*}UCL

[†]NVIDIA

{ucabwja, l.agapito}@ucl.ac.uk jtremblay@nvidia.com

Abstract

Generating novel views of a natural scene, e.g., every-day scenes both indoors and outdoors, from a single view is an under-explored problem, even though it is an organic extension to the object-centric novel view synthesis. Existing diffusion-based approaches focus rather on small camera movements in real scenes or only consider unnatural object-centric scenes, limiting their potential applications in real-world settings. In this paper we move away from these constrained regimes and propose a 3D diffusion model trained with image-only losses on a large-scale dataset of real-world, multi-category, unaligned, and casually acquired videos of everyday scenes. We propose DT-NVS, a 3D-aware diffusion model for generalized novel view synthesis that exploits a transformer-based architecture backbone. We make significant contributions to transformer and self-attention architectures to translate images to 3d representations, and novel camera conditioning strategies to allow training on real-world unaligned datasets. In addition, we introduce a novel training paradigm swapping the role of reference frame between the conditioning image and the sampled noisy input. We evaluate our approach on the 3D task of generalized novel view synthesis from a single input image and show improvements over state-of-the-art 3D aware diffusion models and deterministic approaches, while generating diverse outputs.

1 Introduction

Diffusion models have emerged as a powerful methodology for high-quality 2D image and video generation from multimodal inputs. However, training diffusion models to learn 3D representations for truly 3D-aware generation has not been straightforward. The unique challenge comes from their reliance on large amounts of ground-truth training data which, in the case of 3D scenes, is scarce and costly to acquire. Meanwhile, recent advances in 3D geometry and appearance acquisition from images have resulted in powerful methods such as neural radiance fields (NeRFs) [39], InstantNGP [42] or Gaussian Splatting [31], which can learn 3D implicit scene representations for high-quality new view synthesis from 2D images only, without the need for 3D ground truth. However, these methods are scene-specific, they require a large number of carefully acquired input views with corresponding camera poses, and need costly test-time optimization.

In this paper we focus on the challenging problem of generating novel views of general scenes, from a single input image and using only 2D losses. Numerous approaches have aimed to generalize NeRF to model multiple scenes; some explored using global latents [27], back-projecting features from pre-trained networks [77], or applying generative models such as GANs [7]. More recently, we have witnessed attempts to train 3D diffusion models for this task with 2D-only supervision using implicit representations [2, 62, 30]. However, they require canonicalized datasets with 3D-aligned scenes, assume single-category, object-centric scenes and simplified camera models, or they rely on synthetic data. The recent emergence of large-scale multi-view datasets of casually-captured videos

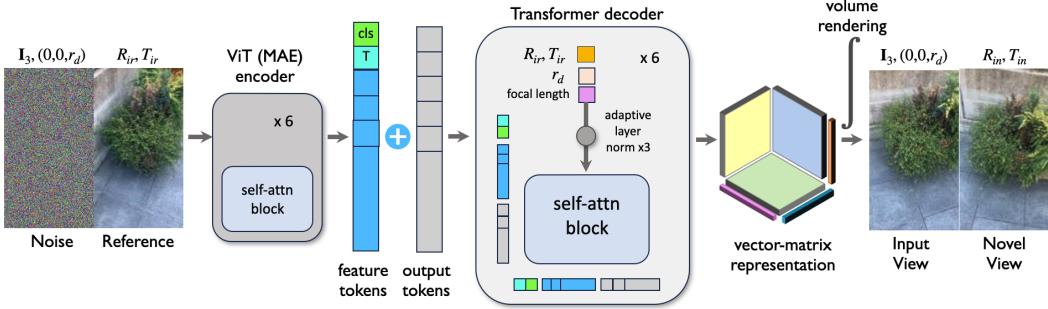


Figure 1: **Architecture:** DT-NVS is a 3D-aware diffusion model that takes noise z_t^i at I_3 and a reference image x^r at reference viewpoints R_{ir}, T_{ir} as input and learns to denoise from c^i and do novel view synthesis on R_{in}, T_{in} . The encoder processes noise z and the reference image x^r separately, and generates feature tokens for each. Our decoder concatenates both feature tokens with its output tokens, then applies self-attention with conditioning tokens on their respective camera parameters (input view for output tokens and feature tokens from z_i^t , and reference viewpoints on features tokens from x^r) We reshape the output tokens into vector-matrix representation, and perform volume rendering. During the training, we supervise the model with denoising loss at input view I_3 and photometric loss on novel view. For inference, we use the predicted output at input view x^i and denoise according to diffusion step t .

of hundreds of object categories such as Co3D [47] or MVImgNet [78] has opened the door for transformer-only architectures also to be proposed for novel-view synthesis [25, 28]. While these methods show promising results, they are deterministic at heart, while the task of novel view synthesis is clearly a generative one.

We propose DT-NVS, a novel view synthesis diffusion model that exploits a transformer-based backbone architecture to predict a radiance field from a single reference image. Unlike many diffusion-based solutions to novel view synthesis, our approach is not limited to masked, object-centric scenes, or 3D aligned scenes, and can be applied to any real-world captures. As such, we propose new self-attention architectures along with novel camera viewpoint conditioning strategies and we introduce a novel training paradigm that switches between sampled noise and reference images to avoid trivial or degraded solutions. We leverage MVImgNet [78], a very large-scale dataset of multi-view videos of indoor and outdoor scenes of hundreds of categories of everyday objects. Our model quantitatively achieves better FID scores than deterministic transformer-based architectures such as NViST [28] and also outperforms non-transformer based diffusion models [2, 62, 1], on both real-world (42% increase) and synthetic (212% increase) datasets. We also provide a detailed ablation study to justify our design choices.

2 Related Work

Diffusion Models: A diffusion model is a generative approach similar to GANs, first proposed by Sohl-Dickstein *et al.* [60], and gained popularity with DDPM [22]. Various techniques have been proposed to enhance the quality of diffusion model outcomes, including Cosine Schedule [12], v-parameterization [50], and classifier-free guidance [23]. DDIM [61] offers a more flexible and efficient sampling process, enabling faster generation of high-quality images with fewer steps. Hang *et al.* [17] proposed a minimum signal-to-noise ratio (Min-SNR) strategy, showing that the model converges faster by using SNR as a weight to the loss function. DiT [44] scales the Diffusion model with Transformers by conditioning diffusion steps with Adaptive Layer Normalization(AdaLN)-Zero [73].

