000 MOVIS: ENHANCING MULTI-OBJECT NOVEL VIEW 001 Synthesis for Indoor Scenes 002 003

Anonymous authors

Paper under double-blind review

ABSTRACT

Repurposing pre-trained diffusion models has been proven to be effective for 012 novel view synthesis (NVS). However, these methods are mostly limited to a single object; directly applying such methods to compositional multi-object scenarios yields inferior results, especially incorrect object placement and inconsistent shape and appearance under novel views. How to enhance and systematically 015 evaluate the cross-view consistency of such models remains under-explored. To 016 address this issue, we propose MOVIS to enhance the structural awareness of the view-conditioned diffusion model for multi-object NVS in terms of model inputs, auxiliary tasks, and training strategy. First, we inject structure-aware features, including depth and object mask, into the denoising U-Net to enhance the model's comprehension of object instances and their spatial relationships. Second, we introduce an auxiliary task requiring the model to simultaneously predict novel view object masks, further improving the model's capability in differentiating and placing objects. Finally, we conduct an in-depth analysis of the diffusion sampling process and carefully devise a structure-guided timestep sampling scheduler during training, which balances the learning of global object placement and fine-grained detail recovery. To systematically evaluate the plausibility of synthesized images, we propose to assess cross-view consistency and novel view object placement alongside existing image-level NVS metrics. Extensive experiments on challenging synthetic and realistic datasets demonstrate that our method exhibits strong generalization capabilities and produces consistent novel view synthesis, highlighting its potential to guide future 3D-aware multi-object NVS tasks.

031 032 033

034

004

010 011

013

014

017

018

019

020

021

022

024

025

026

027

028

029

1 INTRODUCTION

Novel view synthesis (NVS) from a single image is imperative for various applications, including AR/VR, interior designs, robotics, etc. This is highly challenging as it requires understanding complex 037 spatial structures from a single 2D perspective observation while being able to extrapolate consistent 038 and plausible content for unobserved areas. The substantial demands for comprehensive knowledge of the 3D world render it a difficult task, even for humans with rich priors of the 3D environments.

040 Recently, significant progress has been made in the realm of single-object image-to-3D genera-041 tion (Tang et al., 2023b; Liu et al., 2023b; Shi et al., 2023a; Liu et al., 2023a; Shi et al., 2023b; Long 042 et al., 2024) empowered by the advances in 2D diffusion models (Rombach et al., 2022; Ho et al., 043 2020). Among them, one prominent line of research (Liu et al., 2023b; Lin et al., 2023a; Qian et al., 044 2023; Tang et al., 2023a; Weng et al., 2023; Lin et al., 2023b; Liu et al., 2023d; Huang et al., 2023; Chen et al., 2023) has achieved compelling results by building on insights from Zero-1-to-3 (Liu et al., 2023c): repurposing a pre-trained diffusion model as a novel view synthesizer by fine-tuning 046 on large 3D object datasets can provide promising 3D-aware prior for image-to-3D tasks. 047

048 However, these methods are mostly restricted to the single-object level. It remains unclear if this paradigm can be effectively extended to the multi-object level to facilitate more complex tasks like reconstructing an indoor scene. In Fig. 1, we visualize cross-view matching results of directly 051 applying the aforementioned novel view synthesizers (Liu et al., 2023c) in multi-object scenarios, which showcases weak consistency with input views. Specifically, we believe that the lack of 052 structural awareness is the primary reason for the disappearance, distortion, incorrect position, and orientation of objects under novel views. While several works (Sargent et al., 2023; Tung et al., 2024)

077

079

081



Figure 1: Novel view synthesis and cross-view image matching. The first row shows that MOVIS generalizes to different datasets on NVS. We also show visualizations of cross-view consistency compared with Zero-1-to-3 (Liu et al., 2023c) and ground truth by applying image-matching. MOVIS can match a significantly greater number of points, closely aligned with the ground truth.

have explored training on mixed real-world scene datasets, the complexity introduced by multiple
 objects, such as spatial placement, per-instance geometry and appearance, and occlusion relationship,
 makes incorporating such awareness non-trivial.

085 Inspired by the discussion above, our paper seeks to address the question: How to enhance the structural awareness of current diffusion-based novel view synthesizers? We begin by identifying the 087 key challenges in extending single-object methods for multi-object NVS tasks. A multi-object image 088 possesses more complicated structural information than a single-object one. The model must first grasp the hierarchical structure within, which includes both high-level global object placement, e.g., 089 position and orientation, and low-level ones like per-object geometry and appearance. High-level 090 structural information significantly reduces the ambiguity in object composition while low-level 091 details are essential for accurately capturing the characteristics of each object instance. Subsequently, 092 the model needs to retain this hierarchical information captured from the input view while synthesizing 093 novel-view images to ensure cross-view structural consistency. These capabilities are less critical in 094 single-object level NVS tasks due to the reduced ambiguity in one-to-one mapping but are crucial for 095 effective multi-object NVS models. 096

Building on these insights, our technical designs are threefold. We first propose injecting structureaware features, *i.e.*, depth and object mask, from the input view as additional inputs to provide 098 information on both high-level global placement and fine-grained local details. Secondly, we utilize the prediction of novel view object masks as an auxiliary task during training for the model to 100 differentiate object instances, laying a solid foundation for fine-grained geometry and appearance 101 recovery. Finally, through an in-depth analysis of the model's inference process, we highlight the 102 importance of revising the noise timestep sampling schedule, which influences the learning focus 103 in the training process. To be specific, larger timesteps emphasize global placement learning, while 104 smaller timesteps focus on local fine-grained object geometry and appearance recovery. To endow 105 the view-conditioned diffusion model with both capabilities, we propose a structure-guided timestep sampling scheduler that prioritizes larger timesteps in the initial stage, gradually decreasing over time 106 to balance these two conflicting inductive biases. This design is fundamental to our proposed model's 107 effectiveness in addressing the complexity of multi-object level NVS tasks.

108 To systematically assess the plausibility of synthesized novel view images, we additionally evaluate 109 novel-view object mask and cross-view structural consistency apart from the existing NVS metrics. 110 Specifically, we employ image-matching techniques (Wang et al., 2024; Leroy et al., 2024) to compare 111 the input-view image with both the ground-truth and synthesized novel-view images. Cross-view 112 structural consistency evaluates how closely the matching results align, providing a measure of the accuracy in recovering object placement, shape, and appearance. On the other hand, the object 113 mask, as measured by Intersection over Union (IoU), assesses the precision of object placement. 114 Extensive experiments demonstrate that our method excels at multi-object level NVS in indoor 115 scenes, achieving consistent object placement, shape, and appearance. Notably, it exhibits strong 116 generalization capabilities for generating novel views on unseen datasets, including both synthetic 117 ones 3D-FRONT (Fu et al., 2021a), Room-Texture (Luo et al., 2024) and Objaverse (Deitke et al., 118 2023b), as well as the real-world SUNRGB-D (Song et al., 2015). 119

In summary, our paper focuses on enhancing structural awareness of view-conditioned diffusion models, improving the quality and consistency of synthesized images. Our main contributions are:

- We introduce structure-aware features as model inputs and incorporate novel view mask prediction as an auxiliary task during training. This enhances the model's understanding of hierarchical structures in multi-object scenarios, leading to improved NVS performance.
 - 2. We present a novel noise timestep sampling scheduler designed to balance the learning of global object placement and fine-grained detail recovery, which is critical for addressing the increased complexity in multi-object scenarios.
 - 3. We introduce additional metrics to systematically evaluate the novel view structural consistency. Through extensive experiments, our model demonstrates superiority in consistent object placement, geometry, and appearance recovery, showcasing strong generalization capability to unseen datasets.
- 130 131 132

133 134

135

125

126

127

128

129

2 RELATED WORK

2.1 SINGLE OBJECT NVS WITH GENERATIVE MODELS

Synthesizing novel view images for single objects given a single-view image is an extremely ill-posed 136 problem that requires strong priors. With great advances achieved in diffusion models (Ho et al., 137 2020; Rombach et al., 2022), research efforts (Xu et al., 2023; Tang et al., 2023b; Melas-Kyriazi 138 et al., 2023) seek to distill priors (Jain et al., 2022; Poole et al., 2022) learned from Text-to-Image 139 (T2I) diffusion models via image captioning like Li et al. (2023). However, this presents a huge gap 140 between the image and semantics due to the ambiguity of the text, hindering the 3D consistency of 141 these methods. On the other hand, view-conditioned diffusion models like Zero-1-to-3 (Liu et al., 142 2023c) explore an Image-to-Image (I2I) generation paradigm that "teaches" the diffusion model to 143 control viewpoints to synthesize plausible images under novel views, providing a more consistent 144 3D-aware prior. Subsequent work focuses on accelerating the generation speed (Liu et al., 2023b; 145 Tang et al., 2023a), enhancing the view consistency (Chen et al., 2023; Lin et al., 2023b; Weng et al., 2023; Liu et al., 2023d; Huang et al., 2023), or accelerating the training process (Jiang et al., 2023). 146 However, all these methods deal with single and complete object novel view synthesis tasks since 147 they usually fine-tune their model on Objaverse (Deitke et al., 2023b;a), an extensive single-object 148 level dataset, contrary to real images which normally consist of multiple or incomplete objects. The 149 lack of specific model designs for compositional scenes also leads to significant inconsistencies when 150 directly applying them to the multi-object scenarios, as can be seen from Fig. 1.

