# Sparse3D: Distilling Multiview-Consistent Diffusion for Object Reconstruction from Sparse Views

**Zi-Xin Zou[1], Weihao Cheng[2], Yan-Pei Cao[2], Shi-Sheng Huang[3],**
**Ying Shan[2], Song-Hai Zhang[1]**

[1]Tsinghua University  [2]ARC Lab, Tencent PCG  [3]Beijing Normal University

## Abstract

Reconstructing 3D objects from extremely sparse views is a long-standing and challenging problem. While recent techniques employ image diffusion models for generating plausible images at novel viewpoints or for distilling pre-trained diffusion priors into 3D representations using score distillation sampling (SDS), these methods often struggle to simultaneously achieve high-quality, consistent, and detailed results for both novel-view synthesis (NVS) and geometry. In this work, we present *Sparse3D*, a novel 3D reconstruction method tailored for sparse view inputs. Our approach distills robust priors from a multiview-consistent diffusion model to refine a neural radiance field. Specifically, we employ a controller that harnesses epipolar features from input views, guiding a pre-trained diffusion model, such as Stable Diffusion, to produce novel-view images that maintain 3D consistency with the input. By tapping into 2D priors from powerful image diffusion models, our integrated model consistently delivers high-quality results, even when faced with open-world objects. To address the blurriness introduced by conventional SDS, we introduce the category-score distillation sampling (C-SDS) to enhance detail. We conduct experiments on CO3DV2 which is a multi-view dataset of real-world objects. Both quantitative and qualitative evaluations demonstrate that our approach outperforms previous state-of-the-art works on the metrics regarding NVS and geometry reconstruction.

## Introduction

Reconstructing 3D objects from sparse-view images remains a pivotal challenge in the realms of computer graphics and computer vision. This technique has a wide range of applications such as Augmented and Virtual Reality (AR/VR). The advent of the Neural Radiance Field (NeRF) and its subsequent variants has catalyzed significant strides in geometry reconstruction and novel-view synthesis, as delineated in recent studies (Mildenhall et al. 2020; Wang et al. 2021a; Yariv et al. 2021). However, NeRFs exhibit limitations when operating on extremely sparse views, specifically with as few as 2 or 3 images. In these scenarios, the synthesized novel views often suffer in quality due to the limited input observations.

Existing methods for sparse-view reconstruction typically leverage a generalizable NeRF model, pre-trained on multi-view datasets, to infer 3D representations from projected im-
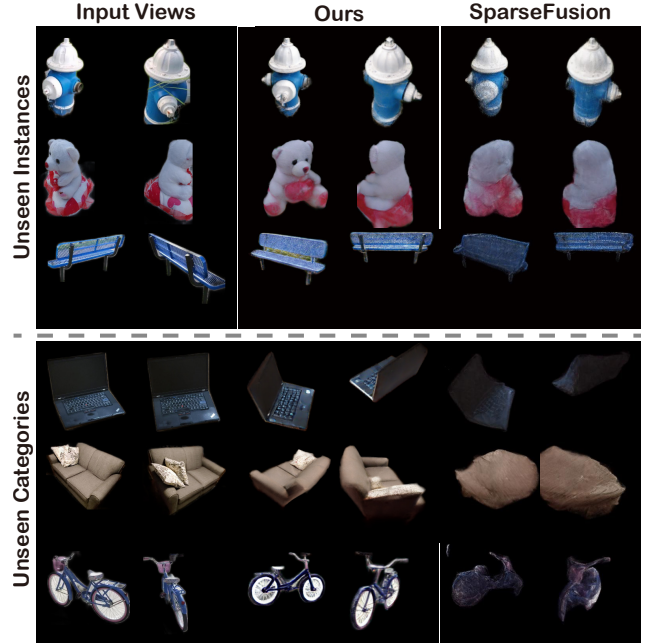


Figure 1: **Novel-view synthesis from two input views using our Sparse3D and SparseFusion.** Our approach can achieve higher-quality images with more details for unseen instances, especially for the unobserved regions of them (e.g., the left face of the teddybear). Furthermore, our approach can generalize to some unseen categories without any further finetuning, while SparseFusion fails.

age features (Yu et al. 2021; Chibane et al. 2021). However, these approaches tend to regress to the mean , failing to produce perceptually sharp outputs, especially in intricate details. To produce plausible results, either in terms of geometry or appearance, from limited observations, several studies have turned to image generation models, such as the diffusion model (Rombach et al. 2022), to "imagine" unseen views based on provided images (Chan et al. 2023; Zhou and Tulsiani 2023). For example, Zero123 (Liu et al. 2023) trains a view-conditioned diffusion model on a large synthetic dataset and achieves impressive results. However, their generated images across different views may not be consistent.

Thus, while these view-conditioned diffusion models can produce satisfactory images, their quality and generalization ability is often constrained by the scarcity of posed image datasets. Large-scale image diffusion models (Ramesh et al. 2021; Saharia et al. 2022; Rombach et al. 2022), which are pre-trained on billions of 2D images (Schuhmann et al. 2022), excel in generating high-quality and diverse images. However, despite the diverse, general capability of such models, in 3D reconstruction tasks, users need to synthesize specific instances that are coherent with user-provided input images. Even with recent model customization methods (Kumari et al. 2023; Ruiz et al. 2023; Gal et al. 2022), they prove unwieldy and often fail to produce the specific concept with sufficient fidelity. Consequently, the potential of merging the capabilities of pre-trained large image diffusion models with viewpoint and appearance perception of specific instances remains an open avenue of exploration.

In contrast to directly generating images at novel views, some recent works explore distilling the priors of pre-trained diffusion models into a NeRF (neural radiance field) framework. This approach facilitates 3D-consistent novel-view synthesis and allows for mesh extraction from the NeRF. Notable works such as DreamFusion (Poole et al. 2023) and SJC (Wang et al. 2023a) employ score distillation sampling (SDS) to harness off-the-shelf diffusion models for text-to-3D generation. However, a persistent challenge with SDS is the production of blurry and oversaturated outputs, attributed to noisy gradients, which in turn compromises the quality of NeRF reconstructions.

In this work, we present *Sparse3D*, a novel 3D reconstruction approach designed to reconstruct high-fidelity 3D objects from sparse and posed input views. Our method hinges on two pivotal components: **(1)** a diffusion model that ensures both multiview consistency and fidelity to user-provided input images, while retaining the powerful generalization capabilities of Stable Diffusion (Rombach et al. 2022), and **(2)** a category-score distillation sampling (C-SDS) strategy. At its core, we distill the priors from our fidelity-preserving, multiview-consistent diffusion model into the NeRF reconstruction using an enhanced category-score distillation sampling. Specifically, for the multiview-consistent diffusion model, we propose to utilize an epipolar controller to guide the off-the-shelf Stable Diffusion model to generate novel-view images which are 3D consistent with the content of input images. Notably, by fully harnessing the 2D priors present in Stable Diffusion, our model exhibits robust generalization capabilities, producing high-quality images even when confronted with open-world, unseen objects. To overcome the problem of blurry, oversaturated, and non-detailed results caused by SDS during NeRF reconstruction, we draw inspiration from VSD (Wang et al. 2023b) and propose a category-score distillation sampling strategy (C-SDS). Additionally, two perception losses between the one-step estimation image from our diffusion model and the rendering image are employed for better results, without incurring much extra computational cost.

We evaluate *Sparse3D* on the Common Object in 3D (CO3DV2) dataset and benchmark it against existing approaches. The results show that our approach outperforms state-of-the-art techniques in terms of the quality of both synthesized novel views and reconstructed geometry. Importantly, *Sparse3D* exhibits superior generalization capabilities, particularly for object categories not present in the training domain.

## Related Works

### Multi-view 3D Reconstruction

Multi-view 3D reconstruction is a long standing problem with impressive works such as traditional Structure-from-Motion (S*f*M) (Schönberger and Frahm 2016) or Multi-view-Stereo (MVS) (Schönberger et al. 2016), and recent learning based approaches (Yao et al. 2018; Yu and Gao 2020). The success of NeRF (Mildenhall et al. 2020; Müller et al. 2022) has led to impressive outcomes in novel-view synthesis and geometric reconstruction. However, these methods still struggle to produce satisfactory results for extremely sparse view scenario. Subsequent works proposed to use regularization (semantic (Jain, Tancik, and Abbeel 2021), frequency (Yang, Pavone, and Wang 2023), geometry and appearance (Niemeyer et al. 2022)) and geometric priors (e.g. depth (Deng et al. 2022; Roessle et al. 2022) or normal (Yu et al. 2022)) but still remain to be inadequate for view generation in unobserved regions, due to the essential lack of scene priors.

### Generalizable Novel-view Synthesis

For generalizable novel-view synthesis using NeRF, some approaches utilize projected features of the sampling points in volumetric rendering (Yu et al. 2021; Wang et al. 2021b; Chibane et al. 2021), or new neural scene representations, such as Light Field Network (Suhail et al. 2022b,a) or Scene Representation Transformer (Sajjadi et al. 2022) for better generalizable novel-view synthesis. Subsequent researches (Kulhánek et al. 2022; Chan et al. 2023; Yoo et al. 2023) propose to further utilize generative models (e.g. VQ-VAE (van den Oord, Vinyals, and Kavukcuoglu 2017) and diffusion model (Rombach et al. 2022)) to generate unseen images. However, these methods didn't have any 3D-aware scene priors, which limits their potential applications. In this paper, we leverage the feature map from generalizable renderer to guide a pre-trained diffusion model to generate multiview-consistent images, and then distill the diffusion prior into NeRF reconstruction for both novel-view synthesis and geometry reconstruction.

### 3D Generation with 2D Diffusion Model

Diffusion denoising probabilistic models has brought a boom of generation task for 2D images and 3D contents in recent years. Inspired by early works which use CLIP embedding (Jain, Tancik, and Abbeel 2021; Wang et al. 2022; Jain et al. 2022) or GAN (Pan et al. 2021) to regularize the NeRF, DreamFusion (Poole et al. 2023) and SJC (Wang et al. 2023a) propose a score distillation sampling (SDS) strategy to guide the NeRF optimization for impressive text-to-3D generation. ProlificDreamer (Wang et al. 2023b) proposes variational score distillation (VSD) for more high-fidelity and diverse text-to-3D generation. Magic3D (Lin

Figure 2: **Overview of Sparse3D**. Our approach consists of two key components: a multiview-consistent diffusion model and a category-score distillation sampling. We utilize epipolar feature map to control Stable Diffusion model to generate image consistent with the content of input images, serving as a multiview-consistent diffusion model. Based on such model, we propose a category-score distillation sampling (C-SDS) strategy to achieve more detailed results during NeRF reconstruction.

et al. 2023) improves the 3D generation quality by a two-stage coarse-to-fine strategy. To generate 3D results consistent with the input image observation, subsequent works leverage textual-inversion (Melas-Kyriazi et al. 2023) or denoised-CLIP loss with depth prior (Tang et al. 2023). When additional geometry prior are available (e.g. point-clouds from Point-E (Nichol et al. 2022)), some works (Seo et al. 2023; Yu et al. 2023) can produce more 3D consistent creation. In additional to lift a pre-trained diffusion model, Zero123 (Liu et al. 2023), SparseFusion (Zhou and Tulsiani 2023) and NerfDiff (Gu et al. 2023) train a viewpoint-conditioned diffusion model and achieve impressive results. Instead of training a diffusion model or directly lifting a pre-trained diffusion model, our approach leverages both the advantages of them to train a multiview-consistent diffusion model, with a category-score distillation sampling to improve the results of SDS for more details.

## Method

Given $N$ input images $\{I_n\}_{n=1}^N$ of an object with corresponding camera poses $\{T_n\}_{n=1}^N$, where $N$ can be as few as 2, our goal is to reconstruct a neural radiance field (NeRF), enabling generalizable novel view synthesis and high-quality surface reconstruction. To realize this goal, we propose Sparse3D, which distills a multiview-consistent diffusion model prior into NeRF representation of an object, using a category-score distillation sampling (C-SDS) strategy. Figure 2 shows the overview of our approach. The multiview-consistent diffusion model extracts epipolar features from sparse input views and uses a control network to guide the Stable Diffusion model generate novel-view images which are faithful to the object shown in the images. A NeRF is then reconstructed with the guidance of the diffusion model. To overcome the blurry problem that occurred in SDS, we propose C-SDS. Benefiting from C-SDS, the gradients conditioned on category prior maintain the optimization with a tightened region of the search space, leading to

more detailed results. Finally, our approach achieves more consistent and high-quality results of novel-view synthesis and geometry reconstruction. We introduce the details of the multiview-consistent diffusion model and the C-SDS based NeRF reconstruction in the following subsections.

## Multiview-Consistent Diffusion Model

Our diffusion model consists of a feature renderer, an epipolar controller and a Stable Diffusion model, where the epipolar controller and the Stable Diffusion model together constitute the noise predictor $\epsilon_\beta$, as shown in Figure 3. The feature renderer $g_\psi$ takes a set of posed images and viewpoint $\pi$ as input, subsequently outputting an epipolar feature map $f_c = g_\psi(\pi, I_1, ..., I_n, T_1, ..., T_n)$, which serves as the input for the epipolar controller. In order to unify pre-trained diffusion model and multiview-consistent perception ability for a specific object, we draw inspiration from ControlNet (Zhang and Agrawala 2023). ControlNet enables image generation controlled by conditional inputs (such as edge maps, segmentation maps, and depth maps). Instead, we use the epipolar feature map to guide a pre-trained diffusion model to generate images consistent with the content of input images from various viewpoints. To align the feature space of feature renderer and controller, we use a convolution layer to map the features before feeding them into the controller.