Neural Radiance Fields: Neural 3D implicit representations were initially proposed in SRN [59] and later used in DVR [43], to train 3D-aware representations without 3D ground-truth. NeRF [39] revolutionized novel view synthesis from collections of posed images. To accelerate NeRF training, grid-based representations have been proposed. [76, 42, 14, 75] EG3D applied triplanes by projecting features into three planes in the context of 3D-aware GANs and TensoRF [9] proposed the vector-matrix representation. Efforts to generalize NeRF include conditioning on global latent vectors [15],

27, 41, 46, 3, 24], associating 2D feature views with target views [77, 67, 10, 47, 19, 65, 11, 26], or supervising NeRF with GAN losses [6, 7, 55, 34, 5].

3D-Aware Diffusion Models: Extending diffusion models to be 3D-aware has proved challenging due to the lack of large 3D ground truth datasets. Previous approaches fine-tuned a pre-trained latent diffusion model [48] on camera parameters [36, 38, 58, 63, 68, 52, 70]. Other approaches [45, 63, 72, 66] regularize NeRF using the pre-trained diffusion models or utilize the pre-trained diffusion models to generate images. Similar to latent diffusion, several approaches first train the 3D-aware GAN and apply the diffusion process on latent vectors [56, 32]. Instead of 2-stage training, there are approaches other methods obtain intermediate features and use them to denoise from the target view [8, 16, 64]. Another line of research involves learning diffusion models in 3D by rendering in 2D or backprojecting 2D features into 3D [30, 1, 2, 62], which usually assume the aligned scenes.

Transformers for 3D Tasks: Vision Transformer (ViT) [13] has been successful in the field of computer vision, and Masked Autoencoder(MAE) [18] learns visual features in self-supervised learning. For 3D, geometry-free methods based on transformer have been explored in [33, 49, 40]. To build 3D implicit representation from Transformer, GINA3D [57] and VQ3D [51] employ the adversarial loss. Other approaches deal with point clouds [29, 71]. LRM [25] and NVIST [28] apply the Transformer for 3D implicit representation in a deterministic way, and DMV [74] extends LRM into a diffusion, but both LRM and DMV do not model the background. DT-NVS is the first method to apply the Transformer with a 3D-aware diffusion model, directly learning from real-world scenes.

3 Background

For each scene \mathcal{S} , we have RGB renderings x^i from camera poses c^i and focal length f (which we assume constant). A diffusion model generates an image using a series of iterative steps, from noisy input to generated content. During this process, the output of the model is fed back as input recursively until convergence. To train such model, a forward process involves gradually adding noise to an image over several steps until it becomes nearly pure noise. Then, during the reverse process, the method aims to reverse the noise addition, transforming the noisy image back into the original image.

Formally, during the forward process, noise is added $\epsilon \sim \mathcal{N}(0, I)$ to x^i to form $z_t^i = \alpha_t x^i + \sigma_t \epsilon$ at each diffusion step $t \in [1, T]$. We assume a variance-preserving scenario, where $\alpha_t^2 + \sigma_t^2 = 1$, and follow the cosine schedule proposed in [12], with $\alpha_t = \cos(\pi t / 2T)$ and $\sigma_t = \sin(\pi t / 2T)$. With the assumptions of a Gaussian distribution and Markov chain, we can represent the marginal as $q(z_t^i | x^i) = \mathcal{N}(z_t^i | \alpha_t x^i, \sigma_t^2 I)$. Given $t > s$, $q(z_t^i | z_s^i) = \mathcal{N}(z_t^i | \alpha_{t|s} z_s^i, \sigma_{t|s}^2 I)$, where $\alpha_{t|s} = \alpha_t / \alpha_s$ and $\sigma_{t|s}^2 = \sigma_t^2 - \alpha_{t|s}^2 \sigma_s^2$.

In the reverse process, $q(z_s^i | z_t^i)$ is intractable; however, there exists a closed form solution when we condition on x^i as $q(z_s^i | z_t^i, x^i) = \mathcal{N}(z_s^i | \mu_{t \rightarrow s}, \sigma_{t \rightarrow s}^2 I)$, where $\mu_{t \rightarrow s} = \alpha_{t|s} \sigma_s^2 / \sigma_t^2 z_t^i + \alpha_s \sigma_{t|s}^2 / \sigma_t^2 x^i$ and $\sigma_{t \rightarrow s}^2 = \sigma_{t|s}^2 \sigma_s^2 / \sigma_t^2$. Applying the result from [22] $s \rightarrow t$, we can represent $p(z_s^i | z_t^i) = q(z_s^i | z_t^i, x^i = \hat{x}^i)$, using \hat{x}^i . Therefore, by estimating \hat{x}^i , we can approximate the posterior $q(z_s^i | z_t^i, x = \hat{x}^i)$.

Since our diffusion model needs to perform novel-view synthesis during training, our diffusion loss cannot be based on the noise (ϵ -parameterization) as DDPM [22], but rather on the ground truth (x -parameterization). To improve convergence we take inspiration from [50, 21] who proposed the v -parameterization, a re-weighting of the ϵ -parameterization: $v_t^i = \alpha_t \hat{\epsilon} - \sigma_t \hat{x}^i$. More concretely, we follow Hang *et al.* [17] who showed that all three parameterizations x , v and ϵ are equivalent by reweighting the x parameterization with signal-to-noise ratio values (SNR) α_t^2 / σ_t^2 .

4 Methodology

4.1 Diffusion in 3D implicit representation

In 3D implicit representations, \mathcal{S} is not directly observable since we only have access to a reference view. Therefore, we add noise ϵ to scene images x^i and train the model \mathcal{F}_θ to denoise them from the camera viewpoint c^i . To be 3D-aware, the model should not only estimate the image x^i but also

predict all of the possible images, x^n , in the scenes from their viewpoints, c^n . Here, we assume that we do not observe the noise level from c^n , so the model estimates \hat{x}^n , which means that the model predicts the ground-truth. For each diffusion step t , we multiply the weight $w(t)=\min(\text{SNR}(t)+1, 5)$, so that it is similar to v -parameterization. We can write the denoising loss $\mathcal{L}_{denoising}$ as below.

$$\mathcal{L}_{denoising} = \mathbb{E}_{\theta, t, x^r} [w(t)(\hat{x}^i - x^i)^2] \quad (1)$$

We get \hat{x}^n by rendering the estimated scene $\hat{\mathcal{S}}$ generated from the model \mathcal{F}_θ which takes the reference image x^r , noisy image z_t^i , and rays from the camera viewpoints c^i and c^r .