151 152 153

2.2 Multi-object 3D reconstruction with single object priors

154 Following the advance in 3D-aware single object generative prior (Liu et al., 2023c; Shi et al., 2023a), 155 a line of research work (Chen et al., 2024b; Dogaru et al., 2024; Chen et al., 2024a) focuses on 156 extending their application to compositional multi-object scenarios. The core idea is to decompose 157 object compositions into individual objects, thereby fully leveraging the powerful generative priors of 158 single-object models. They first break down a multi-object composition into several components via 159 segmentation models like SAM (Kirillov et al., 2023), and then complete every single object with amodal (Ozguroglu et al., 2024; Xu et al., 2024; Zhan et al., 2024) or inpainting (Rombach et al., 2022; 160 Lugmayr et al., 2022) techniques. The object instances are lifted to 3D via image-to-3D models (Tang 161 et al., 2023a; Wu et al., 2024; Liu et al., 2023c; Wang & Shi, 2023) and finally composited into



Figure 2: **Overview of MOVIS.** Our model performs NVS from the input image and relative camera change. We introduce structure-aware features as additional inputs and employ mask prediction as an auxiliary task (Sec. 3.2). The model is trained with a structure-guided timestep sampling scheduler (Sec. 3.3) to balance the learning of global object placement and local detail recovery.

a whole utilizing spatial-aware optimization, 3D bounding box detection (Brazil et al., 2023; Nie et al., 2020) or carefully estimating the metric depth (Ke et al., 2024; Yang et al., 2024b). However, this divide-and-conquer paradigm is limited by the user-specified spatial relations from language prompts (Chen et al., 2024b) and relies heavily on the cascaded modules of detection (Kirillov et al., 2023; Brazil et al., 2023; Ke et al., 2024), completion (Rombach et al., 2022; Lugmayr et al., 2022) and 3D-aware object-level novel view synthesis (NVS) (Liu et al., 2023b; Wang & Shi, 2023) to provide priors for reconstruction. Unlike any of the above, our method aims to build an end-to-end image-conditioned novel view synthesis model that can directly cope with the increased complexity in multi-object compositions, especially in indoor scenes with multiple furniture items.

2.3 SCENE-LEVEL NVS WITH SPARSE VIEW INPUT

Early efforts (Jain et al., 2021; Yu et al., 2021; Wang et al., 2021) attempted to directly perform scenelevel NVS tasks by extracting image features from input-view images and inferring the underlying 3D representation (Mildenhall et al., 2020). With the development of Gaussian Splatting (Kerbl et al., 2023), recent works (Charatan et al., 2023; Chen et al., 2024c) attempt to switch the underlying representation to Gaussian Splatting for efficiency. However, they mainly deal with synthesizing views near input ones with limited generative capabilities to the unseen region. Inspired by the great success of diffusion models (Rombach et al., 2022) and the object-level 3D-aware novel-view synthesizer (Liu et al., 2023c; Wang & Shi, 2023), several recent works have also attempted to perform scene-level NVS tasks by directly conditioning the generative models on a single-view scene image or a monocular dynamic scene video (Van Hoorick et al., 2024). ZeroNVS (Sargent et al., 2023) proposes to train a view-conditioned diffusion model on a mixture of real-world datasets, MegaScenes (Tung et al., 2024) further scales up the training dataset with Internet-level data pairs for stronger generalization capabilities. However, all these works mainly deal with small view-change and simple scenarios in terms of object number, with few adaptations to tackle the multi-object complexity. In this work, we systematically examine the cross-view consistency of NVS by proposing new metrics, and explore the critical designs required to enhance the structural consistency of the view-conditioned diffusion models in the multi-object scenarios.

²¹⁶ 3 METHOD

217 218

In this section, we address the challenge of enhancing the structural awareness of diffusion-based 219 novel view synthesizers for better cross-view consistency in multi-object scenarios. We begin with 220 a brief introduction to diffusion models and view-conditioned diffusion models (Sec. 3.1). Next, 221 we detail the key architectural designs of MOVIS, including how we incorporate structural-aware features as input to improve the model's understanding of hierarchical structure information (Sec. 3.2) 222 and how we introduce novel view mask prediction as an auxiliary task, instructing the model to 223 differentiate the object instances with correct object placement (Sec. 3.2). Finally, we provide an 224 in-depth analysis of the inference process and adopt a structure-guided timestep sampling scheduler 225 (Sec. 3.3) to balance the learning of global object placement and local fine-grained object geometry 226 and appearance recovery. We provide an overview of our view-conditioned diffusion model in Fig. 2. 227

228 229

238

242

3.1 PRELIMINARIES

230 **Diffusion Models** Diffusion models learn to generate images by gradually adding noise to an image 231 (forward process) and recovering the original image from a noisy image (backward process) (Ho 232 et al., 2020). Specifically, in the forward process, Gaussian noise is progressively introduced to the 233 image via $q(\mathbf{x}_t | \mathbf{x}_{t-1})$. Due to the additivity of Gaussian distributions, this iterative process can be 234 written as $q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\alpha_t \mathbf{x}_0, \sigma_t^2 \mathbf{I})$, where α_t and σ_t are designed to converge to $\mathcal{N}(\mathbf{0}, \mathbf{I})$ at the end of the forward process (Kingma et al., 2021; Song et al., 2020b). In the backward process, the 235 model learns to progressively denoise from a noisy image $p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)$. This learning is formulated 236 as learning the noise estimator $\epsilon_{\theta}(\mathbf{x}_t, t)$ following Ho et al. (2020): 237

$$\mathbb{E}[||\epsilon_{\theta}(\alpha_t \mathbf{x}_0 + \sigma_t \epsilon, t) - \epsilon||_2^2], \tag{1}$$

where ϵ is drawn from $\mathcal{N}(\mathbf{0}, \mathbf{I})$ and the timestep t is uniformly sampled from $\mathcal{U}(\mathbf{1}, \mathbf{1000})$. In the inference stage, one can either apply a stochastic (Ho et al., 2020) or a deterministic (Song et al., 2020a) sampler to generate high-quality images via iterative refinement.

View-conditioned Diffusion Models Diffusion models have been recently repurposed as a novel view synthesizer. By training on posed image pairs $\{(\mathbf{x}_0, \hat{\mathbf{x}}_0)\}$ where $\hat{\mathbf{x}}_0 \in \mathbb{R}^{H \times W \times 3}$ denotes the input view image and $\mathbf{x}_0 \in \mathbb{R}^{H \times W \times 3}$ denotes the target view, view-conditioned diffusion models (Watson et al., 2022; Liu et al., 2023c) use the input image $\hat{\mathbf{x}}_0$ and camera pose transformation as conditions to predict the target view image \mathbf{x}_0 from a different viewpoint. Specifically, the learning objective of view-conditioned diffusion models is:

$$\mathbb{E}[||\epsilon_{\theta}(\alpha_t \mathbf{x}_0 + \sigma_t \epsilon, t, C(\hat{\mathbf{x}}_0, R, T)) - \epsilon||_2^2],$$
(2)

249 where R, T represent the relative camera pose transformation between the target view \mathbf{x}_0 and the 250 input view $\hat{\mathbf{x}}_0$. $C(\hat{\mathbf{x}}_0, R, T)$ is the view-conditioned feature, combining the relative camera pose 251 transformation with encoded image features to form a new 'pose-aware' feature map, taking the place 252 of the origin CLIP (Radford et al., 2021) feature embedding. Moreover, input view image x_0 will 253 be concatenated with the noisy image as the input of the denoising U-Net. As discussed in Sec. 2.1, 254 single-image-based NVS is extremely challenging, current methods inherit natural image priors from 255 large-scale pre-training (Rombach et al., 2022) and fine-tune diffusion models on large-scale 3D object datasets like Objaverse (Deitke et al., 2023b) to learn the transformation between objects 256 in the input and novel views given the relative camera pose. Despite their ability to generalize to 257 in-the-wild objects, these view-conditioned diffusion models struggle with multi-object scenarios 258 like multi-furniture indoor scenes due to the scarcity of similar data and increased complexity arising 259 from intricate object compositions. Our method builds on the insight of repurposing the diffusion 260 model as a novel view synthesizer while emphasizing the inherent properties of multi-object scenarios 261 in both model design and training strategy to facilitate multi-object NVS. 262

