

iFusion: Inverting Diffusion for Pose-Free Reconstruction from Sparse Views

Chin-Hsuan Wu^{*1} Yen-Chun Chen² Bolivar Solarte¹ Lu Yuan² Min Sun^{1,3}

¹National Tsing Hua University ²Microsoft ³Amazon

{chinhsuanwu, enrique.solarte.nthu}@gapp.nthu.edu.tw

{yen-chun.chen, luyuan}@microsoft.com sunmin@ee.nthu.edu.tw

[chinhsuanwu.github.io/ifusion](https://github.com/chinhsuanwu/ifusion)



Figure 1. **Demonstration on real-world 3D reconstruction.** With only two casually taken photos **without camera poses**, *iFusion* can reconstruct plausible 3D assets. The top row example is taken from DreamBooth3D [51], and we took photos for the cat statue by ourselves.

Abstract

We present *iFusion*, a novel 3D object reconstruction framework that requires only two views with unknown camera poses. While single-view reconstruction yields visually appealing results, it can deviate significantly from the actual object, especially on unseen sides. Additional views improve reconstruction fidelity but necessitate known camera poses. However, assuming the availability of pose may be unrealistic, and existing pose estimators fail in sparse-view scenarios. To address this, we harness a pre-trained novel view synthesis diffusion model, which embeds implicit knowledge about the geometry and appearance of diverse objects. Our strategy unfolds in three steps: (1) We invert the diffusion model for camera pose estimation instead of synthesizing novel views. (2) The diffusion model is fine-tuned using provided views and estimated poses, turned into a novel view synthesizer tailored for the target object. (3) Leveraging registered views and the fine-tuned diffusion model, we reconstruct the 3D object. Experiments demonstrate strong performance in both pose estimation and novel view synthesis. Moreover, *iFusion* seamlessly integrates with various reconstruction methods and enhances them.

1. Introduction

Reconstructing objects from sparse views poses a significant challenge yet holds paramount importance for various applications, including 3D content creation, augmented reality, virtual reality, and robotics. Recent breakthroughs, guided by pre-trained models, have facilitated visually plausible reconstructions from a single view, without requiring the camera pose [30, 31, 36, 48, 64, 65, 79]. However, the reconstructed assets might not precisely capture the actual objects due to the inherent single-view ambiguity, *e.g.*, the object’s side opposite to the camera can only be imagined. Furthermore, multiple potential 3D structures could correspond to the same input image.

On the other hand, sparse-view methods assume the availability of an accurate camera pose for each view [3, 17, 24, 33, 63, 77, 86]. To meet this requirement, a Structure-from-Motion (SfM) pre-processing, *e.g.*, COLMAP [55], is typically employed. Paradoxically, these methods demand a substantial number of images, usually more than 50 in practice, for reliable pose estimation. Recent learning-based pose estimation [27, 58, 82] and pose-free reconstruction [19, 20] have sought to alleviate this issue. However, they still require a minimum of five input views and are primarily demonstrated on objects with simple 3D geometry

^{*}Part of this work was done as a research intern at Microsoft.

or within a constrained set of object categories. A generic framework for pose-free, sparse-view 3D reconstruction is still lacking, posing a significant obstacle to real-world applications with casually captured photos. We hereby raise the research question: How can one utilize only *extremely sparse views without poses* while maintaining the *reconstruction fidelity* of diverse objects?

The key is a sparse-view pose estimator. Our motivation stems from a recent novel view synthesis diffusion model, namely Zero123 [31], which is pre-trained on the most extensive 3D object dataset to date [8]. Given a reference view image, Zero123 can generate a novel view (query view) from a specified pose (Fig. 2, left). This indicates that the model has learned rich prior knowledge about the geometry and appearance of diverse objects. We thus hypothesize that it can be leveraged for pose estimation, with an intuition that a well-estimated pose fed into Zero123 will produce an image similar to the query view, and vice versa. Next, gradients may be back-propagated to optimize the pose with a proper loss function. Following this idea, we repurpose Zero123 by inverting it to take the two views and estimate the relative camera transformation (Fig. 2, right). More specifically, we adopt an analysis-by-synthesis paradigm [7, 45, 78] that optimizes the transformation by minimizing the difference between the denoised latent visual features, *i.e.*, Zero123’s output image feature map, and the query view’s feature. Empirically, the proposed approach achieves strong pose estimation with as few as 2 views, even outperforming existing approaches’ results with 5 views.

Well-estimated poses also open up a new opportunity. Using the given views registered with poses, a mini-dataset can be constructed to further fine-tune Zero123 and customize the diffusion model for synthesizing the target object’s novel views. Specifically, we can form a set of (reference view, camera pose, query view) triplets from the given sparse views and fine-tune Zero123. To accelerate training and prevent overfitting, we use Low-Rank Adaptation (LoRA) [15] to fine-tune the diffusion model, a recognized technique for customizing diffusion models.¹ Experiments demonstrate that this step significantly improves novel view synthesis, achieving an average increase of **+3.6** in PSNR across two datasets, and is beneficial to the final reconstruction. Note that our approach shares a similar spirit with test-time training [62], test-time adaptation [66], and self-training [56, 74]. Like test-time training and adaptation, we align the model to the test distribution based on test inputs (given views) but without test labels (novel views). Analogous to self-training, we synthesize additional labels (camera poses) using the learning model itself. To the best of our knowledge, the above combination is new for diffusion-based 3D reconstruction.

¹<https://github.com/clonofsimo/lor>

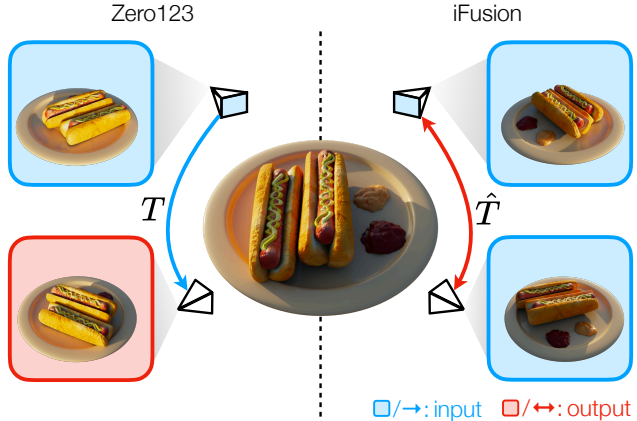


Figure 2. **Zero123 vs. iFusion.** Unlike Zero123 [31] (left), which synthesizes an object’s novel view given an image and a transformation T , *iFusion* (right) instead optimizes an unknown relative transformation \hat{T} from two given views.

To this end, we introduce *iFusion*, a novel framework that reconstructs diverse 3D objects with sparse, pose-free views. First, the pose estimation is achieved by inverting the Zero123 *diffusion* model, as described earlier. With the estimated camera pose, an object-specific improvement on Zero123’s novel view synthesis capability is performed, which can be further utilized as additional reconstruction guidance. Finally, for reconstructing the 3D asset, any differentiable renderer can be plugged in, including NeRFs [38] and the recently proposed 3D Gaussian Splatting [23]. It is noteworthy that our framework does not assume any specific reconstruction pipeline, and experimental results demonstrate that *iFusion* is readily applicable to four different single-view reconstruction methods. Improved geometric fidelity is observed with a significant **+7.2%** increase in volume IoU, showcasing the necessity of additional views for reliable 3D reconstruction.

- Our contributions are summarized as follows:
1. We propose a novel camera pose estimator that significantly outperforms existing methods in terms of both accuracy and required number of input views, while being effective for diverse objects.
 2. A self-training and test-time training inspired fine-tuning stage is innovated. This stage results in a much stronger novel view synthesis diffusion model, which plays a crucial role in guiding the reconstruction process.
 3. For the first time, we escalate diffusion-based single-view reconstruction to multi-view for enhanced fidelity with merely two pose-free images.

