

Free3D: Consistent Novel View Synthesis without 3D Representation

Chuanxia Zheng Andrea Vedaldi

Visual Geometry Group, University of Oxford

{cxzheng, vedaldi}@robots.ox.ac.uk

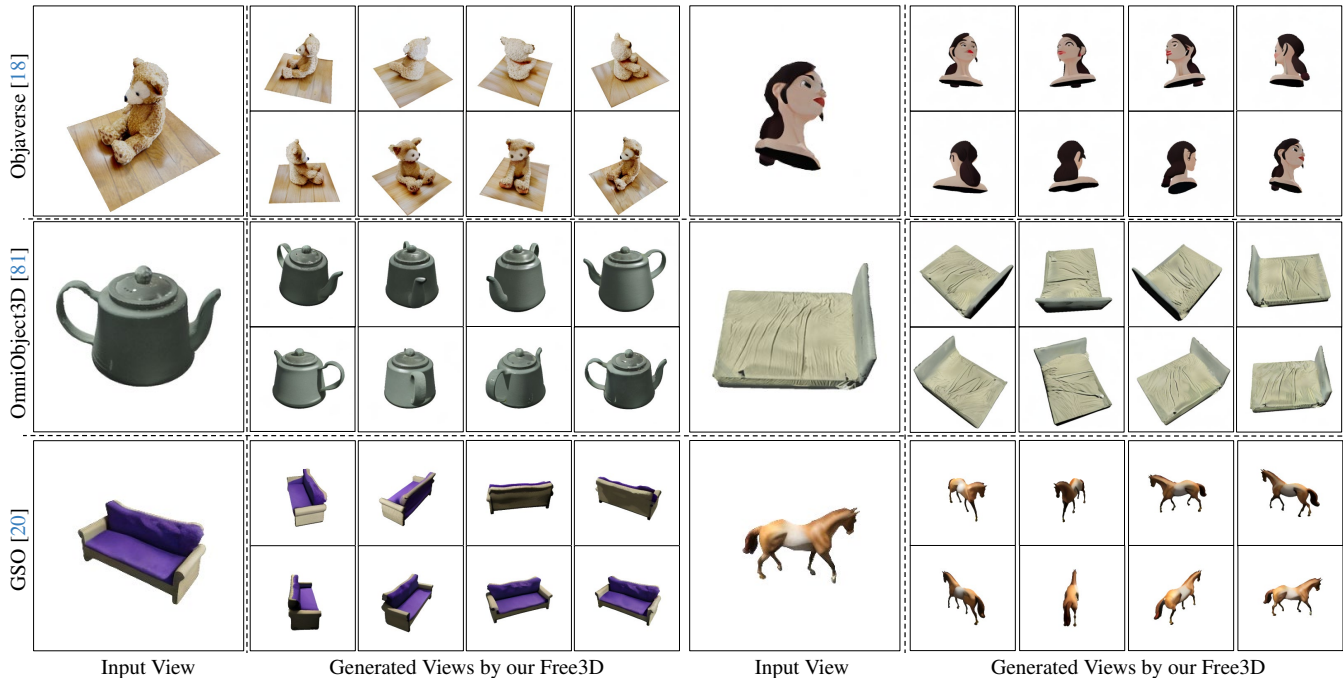


Figure 1. **Free3D novel view synthesis in an open-set setting.** Given a single input view, our approach synthesizes consistent 360° videos accurately and with no explicit 3D representation. Trained on Objaverse only, it generalizes well to new datasets and categories.

Abstract

We introduce *Free3D*, a simple approach designed for open-set novel view synthesis (NVS) from a single image. Similar to *Zero-1-to-3*, we start from a pre-trained 2D image generator for generalization, and fine-tune it for NVS. Compared to recent and concurrent works, we obtain significant improvements without resorting to an explicit 3D representation, which is slow and memory-consuming or training an additional 3D network. We do so by encoding better the target camera pose via a new per-pixel ray conditioning normalization (RCN) layer. The latter injects pose information in the underlying 2D image generator by telling each pixel its specific viewing direction. We also improve multi-view consistency via a light-weight multi-view attention layer and multi-view noise sharing. We train *Free3D* on the *Objaverse* dataset and demonstrate excellent generalization to various new categories in several new datasets, including *OminiObject3D* and *GSO*. We hope our simple

and effective approach will serve as a solid baseline and help future research in NVS with more accuracy pose. The project page is available at <https://chuanxiaz.com/free3d/>.

1. Introduction

Novel view synthesis (NVS) has developed rapidly in the recent past [5, 6, 11, 38, 49], fuelled in part by NeRF [48]. However, many such NVS methods require to optimize a new model from scratch for each scene, and are thus impractical in many applications. Inspired by generative models [25, 39, 67], other NVS methods generalize beyond a single scene to a specific *object category*, or even to *unstructured datasets of objects* [32, 34, 47, 50, 69, 76, 85]. Remarkably, these methods, by utilizing the data prior, can generate new views starting from a single image. However, they trade off some quality in order to achieve generality,

which is still insufficient to build a 3D foundation model. Going beyond, in this work we consider the NVS problem in an *open-set* setting, generalizing beyond a single object, category, and dataset, yet still utilizing a single image as input, and aiming at obtaining high-quality outputs.

There are two primary quality targets: (i) The output must accurately reflect the pose of the target cameras, and (ii), when several views of the same object are generated, they must be multi-view consistent. In order to achieve high 3D accuracy and consistency, recent and concurrent methods [42, 44, 63, 83] use a 3D representation of the reconstructed object, usually on top of a pre-trained 2D generative model [57]. It works well but adds complexity.

In this paper, we introduce Free3D, a simple, flexible, and efficient method that can also achieve consistent NVS results *without the need to rely on a resource-intensive 3D representation*. Zero-1-to-3 [43] is perhaps the best-known example of such a 3D-free NVS system. Like it, Free3D builds upon a pre-trained 2D generative model like Stable Diffusion [57], trained on a large-scale dataset (over 5 billion images [62]), as a data prior. The prior knowledge contained in such a 2D generator is extremely important to be able to ‘guess’ plausible novel views of arbitrary new objects, which is inherently highly ambiguous. However, we show empirically that Zero-1-to-3, has, in practice, poor camera pose control, and, when tasked with generating multiple views, not very consistent. The latter is unavoidable in their design because each view is sampled independently from scratch and, owing to the ambiguity of reconstruction, there is no reason why compatible images would be generated each time. We refer readers to the supplemental videos for a comprehensive comparison.

To mitigate these issues, we first show that much better camera control can be achieved by switching to a different representation of camera pose. Specifically, we introduce a *ray conditioning normalization* (RCN) (Fig. 2 (b)), which tells *each pixel* its viewing direction. This is a *distributed* representation, which should be contrasted to the *concentrated* camera representation used in [43]. They pass the camera elevation, azimuth, and distance as language-like tokens that may be difficult for the network to interpret and utilize [61]. In contrast, with our RCN layer, we show how to effectively incorporate this *per-pixel* ray information in an existing text-to-image diffusion model, which empirically leads to significantly more accurate NVS in our experiments (Tabs. 1 and 2). It is also reminiscent of the design of methods like NeRF [48] and LFNs [66], which also work by processing individual rays.

While RCN leads to more accurate camera control, it *cannot* improve multi-view consistency by itself. To improve the latter, we introduce an additional *pseudo-3D cross-view attention* module (Fig. 2(c)), inspired by video diffusion models [8, 26, 30, 64, 80], that fuses information

across all views instead of processing each view independently. Furthermore, we use *multi-view noise sharing* when generating the different views of the object, which further enhances consistency due to the continuity of the denoising function, reducing aleatoric variations from view to view.

We benchmark Free3D against recent and concurrent state-of-the-art methods [17, 43, 44, 77] on *open-set* NVS. Although the model is trained on only one dataset, it generalizes well to all recent NVS benchmark datasets, including Objaverse [18], OmniObject3D [81], and Google Scanned Object (GSO) [20]. A thorough experimental assessment shows that our approach consistently outperforms existing methods, both quantitatively and qualitatively.

To summarise, with Free3D we make the following contributions: (i) We introduce the *ray conditioning normalization* (RCN) layer and show that representing the camera by utilizing a combination of distributed ray conditioning and concentrated pose tokens significantly improves pose accuracy in NVS. (ii) We show that a small *multi-view attention* module can improve multi-view consistency by exchanging information between views at a low cost. (iii) We find that *multi-view noise sharing* between different views further improves consistency. (iv) We demonstrate empirically that Free3D achieves consistent NVS *without needing a 3D representation* and outperforms the existing state-of-the-art models on both pose accuracy and view consistency.

2. Related Work

Per-Scene NVS. Early NVS works relied on epipolar geometry to interpolate between different views of the same scene [14, 16]. A recent breakthrough was to represent 3D scenes as implicit neural fields, as proposed by SRN [65], DeepSDF [51], NeRF [48] and LFN [66], and further improved in follow-ups [5, 6, 11, 24, 38, 49, 87]. Even so, data efficiency, generalizability, and robustness remain a limitation: *such systems require multiple views to learn the 3D representation from scratch of every single scene*.

To bypass the need for multiple input views, DreamFusion [52] proposed to distill 3D models from a large-scale pre-trained 2D diffusion model [58]. RealFusion [46] extended the latter to single-view image reconstruction by adding the input image as a constraint during distillation. Several follow-up works [53, 54, 61, 70] achieved further improvements to the resolution and quality of the resulting 3D assets. Although these models are open-ended [46], they require lengthy *per-scene* optimisation.