263

263 3.2 MOVIS

Our proposed method extends view-conditioned diffusion models to multi-object level, as illustrated in Fig. 2. The model leverages a pre-trained Stable Diffusion (Rombach et al., 2022) and concatenates the 2D structural information from the input view with a noisy target image as input. Additionally, it integrates a pre-trained image encoder (Oquab et al., 2023) to capture semantic information, which is injected into the network through cross-attention alongside the relative camera pose. Moreover, it predicts novel view mask simultaneously as an auxiliary task to aid global object placement learning. 270 **Structure-Aware Feature Amalgamation** To synthesize plausible images under novel viewpoints, 271 the model must first grasp the compositional structural information from the input view, laying a 272 solid foundation for generation. To address the innate complexity in multi-object scenarios due 273 to the intricate object relationship, we propose to leverage structure-aware features to facilitate 274 model's comprehension. Specifically, we use depth maps and object masks as proxies for image-level structural information. Object masks provide a rough concept of object placement and shape as 275 well as distinguishing distinct object instances, while depth maps encode the rough relative position 276 and shape of the visible objects. Together with input-view images, these conditions provide both global structural information like object placement and local fine-grained details like object shape. 278 Concretely speaking, we normalize the image rendered with object instance IDs of the input view 279 to create a continuous object mask image M. We then replicate the depth map D and object mask 280 image $\hat{\mathbf{M}}$ into three channels to simulate RGB images. These two structural-aware feature images, 281 along with the input image $\hat{\mathbf{x}}_0$, are passed into a VAE to obtain latent features, which will be later 282 concatenated with the noisy target view image x_t as input to the denoising U-Net. Note that both 283 object mask and depth can be obtained with off-the-shelf detectors during the inference stage, such 284 as SAM (Kirillov et al., 2023) and Marigold (Ke et al., 2024). After introducing these additional 285 conditions, the learning objective of MOVIS becomes:

$$\mathbb{E}[||\epsilon_{\theta}(\alpha_t \mathbf{x}_0 + \sigma_t \epsilon, t, C_{SA}(\hat{\mathbf{x}}_0, R, T, \widehat{\mathbf{D}}, \widehat{\mathbf{M}})) - \epsilon||_2^2].$$
(3)

We use $C_{SA}(\cdot)$ as a shorthand for the structure-aware view-conditioned feature throughout the paper.

290 Auxiliary Novel View Mask Prediction Task Input-view depth maps and mask images are intended 291 to help the model indirectly understand the structure of multi-object compositions by incorporating 292 additional structure-aware information into the input. To encourage the model to better grasp overall 293 structure, particularly its ability to generate it, we propose leveraging structural information (*i.e.*, mask 294 image) prediction under the target view as an auxiliary task, providing more direct supervision. Our 295 approach draws inspiration from classifier guidance (Dhariwal & Nichol, 2021), where a classifier 296 $p_{\phi}(y|x_t, t)$ guides the denoising process of image x_t to meet the criterion y via incorporating the gradient $\nabla_{x_t} \log p_{\phi}((y|x_t, t))$ during the inference process as an auxiliary guidance. Similarly, to 297 improve the model's ability to learn compositional structure, particularly in synthesizing novel view 298 plausible object placement (position and orientation), we introduce an auxiliary task during training: 299 predicting object mask images $\mathbf{M}_t \sim p(\mathbf{M}_t | \mathbf{x}_t, t, C_{SA}(\cdot))$ under target view. This prediction is 300 conditioned on the noisy target-view image x_t , timestep t and input-view structure-aware feature 301 $C_{SA}(\cdot)$, derived from the final layer of the denoising U-Net. The supervision could be formulated as: 302

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t, \mathbf{M}_t | t, C_{\mathsf{SA}}(\cdot)) = \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | t, C_{\mathsf{SA}}(\cdot)) + \nabla_{\mathbf{x}_t} \log p(\mathbf{M}_t | \mathbf{x}_t, t, C_{\mathsf{SA}}(\cdot)).$$
(4)

Following Eq. (4), we jointly train the mask predictor and denoising U-Net following:

$$\mathbb{E}[||\epsilon_{\theta}(\alpha_t \mathbf{x}_0 + \sigma_t \epsilon, t, C_{SA}(\cdot)) - \epsilon||_2^2 + \gamma ||\mathbf{M}_{tgt} - \mathbf{M}_t||_2^2],$$
(5)

where we use \mathbf{M}_{tgt} to denote the ground-truth target-view image, and we use the weight $\gamma = 0.1$ to balance the diffusion loss and mask prediction loss.

3.3 STRUCTURE-GUIDED TIMESTEP SAMPLING SCHEDULER

Inspired by previous works (Jiang et al., 2023; Chen, 2023) that identify the importance of different
 scheduling strategies, we first provide an in-depth analysis of the inference process of multi-object
 novel view synthesis, where we adopt a DDIM (Song et al., 2020a) sampler:

$$\mathbf{x}_{t-1} = \sqrt{\alpha_{t-1}} \underbrace{\left(\underbrace{\mathbf{x}_t - \sqrt{1 - \alpha_t} \cdot \mathbf{F}}_{\text{predicted } \mathbf{x}_0} \right)}_{\text{predicted } \mathbf{x}_0} + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \mathbf{F} + \sigma_t \epsilon_t.$$
(6)

317 318 319

287 288

289

303 304

305 306 307

310

311

315 316

We use **F** as a shorthand for $\epsilon_{\theta}(\mathbf{x}_t, t, C_{SA}(\cdot))$ and $\epsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. We examine the predicted \mathbf{x}_0 (as in Eq. (6)) and the predicted mask image \mathbf{M}_t at various timesteps during the inference process as they offer direct visualizations for analysis. These visualized results are presented in Fig. 3.

In Fig. 3, we observe that a blurry image, which indicates the approximate placement of each object, is quickly restored in the early stages (*i.e.*, larger t) of the inference process. This suggests

Figure 3: **Visualization of inference.** The early stage of the denoising process focuses on restoring global object placements, while the prediction of object masks requires a relatively noiseless image to recover fine-grained geometry. This motivates us to seek a balanced timestep sampling scheduler during training. The model trained w/shift yields better mask prediction and thus recovers an image with more details and sharp object boundary. The w/o shift here refers to not shifting the μ value.

that global structural information is prioritized for the model to learn during this stage. Accurate object placements are crucial for synthesizing reasonable novel view images, as incorrect placement predictions indicate a fundamental misunderstanding of the compositional structure. This underscores the importance of training the model with a larger t during the initial training periods, which is even more important for multi-object NVS scenarios considering the increased compositional complexity compared with a single object. Conversely, a mask with a clear boundary is not predicted until a later stage of the sampling process (*i.e.*, smaller t). This is because accurate mask prediction depends heavily on a relatively noiseless image. Therefore, to capture fine-grained geometry and appearance details of objects, it is essential to train the model with a smaller t during later training periods.

Recognizing the importance of timestep t in balancing the learning of global placement information and local fine-grained details, we propose to adjust the original timestep sampling process to:

$$t \sim \mathcal{U}(\mathbf{1}, \mathbf{1000}) \rightarrow t \sim \mathcal{N}(\mu(s), \sigma), \text{ where } \mu(s) = \mu_{\text{local}} + (\mu_{\text{global}} - \mu_{\text{local}}) \cdot \frac{s}{T_s}$$
 (7)

where s denotes the model training iteration, T_s denotes the total number of training steps, $\sigma = 200$ is a constant variance. We sample the timestep t from a Gaussian distribution with mean $\mu(s)$ following a linear decay from a large value $\mu_{global} = 1000$ to a small value $\mu_{local} = 500$. This approach allows the model to initially learn correct global object placement information and gradually turn its focus to refining detailed object geometry in later training stages. In practice, we include a warmup period with 4000 training steps sampling t with a fixed $\mu(s) = \mu_{global}$. After the warmup, we use the linear decay schedule over 2000 steps, and then stabilize the learning for fine-grained details after 6000 steps where we use $\mu(s) = \mu_{local}$.

4 EXPERIMENTS

- 375 4.1 EXPERIMENT SETUP
- 377 We focus on multi-object composite NVS tasks in indoor scenes, with an emphasis on foreground objects, examining novel view structural plausibility regarding object placement, geometry, appear-



ance, and cross-view consistency with input view. This choice stems from the recent advancements in
 object segmentation (Kirillov et al., 2023), while we leave the background modeling for future work.

Datasets. To facilitate the training and evaluation of our proposed method, we curate a scalable synthetic dataset Compositional 3D-FUTURE (C3DF), comprising 100k composites for training and 5k for testing. Each composite is created by composing pre-filtered furniture items from 3D-FUTURE (Fu et al., 2021b) using a heuristic strategy to avoid collision and penetration. Beyond C3DF, we emphasize testing the generalization capability by benchmarking our method on Room-Texture (Luo et al., 2024) and Objaverse (Deitke et al., 2023b). We also evaluate our model on diverse indoor scenes from both the synthetic dataset 3D-FRONT (Fu et al., 2021a) and the real-world dataset SUNRGB-D (Song et al., 2015). Refer to Appx. B.2 for more dataset details.

Baselines. We compare our method against two recent novel view synthesis methods including
Zero-1-to-3 (Liu et al., 2023c) and ZeroNVS (Sargent et al., 2023). The original Zero-1-to-3 is
trained on extensive object-level datasets. Therefore, we also re-train Zero-1-to-3 on our synthetic
dataset C3DF, denoted as Zero-1-to-3[†]. ZeroNVS is trained on a mixture of real-world images with
background, so we use images with backgrounds as its input if possible for a fair comparison.



Figure 4: Qualitative results of NVS and cross-view matching. Our method generates plausible
 novel-view images across various datasets, surpassing baselines regarding object placement, shape,
 and appearance. In cross-view matching, points of the same color indicate correspondences between
 the input and target views. We achieve a higher number of matched points with more precise locations.

