

Distilling Diffusion Models for 3D Reconstruction

Zhizhuo (Z) Zhou

CMU-RI-TR-23-33

July 12



The Robotics Institute
School of Computer Science
Carnegie Mellon University
Pittsburgh, PA

Thesis Committee:

Professor, Shubham Tulsiani *chair*

Professor Jun-Yan Zhu

Professor Deva K. Ramanan

Jason Zhang, *Carnegie Mellon University*

*Submitted in partial fulfillment of the requirements
for the degree of Master of Science in Robotics.*

Copyright © 2023 Zhizhuo (Z) Zhou. All rights reserved.

To my parents.

Abstract

We propose SparseFusion, a sparse view 3D reconstruction approach that unifies recent advances in neural rendering and probabilistic image generation. Existing approaches typically build on neural rendering with re-projected features but fail to generate unseen regions or handle uncertainty under large viewpoint changes. Alternate methods treat this as a (probabilistic) 2D synthesis task, and while they can generate plausible 2D images, they do not infer a consistent underlying 3D. However, we find that this trade-off between 3D consistency and probabilistic image generation does not need to exist. In fact, we show that geometric consistency and generative inference can be complementary in a mode-seeking behavior. By distilling a 3D consistent scene representation from a view-conditioned latent diffusion model, we are able to recover a plausible 3D representation whose renderings are both accurate and realistic. We evaluate our approach across 51 categories in the CO3D dataset and show that it outperforms existing methods, in both distortion and perception metrics, for sparse-view novel view synthesis.

Acknowledgments

The past two years at CMU has been a time of my own reconstruction. I learned a lot. I laughed a lot. I cried sometimes. But most importantly, I grew a lot. And to say that my advisors, peers, and friends helped me is an understatement, they literally pulled me up with them. It really takes a whole community to raise someone. I am so grateful to my community.

I am most grateful to my advisor, Shubham Tulsiani. Without his compassionate kindness, enduring patience, and unwavering support, I would not have made it. I had an incredibly slow start. Many of my projects failed. Despite this, Shubham was there every week, guiding me, through thick and thin. I am eternally grateful. Thank you, Shubham!

I also thank the professors on my committee, Jun-Yan Zhu and Deva Ramanan. I give a special shout out to David Fouhey, who took the naive, undergrad me under his wings and supported my first foray into research.

While research can be a lonely enterprise, my labmates made it a lot more enjoyable. I am thankful to Yen-Chi Cheng and Paritosh Mittal, who patiently answered my questions when I just started dipping my toes in research. I am also thankful to Naveen Venkat and Mayank Agarwal, who provided countless helpful discussions and emotional support from my neighboring table. More recently, I am thankful to Barath Raj for being a titanic role model of hardwork and Hanzhe Hu for being an amazing friend who shared many foolish adventures with me. Moreover, I am thankful to Judy Ye, Jason Zhang, Yehonathan Litman, Homanga Bharadhwaj, Poorvi Hebbar, and Amy Lin for their kind friendships and awesome energy.

In addition to my labmates, I am incredibly thankful to my friends who lit up my heart and made me smile even in the darkest of times. Ziyu (Ken) Liu inspired me to look at where I want to go in life. Gaoyue (Kathy) Zhou made me regret not living in Squirrel Hill. Daohan (Fred) Lu embodied the big D. Jeff Tan rekindled my high-school work ethic. Chonghyuk (Andrew) Song showed me the best calamari. Lihong Jin dragged me through tedious homework assignments. Ruihan (Gary) Gao and Wenyu Xia showed me the joys of running and singing Taylor Swift. Lilly Wang took me to the unforgettable Eras Tour at Ford Field. Himangi Mittal and Anurag Ghosh showed me the fun of working on silly projects.

I am thankful to my friends at ARCC, who gave me many memories to

cherish. I am also thankful to many more members of the community.

Last but not least, I would like to thank my family for always supporting me.

Funding

This work was supported in part by the National Science Foundation Graduate Research Fellowship under Grant No. (DGE1745016, DGE2140739).

Contents

1	Introduction	1
2	Related Work	5
2.1	Prior Works	5
2.2	Concurrent Works	7
3	Background	9
3.1	Denoising Diffusion Models	9
4	Method	11
4.1	Geometry-guided Probabilistic View Synthesis	12
4.2	Epipolar Feature Transformer	12
4.3	View-conditioned Latent Diffusion Model	13
4.4	Extracting 3D Modes via Diffusion Distillation	14
5	Experiments	17
5.1	Experimental Setup	17
5.2	Reconstruction on Real Images	20
5.2.1	Core Subset: 2-view.	20
5.2.2	Core Subset: Varying Views.	21
5.2.3	All Categories: 2-views.	22
5.2.4	Failure Modes.	23
5.3	Additional Analysis	23
5.3.1	Performance Binned by Viewpoint Changes.	23
5.3.2	Importance of Mode Seeking.	23
5.3.3	Ablating Distillation Objective.	25
6	Discussion	27
6.1	Limitations	27
6.2	Ethics and Broader Impact	27
A	Implementation Details	29
A.1	Epipolar Feature Transformer	29
A.2	View-conditioned Diffusion Model	31
A.3	Diffusion Distillation	33

When this dissertation is viewed as a PDF, the page header is a link to this Table of Contents.

List of Figures

1.1	Sparse-view Reconstruction. We present SparseFusion, an approach for 3D reconstruction given a few (<i>e.g.</i> , just two) segmented input images with known relative pose. SparseFusion is able to generate a 3D consistent neural scene representation, enabling us to render novel views and extract the underlying geometry, while being able to generate detailed and plausible structures in uncertain or unobserved regions (<i>e.g.</i> , front of hydrant, teddy’s face, back of laptop, or left side of toybus).	1
4.1	Overview of SparseFusion. SparseFusion comprises of two core components: a view-conditioned latent diffusion model (VLDM) and a diffusion distillation process that optimizes an Instant NGP [20, 39]. We use VLDM to model $p(\mathbf{x} \boldsymbol{\pi}, C)$	11
4.2	View-conditioned Diffusion. We show a diagram of our view-conditioned latent diffusion model. VLDM is conditioned on features \mathbf{y} , which is predicted by EFT.	12
4.3	Diffusion Samples. Given the same input features, the reverse sampling process of diffusion model results in different predictions for unseen regions.	13
4.4	Diffusion Distillation Diagram. We optimize the parameters θ of an Instant NGP network such that rendered images $f_\theta(\boldsymbol{\pi})$ from $\boldsymbol{\pi} \sim \Pi$ are similar to VLDM predictions $\hat{\mathbf{x}}_0$, effectively seeking a mode in $p_\phi(\mathbf{x} \boldsymbol{\pi}, C)$	14
4.5	Mode Seeking Visualization. We show qualitative comparison between a mode-seeking (SparseFusion) and a mean-seeking (VLDM+INGP) objective.	15
5.1	View Synthesis Qualitative Results. We show view synthesis results with 2 input views on donut, hydrant, cake, bench, teddybear, and plant categories. We visualize 2 novel views per instance with PixelNeRF (PN), NerFormer (NF), ViewFormer (VF), EFT, VLDM, and finally, SparseFusion (SF). Corresponding numbers can be found in Table 5.1.	19

5.2	Reconstruction on Diverse Categories. We show SparseFusion reconstructions on a subset of the 51 CO3D categories. We also show a couple of failure modes on the last row. Please see project page for more samples and 360-degree visualizations.	22
5.3	Metrics Binned by Viewpoint Change. We show metrics binned by the angle of query camera to the nearest context view. Results are aggregated from Table 5.1.	24
5.4	Qualitative Results with Pixel Space Loss. Using multi-step denoising and perceptual loss achieves more realistic results.	25
A.1	Epipolar Feature Transformer We show a diagram of EFT. This module processes each query ray independently, using a transformer to aggregate the projected features across views and across possible depths. For each ray, the output is a predicted RGB color (used as a baseline prediction method), and a pixel-aligned feature (used as conditioning in the diffusion model).	31

List of Tables

2.1	Comparison with prior methods. The rows indicate whether each method: 1) has been demonstrated on real world data, 2) works with sparse (2-6) input views, 3) generates geometrically consistent views, 4) generalizes to new scene instances, and 5) hallucinates unseen regions.	6
5.1	Detailed View Synthesis Benchmark. We show 2-view category-specific metrics on 10 CO3D categories from the <i>core subset</i> . We show PSNR \uparrow and LPIPS \downarrow averaged across 10 scenes per category.	20
5.2	View Synthesis on 10 Categories. We benchmark view synthesis results on the 10 categories with 2, 3, and 6 input views.	21
5.3	View Synthesis on 51 Categories. We benchmark novel view synthesis on all CO3D categories with 2 input views.	22
5.4	The Importance of Mode Seeking. We show metrics when EFT and VLDM are naively used to optimize Instant NGP [20] in a mean seeking behavior, versus the mode seeking optimization in SparseFusion. We average across 10 scenes of hydrants with 2 input views. . .	24
5.5	Diffusion Distillation Setup. We show that a combination of multi-step prediction and perceptual loss strikes a balance between all three metrics. (hydrant, 10 scenes, 2 input views)	25
A.1	EFT Configuration. We use default PyTorch hyperparameters for each layer. B is number of rays. M is the number of input views. D is the number of epipolar feature samples along the ray. D is 20.	32
A.2	UNet Parameters. We provide parameters for our UNet.	32
A.3	UNet Blocks. We outline the modules in our denoising UNet.	33