2. Preliminary

iFusion repurposes a novel view synthesizing diffusion model for camera pose prediction. To prepare readers with the necessary backgrounds, we briefly introduce the basics

of Diffusion Models (DM) and how they can be used for novel view synthesis. Next, we summarize a recently popular approach to utilize DM for 3D reconstruction, which we integrate into *iFusion* to allow reconstruction.

Diffusion Models Diffusion models [14, 59, 61] are a class of deep generative models that has become the mainstream approach for high-fidelity visual synthesis. In image generation, they work by “diffusing” an image by adding noise over repeated steps, and then a deep neural network is trained to predict the applied step-wise noise from a corrupted image. This allows the reversion of the diffusion process, thus an image can be generated from a random noise by iterative denoising using the trained noise predicting network. More specifically, Ho *et al.* [14] formulated the diffusion process in the following analytical form:

$$x_t = \sqrt{\alpha_t}x_0 + \sqrt{1 - \alpha_t}\epsilon, \quad t \in [0, 1, \dots, \mathcal{T}], \quad (1)$$

where $\epsilon \sim \mathcal{N}(0, 1)$ denotes the Gaussian noise and hyper-parameter α_t denotes the noise schedule. For the reverse process, the noise predictor is denoted as $\epsilon_\theta(x_t, t)$, where θ is the set of trainable parameters. Instead of directly modeling the RGB pixel values x , a widely used diffusion model, Stable Diffusion (SD),² applies the Latent Diffusion Model (LDM) [52] to model the latent feature maps z . The encoding and reconstruction of images is done via a pre-trained VQ-VAE: $z = \mathcal{E}(x)$, and $x = \mathcal{D}(z)$. Moreover, DM may optionally take conditional inputs c , *e.g.*, texts, bounding box layouts, and depth maps. For instance, the standalone SD takes texts as the condition c and enables text-to-image generation (T2I). Formally, the training loss of the prediction network can be written as:

$$\mathcal{L}(x, c) = \mathbb{E}_{z, \epsilon, t} \left[\|\epsilon - \epsilon_\theta(z_t, t, c)\|_2^2 \right], \quad (2)$$

where $\|\cdot\|_2$ denotes the L2 norm.

Diffusion Models for Novel View Synthesis The original Stable Diffusion was trained on web-scale image-text pairs³ for text-to-image generation. Recently, Liu *et al.* [31] proposed Zero123 to further fine-tune SD on Objaverse [8], a large-scale 3D assets dataset, for object-centric novel view synthesis. Given an image at the reference viewpoint x^r and the reference-to-query transformation $T_{r \rightarrow q} \in \text{SE}(3)$, the model synthesizes the desired query view x^q with condition $c(x^r, T_{r \rightarrow q})$. This is formulated as a DM and shares the same training objective as Eq. (2).

3D Reconstruction via Score Distillation Sampling Recent studies [18, 39, 47, 67] indicated that large-scale pre-trained 2D vision models [50, 52, 54] implicitly encapsulate rich 3D geometric prior. Notably, DreamFusion [47]

introduced the Score Distillation Sampling (SDS) to facilitate 3D generation guided by a pre-trained 2D DM. Let $x = \mathcal{R}_\psi(T)$ be the rendered image at viewpoint $T \in \text{SE}(3)$, where \mathcal{R} is a differentiable renderer parameterized by ψ , *e.g.*, Neural Radiance Fields (NeRFs) [38] or 3D Gaussian Splatting [23]. Given a denoising network ϵ_θ , SDS optimizes the renderer ψ by minimizing the residuals between the predicted noise and the added noise, thereby producing the gradients:

$$\nabla_\psi \mathcal{L}_{SDS}(x, c) = \mathbb{E}_{z, \epsilon, t} \left[(\epsilon_\theta(z_t, t, c) - \epsilon) \frac{\partial z}{\partial \psi} \right]. \quad (3)$$

3. Method

Figure 3 presents an overview of the *iFusion* framework. The key of our pose-free reconstruction framework is the sparse-view pose estimator shown in Fig. 3 (a). By inverting the diffusion model, accurate poses can be estimated. Next, the registered views are leveraged to customized the novel view synthesis model for the target object as in Fig. 3 (b). Finally, 3D reconstruction can be done using the registered views, and the customized diffusion model serves as the guidance, shown in Fig. 3 (c).

3.1. Diffusion as a Pose Estimator

The goal is to recover the relative camera pose $T_{r \rightarrow q}$ from a reference view x^r to the query view x^q , leveraging the pre-trained diffusion model ϵ_θ . Intuitively, a model trained for a task involving camera poses could potentially be used in reverse: to retrieve or estimate the camera pose from given inputs, as evident in Chen *et al.* [7], Park *et al.* [45], Yen-Chen *et al.* [78]. Hence, rather than optimizing DM parameters θ to reconstruct x^q given $c(x^r, T_{r \rightarrow q})$ as in the training stage shown in Eq. (2), we solve the inverse problem by freezing θ and optimizing $\hat{T}_{r \rightarrow q}$ to reconstruct x^q :

$$\hat{T}_{r \rightarrow q} = \underset{T \in \text{SE}(3)}{\text{argmin}} \mathcal{L}(x^q, c(x^r, T)). \quad (4)$$

To minimize Eq. (4), we query a view in its latent space $z_t \sim \mathcal{E}(x^q)$ using Eq. (1), followed by denoising z_t to \hat{z}_{t-1} conditioned on $c(x^r, \hat{T}_{r \rightarrow q})$. Finally, we compute the residuals for backpropagation of the transformation’s gradient $\nabla \hat{T}_{r \rightarrow q}$. To ensure that the estimated pose $\hat{T}_{r \rightarrow q}$ continue to lie on the SE(3) manifold during the gradient-based optimization, we parameterize the pose $T_{r \rightarrow q} = \exp(\xi)$, where $\xi \in \mathbb{R}^6$ is the twist coordinates of the Lie algebra $\mathfrak{se}(3)$ associated with the Lie group SE(3) [60]. Therefore, we reformulate Eq. (4) as follows:

$$\hat{\xi}_{r \rightarrow q} = \underset{\xi \in \mathfrak{se}(3)}{\text{argmin}} \mathcal{L}(x^q, c(x^r, \exp(\xi))). \quad (5)$$

²<https://github.com/CompVis/stable-diffusion>

³<https://laion.ai/blog/laion-aesthetics/>

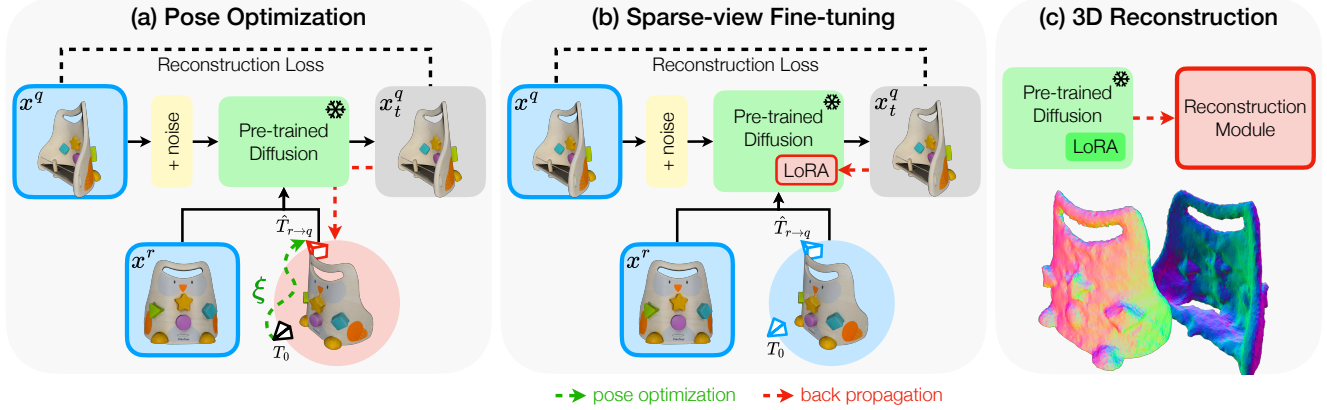


Figure 3. **iFusion framework.** (a) Given as few as two pose-free images (x^r, x^q), we estimate the pose $\hat{T}_{r \rightarrow q}$ from T_0 to optimally reconstruct the input view through the frozen diffusion model. (b) Based on $\hat{T}_{r \rightarrow q}$, we efficiently fine-tune the diffusion model by LoRA [15] to customize the model to synthesize novel views of the given object with enhanced fidelity. (c) Conditioned on $\hat{T}_{r \rightarrow q}$ and the refined diffusion model, we optimize a reconstruction module to perform sparse view 3D reconstruction.

