# DeclutterNeRF: Generative-Free 3D Scene Recovery for Occlusion Removal

Wanzhou Liu[1]  Zhexiao Xiong[1]  Xinyu Li[2]  Nathan Jacobs[1]
[1]Washington University in St. Louis  [2]Georgia Institute of Technology
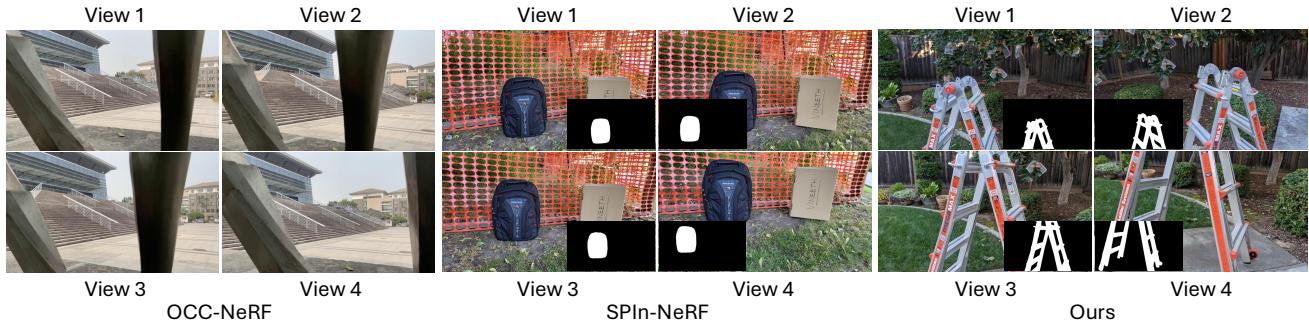{l.wanzhou, x.zhexiao, jacobsn}@wustl.edu, xli3212@gatech.edu

Figure 1. **Comparison of Mainstream Occlusion Removal Datasets.** DeclutterSet is a new dataset reflecting real-world challenges and complexity in occlusion removal. For each dataset, we show four evenly spaced views per scene. As seen in both the RGB images and masks, DeclutterSet exhibits: (i) wider distance distribution, (ii) larger occluded regions, (iii) greater relative motion between viewpoints and occluders, and (iv) more uncertain occluder shapes and mask layouts. In contrast, the OCC-NeRF dataset [58] does not employ masks during selection, limiting it to foreground occlusions and requiring a strict separation between foreground and background, reducing its suitability for complex scenarios. SPIn-NeRF [24] provides limited challenge for cross-view consistency, as it is constrained to small viewpoint variations, keeping occluders and background nearly static across rendered views. A detailed analysis is provided in Sec. 4.1.

## Abstract

Recent novel view synthesis (NVS) techniques, including Neural Radiance Fields (NeRF) and 3D Gaussian Splatting (3DGS) have greatly advanced 3D scene reconstruction with high-quality rendering and realistic detail recovery. Effectively removing occlusions while preserving scene details can further enhance the robustness and applicability of these techniques. However, existing approaches for object and occlusion removal predominantly rely on generative priors, which, despite filling the resulting holes, introduce new artifacts and blurriness. Moreover, existing benchmark datasets for evaluating occlusion removal methods lack realistic complexity and viewpoint variations. To address these issues, we introduce **DeclutterSet**, a novel dataset featuring diverse scenes with pronounced occlusions distributed across foreground, midground, and background, exhibiting substantial relative motion across viewpoints. We further introduce **DeclutterNeRF**, an occlusion removal method free from generative priors. DeclutterNeRF introduces joint multi-view optimization of learnable camera parameters, occlusion annealing regularization, and employs an explainable stochastic structural similarity loss, ensuring high-quality, artifact-free reconstruc- tions from incomplete images. Experiments demonstrate that DeclutterNeRF significantly outperforms state-of-the- art methods on our proposed DeclutterSet, establishing a strong baseline for future research. The code and data are available at *DeclutterNeRF*.

## 1. Introduction

Recent novel view synthesis (NVS) techniques including Neural Radiance Fields (NeRF) [23] and 3D Gaussian Splatting (3DGS) [13] have advanced realistic and efficient 3D scene reconstruction. Removing unwanted objects from rendered scenes would further enhance the flexibility and applicability of these methods for applications in AR, VR, robotics, and autonomous driving [27, 40, 56, 57]. Notably, these real-world scenarios often involve far more complex scene settings than current mainstream occlusion and object removal benchmarks and demand reliable rendering results. This remains a major challenge in 3D reconstruction and calls for a rethinking of existing approaches.

Traditional methods rely on stereo geometry for occlu- sion handling [6, 8, 10, 43, 59]. With the advent of neu- ral view synthesis, filtering-based and optimization-driven

1

techniques have emerged for occlusion selection and removal [32, 33, 58], but their effectiveness remains limited by overly simplified scene assumptions. Recent NeRF and 3DGS approaches have embraced generative models [3–5, 11, 15, 18, 19, 24, 25, 31, 34, 37, 38, 42, 46, 49–51, 55], which can marginally improve reconstruction quality but often introduce significant computational overhead, limiting their practicality. Importantly, most existing methods are developed on OCC-NeRF [58] and SPIn-NeRF [24] datasets, both of which introduce limiting assumptions. As shown in Fig. 1, OCC-NeRF considers only foreground occlusions, while SPIn-NeRF assumes all objects lie on a background plane with minimal relative motion across viewpoints. When these assumptions are violated, *i.e.*, when objects are at different distances or exhibit large motion relative to viewpoints, both generative and non-generative methods struggle, leading to severe artifacts, inconsistent geometry, and unrealistic texture.

To address these limitations, we introduce DeclutterSet, a novel dataset designed to reflect real-world occlusion complexities. Unlike the settings in the previous datasets, DeclutterSet carefully considers the spatial distribution of objects at varying distances, ensuring that occlusions exhibit substantial motion relative to obvious viewpoint changes. By incorporating diverse scenarios where foreground, midground, and background objects shift across views, DeclutterSet provides a more realistic benchmark for evaluating occlusion removal methods.