4.2 Predicting the scene: Transformer, Relative Pose, VM-Representation

Applying the transformer architecture to diffusion models presents unique challenges compared to U-Net, particularly in translating 2D image features into 3D output tokens within each grid. We propose four ways to help generate 3d outputs: conditioning on relative camera poses, a self-attention block for all tokens, randomly swapping the camera reference frame to be aligned with the input view or the reference view, applying dropout to conditioning images and leveraging a vector-matrix representation.

Relative Pose: To facilitate the job of the decoder, we assume that input views c^i always have the identity rotation. Since MVImgNet acquires the camera pose from SfM [54, 53] where different scenes have different scales and alignment, we apply an affine transformation to move the input camera to be at $(0, 0, -r_d)$, with identity rotation, where r_d is the distance between the centre of the coordinate system and the input camera. The reference image x^r transforms accordingly so that it has relative pose R_{ir}, T_{ir} between c^i and c^r .

Conditioning on Camera Parameters: The decoder employs self-attention only, by concatenating feature tokens from the encoder with output tokens which are replicated for grid position, then differentiated by learnable positional embedding. The decoder conditions on camera parameters using adaptive layer normalization (AdaLN) [73]. Unlike DiT [44] which conditions diffusion steps via AdaLN, we condition camera parameters only. Within each attention block, we split the tokens as shown in Figure 2, and apply AdaLN separately following their relative camera pose. This allows the decoder to condition on camera parameters better as shown in Table 1.

Swapping input views and reference views: The model tended to stick to a degraded solution when the noisy images were provided as input throughout the training, as shown in Table 1. To address this, We randomly swap the positions of noisy input images with reference images during the training, effectively placing the reference image with identity rotation. Additionally, we apply dropout to reference images, which regularizes the model and enables it to perform unconditional generation and classifier-free guidance.

VM Representation: Representing the scene as a voxel-grid is computationally expensive. Unlike 2D matrices, there is no Eckart-Young Theorem in tensor decomposition. We adopt the Vector-Matrix Representation (VM Representation) proposed by TensoRF [9]. The idea is to decompose a 3D tensor \mathcal{S} into the summation of three matrices $M^{Y,Z}, M^{Z,X}, M^{X,Y}$ and vectors V^X, V^Y, V^Z with k number of channels.

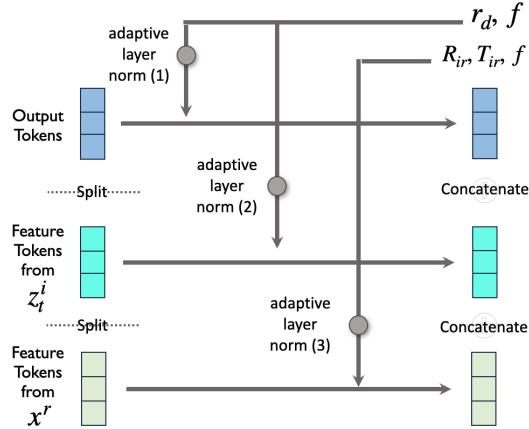


Figure 2: Camera Conditioning: We apply AdaLN separately for output tokens, feature tokens from noisy input image and those from reference image. As input rotation matrices are always identity, so we condition on camera distance r_d and focal length f on output tokens and feature tokens from z_t^i with different embedding MLPs. We condition feature tokens from x^r on relative pose R_{ir}, T_{ir} between input camera pose c^i and reference camera pose c^r .

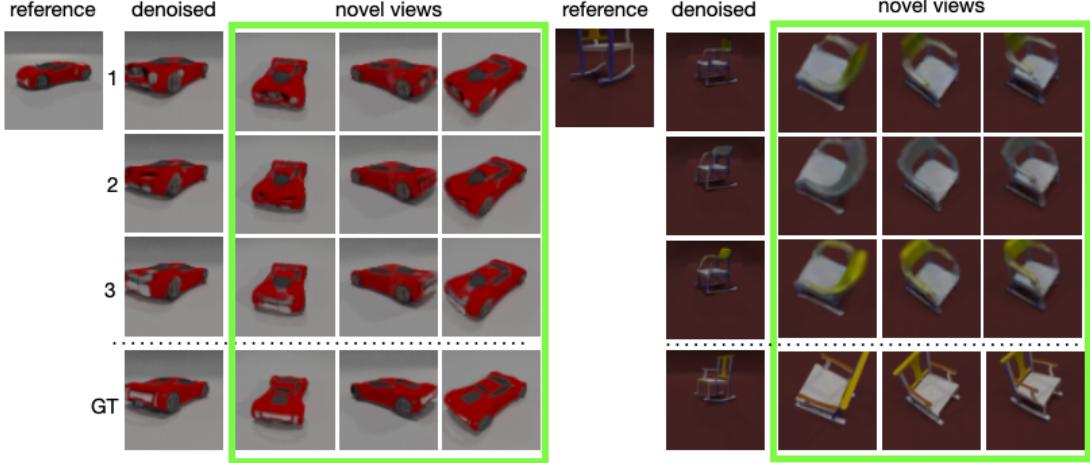


Figure 3: **Sampling results on ShapeNet:** Given reference images, we show different sampling results at input viewpoint (denoised), e.g., identity pose, and novel view synthesis results for objects of the car and chair category from ShapeNet.

$$\hat{\mathcal{S}} = \sum_{r_1=1}^k V_{r_1}^X \circ M_{r_1}^{Y,Z} + \sum_{r_2=1}^k V_{r_2}^Y \circ M_{r_2}^{Z,X} + \sum_{r_3=1}^k V_{r_3}^Z \circ M_{r_3}^{X,Y}. \quad (2)$$

The difference between DT-NVS and TensoRF is that each grid in DT-NVS communicates via attention mechanism within the decoder, leading to a more compact and efficient representation. After the decoder, DT-NVS embeds and reshapes the output tensors to build the VM representation. Finally, points are queried along the ray computed from the camera poses c^i and c^n , both are affine-transformed as c^i is the identity rotation.