Note that Eq. (5) can further be constrained by the inverse transformation defined by the same vector representation, *i.e.*, $T_{q \rightarrow r} = \exp(-\xi)$. We therefore obtain:

$$\hat{\xi}_{r \rightarrow q} = \underset{\xi \in \mathfrak{se}(3)}{\operatorname{argmin}} \mathcal{L}(x^q, c(x^r, \exp(\xi))) + \mathcal{L}(x^r, c(x^q, \exp(-\xi))). \quad (6)$$

In practice, we initialize our optimization from four distinct canonical poses relative to the reference view, *i.e.*, front, left, right, and back, designated as T_0 . This helps reduce the possibility of sticking at a local minima during the optimization. The final estimated camera pose can be denoted as follows:

$$\hat{T}_{r \rightarrow q} = T_0 \cdot \exp(\hat{\xi}_{r \rightarrow q}). \quad (7)$$

Furthermore, taking inspiration from Huang *et al.* [16], instead of sampling the timestep t from a uniform distribution as in training, we linearly decrease t . This adjustment aligns with diffusion models’ coarse-to-fine progressive optimization and has been empirically observed to lead to more stable optimization.

3.2. From Single-View to Multi-View

Even with a fairly accurate estimated pose $\hat{T}_{r \rightarrow q}$, there is still no guarantee that the diffusion model generates the pixel-exact query image x^q . We propose to close the gap by further fine-tuning the DM with the given views and estimated poses. However, due to limited training samples, naively optimizing all trainable parameters θ is inefficient and may jeopardize the pre-trained model. To this end, we incorporate LoRA [15], injecting thin trainable layers ϕ to the attention module in the U-Net ϵ_θ while freezing the pre-trained θ . The objective in Eq. (2) is reformulated as fol-

lows:

$$\mathcal{L}_\phi(x, c) = \mathbb{E}_{z, \epsilon, t} \left[\|\epsilon - \epsilon_{\theta, \phi}(z_t, t, c)\|_2^2 \right], \quad (8)$$

where $(x, c) \in \left\{ \left(x^q, (x^r, \hat{T}_{r \rightarrow q}) \right), \left(x^r, (x^q, \hat{T}_{q \rightarrow r}) \right) \right\}$. In other words, the fine-tuning process adapts the DM to generate the query view x^q from condition $c(x^r, \hat{T}_{r \rightarrow q})$, and vice versa, for a specific object. Empirically, this LoRA fine-tuning effectively customize the DM to generate novel views different from x^r and x^q of the target object, despite the small number of training samples and parameters ϕ , and the inherent noise from the estimated poses.

While the original Zero123 only conditions on a single view, we have multiple images available along with their relative transformations in a sparse-view setting.⁴ This raises the question: How can we better utilize these additional views for improved generation quality? To address this, we employ a simple stochastic conditioning strategy inspired by Watson *et al.* [72]. The key concept is that all given views should collectively shape the final output. More specifically, we randomly sample a registered view as the input condition at each denoising timestep. Empirically, this stochastic multi-view conditioning (MVC) significantly improves the novel view synthesis results compared to naively using the nearest view as the condition. Moreover, the final reconstruction quality is also improved.

3.3. From Sparse Views to 3D Reconstruction

There are two primary lines of existing literature for 3D object reconstruction via diffusion, namely image-based reconstruction [30, 32] and SDS-based generation [29, 47, 48,

⁴We mainly formulate the two-view setting (x^r and x^q). Multi-view settings are achieved via treating all distinct image pairs as query-reference pairs and estimating the pose transform for each pair.

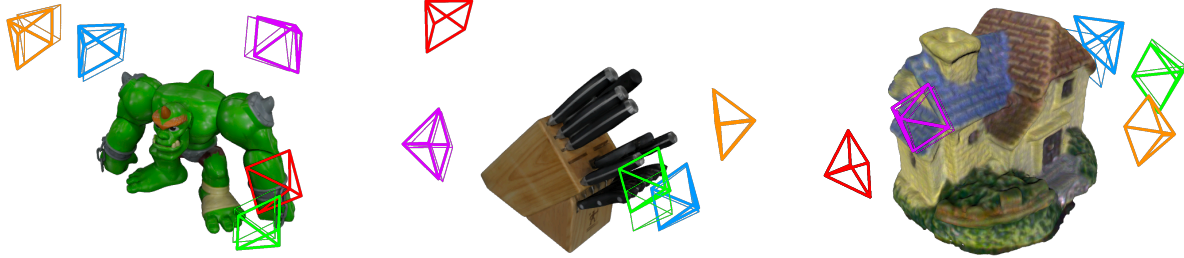


Figure 4. **Qualitative results on pose estimation.** We visualize the predicted poses (thin) alongside the ground truth (bold), using the same color, while the **reference views** are plotted in red. *iFusion* accurately predicts poses even on the opposite side of the **reference view** (red), emphasizing its effectiveness in leveraging the strong prior knowledge embedded in Zero123 [31].

64]. To integrate our proposed technique with the image-based approaches, we may simply generate multi-view images using the fine-tuned model obtained from Eq. (8) with stochastic multi-view conditioning outlined in Sec. 3.2, and then feed them as the training data to the differentiable renderer, *e.g.*, NeRF [38] and NeuS [69]. For SDS-based methods, in addition to Eq. (3), we further incorporate the reconstruction loss on the registered input views:

$$\mathcal{L}_{rec} = \left\| x - \mathcal{R}_{\psi}(\hat{T}) \right\|_2^2, \quad (9)$$

where x is the input image and $\mathcal{R}_{\psi}(\hat{T})$ is the rendered view from viewpoint \hat{T} acquired from Eq. (7). The final objective is the weighted sum of \mathcal{L}_{rec} and \mathcal{L}_{SDS} . For above steps, the LoRA model and MVC are also employed.

4. Experiments

4.1. Experimental Setup

Datasets We conduct experiments using two publicly available object datasets: Google Scanned Object (GSO) [9] and OmniObject3D (OO3D) [73]. We sample 70 instances from each dataset, randomly synthesizing camera poses and rendering observation views. For pose estimation experiments, we render five views per object, accumulating 1,400 views in total with their corresponding camera poses for each dataset. Regarding novel view synthesis and 3D reconstruction experiments, we sample two views from the rendered five with the largest parallax motion around the object to minimize the overlapping between views.

Experiments and Metrics We evaluate our proposed framework on pose estimation, novel view synthesis, and 3D reconstruction. For pose estimation, we report the median error in rotation and translation along with a recall evaluation with a 5° threshold for both, *i.e.*, we consider a true positive only when both rotation and translation errors are within the threshold. Recall results are reported in percentage. Following Liu *et al.* [31], Mildenhall *et al.* [38], we adopt the standard metrics PSNR, SSIM, and LPIPS to

Table 1. **Evaluation results on pose estimation.** *iFusion* achieves significant improvements for all metrics under 2 input views.