Building on our DeclutterSet benchmark, we propose DeclutterNeRF, a straightforward optimization-driven approach that leverages NeRF's inherent cross-view consistency to tackle recovery after occlusion removal. Rather than relying on generative models, we demonstrate that targeted improvements to the classic NeRF framework can achieve superior results for this task. Using SAM [15] for initial occlusion segmentation, our approach focuses on optimizing reconstruction from visible regions with minimal computational overhead. We first observe that occlusion presence alters camera parameter estimation, leading to suboptimal pose reconstruction. Inspired by camera posture estimation methods in 3D reconstruction [9, 16, 48], we incorporate camera parameter optimization as a learnable component, allowing multi-view joint optimization to correct pose shifts and mitigate local minima issues. To ensure stable learning after occlusion removal, where only limited pixels are available for rendering, we propose Occlusion Annealing Regularization (OAR), which reduces the impact of occluded regions, improving training stability and preventing overfitting. Finally, we employ Stochastic Structural Similarity Loss (S3IM) [52] to address the long-tail distribution of background pixels caused by non-fixed occlusion regions, which leads to imbalanced ray sampling. Our experiments demonstrate that these targeted

optimizations enable DeclutterNeRF to significantly outperform both previous optimization-based and generative methods in occlusion removal and recovery tasks, while maintaining computational efficiency. We summarize our contributions as follows:

- We introduce DeclutterSet, a novel occlusion removal dataset with diverse real-world occlusion scenarios, capturing multi-position spatial distributions and viewpoint-dependent changes.
- We propose DeclutterNeRF, a generative-free occlusion removal framework that reconstructs 3D scenes using NeRF's implicit multi-view consistency, ensuring reliable and high-quality results without additional training costs.
- We highlight the impact of occlusion removal on camera pose estimation, incorporate multi-view joint learnable camera parameter optimization, and propose Occlusion Annealing Regularization (OAR) to improve robust rendering progress and stabilize training after occlusion removal, mitigating local minima and overfitting issues.
- We theoretically and experimentally validate the "Unreasonable Effectiveness" of random structural similarity [52], showing its broader applicability in our task.

## 2. Related Work

**Occlusion and Object removal.** Traditional approaches to occlusion removal and object deletion often rely on stereo geometry and multi-view consistency cues, such as disparity maps [12], dense flow fields [53], and synthetic apertures [41]. Later, deep learning techniques leveraging temporal information [7] and optical flow [20] emerged. These methods often have limitations due to restricted camera movement and poor generalization to novel viewpoints.

Recent progress in novel view synthesis is driven by NeRF and 3DGS [13, 23], which offer high-fidelity reconstruction via ray tracing and real-time performance via point-based rendering. Both have been extended to generative-free and optimization-based occlusion removal [32, 33, 58] with simplified assumptions. Concretely, OCC-NeRF [58] removes close-range occluders based on bidirectional depth inconsistency, but assumes all occlusions lie in the foreground, resulting in missing foreground details indiscriminately. RobustNeRF [32] and SpotlessSplats [33] handle transient occlusions by removing outliers that appear sporadically across views, but are not designed for persistent or structured obstacles. In contrast, our approach flexibly removes occlusions across diverse object categories, distributions and varying depth ranges.

**2D & 3D Inpainting.** Early image inpainting techniques restore missing regions via local texture synthesis or structural propagation, using exemplar-based [1] or PDE-driven [2] approaches. With deep learning, 2D inpainting evolved to adopt RGB priors (*e.g.*, LaMa [38]) and explicit depth
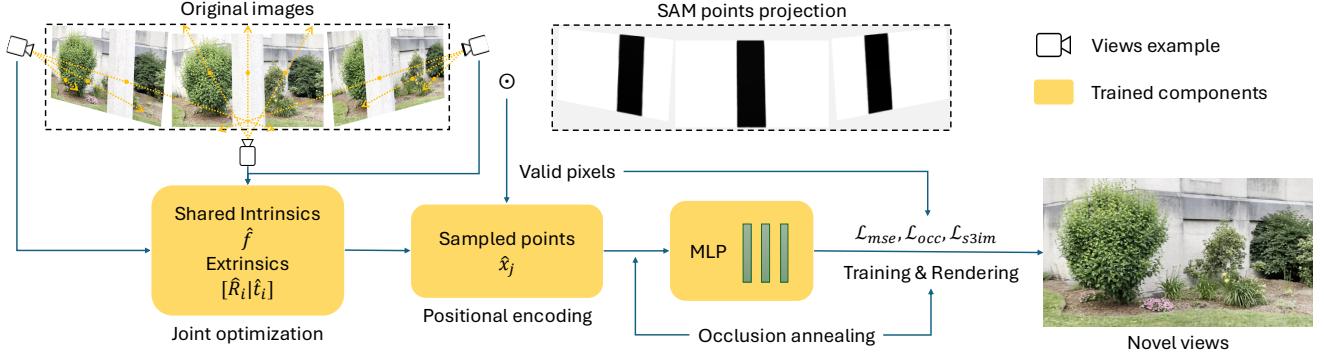
Figure 2. **Overview of Our Optimization Framework.** Our method builds on the NeRF architecture to recover occluded scenes without generative priors. Starting with a single-view SAM segmentation method [55], we propagate occluder masks across views via stereo matching. Camera parameters are jointly optimized with masked photometric supervision to correct occlusion-induced pose errors (Sec. 3.2). To stabilize training and mitigate overfitting to visible regions, we propose Occlusion Annealing Regularization (Sec. 3.3). The Stochastic Structural Similarity loss (Sec. 3.4) enforces global consideration across views and improves reconstruction under long-tail visibility.

cues [19, 24, 42, 50, 55], achieving recovery reconstruction in single and multi-view settings.

Beyond the 2D inpainting methods, diffusion-based techniques have been integrated into 3D reconstruction process with NeRF and 3DGS [4, 5, 11, 18, 21, 34, 37, 46, 49, 51]. MVIP-NeRF [4] leverages diffusion and cross-view distillation to hallucinate missing content, but at the cost of high memory and training time. GScream [46] incorporates depth supervision, which makes it sensitive to the quality of depth estimation. While these generative approaches aim to improve visual fidelity, they are also prone to introducing artifacts, suffer from geometric inconsistencies, and incur significant computational overhead, which limits their scalability in real-world applications.

**3D Reconstruction from Limited Pixels.** Traditional approaches for 3D reconstruction under incomplete observations rely on stereo correspondence [10, 59], image-based priors [6, 43], or local texture synthesis [8]. Segment-based stereo matching improves robustness at object boundaries, while image quilting demonstrates the feasibility of patch-based texture propagation. Despite their contributions, these methods typically involve handcrafted priors and computationally intensive optimization, limiting efficiency and scalability in complex scenes.

The availability of powerful segmentation tools such as SAM [15] and SAM2 [31] has popularized object-level masking in novel view synthesis [11, 25, 46, 55]. This trend amplifies the need for 3D reconstruction from incomplete images, especially when large scene portions are masked out. In this work, we focus on recovering occluded geometry directly from the visible regions, without relying on synthetic content. By leveraging NeRF's cross-view consistency and introducing optimization methods to occlusion scenarios, our method ensures structurally coherent and robust reconstruction under different occlusion scenarios.
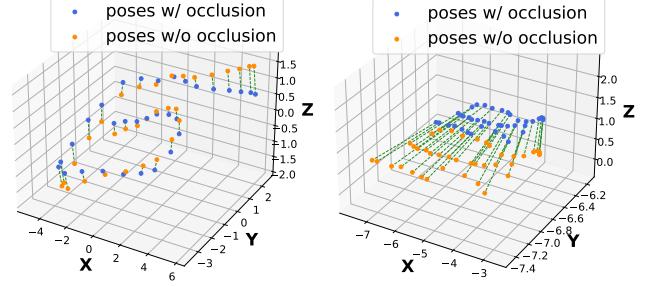


Figure 3. **Visualization of the Impact of Obstacles on Pose Estimation.** Structure-from-motion methods, including the widely used COLMAP [35, 36] and the recently proposed GLOMAP [26], struggle to maintain stable camera pose estimation after occlusion is removed. This is illustrated in the Ladder scene (left) and the Lamp Post scene (right). Green dashed lines connect corresponding samples before and after occlusion removal, highlighting positional shifts. Axes are rotated for clearer visualization.