435

450

462

463

470

471

472

473

474

475

476

477 478 479

Table 1: Quantitative results of multi-object NVS, Object Placement, and Cross-view Consistency. We evaluate on C3DF test set, along with generalization experiments on Room-Texture (Luo et al., 2024) and Objaverse (Deitke et al., 2023b). [†] indicates re-training on C3DF.

136	Dataset	Madha d	Nove	el View Synt	hesis	Placement	Cross-view Consistency	
137		Method	PSNR(↑)	SSIM(†)	$LPIPS(\downarrow)$	IoU(†)	Hit Rate(↑)	Dist(↓)
438	C3DF	ZeroNVS	10.704	0.533	0.481	21.6	1.2	130.3
439		Zero-1-to-3	14.255	0.771	0.302	33.7	5.8	86.9
440		Zero-1-to-3 [†]	14.811	0.794	0.283	34.4	1.6	120.3
441		Ours	17.432	0.825	0.171	58.1	37.0	44.8
442		ZeroNVS	8.217	0.647	0.487	8.2	1.2	140.3
443	Doom Taxtura	Zero-1-to-3	9.860	0.712	0.406	13.9	2.9	104.1
444	Room- Texture	Zero-1-to-3 [†]	8.342	0.657	0.452	13.5	0.5	157.4
445		Ours	10.014	0.718	0.366	24.2	6.1	78.1
446	Objaverse	ZeroNVS	10.557	0.513	0.486	17.3	2.3	126.9
447		Zero-1-to-3	15.850	0.810	0.259	34.7	10.7	80.7
448		Zero-1-to-3 [†]	15.433	0.815	0.273	29.7	1.7	126.7
449		Ours	17.749	0.840	0.169	51.3	50.0	47.2

451 Metrics. We utilize PSNR, SSIM, and LPIPS as metrics for evaluating the quality of Novel View 452 Synthesis. To assess global object *Placement*, we compute the foreground-background IoU with ground-truth masks. Finally, we propose metrics to evaluate Cross-view Consistency with image-453 matching. More specifically, we first apply MASt3R (Leroy et al., 2024) to acquire the image 454 matching between the input-view image and target-view image for both ground truth and model 455 predictions. With the ground-truth matching as references, we compute each method's Hit Rate and 456 the nearest matching distance (Dist.). Hit Rate measures the proportion of predicted matches that 457 align with the ground truth matches. Dist. quantifies the distance between the predicted matching and 458 ground-truth matching in the target view. Please refer to Appx. B.3 for more details about the metrics. 459

460 4.2 **RESULTS AND DISCUSSIONS** 461

Fig. 4 presents qualitative results of multi-object NVS and cross-view matching visualization on different datasets, with quantitative results in Tab. 1. We summarize the following key observations:

- 464 1. Our method realizes the highest PSNR and generates high-quality images under novel views, 465 closely aligned with the ground truth images, especially regarding novel-view object placement 466 (position and orientation), shape, and appearance. In contrast, the baseline models struggle to 467 accurately capture the compositional structure under novel views. For example, in the first row, the 468 red bed is incorrectly oriented in Zero-1-to-3 and is either missing or distorted in other baselines. 469
 - 2. From the visualized cross-view matching results and the metrics in Tab. 1, it is evident that our method significantly outperforms the baseline approaches in *Cross-view Consistency*. It achieves a much higher IoU and *Hit Rate* while exhibiting a considerably lower matching distance. The visualized results are consistent with the metrics, further validating our method's accuracy in capturing cross-view structural consistency, which cannot be reflected by existing NVS metrics.
 - 3. Our model exhibits strong generalization capabilities on unseen datasets, e.g., Room-Texture and Objaverse. We demonstrate more qualitative results, including on 3D-FRONT and SUNRGB-D, in Appx. B.4. We showcase potential applications, including object removal and reconstruction in Appx. B.5. Further discussion about limitation and failure cases are presented in Appx. C.

180	Dataset	Madha d	Nove	Novel View Synthesis		Placement	Cross-view Consistency	
82		Method	$\overline{\text{PSNR}(\uparrow) \text{SSIM}(\uparrow) \text{LPIPS}(\downarrow)}$	IoU(†)	Hit Rate(↑)	Dist(↓)		
183	C3DF	w/o depth	17.080	0.819	0.178	57.2	39.2	45.2
484		w/o mask	16.914	0.818	0.187	54.7	25.4	50.4
485		w/o sch. Ours	16.166 17.432	0.808 0.825	0.212 0.171	49.1 58.1	11.9 37.0	48.6 44.8

Table 2:	Ablation	results	on	C3DF
----------	----------	---------	----	------



Figure 5: Qualitative comparison for ablation study. Removing depth or mask predictions weakens the model's understanding of object placement and existence, exemplified by the missing white cabinet in the first example. Excluding mask predictions or the scheduler reduces the model's ability to learn object placement, as shown by the misoriented brown cabinet in the second example.

4.3 ABLATION STUDY

To verify the efficacy of each component, we perform an ablation study on our key technical designs, including the depth input (w/o depth), mask prediction auxiliary task (w/o mask), and the scheduler (w/o sch. learns with a uniform sampler $t \sim \mathcal{U}(1, 1000)$). Results in Tab. 2 show that the auxiliary mask prediction task and the timestep sampler are the most critical components, significantly affecting all the metrics and the realistic object recovery as demonstrated by the misoriented brown cabinet in the second example from Fig. 5. Without the scheduler, the model produces less accurate object positions, evident both qualitatively and quantitatively. Furthermore, removing depth or mask predictions weakens the model's understanding of spatial relationships and object existence, exemplified by the completely missing white cabinet in the first example. This also shows incorporating structure-aware features as inputs, though seemingly intuitive, offers the most straightforward approach to enhancing the model's structural awareness, particularly given recent advancements in monocular predictors (Kirillov et al., 2023; Ke et al., 2024). We present a more comprehensive discussion on the scheduler strategy in Appx. A.4 and ablations on more datasets are in Appx. B.4.

CONCLUSION

We extend diffusion-based novel view synthesizers to handle multi-object compositions in indoor scenes. Our proposed model generalize well across diverse datasets with more accurate object placement, shape, and appearance, showing a stronger cross-view consistency with input view. The core of our approach lies in integrating structure-aware features as additional inputs, an auxiliary mask prediction task, and a structure-guided timestep sampling scheduler. These components enhance the model's awareness of compositional structure while balancing the learning of global object placement and fine-grained local shape and appearance. Given the prevalence of multi-object compositions in real-world scenes, we believe that our model designs and comprehensive evaluations can offer valuable insights for advancing scene-level NVS models in more complex environments.

540 REFERENCES

542 543 544	Garrick Brazil, Abhinav Kumar, Julian Straub, Nikhila Ravi, Justin Johnson, and Georgia Gkioxari. Omni3d: A large benchmark and model for 3d object detection in the wild. In <i>Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pp. 13154–13164, 2023.
545 546 547	David Charatan, Sizhe Li, Andrea Tagliasacchi, and Vincent Sitzmann. pixelsplat: 3d gaussian splats from image pairs for scalable generalizable 3d reconstruction. In <i>arXiv</i> , 2023.
548	Ting Chen. On the importance of noise scheduling for diffusion models, 2023.
549 550 551 552	Yabo Chen, Jiemin Fang, Yuyang Huang, Taoran Yi, Xiaopeng Zhang, Lingxi Xie, Xinggang Wang, Wenrui Dai, Hongkai Xiong, and Qi Tian. Cascade-zero123: One image to highly consistent 3d with self-prompted nearby views. <i>arXiv preprint arXiv:2312.04424</i> , 2023.
553 554 555	Yixin Chen, Junfeng Ni, Nan Jiang, Yaowei Zhang, Yixin Zhu, and Siyuan Huang. Single-view 3d scene reconstruction with high-fidelity shape and texture. In <i>International Conference on 3D Vision (3DV)</i> , 2024a.
555 557 558 559	Yongwei Chen, Tengfei Wang, Tong Wu, Xingang Pan, Kui Jia, and Ziwei Liu. Comboverse: Compositional 3d assets creation using spatially-aware diffusion guidance. <i>arXiv preprint</i> <i>arXiv:2403.12409</i> , 2024b.
560 561 562	Yuedong Chen, Haofei Xu, Chuanxia Zheng, Bohan Zhuang, Marc Pollefeys, Andreas Geiger, Tat-Jen Cham, and Jianfei Cai. Mvsplat: Efficient 3d gaussian splatting from sparse multi-view images. <i>arXiv preprint arXiv:2403.14627</i> , 2024c.
563 564 565	Blender Online Community. <i>Blender - a 3D modelling and rendering package</i> . Blender Foundation, Stichting Blender Foundation, Amsterdam, 2018. URL http://www.blender.org.
566 567 568	Matt Deitke, Ruoshi Liu, Matthew Wallingford, Huong Ngo, Oscar Michel, Aditya Kusupati, Alan Fan, Christian Laforte, Vikram Voleti, Samir Yitzhak Gadre, et al. Objaverse-xl: A universe of 10m+ 3d objects. <i>arXiv preprint arXiv:2307.05663</i> , 2023a.
570 571 572	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 <i>CVPR</i> , pp. 13142–13153, 2023b.
573 574 575	Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. Advances in Neural Information Processing Systems (NeurIPS), 34:8780–8794, 2021.
576 577	Andreea Dogaru, Mert Özer, and Bernhard Egger. Generalizable 3d scene reconstruction via divide and conquer from a single view. <i>arXiv preprint arXiv:2404.03421</i> , 2024.
578 579 580 581	Huan Fu, Bowen Cai, Lin Gao, Ling-Xiao Zhang, Jiaming Wang, Cao Li, Qixun Zeng, Chengyue Sun, Rongfei Jia, Binqiang Zhao, et al. 3d-front: 3d furnished rooms with layouts and semantics. In <i>International Conference on Computer Vision (ICCV)</i> , pp. 10933–10942, 2021a.
582 583 584	Huan Fu, Rongfei Jia, Lin Gao, Mingming Gong, Binqiang Zhao, Steve Maybank, and Dacheng Tao. 3d-future: 3d furniture shape with texture. <i>International Journal of Computer Vision (IJCV)</i> , 129: 3313–3337, 2021b.
586 587 588	Ming Gui, Johannes S Fischer, Ulrich Prestel, Pingchuan Ma, Dmytro Kotovenko, Olga Grebenkova, Stefan Andreas Baumann, Vincent Tao Hu, and Björn Ommer. Depthfm: Fast monocular depth estimation with flow matching. <i>arXiv preprint arXiv:2403.13788</i> , 2024.
589 590 591	Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. <i>NeurIPS</i> , 33: 6840–6851, 2020.
592 593	Zehuan Huang, Hao Wen, Junting Dong, Yaohui Wang, Yangguang Li, Xinyuan Chen, Yan-Pei Cao, Ding Liang, Yu Qiao, Bo Dai, et al. Epidiff: Enhancing multi-view synthesis via localized epipolar-constrained diffusion. <i>arXiv preprint arXiv:2312.06725</i> , 2023.