Chapter 1

Introduction

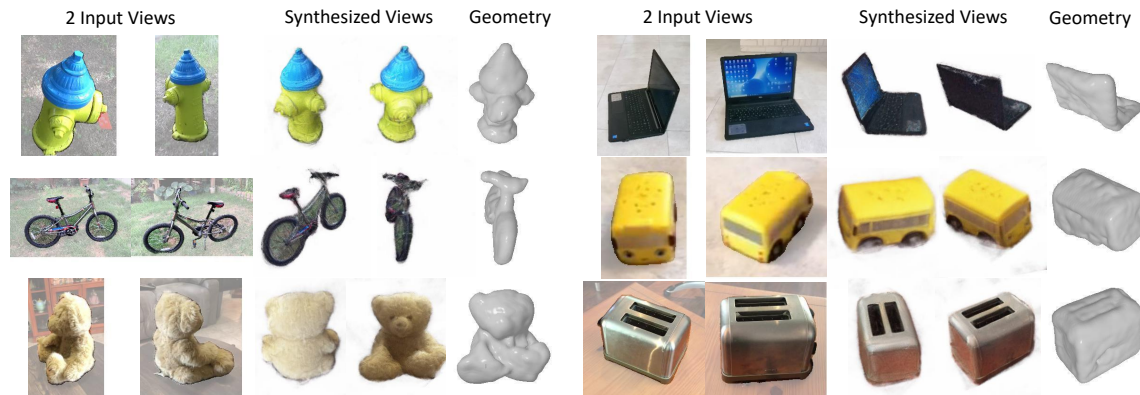


Figure 1.1: **Sparse-view Reconstruction.** We present SparseFusion, an approach for 3D reconstruction given a few (*e.g.*, just two) segmented input images with known relative pose. SparseFusion is able to generate a 3D consistent neural scene representation, enabling us to render novel views and extract the underlying geometry, while being able to generate detailed and plausible structures in uncertain or unobserved regions (*e.g.*, front of hydrant, teddy’s face, back of laptop, or left side of toybus).

Consider the two images of the teddybear shown in Figure 1.1 and try to imagine the underlying 3D object. Relying on the direct visual evidence in these images, you can easily infer that the teddybear is white, has a large head, and has small arms. Even more remarkably, you can imagine beyond the directly visible to estimate a *complete* 3D model of this object *e.g.*, forming a mental model of the teddy’s face with (likely black) eyes even though these were not observed. In this work, we build

1. Introduction

a computational approach that can similarly predict 3D from just a few images – by integrating visual measurements and priors via probabilistic modeling and then seeking likely 3D modes.

A growing number of recent works have studied the related tasks of *sparse-view* 3D reconstruction and novel view synthesis, *i.e.*, inferring 3D representations and/or synthesizing novel views of an object given just a few (typically 2-3) images with known relative camera poses. By leveraging data-driven priors, these approaches can learn to efficiently leverage multi-view cues and infer 3D from sparse views. However, they still yield blurry predictions under large viewpoint changes and cannot hallucinate plausible content in unobserved regions. This is because they do not account for the uncertainty in the outputs *e.g.*, the unobserved nose of a teddybear may be either red or black, but these methods, by reducing inference to independent pixel-wise or point-wise predictions, cannot model such variation.

In this work, we propose to instead model the *distribution* over the possible images given observations from some context views and an arbitrary query viewpoint. Leveraging a geometrically-informed backbone that computes pixel-aligned features in the query view, our approach learns a (conditional) diffusion model that can then infer detailed plausible novel-view images. While this probabilistic image synthesis approach allows the generation of higher quality image outputs, it does not directly yield a 3D representation of underlying the object. In fact, the (independently) sampled outputs for each query view often do not even correspond to a consistent underlying 3D *e.g.*, if the nose of the teddybear is unobserved in context views, one sampled query view may paint it red, while another one black.

To obtain a consistent 3D representation, we propose a *Diffusion Distillation* technique that ‘distills’ the predicted distributions into an instance-specific 3D representation. We note that the conditional diffusion model not only gives us the ability to sample novel-view images but also to (approximately) compute the likelihood of a generated one. Using this insight, we optimize an instance-specific (neural) 3D representation by maximizing the diffusion-based likelihood of its renderings. We show that this leads to a mode-seeking optimization that results in more accurate and realistic renderings, while also recovering a 3D-consistent representation of the underlying object. We demonstrate our approach on over 50 real-world categories from the CO3D dataset and show that our method allows recovering accurate 3D

and novel views given as few as 2 images as input – please see Figure 1.1 for sample results.

1. Introduction

Chapter 2

Related Work

2.1 Prior Works

Instance-specific Reconstruction from Multiple Views. Leveraging Structure-from-Motion [32, 36] to recover camera viewpoints, early Multi-view-Stereo (MVS) [7, 33] methods could recover dense 3D outputs. Recent neural incarnations of these [19, 48, 49] use volumetric rendering to learn a compact neural scene representation. Follow up works [3, 6, 20] seek to make the training and rendering orders of magnitudes faster. However, these methods require many input views, making them impractical for real world applications. While some works [10, 21, 52] seek to reduce the input views required, they still do not make predictions for unseen regions.

Single-view 3D Reconstruction. The ability to predict 3D geometry (and appearance) beyond the visible is a key goal for single-view 3D prediction methods. While these approaches have pursued prediction of different 3D representations *e.g.*, volumetric [4, 8, 12, 42, 50], mesh-based [9, 14], or neural implicit [16, 18, 43] 3D, the use of a single input image fundamentally limits the details that can be predicted. Moreover, these methods do not prioritize view synthesis as a goal. While our approach similarly learns data driven inference, we aim for a more detailed reconstruction and high quality novel-view renderings.

Generalizable View Synthesis from Fewer Views. Novel view synthesis (NVS), while similar to reconstruction, has slightly different roots. Earlier works [40, 55] frame

	Single-instance				Re-projection				Latent			Ours
	NeRF [19]	RegNeRF [21]	VolSDF [49]	NeRS [52]	IBRNet [46]	PixeNeRF [51]	NerFormer [24]	GPNR [38]	LFN [35]	SRT [31]	ViewFormer [15]	SparseFusion
1) Real data	✓	✓	✓	✓	✓	✓	✓	✓	×	✓	✓	✓
2) Sparse-views	×	✓	✓	✓	×	✓	✓	×	✓	✓	✓	✓
3) 3D consistent	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	×	✓
4) Generalization	×	×	×	×	✓	✓	✓	✓	✓	✓	✓	✓
5) Generate unseen	×	×	×	×	×	×	×	×	×	×	✓	✓

Table 2.1: **Comparison with prior methods.** The rows indicate whether each method: 1) has been demonstrated on real world data, 2) works with sparse (2-6) input views, 3) generates geometrically consistent views, 4) generalizes to new scene instances, and 5) hallucinates unseen regions.

NVS as a 2D problem, using deep networks to make predictions from global encodings. Recent approaches combine deep networks with various rendering formulations [31, 34, 35]. Strong performing approaches often leverage re-projected features from input views with volumetric rendering [24, 41, 51] or image based rendering [2, 38, 46]. While feature re-projection methods are 3D consistent, they regress to the mean and fail to produce perceptually sharp outputs. Another line of work [15, 26] revisits NVS as a probabilistic 2D generation task, using newer generative backbones to offer better perceptual quality at the cost of larger distortion and 3D consistency. See Table 2.1 for a comparison of our method against existing approaches.

Diffusion Models. Several works extend upon denoising diffusion models [13, 37] to achieve impressive applications, such as generating images from text [23, 30] and placing foreground objects in different backgrounds [29]. In this work, we leverage this class of models for (probabilistic) novel view synthesis while using geometry-aware features as conditioning. Inspired by the impressive results in DreamFusion [22] which optimized 3D scenes using text-conditioned diffusion models, we propose a view-conditioned diffusion distillation mechanism to similarly extract 3D modes in the sparse view reconstruction task.

2.2 Concurrent Works

Several concurrent works also leverage diffusion models for 3D reconstruction and view synthesis. 3DiM [47] proposes a 2D diffusion approach for image-conditioned novel view synthesis, but does not infer a 3D representation like our approach. Closer to our work, Deng *et al.* [5] uses (pre-trained) 2D diffusion models as guidance for single-view 3D, but obtain coarser reconstructions in this more challenging setting. While we leverage a 2D diffusion model for optimizing 3D, RenderDiffusion [1] learns a diffusion model in 3D space. Concurrently to DreamFusion [22], which inspired our distillation objective, Wang *et al.* [44] provide a different mathematical intuition for a similar objective.

2. Related Work

Chapter 3

Background

3.1 Denoising Diffusion Models

Denoising diffusion probabilistic models [13] approximate a distribution $p(\mathbf{x})$ over real data by reversing a Markov chain of diffusion steps, starting from Gaussian noise at \mathbf{x}_T to a realistic image at $\hat{\mathbf{x}}_0$. See [13] for details.

Forward Process. The forward diffusion process, which incrementally adds noise to a real image \mathbf{x}_0 until the image becomes Gaussian noise \mathbf{x}_T , is defined in Eq. 3.1. Forward variance β is usually defined by a fixed schedule.

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I}) \quad (3.1)$$

Reverse Process. The reverse diffusion process reverses the noise added in the forward process, effectively denoising a noisy image. When we generate a sample from a diffusion model, we apply the reverse process T times from $t = T$ to $t = 1$. The reverse process is defined in Eq. 3.2, where posterior mean $\mu_\phi(\mathbf{x}_t, t)$ is predicted from a network and posterior variance σ^2 follows a fixed schedule (though other works such as [37] also learn σ^2 with a network).

$$p(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mu_\phi(\mathbf{x}_t, t), \sigma^2\mathbf{I}) \quad (3.2)$$

3. Background

Posterior Mean. Prior works [13, 37] have found that parameterizing the neural network to predict $\boldsymbol{\epsilon}$ instead of \boldsymbol{x}_{t-1} or \boldsymbol{x}_0 works better in practice. We write posterior mean in terms of $\boldsymbol{\epsilon}$ in Eq. 3.3 where $\alpha_t = 1 - \beta_t$ and $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$.

$$\mu_\phi(\boldsymbol{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(\boldsymbol{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_\phi(\boldsymbol{x}_t, t) \right) \quad (3.3)$$

As mentioned in the main text, this parametrization leads to a training framework where one adds (time-dependent) noise to a data point \boldsymbol{x}_0 , and then trains the network $\boldsymbol{\epsilon}_\phi$ to predict this noise given the noisy data point \boldsymbol{x}_t .

$$\begin{aligned} \mathcal{L}_{DM} &= \mathbb{E}_{\boldsymbol{x}_0, \boldsymbol{\epsilon}, t} [w_t \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\phi(\boldsymbol{x}_t, t)\|^2] \\ \text{where } \boldsymbol{x}_t &= \sqrt{\bar{\alpha}_t} \boldsymbol{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}; \quad \boldsymbol{\epsilon} \sim \mathcal{N}(0, 1) \end{aligned} \quad (3.4)$$

In this work, we use conditional diffusion models to infer distributions of the form $p(\boldsymbol{x}|\boldsymbol{y})$ by additionally using \boldsymbol{y} as an input for the noise prediction network $\boldsymbol{\epsilon}_\phi(\boldsymbol{x}, \boldsymbol{y}, t)$.