## 3. Method

### 3.1. Preliminaries

**Neural Radiance Fields.** NeRF [23] is an approach to view synthesis, encoding scenes as implicit continuous volumetric functions. Let $\mathbf{x} = (x, y, z)$ denote a 3D point in space and $\mathbf{d} = (\theta, \phi)$ represent a viewing direction. The core of NeRF is a multi-layer perceptron (MLP) $F_\Theta : (\mathbf{x}, \mathbf{d}) \rightarrow (\mathbf{c}, \sigma)$, where $\Theta$ are the parameters of MLP. It maps a 3D location and viewing direction to a color $\mathbf{c} = (r, g, b)$ and volume density $\sigma$. The camera poses $\mathbf{x}$ are primarily derived from the pose estimation tool COLMAP [35, 36].

**Positional Encoding in NeRF.** Directly optimizing over raw inputs $(\mathbf{x}, \mathbf{d})$ makes it difficult for NeRF to capture high-frequency details. To mitigate this, the mapping $F_\Theta$ is decomposed as $F'_\Theta \circ \lambda$, where $\lambda$ encodes inputs into a higher-dimensional space $\mathbb{R}^{2L}$. The positional encoding

$\lambda(\cdot)$ is defined as:

$$\lambda(p) = \big( \sin(2^0 \pi p), \cos(2^0 \pi p), \ldots,$$
$$\sin(2^{L-1} \pi p) \cos(2^{L-1} \pi p) \big) \quad (1)$$

where $L$ is a hyperparameter that controls the highest encoded frequency. The encoding is applied to each component of the 3D position vector $\mathbf{x}$ (normalized to $[-1, 1]$) and the viewing direction vector $\mathbf{d}$ (unit vector in $[-1, 1]$). In most cases, it has $L = 10$ for $\lambda(\mathbf{x})$ and $L = 4$ for $\lambda(\mathbf{d})$.

**Stochastic Structural Similarity (S3IM).** S3IM[52] is a patch-based, stochastic variant of SSIM [47], designed to introduce global structural supervision into NeRF training. Instead of point-wise loss MSE or local supervise SSIM, S3IM computes SSIM between randomly sampled image patches, capturing non-local and cross-view structural consistency. Given rendered radiance fields $\hat{R}$ and ground-truth images $R$, the loss is computed by sampling $M$ patch pairs $P^{(m)}(\hat{C}), P^{(m)}(C)_{m=1}^{M}$ from both rendered images $\hat{C}$ and ground-truth $C$. Each patch is of size $K \times K$, sampled with stride $s = K$. The final S3IM loss is defined as:

$$\text{S3IM}(\hat{\mathcal{R}}, \mathcal{R}) = \frac{1}{M} \sum_{m=1}^{M} \text{SSIM}(\mathcal{P}^{(m)}(\hat{C}), \mathcal{P}^{(m)}(C)) \quad (2)$$

where $\text{SSIM}(\cdot, \cdot)$ denotes the structural similarity between corresponding patches.

### 3.2. Joint Optimization for Camera Parameters

Our concern about the impact of occlusions on camera parameter estimation originates from classical insights in computer vision and graphics [39, 44]. As demonstrated in Fig. 3, occluders can disturb pose estimation, leading to reconstruction degradation. A straightforward solution would be to recalibrate camera parameters after occlusion removal using the cleaner, occlusion-free 2D observations, which typically enhances reconstruction quality. However, for fair comparison and to test robustness, we retain the original camera parameters estimated under occlusions. Our goal is to leverage the occlusion-free setting as a means to further refine these parameters. To this end, we incorporate the camera parameters into our joint optimization framework. With photometric loss as the major supervision, our framework progressively corrects camera poses, resulting in improved reconstruction performance.

Similar work was proposed by [16, 48], which aimed to completely resolve NeRF's dependence on camera parameters. However, these approaches introduced the problem of easily falling into local minima during training. OCC-NeRF [58], which also employs this method, often produces poor reconstruction quality and can only handle small camera position movements due to this issue. Based on our analysis, they sample only one single image per training iteration, although OCC-NeRF utilized a pretrained ResNet

to extract features for the warped feature map, such high-level feature extraction and projection transformation cannot meet NeRF's requirements for fine-grained geometric details, making it difficult to effectively handle subtle differences in camera poses.

The improvement in our method is intuitive and easy to understand. As illustrated in the initial sampling process in Fig. 2, instead of the traditional approach of sampling from a single view, we jointly optimize across all views by uniformly sampling valid pixels from the entire image set, enabling simultaneous refinement of camera parameters for all viewpoints. This approach is well-founded. Firstly, previous results [48] show that when NeRF parameters are trapped in local minima, focal length parameters often deviate significantly from calibrated values. Since focal length is shared across all input views, distributing focal length sampling across all views contributes to its stable optimization. Second, as illustrated in the sampling process in Fig. 2, where multiple intersecting rays are intentionally drawn as a demonstrative example, the volumetric rendering process handles multiple intersecting rays during multi-view optimization. This leverages the advantage of stereoscopic input, where intersecting rays jointly optimize shared parameters, enhancing stability. This is also consistent with the recent prevailing trends in NeRF training methods. Finally, by adjusting the learning rates and implementing a delayed camera optimization strategy, we avoid potential local minima issues during training.

Let $\Theta$ denote the parameters of our MLP, $\phi$ represent the camera parameters, $\phi_t$ represent the current camera parameters at step $t$, $T$ be the total number of the iteration and $I$ be the set of input images. We formulate our joint optimization objective as:

$$\arg \min_{\Theta, \{\phi_t\}_{t=t_c}^{T}} \sum_{t=1}^{T} \mathbb{E}_{B \sim \mathcal{U}(I,b)} \left[ \frac{1}{|B|} \sum_{(i,j) \in B} \mathcal{L}_{\text{photo}}(R_{\Theta, \phi_t}(r_{i,j}), I_{i,j}) \right]$$
$$\text{where } \phi_t = \begin{cases} \phi_0, & \text{if } t < t_c \\ \text{optimized}, & \text{if } t \geq t_c \end{cases}$$
$$(3)$$

where $t_c$ is our delayed camera parameters optimization start step, $B \sim \mathcal{U}(I, b)$ indicates a batch of $B$ rays are used in total, and has uniform $b$ samples from each image. $\mathcal{L}_{\text{photo}}$ is the photometric loss between the rendered images and limited visible ground truth pixels. In our joint optimization framework, it is primarily supervised by $\mathcal{L}_{\text{MSE}}$. The specific weighting of this and other losses is detailed in Sec. 4.1.