4.3 Rendering

Once we build the estimated scene $\hat{\mathcal{S}}$ as a vector-matrix representation, we can render from any viewpoint. From camera poses c^i and c^n , we sample the points q_j along the ray r and query them to $\hat{\mathcal{S}}$ which is composed of vectors V_X, V_Y, V_Z and matrices $M_{Y,Z}, M_{Z,X}, M_{X,Y}$. Query points q_j are projected onto three vectors and matrices, and features are computed via bilinear sampling. We compute the density σ_u at q_u using Equation 2, and for color \mathbf{c}_u , we use shallow MLP to regress the dimensions into 3. Then, we can compute the color of the ray $\hat{C}(r)$ by volume rendering, by first computing the transmittance $T_u = \exp(-\sum_{l=1}^{u-1} \sigma_l \delta_l)$ where δ_l is the distance between adjacent query points.

$$\hat{C}(r) = \sum_{u=1}^N T_u (1 - \exp(-\sigma_u \delta_u)) \mathbf{c}_u \quad (3)$$

4.4 Losses: Photometric, Denoising, and Sampling

We supervise the model with the denoising loss on x^i as shown in Equation 1. In addition, to make the model truly 3d-aware, we also generate novel views by rendering the vector-matrix representation from new viewpoints, and compute a photometric loss.

Using Equation 3, we can render novel views \hat{x}^n from novel viewpoints c^n , as well as the input view \hat{x}^i as shown below.

$$\hat{x}^n, \hat{x}^i = \mathcal{F}_\theta(z_t^i, x^r, R_{ir}, T_{ir}; c^i, c^r, c^n) \quad (4)$$

Table 1: **Quantitative Results on MVImgNet** : DT-NVS achieves similar performance with NViST, and performs better on FID. This table also present our ablation study on MVImgNet Landscape dataset. Note that we compute FID for ours and ϵ -parameterization, as others do not perform well on other metrics. We also present GIBR [1] result, though their resolution is 256×256 .

	MVImgNet (Landscape / Portrait)			
	PSNR↑	SSIM↑	LPIPS↓	FID↓
Ours	20.82 / 20.65	0.62 / 0.63	0.22 / 0.21	37.8 / 14.32
Ours (w/cfg 2.0)	20.91 / 20.85	0.63 / 0.63	0.21 / 0.20	31.82 / 11.68
NViST [28]	21.03 / 21.23	0.62 / 0.62	0.20 / 0.21	39.35 / 16.68
w/ Cross-Attn		did not converge		
w/ noisy input only	18.62	0.41	0.30	—
w/o Encoder	17.32	0.42	0.35	—
w/o Decoder Camera Conditioning	17.19	0.38	0.33	—
w/ ϵ -parameterization (w/cfg 2.0)	20.72	0.61	0.22	44.2
GIBR [1] (higher resolution)	17.96	0.554	0.519	107.3

We employ a photometric loss by computing the L_2 loss that measures the discrepancy between the ground truth x^n and rendered views \hat{x}^n .

$$\mathcal{L}_{photo} = \mathbb{E}_{\theta, t, x^r}[(\hat{x}^n - x^n)^2] \quad (5)$$

We additionally use the LPIPS [79] loss on novel views \hat{x}^n as well as the input view \hat{x}^i . To regularize the model further, we also employ the distortion loss \mathcal{L}_{dist} proposed in MipNeRF360 [4]. The total loss function is as below.

$$\mathcal{L}_{total} = \mathcal{L}_{denoising} + \gamma_1 \mathcal{L}_{photo} + \gamma_2 \mathcal{L}_{LPIPS} + \gamma_3 \mathcal{L}_{dist} \quad (6)$$

Sampling: In 3D implicit representation, we only observe the noise level from the specific viewpoint. As we take the relative-pose-based approach, we denoise from the input viewpoint, which has identity rotation matrix. During the inference time, we use DDIM sampling [61], see Figure 3 for examples of sampling results.

5 Experimental Evaluation

5.1 MVImgNet

Dataset: MVImgNet consists of 6.5M images from 220K scenes across 238 categories, all of them are real world captures, and the camera poses are computed through COLMAP [53, 54], and each scene usually contains around 30 images. The dataset contains both portrait and aspect ratio, and we train the model separately for each aspect ratio. We split train/test for holding out every 50 scenes in alphabetical order. We downsample and center-crop images to 56×32 and 32×56 , we also downscale the point clouds from COLMAP to unit-cube, and change focal length accordingly. See Figure 4 for qualitative results and Figure 5 for depth predictions.

Baselines: We train NViST [28], which is a deterministic approach for training multiple unaligned real-world captures using the MVImgNet dataset. For the diffusion-based model, many approaches assume that scenes are aligned [62, 30, 2], but they would not be appropriate for this dataset which is not aligned. We quantitatively compare with GIBR [1] on MVImgNet as in Table 1, but note that they use a different train/test split, and a different resolutions.

Results: Table 1 presents our results. DT-NVS achieves better performance than NViST on FID, and performs similar to NViST on LPIPS [79] and SSIM [69]. Note that NViST minimizes the L2 distance, and it may contribute to the better results on these metrics. Qualitatively, DT-NVS generates sharper images, and NViST sometimes struggles to estimate the scale of scenes with compared to DT-NVS, see Figure 6. Moreover our method outperforms GIBR on FID [20], although, it is worth noting that their method uses a different train/test split, and uses the different resolutions. GIBR