Category-Specific NVS. Inspired by the success of 2D generation [25, 39, 67], some work [15, 50, 71, 78] built the autoencoder architecture for NVS. Driven by the effectiveness of light field rendering (LFR) [1, 2], other generalizable NVS approaches [59, 60, 68] query a network for colour of different rays. 3DiM [76] then introduced diffusion models

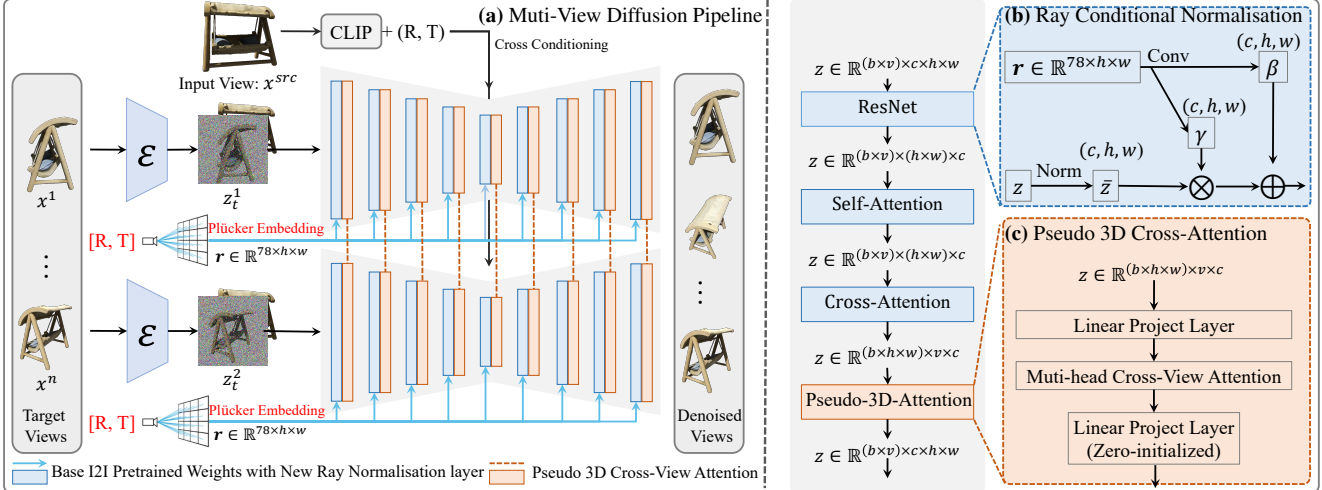


Figure 2. **The overall pipeline of our Free3D.** (a) Given a single source input image, the proposed architecture jointly predicts multiple target views, instead of processing them independently. To achieve a consistent novel view synthesis without the need for 3D representation, (b) we first propose a novel *ray conditional normalization* (RCN) layer, which uses a *per-pixel* oriented camera ray to modulate the latent features, enabling the model’s ability to capture more precise viewpoints. (c) A memory-friendly *pseudo-3D cross-attention* module is introduced to efficiently bridge information across multiple generated views. Note that, here the similarity score is only calculated across multiple views in temporal instead of spatial, resulting in a minimal computational and memory cost.

into NVS. This is followed by RenderDiffusion [3], HoloDiffusion [36], SparseFusion [90], GeNVS [9], Viewset Diffusion [69] and LFD [82]. However, these methods train their models from scratch using pure 3D data. Because such datasets are *not* sufficient for *open-set* generalization, they are limited to one or few categories. Here we start from an open-ended 2D image generator instead.

Open-set NVS. In order to operate in an open-set setting, Zero-1-to-3 [43] built on a 2D image generator trained on a large-scale image dataset [62] and fine-tune it on the Objaverse dataset [18]. While it has good generalizability, Zero-1-to-3 fails to achieve high pose accuracy, and their reconstructions are not very consistent across views. To mitigate these issues, other concurrent approaches either integrate 3D representations into the network [35, 44, 83] or train another auxiliary 3D network [41, 42, 63, 77]. These methods can output high-quality target views, but they are computationally expensive and use significant memory for training. More importantly, these methods do *not* address the issue of pose representation, which we find to be a key bottleneck in NVS. In contrast, our method is *3D-free* and achieves comparable or better NVS quality, due to *ray conditioning normalization*, *multi-view attention*, and *multi-view noise sharing*. A more recent concurrent work is Stable Video Diffusion [7], which also synthesizes multi-view videos. However, they are trained with fixed frames with limited elevation angles. In contrast, our Free3D is trained to handle arbitrary camera viewpoints, and ideally generates 360° spectrum, with flexible-frames for any elevation angles.

3. Method

Our goal is to learn a model Φ that, given as input a source image x^{src} and a sequence of camera poses $\mathcal{P} = \{\mathbf{P}^i\}_{i=1}^N$, $\mathbf{P}^i = (\mathbf{K}^i, \mathbf{R}^i, \mathbf{T}^i)$ comprising intrinsics \mathbf{K}^i , rotation \mathbf{R}^i and translation \mathbf{T}^i , synthesizes corresponding novel views $\{x^i\}_{i=1}^N$ which are accurate and consistent *without relying on an explicit 3D representation*. We approach this problem by generating all the views together, conditioned on the source image x^{src} and a series of camera parameters. In this way, the network has a chance to produce several consistent views together. The *broad motivation* here is that the 1D sequential representation has been used extensively in the 2D generation [13, 23, 56, 73, 74, 75, 89]. Analogously, it should *be able to produce visually consistent video representing the underlying 3D structure using 2D sequential representation*, i.e. *multi-frames* [4, 7]. However, we tackle it in a bidirectional manner [10, 22, 40, 86, 89], instead of sequentially auto-regressive generation.

To realize this goal, we must tackle two challenges: (i) ensuring that the model accurately captures the target view and (ii) ensuring that different views of the object are consistent in geometry and appearance. To achieve these, our framework, illustrated in Fig. 2, extends a 2D generator by injecting at each layer a *ray-conditional normalization* (RCN) layer (Sec. 3.1) as well as *pseudo-3D attention* layer (Sec. 3.2). It also utilizes *multi-view noise sharing* (Sec. 3.2). The former captures pose more accurately, whereas the last two improve consistency. Although several prior NVS works have started from 2D generators and also

focused on improving multi-view consistency, with RCN we are the first (to the best of our knowledge) to use ray conditioning in this setting and obtain high pose accuracy without a 3D representation or an additional 3D network.

3.1. Ray Conditioning Normalization (RCN)

Ray conditioning was originally proposed for NVS by [66]. However, it has *not* been used for NVS models that build on 2D generators like Zero-1-to-3 [43] — in practice, doing so is essential for generalization due to the lack of very large-scale 3D data. Here, we thus develop a method to extend a pretrained 2D generator with ray conditioning.

Ray Conditioning Embedding. Given a target view $\mathbf{P} = (\mathbf{K}, \mathbf{R}, \mathbf{T})$, for each pixel (u, v) in the image, we define the Plücker coordinates $\mathbf{r}_{uv} = \phi(\mathbf{o}, \mathbf{d}_{uv}) = (\mathbf{o} \times \mathbf{d}_{uv}, \mathbf{d}_{uv}) \in \mathbb{R}^6$, of the ray going from the camera center $\mathbf{o} \in \mathbb{R}^3$ through the pixel, while $\mathbf{d}_{uv} = \mathbf{R}^\top (\mathbf{K}^{-1}(u, v, 1)^\top - \mathbf{T}) \in \mathbb{R}^3$ is the ray direction. This encoding was originally introduced by LFN [66]. It is invariant to shifting the camera along the ray, meaning that $\phi(\mathbf{o} + \lambda \mathbf{d}, \mathbf{d}) = ((\mathbf{o} + \lambda \mathbf{d}) \times \mathbf{d}_{uv}, \mathbf{d}_{uv}) = (\mathbf{o} \times \mathbf{d}_{uv}, \mathbf{d}_{uv}) = \phi(\mathbf{o}, \mathbf{d})$, which matches the fact that light propagates in straight lines.

Ray Conditioning Architectures. Our goal is to modify the denoising neural network $\hat{e}(z, t, y)$ so that the conditioning information y includes ray conditioning. We experiment with a number of different architectures to do so, proposing various ray conditioning layers:

- *Concatenation.* Following Zero-1-to-3 [43], a natural choice is to concatenate the noised target z_t^{tgt} , original source embedding z^{src} , and ray conditioning at the input. Then, the input is now $(z_t^{\text{tgt}}, z^{\text{src}}, \mathbf{r})$ instead of z_t^{tgt} alone.
- *Multi-scales concatenation.* We further consider concatenating ray embeddings \mathbf{r} to each intermediate layer in the UNet ϵ_θ . Note that each layer operates at a different resolution, so this amounts to injecting the information \mathbf{r} at different scales.
- *Ray conditioning normalization (RCN).* Finally, we propose to combine the adaptive layer norm [21, 31, 37, 89] with ray conditioning to modulate the image latents. Specifically, the activation latent F_i of the i -th layer in the UNet ϵ_θ is modulated by:

$$\text{ModLN}_{\mathbf{r}}(F_i) = \text{LN}(F_i) \cdot (1 + \gamma) + \beta, \quad (1)$$

where $(\gamma, \beta) = \text{MLP}^{\text{mod}}(\mathbf{r})$ are the learned *scale* and *shift* from ray conditioning \mathbf{r} . Such modulation is applied to each sub-module of UNet ϵ_θ (Fig. 2 (a) & (b)).