594 595 596	Ajay Jain, Matthew Tancik, and Pieter Abbeel. Putting nerf on a diet: Semantically consistent few-shot view synthesis. In <i>International Conference on Computer Vision (ICCV)</i> , pp. 5885–5894, 2021.
597 598 599	Ajay Jain, Ben Mildenhall, Jonathan T Barron, Pieter Abbeel, and Ben Poole. Zero-shot text-guided object generation with dream fields. In <i>CVPR</i> , pp. 867–876, 2022.
600 601 602	Yifan Jiang, Hao Tang, Jen-Hao Rick Chang, Liangchen Song, Zhangyang Wang, and Liangliang Cao. Efficient-3dim: Learning a generalizable single-image novel-view synthesizer in one day. <i>arXiv preprint arXiv:2310.03015</i> , 2023.
603 604 605 606	Bingxin Ke, Anton Obukhov, Shengyu Huang, Nando Metzger, Rodrigo Caye Daudt, and Konrad Schindler. Repurposing diffusion-based image generators for monocular depth estimation. In <i>Conference on Computer Vision and Pattern Recognition (CVPR)</i> , 2024.
607 608	Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. <i>ACM Transactions on Graphics</i> , 2023.
609 610 611	Diederik Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. <i>NeurIPS</i> , 34:21696–21707, 2021.
612 613 614	Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In <i>ICCV</i> , pp. 4015–4026, 2023.
615 616 617	Xin Kong, Shikun Liu, Xiaoyang Lyu, Marwan Taher, Xiaojuan Qi, and Andrew J Davison. Eschernet: A generative model for scalable view synthesis. <i>arXiv preprint arXiv:2402.03908</i> , 2024.
618 619	Vincent Leroy, Yohann Cabon, and Jérôme Revaud. Grounding image matching in 3d with mast3r. arXiv preprint arXiv:2406.09756, 2024.
620 621 622 623	Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In <i>ICML</i> , pp. 19730–19742. PMLR, 2023.
624 625 626	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 <i>CVPR</i> , pp. 300–309, 2023a.
627 628 629 630	Yukang Lin, Haonan Han, Chaoqun Gong, Zunnan Xu, Yachao Zhang, and Xiu Li. Consistent123: One image to highly consistent 3d asset using case-aware diffusion priors. <i>arXiv preprint arXiv:2309.17261</i> , 2023b.
631 632 633	Minghua Liu, Ruoxi Shi, Linghao Chen, Zhuoyang Zhang, Chao Xu, Xinyue Wei, Hansheng Chen, Chong Zeng, Jiayuan Gu, and Hao Su. One-2-3-45++: Fast single image to 3d objects with consistent multi-view generation and 3d diffusion. <i>arXiv preprint arXiv:2311.07885</i> , 2023a.
634 635 636 637	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. <i>arXiv preprint arXiv:2306.16928</i> , 2023b.
638 639 640	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 <i>International Conference on Computer Vision</i> (<i>ICCV</i>), 2023c.
641 642 643 644	Yuan Liu, Cheng Lin, Zijiao Zeng, Xiaoxiao Long, Lingjie Liu, Taku Komura, and Wenping Wang. Syncdreamer: Generating multiview-consistent images from a single-view image. <i>arXiv preprint</i> <i>arXiv:2309.03453</i> , 2023d.
645 646 647	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. In <i>Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pp. 9970–9980, 2024.

648 649 650	Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In CVPR, pp. 11461–11471, 2022.
651 652 653	Rundong Luo, Hong-Xing Yu, and Jiajun Wu. Unsupervised discovery of object-centric neural fields. arXiv preprint arXiv:2402.07376, 2024.
654 655 656	Luke Melas-Kyriazi, Iro Laina, Christian Rupprecht, and Andrea Vedaldi. Realfusion: 360deg reconstruction of any object from a single image. In <i>Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pp. 8446–8455, 2023.
658 659 660	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 <i>European</i> <i>Conference on Computer Vision (ECCV)</i> , 2020.
661 662 663 664	Yinyu Nie, Xiaoguang Han, Shihui Guo, Yujian Zheng, Jian Chang, and Jian Jun Zhang. To- tal3dunderstanding: Joint layout, object pose and mesh reconstruction for indoor scenes from a single image. In <i>Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pp. 55–64, 2020.
665 666 667 668	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. <i>arXiv preprint arXiv:2304.07193</i> , 2023.
669 670 671	Ege Ozguroglu, Ruoshi Liu, Dídac Surís, Dian Chen, Achal Dave, Pavel Tokmakov, and Carl Vondrick. pix2gestalt: Amodal segmentation by synthesizing wholes. In <i>Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pp. 3931–3940, 2024.
672 673 674	Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. <i>arXiv preprint arXiv:2209.14988</i> , 2022.
675 676 677	Guocheng Qian, Jinjie Mai, Abdullah Hamdi, Jian Ren, Aliaksandr Siarohin, Bing Li, Hsin-Ying Lee, Ivan Skorokhodov, Peter Wonka, Sergey Tulyakov, et al. Magic123: One image to high-quality 3d object generation using both 2d and 3d diffusion priors. <i>arXiv preprint arXiv:2306.17843</i> , 2023.
678 679 680	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 <i>ICML</i> , pp. 8748–8763. PMLR, 2021.
682 683	Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High- resolution image synthesis with latent diffusion models. In <i>CVPR</i> , pp. 10684–10695, 2022.
684 685 686 687	Kyle Sargent, Zizhang Li, Tanmay Shah, Charles Herrmann, Hong-Xing Yu, Yunzhi Zhang, Eric Ryan Chan, Dmitry Lagun, Li Fei-Fei, Deqing Sun, et al. Zeronvs: Zero-shot 360-degree view synthesis from a single real image. <i>arXiv preprint arXiv:2310.17994</i> , 2023.
688 689 690	Ruoxi Shi, Hansheng Chen, Zhuoyang Zhang, Minghua Liu, Chao Xu, Xinyue Wei, Linghao Chen, Chong Zeng, and Hao Su. Zero123++: a single image to consistent multi-view diffusion base model. <i>arXiv preprint arXiv:2310.15110</i> , 2023a.
691 692 693	Yichun Shi, Peng Wang, Jianglong Ye, Long Mai, Kejie Li, and Xiao Yang. Mvdream: Multi-view diffusion for 3d generation. <i>arXiv:2308.16512</i> , 2023b.
694 695	Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In Interna- tional Conference on Learning Representations (ICLR), 2020a.
696 697 698 699	Shuran Song, Samuel P Lichtenberg, and Jianxiong Xiao. Sun rgb-d: A rgb-d scene understanding benchmark suite. In <i>Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pp. 567–576, 2015.
700 701	Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. <i>arXiv preprint arXiv:2011.13456</i> , 2020b.