**Feature Renderer.** Previous works acquire the feature map $f_c$ through rendering from Triplane (Gu et al. 2023), 3D Volume (Chan et al. 2023) or epipolar feature transformer (Zhou and Tulsiani 2023). In this paper, we adapt epipolar feature transformer (EFT) following (Zhou and Tulsiani 2023). The EFT, derived from GPNR (Suhail et al. 2022a), learns a network $g_\psi$ to predict color of given ray $r$ from input images. The rendering process primarily involves three transformers, which output attention weights used to blend colors over input views and epipolar lines for the final prediction. We recommend readers to refer to (Zhou and Tulsiani 2023) and (Suhail et al. 2022a) for the details of
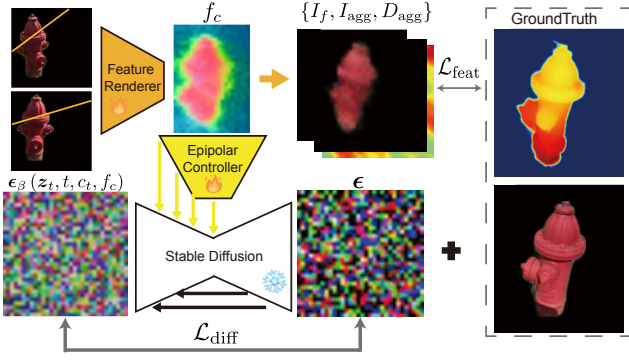
Figure 3: **Multiview-consistent diffusion model.** Our multiview-consistent diffusion model comprises a feature renderer, an epipolar controller and a Stable Diffusion model.

epipolar feature rendering. We implement two modifications to the EFT for improved results: (1) a mask embedding and a relative camera transformation embedding are concatenated with other transformer token features. (2) To enhance generalizability and achieve better geometry-awareness, we also obtain the aggregated color $I_{agg}$ and depth images $D_{agg}$ by attention weights of transformers to compute loss.

**Epipolar Controller.** Given feature maps $f_c$ rendered at arbitrary viewpoints, we propose to learn an epipolar controller to guide a pre-trained diffusion model generate multiview-consistent images with high quality. Our epipolar controller takes epipolar feature map $f_c$ and category text prompt $c_t$ as input, subsequently outputting the latent features that are fused with the latent features of Stable Diffusion. We also employ a convolution layer to align the dimension of feature map and epipolar controller input. Rather than training a new diffusion model, we hope to retain the rich 2D priors from Stable Diffusion. Consequently, we jointly train our epipolar controller and feature renderer, while keeping the parameters of Stable Diffusion fixed. On the one hand, by utilizing the feature map, which contains implicit information about the appearance of the specific object and perception of observation viewpoint, we can control a pre-trained text-to-image diffusion model to generate image consistent with the content of input images from different viewpoints. On the other hand, our diffusion model inherits the high-quality image generation capabilities from Stable Diffusion, and the additional category prior in text domain can also enhance the multiview consistency. Furthermore, these priors also enable our model to generalize to open-world unseen categories.

**Training.** Finally, we jointly train the feature renderer and the epipolar controller by the following objective function:

$$\mathcal{L} = \mathcal{L}_{feat} + \mathcal{L}_{diff} \qquad (1)$$

where $\mathcal{L}_{feat}$ is the loss for feature renderer and $\mathcal{L}_{diff}$ is the loss for epipolar controller.

While the feature map primarily serves as input for the controller in our pipeline, we also supervise it by color images and depth images to enhance its perception of ap-

pearance, observation viewpoints and geometry-awareness. For a query ray $r$ from novel view when given input images, we decode the color $I_f$ from feature map and supervise it using ground-truth color values. Additionally, to improve generalizability and geometry-awareness, we employ a Mean Squared Error (MSE) loss on aggregated color $I_{agg}$ and depth $D_{agg}$. We then formulate the objective function as follows:

$$\mathcal{L}_{feat} = \sum_r ||I_f(r) - I(r)||^2 + ||I_{agg}(r) - I(r)||^2 \\ + ||D_{agg}(r) - D(r)|| \qquad (2)$$

where $I(r)$ and $D(r)$ are ground-truth color and depth image respectively.

The diffusion model learns a conditional noise predictor to estimate denoising score by adding Guassian-noise $\epsilon$ to clean data in $T$ timesteps. We minimize the noise prediction error at randomly sampled timestep $t$. The objective of diffusion model conditioned on text prompt $c_t$ (we use the category name as the conditioned text prompt, e.g. "hydrant") and feature map $f_c$ is given by:

$$\mathcal{L}_{diff} = \mathbb{E}_{\epsilon \sim \mathcal{N}(0,1)} ||\epsilon - \epsilon_\beta(z_t, t, c_t, f_c)||^2 \qquad (3)$$

where $\epsilon_\beta$ is the conditional noise predictor of our diffusion model.

## NeRF Reconstruction with C-SDS

Building on our multiview-consistent diffusion model, we aim to optimize a neural radiance field (NeRF) parameterized with $\theta$, from which more 3D-consistent novel-view synthesis and underlying explicit geometry can be derived. Then to overcome the problem of blurry and non-detailed results in SDS, we propose a category-score distillation sampling (C-SDS) strategy.

**Category-Score Distillation Sampling.** To overcome the problem of SDS, we draw inspiration from VSD (Wang et al. 2023b) and propose a C-SDS for more detailed outcomes as follows:

$$\nabla_\theta \mathcal{L}_{C-SDS}(\theta) \approx \mathbb{E}_{t,\epsilon} \left[ \omega(t) (\epsilon_{mc} - \epsilon_{cat}) \frac{\partial z_t}{\partial x} \frac{\partial x}{\partial \theta} \right] \qquad (4)$$

where $\epsilon_{mc} = \epsilon_\beta(z_t, t, c_t, f_c)$ is the predicted noise by our multiview-consistent diffusion model, $\epsilon_{cat} = \epsilon_{sd}(z_t; t, c_t)$ is the predicted noise by Stable Diffusion conditioned text prompt of category $c_t$. And $\omega(t)$ is a weighting function that depends on the timestep $t$.

Instead of employing a Gaussian noise as SDS does, we replace it with an estimation $\epsilon_{cat}$ incorporating category prior from Stable Diffusion. By providing an approximation of the estimation of the score function of the distribution on rendering images with category prior, our C-SDS can deliver a better gradient with a tightened region of the search space, resulting in more detailed outputs. SDS relies on a high classifier-free guidance (CFG, i.e. 100) to achieve a better convergence, but such high CFG may lead to over-saturation and over-smooth problems (Poole et al. 2023). One reason for the reliance on large CFG is that the pre-trained text-to-image diffusion model may not obtain multiview-consistent
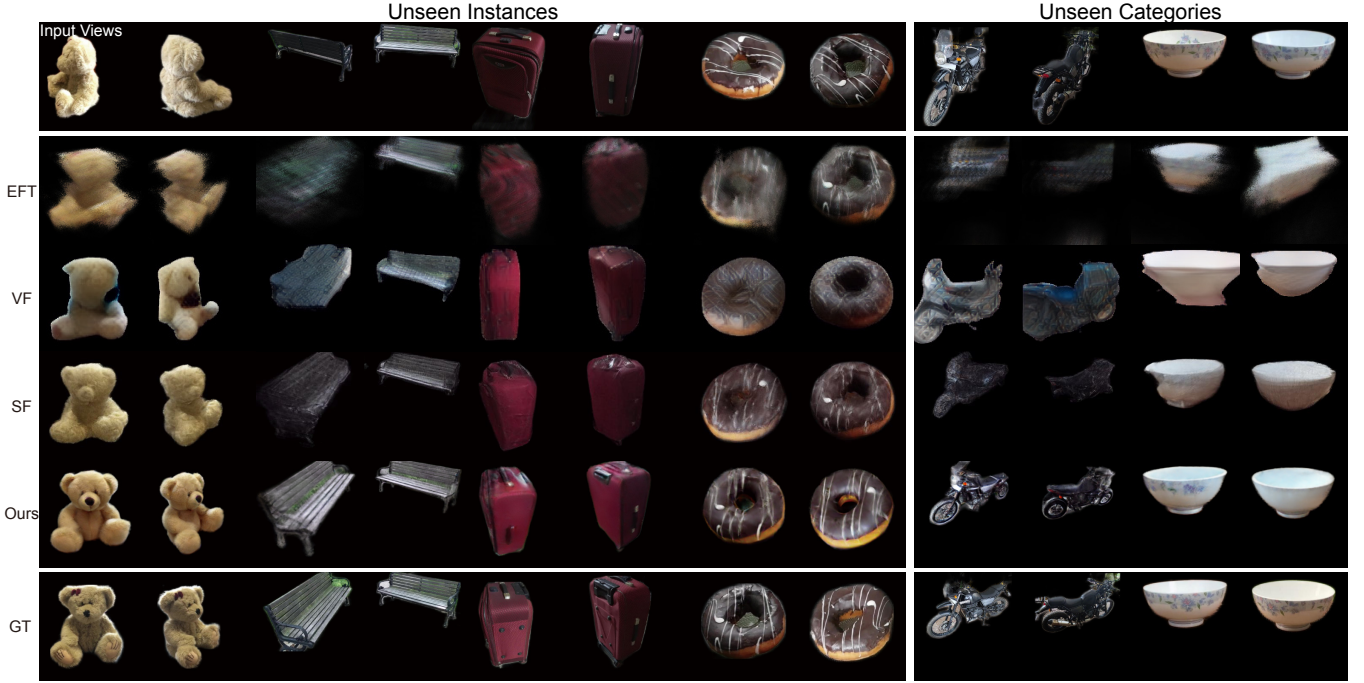
Figure 4: **Qualitative comparison of novel-view synthesis when given 2 input views.** Our approach achieves both high quality and more details of novel-view images compared to the others (e.g., the face of teddybear), whenever with unseen instances and unseen categories.

image generation at novel views, thus providing a noisy estimation across different viewpoints. In our experiment, when using a more multiview-consistent diffusion model, it can work with a small CFG (i.e. 7.5). However, the results still suffer from blurry and non-detailed outputs, as the update gradient is not accurate enough. ProlificDreamer utilizes a low-rank adaption (LoRA) of a pre-trained diffusion model to estimate the score function of the distribution on rendered images. In our experiment, we find that it is hard for LoRA to provide good estimation during our instance-specific optimization. Therefore, our proposed C-SDS offers a simple yet effective way to estimate the score function of the distribution on rendered images for more detailed results.

**One-step Estimation from Diffusion Model.** The predicted noise from diffusion model can be used not only in C-SDS, but also to estimate its one-step denoising image without requiring much extra computation:

$$\boldsymbol{z}_{1\text{step}} = \frac{1}{\sqrt{\bar{\alpha}_t}} \left( \boldsymbol{z}_t - \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}_\beta \left( \boldsymbol{z}_t, t, c_t, f_c \right) \right),$$
$$\boldsymbol{x}_{1\text{step}} = \mathcal{D}(\boldsymbol{z}_{1\text{step}}) \tag{5}$$

where $\mathcal{D}$ is the decoder of Stable Diffusion which decodes latent features to image space. Although its one-step estimation may be blurry and sometimes inaccurate, making it unsuitable for performing pixel-wise loss, we can leverage it to provide an additional regularization term by using perceptual distance. We find that the perception regularization from one-step estimation improves the metrics of results. Specifically, we employ two perceptual losses, which includes LPIPS loss (Zhang et al. 2018) and contextual loss

(Mechrez, Talmi, and Zelnik-Manor 2018) to formulate the perception regularization from one-step estimation image:

$$\mathcal{L}_{\text{perp}} = \lambda_p \mathcal{L}_{\text{lpips}}(I, \boldsymbol{x}_{1step}) + \lambda_c \mathcal{L}_{\text{contextual}}(I, \boldsymbol{x}_{1step}) \tag{6}$$

**Reference Supervision.** In addition to the guidance of diffusion priors at novel views, we use the reference input images $I$ with their masks $M$ to encourage the consistent appearance with the input images:

$$\mathcal{L}_{\text{ref}} = \lambda_r ||(\hat{I} - I) * \hat{M}||_2^2 + \lambda_m ||\hat{M} - M||_2^2 \tag{7}$$

where $\hat{I}$ and $\hat{M}$ are rendering image and mask, respectively.

**Overall Training.** We combine all of the losses, including $\mathcal{L}_{\text{C-SDS}}, \mathcal{L}_{\text{perp}}, \mathcal{L}_{\text{ref}}$, to formulate the objective function of NeRF reconstruction for specific object. Once NeRF reconstruction is complete, we can perform volume rendering for novel-view synthesis, and the underlying mesh can be extracted using Marching Cubes (Lorensen and Cline 1987).