Chapter 4

Method



Figure 4.1: **Overview of SparseFusion.** SparseFusion comprises of two core components: a view-conditioned latent diffusion model (VLDM) and a diffusion distillation process that optimizes an Instant NGP [20, 39]. We use VLDM to model $p(\mathbf{x}|\boldsymbol{\pi}, C)$.

Given sparse-view observations of an object (typically 2-3 images with masked foreground) with known camera viewpoints, our approach aims to infer a (3D) representation capable of synthesizing novel views while also capturing the geometric structure. However, as aspects of the object may be unobserved and its geometry difficult to precisely infer, direct prediction of 3D or novel views leads to implausibly blurry outputs in regions of uncertainty.

To enable plausible and 3D-consistent predictions, we instead take a two step approach as outlined in Figure 4.1. First, we learn a probabilistic view-synthesis model that, using geometry-guided diffusion, can model the *distribution* of images from query views given the sparse-view context (Section 4.1). While this allows the generation of detailed and diverse outputs, the obtained renderings lack 3D consistency. To extract a 3D representation, we propose a 3D neural distillation process that ‘distills’ the predicted view distributions into a 3D mode (Section 4.4).

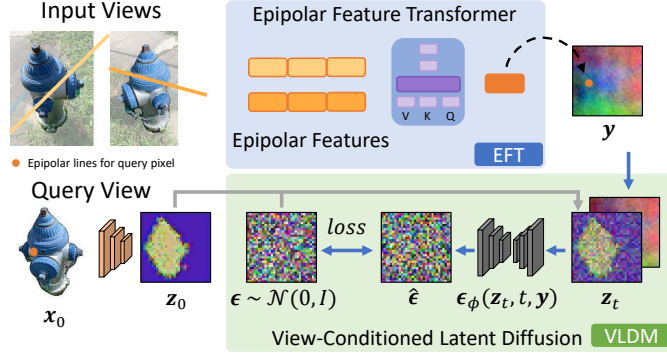


Figure 4.2: **View-conditioned Diffusion.** We show a diagram of our view-conditioned latent diffusion model. VLDM is conditioned on features y , which is predicted by EFT.

4.1 Geometry-guided Probabilistic View Synthesis

Given a target view pose π along with a set of reference images and their relative poses $C \equiv (\mathbf{x}_m, \pi_m)$, we want to model the conditional distribution $p(\mathbf{x}|\pi, C)$, from which we can synthesize an image $\hat{\mathbf{x}}$. We illustrate our approach to modeling this distribution in Figure 4.2. First, we use an epipolar feature transformer (EFT) inspired by [38] as feature extractor to obtain a low resolution feature grid y in the view space of π given the context C . In conjunction, we train a view-conditioned latent diffusion model (VLDM) that models the distribution over novel-view images condition on these geometry-aware features.

4.2 Epipolar Feature Transformer

We build upon GPNR [38] to extract features from context C . GPNR learns a feedforward network, $g_\psi(\mathbf{r}, C)$, that predicts color given a query ray \mathbf{r} by extracting features along its epipolar lines in all context images and aggregating them with transformers. We make several modifications to GPNR to suit our needs. First, we replace the patch projection layer with a ResNet18 [11] convolutional encoder as we found the lightweight patch encodings, while suitable for small baseline view synthesis, are not robust under the sparse-view setting. Furthermore, we modify the last layer

to predict both an RGB value and a feature vector. We denote the RGB branch as g_ψ and the feature branch as h_ψ . We refer to our modified epipolar patch-based feature transformer as **EFT** and present its color branch as a strong baseline.

We train the color branch of the EFT to minimize a simple reconstruction loss in Eq. 4.1, where \mathbf{r} is a query ray sampled from $\boldsymbol{\pi}$, C is the set of reference images and their relative poses, and $I(\mathbf{r})$ is the ground truth pixel value.

$$\mathcal{L}_{EFT} = \sum_{\mathbf{r} \in R(\boldsymbol{\pi})} \|g_\psi(\mathbf{r}, C) - I(\mathbf{r})\|^2 \quad (4.1)$$

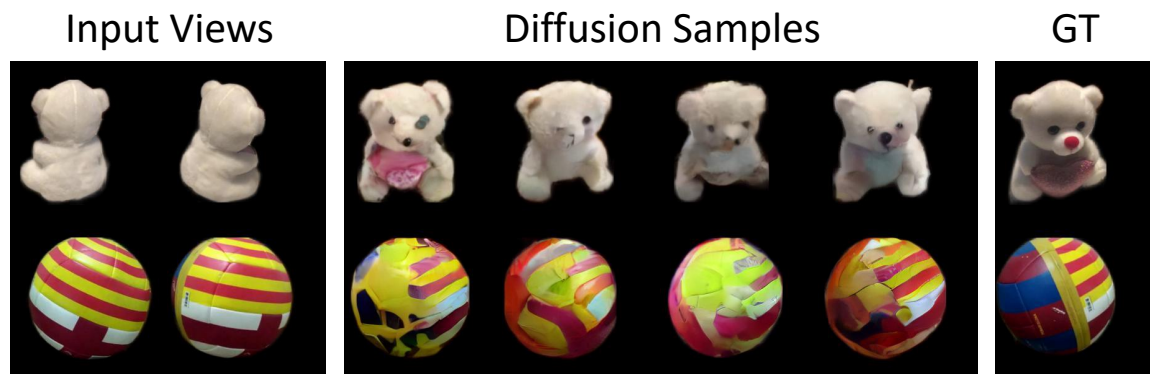


Figure 4.3: **Diffusion Samples**. Given the same input features, the reverse sampling process of diffusion model results in different predictions for unseen regions.

4.3 View-conditioned Latent Diffusion Model

While EFT can directly predict novel views, the pixelwise prediction mechanism does not allow it to model the underlying probability distribution, thus resulting in blurry mean-seeking predictions under uncertainty. To model the distribution over plausible images, we train a view-conditioned diffusion model to estimate $p(\mathbf{x}|\boldsymbol{\pi}, C)$ while using EFT as a geometric feature extractor. Instead of directly modeling the distribution in pixel space, we find it computationally efficient to do so in a lower-resolution latent space $\mathbf{z} = \mathcal{E}(\mathbf{x})$, which can be decoded back to an image as $\mathbf{x} = \mathcal{D}(\mathbf{z})$. Please see the appendix for details.

Given target view $\boldsymbol{\pi}$ and a set of input images C , we extract a 32 by 32 feature grid $\mathbf{y} = h_\psi(\boldsymbol{\pi}, C)$ using the EFT backbone. We train our VLDM to recover ground truth

4. Method

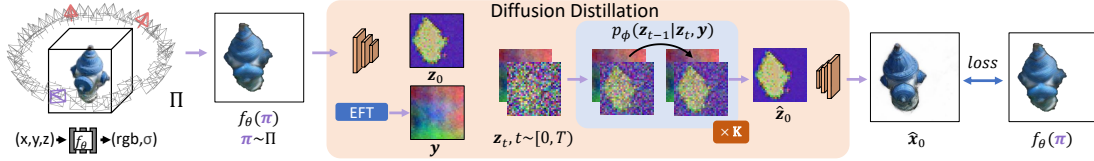


Figure 4.4: **Diffusion Distillation Diagram.** We optimize the parameters θ of an Instant NGP network such that rendered images $f_\theta(\boldsymbol{\pi})$ from $\boldsymbol{\pi} \sim \Pi$ are similar to VLDM predictions $\hat{\mathbf{x}}_0$, effectively seeking a mode in $p_\phi(\mathbf{x}|\boldsymbol{\pi}, C)$.

image latent \mathbf{z}_0 conditioned on \mathbf{y} . Following diffusion model training conventions [13, 23, 37], we optimize a simplified variational lower bound in Eq. 4.2.

$$\mathcal{L}_{VLDM} = \mathbb{E}_{\mathbf{z}, \epsilon \sim \mathcal{N}(0,1), t, \mathbf{y}} [\|\epsilon - \epsilon_\phi(\mathbf{z}_t, t, \mathbf{y})\|^2] \quad (4.2)$$

Figure 4.2 shows a diagram of the training setup. Our VLDM model allows us to approximate $p(\mathbf{x}|\boldsymbol{\pi}, C)$, and enables drawing multiple sample predictions. In Figure 4.3, we see variations in VLDM predictions. Nevertheless, all predictions are plausible explanations for the target view given that majority of it is unseen.

4.4 Extracting 3D Modes via Diffusion Distillation

While the proposed VLDM gives us the ability to hallucinate unseen regions and make realistic predictions under uncertainty, it does not output a 3D representation. In fact, as it models the distribution over images, the views sampled from the VLDM do not (and should not!) necessarily correspond to a single underlying 3D interpretation. How can we then obtain an output 3D representation while preserving the high-quality of renderings?