### 3.3. Occlusion Annealing Regularization (OAR)

In reconstruction after occlusion removal, the most significant issue arises from the variability of the visible region due to the non-fixed distribution of the occlusion. This results in two effects: 1) some areas being underfitted because they are not adequately visible across multiple views,

and 2) other non-occluded regions may be overfitted due to being fully visible in every input image. This imbalance leads to underfitting in regions that are sparsely visible and overfitting in consistently visible ones, often causing rendering artifacts within the same view. Our objective is to reduce the impact of overfitting on the final rendering effect while ensuring sufficient generalization across views. Given that neural networks tend to learn low-frequency features [23, 30], and inspired by regularization work in neural rendering across different frequencies [17, 28, 54], we first use lower-dimensional pose embeddings to learn consistent low-frequency scene features from various perspectives, then gradually increase the dimension to standard frequency encoding. This simple approach leads to blurred boundaries where occlusions exist. Additionally, it delays the training convergence in areas with higher visibility frequencies, contributing to the generation of consistent rendering effects.

Due to limited multi-view supervision, artifacts in occluded regions frequently manifest as floaters near the camera, where rendering is more sensitive to density misestimation. While prior work [54] penalizes near-camera rays to suppress floaters, we find that early-stage low-frequency training can cause unstable feature clustering in these regions, making global penalties detrimental and prone to collapse. To mitigate this, we propose Occlusion Annealing Regularization (OAR), which gradually introduces occlusion loss during frequency ramp-up, stabilizing training under view redundancy.

The position and direction encodings at iteration $t$ are represented as:

$$
\begin{aligned}
\mathbf{e}_{\mathrm{pos}}(t) &= \mathbf{x} \odot \mathbf{m}_{\mathrm{pos}}(t), \\
\mathbf{e}_{\mathrm{dir}}(t) &= \mathbf{d} \odot \mathbf{m}_{\mathrm{dir}}(t),
\end{aligned}
\tag{4}
$$

where $\mathbf{x}$ and $\mathbf{d}$ are the original position and direction encodings, and $\mathbf{m}_{\mathrm{pos}}(t)$ and $\mathbf{m}_{\mathrm{dir}}(t)$ are frequency masks that depend on the current iteration $t$.

The masks are defined as:

$$
\mathbf{m}_{\mathrm{pos,dir}}(t) = \begin{cases} 1, & \text{if } f \le f_{\max}(t), \\ 0, & \text{otherwise}, \end{cases}
\tag{5}
$$

where $f$ is the frequency of each encoding dimension, and $f_{\max}(t)$ is the maximum allowed frequency at iteration $t$, which increases linearly from 0 to the maximum frequency over the course of training. Through the masks, low-frequency and high-frequency information is progressively exposed to the network.

This progressive exposure is synchronized with an annealed occlusion loss weight, defined via a cosine schedule between iterations $t_{\mathrm{start}}$ and $t_{\mathrm{end}}$, ensuring a smooth transi-
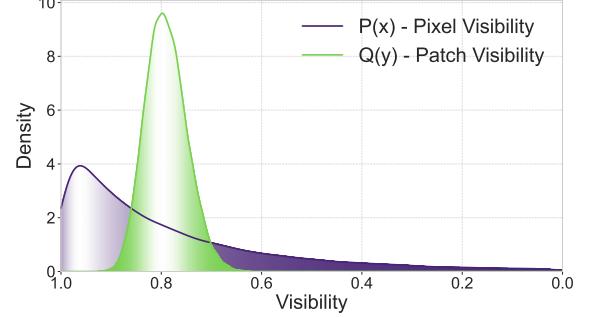


Figure 4. **Visualization of Sampling Distribution.** For a demonstration of the principle of global patched S3IM, the distribution of pixels exhibits a marked imbalance. This issue can be addressed through our patch reorganization. The distribution of each patch becomes more concentrated and uniform, eliminating the regional long-tail distribution of pixels and promoting stable model iteration. Darker regions indicate more extreme long-tailed visibility, which require targeted optimization.

tion toward full supervision as frequency increases:

$$
w_{\mathrm{occ}}(t) = \begin{cases} 0, & \text{if } t < t_{\mathrm{start}}, \\ \frac{w_{\mathrm{full}}}{2}\left(1 + \cos\left(\pi \frac{t_{\mathrm{end}} - t}{t_{\mathrm{end}} - t_{\mathrm{start}}}\right)\right), & \text{if } t_{\mathrm{start}} < t < t_{\mathrm{end}}, \\ w_{\mathrm{full}}, & \text{if } t \ge t_{\mathrm{end}}. \end{cases}
\tag{6}
$$

Here, $w_{\mathrm{full}}$ is the full occlusion loss weight, $t_{\mathrm{start}}$ is the iteration to start introducing the occlusion loss, and corresponding $t_{\mathrm{end}}$ is the iteration when the full weight is reached.

The occlusion loss is then calculated as:

$$
\mathcal{L}_{\mathrm{occ}}(t) = w_{\mathrm{occ}}(t) \cdot \mathcal{L}_{\mathrm{occ\_base}}
\tag{7}
$$

$$
\mathcal{L}_{occ\_base} = \frac{\boldsymbol{\sigma}_K^{\mathbf{T}} \cdot \mathbf{m}_K}{K} = \frac{1}{K}\sum_K \sigma_k \cdot m_k,
\tag{8}
$$

where $\mathbf{m}_k$ is a binary mask vector and $\boldsymbol{\sigma}_K$ denotes the density values of the $K$ sampled points along the ray. The frequency regularization end ($t_{\mathrm{freq\_end}}$) and occlusion annealing are connected through $\lambda$:

$$
t_{\mathrm{end}} = \frac{t_{\mathrm{freq\_end}}}{\lambda}
\tag{9}
$$

This coordination between frequency and occlusion schedules promotes stable learning early on while effectively penalizing artifacts later in training.

### 3.4. S3IM in Occluded Long-Tailed Visibility

Occluded scenes often exhibit a long-tailed distribution of pixel visibility, *i.e.*, most pixels appear frequently across views, while a minority are rarely observed. This imbalance hampers stable training and leads to biased reconstructions, as frequently visible pixels dominate the optimization signal. We employ a patch-based global stochastic structural

5

similarity method [52] to address this issue and validate its effectiveness in our task. As described in Sec. 3.1 with all the notations, S3IM lies within $[-1, 1]$ and is positively correlated with image quality, so its loss definition is:

$$\mathcal{L}_{S3IM}(\Theta, \mathcal{R}) = 1 - S3IM(\hat{\mathcal{R}}, \mathcal{R})$$
$$= 1 - \frac{1}{M} \sum_{m=1}^{M} SSIM(\mathcal{P}^{(m)}(\hat{C}), \mathcal{P}^{(m)}(C)). \quad (10)$$