Interestingly, while RCN works best, we show empirically (Tab. 3) that all such architectures lead to strong improvements. This confirms the importance of ray conditioning for the significant improvements, with the novel RCN emerging as the frontrunner in this progression.

Discussion. RCN differs significantly from Zero-1-to-3 [43] and follow-ups [33, 42, 44, 77, 83] which coded \mathbf{P} as *global* tokens. Ray-Cond [12] and LFD [82] also applied ray conditioning for NVS in GAN and Diffusion, respectively, but just concatenated it as additional channels. Besides, they are only trained for *category-specific* NVS. Similar to *ray-casting* [45, 48], RCN considers one ray per pixel. However, it eschews the computation of hundreds of samples per ray. Because of this, Free3D dramatically reduces the rendering time and memory consumption compared to the concurrent works of [42, 44, 83].

3.2. View Consistent Rendering

Given a source image x^{src} , our goal is to render a series of consistent novel views x^i . While the camera encoding technique of Sec. 3.1 significantly enhances the pose accuracy of our model, if images are sampled independently, they will almost never be visually consistent due to the intrinsic ambiguity of the reconstruction task. This is visible as temporal flickering in the rendered 3D video (as shown in supplemental videos). To remove or at least greatly mitigate this problem, we propose to sample images jointly as well as to share noise among them.

Multi-view attention. We first adapt the frame attention, a well-established method in video diffusion models [8, 26, 30, 64, 80], to capture temporal dependencies across different views. As shown in Fig. 2(c), given a 5D latent $z \in \mathbb{R}^{B \times v \times c \times h \times w}$, we initially reshape it to $z \in \mathbb{R}^{(B \times h \times w) \times v \times c}$, resulting in *batch* \times *height* \times *width* sequences at the length of *views*. Subsequently, this reshaped latent is passed through the pseudo-3D attention module to calculate the similarity across different views. Since this attention layer operates across views but separately for each spatial location, the computational and memory costs are quite low (Tab. 1). Similar to ray conditioning, we inject multi-view attention at each level of the UNet \hat{e}_θ .

Multi-view noise sharing. At test time, views are sampled using the backward process Eq. (A.2) starting from a normally-distributed, yet randomly-sampled noise vector x_T . The idea is that the noise distribution is a representation of the data distribution, and models aleatoric uncertainty. In our case, while we are interested in generating *multiple* and *diverse* solutions, we still wish the different views to be *internally consistent* for each 360° rendered video.

In order to reduce the variance between different views, we propose to start sampling each view from the *same* noise vector x_T . The network ϵ_θ still generates different views because it is conditioned on the camera parameters. It can also generate different reconstructions by resampling x_T . However, the variance *between views* is thus reduced. Noise sharing can be justified by noting that the network $\hat{e}_\theta(z_t, t, y)$ is a continuous function of both z_t and y [79].

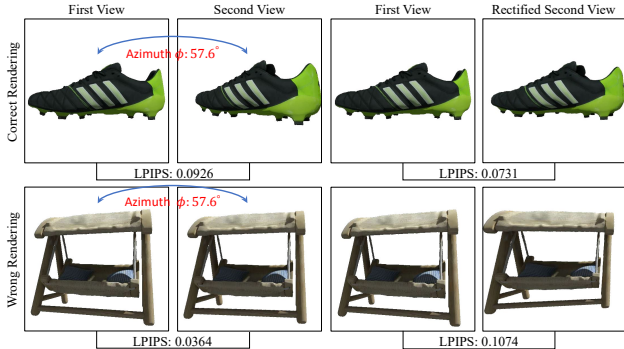


Figure 3. **Perceptual Path Length Consistency (PPLC)**. To partly compensate for the viewpoint change, the second image is rectified w.r.t. the first before comparison. To illustrate the importance of using rectification, the figure shows two objects in a large azimuth $\phi : 57.6^\circ$. The top row shows the left an ideally-rendered image pair, which however attains a large LPIPS loss due to the viewpoint change. To the right, rectification reduces this score. The bottom row shows the opposite, where a pair of incorrectly rendered views has its LPIPS loss increased by rectification.

Besides, different from [72], we do not specially set the noise schedule in different time steps t . More complex designs and training strategies have the potential to improve performance but are not the focus of this work.

3.3. Learning formulation

The learning objective is given by

$$\mathcal{L} = \mathbb{E}_{(Z_0, z^{\text{src}}, \mathbf{P}), \epsilon, t} [\|\epsilon - \epsilon_\theta(Z_t, t, y)\|_2^2], \quad (2)$$

where each training sample $(Z_0, z^{\text{src}}, \mathbf{P})$ consist of N encoded target views $Z_0 = \{\mathcal{E}(x^i)\}_{i=1}^N$, the encoded source view $z^{\text{src}} = \mathcal{E}(x^{\text{src}})$ and the viewpoints \mathbf{P} . The network is conditioned on the source view and cameras, *i.e.*, $y = (z^{\text{src}}, \mathbf{P})$ and, as in Fig. 2, estimates the noise for all target views jointly. This information is incorporated as discussed above via ray conditioning modulation; following [43], we also concatenate the input image code z^{src} to z_t along the channel dimension and add the CLIP [55] encoding.

3.4. Perceptual Path Length Consistency

To quantify the consistency across different novel views, recent methods [43, 76, 77, 83] trained NeRF models from the sampled views and evaluated them on the remaining views. However, this requires training a NeRF for each test instance, which is unfeasible for large-scale testing required for proper open-set NVS evaluation. Here, we propose to use instead the pairwise perceptual distance [88] to measure the consistency between generated views. In particular, we first subdivide the 360° rendering path into linear segments and then calculate the LPIPS score between two neighboring generated images x^i and x^{i+1} . Naturally,

Method	Objaverse[18]				
	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	PPLC \downarrow	Time \downarrow
Zero-1-to-3[43]	0.8462	0.0938	1.52	18.84	3s/44s
Zero123-XL[17]	0.8339	0.1098	1.67	25.61	3s/44s
SyncDreamer[44]	0.8063	0.1910	7.57	16.32	25s/77s
Consistent123[77]	0.8530	0.0913	1.48	17.89	4s/63s
Ours (Free3D)	0.8620	0.0784	1.21	10.82	3s/52s

Table 1. **Comparison with SoTA methods** on all 7,729 objects in Objaverse test-set. Recent works like [17, 43, 76, 77] were originally evaluated using a subset of this data only due to the cost of training additional 3D models. Unlike them, we directly evaluate the model on whole test-set without using additional 3D networks. The inference time is for rendering a single target view and a 360° video, respectively.

these images differ due to the different viewpoints (Fig. 3 to row). We partly compensate for the viewpoint change by rectifying [27] the second image w.r.t. the first. We then define Perceptual Path Length Consistency (PPLC) of one rendered sequence $\{x^i\}_{i=1}^N$ as follows:

$$l_{\text{pplc}} = \mathbb{E} \left[\frac{1}{\phi^2} \|\mathcal{F}(\text{Rect}(x^i)) - \mathcal{F}(\text{Rect}(x^{i+1}))\|_2^2 \right], \quad (3)$$

where ϕ is the degree between views x^i and x^{i+1} , which is set as $\phi = 2\pi/50$, with the azimuth 7.2° in all our video rendering. \mathcal{F} is a pre-trained network to ensure the metric matches with human perceptual similarity judgment.

4. Experiments

We compare Free3D to state-of-the-art open-ended single-image NVS methods and also ablate our design choices.

4.1. Experimental Details

Datasets. For fairness, our model is trained using the exact same protocol as Zero-1-to-3 [43]. They render multiple views for 772,870 objects from Objaverse dataset [18]. We use the identical test split as they do. To assess how well our mode generalizes to other datasets, and how it compares to other models, we consider two more datasets: OmniObject3D [81] and Google Scanned Objects (GSO) [20], which contain real-life scanned objects. Since we do not use these datasets for training at all, we use the *entirety* of OmniObject3D and GSO objects for evaluation (6,000 and 1,030 objects, respectively).