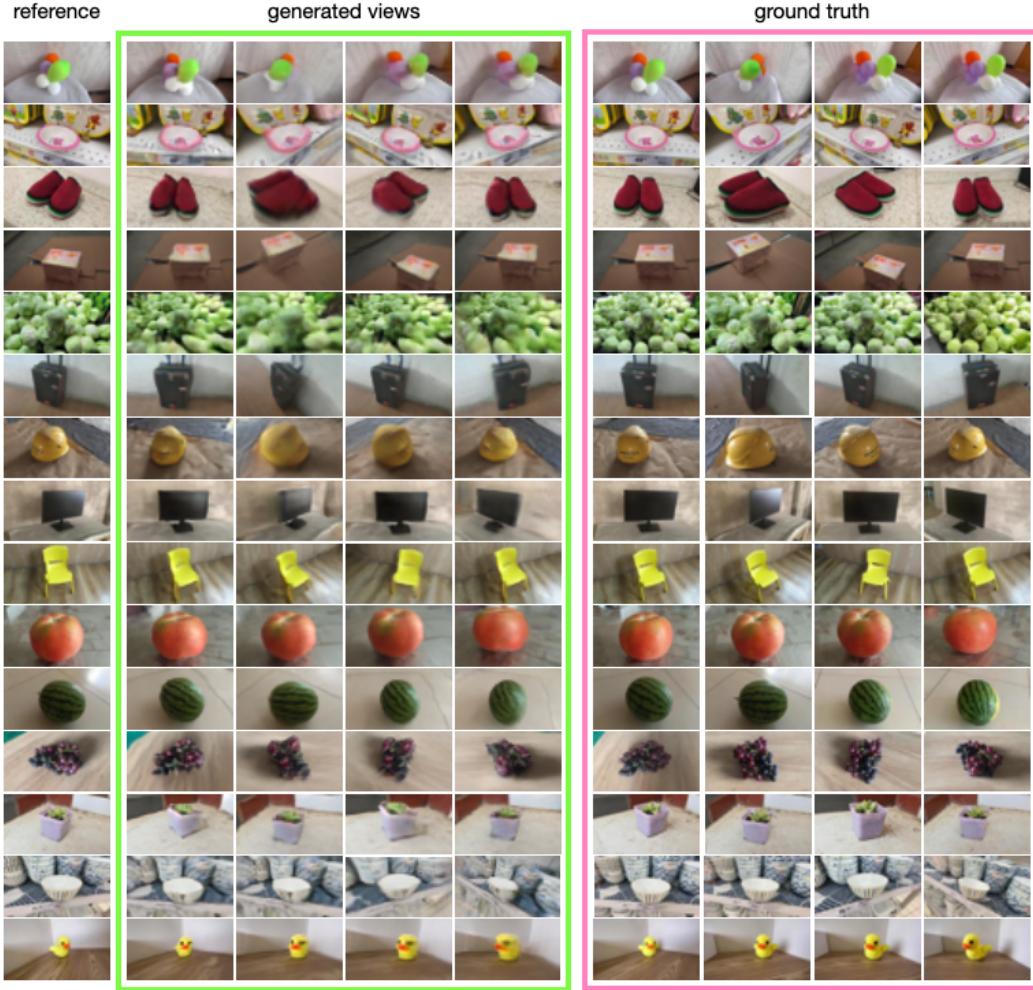


Figure 4: Qualitative Results on MVImgNet: We show the capabilities of DT-NVS on test images from unknown scenes of MVImgNet [78]. The model takes a reference image as input and denoises at input viewpoint during sampling. We show novel view synthesis results for three different viewpoints after the denoising step.

cropped and downsampled images to 90×90 or 256×256 , while DT-NVS pre-processes to 32×56 or 56×32 .

5.2 Ablation Study

In order to develop DT-NVS we made different design choices, such as not using cross-attention. The lower part of Table 1 presents different design variations and we discuss these variations as follow:

Cross-Attention based model: NViST [28] and LRM [25], which are both deterministic approaches, employ cross-attention and self-attention architectures in their decoders. The primary difference between self-attention and cross-attention is that feature tokens (output from the encoder) from noisy or reference images are not updated in the decoder, leading to smaller attention matrices. However, our training did not converge when we applied cross and self-attention approach to the decoder.

Fixing noisy images at input view: During the inference, we denoise from the input view (*i.e.*, identity rotation). If we always feed the noisy images with identity rotation during training, the model tends to find a trivial solution instead of learning the 3D structure from multiple views, which leads to lower performance. Randomly swapping the reference frame between the reference image and the sampled noisy input image regularizes the model for better performance, as shown in Table 1.

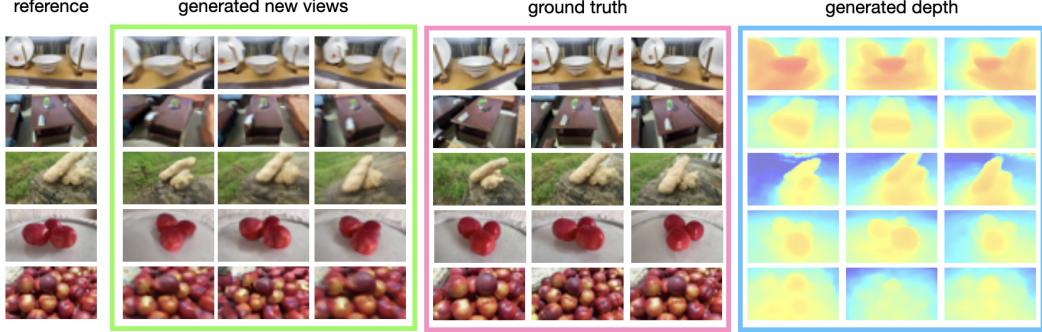


Figure 5: **Qualitative Results with depth prediction:** We show reference images, generated new views and the generated depth maps for a variety of scenes from MVImgNet, as well as a comparison with the ground truth views.

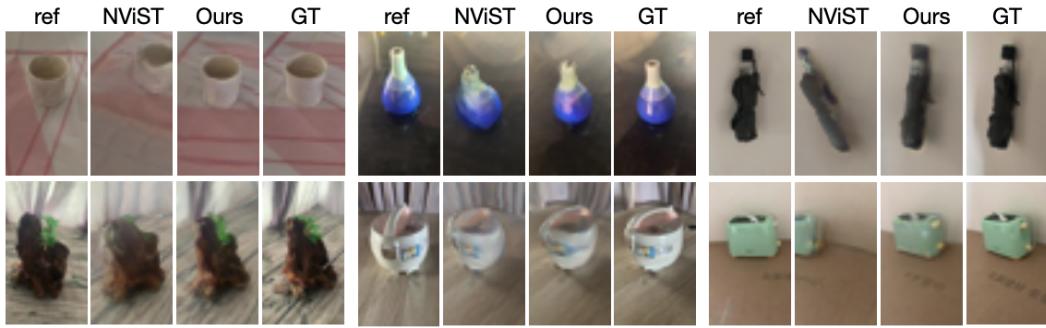


Figure 6: **Qualitative Comparison:** We compare DT-NVS with NViST [28]. For NViST, we provide the same reference images as input images. NViST sometimes fails to understand the scale and pose, and also generates more blurry images.