3D Inference as Neural Mode Seeking. Our key insight is that the VLDM model not only allows us to sample plausible novel views, but the modeled distribution also gives us a mechanism to approximate the likelihood of a generated novel view. Building on this insight, we propose to distill the VLDM predictions to obtain an instance-specific 3D neural scene representation f_θ , such as NeRF [19] or Instant NGP (INGP) [20]. Intuitively, we want to arrive at a solution for f_θ such that its

renderings $\mathbf{x} \equiv f_\theta(\boldsymbol{\pi})$ from arbitrary viewpoints $\boldsymbol{\pi}$ are likely under the conditional distribution modeled by the VLDM $p_\phi(\mathbf{x}|\boldsymbol{\pi}, C)$:

$$\min_{\theta} \mathbb{E}_{\boldsymbol{\pi} \sim \Pi} -\log p_\phi(f_\theta(\boldsymbol{\pi})|\boldsymbol{\pi}, C) \quad (4.3)$$

where we minimize the negative log-likelihood for images rendered with f_θ over cameras sampled from a prior camera distribution Π (constructed by assuming a circular camera trajectory and that all cameras look at a common center). We term this process as ‘neural mode seeking’ as it encourages a representation which maximizes likelihood as opposed to minimizing distance to samples (mean seeking).

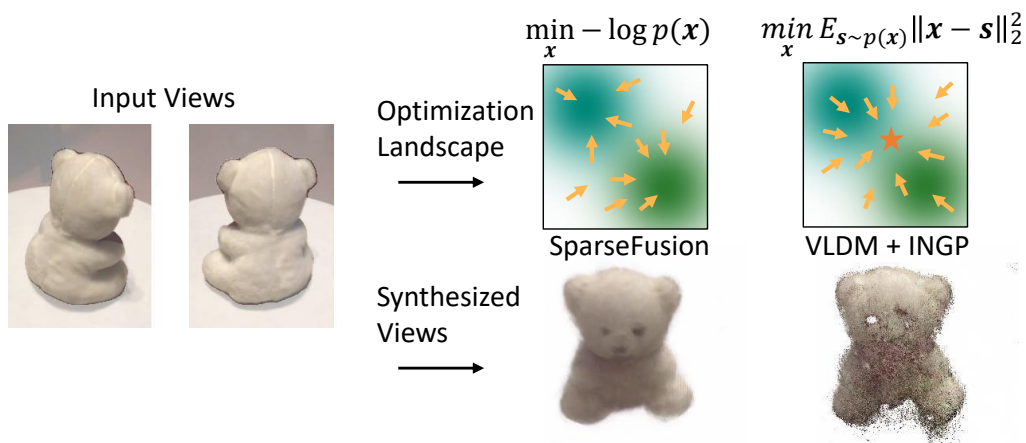


Figure 4.5: **Mode Seeking Visualization.** We show qualitative comparison between a mode-seeking (SparseFusion) and a mean-seeking (VLDM+INGP) objective.

Neural Mode Seeking via Diffusion Distillation. Given a learned diffusion model, the reconstruction objective yields a bound on the log-likelihood of a data point \mathbf{x} . This approximation yields a simple mechanism for computing the likelihood of a (rendered) image $f_\theta(\boldsymbol{\pi})$ to be used in the mode-seeking optimization (Eq. 4.3):

$$-\log p_\phi(\mathbf{x}_0) \approx \mathbb{E}_{\epsilon, t} [w_t \|\mathbf{z}_0 - \hat{\mathbf{z}}_{0,t}\|^2] + C \quad (4.4)$$

where $\mathbf{z}_0 = \mathcal{E}(f_\theta(\boldsymbol{\pi}))$ is the latent of the rendered image, $t \sim (0, T]$, and $\hat{\mathbf{z}}_{0,t}$ is the predicted latent. Intuitively, this objective implies that if, after adding noise to obtain \mathbf{z}_t from \mathbf{z}_0 , the denoising diffusion model predicts $\hat{\mathbf{z}}_0$ close to the original input, one has reached a mode under $p_\phi(\mathbf{z})$. We visualize the behavior of mode seeking versus

4. Method

mean seeking in Figure 4.5.

Multi-step Denoising and Image-space Reconstruction. In practice, we make three modifications to the single-step objective in Eq. 4.4 for better performance: 1) taking loss in pixel space instead of latent space *i.e.*, using \mathbf{x}_0 instead of \mathbf{z}_0 , 2) using perceptual distance [54] in addition to the pixelwise distance, and 3) performing multi-step denoising. Instead of directly predicting $\hat{\mathbf{z}}_{0,t}$, we adaptively use multiple time-steps (up to 50 steps) $\mathcal{T} = (t_1, \dots, t_k, t)$, and successively predict $\hat{\mathbf{z}}_{t_{k-1}, t_k}$ (via [17]) *i.e.*, predict a denoised estimate for time t_{k-1} given a sample from time t_k . We denote this reconstruction as $\hat{\mathbf{z}}_{0,\mathcal{T}}$ to highlight the multiple-step reconstruction. We express our final objective for optimizing for neural mode seeking with view-conditioned diffusion models as:

$$\mathcal{L} = \mathbb{E}_{\boldsymbol{\pi}, \epsilon, t} [w_t \|f_\theta(\boldsymbol{\pi}) - \hat{\mathbf{x}}_{0,\mathcal{T}}\|^2 + \mathcal{L}_{Perp}(f_\theta(\boldsymbol{\pi}), \hat{\mathbf{x}}_{0,\mathcal{T}})] \quad (4.5)$$

where $\hat{\mathbf{x}}_{0,\mathcal{T}} = \mathcal{D}(\hat{\mathbf{z}}_{0,\mathcal{T}})$, and $\hat{\mathbf{z}}_{0,\mathcal{T}}$ is the multi-step reconstruction from \mathbf{z}_t – which is obtained by adding noise to $\mathbf{z}_0 = \mathcal{E}(f_\theta(\boldsymbol{\pi}))$. While $\hat{\mathbf{z}}$ in the above objective does (indirectly) depend on the neural representation f_θ , we follow [22] in ignoring this dependence when computing parameter gradients (see [44] for a justification). We outline the multi-step denoising diffusion distillation in Figure 4.4.

Chapter 5

Experiments

We demonstrate our approach on a challenging real world multi-view dataset CO3Dv2 [24], across 51 diverse categories. First, we compare SparseFusion against prior works, highlighting the benefit of our approach in sparse view settings. Then, we show the importance of diffusion distillation and its probabilistic mode-seeking formulation.

5.1 Experimental Setup

Dataset. We perform experiments on CO3Dv2 [24], a multi-view dataset of real world objects annotated with relative camera poses and foreground masks. We use the specified *fewview-train* and *fewview-dev* splits for training and evaluation. Since SparseFusion optimizes an instance-specific Instant NGP, it is computationally prohibitive to evaluate on all evaluation scenes. Instead, we perform most experiments on a **core subset** of 10 categories proposed by [24], evaluating 10 scenes per category. Furthermore, we demonstrate that SparseFusion extends to diverse categories by evaluating 5 scenes per category across 51 categories.

Baselines. We compare SparseFusion against current state-of-the-art methods. We first compare against *PixelNeRF* [51], a feature re-projection method. We adapt *PixelNeRF* to CO3Dv2 dataset and train category-specific models on the 10 categories of the *core subset*, each for 300k steps. We also compare against *NerFormer* [24], another feature re-projection method. We use category-specific models provided by

5. Experiments

the authors for all 51 categories. Moreover, we compare against *ViewFormer*¹ [15], an autoregressive image generation method, using models provided by the authors. Lastly, we present components of SparseFusion, EFT and VLDM, as strong baselines.

Metrics. We report standard image metrics PSNR, SSIM, and LPIPS [54]. We recognize that no metric is perfect for ambiguous cases of novel view synthesis; PSNR derives from pixelwise MSE and favors mean color prediction while SSIM and LPIPS favor perceptual agreement.

Implementation Details. For EFT, we use a ResNet18 [11] backbone and three groups of transformer encoders with 4 layers each. We use 256 hidden dimensions for all layers. For VLDM, we freeze the VAE from [27] that encodes 256x256 images to 32x32 latents with channel dimension of 4. We construct a 400M parameter denoising UNet similar to [28, 30] for probabilistic modeling. We jointly train category-specific EFT and VLDM models, using Eq. 4.1 and Eq. 4.2, across all categories in CO3Dv2. We use a batch size of 2 and train for 100K iterations.

For diffusion distillation, we use a PyTorch implementation of Instant NGP [20, 39]. Due to memory constraints, we render images at 128x128 and upsample to 256x256 before performing diffusion distillation. For each instance, we optimize Instant NGP for 3,000 steps. During the first 1,000 steps, we optimize rendering loss on input images and predicted EFT images from a circular camera trajectory to initialize a rough volume. During the next 2,000 steps, we perform diffusion distillation. Reconstructing a single instance takes roughly an hour on an A5000 gpu.

¹Only category-agnostic CO3Dv1 weights are compatible with our evaluation. We use the 10-category weights for our *core subset* experiments and all-category weights for our all category experiments. Despite this difference, the comparative results of ViewFormer against our baselines are consistent with the comparisons reported in their original paper.