## Experiment

In this section, we conduct qualitative and quantitative evaluation of our approach on 3D object dataset, CO3Dv2 dataset (Reizenstein et al. 2021), to demonstrate its effectiveness. CO3Dv2 dataset is a real-world dataset, which contains 51 common object categories encountered in daily life. We first show the superior quality of ours on novel-view synthesis and 3D reconstruction for unseen object instances in category-specific scenario with varying number of input, and then show its out-of-domain generalization ability for unseen categories.

| | Unseen Instances - 2 views | | | | | |
|---|---|---|---|---|---|---|
| | PSNR ↑ | SSIM ↑ | LPIPS ↓ | FID ↓ | CLIP ↑ | DISTS ↓ |
| PixelNeRF | 15.33 | 0.29 | 0.59 | 371.23 | 0.83 | 0.44 |
| EFT | **21.28** | 0.69 | 0.34 | 293.36 | 0.87 | 0.33 |
| ViewFormer | 18.42 | 0.71 | 0.29 | 248.23 | 0.82 | 0.29 |
| SparseFusion | **21.28** | 0.76 | 0.23 | 187.22 | 0.91 | 0.26 |
| Ours | 20.95 | **0.77** | **0.22** | **147.65** | **0.93** | **0.23** |
| | Unseen Instance - 3 views | | | | | |
| | PSNR ↑ | SSIM ↑ | LPIPS ↓ | FID ↓ | CLIP ↑ | DISTS ↓ |
| PixelNeRF | 15.50 | 0.31 | 0.58 | 363.68 | 0.83 | 0.43 |
| EFT | **22.62** | 0.74 | 0.29 | 242.87 | 0.89 | 0.30 |
| ViewFormer | 18.91 | 0.72 | 0.28 | 240.21 | 0.87 | 0.29 |
| SparseFusion | 22.31 | 0.78 | 0.22 | 175.02 | 0.92 | 0.24 |
| Ours | 22.06 | **0.79** | **0.20** | **134.22** | **0.94** | **0.21** |
| | Unseen Instances - 6 views | | | | | |
| | PSNR ↑ | SSIM ↑ | LPIPS ↓ | FID ↓ | CLIP ↑ | DISTS ↓ |
| PixelNeRF | 15.65 | 0.33 | 0.55 | 344.58 | 0.85 | 0.42 |
| EFT | **24.47** | 0.80 | 0.23 | 161.78 | 0.93 | 0.25 |
| ViewFormer | 19.77 | 0.74 | 0.27 | 232.30 | 0.89 | 0.28 |
| SparseFusion | 23.69 | 0.80 | 0.20 | 154.20 | 0.93 | 0.22 |
| Ours | 23.92 | **0.82** | **0.18** | **116.10** | **0.95** | **0.19** |
| | Unseen Categories - 2 views | | | | | |
| | PSNR ↑ | SSIM ↑ | LPIPS ↓ | FID ↓ | CLIP ↑ | DISTS ↓ |
| PixelNeRF | 14.82 | 0.31 | 0.50 | 314.45 | 0.81 | 0.44 |
| EFT | **19.31** | 0.56 | 0.41 | 318.64 | 0.87 | 0.38 |
| ViewFormer | 15.43 | 0.63 | 0.34 | 301.19 | 0.85 | 0.36 |
| SparseFusion | 18.83 | 0.70 | 0.28 | 290.45 | 0.88 | 0.34 |
| Ours | 18.83 | **0.72** | **0.23** | **164.30** | **0.93** | **0.26** |

Table 1: **Quantitative comparisons of novel-view synthesis.** We evaluate methods on unseen instances with varying numbers of input images, such as 2, 3, 6 (the top three blocks), and on unseen categories with 2 inputs views (the bottom block). We report the average results across categories for each blocks.



Figure 5: **Geometry reconstruction using SparseFusion and Ours.** The last column shows that ground-truth point cloud.

| | Unseen Instances | | | | | | Unseen Categories | |
|---|---|---|---|---|---|---|---|---|
| Views Num. | 2 | | 3 | | 6 | | 2 | |
| | CD ↓ | F-score ↑ | CD ↓ | F-score ↑ | CD ↓ | F-score ↑ | CD ↓ | F-score ↑ |
| SparseFusion | 0.27 | 0.23 | 0.26 | 0.23 | 0.24 | 0.25 | 0.37 | 0.18 |
| Ours | **0.21** | **0.32** | **0.19** | **0.38** | **0.16** | **0.45** | **0.27** | **0.28** |

Table 2: **Quantitative comparison of geometry reconstruction.** Since other baselines only produce image at novel views without 3D representation, we only report the results of ours and SparseFusion.

**Implementation details.** For feature renderer, we follow SparseFusion (Zhou and Tulsiani 2023) to use three groups of transformer encoders with four 256-dimensional layers to aggregate epipolar features. For multiview-consistent model, we adopt the Stable Diffusion model v1.5 as our priors. For NeRF reconstruction, we adapt the threestudio (Guo et al. 2023), which is an unified framework for 3D content creation from various input, to implement the NeRF reconstruction for specific object. We set the weights of the losses with $\lambda_p = 100$, $\lambda_c = 10$, $\lambda_r = 1000$ and $\lambda_m = 50$. NeRF optimization runs for 10,000 steps, which takes about an hour on an A100 GPU.

## Experimental Settings

**Dataset.** We follow the *fewview-train* and *fewview-dev* splits provided by CO3Dv2 dataset (Reizenstein et al. 2021) for training and evaluation purposes, respectively. For the evaluation of unseen object instances within the same categories, we use the core-subset with 10 categories to train category-specific diffusion model for each category. To assess the out-of-domain generalization ability on unseen categories, we select 10 categories for evaluation and use the remaining 41 categories together for training. Due to the hour-long of computation time required for our method, we evaluate only the first 10 object instances of each test split.

**Baselines.** We compare our approach with previous state-of-the-art baselines, including PixelNeRF (Yu et al. 2021), ViewFormer (Kulhánek et al. 2022), EFT and SparseFusion (Zhou and Tulsiani 2023). PixelNeRF and EFT ar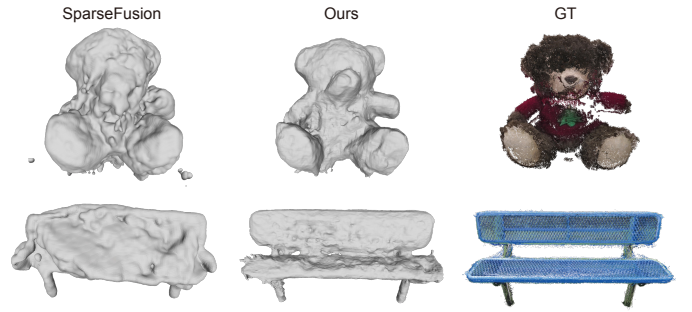e regression-based methods that deduce image at novel view by projection feature, where EFT is a adapted from GPNR for sparse views settings by (Zhou and Tulsiani 2023). ViewFormer is a generative model which employs a VQ-VAE codebook and a transformer module for image generation. Unlike the other methods that directly obtain novel-view synthesis with a single feed-forward pass, SparseFusion is a most relevant baseline to our approach, as it distills the diffusion model prior into NeRF reconstruction.

**Metrics.** We adopt several popular image quality assessment (IQA) to evaluate the quality of novel-view synthesis, including PSNR, SSIM, LPIPS (Zhang et al. 2018), FID (Heusel et al. 2017) and DISTS (Ding et al. 2022). Additionally, since our method can generate plausible results for unobserved regions, the evaluation between them and GT images may not be fair. Thus, we also adopt CLIP embedding similarity (Radford et al. 2021) of generated images with input images. Additionally, we evaluate the most commonly used 3D reconstruction quality metrics, including Chamfer Distance and F-score.

## Qualitative and Quantitative Evaluation

**Unseen Instances: 2 Views.** We first evaluate our approach with extremely sparse views (i.e. 2 views) for unseen object instances within the same categories. Table 1 demonstrates the quantitative comparison of ours and other baselines, with metrics averaged across 10 categories. We can observe that our method outperforms the others on most metrics of image quality except PSNR. Although ours has a slightly lower PSNR compared to the others, due to its formulation of pixel-wise MSE which favors mean color rendering results, our approach outperforms all of the others in perception metrics (e.g. LPIPS, FID, etc.). As the qualitative results shown in Figure 4, benefiting from two proposed key components,

Figure 6: **Effect of Stable Diffusion priors.** (a) diffusion model from SparseFusion; (b) ours diffusion model with Stable Diffusion priors.
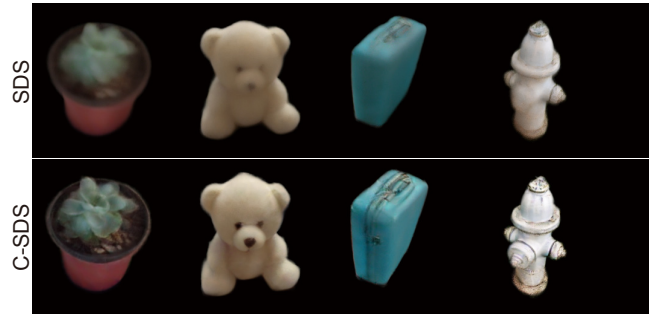


Figure 7: **Effect of C-SDS to the quality of NVS from NeRF reconstruction.** We can find the results of SDS are blurry and non-detailed in unobserved regions, while ours can generate more details with the same diffusion model.

our approach achieves both high-quality and more detailed results with 3D consistency. In addition to novel-view synthesis, we evaluate the quality of geometry reconstruction by extracting underlying mesh from NeRF. It should be noted that we only compare ours with SparseFusion, while the others (PixelNeRF, EFT and ViewFormer) are without 3D representation. From Table 2, we can find that our approach significantly outperforms SparseFusion with a wide margin, which demonstrates the ours method's superiority. Figure 5 also shows the mesh extracted from NeRF, where our results achieve sharper geometry with more details.

**Unseen Instances: Varying Views.** It's obviously that as the number of input views increases, the results of novel-view synthesis and geometry reconstruction improve. Table 1 and Table 2 also show the comparison of novel-view synthesis and geometry reconstruction on 3 and 6 input views, which demonstrates that our approach consistently outperforms the others with varying input views. More detailed evaluation results for each category and more qualitative results of novel-view synthesis and explicit geometry can be found in supplementary materials.

**Unseen Categories.** We conduct an experiment to evaluate the generalization ability to unseen categories between ours and the other baselines. Table 1 and Table 2 show the quantitative results of novel-view synthesis and geometry reconstruction. When confronted with the unseen categories that are out of the training domain, the performance of the other methods have a significant drop, while ours still maintains good performance, achieving the best results among them. The priors from Stable Diffusion enable our diffusion model to faithfully generate image of unseen categories. The last two columns of Figure 4 shows the novel-view synthesis of these methods. Our approach still can achieve high-quality images with more details, while the others are blurry and somewhat meaningless. More evaluation on unseen categories (e.g. with varying views) can be found in supplementary materials.

### Ablation Studies

**Stable Diffusion Priors.** To evaluate the effect of Stable Diffusion priors, we compare ours and SparseFusion in directly generating novel view images without performing

NeRF reconstruction, as shown in Figure 6. In unseen instances scenario, the diffusion model of SparseFusion can generate image at novel viewpoints consistent with the appearance of input images in a certain way (e.g. the blue hydrant with white head), but fails to achieve high-quality image generation. When the feature map is not reliable in some views, SparseFusion also fails to generate a multiview-consistent image (e.g. the bench). However, our diffusion model can achieve higher-quality of image generation that is more multiview-consistent regarding input images. In unseen categories scenario, the diffusion model of SparseFusion fails to generate meaningful images, while our method can be generalized to the these objects, benefiting from the Stable Diffusion priors (the last two columns in Figure 6).

**C-SDS.** We also investigate the effect of our distillation strategy on the quality of NeRF reconstruction. We also implement a version of using SDS, denoted as Ours$^\dagger$. When using our multiview-consistent diffusion model, which can provide a more accurate gradient update direction, thus there is no need of a large CFG, but it's still not enough for detailed results. In our experiment with setting the CFG value of SDS as 7.5, it can achieve plausible results with successfully convergence, but the blur problem is still unsolved, as shown the first row of Figure 7. When applying our proposed C-SDS with the same CFG, it's evident that the results show more details, which demonstrates the effectiveness of method. The quantitative and more results of ablation studies can be found in supplementary materials.

### Conclusion

In this paper, we introduce Sparse3D, a new approach to reconstruct high-quality 3D objects from sparse input views with camera poses. We utilize an epipolar controller to guide a pre-trained diffusion model to generate high-quality images which are 3D consistent with the content of input images, leading to a multiview-consistent diffusion model. Then, we distills the diffusion priors into NeRF optimization in a better way by using a category-score distillation sampling (C-SDS) strategy, resulting more detailed results. Experiments demonstrate that our approach can achieve the state-of-the-art results with higher-quality and more details, even when confronted with open-world, unseen objects.

# References

Chan, E. R.; Nagano, K.; Chan, M. A.; Bergman, A. W.; Park, J. J.; Levy, A.; Aittala, M.; Mello, S. D.; Karras, T.; and Wetzstein, G. 2023. Generative Novel View Synthesis with 3D-Aware Diffusion Models. *CoRR*, abs/2304.02602.