Metrics. We follow [43, 44, 77] and assess the NVS quality by comparing the generated images and the ground-truth views at different levels of granularity, including *pixel*-level Peak Signal-to-Noise Ration (PSNR), *patch*-level Structure SIMilarity index (SSIM) and *feature*-level Learned Percep-

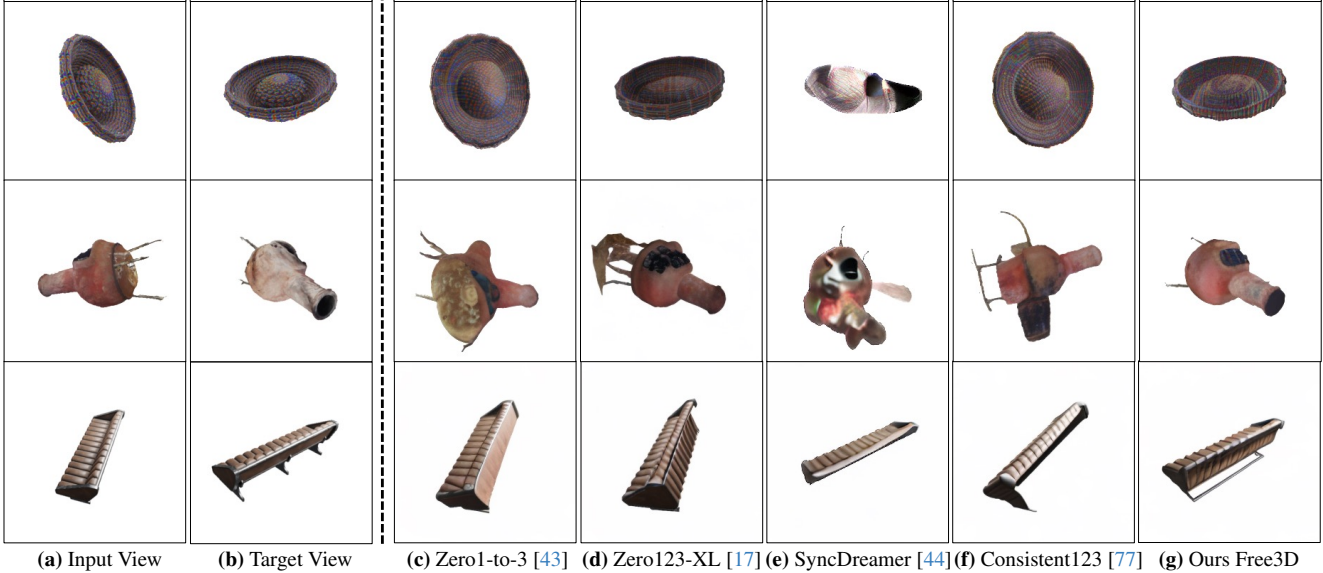


Figure 4. **Qualitative comparisons** on Objaverse. Given a target pose, our Free3D significantly improves the accuracy of the generated pose compared to existing state-of-the-art methods. Note that Zero123-XL [17] is trained on the much larger Objaverse-XL dataset [17], which contains 10 million 3D objects. More comparisons are provided in the supplement Figs. C.1 and C.2.

tual Image Patch Similarity (LPIPS) [88]¹, and *dataset-level* Fréchet Inception Distance (FID) [28]². Some of these metrics look at the statistics of individually-reconstructed images rather than exact reconstructions, which is the correct approach given that the reconstruction task is inherently ambiguous; however, for the same reason, they are *not suitable* for assessing multi-view consistency. Instead, we use PPLC (Sec. 3.4) to evaluate consistency. To do so, we generate 50-frame videos along a pre-defined circular camera trajectory for each object and then compute the PPLC score between neighbouring frames.

Baselines. We compare Free3D to the state-of-the-art Zero-1-to-3 [43]_{ICCV’2023} and the follow-ups Zero123-XL [17]_{NeurIPS’2023}, SyncDreamer [44]_{arXiv’2023} and Consistent123 [77]_{arXiv’2023}, which are built on Zero-1-to-3. While noticed other concurrent works in arXiv for NVS [35, 83, 84], but, as no codes are available yet, we could not re-implement and re-train all of them for comparison.

4.2. Assessing Quality

Quantitative comparison. We first evaluate Free3D and other methods on the Objaverse dataset [18, 43], where the model is trained. Unlike previous work [43, 44, 77] which consider only a subset of the test data due to expensive post-processing, we use the *entire test set* of 7,729 3D objects. Quantitative results in Tab. 1 show that Free3D outperforms state-of-the-art models. This includes Zero123-XL [17] and SyncDreamer [44], which are trained on a

much larger dataset and with explicit 3D volume representation, respectively. While the concurrent Consistent123 [77] also utilizes a form of multi-view diffusion with cross-view attention, our Free3D significantly improves the quantitative results (16% relative improvement on LPIPS) on the large evaluation dataset. This suggests that our RCN layer is the primary reason for the observed improvements.

Qualitative comparison. Qualitative results on various categories are visualized in Fig. 4. Free3D achieves better results even under challenging viewpoints with very different categories. SyncDreamer [44] uses an *explicit volumetric representation* representation, but can only generate views with fixed elevation (30°), leading to worse results on synthesizing arbitrary target views. While the concurrent Consistent123 [77] aims to improve rendering consistency with a version of multi-view attention, they cannot directly improve pose accuracy. The Zero123-XL [17] learns to capture pose better than the baseline [43] by training on the larger Objaverse-XL, but the pose is still *not* very accurate. Because of the RCN layer, Free3D shows no such pose errors and results in better images.

4.3. Assessing Generalization

In Tab. 2 and Fig. 5, we validate the ability of Free3D to generalize to datasets not seen during training. This includes the OmniObject3D and the GSO datasets, with 6,000 scanned objects in 190 categories and 1,030 scanned objects in 17 categories, respectively. For this result, we directly test all trained models without any fine-tuning.

¹<https://github.com/richzhang/PerceptualSimilarity> “squeeze”-net.

²<https://github.com/GaParmar/clean-fid> “clip_vit_b_32”-net.

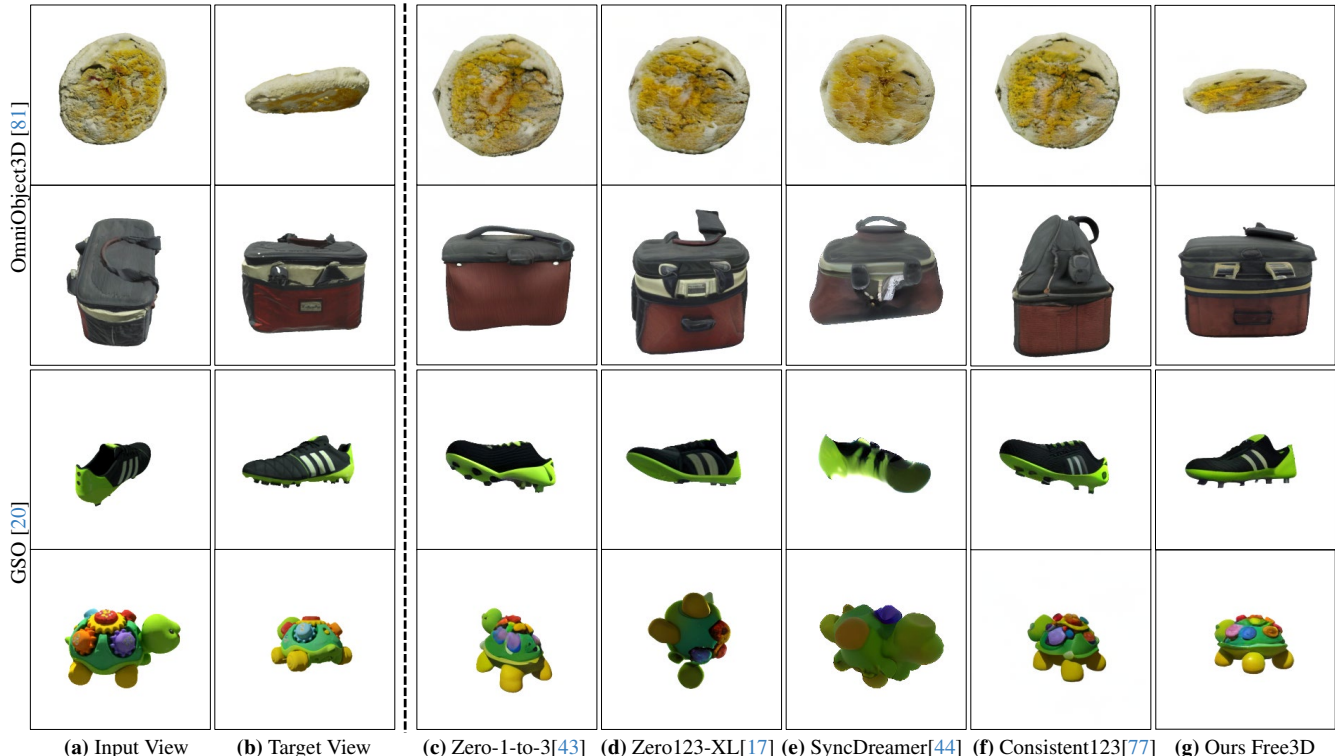


Figure 5. **Qualitative comparisons** on OmniObject3D (top two rows) and GSO (bottom two rows) dataset. Interestingly, exciting methods cannot deal with unconventional objects, such as the “pie” in the first row, while our Free3D is still robust for such a challenging scenario. More comparisons are provided in supplemental Figs. C.3 and C.4.

Method	OmniObject3D [81]					GSO [20]				
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	PPLC \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	PPLC \downarrow
Zero-1-to-3 [43]	16.84	0.7813	0.1321	1.73	24.58	19.65	0.8501	0.0758	3.24	33.15
Zero123-XL [17]	17.11	0.7818	0.1291	1.51	21.33	20.43	0.8589	0.0706	3.23	28.03
SyncDreamer [44]	17.00	0.7941	0.1442	6.58	11.49	14.72	0.7835	0.1533	8.65	9.42
Consistent123 [77]*	17.13	0.7821	0.1255	1.55	18.02	20.11	0.8553	0.0716	3.24	20.08
Ours (Free3D)	18.23	0.8090	0.0996	1.34	8.67	21.13	0.8686	0.0619	2.85	9.10

Table 2. **Generalizable results on unseen datasets**, including OmniObject3D [81] and GSO [20], with 6,000 and 1,030 3D instances, respectively. Note that, although Zero123-XL [17] is trained on a larger dataset, and shows better generalizability, the proposed Free3D still significantly outperforms it with *precise* pose estimation for target views.