Dataset	Method	Rot. ↓	Trans. ↓	Recall ↑
GSO [9]	RelPose++ [27]	109.89	90.58	0.21
	FORGE [19]	111.15	88.01	0.00
	<i>iFusion</i>	2.29	2.22	74.79
OO3D [73]	RelPose++	108.83	90.83	0.00
	FORGE	107.82	87.21	0.00
	<i>iFusion</i>	2.97	2.80	69.29

evaluate novel view synthesis results. For 3D reconstruction, we report Chamfer Distances and volumetric IoU between ground truth shapes and reconstructed ones.

4.2. Experimental Result

Pose Estimation We compare our proposed *iFusion* with RelPose++ [27] and FORGE [19] for pose estimation given two views of each object. Quantitative and qualitative results are depicted in Table 1 and Fig. 4, respectively. Table 1 verifies the effectiveness of our proposed solution over the baselines with significant improvements for all metrics. We found that by leveraging the diffusion model [31], *iFusion* excels at handling diverse objects thanks to its strong prior knowledge learned during pre-training, whereas RelPose++ and FORGE fall short due to their smaller training dataset with limited object diversity. Based on the qualitative results presented in Fig. 4, we corroborate the benefits of our proposed solution in estimating the pose between two given views. We consistently find that our solution estimates accurate camera poses even with minimal overlapping. This is evident in Fig. 4, where all samples show several cameras on the opposite sides to the camera reference (red camera) and *iFusion* still achieves accurate estimations. Notably, COLMAP [55] cannot serve as a baseline in our evaluation due to the structural limitations of Structure-from-Motion, which requires a large number of views for optimization.



Figure 5. **Qualitative examples on novel view synthesis.** *iFusion* takes two unposed images and Zero123 [31] only conditions on the first view. We observe that *iFusion* effectively leverages the additional images without camera poses and generates more faithful images.

Table 2. **Novel view synthesis results.** *iFusion* performed significantly better than the original Zero123 and 3D-based methods.

Dataset	Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
GSO [9]	FORGE [19]	10.45	0.673	0.449
	LEAP [20]	12.51	0.751	0.312
	Zero123 [31]	15.40	0.788	0.184
	<i>iFusion</i>	18.73	0.836	0.121
OO3D [73]	FORGE	10.48	0.684	0.447
	LEAP	12.63	0.759	0.305
	Zero123	15.84	0.801	0.184
	<i>iFusion</i>	19.78	0.851	0.117

Novel View Synthesis Table 2 shows our novel view synthesis comparison against 2D-based Zero123⁵ and 3D-based methods, *i.e.*, FORGE and LEAP [20]. It is observed that both the 3D-based methods do not perform well under extremely few-view scenarios. Moreover, *iFusion* significantly outperforms all methods on all metrics. Figure 5 includes qualitative examples to demonstrate *iFusion*’s advantage in novel view synthesis. We observe that images generated by Zero123, although mostly visually plausible, do not faithfully represent the actual objects, especially

⁵By default, we use Zero123-XL for all modules that require Zero123.

those with complex geometry. In contrast, our *iFusion* improves novel views’ image fidelity by conditioning on an additional pose-free view.

3D Reconstruction We showcase the efficacy of the *iFusion* framework in 3D reconstruction by integrating it with various existing reconstruction methods. Specifically, One-2-3-45 [30] represents image-based methods, which directly regresses SDFs from the generated multi-view images; on the other hand, Zero123-SDS [31], Magic123 [48], and DreamGaussian [64] are SDS-based approaches. For completeness, Zero123-SDS trains Instant-NGP [41] via Zero123-guided SDS. Magic123 combines Zero123 and SD for improved quality.⁶ DreamGaussian leverages the recent 3D Gaussian Splatting renderer [23]. As illustrated in Table 3 and Fig. 6, the incorporation of *iFusion* enhances the performance of all reconstruction modules by a large margin. In addition, *iFusion* clearly outperforms other non-optimization-based methods Point-E [42] and Shape-E [21], which are trained on a large-scale private dataset. To conclude, when faithful reconstruction is desired, *iFusion* is extremely beneficial, requiring very few additional view that can be casually captured without knowing the camera poses.

⁶The implementations of Zero123-SDS and Magic123 are adopted from threestudio: <https://github.com/threestudio-project/threestudio>.

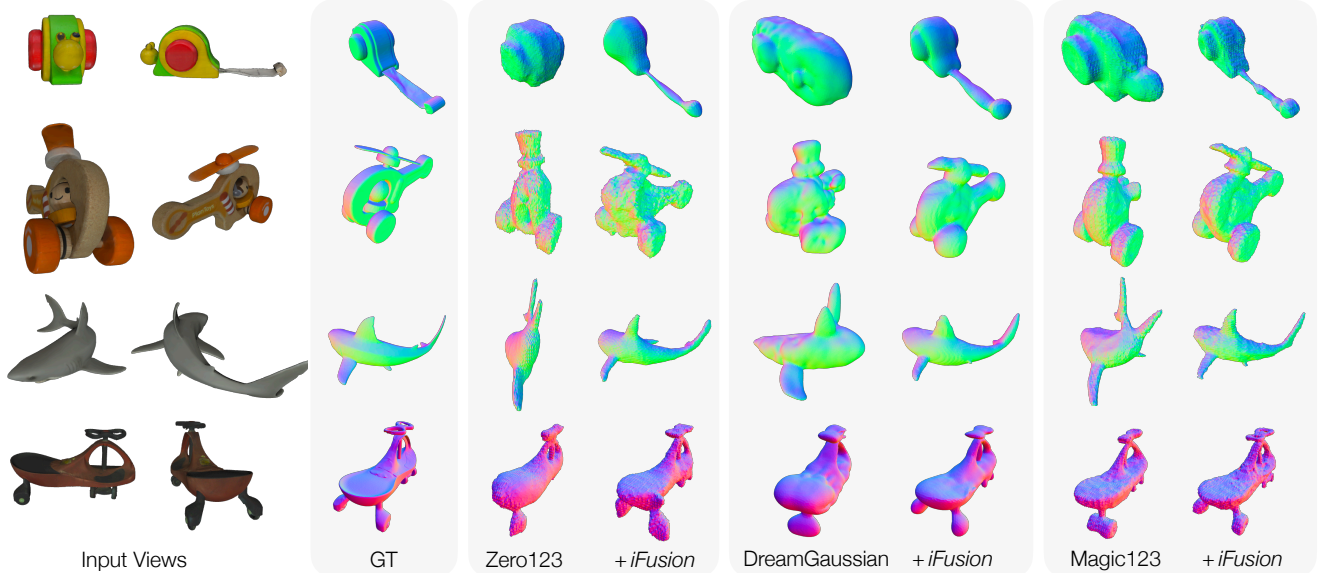


Figure 6. **Qualitative comparison of surface reconstruction.** It is clear that *iFusion* significantly enhances existing reconstruction methods including Zero123-SDS [31], DreamGaussian [64], and Magic123 [48], by adding an additional view without the camera pose.

Table 3. **Evaluation results on 3D reconstruction.** Strong single-view reconstruction baselines are improved by *iFusion* consistently.

Method	GSO [9]		OO3D [73]	
	Chamfer Dist. ($\times 10^3$) \downarrow	Volume IoU (%) \uparrow	Chamfer Dist. ($\times 10^3$) \downarrow	Volume IoU (%) \uparrow
Point-E [42]	6.414	18.92	6.766	19.83
Shape-E [21]	5.839	29.00	6.086	29.02
One-2-3-45 [30]	7.173	28.77	5.424	43.75
+ <i>iFusion</i>	<u>6.359</u>	<u>31.68</u>	<u>4.739</u>	<u>48.32</u>
Zero123-SDS [31]	6.456	33.63	5.676	45.90
+ <i>iFusion</i>	<u>4.178</u>	<u>39.73</u>	<u>3.293</u>	<u>56.36</u>
DreamGaussian [64]	4.728	35.35	4.298	44.35
+ <i>iFusion</i>	<u>3.977</u>	<u>42.07</u>	<u>2.947</u>	<u>57.58</u>
Magic123 [48]	4.839	39.46	3.842	53.69
+ <i>iFusion</i>	3.076	46.70	2.682	60.31

4.3. Ablation Study

Pose Estimation We first validate whether the use of more poses for initialization, namely T_0 in Eq. (7), leads to more accurate camera pose estimation, and it is confirmed in Table 4. The reported computation time was measured on a single Nvidia 3090 GPU. Based on Table 4, we employed $n = 4$ initial poses for a better trade-off between speed and accuracy for all experiments unless otherwise specified. Additionally, we observed that linearly annealing the timestep t lead to significantly more accurate pose estimation, as demonstrated in Table 5.