ϵ -parameterization: DDPM [22] proposed using ϵ -parameterization, and Hang *et al.* [17] suggested applying the minimum Signal Noise Ratio (SNR) to x -parameterization, resulting in a model that converges similarly to ϵ -parameterization. Here, we use the weight $\min(\text{SNR}(t) + 1, 5)$ for x -parameterization, which is akin to v -parameterization. In our ablation study, we compare models using the weights $\min(\text{SNR}(t), 5)$ and $\min(\text{SNR}(t) + 1, 5)$, corresponding to ϵ -parameterization and v -parameterization, respectively. We found that with the weight $\min(\text{SNR}(t), 5)$, the model performance degrades in terms of FID score.

Without Encoder: Since we employ self-attention in our Transformers, we conducted an ablation study to determine the role of the encoder. Without the encoder, the model’s performance significantly drops.

No Camera Conditioning on Decoder: To test the effectiveness of our camera conditioning approach depicted in Figure 2, we replaced our decoder with ViT. In this setup, the decoder needs to infer the camera pose from feature image tokens alone. As shown in Table 1, the model fails to train properly without conditioning on camera parameters.

5.3 ShapeNet

Dataset: We use ShapeNet renderings from [2] to validate our model. The dataset includes three categories (cars, chairs and planes) with each category containing 3,200 scenes, divided into training (2700) and testing (500) sets. The dataset assumes an object-centric scene, meaning the objects are always located at the center of the coordinate system. It also adopts the simplified camera model which always points toward the center of the coordinate system. Additionally, all 3D objects are aligned, sharing the same reference frame, *e.g.*, plane nose aligned with the positive x-axis. This

Table 2: **Quantitative Results on ShapeNet:** DT-NVS performs similar or better than baseline models on all ShapeNet categories. Our model performs much better than other diffusion-based models (RD, GIBR, VSD) in terms of FID score, and achieves similar FID score compared to EG3D which is 3D-aware GAN [7]. We denote split from RD and GIBR with ¹ and ².

	Cars	Planes	Chairs
	PSNR↑/SSIM↑/LPIPS↓/FID↓	PSNR↑/SSIM↑/LPIPS↓/FID↓	PSNR↑/SSIM↑/LPIPS↓/FID↓
Ours	27.64 / 0.87 / 0.12 / 14.9	26.9 / 0.88 / 0.15 / 21.8	28.75 / 0.88 / 0.12 / 15.6
RD ¹ [2]	25.4 / 0.81 / - / 46.5	26.3 / 0.83 / - / 53.3	26.6 / 0.83 / - / 47.8
GIBR ² [1]	29.74 / 0.90 / 0.14 / 90.1	(N/A)	(N/A)
VSD ² [62]	28.00 / 0.87 / 0.17 / 56.0	(N/A)	(N/A)
EG3D ¹ [7]	21.8 / 0.71 / - / 17.9	25.0 / 0.80 / - / 20.9	25.5 / 0.80 / - / 14.2

setting differs significantly from MVImgNet [78], where all scenes are not aligned and assume a 6-degree-of-freedom camera model.

Baselines and Results: Even with aligned scenes, we adopt the relative-pose-based approach, meaning we do not exploit the 3D alignment of objects in this dataset. We follow the training protocol from RenderDiffusion and compare our results with RenderDiffusion, GIBR, Viewset Diffusion [62], and the 3D-aware GAN EG3D [7]. As shown in Table 2, DT-NVS performs similarly or better across most metrics. Notably, DT-NVS excels in terms of the FID score compared to other diffusion models, achieving a level better or comparable to EG3D.

Diverse Sampling: We validate our model’s capability of generating multiple outputs for occluded regions from a single image, as shown in Figure 3. We deliberately choose viewpoints that are not visible in the reference image and denoise from there. The model demonstrates that it can generate multiple plausible outputs, maintaining consistency across viewpoints.

5.4 Implementation Details

We train our model using 2 A100-40GB GPUs for both MVImgNet and ShapeNet datasets, using the same architecture with both. Training MVImgNet takes 5 days (700,000 iterations) for both landscape and portrait, while training ShapeNet takes 2 days (400,000 iterations). We use a batch size of 44 for MVImgNet and 26 for ShapeNet. We use AdamW [37] optimizer and the learning rate is set to 2e-4 with cosine decay, and a warm-up period of 50,000 iterations is applied for both datasets.

6 Conclusion

We have introduced DT-NVS, a transformer-based 3D-aware diffusion model trained on real-world multiview images. Our evaluation demonstrates that DT-NVS outperforms baseline models, particularly in terms of FID score. We show qualitative results on a challenging real-world dataset of casually captured videos of everyday scenes (MVImgNet) that demonstrate the ability of DT-NVS to generate novel views of complex scenes of a large variety of object categories with very different backgrounds. A detailed ablation study is provided to illustrate the successful training of this model. An interesting future direction includes extending this model to more challenging outdoor scenes, higher resolution or associating with other modalities.

Limitation: We needed to downsample images significantly to train this model, resulting in a loss of output quality. This was because we only had access to two A100 GPUs. The model also occasionally struggles with outdoor scenes, partly due to the VM representation. Recent diffusion approaches, like flow models [35], could not be employed. This is because they predict velocity (defined as "noise - ground truth"), but we cannot observe this noise for novel view synthesis. Instead, we apply a mathematical equivalent of v -parameterization when we supervise the model with an L2 loss and other auxiliary losses after volume rendering.