Quantitative comparison. In Tab. 2, Free3D, which uses the same training set as Zero-1-to-3 [43], outperforms the baseline and all state-of-the-art variants, including Zero123-XL [17], which is trained on a larger 3D dataset, and SyncDreamer [44], which utilizes a heavier 3D volumetric representation. SyncDreamer improves the baseline on OmniObject3D, which has small elevation views change, while achieving worse results on GSO with random elevation angles. The concurrent Consistent123 [77] is also trained multiple views jointly, with additional cross-view attention. However, their relative improvement is limited. Table 2 shows that Free3D achieves

very substantial improvements on all instantiations on both OmniObject3D and GSO datasets. This further indicates that ray conditioning can successfully improve the pose accuracy of the target view.

Qualitative comparison. A qualitative comparison is given in Fig. 5 (more in supplemental Figs. C.3 and C.4). Though Free3D is only trained on objaverse dataset, it works quite well in the open-set setting. Moreover, it shows significantly better results than all state-of-the-art models, and even better than the Zero123-XL [17], a concurrent model trained on much larger datasets.

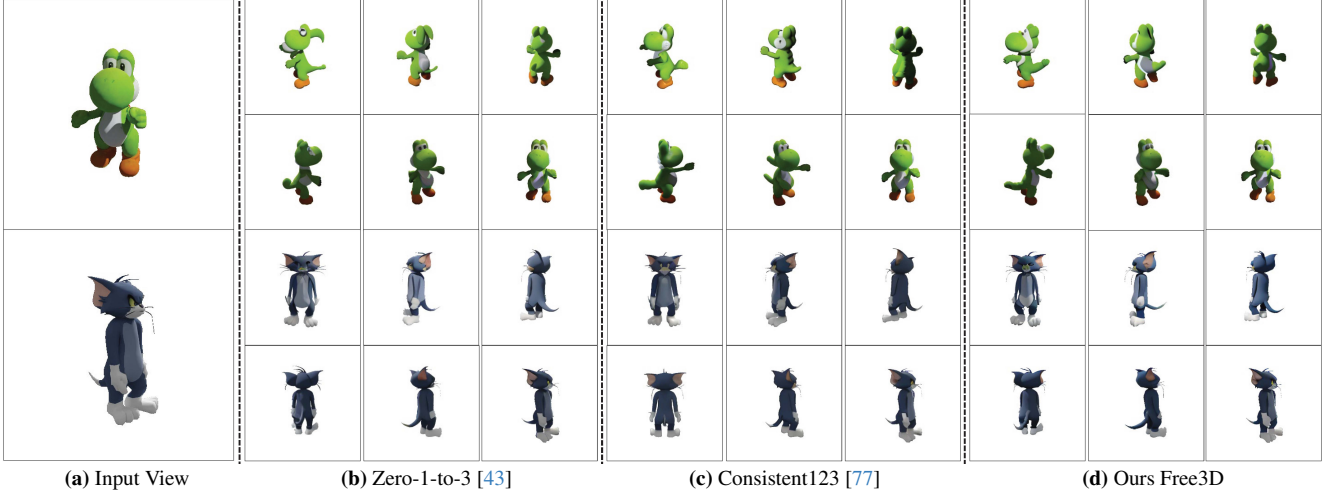


Figure 6. **Qualitative new-view synthesis comparisons.** Zero-1-to-3 [43] generates diverse details (e.g. various ears and tails) for different views in one sampling. Consistent123 [77] improves it through the multi-view diffusion with the cross-view attention. However, it still requires training additional NeRF for 3D reconstruction. For more comparisons, see the supplemental videos for better visualisation.

Method	Objaverse [18]					GSO [20]				
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	PPLC \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	PPLC \downarrow
A Baseline Zero-1-to-3 [43]	19.65	0.8462	0.0938	1.52	22.10	19.65	0.8501	0.0758	3.24	33.15
B + Input Ray Embeddings	20.21	0.8550	0.0858	1.26	16.63	20.49	0.8617	0.0677	3.01	22.08
C + Multi-Scale Ray Emb.	20.56	0.8609	0.0797	1.30	15.94	20.50	0.8615	0.0667	2.98	20.09
D + RCN	20.78	0.8620	0.0784	1.21	15.67	21.13	0.8686	0.0619	2.85	18.48
E + Pseudo-3D attention	20.81	0.8620	0.0781	1.25	14.76	21.20	0.8697	0.0617	2.86	17.39
F E + noise sharing	—	—	—	—	11.39	—	—	—	—	9.10

Table 3. **Ablations** of Free3D design choices. Here, to reduce the computational cost and testing time, we only evaluate 2,000 instances in Objaverse dataset, instead of running the whole dataset with 7,729 for 50-frame 360° video rendering. Therefore, the PPLC score is slight different from the values reported in Tab. 1.

4.4. Assessing 3D consistency

To assess 3D consistency, we render a 360° video with the fixed elevation angle $\theta := 0^\circ$ and 50 azimuth angles uniformly sampled in $[0^\circ, 360^\circ]$. The quantitative results are reported in Tabs. 1 and 2 using the proposed PPLC score, and the qualitative results are shown in Figs. 1 and 6. Consistent123 [77] still requires training additional NeRF to achieve 3D consistent video, while the directly rendered videos are still flickering. SyncDreamer [44] obviously improves view consistency by using an explicit 3D volume representation. However, it is trained only on a fixed elevation angle and can generate only 16 fixed frames for a video in its code. Interestingly, Free3D, by using *multi-view attention* and *multi-view noise sharing* is nearly as effective (while being much cheaper). This is also partly due to ray conditioning, which can capture more precise poses of the target view and thus reduce ambiguity. The additional results are provided as videos in the supplemental materials.

4.5. Ablation Study

We run a number of ablations to analyse Free3D, testing also the variants described in Sec. 3. Results are shown in Tab. 3 and discussed in detail next.

Our baseline configuration (denoted A) is the same as Zero-1-to-3 [43], which is derived from the SD model [57] by replacing text conditioning with source image and target pose conditioning. Then, in B we extend this model with ray conditioning by concatenating the ray embeddings to the source image x^{src} . This alone improves the performance dramatically in both Objaverse and GSO. In C we test injecting ray conditioning into each level of the diffusion UNet ϵ_θ , but the generalizable performance remains similar to B in the unseen GSO dataset. In D, we test instead the RCN layer, which results in further improvements on not only Objaverse dataset, but also on the new GSO. In E, we add *multi-view attention* to exchange information between different frames. As expected, this does *not* improve

the metrics that measure the quality of individual views, but it slightly improves the PPLC score, measuring consistency. In \mathbb{F} , we further add *multi-view noise sharing*, which significantly enhances the consistency between different views, while preserving the quality of single-view rendering. Rendered videos are provided in the supplement.

5. Conclusion

We have introduced Free3D, an open-set single-view NVS method with state-of-the-art performance on various categories, yet bypass the requirement of building on heavy 3D representation or training additional auxiliary 3D models. It is a simple approach that (i) obtains data prior from an off-the-shelf pre-trained 2D image generator, (ii) injects ray conditioning utilizing the new RCN layer to accurately code for the target pose, and (iii) combines that with multi-view attention and noise sharing to improve multi-view consistency. Experimental results show that Free3D significantly outperforms recent and concurrent state-of-the-art NVS models without incurring the cost of utilizing a 3D representation. We hope that Free3D will serve as a new strong baseline for single-image NVS and inspire future research in this area.

Acknowledgements. This research is supported by ERC-CoG UNION 101001212. Many thanks to Stanislaw Szymanowicz, Edgar Sucar, and Luke Melas-Kyriazi of VGG for insightful discussions and Ruining Li, Eldar Insafutdinov, and Yash Bhalgat of VGG for their helpful feedback. We would also like to thank the authors of Zero-1-to-3 [43] and Objaverse-XL [17] for their helpful discussions.

Ethics. We use the Objaverse [18], OmniObject3D [81], and Google Scanned Object datasets (GSO) [20] following their terms and conditions. These datasets contain synthetic or scanned 3D objects, but, as far as we could determine, no personal data. For further details on ethics, data protection, and copyright please see <https://www.robots.ox.ac.uk/~vedaldi/research/union/ethics.html>.