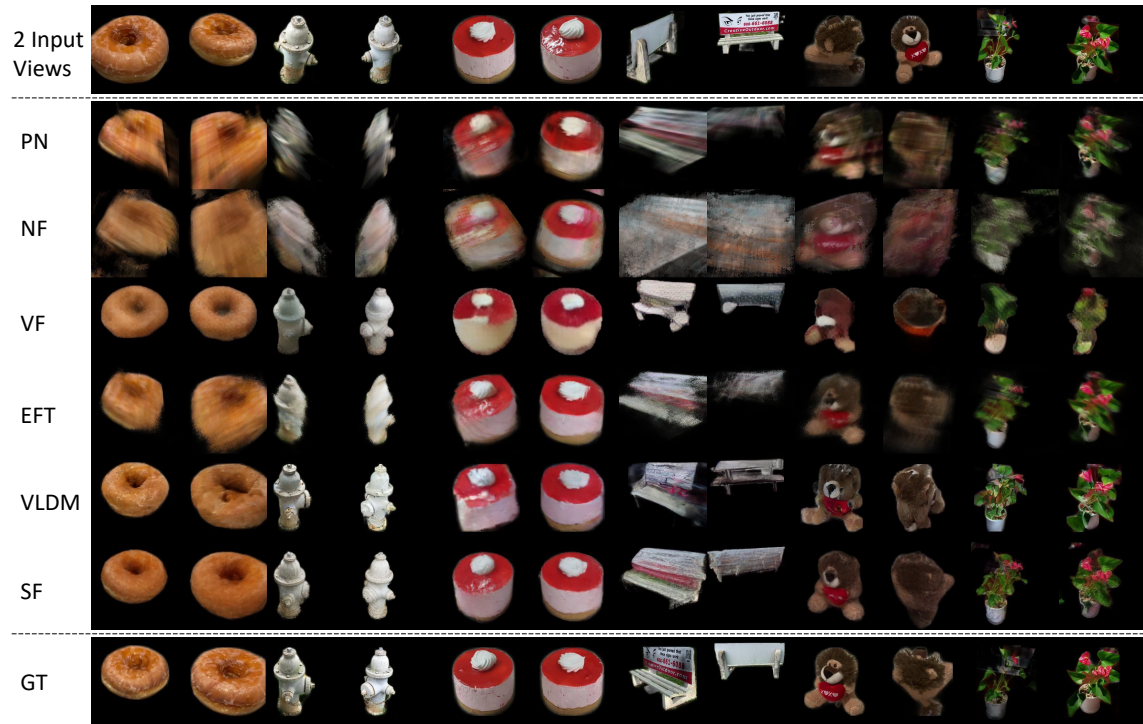


Figure 5.1: **View Synthesis Qualitative Results.** We show view synthesis results with 2 input views on donut, hydrant, cake, bench, teddybear, and plant categories. We visualize 2 novel views per instance with PixelNeRF (PN), NerFormer (NF), ViewFormer (VF), EFT, VLDM, and finally, SparseFusion (SF). Corresponding numbers can be found in Table 5.1.

5.2 Reconstruction on Real Images

Table 5.1: **Detailed View Synthesis Benchmark.** We show 2-view category-specific metrics on 10 CO3D categories from the *core subset*. We show PSNR \uparrow and LPIPS \downarrow averaged across 10 scenes per category.

	Domt		Apple		Hydrant		Vase		Cake		Ball		Bench		Suitcase		Teddybear		Plant	
	PSNR	LPIPS	PSNR	LPIPS	PSNR	LPIPS	PSNR	LPIPS	PSNR	LPIPS	PSNR	LPIPS	PSNR	LPIPS	PSNR	LPIPS	PSNR	LPIPS	PSNR	LPIPS
PixelNeRF [51]	20.9	0.30	20.0	0.35	19.0	0.27	21.3	0.26	18.3	0.37	18.5	0.36	17.7	0.35	21.7	0.30	18.5	0.35	19.3	0.36
NerFormer [24]	20.3	0.34	19.5	0.33	18.2	0.30	17.7	0.34	16.9	0.44	16.8	0.35	15.9	0.44	20.0	0.39	15.8	0.43	17.8	0.45
ViewFormer ¹ [15]	19.3	0.29	20.1	0.26	17.5	0.22	20.4	0.21	17.3	0.33	18.3	0.31	16.4	0.30	21.0	0.26	15.5	0.32	17.8	0.31
EFT	21.5	0.31	22.0	0.29	21.6	0.22	21.1	0.25	19.9	0.33	21.4	0.29	17.8	0.34	23.0	0.26	19.8	0.30	20.4	0.31
VLDM	20.1	0.25	21.3	0.22	20.1	0.18	20.2	0.20	18.9	0.30	20.3	0.25	16.6	0.29	21.3	0.23	17.9	0.27	18.9	0.27
SparseFusion	22.8	0.22	22.8	0.20	22.3	0.16	22.8	0.18	20.8	0.28	22.4	0.22	16.7	0.28	22.2	0.22	20.6	0.24	20.0	0.25

5.2.1 Core Subset: 2-view.

We show 2-view category-specific reconstruction results for the 10 *core subset* categories. We evaluate metrics on the first 10 scenes of each category. For each scene, we load 32 linearly spaced views, from which we randomly sample two input views and evaluate on the remaining 30 unseen views. The input and evaluation views are held constant across methods. We report category-specific PSNR and LPIPS in Table 5.1. We show qualitative comparisons in Figure 5.1.

SparseFusion outperforms all other methods in LPIPS, only losing out in PSNR for 3 categories. Despite PSNR favoring mean predicting methods, SparseFusion achieves higher PSNR in 7 categories. The strong performance of SparseFusion is reflected in the qualitative comparison. Existing methods either predict a blurry view for unseen regions or a perceptually reasonable view that disregards 3D consistency. SparseFusion predicts views that are both perceptually reasonable and geometrically consistent.

Table 5.2: **View Synthesis on 10 Categories.** We benchmark view synthesis results on the 10 categories with 2, 3, and 6 input views.

	2 Views			3 Views			6 Views		
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
PixelNeRF [51]	19.52	0.667	0.327	20.67	0.712	0.293	22.47	0.776	0.241
NerFormer [24]	17.88	0.598	0.382	18.54	0.618	0.367	19.99	0.661	0.332
ViewFormer ¹ [15]	18.37	0.697	0.282	18.91	0.704	0.275	19.72	0.717	0.266
EFT	20.85	0.680	0.289	22.71	0.747	0.262	24.57	0.804	0.210
VLDM	19.55	0.711	0.247	20.85	0.737	0.225	22.35	0.768	0.201
SparseFusion	21.34	0.752	0.225	22.35	0.766	0.216	23.74	0.791	0.200

5.2.2 Core Subset: Varying Views.

We examine performance of the different methods as we increase the number of input views. As the number of input views increases, more regions are observed, giving an advantage to methods that explicitly use feature re-projection. We evaluate 2, 3, and 6 view reconstruction on the *core subset* categories and show PSNR, SSIM, and LPIPS in Table 5.2.

We see feature re-projection methods improve drastically with more input views as the need for hallucination of unseen regions decreases. EFT outperforms SparseFusion in PSNR for the 3-view and 6-view settings. However, SparseFusion remains competitive in PSNR while being better in LPIPS. SSIM results further underscore the advantage of SparseFusion with sparse (2, 3) input views. Moreover, SparseFusion outperforms all current state-of-the-art methods in all three metrics for 2, 3, and 6 view reconstruction.

5. Experiments

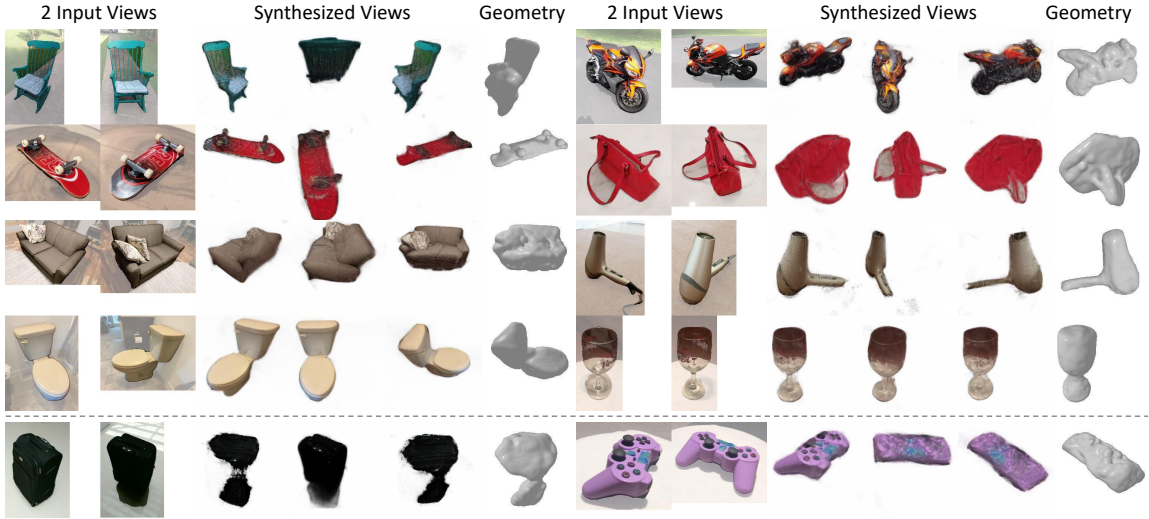


Figure 5.2: **Reconstruction on Diverse Categories.** We show SparseFusion reconstructions on a subset of the 51 CO3D categories. We also show a couple of failure modes on the last row. Please see project page for more samples and 360-degree visualizations.

Table 5.3: **View Synthesis on 51 Categories.** We benchmark novel view synthesis on all CO3D categories with 2 input views.

	2 Views		
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
NerFormer [24]	18.44	0.614	0.365
ViewFormer ¹ [15]	18.91	0.718	0.265
EFT	21.44	0.719	0.281
VLDL	19.85	0.732	0.229
SparseFusion	21.20	0.756	0.223

5.2.3 All Categories: 2-views.

We compare against NerFormer and ViewFormer across all 51 categories to demonstrate SparseFusion’s performance on diverse categories. We evaluate with 2 random input views on the first 5 scenes of each category for all 51 categories and report the averaged metrics in Table 5.3. While EFT edges out in PSNR, SparseFusion achieves better SSIM and LPIPS. Existing methods, NerFormer and ViewFormer perform

significantly worse. We show qualitative results of SparseFusion on diverse categories in Figure 5.2 where, in addition to 3 synthesized novel views, we also visualize the underlying geometry by extracting an iso-surface via marching cubes.