This involves a patch-based stochastic structural similarity SSIM [47] but in a global range. To better understand how this loss formulation mitigates long-tailed visibility issues, we begin by characterizing the underlying pixel visibility distribution. In the illustration for occluded scenes, pixel visibility follows a long-tailed pattern:

$$P(x) \approx \frac{1}{(x_{max} - x + 1)^{\alpha}}, \quad \alpha > 1 \quad (11)$$

where $x$ represents pixel visibility, the number of views in which a pixel is visible, and $x_{max}$ corresponds to the maximum visibility. S3IM randomly samples rays across all views and group them into $K \times K$ patches. Each patch's visibility is defined as:

$$\text{vis}(p) = \frac{1}{K^2} \sum_{i=1}^{K^2} \text{vis}(pixel_i) \quad (12)$$

which converts the per-pixel visibility distribution $P(x)$ into a patch-level distribution $Q(y)$, where

$$Q(y) = P(y = \frac{1}{K^2} \sum_{i=1}^{K^2} X_i), \quad X_i \sim P(x) \quad (13)$$

Compared with $x$ and $x_{max}$, $y$ is the average visibility of a patch. By aggregating over both high and low visibility pixels, each patch visibility becomes naturally moderated:

$$\min_{i \in p} \text{vis}(pixel_i) \le \text{vis}(p) \le \max_{i \in p} \text{vis}(pixel_i) \quad (14)$$

This mixing effect shortens the tail of the visibility distribution (as visualized in Fig. 4), resulting in more centralized gradients during optimization. By balancing supervision across visibility levels, it enhances stability and improves reconstruction in sparsely observed regions.

## 4. Experiments

### 4.1. Experimental Setup

**Dataset.** Due to the novelty of the occlusion removal and reconstruction problem and the limited availability of existing datasets, we follow the pattern of NeRF, which uses 8 scenes from the LLFF dataset [22], and create a dataset
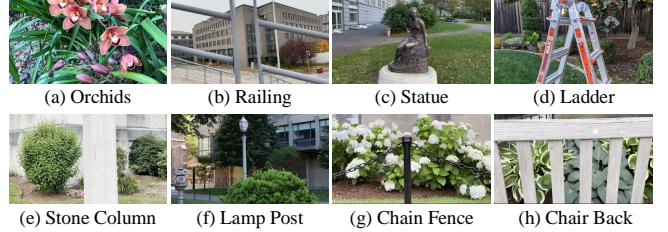


(a) Orchids  (b) Railing  (c) Statue  (d) Ladder

(e) Stone Column  (f) Lamp Post  (g) Chain Fence  (h) Chair Back

Figure 5. **The DeclutterSet.** (a)Orchids and (f)Lamp Post illustrate occluders at different distances: in (a), both the buds and flowers lie on the same near-depth plane close to the camera, while in (f), the occluding object is situated farther away in the mid-background; (b)Railing and (c)Statue resemble traditional occlusion and object removal settings commonly found in existing benchmarks; (e)Stone Column and (g)Chain Fence exhibit occlusions that scatter across different image regions as the viewpoint shifts; (d)Ladder and (h)Chair Back feature larger, irregularly shaped occluders and more pronounced viewpoint variations, posing further challenges to cross-view consistency and geometry recovery. Further details are provided in the supplementary material.

comprising 8 occluded scenes. As shown in Fig. 5, DeclutterSet comprises eight occluded scenes, including four sourced from existing benchmarks and four newly captured, ensuring a balanced distribution of occlusion types and scene layouts. Specifically, (a) Orchids is taken from the classic LLFF dataset, (b) Railing is from the OCC-NeRF dataset, and (c) Statue and (d) Ladder are from dataset IBRNet [45] which currently has become the mainstream data in object removal. For the data we constructed, (e) to (h), each scene consists of approximately 30 images captured using a Canon R6 Mark II or an iPhone 12 Pro. Following the mainstream approach, we created the test set by holding out $1/8$ of the images. The details for mask annotation and propagation are described in the supplementary material.

**Baselines & Metrics.** We compare our method against both generative and non-generative state-of-the-art approaches for object and occlusion removal in NVS. Specifically, OCC-NeRF [58] serves as the generative-free baseline, while SPIn-NeRF [24] and MVIP-NeRF [4] represent the generative baselines. We provide both qualitative and quantitative evaluations for the rendering results. For qualitative analysis, we include visualizations across all scenes in the main paper and the supplementary material. For quantitative evaluation, we report standard NeRF reconstruction metrics, PSNR, SSIM, and LPIPS, computed over non-occluded pixels only.

**Parameters.** Similar to most learning-based 3D reconstruction methods, DeclutterNeRF is also influenced by hyperparameters. We set the termination of frequency regularization at 10% of the total iterations, begin camera parameter optimization at 20% of the total iterations, and set the Occlusion Annealing Regularization $\lambda$ to 100. Our $\mathcal{L}_{photo}$ consists of three mainstream loss functions: $\mathcal{L}_{mse}$, $\mathcal{L}_{occ}$,

Table 1. **Quantitative Comparisons With the Generative-Free Baseline.** Due to the lack of evaluation code or the inability to reconstruct all scenes, generative baselines are excluded. OCC-NeRF is selected as the only non-generative method that supports reconstruction on DeclutterSet. Under its original parameter settings, OCC-NeRF underperforms due to its rigid near-range removal strategy and lack of post-removal refinement. In contrast, DeclutterNeRF yields superior performance across all metrics and scenes, demonstrating improved robustness to complex scenarios. On average, our method improves PSNR by 68.4%, SSIM by 238.0%, and reduces LPIPS by 54.8%.