Sparse-view Fine-tuning Table 6 assesses the efficacy of the proposed fine-tuning stage for object-specific novel view synthesis. Upon examining row (a), *i.e.*, Zero123, alongside row (b), it is evident that the performance is

boosted by incorporating the additional view and an accurately estimated pose. Row (c) highlights the substantial improvement from the stochastic re-sampling of multi-view conditions at each timestep, providing more robust outcomes than row (b). Moreover, the multi-view fine-tuning with LoRA in row (d) significantly enhances performance by improving the understanding of the target object. Finally, row (e) underscores the potential for achieving higher-quality synthesis by incorporating more views. All are achieved with self-estimated camera poses.

3D Reconstruction We validate the proposed components contributing to reconstruction in Table 7, using DreamGaussian as the reconstruction module on the OO3D dataset. The results in rows (a) and (b) distinctly illustrate that adding an extra view with an estimated pose and su-

Table 4. Ablation of the **number of initial poses** for pose estimation on GSO [9].

	n poses	Recall \uparrow			Time (s) \downarrow
		5°	10°	20°	
(a)	1	33.07	36.21	38.36	22.30
(b)	2	60.57	69.14	73.07	38.51
(c)	4	74.79	84.29	88.57	70.59
(d)	8	78.21	88.93	92.43	133.73

Table 5. Ablation of **t annealing** for pose estimation on GSO [9].

	n poses	t annealing	Recall \uparrow		
			5°	10°	20°
(a)	4	-	48.61	56.67	61.39
(b)	4	\checkmark	74.79	84.29	88.57

pervising with reconstruction loss significantly enhance the single-view baseline. Incorporating stochastic multi-view conditioning (MVC) further improves the performance, as evident in row (c). Finally, fine-tuning via LoRA demonstrates an additional improvement in customizing the model for faithful reconstruction of the given object.

5. Related Work

Few-shot NeRFs Neural Radiance Fields (NeRFs) [38] have revolutionized 3D modeling with its powerful representations and high-fidelity render quality, but struggling under insufficient views. Follow-up works introduced regularizations to stabilize training [24, 43, 77], or prior models for auxiliary 3D reasoning [5, 17, 70, 79]. Nevertheless, the dependency on precise camera poses remains an issue, as Lin *et al.* [28] showed that inaccurate poses, which often arise in pose estimation using a limited number of views, lead to degraded performance.

Diffusion for 3D Generation Diffusion models [14, 59, 61] have emerged as the leading visual generative models. They generate visually plausible images from various input conditions [10, 11, 26, 37, 75, 76] and customize or edit existing photos with diverse controlling signals [2, 12, 49, 53, 83, 84]. Promising results have been achieved in 3D generation as well, spanning various representations such as point-clouds [35, 80, 85], voxel grids [40, 85], and tri-planes [1, 13, 57]; however, they are constrained by the limited diversity of 3D datasets, *e.g.*, ShapeNet [4]. To overcome the data scarcity, researchers utilize pre-trained 2D diffusion models [52, 54] for text-to-3D generation [6, 29, 47, 71], and further extend them for single-view reconstruction [30, 31, 36, 48, 64, 65], where the diffusion model “*dreams up*” unobserved views. However, single-view methods diverge from real-world recon-

Table 6. **Ablation of novel view synthesis** on GSO [9]. Multi-view conditioning and LoRA [15] finetuning are validated. Increased views also improve the scores.

	n views	Strategy	LoRA	PSNR \uparrow	LPIPS \downarrow
(a)	1	-	-	15.40	0.184
(b)	2	closest-view	-	16.19	0.169
(c)	2	multi-view	-	17.30	0.149
(d)	2	multi-view	\checkmark	18.73	0.121
(e)	4	multi-view	\checkmark	21.32	0.092

Table 7. **Ablation of 3D reconstruction** on OO3D [73] based on DreamGaussian [64]. MVC + LoRA achieves the best result.

	n views	MVC	LoRA	Chamfer Dist. \downarrow	IoU \uparrow
(a)	1	-	-	4.298	44.35
(b)	2	-	-	3.427	53.04
(c)	2	\checkmark	-	3.241	54.16
(d)	2	\checkmark	\checkmark	2.947	57.58

*Chamfer distance measured by $\times 10^3$ and IoU in (%)

struction scenarios — the target object needs to be accurately reconstructed, not over-imagined. Although several methods propose to include additional views, accurate camera poses are still assumed [3, 22, 63, 86].

Pose-free Reconstruction To recover the unknown camera poses from sparse views, recent studies have explored learnable pose estimation, either by directly regressing the pose [19, 27, 82] or through iterative refinement [58, 68]. The estimated poses can then be utilized for reconstruction [19, 28, 81]. Notably, FORGE [19] combines the two stages to achieve pose-free reconstruction but lacks robustness for intricate geometry and is sensitive to lighting. A recent follow-up, LEAP [20], eliminates pose estimation by employing DINOv2 [44] for feature mapping, showing improved generalization but struggling at unseen regions. In contrast, our solution excels in these scenarios, empirically proving its value in extreme few-shot situations.

6. Conclusion

We propose *iFusion*, a framework that reconstructs 3D objects without requiring poses, leveraging a large-scale pre-trained diffusion model as a prior. Given a few unposed images, we begin with inverting the diffusion for gradient-based pose optimization. The estimated poses, in turn, enhance the novel view synthesis diffusion model through multi-view fine-tuning and conditioning. Finally, by combining the estimated poses and the refined diffusion model, we demonstrate how *iFusion* achieves pose-free reconstruction. Experimental results show that our solution outperforms strong baselines on three key tasks: pose estimation, novel view synthesis, and 3D reconstruction.

Acknowledgement

This work is supported in part by the National Science and Technology Council (NSTC 111-2634-F-002-022). The views and opinions expressed in this paper are solely those of the authors and do not necessarily represent the official policies or positions of their affiliations or funding agencies.