5.2.4 Failure Modes.

We show failure modes on the bottom row of Figure 5.2. On the bottom left, SparseFusion fails to reconstruct a good geometry for the black suitcase. As Instant NGP is trained to output a default black color for the background, the neural representation sometimes fails to disambiguate black foreground from black background. On the bottom right, we see SparseFusion propagating a dataset bias for the category, remote. Since most remote images are TV remotes, SparseFusion attempts to make the video game controller a TV remote.

5.3 Additional Analysis

5.3.1 Performance Binned by Viewpoint Changes.

We investigate the relationship between magnitude of viewpoint change and reconstruction performance. We analyze SparseFusion, EFT, and PixelNeRF results on the *core subset* and visualize PSNR and LPIPS binned by angle in degrees to the nearest context view in Figure 5.3. We show that for small viewpoint changes, SparseFusion performs better in LPIPS and competitively in PSNR against EFT. As viewpoint change increases, feature re-projection methods fall off quite fast while SparseFusion remains more robust and performs relatively better.

5.3.2 Importance of Mode Seeking.

We compare the diffusion distillation formulation against a naive method to obtain a neural representation given a view synthesis method (VLDM or EFT). Concretely, we obtain several rendered samples ($\{\hat{I}, \hat{\pi}\}$) from the base view synthesis method given the context views C , and simply train an INGP to fit a 3D representation to these.

We present the results in Table 5.4, and see no significant change when we fit INGP to EFT renderings because EFT predicts consistent mean outputs. However,

5. Experiments

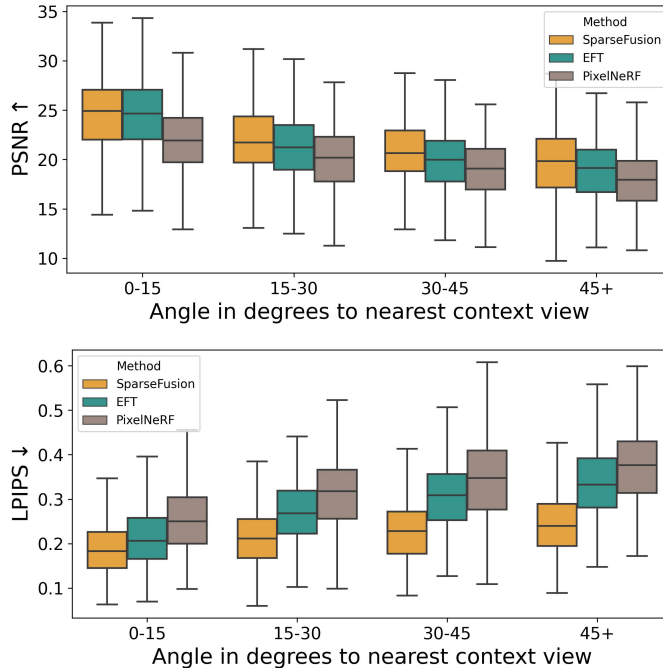


Figure 5.3: **Metrics Binned by Viewpoint Change.** We show metrics binned by the angle of query camera to the nearest context view. Results are aggregated from Table 5.1.

Table 5.4: **The Importance of Mode Seeking.** We show metrics when EFT and VLDM are naively used to optimize Instant NGP [20] in a mean seeking behavior, versus the mode seeking optimization in SparseFusion. We average across 10 scenes of hydrants with 2 input views.

Backbone	Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
EFT	base	21.58	0.732	0.224
	base w/ INGP	21.57	0.780	0.219
VLDM	base	20.05	0.776	0.178
	base w/ INGP	20.61	0.753	0.230
	SparseFusion	22.35	0.817	0.153

when we fit INGP to VLDM predictions, we see that perceptual quality decreases. We show a qualitative example in Figure 4.5 and also illustrate a toy 2D scenario which explains this drop due to mean seeking where averaging over conflicting samples leads to a poor reconstruction. However, when we optimize INGP using the diffusion distillation objective, all metrics improve, underscoring the importance of our proposed mode seeking optimization.

Table 5.5: **Diffusion Distillation Setup.** We show that a combination of multi-step prediction and perceptual loss strikes a balance between all three metrics. (hydrant, 10 scenes, 2 input views)

Loss Space	Denoising Steps	Perceptual Loss	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
Latent	Single	No	22.25	0.720	0.211
		Yes	22.15	0.770	0.187
	Multiple	No	21.92	0.744	0.211
		Yes	22.03	0.781	0.170
Pixel	Single	No	22.13	0.792	0.208
		Yes	22.49	0.826	0.169
	Multiple	No	22.36	0.797	0.200
		Yes	22.35	0.817	0.153

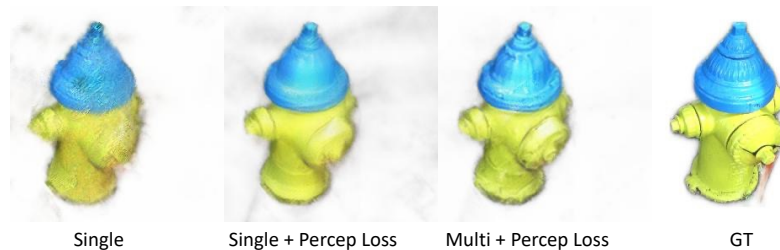


Figure 5.4: **Qualitative Results with Pixel Space Loss.** Using multi-step denoising and perceptual loss achieves more realistic results.

5.3.3 Ablating Distillation Objective.

We examine performance across various distillation design choices in Table 5.5. We observe that for all methods, PSNR remains relatively similar. However, computing loss in pixel space and additionally using perceptual loss improves both SSIM and LPIPS. Moreover, the multi-step denoising leads to the best perceptual results. While single-step denoising with perceptual loss achieves better PSNR and SSIM by a small margin, qualitative results in Figure 5.4 show that the predicted texture is smooth and unrealistic.

5. Experiments

Chapter 6

Discussion

6.1 Limitations

We presented an approach for inferring 3D neural representations from sparse-view observations. Unlike prior methods that struggled to deal with uncertainty, our approach allowed predicting 3D-consistent representations with plausible and realistic outputs even in unobserved regions. While we believe our work represents a significant step forward in recovering detailed 3D from casually captured images, a few challenges still remain. A key limitation of our work (as well as prior methods) is the reliance on known (relative) camera poses across the observations, and while there have been recent promising advances [25, 53], this remains a challenging task in general. Additionally, our approach requires optimizing instance-specific neural fields and is computationally expensive. Finally, while our work introduced the view-conditioned diffusion distillation in context of sparse-view reconstruction, we believe even single-view 3D prediction approaches can benefit from leveraging similar objectives.

6.2 Ethics and Broader Impact

Compared to existing novel view synthesis methods, SparseFusion is more computationally expensive. This poses a hardware limitation for potential downstream tasks and may also increase carbon emissions. Additionally, SparseFusion relies on

6. Discussion

view-conditioned latent diffusion models (VLDM), which are trained on multi-view datasets. VLDMs are good at representing their training data, potentially learning harmful biases that will propagate to reconstructed 3D scenes. While our current use case for reconstructing static objects from CO3D categories does not present ethical concerns, adapting SparseFusion to humans or animals requires more thorough examination of bias present in the training data.

Appendix A

Implementation Details

A.1 Epipolar Feature Transformer

Overview. Epipolar feature transformer is a feed-forward network that first gathers features along the epipolar lines of input images before aggregating them through a series of transformers. EFT is inspired by the GPNR approach by Suhail *et al.* [38], but we modify the feature extractor backbone to better suit the sparse-view setup and additionally use epipolar features for conditional diffusion. We describe our implementation below.

Notation: Let g_ψ be the RGB branch and h_ψ be the feature branch.

Inputs: $C \equiv (\mathbf{x}_m, \boldsymbol{\pi}_m)$, a set of input images with known camera poses and a query pose $\boldsymbol{\pi}$ – note that the poses are w.r.t. an arbitrary world-coordinate system and we only use their relative configuration.

Outputs: an RGB image \mathbf{x} and a feature grid \mathbf{y} corresponding to the query viewpoint $\boldsymbol{\pi}$.

Feature Extractor Backbone. Given input views $C \equiv (\mathbf{x}_m, \boldsymbol{\pi}_m)$ where \mathbf{x}_m is the m^{th} masked (black background) input image of shape (256, 256, 3). We use ResNet18 [11] as our backbone to extract pixel-aligned features by concatenating intermediate features from the first 4 layer groups of ResNet18, using bilinear upsampling to ensure all features are 128 by 128. For each image \mathbf{x}_m , we arrive at a feature grid of shape

(128, 128, 512).

Epipolar Points Projector. Given a query camera π , each pixel in its image plane corresponds to some ray. Our Epipolar Transformer seeks to infer per-pixel colors or features, and does so by processing each ray using the multi-view projections of points along it. For each ray \mathbf{r} (parameterized by its origin and direction), we project 20 points along the ray direction with depth values linearly spaced between z_{near} and z_{far} . We set z_{near} to $s - 5$ and z_{far} to $s + 5$ where s is the average distance from scene cameras to origin computed per scene. The 20 points, with shape (20, 3), are then projected into the screen space of each of the m input cameras, giving us epipolar points with shape (M, 20, 2). We use bilinear sampling to sample image features at the epipolar points, giving us combined epipolar features of shape (M, 20, 512) per ray. This becomes the input to our epipolar feature transformer.

Epipolar Feature Transformer. EFT aggregates the epipolar features from a single ray with a series of three transformers to predict an RGB pixel color and a 256-dimension feature. We visualize the EFT in Figure A.1. We show details of the transformers in Table A.1. All transformer encoders have hidden and output dimensions of 256. Both the depth aggregator and view aggregator transformers are followed by a weighted average operation, where the output features from the transformers are multiplied by a weight, which sums to 1 along the sequence length dimension. The relative weights are predicted by a linear layer before passing through softmax. This effectively performs weighted averaging along the sequence dimension.