Chibane, J.; Bansal, A.; Lazova, V.; and Pons-Moll, G. 2021. Stereo Radiance Fields (SRF): Learning View Synthesis from Sparse Views of Novel Scenes. In *IEEE (CVPR)*.

Deng, K.; Liu, A.; Zhu, J.; and Ramanan, D. 2022. Depth-supervised NeRF: Fewer Views and Faster Training for Free. In *IEEE CVPR*, 12872–12881.

Ding, K.; Ma, K.; Wang, S.; and Simoncelli, E. P. 2022. Image Quality Assessment: Unifying Structure and Texture Similarity. *IEEE TPAMI.*, 44(5): 2567–2581.

Gal, R.; Alaluf, Y.; Atzmon, Y.; Patashnik, O.; Bermano, A. H.; Chechik, G.; and Cohen-Or, D. 2022. An image is worth one word: Personalizing text-to-image generation using textual inversion. *arXiv preprint arXiv:2208.01618*.

Gu, J.; Trevithick, A.; Lin, K.; Susskind, J. M.; Theobalt, C.; Liu, L.; and Ramamoorthi, R. 2023. NerfDiff: Single-image View Synthesis with NeRF-guided Distillation from 3D-aware Diffusion. *CoRR*, abs/2302.10109.

Guo, Y.-C.; Liu, Y.-T.; Shao, R.; Laforte, C.; Voleti, V.; Luo, G.; Chen, C.-H.; Zou, Z.-X.; Wang, C.; Cao, Y.-P.; and Zhang, S.-H. 2023. threestudio: A unified framework for 3D content generation. https://github.com/threestudio-project/threestudio.

Heusel, M.; Ramsauer, H.; Unterthiner, T.; Nessler, B.; and Hochreiter, S. 2017. GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium. In *NeurIPS*, 6626–6637.

Jain, A.; Mildenhall, B.; Barron, J. T.; Abbeel, P.; and Poole, B. 2022. Zero-Shot Text-Guided Object Generation with Dream Fields. In *IEEE CVPR*, 857–866.

Jain, A.; Tancik, M.; and Abbeel, P. 2021. Putting NeRF on a Diet: Semantically Consistent Few-Shot View Synthesis. In *ICCV*, 5865–5874.

Kulhánek, J.; Derner, E.; Sattler, T.; and Babuska, R. 2022. ViewFormer: NeRF-Free Neural Rendering from Few Images Using Transformers. In *ECCV*, volume 13675, 198–216.

Kumari, N.; Zhang, B.; Zhang, R.; Shechtman, E.; and Zhu, J.-Y. 2023. Multi-concept customization of text-to-image diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 1931–1941.

Lin, C.-H.; Gao, J.; Tang, L.; Takikawa, T.; Zeng, X.; Huang, X.; Kreis, K.; Fidler, S.; Liu, M.-Y.; and Lin, T.-Y. 2023. Magic3D: High-Resolution Text-to-3D Content Creation. In *IEEE CVPR*.

Liu, R.; Wu, R.; Hoorick, B. V.; Tokmakov, P.; Zakharov, S.; and Vondrick, C. 2023. Zero-1-to-3: Zero-shot One Image to 3D Object. arXiv:2303.11328.

Lorensen, W. E.; and Cline, H. E. 1987. Marching cubes: A high resolution 3D surface construction algorithm. In Stone, M. C., ed., *SIGGRAPH*, 163–169.

Mechrez, R.; Talmi, I.; and Zelnik-Manor, L. 2018. The Contextual Loss for Image Transformation with Non-aligned Data. In *ECCV*, volume 11218, 800–815.

Melas-Kyriazi, L.; Rupprecht, C.; Laina, I.; and Vedaldi, A. 2023. RealFusion: 360 Reconstruction of Any Object from a Single Image. In *IEEE CVPR*.

Mildenhall, B.; Srinivasan, P. P.; Tancik, M.; Barron, J. T.; Ramamoorthi, R.; and Ng, R. 2020. NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis. In *ECCV*, 405–421.

Müller, T.; Evans, A.; Schied, C.; and Keller, A. 2022. Instant neural graphics primitives with a multiresolution hash encoding. *ACM TOG*, 41(4): 102:1–102:15.

Nichol, A.; Jun, H.; Dhariwal, P.; Mishkin, P.; and Chen, M. 2022. Point-E: A System for Generating 3D Point Clouds from Complex Prompts. abs/2212.08751.

Niemeyer, M.; Barron, J. T.; Mildenhall, B.; Sajjadi, M. S. M.; Geiger, A.; and Radwan, N. 2022. RegNeRF: Regularizing Neural Radiance Fields for View Synthesis from Sparse Inputs. In *IEEE CVPR*, 5470–5480.

Pan, X.; Dai, B.; Liu, Z.; Loy, C. C.; and Luo, P. 2021. Do 2D GANs Know 3D Shape? Unsupervised 3D Shape Reconstruction from 2D Image GANs. In *ICLR*.

Poole, B.; Jain, A.; Barron, J. T.; and Mildenhall, B. 2023. DreamFusion: Text-to-3D using 2D Diffusion. In *ICLR*.

Radford, A.; Kim, J. W.; Hallacy, C.; Ramesh, A.; Goh, G.; Agarwal, S.; Sastry, G.; Askell, A.; Mishkin, P.; Clark, J.; Krueger, G.; and Sutskever, I. 2021. Learning Transferable Visual Models From Natural Language Supervision. In *ICML*, volume 139, 8748–8763.

Ramesh, A.; Pavlov, M.; Goh, G.; Gray, S.; Voss, C.; Radford, A.; Chen, M.; and Sutskever, I. 2021. Zero-Shot Text-to-Image Generation. In *ICML*, volume 139, 8821–8831.

Reizenstein, J.; Shapovalov, R.; Henzler, P.; Sbordone, L.; Labatut, P.; and Novotny, D. 2021. Common Objects in 3D: Large-Scale Learning and Evaluation of Real-life 3D Category Reconstruction. In *ICCV*.

Roessle, B.; Barron, J. T.; Mildenhall, B.; Srinivasan, P. P.; and Nießner, M. 2022. Dense Depth Priors for Neural Radiance Fields from Sparse Input Views. In *IEEE CVPR*, 12882–12891.

Rombach, R.; Blattmann, A.; Lorenz, D.; Esser, P.; and Ommer, B. 2022. High-Resolution Image Synthesis with Latent Diffusion Models. In *IEEE CVPR*, 10674–10685.

Ruiz, N.; Li, Y.; Jampani, V.; Pritch, Y.; Rubinstein, M.; and Aberman, K. 2023. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 22500–22510.

Saharia, C.; Chan, W.; Saxena, S.; Li, L.; Whang, J.; Denton, E. L.; Ghasemipour, S. K. S.; Lopes, R. G.; Ayan, B. K.; Salimans, T.; Ho, J.; Fleet, D. J.; and Norouzi, M. 2022. Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding. In *NeurIPS*.

Sajjadi, M. S. M.; Meyer, H.; Pot, E.; Bergmann, U.; Greff, K.; Radwan, N.; Vora, S.; Lucic, M.; Duckworth, D.; Dosovitskiy, A.; Uszkoreit, J.; Funkhouser, T. A.; and Tagliasacchi, A. 2022. Scene Representation Transformer: Geometry-Free Novel View Synthesis Through Set-Latent Scene Representations. In *IEEE CVPR*, 6219–6228.

Schönberger, J. L.; and Frahm, J. 2016. Structure-from-Motion Revisited. In *IEEE CVPR*, 4104–4113.

Schönberger, J. L.; Zheng, E.; Frahm, J.; and Pollefeys, M. 2016. Pixelwise View Selection for Unstructured Multi-View Stereo. In *ECCV*, volume 9907, 501–518.

Schuhmann, C.; Beaumont, R.; Vencu, R.; Gordon, C.; Wightman, R.; Cherti, M.; Coombes, T.; Katta, A.; Mullis, C.; Wortsman, M.; Schramowski, P.; Kundurthy, S.; Crowson, K.; Schmidt, L.; Kaczmarczyk, R.; and Jitsev, J. 2022. LAION-5B: An open large-scale dataset for training next generation image-text models. In *NeurIPS*.

Seo, J.; Jang, W.; Kwak, M.; Ko, J.; Kim, H.; Kim, J.; Kim, J.; Lee, J.; and Kim, S. 2023. Let 2D Diffusion Model Know 3D-Consistency for Robust Text-to-3D Generation. abs/2303.07937.

Suhail, M.; Esteves, C.; Sigal, L.; and Makadia, A. 2022a. Generalizable Patch-Based Neural Rendering. In *ECCV*, volume 13692, 156–174.

Suhail, M.; Esteves, C.; Sigal, L.; and Makadia, A. 2022b. Light Field Neural Rendering. In *IEEE CVPR*, 8259–8269.

Tang, J.; Wang, T.; Zhang, B.; Zhang, T.; Yi, R.; Ma, L.; and Chen, D. 2023. Make-It-3D: High-Fidelity 3D Creation from A Single Image with Diffusion Prior. *arXiv preprint arXiv:2303.14184*.

van den Oord, A.; Vinyals, O.; and Kavukcuoglu, K. 2017. Neural Discrete Representation Learning. In *NeurIPS*, 6306–6315.

Wang, C.; Chai, M.; He, M.; Chen, D.; and Liao, J. 2022. CLIP-NeRF: Text-and-Image Driven Manipulation of Neural Radiance Fields. In *IEEE CVPR*, 3825–3834.

Wang, H.; Du, X.; Li, J.; Yeh, R. A.; and Shakhnarovich, G. 2023a. Score Jacobian Chaining: Lifting Pretrained 2D Diffusion Models for 3D Generation. *IEEE CVPR*.

Wang, P.; Liu, L.; Liu, Y.; Theobalt, C.; Komura, T.; and Wang, W. 2021a. NeuS: Learning Neural Implicit Surfaces by Volume Rendering for Multi-view Reconstruction. In *NeurIPS*, 27171–27183.

Wang, Q.; Wang, Z.; Genova, K.; Srinivasan, P. P.; Zhou, H.; Barron, J. T.; Martin-Brualla, R.; Snavely, N.; and Funkhouser, T. A. 2021b. IBRNet: Learning Multi-View Image-Based Rendering. In *IEEE CVPR*, 4690–4699.

Wang, Z.; Lu, C.; Wang, Y.; Bao, F.; Li, C.; Su, H.; and Zhu, J. 2023b. ProlificDreamer: High-Fidelity and Diverse Text-to-3D Generation with Variational Score Distillation. *CoRR*, abs/2305.16213.

Yang, J.; Pavone, M.; and Wang, Y. 2023. FreeNeRF: Improving Few-shot Neural Rendering with Free Frequency Regularization. In *IEEE CVPR*.

Yao, Y.; Luo, Z.; Li, S.; Fang, T.; and Quan, L. 2018. MVS-Net: Depth Inference for Unstructured Multi-view Stereo. In *ECCV*, 785–801.

Yariv, L.; Gu, J.; Kasten, Y.; and Lipman, Y. 2021. Volume Rendering of Neural Implicit Surfaces. In *NeurIPS*.

Yoo, P.; Guo, J.; Matsuo, Y.; and Gu, S. S. 2023. DreamSparse: Escaping from Plato's Cave with 2D Frozen Diffusion Model Given Sparse Views. *CoRR*, abs/2306.03414.

Yu, A.; Ye, V.; Tancik, M.; and Kanazawa, A. 2021. pixelNeRF: Neural Radiance Fields From One or Few Images. In *IEEE CVPR*, 4578–4587.

Yu, C.; Zhou, Q.; Li, J.; Zhang, Z.; Wang, Z.; and Wang, F. 2023. Points-to-3D: Bridging the Gap between Sparse Points and Shape-Controllable Text-to-3D Generation. abs/2307.13908.

Yu, Z.; and Gao, S. 2020. Fast-MVSNet: Sparse-to-Dense Multi-View Stereo With Learned Propagation and Gauss-Newton Refinement. In *IEEE CVPR*, 1946–1955.

Yu, Z.; Peng, S.; Niemeyer, M.; Sattler, T.; and Geiger, A. 2022. MonoSDF: Exploring Monocular Geometric Cues for Neural Implicit Surface Reconstruction. In *NeurIPS*.

Zhang, L.; and Agrawala, M. 2023. Adding Conditional Control to Text-to-Image Diffusion Models. arXiv:2302.05543.

Zhang, R.; Isola, P.; Efros, A. A.; Shechtman, E.; and Wang, O. 2018. The Unreasonable Effectiveness of Deep Features as a Perceptual Metric. In *IEEE CVPR*, 586–595.

Zhou, Z.; and Tulsiani, S. 2023. SparseFusion: Distilling View-conditioned Diffusion for 3D Reconstruction. In *CVPR*.

# Supplementary Materials:
# Sparse3D: Distilling Multiview-Consistent Diffusion for Object Reconstruction from Sparse Views

## Implementation Details

### Feature Renderer

Since we adapt the Epipolar Feature Transformer (EFT) from SparseFusion (Zhou and Tulsiani 2023) to feature renderer, we firstly provide a review of EFT and introduce the difference between them.