References

- [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. 8
- [2] Omri Avrahami, Thomas Hayes, Oran Gafni, Sonal Gupta, Yaniv Taigman, Devi Parikh, Dani Lischinski, Ohad Fried, and Xi Yin. Spatext: Spatio-textual representation for controllable image generation. In *CVPR*, 2023. 8
- [3] Eric R. Chan, Koki Nagano, Matthew A. Chan, Alexander W. Bergman, Jeong Joon Park, Axel Levy, Miika Aittala, Shalini De Mello, Tero Karras, and Gordon Wetzstein. GeNVS: Generative novel view synthesis with 3D-aware diffusion models. In *ICCV*, 2023. 1, 8
- [4] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, *et al.* Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015. 8
- [5] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. Mvsnr: Fast generalizable radiance field reconstruction from multi-view stereo. In *CVPR*, 2021. 8
- [6] Rui Chen, Yongwei Chen, Ningxin Jiao, and Kui Jia. Fantasia3d: Disentangling geometry and appearance for high-quality text-to-3d content creation. In *ICCV*, 2023. 8
- [7] Xu Chen, Zijian Dong, Jie Song, Andreas Geiger, and Otmar Hilliges. Category level object pose estimation via neural analysis-by-synthesis. In *ECCV*, 2020. 2, 3
- [8] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3d objects. In *CVPR*, 2023. 2, 3, 13
- [9] Laura Downs, Anthony Francis, Nate Koenig, Brandon Kinman, Ryan Hickman, Krista Reymann, Thomas B McHugh, and Vincent Vanhoucke. Google scanned objects: A high-quality dataset of 3d scanned household items. In *ICRA*, 2022. 5, 6, 7, 8
- [10] Wan-Cyuan Fan, Yen-Chun Chen, Dongdong Chen, Yu Cheng, Lu Yuan, and Yu-Chiang Frank Wang. Frido: Feature pyramid diffusion for complex scene image synthesis. In *AAAI*, 2023. 8
- [11] Oran Gafni, Adam Polyak, Oron Ashual, Shelly Sheynin, Devi Parikh, and Yaniv Taigman. Make-a-scene: Scene-based text-to-image generation with human priors. In *ECCV*, 2022. 8
- [12] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit Haim Bermano, Gal Chechik, and Daniel Cohen-or. An image is worth one word: Personalizing text-to-image generation using textual inversion. In *ICLR*, 2023. 8
- [13] Jiatao Gu, Alex Trevithick, Kai-En Lin, Joshua M Susskind, Christian Theobalt, Lingjie Liu, and Ravi Ramamoorthi. Nerfdiff: Single-image view synthesis with nerf-guided distillation from 3d-aware diffusion. In *ICML*, 2023. 8
- [14] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020. 3, 8
- [15] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *ICLR*, 2022. 2, 4, 8
- [16] Yukun Huang, Jianan Wang, Yukai Shi, Xianbiao Qi, Zheng-Jun Zha, and Lei Zhang. Dreamtime: An improved optimization strategy for text-to-3d content creation. *arXiv preprint arXiv:2306.12422*, 2023. 4
- [17] Ajay Jain, Matthew Tancik, and Pieter Abbeel. Putting nerf on a diet: Semantically consistent few-shot view synthesis. In *ICCV*, 2021. 1, 8
- [18] Ajay Jain, Ben Mildenhall, Jonathan T. Barron, Pieter Abbeel, and Ben Poole. Zero-shot text-guided object generation with dream fields. In *CVPR*, 2022. 3
- [19] Hanwen Jiang, Zhenyu Jiang, Kristen Grauman, and Yuke Zhu. Few-view object reconstruction with unknown categories and camera poses. *arXiv preprint arXiv:2212.04492*, 2022. 1, 5, 6, 8
- [20] Hanwen Jiang, Zhenyu Jiang, Yue Zhao, and Qixing Huang. Leap: Liberate sparse-view 3d modeling from camera poses. *arXiv preprint arXiv:2310.01410*, 2023. 1, 6, 8
- [21] Heewoo Jun and Alex Nichol. Shap-e: Generating conditional 3d implicit functions. *arXiv preprint arXiv:2305.02463*, 2023. 6, 7
- [22] Animesh Karnawar, Andrea Vedaldi, David Novotny, and Niloy J Mitra. Holodiffusion: Training a 3d diffusion model using 2d images. In *CVPR*, 2023. 8
- [23] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM TOG*, 2023. 2, 3, 6
- [24] Mijeong Kim, Seonguk Seo, and Bohyung Han. Infonerf: Ray entropy minimization for few-shot neural volume rendering. In *CVPR*, 2022. 1, 8
- [25] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 12
- [26] Yuheng Li, Haotian Liu, Qingyang Wu, Fangzhou Mu, Jianwei Yang, Jianfeng Gao, Chunyuan Li, and Yong Jae Lee. Gligen: Open-set grounded text-to-image generation. In *CVPR*, 2023. 8
- [27] Amy Lin, Jason Y Zhang, Deva Ramanan, and Shubham Tulsiani. Relpose++: Recovering 6d poses from sparse-view observations. *arXiv preprint arXiv:2305.04926*, 2023. 1, 5, 8, 12
- [28] Chen-Hsuan Lin, Wei-Chiu Ma, Antonio Torralba, and Simon Lucey. Barf: Bundle-adjusting neural radiance fields. In *ICCV*, 2021. 8
- [29] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3d: High-resolution text-to-3d content creation. In *CVPR*, 2023. 4, 8

- [30] Minghua Liu, Chao Xu, Haian Jin, Linghao Chen, Zexiang Xu, Hao Su, *et al.* One-2-3-45: Any single image to 3d mesh in 45 seconds without per-shape optimization. In *NeurIPS*, 2023. 1, 4, 6, 7, 8, 12
- [31] Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick. Zero-1-to-3: Zero-shot one image to 3d object. In *ICCV*, 2023. 1, 2, 3, 5, 6, 7, 8, 12, 13, 14
- [32] Yuan Liu, Cheng Lin, Zijiao Zeng, Xiaoxiao Long, Lingjie Liu, Taku Komura, and Wenping Wang. Syncdreamer: Learning to generate multiview-consistent images from a single-view image. *arXiv preprint arXiv:2309.03453*, 2023. 4, 12
- [33] Xiaoxiao Long, Cheng Lin, Peng Wang, Taku Komura, and Wenping Wang. Sparseneus: Fast generalizable neural surface reconstruction from sparse views. In *ECCV*, 2022. 1
- [34] Xiaoxiao Long, Yuan-Chen Guo, Cheng Lin, Yuan Liu, Zhiyang Dou, Lingjie Liu, Yuexin Ma, Song-Hai Zhang, Marc Habermann, Christian Theobalt, *et al.* Wonder3d: Single image to 3d using cross-domain diffusion. *arXiv preprint arXiv:2310.15008*, 2023. 12
- [35] Shitong Luo and Wei Hu. Diffusion probabilistic models for 3d point cloud generation. In *CVPR*, 2021. 8
- [36] Luke Melas-Kyriazi, Christian Rupprecht, Iro Laina, and Andrea Vedaldi. Realfusion: 360 reconstruction of any object from a single image. In *CVPR*, 2023. 1, 8
- [37] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. SDEdit: Guided image synthesis and editing with stochastic differential equations. In *ICLR*, 2022. 8
- [38] 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, 3, 5, 8, 12
- [39] Nasir Mohammad Khalid, Tianhao Xie, Eugene Belilovsky, and Tiberiu Popa. Clip-mesh: Generating textured meshes from text using pretrained image-text models. In *SIGGRAPH Asia*, 2022. 3
- [40] Norman Müller, Yawar Siddiqui, Lorenzo Porzi, Samuel Rota Buló, Peter Kotschieder, and Matthias Nießner. Diffrf: Rendering-guided 3d radiance field diffusion. In *CVPR*, 2023. 8
- [41] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM TOG*, 2022. 6
- [42] Alex Nichol, Heewoo Jun, Prafulla Dhariwal, Pamela Mishkin, and Mark Chen. Point-e: A system for generating 3d point clouds from complex prompts. *arXiv preprint arXiv:2212.08751*, 2022. 6, 7
- [43] Michael Niemeyer, Jonathan T. Barron, Ben Mildenhall, Mehdi S. M. Sajjadi, Andreas Geiger, and Noha Radwan. Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs. In *CVPR*, 2022. 8
- [44] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, *et al.* Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 8
- [45] Keunhong Park, Arsalan Mousavian, Yu Xiang, and Dieter Fox. Latentfusion: End-to-end differentiable reconstruction and rendering for unseen object pose estimation. In *CVPR*, 2020. 2, 3
- [46] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, *et al.* Pytorch: An imperative style, high-performance deep learning library. In *NeurIPS*, 2019. 12
- [47] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. In *ICLR*, 2023. 3, 4, 8
- [48] Guocheng Qian, Jinjie Mai, Abdullah Hamdi, Jian Ren, Aliaksandr Siarohin, Bing Li, Hsin-Ying Lee, Ivan Skokhodov, Peter Wonka, Sergey Tulyakov, and Bernard Ghanem. Magic123: One image to high-quality 3d object generation using both 2d and 3d diffusion priors. *arXiv preprint arXiv:2306.17843*, 2023. 1, 4, 6, 7, 8, 12, 14
- [49] Can Qin, Shu Zhang, Ning Yu, Yihao Feng, Xinyi Yang, Yingbo Zhou, Huan Wang, Juan Carlos Niebles, Caiming Xiong, Silvio Savarese, Stefano Ermon, Yun Fu, and Ran Xu. Unicontrol: A unified diffusion model for controllable visual generation in the wild. In *NeurIPS*, 2023. 8
- [50] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, *et al.* Learning transferable visual models from natural language supervision. In *ICML*, 2021. 3
- [51] Amit Raj, Srinivas Kaza, Ben Poole, Michael Niemeyer, Ben Mildenhall, Nataniel Ruiz, Shiran Zada, Kfir Aberman, Michael Rubenstein, Jonathan Barron, Yuanzhen Li, and Varun Jampani. Dreambooth3d: Subject-driven text-to-3d generation. In *ICCV*, 2023. 1
- [52] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 3, 8
- [53] 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. In *CVPR*, 2023. 8
- [54] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, *et al.* Photorealistic text-to-image diffusion models with deep language understanding. In *NeurIPS*, 2022. 3, 8
- [55] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *CVPR*, 2016. 1, 5
- [56] H. Scudder. Probability of error of some adaptive pattern-recognition machines. *IEEE Transactions on Information Theory*, 11(3):363–371, 1965. 2
- [57] J Ryan Shue, Eric Ryan Chan, Ryan Po, Zachary Ankner, Jiajun Wu, and Gordon Wetzstein. 3d neural field generation using triplane diffusion. In *CVPR*, 2023. 8
- [58] Samarth Sinha, Jason Y Zhang, Andrea Tagliasacchi, Igor Gilitschenski, and David B Lindell. SparsePose: Sparse-view camera pose regression and refinement. In *CVPR*, 2023. 1, 8