The inputs to the transformer are the sampled features concatenated with additional ray and depth encodings. Given a point along the query ray \mathbf{r}_q at depth d , we denote by \mathbf{p}_{md} its projection in the m^{th} context view. In addition to the pixel-aligned feature \mathbf{f}_{md} (described in previous paragraph), we also concatenate encodings of the query ray \mathbf{r}_q , the depth \mathbf{d} , and the ray \mathbf{r}_{md} connecting the m^{th} camera center to the 3D point. We use plucker coordinates to represent each ray, and compute harmonic embeddings for each to $(\mathbf{r}_q, \mathbf{r}_{md}, \mathbf{d})$ (using 6 harmonic functions) before concatenating them with \mathbf{f}_{md} to form the input tokens to the transformer.

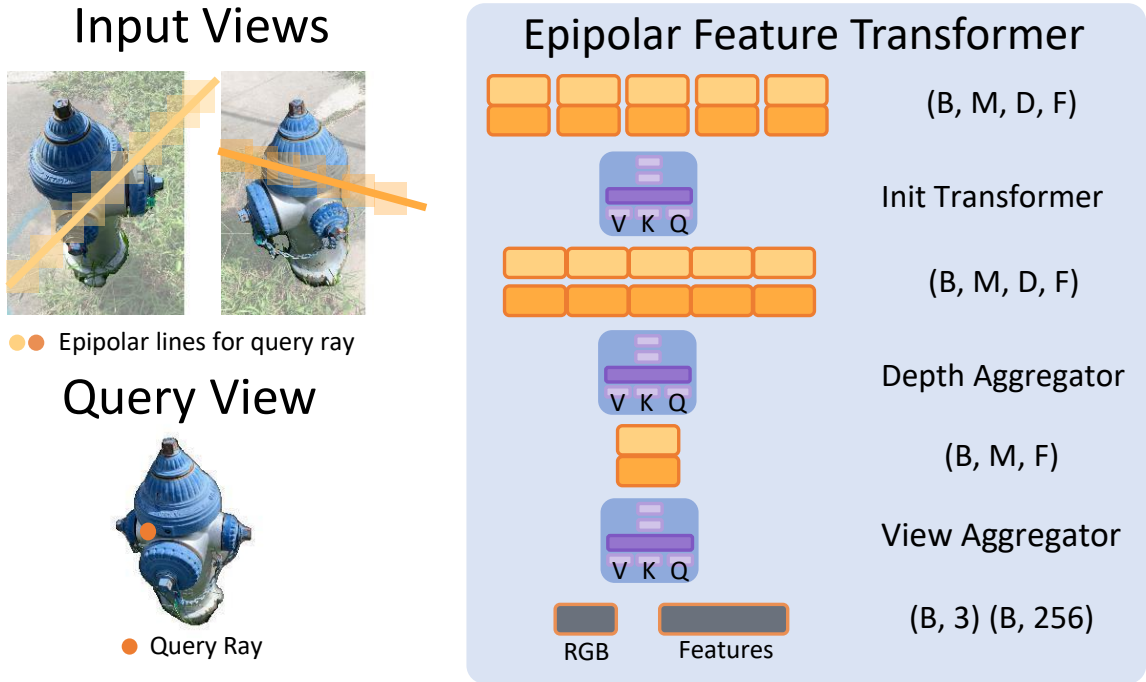


Figure A.1: **Epipolar Feature Transformer** We show a diagram of EFT. This module processes each query ray independently, using a transformer to aggregate the projected features across views and across possible depths. For each ray, the output is a predicted RGB color (used as a baseline prediction method), and a pixel-aligned feature (used as conditioning in the diffusion model).

Training Procedure. We can train the color branch of EFT as a standalone novel view synthesis baseline. In our work, EFT is jointly trained with VLDM. Please see supplementary Section A.2 for details.

A.2 View-conditioned Diffusion Model

Overview. View-conditioned diffusion model is a latent diffusion model that conditions on a pixel-aligned feature grid \mathbf{y} .

Notation: Let ϵ_ϕ be the denoising UNet, \mathcal{E} be the VAE encoder, and \mathcal{D} be the VAE decoder.

VAE. We use the VAE from Stable Diffusion [27]. We use the provided v1-3 weights and keep the VAE frozen for all experiments. We use $(256, 256, 3)$ RGB images as

A. Implementation Details

Table A.1: **EFT Configuration.** We use default PyTorch hyperparameters for each layer. B is number of rays. M is the number of input views. D is the number of epipolar feature samples along the ray. D is 20.

Transformer	Layers	Sequence Dims / Dim	Output Shape
Init Transformer	Transformer Encoder x4	M	(B, M, D, 256)
Depth Aggregator	Transformer Encoder x4	D	(B, M, D, 256)
	Linear + Softmax	D	(B, M, D, 1)
	Weighted Average		(B, M, 256)
View Aggregator	Transformer Encoder x4	M	(B, M, 256)
	Linear + Softmax	M	(B, M, 1)
	Weighted Average		(B, 256)
Color Branch	Linear		(B, 3)

input, and the VAE encodes them into latents of shape (32, 32, 4). We refer readers to [27] for more details.

Denosing UNet. Our 400M parameter UNet roughly follows [30]. We construct our UNet using code from [45] with the parameters in Table A.2.

Table A.2: **UNet Parameters.** We provide parameters for our UNet.

Parameter	Value
channels	4
dim	256
dim_mults	(1,2,4,4)
num_resnet_blocks	(2,2,2,2)
layer_attns	(False, False, False, True)
cond_images_channels	256

The UNet comprises of 4 down-sampling blocks, a middle block, and 4 up-sampling blocks. We show the input and output shape for the modules of the UNet in Table A.3. We refer readers to [45] for UNet details. We disable all text conditioning and cross attention mechanisms; instead, we concatenate EFT features, \mathbf{y} , with image latents, \mathbf{z}_t . These EFT features are computed for the of 32×32 rays corresponding to the patch centers.

Training Procedure. We train with batch size of 2, randomly chosen number of input views between 2-5, and learning rate of 5e-5 using Adam optimizer with default

Table A.3: **UNet Blocks.** We outline the modules in our denoising UNet.

Modules	Block	Output Shape
Input		(B, 260, 32, 32)
Init. Conv	InitBlock	(B, 256, 32, 32)
Down 1	DownBlock	(B, 256, 16, 16)
Down 2	DownBlock	(B, 512, 8, 8)
Down 3	DownBlock	(B, 1024, 4, 4)
Down 4	DownBlock Self-attention	(B, 1024, 4, 4) (B, 1024, 4, 4)
Middle	MiddleBlock	(B, 1024, 4, 4)
Up 1	UpBlock Self-attention	(B, 1024, 8, 8) (B, 1024, 8, 8)
Up 2	UpBlock	(B, 512, 16, 16)
Up 3	UpBlock	(B, 256, 32, 32)
Up 4	UpBlock	(B, 256, 32, 32)
Final Conv.	Conv2D	(B, 4, 32, 32)

hyperparameters for 100K steps. We optimize both the UNet weights and also the EFT weights. We optimize the UNet and feature branch of EFT with the simplified variational lower bound [13]. We optimize the color branch of EFT with pixel-wise reconstruction loss.

A.3 Diffusion Distillation

Overview. We optimize a 3D neural scene representation, Instant NGP [20, 39], with our VLDM.

Notation: Let f_θ be the volumetric Instant NGP renderer, $p_\phi(\mathbf{z}_{0:T}|\boldsymbol{\pi}, C)$ be the multi-step denoising process that estimates $\hat{\mathbf{z}}_0$. Let Π be an instance-specific camera distribution.

Instant NGP. We use the PyTorch Instant NGP implementation from [39]. We set scene bounds to 4 with desired hashgrid resolution of 8,192. We use a small 3 layer MLP with hidden dimension of 64 to predict RGB and density. We do not use

A. Implementation Details

view direction as input.

Camera Distribution. Given a set of input cameras $C_I \equiv (\boldsymbol{\pi}_m)$ and a query camera $\boldsymbol{\pi}_q$, we first find the look-at point P_{at} by finding the nearest point to all $m + 1$ rays originating from camera centers. Then, we fit a circle O in 3D space with center being the mean of all camera centers. Let the normal of circle O be \mathbf{n} . To sample a camera, we first sample a point P_i on O and jitter the angle between $\overline{P_{at}P_i}$ and \mathbf{n} by $\mathcal{N}(0, 0.17)$ radians to get jittered point P'_i . We then construct a camera $\boldsymbol{\pi}$ with center P'_i looking at P_{at} .

Multi-step Diffusion Guidance. Given a rendered image \mathbf{x}_0 , we encode it to obtain \mathbf{z}_0 . Then, we uniformly sample $t \sim (0, T]$ and construct a noisy image latent \mathbf{z}_t . We perform multi-step denoising to obtain $\hat{\mathbf{z}}_0$ by iteratively sampling $\hat{\mathbf{z}}_{t_{k-1}} \sim p_\phi(\mathbf{z}_{t_{k-1}} | \hat{\mathbf{z}}_{t_k}, y)$ on an interval of time steps $\mathcal{T} = (t_1, \dots, t_k, t)$ using a linear multi-step method [17]. We construct \mathcal{T} by linearly spacing $k + 1$ time steps between $(0, t]$. We define k with a simple scheduler:

$$k + 1 = \begin{cases} \frac{100t}{T}, & \text{if } t \leq \frac{T}{2} \\ 50, & \text{if } t > \frac{T}{2} \end{cases} \quad (\text{A.1})$$

Finally, given $\hat{\mathbf{z}}_0$, we get the predicted image $\hat{\mathbf{x}}_0 = \mathcal{D}(\hat{\mathbf{z}}_0)$. We do not compute gradients through multi-step diffusion and treat $\hat{\mathbf{x}}_0$ as a detached tensor.