References

- [1] Titas Auciukevicius, Fabian Manhardt, Federico Tombari, and Paul Henderson. Denoising diffusion via image-based rendering. *arXiv preprint arXiv:2402.03445*, 2024.
- [2] Titas Auciukevičius, Zexiang Xu, Matthew Fisher, Paul Henderson, Hakan Bilen, Niloy J Mitra, and Paul Guerrero. Renderdiffusion: Image diffusion for 3d reconstruction, inpainting and generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12608–12618, 2023.
- [3] ShahRukh Athar, Zexiang Xu, Kalyan Sunkavalli, Eli Shechtman, and Zhixin Shu. Rignerf: Fully controllable neural 3d portraits. In *Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [4] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5470–5479, 2022.
- [5] Shengqu Cai, Anton Obukhov, Dengxin Dai, and Luc 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.
- [6] Eric Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *arXiv*, 2020.
- [7] Eric R. Chan, Connor Z. Lin, Matthew A. Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas Guibas, Jonathan Tremblay, Sameh Khamis, Tero Karras, and Gordon Wetzstein. Efficient geometry-aware 3D generative adversarial networks. In *arXiv*, 2021.
- [8] Eric R Chan, Koki Nagano, Matthew A Chan, Alexander W Bergman, Jeong Joon Park, Axel Levy, Miika Aittala, Shalini De Mello, Tero Karras, and Gordon Wetzstein. Genvs: Generative novel view synthesis with 3d-aware diffusion models. *arXiv preprint arXiv:2311.10709*, 2023.
- [9] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. In *European Conference on Computer Vision*, pages 333–350. Springer, 2022.
- [10] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 14124–14133, 2021.
- [11] Yuedong Chen, Haofei Xu, Qianyi Wu, Chuanxia Zheng, Tat-Jen Cham, and Jianfei Cai. Explicit correspondence matching for generalizable neural radiance fields. *arXiv preprint arXiv:2304.12294*, 2023.
- [12] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021.
- [13] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [14] Sara Fridovich-Keil, Giacomo Meanti, Frederik Rahbæk Warburg, Benjamin Recht, and Angjoo Kanazawa. K-planes: Explicit radiance fields in space, time, and appearance. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12479–12488, 2023.
- [15] Guy Gafni, Justus Thies, Michael Zollhöfer, and Matthias Nießner. Dynamic neural radiance fields for monocular 4d facial avatar reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8649–8658, June 2021.

- [16] Jiatao Gu, Alex Trevithick, Kai-En Lin, Joshua M Susskind, Christian Theobalt, Lingjie Liu, and Ravi Ramamoorthi. Nerfdiff: Single-image view synthesis with nerf-guided distillation from 3d-aware diffusion. In *Proceedings of the International Conference on Machine Learning*, 2023.
- [17] Tiankai Hang, Shuyang Gu, Chen Li, Jianmin Bao, Dong Chen, Han Hu, Xin Geng, and Baining Guo. Efficient diffusion training via min-snr weighting strategy. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7441–7451, 2023.
- [18] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16000–16009, 2022.
- [19] Philipp Henzler, Jeremy Reizenstein, Patrick Labatut, Roman Shapovalov, Tobias Ritschel, Andrea Vedaldi, and David Novotny. Unsupervised learning of 3d object categories from videos in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4700–4709, June 2021.
- [20] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *arXiv preprint arXiv:1706.08500*, 2017.
- [21] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022.
- [22] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- [23] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.
- [24] Yang Hong, Bo Peng, Haiyao Xiao, Ligang Liu, and Juyong Zhang. Headnerf: A real-time nerf-based parametric head model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20374–20384, 2022.
- [25] Yicong Hong, Kai Zhang, Jiuxiang Gu, Sai Bi, Yang Zhou, Difan Liu, Feng Liu, Kalyan Sunkavalli, Trung Bui, and Hao Tan. Lrm: Large reconstruction model for single image to 3d. *arXiv preprint arXiv:2311.04400*, 2023.
- [26] Muhammad Zubair Irshad, Sergey Zakharov, Katherine Liu, Vitor Guizilini, Thomas Kollar, Adrien Gaidon, Zsolt Kira, and Rares Ambrus. Neo 360: Neural fields for sparse view synthesis of outdoor scenes. 2023.
- [27] Wonbong Jang and Lourdes Agapito. Codenerf: Disentangled neural radiance fields for object categories. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12949–12958, 2021.
- [28] Wonbong Jang and Lourdes Agapito. Nvist: In the wild new view synthesis from a single image with transformers. *arXiv preprint arXiv:2312.08568*, 2023.
- [29] Heewoo Jun and Alex Nichol. Shap-e: Generating conditional 3d implicit functions. *arXiv preprint arXiv:2305.02463*, 2023.
- [30] Animesh Karnewar, Andrea Vedaldi, David Novotny, and Niloy J Mitra. Holodiffusion: Training a 3d diffusion model using 2d images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18423–18433, 2023.
- [31] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4), July 2023.

- [32] Seung Wook Kim, Bradley Brown, Kangxue Yin, Karsten Kreis, Katja Schwarz, Daiqing Li, Robin Rombach, Antonio Torralba, and Sanja Fidler. Neuralfield-ldm: Scene generation with hierarchical latent diffusion models. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.
- [33] Jonáš Kulhánek, Erik Derner, Torsten Sattler, and Robert Babuška. Viewformer: Nerf-free neural rendering from few images using transformers. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XV*, pages 198–216. Springer, 2022.
- [34] Eric-Tuan Le, Edward Bartrum, and Iasonas Kokkinos. Stylemorph: Disentangled 3d-aware image synthesis with a 3d morphable stylegan. In *The Eleventh International Conference on Learning Representations*, 2022.
- [35] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.
- [36] Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick. Zero-1-to-3: Zero-shot one image to 3d object. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9298–9309, 2023.
- [37] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019.
- [38] Luke Melas-Kyriazi, Iro Laina, Christian Rupprecht, and Andrea Vedaldi. Realfusion: 360deg reconstruction of any object from a single image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023.
- [39] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021.
- [40] Takeru Miyato, Bernhard Jaeger, Max Welling, and Andreas Geiger. Gta: A geometry-aware attention mechanism for multi-view transformers. In *International Conference on Learning Representations (ICLR)*, 2024.
- [41] Norman Müller, Andrea Simonelli, Lorenzo Porzi, Samuel Rota Bulò, Matthias Nießner, and Peter Kontschieder. Autorf: Learning 3d object radiance fields from single view observations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3971–3980, 2022.
- [42] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)*, 41(4):1–15, 2022.
- [43] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3504–3515. IEEE, 2020.
- [44] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4195–4205, 2023.
- [45] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. In *International Conference on Learning Representations*, 2022.
- [46] Daniel Rebain, Mark Matthews, Kwang Moo Yi, Dmitry Lagun, and Andrea Tagliasacchi. Lolnerf: Learn from one look. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1558–1567, 2022.
- [47] Jeremy Reizenstein, Roman Shapovalov, Philipp Henzler, Luca Sbordone, Patrick Labatut, and David Novotny. Common objects in 3d: Large-scale learning and evaluation of real-life 3d category reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10901–10911, 2021.