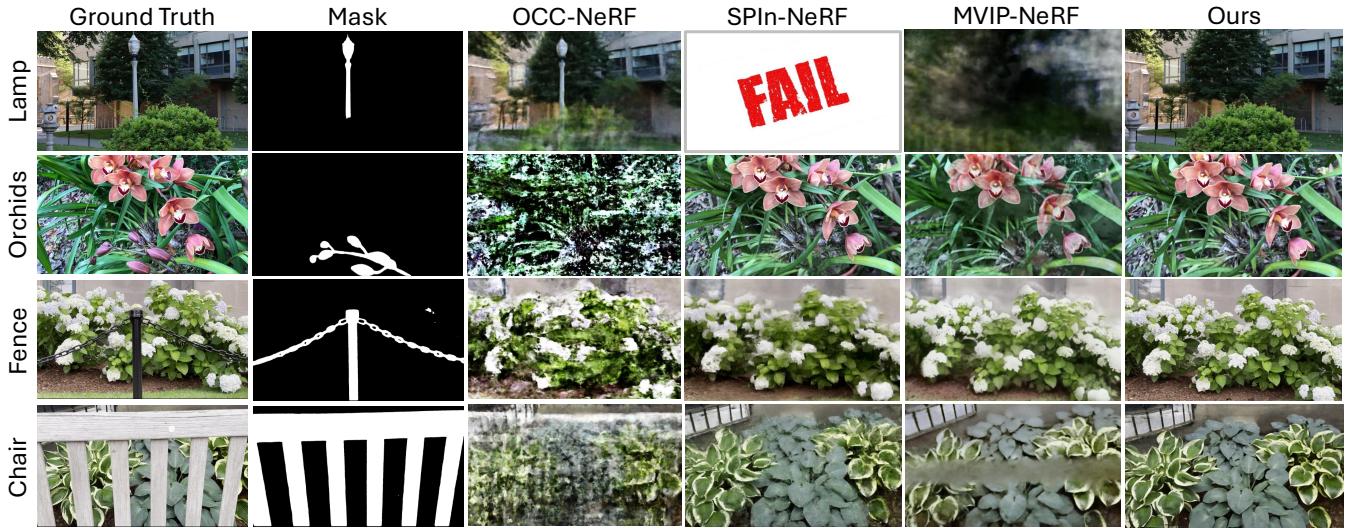| | PSNR↑ | | SSIM↑ | | LPIPS↓ | |
|---|---|---|---|---|---|---|
| Scene | OCC-NeRF | Ours | OCC-NeRF | Ours | OCC-NeRF | Ours |
| (a) Orchids | 10.76 | 20.86 | 0.213 | 0.894 | 0.371 | 0.130 |
| (b) Railing | 14.00 | 23.12 | 0.324 | 0.860 | 0.457 | 0.241 |
| (c) Statue | 16.13 | 24.87 | 0.197 | 0.902 | 0.502 | 0.135 |
| (d) Ladder | 13.17 | 21.37 | 0.074 | 0.656 | 0.534 | 0.352 |
| (e) Stone Column | 14.73 | 21.73 | 0.240 | 0.864 | 0.435 | 0.229 |
| (f) Lamp Post | 15.62 | 22.67 | 0.581 | 0.903 | 0.403 | 0.240 |
| (g) Chain Fence | 11.90 | 23.71 | 0.294 | 0.927 | 0.389 | 0.125 |
| (h) Chair Back | 10.59 | 21.67 | 0.118 | 0.887 | 0.432 | 0.143 |
| Average | 13.36 | 22.50 | 0.255 | 0.862 | 0.440 | 0.199 |



Figure 6. **Qualitative Comparisons With Baselines.** We compare DeclutterNeRF against OCC-NeRF, SPIn-NeRF and MVIP-NeRF across both standard and newly collected datasets. Our method reliably removes occluders at varying depths and achieves photorealistic reconstruction. In contrast, OCC-NeRF struggles in close-range scenes due to its distant-only rendering assumptions. SPIn-NeRF and MVIP-NeRF, designed on previous benchmarks, frequently suffer from inconsistent floaters and hallucinated artifacts when occluders shift their relative positions across views—a mode exposed by our DeclutterSet but overlooked in prior benchmarks.

and $\mathcal{L}_{s3im}$. The weights of $L_{occ}$ and $L_{s3im}$ are set to 0.01. More details can be found in the supplementary material.

## 4.2. Comparison Results

**Qualitative Evaluation.** Figure 6 compares our method's rendering results with OCC-NeRF, SPIn-NeRF, and MVIP-NeRF. The mask represents the occluding objects we aim to remove. Our pipeline enables selective removal of occlusions with minimal manual intervention, including the distant lamp and the unopened orchid buds in the foreground, while achieving photorealistic reconstructions. In contrast, OCC-NeRF only handles distant scenes adequately, often

removing desired nearby objects and failing to reconstruct close-range details. For datasets like LLFF, which primarily consist of close-range scenes, OCC-NeRF's performance is significantly limited. Even for distant scenes, OCC-NeRF's depth warping strategy, impedes the optimization process, causing the model to struggle with complex geometries and leading to poor reconstruction quality that often appears smeared. Our joint camera parameter optimization strategy effectively avoids local minima traps and leverages this optimization to achieve high quality reconstructions.

The performance of SPIn-NeRF and MVIP-NeRF on our DeclutterSet also shows notable differences compared to

our model. SPIn-NeRF's heavy reliance on COLMAP and pre-rendered depth priors impacts its reconstruction capability, often leading to failure rendering when depth information cannot be accurately recovered, particularly in distant scenes. More failure cases, parameter settings, and detailed analysis can be found in the supplementary material. Even MVIP-NeRF, despite claiming independence from depth priors, struggles with outdoor distant scenes. Moreover, in relatively simple scenes, these generative methods are highly prone to overfitting. As training progresses, reconstruction quality plateaus while artifacts increase, degrading overall results. Our method effectively avoids overfitting and underfitting issues caused by varying exposure levels in occluded regions, and it suppresses artifact generation, achieving optimal reconstruction results.

**Quantitative Evaluation.** We quantitatively evaluate our method against OCC-NeRF, the only generative-free baseline capable of handling all scenes in DeclutterSet. As shown in Table 1, DeclutterNeRF achieves consistent and significant improvements across all standard NeRF metrics. The performance gap is especially pronounced in challenging scenes such as Orchids, Chain Fence, and Chair Back, where OCC-NeRF's fixed near-range removal strategy struggles to adapt to varying occlusion depths and scene complexity. This is particularly evident in the Orchids scene, as corroborated by Fig. 6, where none of the flowers are retained in OCC-NeRF's output—resulting in one of the lowest PSNR scores. In contrast, our method contributes to more robust, artifact-free reconstruction under diverse occlusion settings.

Beyond accuracy, DeclutterNeRF also demonstrates high practical efficiency, enabled by our multi-stage architectural improvements. It completes training in under 10 hours on a single NVIDIA RTX 4090 GPU, with memory consumption kept below 10 GB. In comparison, OCC-NeRF requires over 30 hours of training, while the diffusion and distillation-based learning MVIP-NeRF demands more than 100 GB of GPU memory. This level of efficiency makes DeclutterNeRF more accessible and better suited for broader adoption and large-scale experimentation.

### 4.3. Ablation Studies

We demonstrate the impact of ablation and the gradual introduction of each component in Table 2. Initially, we train (i) a masked NeRF, which simply uses the non-occluded mask areas for training. Subsequently, we introduce (ii) joint optimization for camera parameters, which improves all the rendering metrics. Notably, while the use of (iii) occlusion annealing regularization (OAR) results in a slight regression in rendering metrics, it addresses the main issues of artifacts and incomplete rendering. We show this in Fig. 7, which observably enhances the actual visual effect. Finally, we introduced the S3IM loss to further address the

Table 2. **Ablation Studies.** Metrics evaluation through ablation and gradual introduction of each module in our framework. Although the introduction of OAR leads to a slight drop in overall quantitative accuracy, it plays a critical role in cross-view generalization and is essential for successful reconstruction after occlusion removal. This effect is further illustrated in Fig 7.