References

- [1] Edward H Adelson and John YA Wang. Single lens stereo with a plenoptic camera. *IEEE transactions on pattern analysis and machine intelligence (TPAMI)*, 14(2):99–106, 1992. 2
- [2] Edward H Adelson, James R Bergen, et al. The plenoptic function and the elements of early vision. *Computational models of visual processing*, 1(2):3–20, 1991. 2
- [3] Titas Anciukevičius, 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 (CVPR)*, pages 12608–12618, 2023. 3
- [4] Yutong Bai, Xinyang Geng, Karttikeya Mangalam, Amir Bar, Alan Yuille, Trevor Darrell, Jitendra Malik, and Alexei A Efros. Sequential modeling enables scalable learning for large vision models. *arXiv preprint arXiv:2311.15127*, 2023. 3
- [5] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5855–5864, 2021. 1, 2
- [6] 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 (CVPR)*, pages 5470–5479, 2022. 1, 2
- [7] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023. 3
- [8] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 22563–22575, 2023. 2, 4
- [9] 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. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2023. 3
- [10] Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T Freeman. Maskgit: Masked generative image transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11315–11325, 2022. 3
- [11] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. In *European Conference on Computer Vision (ECCV)*, pages 333–350. Springer, 2022. 1, 2

- [12] Eric Ming Chen, Sidhanth Holalkere, Ruyu Yan, Kai Zhang, and Abe Davis. Ray conditioning: Trading photo-realism for photo-consistency in multi-view image generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023. 4
- [13] Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya Sutskever. Generative pre-training from pixels. In *International conference on machine learning (ICML)*, pages 1691–1703. PMLR, 2020. 3
- [14] Shenchang Eric Chen and Lance Williams. View interpolation for image synthesis. In *Proceedings of the 20th Annual Conference on Computer Graphics and Interactive Techniques (SIGGRAPH)*, page 279–288, New York, NY, USA, 1993. Association for Computing Machinery. 2
- [15] 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. 2
- [16] Paul E Debevec, Camillo J Taylor, and Jitendra Malik. Modeling and rendering architecture from photographs: A hybrid geometry-and image-based approach. In *Proceedings of the 23th Annual Conference on Computer Graphics and Interactive Techniques (SIGGRAPH)*, 1996. 2
- [17] Matt Deitke, Ruoshi Liu, Matthew Wallingford, Huong Ngo, Oscar Michel, Aditya Kusupati, Alan Fan, Christian Laforte, Vikram Voleti, Samir Yitzhak Gadre, Eli VanderBilt, Anirudha Kembhavi, Carl Vondrick, Georgia Gkioxari, Kiana Ehsani, Ludwig Schmidt, and Ali Farhadi. Objaverse-xl: A universe of 10m+ 3d objects. *Advances in Neural Information Processing Systems (NeurIPS)*, 2023. 2, 5, 6, 7, 9, 3, 4
- [18] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Anirudha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3d objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13142–13153, 2023. 1, 2, 3, 5, 6, 8, 9
- [19] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems (NeurIPS)*, 34:8780–8794, 2021. 1
- [20] Laura Downs, Anthony Francis, Nate Koenig, Brandon Kinman, Ryan Hickman, Krista Reymann, Thomas B McHugh, and Vincent Vanhoucke. Google scanned objects: A high-quality dataset of 3d scanned household items. In *2022 International Conference on Robotics and Automation (ICRA)*, pages 2553–2560. IEEE, 2022. 1, 2, 5, 7, 8, 9
- [21] Vincent Dumoulin, Jonathon Shlens, and Manjunath Kudlur. A learned representation for artistic style. In *International Conference on Learning Representations (ICLR)*, 2016. 4
- [22] Patrick Esser, Robin Rombach, Andreas Blattmann, and Bjorn Ommer. Imagebart: Bidirectional context with multinomial diffusion for autoregressive image synthesis. *Advances in neural information processing systems (NeurIPS)*, 34:3518–3532, 2021. 3
- [23] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (CVPR)*, pages 12873–12883, 2021. 3
- [24] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5501–5510, 2022. 2
- [25] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014. 1, 2
- [26] Yuwei Guo, Ceyuan Yang, Anyi Rao, Yaohui Wang, Yu Qiao, Dahua Lin, and Bo Dai. Animatediff: Animate your personalized text-to-image diffusion models without specific tuning. *arXiv preprint arXiv:2307.04725*, 2023. 2, 4
- [27] Richard Hartley and Andrew Zisserman. *Multiple view geometry in computer vision*. Cambridge university press, 2003. 5
- [28] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Proceedings of the 31st International Conference on Neural Information Processing Systems (NeurIPS)*, pages 6626–6637, 2017. 6
- [29] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems (NeurIPS)*, 33:6840–6851, 2020. 1
- [30] 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. 2, 4
- [31] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *Proceedings of the IEEE international conference on computer vision (ICCV)*, pages 1501–1510, 2017. 4
- [32] Zixuan Huang, Stefan Stojanov, Anh Thai, Varun Jampani, and James M Rehg. Planes vs. chairs: Category-guided 3d shape learning without any 3d cues. In *European Conference on Computer Vision*, pages 727–744. Springer, 2022. 1
- [33] Yifan Jiang, Hao Tang, Jen-Hao Rick Chang, Liangchen Song, Zhangyang Wang, and Liangliang Cao. Efficient-3dim: Learning a generalizable single-image novel-view synthesizer in one day. *arXiv preprint arXiv:2310.03015*, 2023. 4
- [34] Angjoo Kanazawa, Shubham Tulsiani, Alexei A Efros, and Jitendra Malik. Learning category-specific mesh reconstruction from image collections. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 371–386, 2018. 1
- [35] Yash Kant, Aliaksandr Siarohin, Michael Vasilkovsky, Riza Alp Guler, Jian Ren, Sergey Tulyakov, and Igor Gilitschenski. invs: Repurposing diffusion inpainters for novel view synthesis. *arXiv preprint arXiv:2310.16167*, 2023. 3, 6
- [36] 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 (CVPR)*, pages 18423–18433, 2023. 3
- [37] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019. 4
- [38] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics (TOG)*, 42(4):1–14, 2023. 1, 2
- [39] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2014. 1, 2
- [40] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*, 2019. 3
- [41] Minghua Liu, Ruoxi Shi, Linghao Chen, Zhuoyang Zhang, Chao Xu, Xinyue Wei, Hansheng Chen, Chong Zeng, Jiayuan Gu, and Hao Su. One-2-3-45++: Fast single image to 3d objects with consistent multi-view generation and 3d diffusion. *arXiv preprint arXiv:2311.07885*, 2023. 3
- [42] Minghua Liu, Chao Xu, Haian Jin, Linghao Chen, Zexiang Xu, Hao Su, et al. One-2-3-45: Any single image to 3d mesh in 45 seconds without per-shape optimization. *arXiv preprint arXiv:2306.16928*, 2023. 2, 3, 4
- [43] 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 (ICCV)*, pages 9298–9309, 2023. 2, 3, 4, 5, 6, 7, 8, 9, 1
- [44] Yuan Liu, Cheng Lin, Zijiao Zeng, Xiaoxiao Long, Lingjie Liu, Taku Komura, and Wenping Wang. Syncdreamer: Learning to generate multiview-consistent images from a single-view image. *arXiv preprint arXiv:2309.03453*, 2023. 2, 3, 4, 5, 6, 7, 8, 1
- [45] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural volumes: learning dynamic renderable volumes from images. *ACM Transactions on Graphics (TOG)*, 38(4):1–14, 2019. 4
- [46] 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 (CVPR)*, pages 8446–8455, 2023. 2
- [47] Luke Melas-Kyriazi, Christian Rupprecht, and Andrea Vedaldi. Pc2: Projection-conditioned point cloud diffusion for single-image 3d reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12923–12932, 2023. 1
- [48] B Mildenhall, PP Srinivasan, M Tancik, JT Barron, R Ramamoorthi, and R Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *Proceedings of the European conference on computer vision (ECCV)*, 2020. 1, 2, 4
- [49] 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. 1, 2
- [50] Eunbyung Park, Jimei Yang, Ersin Yumer, Duygu Ceylan, and Alexander C Berg. Transformation-grounded image generation network for novel 3d view synthesis. In *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, pages 3500–3509, 2017. 1, 2
- [51] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (CVPR)*, pages 165–174, 2019. 2
- [52] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. In *The Eleventh International Conference on Learning Representations (ICLR)*, 2023. 2
- [53] Senthil Purushwalkam and Nikhil Naik. Conrad: Image constrained radiance fields for 3d generation from a single image. *Advances in Neural Information Processing Systems (NeurIPS)*, 2023. 2
- [54] Guocheng Qian, Jinjie Mai, Abdullah Hamdi, Jian Ren, Aliaksandr Siarohin, Bing Li, Hsin-Ying Lee, Ivan Skokhodov, Peter Wonka, Sergey Tulyakov, et al. Magic123: One image to high-quality 3d object generation using both 2d and 3d diffusion priors. *arXiv preprint arXiv:2306.17843*, 2023. 2
- [55] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning (ICML)*, pages 8748–8763. PMLR, 2021. 5
- [56] Ali Razavi, Aaron Van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images with vq-vae-2. *Advances in neural information processing systems (NeurIPS)*, 32, 2019. 3
- [57] 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, 2022. 2, 8, 1
- [58] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems (NeurIPS)*, 35:36479–36494, 2022. 2
- [59] Mehdi SM Sajjadi, Daniel Duckworth, Aravindh Mahendran, Sjoerd van Steenkiste, Filip Pavetic, Mario Lucic, Leonidas J Guibas, Klaus Greff, and Thomas Kipf. Object scene representation transformer. *Advances in Neural Information Processing Systems (NeurIPS)*, 35:9512–9524, 2022. 2
- [60] 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 (CVPR)*, pages 6229–6238, 2022. 2
- [61] Kyle Sargent, Zizhang Li, Tanmay Shah, Charles Herrmann, Hong-Xing Yu, Yunzhi Zhang, Eric Ryan Chan, Dmitry Lagun, Li Fei-Fei, Deqing Sun, et al. Zeronvs: Zero-shot 360-degree view synthesis from a single real image. *arXiv preprint arXiv:2310.17994*, 2023. 2
- [62] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems (NeurIPS)*, 35:25278–25294, 2022. 2, 3, 1
- [63] Yichun Shi, Peng Wang, Jianglong Ye, Mai Long, Kejie Li, and Xiao Yang. Mvdream: Multi-view diffusion for 3d generation. *arXiv preprint arXiv:2308.16512*, 2023. 2, 3
- [64] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, et al. Make-a-video: Text-to-video generation without text-video data. In *The Eleventh International Conference on Learning Representations (ICLR)*, 2022. 2, 4
- [65] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. *Advances in Neural Information Processing Systems (NeurIPS)*, 32, 2019. 2
- [66] Vincent Sitzmann, Semon Rezchikov, Bill Freeman, Josh Tenenbaum, and Fredo Durand. Light field networks: Neural scene representations with single-evaluation rendering. *Advances in Neural Information Processing Systems (NeurIPS)*, 34:19313–19325, 2021. 2, 4, 1
- [67] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning (ICML)*, pages 2256–2265. PMLR, 2015. 1, 2
- [68] Mohammed Suhail, Carlos Esteves, Leonid Sigal, and Ameesh Makadia. Light field neural rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8269–8279, 2022. 2
- [69] Stanislaw Szymanowicz, Christian Rupprecht, and Andrea Vedaldi. Viewset diffusion: (0-)image-conditioned 3D generative models from 2D data. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2023. 1, 3
- [70] Junshu Tang, Tengfei Wang, Bo Zhang, Ting Zhang, Ran Yi, Lizhuang Ma, and Dong Chen. Make-it-3d: High-fidelity 3d creation from a single image with diffusion prior. *arXiv preprint arXiv:2303.14184*, 2023. 2
- [71] Maxim Tatarchenko, Alexey Dosovitskiy, and Thomas Brox. Multi-view 3d models from single images with a convolutional network. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part VII 14*, pages 322–337. Springer, 2016. 2
- [72] Hung-Yu Tseng, Qinbo Li, Changil Kim, Suhub Alsisan, Jia-Bin Huang, and Johannes Kopf. Consistent view synthesis with pose-guided diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16773–16783, 2023. 5
- [73] Aaron Van den Oord, Nal Kalchbrenner, Lasse Espeholt, Oriol Vinyals, Alex Graves, et al. Conditional image generation with pixelcnn decoders. *Advances in neural information processing systems (NeurIPS)*, 29, 2016. 3
- [74] Aäron Van Den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural networks. In *International conference on machine learning (ICML)*, pages 1747–1756. PMLR, 2016. 3
- [75] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in neural information processing systems (NeurIPS)*, 30, 2017. 3
- [76] Daniel Watson, William Chan, Ricardo Martin Brullalla, Jonathan Ho, Andrea Tagliasacchi, and Mohammad Norouzi. Novel view synthesis with diffusion models. In *The Eleventh International Conference on Learning Representations (ICLR)*, 2023. 1, 2, 5
- [77] Haohan Weng, Tianyu Yang, Jianan Wang, Yu Li, Tong Zhang, CL Chen, and Lei Zhang. Consistent123: Improve consistency for one image to 3d object synthesis. *arXiv preprint arXiv:2310.08092*, 2023. 2, 3, 4, 5, 6, 7, 8
- [78] Olivia Wiles, Georgia Gkioxari, Richard Szeliski, and Justin Johnson. Synsin: End-to-end view synthesis from a single image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7467–7477, 2020. 2
- [79] Qiucheng Wu, Yujian Liu, Handong Zhao, Ajinkya Kale, Trung Bui, Tong Yu, Zhe Lin, Yang Zhang, and Shiyu Chang. Uncovering the disentanglement capability in text-to-image diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1900–1910, 2023. 4
- [80] Ruiqi Wu, Liangyu Chen, Tong Yang, Chunle Guo, Chongyi Li, and Xiangyu Zhang. Lamp: Learn a motion pattern for few-shot-based video generation. *arXiv preprint arXiv:2310.10769*, 2023. 2, 4
- [81] Tong Wu, Jiarui Zhang, Xiao Fu, Yuxin Wang, Jiawei Ren, Liang Pan, Wayne Wu, Lei Yang, Jiaqi Wang, Chen Qian, et al. Omniobject3d: Large-vocabulary 3d object dataset for realistic perception, reconstruction and generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 803–814, 2023. 1, 2, 5, 7, 9
- [82] Yifeng Xiong, Haoyu Ma, Shanlin Sun, Kun Han, and Xi-aohui Xie. Light field diffusion for single-view novel view synthesis. *arXiv preprint arXiv:2309.11525*, 2023. 3, 4
- [83] Jiayu Yang, Ziang Cheng, Yunfei Duan, Pan Ji, and Hongdong Li. Consistnet: Enforcing 3d consistency for multi-view images diffusion. *arXiv preprint arXiv:2310.10343*, 2023. 2, 3, 4, 5, 6