- [59] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *ICML*, 2015. 3, 8
- [60] Joan Sola, Jeremie Deray, and Dinesh Atchuthan. A micro lie theory for state estimation in robotics. *arXiv preprint arXiv:1812.01537*, 2018. 3
- [61] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. In *NeurIPS*, 2019. 3, 8
- [62] Yu Sun, Xiaolong Wang, Zhuang Liu, John Miller, Alexei Efros, and Moritz Hardt. Test-time training with self-supervision for generalization under distribution shifts. In *ICML*, 2020. 2
- [63] Stanislaw Szymanowicz, Christian Rupprecht, and Andrea Vedaldi. Viewset diffusion:(0-) image-conditioned 3d generative models from 2d data. In *ICCV*, 2023. 1, 8
- [64] Jiaxiang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang Zeng. Dreamgaussian: Generative gaussian splatting for efficient 3d content creation. *arXiv preprint arXiv:2309.16653*, 2023. 1, 5, 6, 7, 8, 12, 14
- [65] Junshu Tang, Tengfei Wang, Bo Zhang, Ting Zhang, Ran Yi, Lizhuang Ma, and Dong Chen. Make-it-3d: High-fidelity 3d creation from a single image with diffusion prior. In *ICCV*, 2023. 1, 8
- [66] Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen, and Trevor Darrell. Tent: Fully test-time adaptation by entropy minimization. In *ICLR*, 2021. 2
- [67] 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. 3
- [68] Jianyuan Wang, Christian Rupprecht, and David Novotny. PoseDiffusion: Solving pose estimation via diffusion-aided bundle adjustment. In *ICCV*, 2023. 8
- [69] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. In *NeurIPS*, 2021. 5
- [70] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul Srinivasan, Howard Zhou, Jonathan T. Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. Ibrnet: Learning multi-view image-based rendering. In *CVPR*, 2021. 8
- [71] Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation. In *NeurIPS*, 2023. 8
- [72] Daniel Watson, William Chan, Ricardo Martin-Brualla, Jonathan Ho, Andrea Tagliasacchi, and Mohammad Norouzi. Novel view synthesis with diffusion models. In *ICLR*, 2023. 4
- [73] Tong Wu, Jiarui Zhang, Xiao Fu, Yuxin Wang, Jiawei Ren, Liang Pan, Wayne Wu, Lei Yang, Jiaqi Wang, Chen Qian, *et al.* Omniobject3d: Large-vocabulary 3d object dataset for realistic perception, reconstruction and generation. In *CVPR*, 2023. 5, 6, 7, 8
- [74] Qizhe Xie, Minh-Thang Luong, Eduard Hovy, and Quoc V Le. Self-training with noisy student improves imagenet classification. In *CVPR*, 2020. 2
- [75] Han Xue, Zhiwu Huang, Qianru Sun, Li Song, and Wenjun Zhang. Freestyle layout-to-image synthesis. In *CVPR*, 2023. 8
- [76] Binbin Yang, Yi Luo, Ziliang Chen, Guangrun Wang, Xiaodan Liang, and Liang Lin. Law-diffusion: Complex scene generation by diffusion with layouts. In *ICCV*, 2023. 8
- [77] Jiawei Yang, Marco Pavone, and Yue Wang. Freenerf: Improving few-shot neural rendering with free frequency regularization. In *CVPR*, 2023. 1, 8
- [78] Lin Yen-Chen, Pete Florence, Jonathan T. Barron, Alberto Rodriguez, Phillip Isola, and Tsung-Yi Lin. iNeRF: Inverting neural radiance fields for pose estimation. In *IROS*, 2021. 2, 3
- [79] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelNeRF: Neural radiance fields from one or few images. In *CVPR*, 2021. 1, 8
- [80] Xiaohui Zeng, Arash Vahdat, Francis Williams, Zan Gojcic, Or Litany, Sanja Fidler, and Karsten Kreis. Lion: Latent point diffusion models for 3d shape generation. In *NeurIPS*, 2022. 8
- [81] Jason Zhang, Gengshan Yang, Shubham Tulsiani, and Deva Ramanan. Ners: Neural reflectance surfaces for sparse-view 3d reconstruction in the wild. In *NeurIPS*, 2021. 8
- [82] Jason Y Zhang, Deva Ramanan, and Shubham Tulsiani. Rel-pose: Predicting probabilistic relative rotation for single objects in the wild. In *ECCV*, 2022. 1, 8
- [83] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *ICCV*, 2023. 8
- [84] Shihao Zhao, Dongdong Chen, Yen-Chun Chen, Jianmin Bao, Shaozhe Hao, Lu Yuan, and Kwan-Yee K Wong. Uni-controlnet: All-in-one control to text-to-image diffusion models. In *NeurIPS*, 2023. 8
- [85] Linqi Zhou, Yilun Du, and Jiajun Wu. 3d shape generation and completion through point-voxel diffusion. In *ICCV*, 2021. 8
- [86] Zhizhuo Zhou and Shubham Tulsiani. Sparsefusion: Distilling view-conditioned diffusion for 3d reconstruction. In *CVPR*, 2023. 1, 8

Appendix

We provide implementation details including hyper-parameters and dataset, more qualitative examples, and limitations and future directions as appendices.

A. Implementation Details

A.1. Hyper-parameters

We employ Adam [25] as the optimizer for both pose estimation and sparse-view fine-tuning. In the case of pose estimation, we optimize the initial poses for 100 steps with an initial learning rate of 0.1. The learning rate is dynamically reduced if the L2 loss stops decreasing, handled by the *ReduceLROnPlateau* scheduler from PyTorch [46]. Specifically, we set the reduction factor to 0.6 and the patience to 10. Afterward, in the sparse-view fine-tuning stage, the model is fine-tuned for 30 steps, with the learning rate annealed from 10^{-3} to 10^{-4} and the rank of the injected LoRA parameter set to 12. This process takes approximately 30 seconds on a single Nvidia 3090 GPU with a batch size of 16. For 3D reconstruction, we follow the default hyper-parameters of each reconstruction method, *i.e.*, One2345 [30], Zero123-SDS [31], Magic123 [48], and DreamGaussian [64], when combining with *iFusion*. Please refer to their official implementations for details.