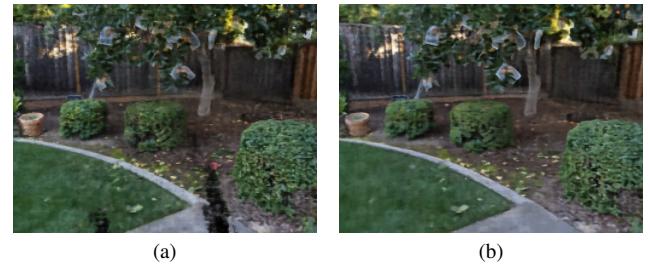| Method | PSNR↑ | SSIM↑ | LPIPS↓ |
|---|---|---|---|
| (i) (masked NeRF) | 19.59 | 0.56 | 0.38 |
| (ii) (+camera opt.) | 21.01 | **0.67** | **0.28** |
| (iii) (+OAR) | 20.89 | 0.61 | 0.29 |
| (iv) (+s3im loss) | **21.37** | 0.656 | 0.352 |



| (a) | (b) |

Figure 7. **Visual Ablation Studies for Occlusion Annealing Regularization (OAR).** Visual effects between introducing (iii) OAR on the Ladder scene: (a) before, (b) after.

reconstruction costs caused by occlusions.

### 4.4. Limitations

Our experiments assume that occluded regions are at least partially visible from other viewpoints, as we have no generative priors to reconstruct unseen parts of the scene. In cases where occlusions completely hide target content from all views, generative priors remain necessary. However, we believe that future generative approaches can build upon our framework—first maximizing reconstruction from observable data, then refining the remaining gaps through targeted generation. This layered strategy promises both efficiency and consistency for occlusion-aware scene recovery.

### 5. Conclusion

We introduced DeclutterSet, a dataset designed to reflect the real-world complexity of occlusions with diverse object layouts and viewpoint variations, addressing critical limitations of existing benchmarks. Based on this, we proposed DeclutterNeRF, a generative-free framework that leverages NeRF's multi-view consistency, joint camera optimization, occlusion annealing regularization, and stochastic structural similarity loss. Our method achieves state-of-the-art performance on this specific task with minimal computational overhead. We hope this work offers a broader perspective on occlusion and object removal and serves as a foundation for future research, whether generative or optimization-based, in robust and efficient 3D scene reconstruction.

# References

[1] Coloma Ballester, Marcelo Bertalmio, Vicent Caselles, Guillermo Sapiro, and Joan Verdera. Filling-in by joint interpolation of vector fields and gray levels. *IEEE transactions on image processing*, 10(8):1200–1211, 2001. 2

[2] Marcelo Bertalmio, Guillermo Sapiro, Vincent Caselles, and Coloma Ballester. Image inpainting. In *Proceedings of the 27th Annual Conference on Computer Graphics and Interactive Techniques*, page 417–424, USA, 2000. ACM Press/Addison-Wesley Publishing Co. 2

[3] Chenjie Cao, Chaohui Yu, Yanwei Fu, Fan Wang, and Xiangyang Xue. Mvinpainter: Learning multi-view consistent inpainting to bridge 2d and 3d editing. *ArXiv*, abs/2408.08000, 2024. 2

[4] Honghua Chen, Chen Change Loy, and Xingang Pan. Mvipnerf: Multi-view 3d inpainting on nerf scenes via diffusion prior. In *CVPR*, 2024. 3, 6

[5] Jiafu Chen, Tianyi Chu, Jiakai Sun, Wei Xing, and Lei Zhao. Single-mask inpainting for voxel-based neural radiance fields. In *European Conference on Computer Vision*, 2024. 2, 3

[6] Ming-Ming Cheng, Fang-Lue Zhang, Niloy J Mitra, Xiaolei Huang, and Shi-Min Hu. Repfinder: finding approximately repeated scene elements for image editing. *ACM transactions on graphics (TOG)*, 29(4):1–8, 2010. 1, 3

[7] Chen Du, Byeongkeun Kang, Zheng Xu, Ji Dai, and Truong Nguyen. Accurate and efficient video de-fencing using convolutional neural networks and temporal information. In *2018 IEEE International Conference on Multimedia and Expo (ICME)*, pages 1–6. IEEE, 2018. 2

[8] Alexei A. Efros and William T. Freeman. Image quilting for texture synthesis and transfer. In *Proceedings of the 28th Annual Conference on Computer Graphics and Interactive Techniques*, page 341–346, New York, NY, USA, 2001. Association for Computing Machinery. 1, 3

[9] Yang Fu, Sifei Liu, Amey Kulkarni, Jan Kautz, Alexei A. Efros, and Xiaolong Wang. Colmap-free 3d gaussian splatting. *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 20796–20805, 2023. 2

[10] Heiko Hirschmuller. Stereo processing by semiglobal matching and mutual information. *IEEE Transactions on pattern analysis and machine intelligence*, 30(2):328–341, 2007. 1, 3

[11] Sheng-Yu Huang, Zi-Ting Chou, and Yu-Chiang Frank Wang. 3d gaussian inpainting with depth-guided cross-view consistency. *ArXiv*, abs/2502.11801, 2025. 2, 3

[12] Sankaraganesh Jonna, Sukla Satapathy, and Rajiv R Sahay. Stereo image de-fencing using smartphones. In *2017 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 1792–1796. IEEE, 2017. 2

[13] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Trans. Graph.*, 42(4):139–1, 2023. 1, 2

[14] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014. 1

[15] 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 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026, 2023. 2, 3

[16] Axel Levy, Mark J. Matthews, Matan Sela, Gordon Wetzstein, and Dmitry Lagun. Melon: Nerf with unposed images using equivalence class estimation. *ArXiv*, abs/2303.08096, 2023. 2, 4

[17] Chen-Hsuan Lin, Wei-Chiu Ma, Antonio Torralba, and Simon Lucey. Barf: Bundle-adjusting neural radiance fields. In *IEEE International Conference on Computer Vision (ICCV)*, 2021. 5

[18] Chieh Hubert Lin, Changil Kim, Jia-Bin Huang, Qinbo Li, Chih-Yao Ma, Johannes Kopf, Ming-Hsuan Yang, and Hung-Yu Tseng. Taming latent diffusion model for neural radiance field inpainting. *ArXiv*, abs/2404.09995, 2024. 2, 3

[19] Hao-Kang Liu, I Shen, Bing-Yu Chen, et al. Nerf-in: Free-form nerf inpainting with rgb-d priors. *arXiv preprint arXiv:2206.04901*, 2022. 2, 3

[20] Yu-Lun Liu, Wei-Sheng Lai, Ming-Hsuan Yang, Yung-Yu Chuang, and Jia-Bin Huang. Learning to see through obstructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14215–14224, 2020. 2

[21] Zhiheng Liu, Ouyang Hao, Qiuyu Wang, Ka Leong Cheng, Jie Xiao, Kai Zhu, Nan Xue, Yu Liu, Yujun Shen, and Yang Cao. Infusion: Inpainting 3d gaussians via learning depth completion from diffusion prior. *ArXiv*, abs/2404.11613, 2024. 3

[22] Ben Mildenhall, Pratul P. Srinivasan, Rodrigo Ortiz-Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi, Ren Ng, and Abhishek Kar. Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Transactions on Graphics (TOG)*, 2019. 6