- 702 Jiaxiang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang Zeng. Dreamgaussian: Generative 703 gaussian splatting for efficient 3d content creation. arXiv preprint arXiv:2309.16653, 2023a. 704 705 Junshu Tang, Tengfei Wang, Bo Zhang, Ting Zhang, Ran Yi, Lizhuang Ma, and Dong Chen. Make-it-3d: High-fidelity 3d creation from a single image with diffusion prior. arXiv preprint 706 arXiv:2303.14184, 2023b. 707 708 Joseph Tung, Gene Chou, Ruojin Cai, Guandao Yang, Kai Zhang, Gordon Wetzstein, Bharath 709 Hariharan, and Noah Snavely. Megascenes: Scene-level view synthesis at scale. arXiv preprint 710 arXiv:2406.11819, 2024. 711 Basile Van Hoorick, Pavel Tokmakov, Simon Stent, Jie Li, and Carl Vondrick. Tracking through 712 containers and occluders in the wild. In Conference on Computer Vision and Pattern Recognition 713 (CVPR), pp. 13802–13812, 2023. 714 715 Basile Van Hoorick, Rundi Wu, Ege Ozguroglu, Kyle Sargent, Ruoshi Liu, Pavel Tokmakov, Achal 716 Dave, Changxi Zheng, and Carl Vondrick. Generative camera dolly: Extreme monocular dynamic novel view synthesis. arXiv preprint arXiv:2405.14868, 2024. 717 718 Peng Wang and Yichun Shi. Imagedream: Image-prompt multi-view diffusion for 3d generation. 719 arXiv preprint arXiv:2312.02201, 2023. 720 Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul Srinivasan, Howard Zhou, Jonathan T. Barron, 721 Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. Ibrnet: Learning multi-view 722 image-based rendering. In Conference on Computer Vision and Pattern Recognition (CVPR), 2021. 723 724 Shuzhe Wang, Vincent Leroy, Yohann Cabon, Boris Chidlovskii, and Jerome Revaud. Dust3r: 725 Geometric 3d vision made easy. In Conference on Computer Vision and Pattern Recognition 726 (CVPR), pp. 20697-20709, 2024. 727 Daniel Watson, William Chan, Ricardo Martin-Brualla, Jonathan Ho, Andrea Tagliasacchi, and Mo-728 hammad Norouzi. Novel view synthesis with diffusion models. arXiv preprint arXiv:2210.04628, 729 2022. 730 731 Haohan Weng, Tianyu Yang, Jianan Wang, Yu Li, Tong Zhang, CL Chen, and Lei Zhang. Consistent123: Improve consistency for one image to 3d object synthesis. arXiv preprint 732 arXiv:2310.08092, 2023. 733 734 Sanghyun Woo, Shoubhik Debnath, Ronghang Hu, Xinlei Chen, Zhuang Liu, In So Kweon, and 735 Saining Xie. Convnext v2: Co-designing and scaling convnets with masked autoencoders. In 736 Conference on Computer Vision and Pattern Recognition (CVPR), pp. 16133–16142, 2023. 737 Kailu Wu, Fangfu Liu, Zhihan Cai, Runjie Yan, Hanyang Wang, Yating Hu, Yueqi Duan, and 738 Kaisheng Ma. Unique3d: High-quality and efficient 3d mesh generation from a single image. 739 arXiv preprint arXiv:2405.20343, 2024. 740 741 Dejia Xu, Yifan Jiang, Peihao Wang, Zhiwen Fan, Yi Wang, and Zhangyang Wang. Neurallift-360: 742 Lifting an in-the-wild 2d photo to a 3d object with 360deg views. In CVPR, pp. 4479–4489, 2023. 743 Katherine Xu, Lingzhi Zhang, and Jianbo Shi. Amodal completion via progressive mixed context 744 diffusion. In Conference on Computer Vision and Pattern Recognition (CVPR), pp. 9099–9109, 745 2024. 746 747 Jiayu Yang, Ziang Cheng, Yunfei Duan, Pan Ji, and Hongdong Li. Consistnet: Enforcing 3d 748 consistency for multi-view images diffusion. In Conference on Computer Vision and Pattern Recognition (CVPR), pp. 7079–7088, 2024a. 749 750 Lihe Yang, Bingyi Kang, Zilong Huang, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth 751 anything: Unleashing the power of large-scale unlabeled data. In Conference on Computer Vision 752 and Pattern Recognition (CVPR), pp. 10371-10381, 2024b. 753 Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields 754 from one or few images. In Conference on Computer Vision and Pattern Recognition (CVPR), pp. 755 4578-4587, 2021.
 - 14

756 757 758 750	Guanqi Zhan, Chuanxia Zheng, Weidi Xie, and Andrew Zisserman. Amodal ground truth and completion in the wild. In <i>Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pp. 28003–28013, 2024.
759 760 761	Xiaohang Zhan, Xingang Pan, Bo Dai, Ziwei Liu, Dahua Lin, and Chen Change Loy. Self-supervised scene de-occlusion. In <i>Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pp. 2704, 2702, 2000.
762	3784-3792, 2020.
763	Yan Zhu, Yuandong Tian, Dimitris Metaxas, and Piotr Dollár, Semantic amodal segmentation. In
764	Conference on Computer Vision and Pattern Recognition (CVPR), pp. 1464–1472, 2017.
765	
766	
767	
768	
769	
770	
771	
772	
773	
774	
775	
776	
777	
778	
779	
780	
781	
782	
783	
784	
785	
700	
700	
790	
709	
791	
792	
793	
794	
795	
796	
797	
798	
799	
800	
801	
802	
803	
804	
805	
806	
807	
808	
809	

⁸¹⁰ A MODEL DETAILS

812 A.1 DINO PATCH FEATURE AND CAMERA VIEW EMBEDDING 813

814 The original image encoder of Stable Diffusion is CLIP, which excels at aligning images with text. Other image encoders like DINO-v2 (Oquab et al., 2023) or ConvNeXtv2 (Woo et al., 2023) 815 may provide denser image features that may benefit generation tasks as mentioned by previous 816 works (Jiang et al., 2023; Kong et al., 2024). Therefore, we opt to use the DINO feature instead of 817 the original CLIP feature in our network following Jiang et al. (2023). To inject the DINO patch 818 feature into our network, we encode the input view image using DINO-v2 (Oquab et al., 2023) "norm 819 patchtokens", whose shape dimension is (b, 16, 16, 1024). We will simply flatten it into (b, 256, 1024)820 to apply cross-attention, and b means batch size here. 821

As for the camera view embedding, we choose to embed it using a 6 degrees of freedom (6DoF) 822 representation. To be specific, let E_i be the extrinsic matrix under the input view and E_i be the 823 extrinsic matrix under the output view, we represent relative camera pose change as $E_i^{-1}E_i$. We will 824 also flatten it into 16 dimensions to concatenate it to the image feature. Afterwards, we will replicate 825 the 16-dimension embedding 256 times to concatenate the embedding to every channel of the DINO 826 feature map. A projection layer will later be employed to project the feature map into (b, 256, 768)827 to match the dimension of the CLIP encoder, which was originally used by Stable Diffusion so 828 that we can fine-tune the pre-trained checkpoint. It is worth noting that we also tried other novel 829 view synthesizer's camera embedding like Zero-1-to-3 (Liu et al., 2023c) using a 3DoF spherical 830 coordinates in early experiments, but we found that it does not make much of a difference.

831 832

833

A.2 DEPTH AND MASK CONDITION

In this section, we will explain how input view depth and mask are incorporated as additional conditioning inputs. For depth maps, regions with infinite depth values are assigned a value equal to twice the maximum finite depth value in the rest of the image. After this adjustment, we apply a normalization technique to scale the depth values to the range of [-1, 1], enabling the use of the same VAE architecture as for images.

For mask images, we assign unique values to different object instances in the input view. For
instance, if there are four objects in the multi-object composite, they will be labeled as 1, 2, 3, and 4,
respectively, while the background will be assigned a value of 0. The same normalization technique
used for depth maps is applied to these mask images. These mask images, like all other inputs, are
processed by the VAE, with all images set to a resolution of 256 × 256.

844 845

A.3 SUPERVISION FOR AUXILIARY MASK PREDICTION TASK

To implement the auxiliary mask prediction task, we encode the output view mask images into the same latent space as the input view mask images. Object instances viewed from different angles will be assigned the same value, which is ensured during the curation of our compositional dataset. Supervision is directly applied to the latent mask features extracted from the final layer of the denoising U-Net. Only the input view mask images are required during inference, simplifying the process while preserving consistency across views.

852 853

854

A.4 TIMESTEP SCHEDULER

Though we finally employed a linearly declining strategy, we experimented with several alternatives. Specifically, we tested linearly declining the mean of the Gaussian distribution (LDC), linearly increasing the mean after a sudden drop (LIND), and keeping the mean constant (KMS). These strategies are illustrated in Fig. A1. The metrics on the test set of our C3DF are provided in Tab. A1, with some visual comparisons in Fig. A2. *w/o sch.* in Tab. A1 refers to applying a uniform sampler, same as the one in the main paper. From the results, we observe that LDC achieves slightly better performance than LIND and KMS, largely outperforming *w/o sch.*

However, we observed significant visual artifacts such as weird colors and extremely blurry mask
 images when combining the auxiliary mask prediction task with the KMS sampling strategy, as
 demonstrated in Fig. A2. For example, the bed in the second example possesses unclear object

Mean 4000 6000 8000 Steps

Table A1: **Ablation on different strategies.** Incorporating sampling strategies significantly improves the model performance, while the linear decline (LDC) achieves the best.

Dataat		Novel View Synthesis				
Dataset	Method	PSNR(↑)	SSIM(†)	$LPIPS(\downarrow)$		
	w/o sch.	16.166	0.808	0.212		
C2DE	KMS	17.148	0.823	0.175		
CODF	LIND	17.279	0.824	0.172		
	LDC	17.432	0.825	0.171		

Figure A1: Illustration of different timestep sampling strategies.