**Review of EFT.** EFT, which is derived from GPNR, is a feed-forward network that aggregates the features along the epipolar lines of input images and then aggregates the features of aggregated epipolar features from different input views. Firstly, EFT employs a ResNet18 (He et al. 2016) as image feature extractor backbone to obtain the features by concatenating intermediate features from the first 4 layers. Then, for a ray $r$ casting from a query camera viewpoint $\pi$, the EFT uniformly samples $N$ points along the ray direction between the near $d_{near}$ and far $d_{far}$. The initial features $f_0$ of all sampled points can be concatenated by: (1) projected features tri-linear interpolated from input view images; (2) depths embedding; (3) Plücker coordinates embedding. Afterward, it employs three transformer modules $T_1, T_2, T_3$ to achieve aggregated features for the final feature map and color image calculation. The first transformer module is used to combine information from different views. Then an epipolar transformer and a view transformer are used to aggregate features along epipolar lines and different views to achieve feature map. The process of aggregation is calculated by attention weights $\alpha_k^m$ and $\beta_k$ of transformers and the final feature map $f_c$ can be calculated by:

$$f_c = \sum_{k=1}^{K} \beta_k \left( \sum_{m=1}^{M} \alpha_k^m f_k^m \right) \tag{1}$$

where $K$ is the number of input views, $M$ is the number of points sampled along the epipolar lines and $f_k^m$ is the combined features of $k$-th sampled point on $m$-th input image from the output of the first transformer. Then, a color image $I_f$ can be decoded by an additional linear layer.

**Additional Features Embedding.** In addition to the embedding used in EFT, we incorporate a mask embedding and relative camera transformation embedding. Since we focus on object reconstruction, a supplementary information indicating which sampled point's projection positions are in-side or out-side the object is beneficial for the transformer module to pay more attention to the in-side features. Further-

more, the input views of our approach are extremely sparse and there are almost no overlapping areas between two input images, (e.g., the two inputs show the front and the back of an object). In this case, it degrades to a single input situation, and information across different epipolar cannot be effectively combined. For example, the attention weights $\beta_k$ may have higher values for the other side of the input image, leading to worse view perception. Thus, a relative camera transformation can let the model learn to pay more attention to the nearer input view.

**Aggregated RGB and Depth.** SparseFusion (Zhou and Tulsiani 2023) trains EFT by a loss between a decoder color image $I_f$ and groundtruth $I$. For better generalization ability, we additionally adopt an aggregated color image $I_{agg}$ as GPNR does. And in order to improve the geometry-awareness of our feature renderer, we formulate an aggregated depth $D_{agg}$ with the similar to aggregated color as:

$$I_{agg} = \sum_{k=1}^{K} \beta_k \left( \sum_{m=1}^{M} \alpha_k^m I_k^m \right)$$

$$D_{agg} = \sum_{k=1}^{K} \beta_k \left( \sum_{m=1}^{M} \alpha_k^m d_k^m \right) \tag{2}$$

where $I_k^m$ is the rgb value of $k$-th sampled point projected on $m$-th input image, and $d_k^m$ is the sampled depth of this sampled point at rendering viewpoint. $\alpha_k^m$ and $\beta_k^m$ are the same as Equation . We supervise the $I_{agg}$ and $D_{agg}$ by the groundtruth color image and depth. These can make the transformer modules of feature render to output attention weights with better geometry awareness and generalization ability, which helps for extremely sparse views.

In our implementation, we follow (Zhou and Tulsiani 2023) to sample $N = 20$ points along each ray, and set $d_{near} = s - 5$ and $d_{far} = s + 5$ during training stage, where $s$ is the average distance from scene cameras to origin computed per scene. Finally, the feature map is rendered at the resolution of $32 \times 32$ with 256 dimensions for efficiency.

### Multiview-Consistent Diffusion model

We implement our diffusion model by using the Diffusers (von Platen et al. 2022), which is a go-to library for the state-of-the-art pre-trained diffusion models. We choose Stable Diffusion (Rombach et al. 2022) model as powerful
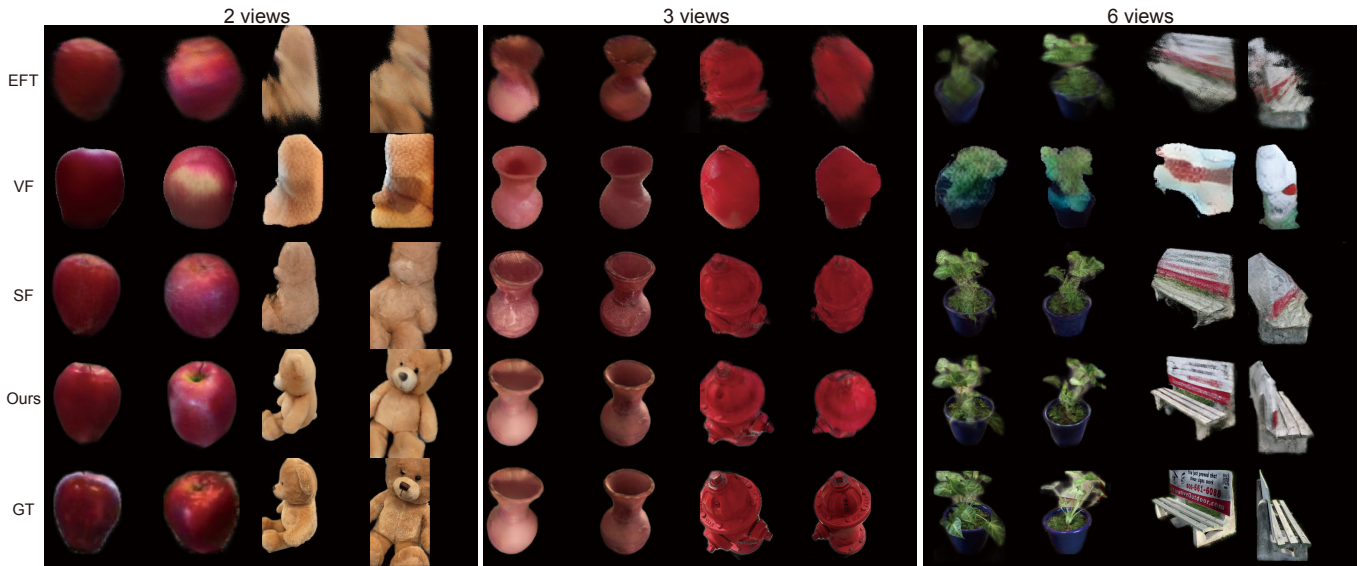
Figure 1: **Qualitative comparison of novel-view synthesis on unseen instances with a varying number of input views.**

well-studied 2D priors, which includes an encoder $\mathcal{E}(x)$, a decoder $\mathcal{D}(z)$ and an UNet $\mathcal{U}(z)$. The encoder and the decoder are employed to transfer between pixel space $x$ and latent space $z$, and the diffusion process with UNet is employed in latent space. Specifically, we adopt the network architecture and checkpoint weights from Stable-Diffusion-v1-5 (Rombach et al. 2022). The epipolar controller is initialized with the same architecture and weights as the encoder blocks and mid-blocks of stable diffusion UNet. And we employ a convolution layer to align the dimension of the feature map and epipolar controller input.

### NeRF Reconstruction

**Rendering from NeRF and Epipolar Features.** During NeRF reconstruction, we employ two renderers for the purpose of rendering images from NeRF representation and rendering epipolar features from input images, respectively. When sampling a camera viewpoint at novel viewpoints, we both render an image by the NeRF renderer and an epipolar feature map by the feature renderer. Our multiview-consistent diffusion model takes the rendered feature map as input to guide the NeRF representation through the back-propagation of differentiable volume rendering. Different from the training stage of the feature renderer, we employ the same near and far values of the feature renderer with NeRF renderer, which are calculated by the intersections between occupancy grids and ray.

**Scene Representation and Rendering.** For faster rendering and optimization, we utilize the Instance-NGP (Müller et al. 2022) as position encoding, a light-weight MLP with one hidden layer can output density $\sigma$ and color $c$. We also implement a progressive coarse-to-fine training strategy similar to (Li et al. 2023b). During rendering from NeRF representation, we also employ an occupancy grid transmittance estimator during optimization to skip the empty spaces, which can reduce the cost of memory for higher res-

olution. The NeRF representation is initialized as a Gaussian sphere for better convergence. During volume rendering from NeRF, we sample $N = 512$ points along the ray implemented by nerfacc (Li et al. 2023a) for acceleration.

**Additional Regularization.** In addition to the losses introduced in our paper, we employ three geometry regularization terms on the NeRF reconstruction, which are widely used in other works (Tang et al. 2023; Melas-Kyriazi et al. 2023), including orientation loss, entropy loss and sparsity loss. The orientation loss is proposed by (Verbin et al. 2022) which acts as a penalty on "foggy" surfaces, and the other two are following:

$$\mathcal{L}_{entropy} = w \cdot log_2(w) - (1-w) \cdot log_2(1-w)$$
$$\mathcal{L}_{sparsity} = \sqrt{w^2 + 0.01} \qquad (3)$$

where $w$ is the cumulative sum of the density.

## More Results of Individual Categories

### Novel-view Synthesis

**Unseen Instances.** To evaluate the performance on unseen object instances within the same categories, we conduct the experiment on 10 categories of CO3D dataset, including donut, apple, hydrant, vase, cake, ball, bench, suitcase, teddybear and plant. Specifically, we train the models for each categories using *fewview-train* split, and evaluate them on *fewview-test* split of each category with varying number of input images (2, 3 and 6), respectively. We compare our method with PixelNeRF (PN), ViewFormer (VF), EFT and SparseFusion (SF). Table 3 demonstrates the detailed quantitative results of novel-view synthesis of each category. Since the metrics of PSNR and SSIM have no perception ability, we pay more attention to the other perception metrics (LPIPS, FID, CLIP and DISTS), which are more suitable to our evaluation. SparseFusion achieves comparable results in most categories with simple geometry or appearance to
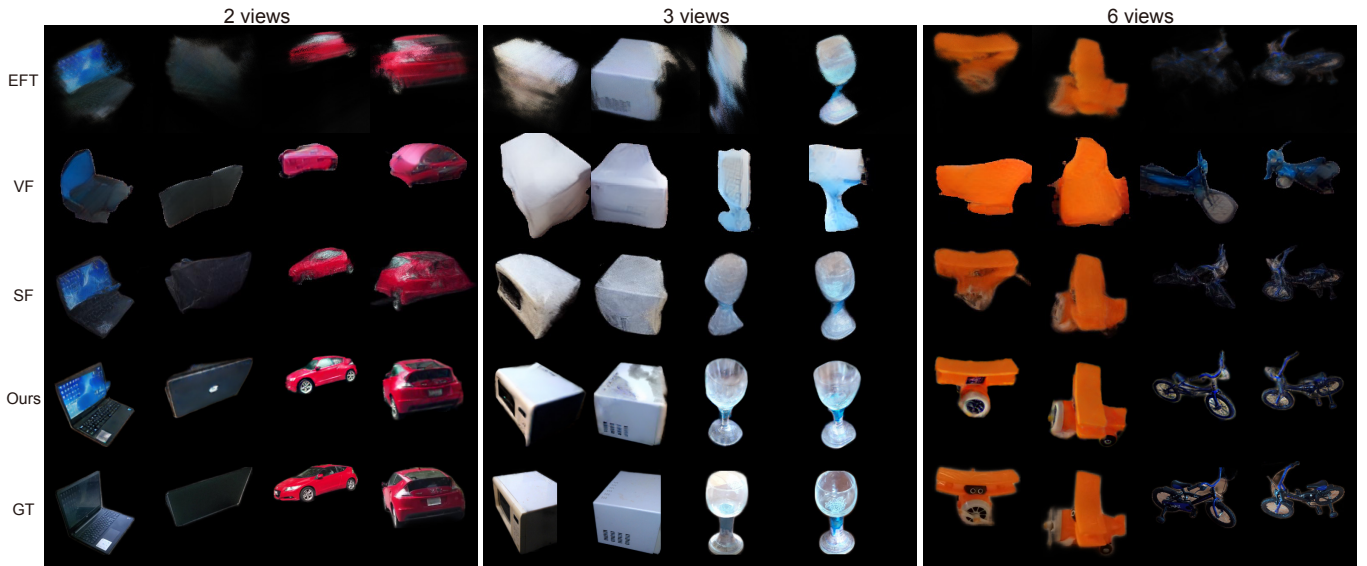
Figure 2: **Qualitative comparison of novel-view synthesis on unseen categories with a varying number of input views.**

ours, even though few of the results among them slightly surpass our method (e.g., donut and apple). However, it has poor performance on those with more complex geometry or appearance (e.g., bench), while our approach can produce higher-quality results on them. Along with the number of input views increasing, our approach outperforms the others in almost all categories (e.g., except ball with 6 views), consistently with varying numbers of input views. As Figure 1 illustrates some visual results of novel-view synthesis with varying numbers of input views, our approach achieves both high quality and more details than all of the other baselines.