- [84] Jianglong Ye, Peng Wang, Kejie Li, Yichun Shi, and Heng Wang. Consistent-1-to-3: Consistent image to 3d view synthesis via geometry-aware diffusion models. In *Proceedings of the International Conference on 3D Vision (3DV)*, 2024. 6
- [85] 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 (CVPR)*, pages 4578–4587, 2021. 1
- [86] Lijun Yu, Yong Cheng, Kihyuk Sohn, José Lezama, Han Zhang, Huiwen Chang, Alexander G Hauptmann, Ming-Hsuan Yang, Yuan Hao, Irfan Essa, et al. Magvit: Masked generative video transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10459–10469, 2023. 3
- [87] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*, 2020. 2
- [88] 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 (CVPR)*, pages 586–595, 2018. 5, 6
- [89] Chuanxia Zheng, Tung-Long Vuong, Jianfei Cai, and Dinh Phung. Movq: Modulating quantized vectors for high-fidelity image generation. *Advances in Neural Information Processing Systems (NeurIPS)*, 35:23412–23425, 2022. 3, 4
- [90] Zhizhuo Zhou and Shubham Tulsiani. Sparsefusion: Distilling view-conditioned diffusion for 3d reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12588–12597, 2023. 3

Free3D: Consistent Novel View Synthesis without 3D Representation

Supplementary Material

The supplementary materials are organized as follows:

- A video to illuminate our work and the rendered videos.
- Introduction for the baseline diffusion model.
- Experiment details.
- Results for more single view NVS.

A. Background: Diffusion Generators

In order to achieve sufficient generalization to operate in an *open-set* category setting, Free3D builds on a pre-trained 2D image generation, and specifically Stable Diffusion (SD) [57]. SD is a Latent Diffusion Model (LDM) trained on billions of text-image pairs from LAION-5B [62]. It consists of two stages. The first stage embeds the given image $x_0 \in \mathbb{R}^{H \times W \times 3}$ in a latent space $z \in \mathbb{R}^{\frac{H}{f} \times \frac{W}{f} \times c}$ through an autoencoder $\mathcal{E} : x \mapsto z$, paired with a decoder $\mathcal{D} : z \mapsto x$, which reconstructs the image ($x = \mathcal{D} \circ \mathcal{E}(x)$). The second stage uses diffusion to model the distribution $p(z|y)$ over such latent codes, where y lumps any conditioning information (e.g., text, image, or viewpoint). Diffusion involves a forward noising process that gradually perturbs the given latent $z_0 = z$ by adding the Gaussian noise ϵ in a Markovian fashion:

$$z_t = \sqrt{\bar{\alpha}_t} z_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, \quad \epsilon \sim \mathcal{N}(0, I), \quad (\text{A.1})$$

producing a sequence $z_t, t = 1, \dots, T$, and $\bar{\alpha}_t := \prod_{s=1}^t \alpha_s$, $\alpha_t := 1 - \beta_t$ denote the noise strength at different steps. $\{\beta_t\}_{t=1}^T$ is a pre-defined variance schedule. Ultimately, $p(z_T|y)$ is approximately normal; we can thus easily sample z_T , and then go back to z_0 via the backward denoising process using the predicted noise:

$$z_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(z_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(z_t, t, y) \right) + \sigma_t \epsilon, \quad (\text{A.2})$$

where ϵ_θ is typically an UNet [19], and $\{\sigma_t\}_{t=1}^T$ is another control of noise ϵ , which is also a pre-defined schedule corresponding to the schedule β_t and introduces uncertainty for the synthesis of different views. Similar to the vanilla DDPM [29], SD uses the following training objective to optimize the UNet ϵ_θ :

$$\mathcal{L} = \mathbb{E}_{z_0, y, \epsilon \sim \mathcal{N}(0, I), t} [\|\epsilon - \epsilon_\theta(z_t, t, y)\|_2^2], \quad (\text{A.3})$$

B. Experiment Details

The Stable Diffusion (SD), originally trained for text-to-images generation, requires adaptation to suit image-conditional NVS tasks. Following Zero-1-to-3 [43], we utilize the image-to-image Stable Diffusion checkpoints³. Our

baseline code is built upon the Zero-1-to-3 [43]⁴. Hyper-parameters are configured in accordance with the default settings of the baseline code. The *ray conditioning normalisation* (RCN) is incorporated into each ResNet block within the diffusion Unet ϵ_θ , while the *pseudo-3D cross-attention* is introduced after the original CLIP-conditional cross-attention layer (as illustrated in Fig. 2).