Figure A2: **Comparison of different strategies**. The predicted images and mask images under novel views using different strategies are visualized. We can observe that images predicted by the KMS strategy possess weird and blurry color while LDC strategy seems to be slightly better than LIND.

boundaries and distorted texture. We believe this is due to KMS focusing primarily on denoising at larger timesteps, which provides limited guidance for recovering mask images and refining finegrained geometry and appearance. Consequently, without a dedicated period for denoising smaller timesteps, the per-object shape and appearance appear distorted and unrealistic.

B EXPERIMENT DETAILS

914 B.1 IMPLEMENTATION DETAILS

We solely utilize the data from C3DF as the training source for our model. The training process takes
 around 2 days on 8 NVIDIA A100 (80G) GPUs, employing a batch size of 172 per GPU. The exact training steps are 8,000 steps. During the inference process, we apply 50 DDIM steps and set the

guidance scale to 3.0. We use DepthFM (Gui et al., 2024) and SAM (Kirillov et al., 2023) to extract the depth maps and object masks when they are not available, as well as for all real-world images.

920 921 922

B.2 DATASETS

C3DF We use the furniture models from the 3D-FUTURE dataset (Fu et al., 2021b) to create our synthetic multi-object compositional data. The 3D-FUTURE dataset contains 9,992 detailed 3D furniture models with high-resolution textures and labels. Following previous work Chen et al. (2024a), we categorize the furniture into seven groups: bed, bookshelf, cabinet, chair, nightstand, sofa, and table. To ensure unbiased evaluation, we further split the furniture into distinct training and test sets, ensuring that none of the test set items are seen during training.

929 After filtering the furniture, we first determine the number of pieces to include in each composite, 930 which is randomly selected to be between 3 and 6. Next, we establish a probability distribution based 931 on the different types of furniture items and sample each piece according to this distribution. To 932 prevent collisions and penetration between furniture items, we employ a heuristic strategy. Specifi-933 cally, for each furniture item to be added, we apply a random scale adjustment within the range of 934 [0.95, 1.05], as the inherent scale of the furniture models accurately reflects real-world sizes. We also 935 rotate each model by a random angle to introduce additional variability. Once these adjustments are 936 complete, we begin placing the furniture items in the scene. The first item is positioned at the center 937 of the scene at coordinates (0, 0, 0). Subsequent objects are added one by one, initially placed at the same central location. Since this results in inevitable collisions, we randomly sample a direction 938 and gradually move the newly added item along this vector until there is no intersection between 939 the bounding boxes of the objects. By following these steps, we generate a substantial number of 940 multiple furniture items composites, ultimately creating a training set of 100,000 composites and a 941 test set of 5,000 to evaluate the capabilities of our network. 942

After placing all the furniture items, we render multi-view images to facilitate training, using
Blender (Community, 2018) as our renderer due to its high-quality output. We first normalize each
composite along its longest axis. To simulate real-world camera poses and capture meaningful
multi-object compositions, we employ the following method for sampling camera views.

Cameras are randomly sampled using spherical coordinates, with a radius range of [1.3, 1.7] and an elevation angle range of $[2^{\circ}, 40^{\circ}]$. There are no constraints on the azimuth angle, allowing the camera to rotate freely around multiple objects. The chosen ranges for the radius and elevation angles are empirical. In addition to determining the camera positions, we establish a "look-at" point to compute the camera pose. This point is randomly selected on a spherical shell with a radius range of [0.01, 0.2].

To enhance the model's compositional structural awareness, we also render depth maps and instance
masks (both occluded and unoccluded) from 12 different viewpoints. The unoccluded instance mask
ensures that if one object is blocked by another, the complete amodal mask of the occluded object
is still provided, regardless of any obstructions. Although these unoccluded instance masks are not
currently necessary for our network, we render them for potential future use.

958 959

960

961

962

963

964 965 966

Objaverse To evaluate our network's generalization capability, we create a small dataset comprising 300 composites sourced from Objaverse (Deitke et al., 2023b). Specifically, we utilize the provided LVIS annotations to select categories that are commonly found in indoor environments, such as beds, chairs, sofas, dressers, tables, and others. Since the meshes from Objaverse vary in scale, we rescale each object based on reference object scales from the 3D-FUTURE dataset (Fu et al., 2021b). The composition and rendering processes follow the same strategy employed in C3DF.

501						
968		C3DFS	Room-Texture	Objaverse	SUNRGB-D	3D-FRONT
969	denth	.(~		~	~
970	mask	v ./		v	~	~
971	mask	v	v	v	~	~

Inference Details Since our model requires input-view depth map and mask images as additional inputs, we need to use DepthFM (Gui et al., 2024) and SAM (Kirillov et al., 2023) to extract the depth maps and object masks when they are not available, as well as for all real-world images. We show whether all the used datasets have provided depth maps and mask images in Tab. A2. '×' means they do not provide such conditions while '√' means they do provide such conditions.

977 978

979 980

981

982

983

984

B.3 METRICS

Intersection over Union (IoU) Since all baseline methods do not possess the concept of every object instance, we compute a foreground-background IoU for comparison. This metric can provide a rough concept of the overall placement alignment with ground truth images. We extract the foreground object mask by converting the generated image to grayscale (I_L) . Given that the generated image has a white background, we compute the foreground mask \mathbf{M} as $\mathbf{M} = I_L < \beta_{th}$, where β_{th} is a threshold that is fixed as 250.

985 986 987

988 Cross-view Matching As outlined in the main paper, we introduce two metrics to systematically
 989 evaluate cross-view consistency with the input view: Hit Rate and Nearest Matching Distance.
 990 Since direct assessment of cross-view consistency is not feasible by merely evaluating the success
 991 matches between each method's predicted novel view images and the input view image, we opt to
 992 how far the predicted matches deviate from the ground-truth matches.

We first compute ground-truth matching points and every model's matching points using MASt3R (Leroy et al., 2024) upon the input view image and the output view image (ground truth or predicted). Each matching pair is represented as a four-element tuple (x^0, y^0, x^1, y^1) , where (x^0, y^0) corresponds to the point on the input-view image, and (x^1, y^1) corresponds to the point on the output-view image.

For each ground-truth matching pair $(\mathbf{x}_{gt}^0, \mathbf{y}_{gt}^0, \mathbf{x}_{gt}^1, \mathbf{y}_{gt}^1)$, we find the nearest predicted matching pair in each model's results, denoted as $(\mathbf{x}^0, \mathbf{y}^0, \mathbf{x}^1, \mathbf{y}^1)$, based on the Euclidean distance between points in the input view image. If both $\mathbf{L}_2||(\mathbf{x}_{gt}^0, \mathbf{y}_{gt}^0), (\mathbf{x}^0, \mathbf{y}^0)||$ and $\mathbf{L}_2||(\mathbf{x}_{gt}^1, \mathbf{y}_{gt}^1), (\mathbf{x}^1, \mathbf{y}^1)||$ is smaller than a fixed threshold 20, the match is considered a successful hit. The **Hit Rate** is then calculated as the ratio of successful hits to the total number of ground-truth matches.

For Nearest Matching Distance, we examine whether $\mathbf{L}_2||(\mathbf{x}_{gt}^0, \mathbf{y}_{gt}^0), (\mathbf{x}^0, \mathbf{y}^0)||$ is within the threshold. For those passing this check, we compute the mean distance $\mathbf{L}_2||(\mathbf{x}_{gt}^1, \mathbf{y}_{gt}^1), (\mathbf{x}^1, \mathbf{y}^1)||$ as the **Nearest Matching Distance**, averaging over all successful hits. A detailed pseudo-code explanation can be found in Alg. 1 and Alg. 2.

1007 1008

Algorithm 1 Hit Rate Computation 1009 1: // Obtain image-matching pairs using MASt3R and save in a list 1010 2: $Pairs_{gt} = MASt3R(GT)$ 1011 3: $Pairs_{ours} = MASt3R(Ours)$ 1012 4: // Each element in the list is a four-element tuple $\mathbf{p} = (\mathbf{x}^0, \mathbf{y}^0, \mathbf{x}^1, \mathbf{y}^1)$ 5: // $(\mathbf{x}^0, \mathbf{y}^0)$ refers to the point in the input view image and $(\mathbf{x}^1, \mathbf{y}^1)$ the point in output view image 1013 1014 6: Hits = 01015 7: For p_{ot} in Pairs_{gt} 1016 // p_{ours}^i is the *i*-th element of Pairs_{ours} // p[:2] refers to the first two element in the tuple and p[2:] the last two 8: 1017 9: 1018 $i^{\star} = \arg\min(\mathbf{L}_2(\mathbf{p}_{\mathsf{gt}}[:2], \mathbf{p}_{\mathsf{ours}}^i[:2]))$ 10: IF $L_2(p_{gt}[:2], p_{ours}^{i^*}[:2]) < 20$ and $L_2(p_{gt}[2:], p_{ours}^{i^*}[2:]) < 20$ 11: 1020 // Successfully hit one, delete it from gt pairs and ours pairs 12: 1021 $\text{Hits} \gets \text{Hits} + 1$ 13: 1022 **POP**(Pairs_{ours}, p_{ours}^{i}) 14: 1023 15: return Hits/len(Pairs_{gt}) 1024 1025