A.2. Dataset Collection

We use Pyrender to render images for evaluation.⁷ Following Liu *et al.* [31], the transformation is defined using the spherical coordinate system with θ , ϕ , and r representing the elevation angle, azimuth angle, and distance towards the center, respectively. In practice, we sample camera viewpoints on the unit sphere with $\theta \in [\pi/4, 3\pi/4]$, $\phi \in [0, 2\pi]$ and r is uniformly sampled in the interval of $[1.2, 2.0]$. The field of view of the perspective camera is set to 49.1° . All images are rendered in the resolution of 512×512 with transparent background.

B. Qualitative Results

To further corroborate the effectiveness of our proposed pose estimation strategy described in Sec. 3.1, we present additional qualitative visualization in Fig. 8. These results complement the findings presented in Fig. 4 of our main manuscript. They also support our assumption that the learned understanding of diverse objects in Zero123 [31] can be leveraged for other tasks, such as pose estimation. Moreover, examples illustrating the single-view ambiguity, taken from the Blender dataset [38], are shown in Fig. 7. These instances motivated us to fine-tune and condition the model with registered multi-views.

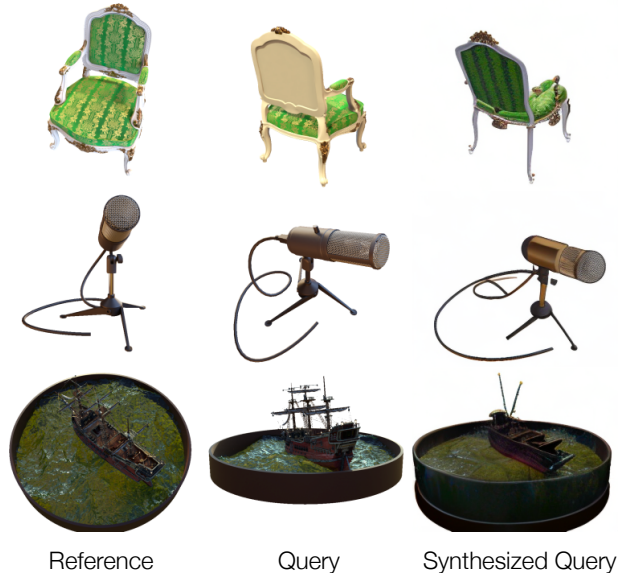


Figure 7. **Single-view ambiguity.** We show the reference view, query view, and the synthesized query view given $c(x^r, T_q^r)$. It is observed that, while the model can generate reasonable novel views, there is a gap between the model’s understanding and the actual object, arising from single-view ambiguity. This prompts us to condition the model with additional views to mitigate this issue.

In Fig. 9, we showcase additional comparisons on 3D reconstruction. These results complement Fig. 6 of our main manuscript, underscoring the efficacy of our proposed framework in achieving faithful reconstruction by considering one extra view without requiring known camera poses.

C. Limitations and Future Works

While our methods deliver highly accurate camera poses, our pose estimation run time is higher than feed-forward-based methods, *e.g.*, RelPose++ [27]. This is attributed to the optimization nature of our approach, which involves back-propagation for updating the poses. Moreover, when we fine-tune Zero123 [31] on estimated poses and additional input views, it is worth noting that Zero123, originally adapted from the 2D-based Stable Diffusion (SD), lacks complete 3D awareness. This structural limitation prevents it from generating multi-views with consistency. However, our framework holds potential for integration with other diffusion-based novel view synthesizers [32, 34] that enforce consistency by incorporating 3D-aware modules onto SD.

⁷<https://github.com/mmatl/pyrender>



Figure 8. **More qualitative results on pose estimation.** The predicted poses (thin) and their corresponding ground truth (bold), are plotted in the same color, while the **reference views** are plotted in red. We confirm that *iFusion* effectively exploits the robust understanding of diverse objects in Zero123 [31] acquired from Objaverse [8].

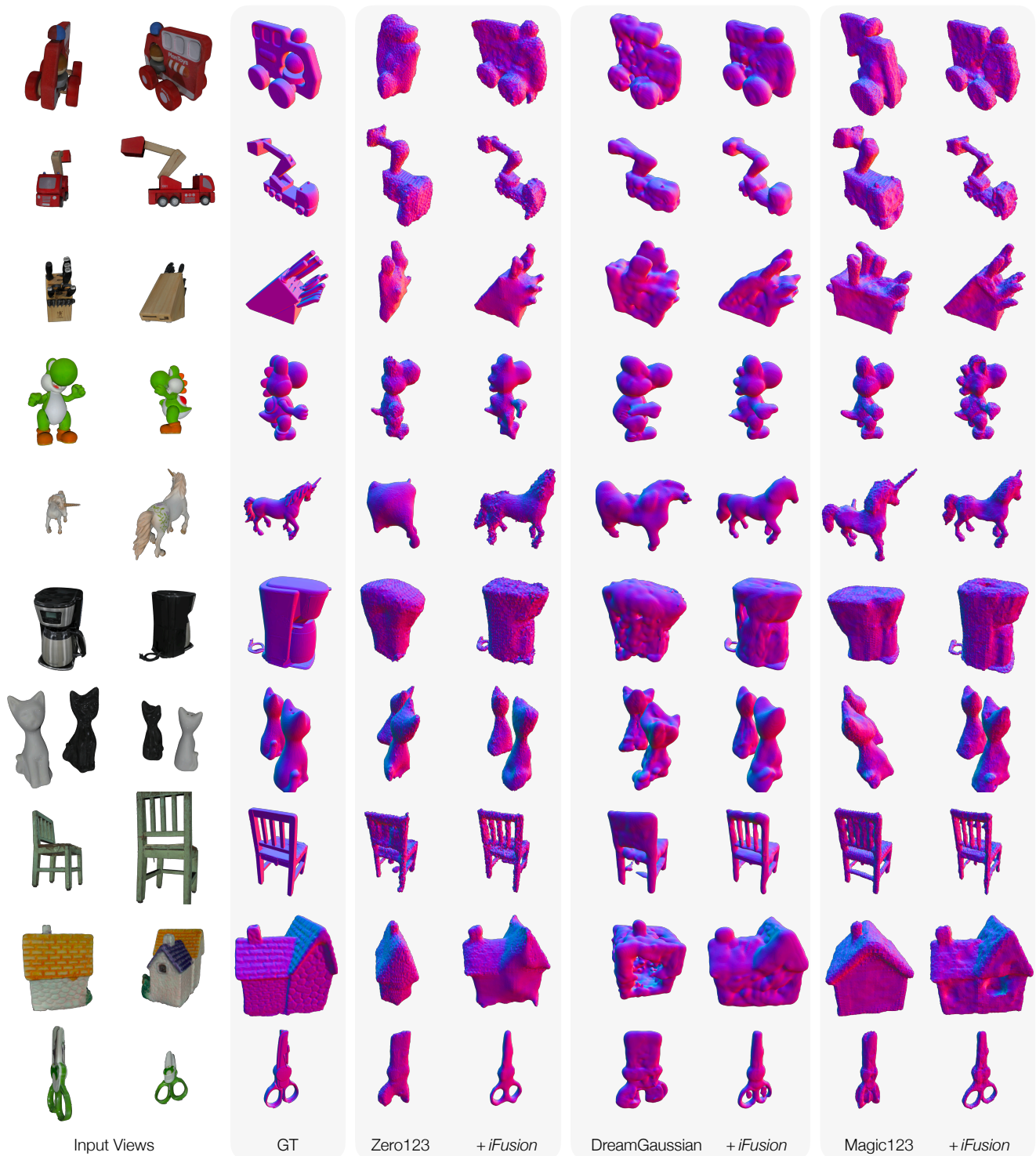


Figure 9. **More qualitative comparisons on surface reconstruction.** We integrate *iFusion* with Zero123-SDS [31], DreamGaussian [64], and Magic123 [48] to perform pose-free reconstruction given sparse views. The results indicate that our method operates as an effective add-on, consistently enhancing existing single-view reconstruction methods.