**Unseen Categories.** To evaluate the generalization ability to unseen categories which are out of the training domain, we conduct the experiment on 10 categories, including bicycle, car, couch, laptop, microwave, motorcycle, bowl, toyplane, tv, wineglass. Specifically, we train only *one* model on the other 41 categories together, and evaluate them on *fewview-test* split of 10 categories, also with a varying number of input images (2, 3, and 6), respectively. We can find that when confronted with objects with unseen categories, our approach significantly outperforms the other methods across all test categories. This is mainly due to the prior from Stable Diffusion and generalizable feature renderer. The generalizable feature map can provide a view and appearance perception of the object, and the priors from stable diffusion with the text prompt of the category's name can successfully generate a plausible image regarding the specific input images, resulting in better generalization ability to unseen categories than the others. Figure 2 also illustrates some visual results of novel-view synthesis with varying numbers of input views. Ours still outperforms the others with higher quality, which demonstrates that our approach can be generalized to unseen categories very well, without further fine-tuning.

## Geometry Reconstruction

In addition to the evaluation of novel-view synthesis, we evaluate the performance of geometry reconstruction, since explicit geometry is also essential to many downstream applications. We only compare ours with SparseFusion for geometry reconstruction evaluation, due to the lack of 3D representation from the other methods. Table 5 shows the detailed quantitative results of geometry reconstruction of each category for both unseen instances and unseen categories, with varying numbers of input views. It demonstrates that our approach outperforms the SparseFusion with a wide margin on both Chamfer Distance and F-score, with all experiment settings. Figure 3 also shows the qualitative comparison between them, which demonstrates our approach can recover more detailed geometry. We find that SparseFusion can achieve approximate geometric shapes, where we can recognize it in some way, when evaluating on unseen instances. However, in unseen categories scenarios, the results of SparseFusion are much worse, and even some of them are hard to recognize the object from the shape. In contrast, our approach achieves much better reconstruction results with more details and sharper geometry. And even can make up the missing regions in ground-truth point cloud (the ball in the last row in Figure 3). Our approach is based on NeRF representation, which is not suitable for geometry reconstruction well. Thus, some ways to improve the quality of geometry reconstruction may be to utilize Signed Distance Function (SDF) field by NeuS (Wang et al. 2021), or transfer NeRF representation to tetrahedral SDF grid (Shen et al. 2021) with coarse-to-fine stage refinement, and we leave it as future work.

## Quantitative Results of Ablative Analysis

**Stable Diffusion Priors.** Table 1 shows the quantitative results of novel-view synthesis between the diffusion model of

Figure 3: **Qualitative comparison of geometry reconstruction with SparseFusion.**

| | Unseen Instances | | | | | |
| | PSNR ↑ | SSIM ↑ | LPIPS ↓ | FID ↓ | CLIP ↑ | DISTS ↓ |
|---|---|---|---|---|---|---|
| SF | **20.93** | 0.75 | 0.22 | 145.95 | **0.93** | **0.21** |
| Ours | 20.51 | **0.77** | **0.21** | **131.99** | **0.93** | 0.23 |
| | Unseen Categories | | | | | |
| | PSNR ↑ | SSIM ↑ | LPIPS ↓ | FID ↓ | CLIP ↑ | DISTS ↓ |
| SF | 18.17 | 0.68 | 0.29 | 257.19 | 0.90 | 0.27 |
| Ours | **18.78** | **0.72** | **0.24** | **138.29** | **0.93** | **0.25** |

Table 1: **Quantitative evaluation between the diffusion model of SparseFusion and Ours.**

| | PSNR ↑ | SSIM ↑ | LPIPS ↓ | FID ↓ | CLIP ↑ | DISTS ↓ | CD ↓ | F-score ↑ |
|---|---|---|---|---|---|---|---|---|
| SDS | **22.35** | **0.79** | 0.23 | 175.09 | 0.91 | 0.27 | 0.21 | 0.38 |
| C-SDS | 22.31 | **0.79** | **0.20** | **132.65** | **0.94** | **0.21** | **0.19** | **0.39** |

Table 2: **Quantitative evaluation between SDS and C-SDS on novel-view synthesis and geometry reconstruction.**



Figure 4: **Failure cases.** (a) extremely partial observation; (b) Janus problem.

SparseFusion and ours. We report the average of each metric across all categories and varying numbers of inputs. In unseen instances experiment, SparseFusion achieves comparable results to ours. However, in the unseen categories experiment, our approach significantly outperforms it by a large margin. The Stable Diffusion priors contain the features of objects from categories that may not have been present in our training domain, but are learned from large-scale images. This enables our diffusion model to generalize to unseen categories. Moreover, the Stable Diffusion priors also encompass the distribution of high-quality image generation and an additional category prior in the text domain. This

further assists our diffusion model in achieving both higher quality image generation and better multiview consistency.

**C-SDS.** Table 2 shows the quantitative results of novel-view synthesis and geometry reconstruction by using SDS and C-SDS, on unseen instances with varying numbers of inputs. In addition to the problem of blurry images shown in the figure in our paper, The blurring issue of SDS leads to poorer performance on most metrics for novel-view synthesis evaluation, with the exception of PSNR and SSIM.

When utilizing the same diffusion model, our C-SDS is able to achieve more details from diffusion model priors, resulting in better performance compared to SDS. Additionally, we have observed that our C-SDS has minimal impact on the geometry reconstruction quality, with SDS and C-SDS achieving comparable results in this aspect. Furthermore, we have found that the geometry reconstruction results achieved using SDS are significantly superior to those of SparseFusion (with 0.26 Chamfer Distance and 0.24 F-score). This further demonstrates that our multiview-consistent diffusion model has a heightened sense of geometry awareness.

## Failure Cases and Limitations

While our method has demonstrated promising results, there are still some limitations to its effectiveness. Figure 4 shows some failure cases of our approach, particularly when all input images are too close to the object and only contain a small portion of it, making it difficult for our approach to achieve satisfactory overall results (Figure 4 (a)). Additionally, the Janus problem still occasionally occurs in our results, as depicted in Figure 4(b) where the back of the car shows the appearance of the side. Furthermore, our approach heavily relies on accurate camera poses, which can be challenging to estimate directly from extremely sparse views, resulting in noisy estimates. As such, improving the method to handle noisy camera inputs is an interesting topic for future work.

## Supplementary Video

We provide a video attached with supplementary materials to show the 360-degree visualizations of the other baselines and ours. Please refer to it for more details.

## References

He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep Residual Learning for Image Recognition. In *IEEE CVPR*, 770–778.

Li, R.; Gao, H.; Tancik, M.; and Kanazawa, A. 2023a. NerfAcc: Efficient Sampling Accelerates NeRFs. *CoRR*, abs/2305.04966.

Li, Z.; Müller, T.; Evans, A.; Taylor, R. H.; Unberath, M.; Liu, M.-Y.; and Lin, C.-H. 2023b. Neuralangelo: High-Fidelity Neural Surface Reconstruction. In *IEEE CVPR*.

Melas-Kyriazi, L.; Rupprecht, C.; Laina, I.; and Vedaldi, A. 2023. RealFusion: 360 Reconstruction of Any Object from a Single Image. In *IEEE CVPR*.

Müller, T.; Evans, A.; Schied, C.; and Keller, A. 2022. Instant neural graphics primitives with a multiresolution hash encoding. *ACM TOG*, 41(4): 102:1–102:15.

Rombach, R.; Blattmann, A.; Lorenz, D.; Esser, P.; and Ommer, B. 2022. High-Resolution Image Synthesis with Latent Diffusion Models. In *IEEE CVPR*, 10674–10685.

Shen, T.; Gao, J.; Yin, K.; Liu, M.-Y.; and Fidler, S. 2021. Deep Marching Tetrahedra: a Hybrid Representation for High-Resolution 3D Shape Synthesis. In *NeurIPS*.

Tang, J.; Wang, T.; Zhang, B.; Zhang, T.; Yi, R.; Ma, L.; and Chen, D. 2023. Make-It-3D: High-Fidelity 3D Creation from A Single Image with Diffusion Prior. *arXiv preprint arXiv:2303.14184*.

Verbin, D.; Hedman, P.; Mildenhall, B.; Zickler, T. E.; Barron, J. T.; and Srinivasan, P. P. 2022. Ref-NeRF: Structured View-Dependent Appearance for Neural Radiance Fields. In *IEEE CVPR*, 5481–5490.

von Platen, P.; Patil, S.; Lozhkov, A.; Cuenca, P.; Lambert, N.; Rasul, K.; Davaadorj, M.; and Wolf, T. 2022. Diffusers: State-of-the-art diffusion models. https://github.com/huggingface/diffusers.

Wang, P.; Liu, L.; Liu, Y.; Theobalt, C.; Komura, T.; and Wang, W. 2021. NeuS: Learning Neural Implicit Surfaces by Volume Rendering for Multi-view Reconstruction. In *NeurIPS*, 27171–27183.

Zhou, Z.; and Tulsiani, S. 2023. SparseFusion: Distilling View-conditioned Diffusion for 3D Reconstruction. In *CVPR*.

Table 3 — **2 Views**

| | Donut LPIPS↓ | Donut FID↓ | Apple LPIPS↓ | Apple FID↓ | Hydrant LPIPS↓ | Hydrant FID↓ | Vase LPIPS↓ | Vase FID↓ | Cake LPIPS↓ | Cake FID↓ | Ball LPIPS↓ | Ball FID↓ | Bench LPIPS↓ | Bench FID↓ | Suitcase LPIPS↓ | Suitcase FID↓ | Teddybear LPIPS↓ | Teddybear FID↓ | Plant LPIPS↓ | Plant FID↓ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PN | 0.67 | 374.06 | 0.65 | 306.26 | 0.49 | 355.23 | 0.44 | 314.34 | 0.64 | 420.66 | 0.71 | 360.53 | 0.62 | 390.23 | 0.49 | 382.97 | 0.65 | 417.98 | 0.55 | 390.05 |
| EFT | 0.32 | 227.03 | 0.34 | 221.76 | 0.24 | 275.57 | 0.27 | 287.28 | 0.40 | 342.09 | 0.37 | 240.34 | 0.41 | 361.97 | 0.31 | 309.97 | 0.35 | 327.01 | 0.38 | 340.54 |
| VF | 0.29 | 197.35 | 0.26 | 128.33 | 0.23 | 232.44 | 0.22 | 188.77 | 0.33 | 289.76 | 0.32 | 226.12 | 0.31 | 349.23 | 0.27 | 286.24 | 0.33 | 286.78 | 0.31 | 297.25 |
| SF | **0.22** | **114.05** | **0.20** | 66.15 | 0.16 | 153.11 | **0.18** | **140.14** | 0.28 | 243.81 | **0.24** | 116.37 | 0.29 | 350.96 | 0.23 | 254.58 | 0.25 | 196.95 | 0.26 | 236.05 |
| Ours | 0.24 | 123.39 | 0.21 | **58.99** | **0.15** | **126.90** | **0.18** | 148.01 | 0.30 | 237.69 | **0.24** | **109.29** | **0.25** | **205.70** | **0.19** | **177.85** | **0.23** | **119.87** | **0.24** | **168.77** |

| | Donut CLIP↑ | Donut DISTS↓ | Apple CLIP↑ | Apple DISTS↓ | Hydrant CLIP↑ | Hydrant DISTS↓ | Vase CLIP↑ | Vase DISTS↓ | Cake CLIP↑ | Cake DISTS↓ | Ball CLIP↑ | Ball DISTS↓ | Bench CLIP↑ | Bench DISTS↓ | Suitcase CLIP↑ | Suitcase DISTS↓ | Teddybear CLIP↑ | Teddybear DISTS↓ | Plant CLIP↑ | Plant DISTS↓ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PN | 0.80 | 0.48 | 0.80 | 0.46 | 0.83 | 0.42 | 0.84 | 0.41 | 0.82 | 0.43 | 0.81 | 0.44 | 0.87 | 0.42 | 0.86 | 0.42 | 0.83 | 0.43 | 0.82 | 0.47 |
| EFT | 0.89 | 0.30 | 0.88 | 0.31 | 0.85 | 0.33 | 0.89 | 0.32 | 0.87 | 0.34 | 0.86 | 0.30 | 0.87 | 0.38 | 0.87 | 0.34 | 0.87 | 0.34 | 0.86 | 0.37 |
| VF | 0.90 | 0.28 | 0.91 | 0.25 | 0.85 | 0.29 | 0.86 | 0.28 | 0.88 | 0.30 | 0.86 | 0.28 | 0.87 | 0.33 | 0.87 | 0.29 | 0.84 | 0.30 | 0.84 | 0.32 |
| SF | **0.93** | **0.23** | 0.94 | 0.21 | 0.91 | 0.23 | 0.92 | **0.24** | **0.91** | 0.27 | **0.92** | **0.21** | 0.89 | 0.37 | 0.89 | 0.31 | 0.90 | 0.24 | 0.91 | 0.28 |
| Ours | **0.93** | **0.23** | **0.96** | **0.20** | **0.93** | **0.19** | **0.93** | **0.24** | **0.91** | 0.26 | **0.92** | 0.22 | **0.92** | 0.25 | **0.92** | 0.22 | **0.93** | **0.22** | **0.93** | **0.24** |