Distillation Details. We perform 3,000 steps of distillation, optimizing weights of the MLP θ with Adam optimizer and learning rate 5e-4. During each step of diffusion distillation, we sample $\boldsymbol{\pi} \sim \Pi$ and render an image $\mathbf{x}_0 = f_\theta(\boldsymbol{\pi})$. For the first 1,000 steps, we compute rendering loss between $f_\theta(\boldsymbol{\pi})$ and $g_\psi(\boldsymbol{\pi} | C)$. During the remaining steps, we compute loss between $f_\theta(\boldsymbol{\pi})$ and $\hat{\mathbf{x}}_0$ and use weighting $w_t = 1 - \bar{\alpha}_t$. To avoid out-of-memory error, we render images at reduced resolution (128, 128) and apply bilinear up-sampling before performing multi-step diffusion. In addition, we compute rendering loss between $f_\theta(\boldsymbol{\pi}_m)$ and \mathbf{x}_m on all m input images. Optimizing a single scene takes roughly 1 hour on an A5000 GPU.

Bibliography

- [1] 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 *CVPR*, 2023. [2.2](#)
- [2] Ang Cao, Chris Rockwell, and Justin Johnson. Fwd: Real-time novel view synthesis with forward warping and depth. In *CVPR*, 2022. [2.1](#)
- [3] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. In *ECCV*, 2022. [2.1](#)
- [4] Christopher B Choy, Danfei Xu, JunYoung Gwak, Kevin Chen, and Silvio Savarese. 3d-r2n2: A unified approach for single and multi-view 3d object reconstruction. In *ECCV*, 2016. [2.1](#)
- [5] Congyue Deng, Chiyu Jiang, Charles R Qi, Xinchun Yan, Yin Zhou, Leonidas Guibas, Dragomir Anguelov, et al. Nerdi: Single-view nerf synthesis with language-guided diffusion as general image priors. In *CVPR*, 2023. [2.2](#)
- [6] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *CVPR*, 2022. [2.1](#)
- [7] Yasutaka Furukawa and Jean Ponce. Accurate, dense, and robust multiview stereopsis. *TPAMI*, 2009. [2.1](#)
- [8] Rohit Girdhar, David F Fouhey, Mikel Rodriguez, and Abhinav Gupta. Learning a predictable and generative vector representation for objects. In *ECCV*, 2016. [2.1](#)
- [9] Georgia Gkioxari, Jitendra Malik, and Justin Johnson. Mesh r-cnn. In *ICCV*, 2019. [2.1](#)
- [10] Shubham Goel, Georgia Gkioxari, and Jitendra Malik. Differentiable stereopsis: Meshes from multiple views using differentiable rendering. In *CVPR*, 2022. [2.1](#)
- [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. [4.2](#), [5.1](#), [A.1](#)
- [12] Philipp Henzler, Niloy J Mitra, and Tobias Ritschel. Escaping plato’s cave: 3d

- shape from adversarial rendering. In *ICCV*, 2019. 2.1
- [13] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *NeurIPS*, 2020. 2.1, 3.1, 3.1, 4.3, A.2
- [14] Angjoo Kanazawa, Shubham Tulsiani, Alexei A Efros, and Jitendra Malik. Learning category-specific mesh reconstruction from image collections. In *ECCV*, 2018. 2.1
- [15] Jonáš Kulháněk, Erik Derner, Torsten Sattler, and Robert Babuška. Viewformer: Nerf-free neural rendering from few images using transformers. In *ECCV*, 2022. ??, 2.1, 5.1, 5.1, 5.2, 5.3
- [16] Chen-Hsuan Lin, Chaoyang Wang, and Simon Lucey. Sdf-srn: Learning signed distance 3d object reconstruction from static images. In *NeurIPS*, 2020. 2.1
- [17] Luping Liu, Yi Ren, Zhijie Lin, and Zhou Zhao. Pseudo numerical methods for diffusion models on manifolds. In *ICLR*, 2022. 4.4, A.3
- [18] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *CVPR*, 2019. 2.1
- [19] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020. 2.1, ??, 4.4
- [20] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Trans. Graph.*, 2022. (document), 2.1, 4.1, 4.4, 5.1, 5.4, A.3
- [21] Michael Niemeyer, Jonathan T Barron, Ben Mildenhall, Mehdi SM Sajjadi, Andreas Geiger, and Noha Radwan. Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs. In *CVPR*, 2022. 2.1, ??
- [22] Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. In *ICLR*, 2023. 2.1, 2.2, 4.4
- [23] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *URL: [https://doi.org/10.48550/arXiv, 2204](https://doi.org/10.48550/arXiv.2204.2024)*, 2022. 2.1, 4.3
- [24] Jeremy Reizenstein, Roman Shapovalov, Philipp Henzler, Luca Sbordone, Patrick Labatut, and David Novotny. Common objects in 3d: Large-scale learning and evaluation of real-life 3d category reconstruction. In *ICCV*, 2021. ??, 2.1, 5, 5.1, 5.1, 5.1, 5.2, 5.3
- [25] Chris Rockwell, Justin Johnson, and David F. Fouhey. The 8-point algorithm as an inductive bias for relative pose prediction by vits. In *3DV*, 2022. 6.1
- [26] Robin Rombach, Patrick Esser, and Björn Ommer. Geometry-free view synthesis:

- Transformers and no 3d priors. In *ICCV*, 2021. [2.1](#)
- [27] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. [5.1](#), [A.2](#)
- [28] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, 2015. [5.1](#)
- [29] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. *arXiv preprint arxiv:2208.12242*, 2022. [2.1](#)
- [30] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. In *NeurIPS*, 2022. [2.1](#), [5.1](#), [A.2](#)
- [31] 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 *CVPR*, 2022. [??](#), [2.1](#)
- [32] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *CVPR*, 2016. [2.1](#)
- [33] Steven M Seitz, Brian Curless, James Diebel, Daniel Scharstein, and Richard Szeliski. A comparison and evaluation of multi-view stereo reconstruction algorithms. In *CVPR*, 2006. [2.1](#)
- [34] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. *NeurIPS*, 2019. [2.1](#)
- [35] Vincent Sitzmann, Semon Rezkikov, Bill Freeman, Josh Tenenbaum, and Fredo Durand. Light field networks: Neural scene representations with single-evaluation rendering. In *NeurIPS*, 2021. [??](#), [2.1](#)
- [36] Noah Snavely, Steven M Seitz, and Richard Szeliski. Modeling the world from internet photo collections. *IJCV*, 2008. [2.1](#)
- [37] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *ICLR*, 2021. [2.1](#), [3.1](#), [3.1](#), [4.3](#)
- [38] Mohammed Suhail, Carlos Esteves, Leonid Sigal, and Ameesh Makadia. Generalizable patch-based neural rendering. In *ECCV*, 2022. [??](#), [2.1](#), [4.1](#), [4.2](#), [A.1](#)
- [39] Jiaxiang Tang. Torch-ngp: a pytorch implementation of instant-ngp, 2022.

- <https://github.com/ashawkey/torch-ngp>. ([document](#)), [4.1](#), [5.1](#), [A.3](#), [A.3](#)
- [40] Maxim Tatarchenko, Alexey Dosovitskiy, and Thomas Brox. Multi-view 3d models from single images with a convolutional network. In *ECCV*, 2016. [2.1](#)
 - [41] Alex Trevithick and Bo Yang. Grf: Learning a general radiance field for 3d representation and rendering. In *ICCV*, 2021. [2.1](#)
 - [42] Shubham Tulsiani, Tinghui Zhou, Alexei A Efros, and Jitendra Malik. Multi-view supervision for single-view reconstruction via differentiable ray consistency. In *CVPR*, 2017. [2.1](#)
 - [43] Kalyan Alwala Vasudev, Abhinav Gupta, and Shubham Tulsiani. Pre-train, self-train, distill: A simple recipe for supersizing 3d reconstruction. In *CVPR*, 2022. [2.1](#)
 - [44] Haochen Wang, Xiaodan Du, Jiahao Li, Raymond A Yeh, and Greg Shakhnarovich. Score jacobian chaining: Lifting pretrained 2d diffusion models for 3d generation. In *CVPR*, 2023. [2.2](#), [4.4](#)
 - [45] Phil Wang. Implementation of imagen, google’s text-to-image neural network, in pytorch, 2022. <https://github.com/lucidrains/imagen-pytorch>. [A.2](#), [A.2](#)
 - [46] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul P Srinivasan, Howard Zhou, Jonathan T Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. Ibrnet: Learning multi-view image-based rendering. In *CVPR*, 2021. [??](#), [2.1](#)
 - [47] Daniel Watson, William Chan, Ricardo Martin-Brualla, Jonathan Ho, Andrea Tagliasacchi, and Mohammad Norouzi. Novel view synthesis with diffusion models. In *ICLR*, 2023. [2.2](#)
 - [48] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Basri Ronen, and Yaron Lipman. Multiview neural surface reconstruction by disentangling geometry and appearance. In *NeurIPS*, 2020. [2.1](#)
 - [49] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces. In *NeurIPS*, 2021. [2.1](#), [??](#)
 - [50] Yufei Ye, Shubham Tulsiani, and Abhinav Gupta. Shelf-supervised mesh prediction in the wild. In *CVPR*, 2021. [2.1](#)
 - [51] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *CVPR*, 2021. [??](#), [2.1](#), [5.1](#), [5.1](#), [5.2](#)
 - [52] Jason Zhang, Gengshan Yang, Shubham Tulsiani, and Deva Ramanan. Ners: Neural reflectance surfaces for sparse-view 3d reconstruction in the wild. *NeurIPS*, 2021. [2.1](#), [??](#)
 - [53] Jason Y. Zhang, Deva Ramanan, and Shubham Tulsiani. RelPose: Predicting probabilistic relative rotation for single objects in the wild. In *ECCV*, 2022. [6.1](#)

- [54] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. [4.4](#), [5.1](#)
- [55] Tinghui Zhou, Shubham Tulsiani, Weilun Sun, Jitendra Malik, and Alexei A Efros. View synthesis by appearance flow. In *ECCV*, 2016. [2.1](#)