- [48] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10684–10695. IEEE, 2022.
- [49] Mehdi SM Sajjadi, Henning Meyer, Etienne Pot, Urs Bergmann, Klaus Greff, Noha Radwan, Suhani Vora, Mario Lučić, Daniel Duckworth, Alexey Dosovitskiy, et al. Scene representation transformer: Geometry-free novel view synthesis through set-latent scene representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6229–6238, 2022.
- [50] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. *arXiv preprint arXiv:2202.00512*, 2022.
- [51] Kyle Sargent, Jing Yu Koh, Han Zhang, Huiwen Chang, Charles Herrmann, Pratul Srinivasan, Jiajun Wu, and Deqing Sun. Vq3d: Learning a 3d-aware generative model on imagenet. *arXiv preprint arXiv:2302.06833*, 2023.
- [52] Kyle Sargent, Zizhang Li, Tanmay Shah, Charles Herrmann, Hong-Xing Yu, Yunzhi Zhang, Eric Ryan Chan, Dmitry Lagun, Li Fei-Fei, Deqing Sun, and Jiajun Wu. ZeroNVS: Zero-shot 360-degree view synthesis from a single real image. *arXiv preprint arXiv:2310.17994*, 2023.
- [53] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [54] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise view selection for unstructured multi-view stereo. In *European Conference on Computer Vision (ECCV)*, 2016.
- [55] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. Graf: Generative radiance fields for 3d-aware image synthesis. *Advances in Neural Information Processing Systems*, 33:20154–20166, 2020.
- [56] Katja Schwarz, Seung Wook Kim, Jun Gao, Sanja Fidler, Andreas Geiger, and Karsten Kreis. Wildfusion: Learning 3d-aware latent diffusion models in view space. In *International Conference on Learning Representations (ICLR)*, 2024.
- [57] Bokui Shen, Xincheng Yan, Charles R Qi, Mahyar Najibi, Boyang Deng, Leonidas Guibas, Yin Zhou, and Dragomir Anguelov. Gina-3d: Learning to generate implicit neural assets in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4913–4926, 2023.
- [58] Yichun Shi, Peng Wang, Jianglong Ye, Mai Long, Kejie Li, and Xiao Yang. Mvdream: Multi-view diffusion for 3d generation. *arXiv preprint arXiv:2312.02201*, 2023.
- [59] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. *Advances in Neural Information Processing Systems*, 32, 2019.
- [60] Jascha Sohl-Dickstein, Eric A Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *Proceedings of the 32nd International Conference on Machine Learning*, pages 2256–2265. PMLR, 2015.
- [61] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020.
- [62] Stanislaw Szymanowicz, Christian Rupprecht, and Andrea Vedaldi. Viewset diffusion: (0-)image-conditioned 3d generative models from 2d data. *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023.
- [63] Jiaxiang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang Zeng. Dreamgaussian: Generative gaussian splatting for efficient 3d content creation. *arXiv preprint arXiv:2309.16653*, 2023.

- [64] Ayush Tewari, Tianwei Yin, George Cazenavette, Semon Rezhikov, Joshua B. Tenenbaum, Frédo Durand, William T. Freeman, and Vincent Sitzmann. Diffusion with forward models: Solving stochastic inverse problems without direct supervision. In *arXiv*, 2023.
- [65] Alex Trevithick and Bo Yang. Grf: Learning a general radiance field for 3d representation and rendering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15182–15192, 2021.
- [66] Hengyi Wang, Jingwen Wang, and Lourdes Agapito. Morpheus: Neural dynamic 360 $\{\backslash \deg\}$ surface reconstruction from monocular rgb-d video. *arXiv preprint arXiv:2312.00778*, 2023.
- [67] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul P Srinivasan, Howard Zhou, Jonathan T Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. Ibrnet: Learning multi-view image-based rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4690–4699, 2021.
- [68] Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Pro-lificdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation. *Advances in Neural Information Processing Systems*, 2023.
- [69] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004.
- [70] Daniel Watson, William Chan, Ricardo Martin-Brualla, Jonathan Ho, Andrea Tagliasacchi, and Mohammad Norouzi. Novel view synthesis with diffusion models. *arXiv preprint arXiv:2210.04628*, 2022.
- [71] Chao-Yuan Wu, Justin Johnson, Jitendra Malik, Christoph Feichtenhofer, and Georgia Gkioxari. Multiview compressive coding for 3D reconstruction. *arXiv:2301.08247*, 2023.
- [72] Rundi Wu, Ben Mildenhall, Philipp Henzler, Keunhong Park, Ruiqi Gao, Daniel Watson, Pratul P. Srinivasan, Dor Verbin, Jonathan T. Barron, Ben Poole, and Aleksander Holynski. Reconfusion: 3d reconstruction with diffusion priors. *arXiv*, 2023.
- [73] Jingjing Xu, Xu Sun, Zhiyuan Zhang, Guangxiang Zhao, and Junyang Lin. Understanding and improving layer normalization. In *Advances in Neural Information Processing Systems. NeurIPS*, 2019.
- [74] Yinghao Xu, Hao Tan, Fujun Luan, Sai Bi, Peng Wang, Jiahao Li, Zifan Shi, Kalyan Sunkavalli, Gordon Wetzstein, Zexiang Xu, et al. Dmv3d: Denoising multi-view diffusion using 3d large reconstruction model. *arXiv preprint arXiv:2311.09217*, 2023.
- [75] Brent Yi, Weijia Zeng, Sam Buchanan, and Yi Ma. Canonical factors for hybrid neural fields. In *International Conference on Computer Vision (ICCV)*, 2023.
- [76] Alex Yu, Sara Fridovich-Keil, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. *arXiv preprint arXiv:2112.05131*, 2021.
- [77] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4578–4587, 2021.
- [78] Xianggang Yu, Mutian Xu, Yidan Zhang, Haolin Liu, Chongjie Ye, Yushuang Wu, Zizheng Yan, Chenming Zhu, Zhangyang Xiong, Tianyou Liang, et al. Mvimgnet: A large-scale dataset of multi-view images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9150–9161, 2023.
- [79] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018.