Instead of directly providing $r_{uv} = (\mathbf{o} \times \mathbf{d}_{uv}, \mathbf{d}_{uv})$ to the network for modulating the features, we embed them into higher-dimensional features, following the approach of NeRF [48] and LFN [66]. In particular, we employ the element-wise mapping $\mathbf{r} \mapsto [\mathbf{r}, \sin(f_1 \pi \mathbf{r}), \cos(f_1 \pi \mathbf{r}), \dots, \sin(f_K \pi \mathbf{r}), \cos(f_K \pi \mathbf{r})]$, where K is the number of Fourier bands, and f_k is equally spaced to the sampling rate. In all experiments, K is set as 6, leading to $78 = 2 \times 3 \times K_o + 3 + 2 \times 3 \times K_d + 3 = (6 + 6) \times 6 + 6$ dimensional features (as depicted in Fig. 2(a)).

Training Details. Our model was trained on $4 \times$ A40 48GB GPUs in two stages: **i)** We first finetuned the model with RCN, utilizing a batch size of 256 for 3 days on random camera viewpoints, enhancing the pose accuracy for target views. **ii)** Subsequently, the pseudo-3D cross-attention was finetuned on the 4 nearest views, employing a batch size of 192 for 2 days. In the second stage, different views from one instance were perturbed by adding noise from the same time step t .

In an alternative approach during the first stage, we initially attempted to jointly train the pseudo-3D cross-attention with random camera viewpoints. However, the performance is worse than the configuration \mathbb{D} . We believe this is because the camera viewpoints have a large gap along these random views in the rendered datasets, making it harder to calculate the similarity across these frames. In all experiments, we use AdamW with a learning rate of 10^{-5} for the old parameters in the original diffusion Unet ϵ_θ and a $10 \times$ larger learning rate for new parameters, namely the parameters for RCN and Pseudo-3D cross-attention.

Inference Details. At the testing phase, we configure the diffusion model with a sampling step set to $T = 50$. The computational time for rendering a novel view using our proposed Free3D is approximately 3 seconds, utilizing an A6000 GPU. For a fair comparison, all models are evaluated on the same A6000 GPU employing the same batch size of 4. This batch size is chosen due to the operational constraints of syndreamer [44], which can only run such a

³<https://huggingface.co/spaces/lambdalabs/stable-diffusion-image-variations>

⁴<https://github.com/cvlab-columbia/zero123>

small size. Additionally, we also utilize the CFG with a scale $s = 3$ to guide the rendering for each target view.

360° Video Rendering. To render a 360° video, we establish a circle trajectory by uniformly subdividing the azimuth ϕ into discrete intervals of $\frac{2\pi}{50} = 7.2^\circ$, while the elevation θ and the distance z remain fixed. For each 3D instance, we replicate the same latent variable z_T over 50 frames, which can minimize temporal flickering across different views. Additionally, we also set the parameter σ_t in Eq. (A.2) to zero, thereby further mitigating uncertainty introduced by varying noise patterns.

C. More Visual Results

More results on Objaverse NVS. In Figs. C.1 and C.2, we present more visual comparisons on Objaverse datasets [18] that given one input image and the target viewpoint, all models render the target novel view. This is an extension of Fig. 4 in the main paper.

Here, all examples shown come from the corresponding test-set following the split, as in Zero-1-to-3 [43]. These examples are good evidence that our Free3D is suitable for *open-set* categories NVS, where it can generate semantically reasonable content with visually realistic appearances across various categories. More importantly, compared to existing state-of-the-art methods, the Free3D provides better results with a more precise pose for the target novel view. This observation suggests that the RCN is able to provide better viewpoint perception for the NVS.

More results on OmniObject3D and GSO NVS. In Figs. C.3 and C.4, we show additional comparison results on OmniObject3D [81] and GSO [20] datasets, respectively. This is an extension of Fig. 5 in the main paper, which demonstrates the generalizability of our Free3D on unseen datasets encompassing various categories.

As can be seen from these results, although the baseline Zero-1-to-3 [43] provides visually realistic appearances for all objects, the content is *not* always reasonable, and the pose is inaccurate in many cases. This indicates the global language token embedding with elevation θ , azimuth ϕ , and distance z is *not* so precise for the network to interpret and utilize the camera viewpoints. While the Zero123-XL [17] and consistent123 [77] enhance the quality by training on a larger dataset and employing multi-view diffusion, respectively, they do *not* directly deal with the camera pose perception. In contrast, our Free3D leverages the *per-pixel* ray conditioning as well as the modulating, which significantly improves the pose perception accuracy.



(a) Input View (b) Target View (c) Zero123[43] (d) Zero123-XL[17] (e) SyncDreamer[44] (f) Consistent123[77] (g) Ours Free3D

Figure C.1. **Qualitative comparisons on Objaverse dataset.** Given the exact target pose, the proposed Free3D significantly improves the pose precision compared to existing state-of-the-art methods.

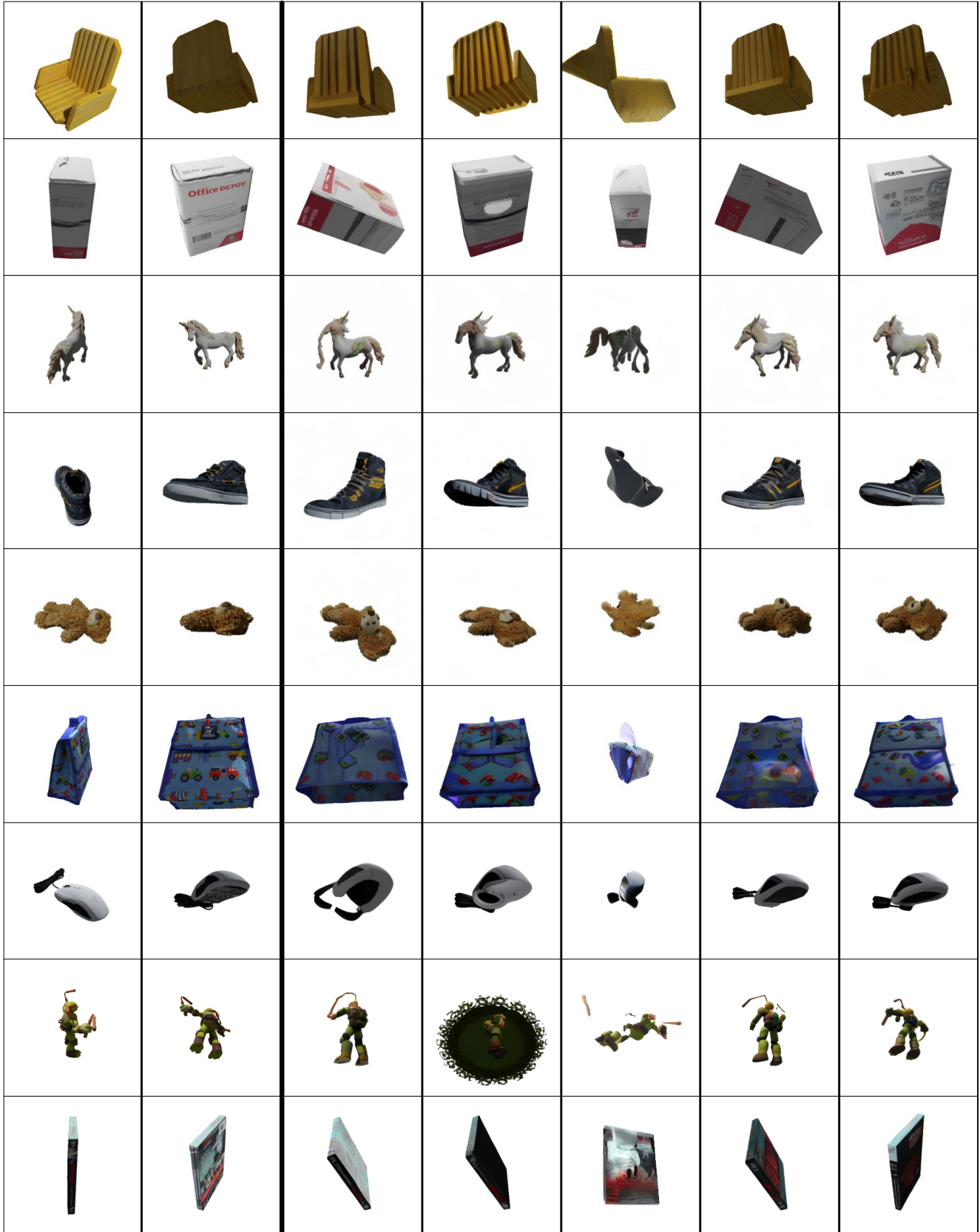


(a) Input View (b) Target View (c) Zero123[43] (d) Zero123-XL[17] (e) SyncDreamer[44] (f) Consistent123[77] (g) Ours Free3D

Figure C.2. **Qualitative comparisons on Objaverse dataset.** Given the exact target pose, the proposed Free3D significantly improves the pose precision compared to existing state-of-the-art methods.



Figure C.3. **Qualitative comparisons on OmniObject3D dataset.** Given the exact target pose, the proposed Free3D significantly improves the pose precision compared to existing state-of-the-art methods.



(a) Input View (b) Target View (c) Zero123[43] (d) Zero123-XL[17] (e) SyncDreamer[44] (f) Consistent123[77] (g) Ours Free3D

Figure C.4. **Qualitative comparisons on GSO dataset.** Given the exact target pose, the proposed Free3D significantly improves the pose precision compared to existing state-of-the-art methods.