Table 3 — **3 Views**

| | Donut LPIPS↓ | Donut FID↓ | Apple LPIPS↓ | Apple FID↓ | Hydrant LPIPS↓ | Hydrant FID↓ | Vase LPIPS↓ | Vase FID↓ | Cake LPIPS↓ | Cake FID↓ | Ball LPIPS↓ | Ball FID↓ | Bench LPIPS↓ | Bench FID↓ | Suitcase LPIPS↓ | Suitcase FID↓ | Teddybear LPIPS↓ | Teddybear FID↓ | Plant LPIPS↓ | Plant FID↓ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PN | 0.66 | 368.98 | 0.65 | 308.78 | 0.49 | 345.07 | 0.43 | 298.97 | 0.62 | 416.85 | 0.69 | 363.16 | 0.61 | 381.76 | 0.47 | 377.21 | 0.63 | 410.14 | 0.53 | 365.95 |
| EFT | 0.28 | 181.47 | 0.29 | 163.61 | 0.21 | 230.53 | 0.24 | 238.92 | 0.35 | 287.92 | 0.31 | 180.23 | 0.37 | 331.04 | 0.27 | 270.05 | 0.29 | 248.43 | 0.33 | 296.55 |
| VF | 0.29 | 196.07 | 0.25 | 117.17 | 0.22 | 224.68 | 0.21 | 187.42 | 0.33 | 280.81 | 0.31 | 216.30 | 0.30 | 344.33 | 0.26 | 280.67 | 0.32 | 269.22 | 0.31 | 285.39 |
| SF | **0.22** | 117.39 | **0.19** | 57.36 | 0.15 | 149.60 | **0.18** | 141.91 | **0.27** | 233.56 | 0.23 | 98.83 | 0.27 | 332.04 | 0.21 | 224.09 | 0.22 | 167.79 | 0.25 | 227.67 |
| Ours | **0.22** | **110.30** | 0.20 | **56.31** | **0.14** | **117.87** | 0.17 | **129.87** | **0.27** | **209.09** | **0.22** | **93.80** | **0.22** | **164.42** | 0.17 | 149.40 | **0.20** | **106.93** | **0.23** | **161.20** |

| | Donut CLIP↑ | Donut DISTS↓ | Apple CLIP↑ | Apple DISTS↓ | Hydrant CLIP↑ | Hydrant DISTS↓ | Vase CLIP↑ | Vase DISTS↓ | Cake CLIP↑ | Cake DISTS↓ | Ball CLIP↑ | Ball DISTS↓ | Bench CLIP↑ | Bench DISTS↓ | Suitcase CLIP↑ | Suitcase DISTS↓ | Teddybear CLIP↑ | Teddybear DISTS↓ | Plant CLIP↑ | Plant DISTS↓ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PN | 0.81 | 0.47 | 0.81 | 0.45 | 0.83 | 0.41 | 0.85 | 0.40 | 0.83 | 0.43 | 0.82 | 0.43 | 0.87 | 0.40 | 0.86 | 0.42 | 0.84 | 0.42 | 0.83 | 0.46 |
| EFT | 0.91 | 0.28 | 0.90 | 0.28 | 0.88 | 0.30 | 0.91 | 0.29 | 0.89 | 0.31 | 0.89 | 0.27 | 0.87 | 0.35 | 0.88 | 0.31 | 0.90 | 0.30 | 0.89 | 0.34 |
| VF | 0.91 | 0.27 | 0.92 | 0.24 | 0.86 | 0.29 | 0.86 | 0.28 | 0.88 | 0.29 | 0.88 | 0.27 | 0.87 | 0.32 | 0.87 | 0.29 | 0.85 | 0.29 | 0.85 | 0.31 |
| SF | 0.93 | 0.22 | 0.94 | 0.19 | 0.92 | 0.22 | 0.92 | **0.24** | **0.92** | 0.26 | **0.94** | **0.20** | 0.89 | 0.33 | 0.90 | 0.27 | 0.92 | 0.22 | 0.92 | 0.26 |
| Ours | **0.94** | **0.21** | **0.96** | **0.18** | **0.94** | **0.18** | **0.93** | 0.23 | **0.92** | **0.24** | **0.94** | **0.20** | **0.92** | 0.23 | **0.93** | 0.21 | **0.93** | **0.21** | **0.94** | **0.20** |

Table 3 — **6 Views**

| | Donut LPIPS↓ | Donut FID↓ | Apple LPIPS↓ | Apple FID↓ | Hydrant LPIPS↓ | Hydrant FID↓ | Vase LPIPS↓ | Vase FID↓ | Cake LPIPS↓ | Cake FID↓ | Ball LPIPS↓ | Ball FID↓ | Bench LPIPS↓ | Bench FID↓ | Suitcase LPIPS↓ | Suitcase FID↓ | Teddybear LPIPS↓ | Teddybear FID↓ | Plant LPIPS↓ | Plant FID↓ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PN | 0.64 | 355.81 | 0.64 | 291.21 | 0.46 | 312.17 | 0.40 | 274.33 | 0.60 | 397.98 | 0.68 | 350.32 | 0.57 | 375.22 | 0.43 | 359.19 | 0.60 | 407.98 | 0.48 | 321.57 |
| EFT | 0.22 | 110.54 | 0.22 | 69.22 | 0.15 | 147.47 | 0.20 | 198.95 | 0.28 | 198.95 | 0.24 | 110.57 | 0.30 | 259.92 | 0.21 | 191.85 | 0.22 | 159.42 | 0.26 | 204.50 |
| VF | 0.28 | 194.34 | 0.24 | 106.93 | 0.21 | 212.40 | 0.21 | 189.51 | 0.32 | 273.57 | 0.29 | 207.04 | 0.29 | 338.11 | 0.25 | 274.26 | 0.30 | 251.17 | 0.30 | 275.71 |
| SF | 0.20 | 113.81 | **0.17** | 54.37 | 0.14 | 137.18 | 0.17 | 137.46 | 0.26 | 214.34 | 0.21 | **79.95** | 0.21 | 181.71 | 0.19 | 181.71 | 0.21 | 143.68 | 0.23 | 198.32 |
| Ours | **0.19** | **93.80** | **0.17** | **53.41** | **0.12** | **106.89** | **0.16** | **127.75** | **0.23** | **166.95** | **0.20** | 82.32 | **0.20** | **168.25** | **0.15** | **133.26** | **0.18** | **89.66** | **0.20** | **138.67** |

| | Donut CLIP↑ | Donut DISTS↓ | Apple CLIP↑ | Apple DISTS↓ | Hydrant CLIP↑ | Hydrant DISTS↓ | Vase CLIP↑ | Vase DISTS↓ | Cake CLIP↑ | Cake DISTS↓ | Ball CLIP↑ | Ball DISTS↓ | Bench CLIP↑ | Bench DISTS↓ | Suitcase CLIP↑ | Suitcase DISTS↓ | Teddybear CLIP↑ | Teddybear DISTS↓ | Plant CLIP↑ | Plant DISTS↓ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PN | 0.82 | 0.47 | 0.83 | 0.43 | 0.84 | 0.39 | 0.86 | 0.39 | 0.85 | 0.41 | 0.84 | 0.43 | 0.89 | 0.39 | 0.88 | 0.40 | 0.86 | 0.41 | 0.84 | 0.45 |
| EFT | 0.94 | 0.23 | 0.94 | 0.22 | 0.92 | 0.24 | 0.93 | 0.26 | 0.93 | 0.26 | 0.93 | 0.22 | 0.90 | 0.32 | 0.92 | 0.26 | 0.94 | 0.24 | 0.92 | 0.28 |
| VF | 0.91 | 0.27 | 0.93 | 0.23 | 0.88 | 0.28 | 0.88 | 0.28 | 0.89 | 0.29 | 0.90 | 0.26 | 0.88 | 0.31 | 0.89 | 0.28 | 0.87 | 0.28 | 0.86 | 0.30 |
| SF | 0.94 | 0.21 | 0.95 | 0.17 | 0.93 | 0.21 | 0.93 | 0.23 | 0.93 | 0.25 | **0.95** | **0.18** | 0.91 | 0.28 | 0.92 | 0.24 | 0.93 | 0.20 | 0.93 | 0.24 |
| Ours | **0.95** | **0.20** | **0.97** | **0.16** | **0.95** | **0.17** | **0.94** | **0.22** | **0.94** | **0.22** | **0.95** | 0.19 | **0.94** | **0.21** | **0.95** | **0.18** | **0.95** | **0.18** | **0.95** | **0.20** |

Table 3: **Quantitative comparison of novel-view synthesis on unseen instances for each category, with varying number of input views (2, 3 and 6).**

Table 4 — **2 Views**

| | Bicycle LPIPS↓ | Bicycle FID↓ | Car LPIPS↓ | Car FID↓ | Couch LPIPS↓ | Couch FID↓ | Laptop LPIPS↓ | Laptop FID↓ | Microwave LPIPS↓ | Microwave FID↓ | Motorcycle LPIPS↓ | Motorcycle FID↓ | Bowl LPIPS↓ | Bowl FID↓ | Toyplane LPIPS↓ | Toyplane FID↓ | TV LPIPS↓ | TV FID↓ | Wineglass LPIPS↓ | Wineglass FID↓ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PN | 0.56 | 351.03 | 0.52 | 310.50 | 0.51 | 369.59 | 0.49 | 320.08 | 0.56 | 342.70 | 0.50 | 376.82 | 0.56 | 390.23 | 0.50 | 377.76 | 0.53 | 344.72 | 0.50 | 314.45 |
| EFT | 0.42 | 319.14 | 0.38 | 280.12 | 0.50 | 323.11 | 0.42 | 312.41 | 0.49 | 329.69 | 0.44 | 360.55 | 0.37 | 283.81 | 0.34 | 315.92 | 0.46 | 319.86 | 0.30 | 308.91 |
| VF | 0.32 | 301.07 | 0.31 | 300.26 | 0.40 | 358.38 | 0.37 | 281.72 | 0.38 | 328.33 | 0.36 | 346.32 | 0.28 | 237.76 | 0.33 | 316.59 | 0.37 | 321.62 | 0.25 | 219.90 |
| SF | 0.29 | 329.32 | 0.26 | 278.98 | 0.37 | 311.19 | 0.32 | 318.97 | 0.37 | 317.61 | 0.30 | 365.12 | 0.21 | 196.06 | 0.22 | 265.71 | 0.30 | 317.07 | 0.18 | 204.46 |
| Ours | **0.23** | **162.88** | **0.19** | **80.67** | **0.32** | **270.68** | **0.23** | **130.68** | **0.31** | **214.97** | **0.26** | **206.15** | **0.16** | **97.94** | **0.19** | **168.73** | **0.28** | **236.83** | **0.16** | **73.51** |

| | Bicycle CLIP↑ | Bicycle DISTS↓ | Car CLIP↑ | Car DISTS↓ | Couch CLIP↑ | Couch DISTS↓ | Laptop CLIP↑ | Laptop DISTS↓ | Microwave CLIP↑ | Microwave DISTS↓ | Motorcycle CLIP↑ | Motorcycle DISTS↓ | Bowl CLIP↑ | Bowl DISTS↓ | Toyplane CLIP↑ | Toyplane DISTS↓ | TV CLIP↑ | TV DISTS↓ | Wineglass CLIP↑ | Wineglass DISTS↓ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PN | 0.86 | 0.47 | 0.85 | 0.45 | 0.90 | 0.42 | 0.84 | 0.45 | 0.87 | 0.40 | 0.86 | 0.50 | 0.83 | 0.43 | 0.79 | 0.51 | 0.89 | 0.41 | 0.81 | 0.44 |
| EFT | 0.86 | 0.43 | 0.85 | 0.38 | 0.90 | 0.38 | 0.84 | 0.37 | 0.88 | 0.36 | 0.87 | 0.43 | 0.88 | 0.32 | 0.83 | 0.38 | 0.90 | 0.36 | 0.85 | 0.35 |
| VF | 0.85 | 0.41 | 0.81 | 0.37 | 0.86 | 0.39 | 0.84 | 0.36 | 0.86 | 0.35 | 0.85 | 0.37 | 0.90 | 0.29 | 0.85 | 0.39 | 0.87 | 0.34 | 0.85 | 0.32 |
| SF | 0.86 | 0.43 | 0.84 | 0.38 | 0.90 | **0.35** | 0.86 | 0.35 | 0.88 | 0.35 | 0.86 | 0.40 | 0.92 | 0.24 | 0.86 | 0.30 | 0.91 | **0.30** | 0.88 | 0.28 |
| Ours | **0.93** | **0.30** | **0.93** | **0.21** | **0.92** | 0.37 | **0.91** | **0.23** | **0.92** | **0.29** | **0.92** | **0.28** | **0.96** | **0.19** | **0.92** | **0.23** | **0.93** | **0.30** | **0.94** | **0.20** |