Alg	orithm 2 Nearest Matching Distance Computation
1:	// The notations are the same as the one in Alg. 1
2:	$Pairs_{gt} = MASt3R(GT)$
3:	$Pairs_{ours} = MASt3R(Ours)$
4:	Dist = EmptyList()
5:	For p _{gt} in Pairs _{gt}
6:	$i^{\star} = \arg\min_{i}(\mathbf{L}_{2}(p_{gt}[:2], p_{ours}^{i}[:2]))$
7:	IF $L_2(p_{gt}[:2], p_{ours}^{i^*}[:2]) < 20$
8:	Append (Dist, $\mathbf{L}_2(p_{gt}[2:], p_{ours}^{i^*}[2:]))$
9:	POP (Pairs _{ours} , $p_{ours}^{i^{\star}}$)
10:	return Mean(Dist)



1055 Figure A3: Visualized comparison with baselines. Our method synthesizes more consistent novel view images and can even hallucinate objects that exceed the edge of image as shown in the first row. 1056 Conversely, baselines may predict unclear object boundary and omit objects under novel views. 1057

1060 **B.4 RESULTS**

We show more visualized results of our own methods along with ground truth on C3DF in Fig. A11, 1062 on Objaverse (Deitke et al., 2023b) in Fig. A12, and on Room-Texture (Luo et al., 2024) in Fig. A13. 1063 More visualized comparisons with baselines on SUNRGB-D (Song et al., 2015) and 3D-FRONT 1064 Fu et al. (2021a) are shown in Fig. A3. A more complete ablation study on other datasets including Objaverse and Room-Texture is shown in Tab. A3. Some continuous rotation examples on SUNRGB-1066 D are shown in Fig. A4, on 3D-FRONT are shown in Fig. A5, and more cross-view matching results 1067 without ground-truth pairs as reference are shown in Fig. A6. 1068

1069

1058 1059

1061

B.5 APPLICATIONS 1070

1071 **Object Removal** Since we can predict mask images under novel views, we can support simple 1072 image editing tasks like novel view object removal by simply setting a threshold value in the mask 1073 image and mask out corresponding pixels to achieve object removal. An example is shown in Fig. A8. 1074

1075

1076 **Reconstruction** The capability to synthesize novel view images that are consistent with the input 1077 view image demonstrates that the model possesses 3D-awareness, which can assist downstream tasks such as reconstruction. We leverage an off-the-shelf reconstruction method DUSt3R (Wang et al., 1078 2024) using the input-view image and novel view images predicted by our method. Two visualized 1079 examples are shown in Fig. A9.



Figure A4: Continuous rotation examples on SUNRGB-D. We rotate the camera around the multi-object composites, successfully synthesizing plausible novel-view images across a wide range of camera pose variations. 1108



Figure A5: Continuous rotation examples on 3D-FRONT.

1129 B.6 MUTUAL OCCLUSION

1106

1107

1127 1128

1130

1131 In multi-object compositions, mutual occlusion between objects is very common. Although we did not specifically design the method to make the model aware of mutual occlusion, the model has 1132 learned some understanding of these occlusion relationships. A series of research efforts (Van Hoorick 1133 et al., 2023; Ozguroglu et al., 2024; Xu et al., 2024; Zhan et al., 2024; Zhu et al., 2017; Zhan et al.,



Figure A6: **Visualized cross-view matching results**. Since we do not have ground truth image for 3D-FRONT and SUNRGB-D, we only visualize cross-view matching results using our predicted images. But we can still observe a strong cross-view consistency from the accurate matching results.

Table A3:	Ablation	study	on	various	datasets.
-----------	----------	-------	----	---------	-----------

1159			Nove	el View Synt	hesis	Placement	Cross-view Consistency	
1160	Dataset	Method	PSNR(↑)	SSIM(†)	$LPIPS(\downarrow)$	IoU(↑)	Hit Rate(↑)	Dist(↓)
1162		w/o depth	9.829	0.705	0.365	25.7	5.5	75.3
1163	Room-Texture	w/o mask	9.576	0.699	0.384	24.2	2.7	92.2
1164		w/o sch.	9.173	0.689	0.392	22.4	2.3	88.6
1165		Ours	10.014	0.718	0.366	24.2	6.1	78.1
1166		w/o depth	17.457	0.835	0.178	50.5	23.0	52.6
1167	Objaverse	w/o mask	17.176	0.834	0.187	47.3	11.1	57.1
1168		w/o sch.	16.642	0.825	0.210	43.2	6.3	55.0
1169		Ours	17.749	0.840	0.169	51.3	50.0	47.2

1154

1155

1156 1157 1158

2020) specifically focus on addressing mutual occlusion relationships by predicting the amodal masks or synthesizing amodal appearance, but these models typically do not consider scenarios involving camera view change. Moreover, there may not be a well-established metric to measure how well the model understands mutual occlusion from novel viewpoints. We provide a simple experiment and discussion in this section to illustrate model's comprehension of mutual occlusion.

First, in the context of novel view synthesis, the comprehension of occlusion relationships can be divided into two parts. The first is the ability to synthesize parts that were occluded in the input view. The second is the ability to synthesize new occlusion relationships under the novel view. We show several examples of synthesizing occluded parts and synthesizing new occlusions in Fig. A7. We believe this capability is learned in a data-driven way since the multi-object composites are physically plausible regarding these occlusion relationships.

Secondly, we now propose a new metric to evaluate the capability of understanding mutual occlusion under this setting. We first use visible mask and amodal mask in the input-view image to determine how heavily an object is occluded:

1. If an object's visible mask is exactly its full mask, there exists no occlusion.

1187 2. If an object's visible mask is more than 70% of its full mask, the object is occluded.

3. If an object's visible mask is less than 70% of its full mask, the object is heavily occluded.



Figure A7: Occlusion Synthesis Capability. Our proposed method can synthesize new occlusion relationship under novel views as shown in the highlighted area of sofa or cabinet in (a). Our method can also hallucinate occluded parts as shown in the highlighted area of chairs in (b).



Figure A8: Object Removal Example. We can remove an object under novel views by setting a threshold to the predicted mask image and delete corresponding pixels.

Afterward, we segment the predicted view image with ground truth per-object visible mask. We calculate the specific region's PSNR, SSIM, and LPIPS metrics as shown in Tab. A4. It can reflect how well our model and baseline models are at synthesizing novel view plausible images that are originally occluded under the input view. There are 10903 fully visible objects, 6058 occluded objects, and 2215 heavily occluded objects. This experiment is conducted on our own C3DF.

Table At. Evaluation on objects with varying catches of occlusion.	Table A4: F	Evaluation on	objects ^v	with varving	extents of	occlusion.
--------------------------------------------------------------------	--------------------	---------------	----------------------	--------------	------------	------------

Mathad	Visible			Occluded			Heavily Occluded		
Method	PSNR(↑)	SSIM(†)	$LPIPS(\downarrow)$	PSNR(↑)	SSIM(†)	$LPIPS(\downarrow)$	PSNR(↑)	SSIM(†)	$LPIPS(\downarrow)$
Ours	11.45	0.56	0.13	11.33	0.55	0.14	10.57	0.55	0.14
Zero-1-to-3	9.46	0.54	0.16	9.33	0.52	0.17	9.00	0.53	0.16
Zero-1-to-3 [†]	9.68	0.55	0.14	9.54	0.52	0.15	9.26	0.53	0.15

FAILURE CASES AND LIMITATIONS С

Failure Cases We showcase two failure cases in Fig. A10. We can observe that delicate structure and texture like colorful pillows on the sofa or slim legs of chairs are hard for our model to learn. Though object placement is approximately accurate, more fine-grained consistency is not quite ideal



Figure A9: **Reconstruction results using DUSt3R.** We rotate our camera around the multi-object composite and use the predicted images along with the input-view image for reconstruction.



Figure A10: **Failure Cases**. It is hard for our model to learn extremely fine-grained consistency on objects with delicate structure and texture.

in these cases. We believe training with a higher resolution and incorporating epipolar constraints
 will mitigate this problem in the future.

Limitations We identify two limitations of our work. Firstly, though we achieve stronger cross-view consistency with the input view image, our model does not guarantee the multi-view consistency between our synthesized images. It is plausible to synthesize any results in areas with ambiguity, leading to potential multi-view inconsistency in our model. We believe incorporating multi-view awareness techniques Shi et al. (2023b); Wang & Shi (2023); Shi et al. (2023a); Kong et al. (2024); Liu et al. (2023d); Yang et al. (2024a) can mitigate this problem. Secondly, we do not model background texture in our framework due to difficulty of realistically mimicking real-world background texture, making it less convenient to directly apply our method to in-the-wild images. We believe training on more realistic data with background in the future can make our model more convenient to use.

D POTENTIAL NEGATIVE IMPACT

The use of diffusion models to generate compositional assets can raise ethical concerns, especially if used to create realistic yet fake environments. This could be exploited for misinformation or deceptive purposes, potentially leading to trust issues and societal harm. Additionally, hallucinations from diffusion generation models can produce misleading or false information within generated images. This is particularly concerning in applications where accuracy and fidelity to the real world are critical.



Figure A11: More visualized results on C3DFS dataset.



Under review as a conference paper at ICLR 2025





Figure A13: More visualized results on Room-Texture dataset.