Table 4 — **3 Views**

| | Bicycle LPIPS↓ | Bicycle FID↓ | Car LPIPS↓ | Car FID↓ | Couch LPIPS↓ | Couch FID↓ | Laptop LPIPS↓ | Laptop FID↓ | Microwave LPIPS↓ | Microwave FID↓ | Motorcycle LPIPS↓ | Motorcycle FID↓ | Bowl LPIPS↓ | Bowl FID↓ | Toyplane LPIPS↓ | Toyplane FID↓ | TV LPIPS↓ | TV FID↓ | Wineglass LPIPS↓ | Wineglass FID↓ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PN | 0.46 | 316.75 | 0.51 | 291.66 | 0.54 | 358.22 | 0.52 | 308.25 | 0.53 | 336.86 | 0.52 | 372.39 | 0.54 | 366.09 | 0.49 | 363.66 | 0.52 | 286.14 | 0.52 | 305.69 |
| EFT | 0.38 | 261.54 | 0.34 | 232.67 | 0.48 | 303.02 | 0.39 | 273.40 | 0.46 | 292.85 | 0.40 | 331.59 | 0.30 | 229.57 | 0.29 | 271.01 | 0.42 | 286.14 | 0.23 | 225.28 |
| VF | 0.31 | 289.28 | 0.30 | 286.54 | 0.40 | 353.42 | 0.37 | 266.16 | 0.38 | 317.59 | 0.36 | 347.63 | 0.26 | 216.74 | 0.33 | 314.39 | 0.36 | 314.44 | 0.24 | 202.62 |
| SF | 0.27 | 320.25 | 0.24 | 257.96 | 0.36 | 303.58 | 0.32 | 330.44 | 0.36 | 300.86 | 0.29 | 362.59 | 0.18 | 162.65 | 0.21 | 254.88 | 0.31 | 328.60 | 0.16 | 166.70 |
| Ours | **0.22** | **153.23** | **0.17** | **69.97** | **0.32** | **265.84** | **0.22** | **116.32** | **0.30** | **202.69** | **0.24** | **198.09** | **0.14** | **85.80** | **0.14** | **149.40** | **0.26** | **225.45** | **0.14** | **74.39** |

| | Bicycle CLIP↑ | Bicycle DISTS↓ | Car CLIP↑ | Car DISTS↓ | Couch CLIP↑ | Couch DISTS↓ | Laptop CLIP↑ | Laptop DISTS↓ | Microwave CLIP↑ | Microwave DISTS↓ | Motorcycle CLIP↑ | Motorcycle DISTS↓ | Bowl CLIP↑ | Bowl DISTS↓ | Toyplane CLIP↑ | Toyplane DISTS↓ | TV CLIP↑ | TV DISTS↓ | Wineglass CLIP↑ | Wineglass DISTS↓ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PN | 0.87 | 0.48 | 0.85 | 0.44 | 0.92 | 0.38 | 0.86 | 0.42 | 0.82 | 0.41 | 0.87 | 0.47 | 0.84 | 0.42 | 0.80 | 0.51 | 0.90 | 0.40 | 0.82 | 0.44 |
| EFT | 0.89 | 0.40 | 0.88 | 0.35 | 0.92 | **0.36** | 0.87 | 0.34 | 0.90 | 0.35 | 0.88 | 0.41 | 0.90 | 0.29 | 0.86 | 0.35 | 0.92 | 0.34 | 0.88 | 0.31 |
| VF | 0.86 | 0.39 | 0.82 | 0.36 | 0.87 | 0.38 | 0.86 | 0.35 | 0.88 | 0.35 | 0.85 | 0.36 | 0.91 | 0.27 | 0.85 | 0.39 | 0.89 | 0.33 | 0.86 | 0.31 |
| SF | 0.87 | 0.39 | 0.86 | 0.35 | 0.91 | **0.36** | 0.87 | 0.34 | 0.89 | 0.35 | 0.86 | 0.38 | 0.93 | 0.22 | 0.87 | 0.29 | 0.92 | 0.37 | 0.91 | 0.26 |
| Ours | **0.94** | **0.27** | **0.93** | **0.19** | **0.93** | 0.37 | **0.92** | **0.22** | **0.93** | **0.29** | **0.93** | **0.26** | **0.96** | **0.17** | **0.93** | **0.20** | **0.94** | **0.28** | **0.95** | **0.19** |

Table 4 — **6 Views**

| | Bicycle LPIPS↓ | Bicycle FID↓ | Car LPIPS↓ | Car FID↓ | Couch LPIPS↓ | Couch FID↓ | Laptop LPIPS↓ | Laptop FID↓ | Microwave LPIPS↓ | Microwave FID↓ | Motorcycle LPIPS↓ | Motorcycle FID↓ | Bowl LPIPS↓ | Bowl FID↓ | Toyplane LPIPS↓ | Toyplane FID↓ | TV LPIPS↓ | TV FID↓ | Wineglass LPIPS↓ | Wineglass FID↓ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PN | 0.45 | 305.20 | 0.48 | 272.57 | 0.51 | 346.37 | 0.47 | 270.74 | 0.53 | 319.84 | 0.47 | 341.82 | 0.50 | 313.18 | 0.53 | 314.44 | 0.48 | 310.93 | 0.47 | 283.57 |
| EFT | 0.30 | 226.22 | 0.27 | 167.28 | 0.41 | 249.15 | 0.29 | 189.89 | 0.40 | 219.17 | 0.32 | 269.31 | 0.21 | 182.22 | 0.21 | 138.38 | 0.36 | 216.08 | 0.18 | 168.32 |
| VF | 0.30 | 267.21 | 0.28 | 270.05 | 0.37 | 336.73 | 0.35 | 250.42 | 0.36 | 285.93 | 0.33 | 331.89 | 0.24 | 183.85 | 0.30 | 288.68 | 0.34 | 308.65 | 0.20 | 185.90 |
| SF | 0.26 | 304.24 | 0.22 | 229.48 | 0.33 | 286.58 | 0.29 | 297.36 | 0.35 | 264.92 | 0.27 | 341.74 | 0.17 | 137.22 | 0.19 | 205.89 | 0.29 | 281.12 | 0.16 | 160.85 |
| Ours | **0.19** | **136.41** | **0.14** | **62.36** | **0.30** | **217.39** | **0.19** | **85.14** | **0.29** | **180.00** | **0.21** | **156.93** | **0.11** | **69.05** | **0.13** | **107.69** | **0.24** | **177.08** | **0.13** | **69.20** |

| | Bicycle CLIP↑ | Bicycle DISTS↓ | Car CLIP↑ | Car DISTS↓ | Couch CLIP↑ | Couch DISTS↓ | Laptop CLIP↑ | Laptop DISTS↓ | Microwave CLIP↑ | Microwave DISTS↓ | Motorcycle CLIP↑ | Motorcycle DISTS↓ | Bowl CLIP↑ | Bowl DISTS↓ | Toyplane CLIP↑ | Toyplane DISTS↓ | TV CLIP↑ | TV DISTS↓ | Wineglass CLIP↑ | Wineglass DISTS↓ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PN | 0.88 | 0.44 | 0.82 | 0.42 | 0.93 | 0.36 | 0.89 | 0.41 | 0.89 | 0.38 | 0.89 | 0.46 | 0.86 | 0.40 | 0.82 | 0.48 | 0.93 | 0.40 | 0.83 | 0.43 |
| EFT | 0.91 | 0.34 | 0.91 | 0.31 | 0.93 | 0.32 | 0.92 | 0.28 | 0.92 | 0.30 | 0.90 | 0.36 | 0.94 | 0.23 | 0.91 | 0.28 | 0.95 | 0.30 | 0.92 | 0.26 |
| VF | 0.88 | 0.37 | 0.83 | 0.33 | 0.89 | 0.37 | 0.90 | 0.33 | 0.89 | 0.32 | 0.86 | 0.34 | 0.92 | 0.25 | 0.88 | 0.35 | 0.91 | 0.32 | 0.88 | 0.28 |
| SF | 0.88 | 0.35 | 0.88 | 0.31 | 0.92 | 0.34 | 0.90 | 0.30 | 0.91 | 0.31 | 0.88 | 0.35 | 0.94 | 0.21 | 0.90 | 0.26 | 0.93 | 0.30 | 0.93 | 0.26 |
| Ours | **0.95** | **0.23** | **0.95** | **0.18** | **0.94** | 0.30 | **0.94** | **0.18** | **0.94** | **0.26** | **0.94** | **0.22** | **0.97** | **0.14** | **0.95** | **0.17** | **0.96** | **0.24** | **0.95** | **0.18** |

Table 4: **Quantitative comparison of novel-view synthesis on unseen categories for each category, with varying number of input views (2, 3 and 6).**

| | Unseen Instances | | | | | | | | | | | | | | | | | | | |
| | 2 Views | | | | | | | | | | | | | | | | | | | |
| | Donut | | Apple | | Hydrant | | Vase | | Cake | | Ball | | Bench | | Suitcase | | Teddybear | | Plant | |
| | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ |
| SF | 0.34 | 0.18 | 0.24 | 0.27 | 0.15 | 0.27 | 0.14 | 0.32 | 0.40 | 0.16 | 0.33 | 0.25 | 0.41 | 0.19 | 0.28 | 0.21 | 0.20 | 0.20 | **0.21** | 0.20 |
| Ours | **0.27** | **0.26** | **0.15** | **0.42** | **0.13** | **0.45** | **0.12** | **0.42** | **0.35** | **0.22** | **0.27** | **0.32** | **0.21** | **0.27** | **0.17** | **0.39** | **0.15** | **0.39** | 0.24 | **0.26** |
| | 3 Views | | | | | | | | | | | | | | | | | | | |
| | Donut | | Apple | | Hydrant | | Vase | | Cake | | Ball | | Bench | | Suitcase | | Teddybear | | Plant | |
| | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ |
| SF | 0.34 | 0.19 | 0.25 | 0.31 | 0.15 | 0.25 | **0.12** | 0.34 | 0.40 | 0.15 | 0.33 | 0.26 | 0.37 | 0.18 | 0.24 | 0.24 | 0.18 | 0.22 | **0.21** | 0.18 |
| Ours | **0.24** | **0.31** | **0.14** | **0.45** | **0.11** | **0.49** | 0.17 | **0.45** | **0.28** | **0.28** | **0.26** | **0.36** | **0.17** | **0.29** | **0.16** | **0.43** | **0.13** | **0.47** | 0.22 | **0.29** |
| | 6 Views | | | | | | | | | | | | | | | | | | | |
| | Donut | | Apple | | Hydrant | | Vase | | Cake | | Ball | | Bench | | Suitcase | | Teddybear | | Plant | |
| | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ |
| SF | 0.32 | 0.22 | 0.26 | 0.30 | 0.13 | 0.26 | 0.14 | 0.36 | 0.38 | 0.18 | 0.33 | 0.27 | 0.26 | 0.20 | 0.23 | 0.28 | 0.19 | 0.23 | 0.21 | 0.18 |
| Ours | **0.22** | **0.38** | **0.13** | **0.52** | **0.09** | **0.59** | **0.12** | **0.48** | **0.25** | **0.34** | **0.25** | **0.42** | **0.13** | **0.39** | **0.14** | **0.50** | **0.11** | **0.56** | **0.18** | **0.36** |
| | Unseen Categories | | | | | | | | | | | | | | | | | | | |
| | 2 Views | | | | | | | | | | | | | | | | | | | |
| | Bicycle | | Car | | Couch | | Laptop | | Microwave | | Motorcycle | | Bowl | | Toyplane | | TV | | Wineglass | |
| | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ |
| SF | 0.44 | 0.16 | 0.35 | 0.23 | 0.66 | 0.09 | 0.45 | 0.15 | 0.52 | 0.14 | 0.24 | 0.24 | 0.35 | 0.19 | 0.22 | 0.23 | 0.36 | 0.15 | 0.18 | 0.23 |
| Ours | **0.32** | **0.28** | **0.24** | **0.33** | **0.53** | **0.16** | **0.27** | **0.20** | **0.28** | **0.26** | **0.23** | **0.27** | **0.18** | **0.33** | **0.19** | **0.32** | **0.31** | **0.27** | **0.14** | **0.36** |
| | 3 Views | | | | | | | | | | | | | | | | | | | |
| | Bicycle | | Car | | Couch | | Laptop | | Microwave | | Motorcycle | | Bowl | | Toyplane | | TV | | Wineglass | |
| | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ |
| SF | 0.41 | 0.16 | **0.25** | 0.28 | **0.65** | 0.11 | 0.43 | 0.16 | 0.49 | 0.16 | 0.21 | 0.27 | 0.30 | 0.22 | 0.20 | 0.22 | 0.60 | 0.15 | 0.14 | 0.24 |
| Our | **0.27** | **0.33** | 0.42 | **0.39** | 0.75 | **0.17** | **0.20** | **0.25** | **0.30** | **0.25** | **0.15** | **0.35** | **0.15** | **0.39** | **0.15** | **0.38** | **0.31** | **0.31** | **0.11** | **0.42** |
| | 6 Views | | | | | | | | | | | | | | | | | | | |
| | Bicycle | | Car | | Couch | | Laptop | | Microwave | | Motorcycle | | Bowl | | Toyplane | | TV | | Wineglass | |
| | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ | CD↓ | F-score↑ |
| SF | 0.33 | 0.16 | **0.22** | 0.34 | **0.61** | 0.14 | 0.31 | 0.23 | 0.43 | 0.19 | 0.15 | 0.30 | 0.23 | 0.26 | 0.17 | 0.28 | 0.43 | 0.17 | 0.15 | 0.26 |
| Our | **0.19** | **0.42** | 0.37 | **0.47** | 0.63 | **0.24** | **0.15** | **0.34** | **0.29** | **0.29** | **0.09** | **0.44** | **0.14** | **0.45** | **0.12** | **0.46** | **0.22** | **0.38** | **0.11** | **0.44** |

Table 5: **Quantitative comparison of geometry reconstruction both on unseen instances and unseen categories for each category, with varying number of input views (2, 3 and 6).**