[23] 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. 1, 2, 3, 5

[24] Ashkan Mirzaei, Tristan Aumentado-Armstrong, Konstantinos G. Derpanis, Jonathan Kelly, Marcus A. Brubaker, Igor Gilitschenski, and Alex Levinshtein. SPIn-NeRF: Multiview segmentation and perceptual inpainting with neural radiance fields. In *CVPR*, 2023. 1, 2, 3, 6

[25] Jingcheng Ni, Weiguang Zhao, Daniel Wang, Ziyao Zeng, Chenyu You, Alex Wong, and Kaizhu Huang. Efficient interactive 3d multi-object removal. *ArXiv*, abs/2501.17636, 2025. 2, 3

[26] Linfei Pan, Daniel Barath, Marc Pollefeys, and Johannes Lutz Schönberger. Global Structure-from-Motion Revisited. In *European Conference on Computer Vision (ECCV)*, 2024. 3

[27] Mingjie Pan, Jiaming Liu, Renrui Zhang, Peixiang Huang, Xiaoqi Li, Hongwei Xie, Bing Wang, Li Liu, and Shanghang Zhang. Renderocc: Vision-centric 3d occupancy pre-

diction with 2d rendering supervision. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 12404–12411. IEEE, 2024. 1

[28] Keunhong Park, Utkarsh Sinha, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Steven M. Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. *ICCV*, 2021. 5

[29] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. *ArXiv*, abs/1912.01703, 2019. 1

[30] Nasim Rahaman, Aristide Baratin, Devansh Arpit, Felix Draxler, Min Lin, Fred Hamprecht, Yoshua Bengio, and Aaron Courville. On the spectral bias of neural networks. In *International conference on machine learning*, pages 5301–5310. PMLR, 2019. 5

[31] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, Eric Mintun, Junting Pan, Kalyan Vasudev Alwala, Nicolas Carion, Chao-Yuan Wu, Ross Girshick, Piotr Dollár, and Christoph Feichtenhofer. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024. 2, 3

[32] Sara Sabour, Suhani Vora, Daniel Duckworth, Ivan Krasin, David J. Fleet, and Andrea Tagliasacchi. Robustnerf: Ignoring distractors with robust losses. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 20626–20636, 2023. 2

[33] Sara Sabour, Lily Goli, George Kopanas, Mark J. Matthews, Dmitry Lagun, Leonidas J. Guibas, Alec Jacobson, David J. Fleet, and Andrea Tagliasacchi. Spotlesssplats: Ignoring distractors in 3d gaussian splatting. *ArXiv*, abs/2406.20055, 2024. 2

[34] Ahmad Salimi, Tristan Aumentado-Armstrong, Marcus A. Brubaker, and Konstantinos G. Derpanis. Geometry-aware diffusion models for multiview scene inpainting. *ArXiv*, abs/2502.13335, 2025. 2, 3

[35] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 3

[36] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise view selection for unstructured multi-view stereo. In *European Conference on Computer Vision (ECCV)*, 2016. 3

[37] Zhihao Shi, Dong Huo, Yuhongze Zhou, Kejia Yin, Yan Min, Juwei Lu, and Xinxin Zuo. Imfine: 3d inpainting via geometry-guided multi-view refinement. 2025. 2, 3

[38] Roman Suvorov, Elizaveta Logacheva, Anton Mashikhin, Anastasia Remizova, Arsenii Ashukha, Aleksei Silvestrov, Naejin Kong, Harshith Goka, Kiwoong Park, and Victor Lempitsky. Resolution-robust large mask inpainting with fourier convolutions. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 2149–2159, 2022. 2

[39] Richard Szeliski and Philip HS Torr. Geometrically constrained structure from motion: Points on planes. In *European Workshop on 3D Structure from Multiple Images of Large-Scale Environments*, pages 171–186. Springer, 1998. 4

[40] Adam Tonderski, Carl Lindström, Georg Hess, William Ljungbergh, Lennart Svensson, and Christoffer Petersson. Neurad: Neural rendering for autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14895–14904, 2024. 1

[41] Vaibhav Vaish, Marc Levoy, Richard Szeliski, C Lawrence Zitnick, and Sing Bing Kang. Reconstructing occluded surfaces using synthetic apertures: Stereo, focus and robust measures. In *2006 IEEE computer society conference on computer vision and pattern recognition (CVPR'06)*, pages 2331–2338. IEEE, 2006. 2

[42] Dongqing Wang, Tong Zhang, Alaa Abboud, and Sabine Süsstrunk. Innerf360: Text-guided 3d-consistent object inpainting on 360-degree neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12677–12686, 2024. 2, 3

[43] Jiaping Wang, Shuang Zhao, Xin Tong, John Snyder, and Baining Guo. Modeling anisotropic surface reflectance with example-based microfacet synthesis. In *ACM SIGGRAPH 2008 papers*, pages 1–9. Association for Computing Machinery, 2008. 1, 3

[44] John YA Wang and Edward H Adelson. Representing moving images with layers. *IEEE transactions on image processing*, 3(5):625–638, 1994. 4

[45] 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. 6

[46] Yuxin Wang, Qianyi Wu, Guofeng Zhang, and Dan Xu. Gscream: Learning 3d geometry and feature consistent gaussian splatting for object removal. In *European Conference on Computer Vision*, 2024. 2, 3

[47] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. 4, 6

[48] Zirui Wang, Shangzhe Wu, Weidi Xie, Min Chen, and Victor Adrian Prisacariu. Nerf–: Neural radiance fields without known camera parameters. *arXiv preprint arXiv:2102.07064*, 2021. 2, 4, 1

[49] Ethan Weber, Aleksander Holynski, Varun Jampani, Saurabh Saxena, Noah Snavely, Abhishek Kar, and Angjoo Kanazawa. Nerfiller: Completing scenes via generative 3d inpainting. In *CVPR*, 2024. 2, 3

[50] Silvan Weder, Guillermo Garcia-Hernando, Áron Monszpart, Marc Pollefeys, Gabriel Brostow, Michael Firman, and Sara Vicente. Removing objects from neural radiance fields. In *CVPR*, 2023. 3

[51] Chung-Ho Wu, Yang-Jung Chen, Ying-Huan Chen, Jie-Ying Lee, Bo-Hsu Ke, Chun-Wei Tuan Mu, Yi-Chuan Huang, Chin-Yang Lin, Min-Hung Chen, Yen-Yu Lin, and Yu-

Lun Liu. Aurafusion360: Augmented unseen region alignment for reference-based 360° unbounded scene inpainting. *ArXiv*, abs/2502.05176, 2025. 2, 3

[52] Zeke Xie, Xindi Yang, Yujie Yang, Qi Sun, Yixiang Jiang, Haoran Wang, Yunfeng Cai, and Mingming Sun. S3im: Stochastic structural similarity and its unreasonable effectiveness for neural fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 18024–18034, 2023. 2, 4, 6

[53] Tianfan Xue, Michael Rubinstein, Ce Liu, and William T Freeman. A computational approach for obstruction-free photography. *ACM Transactions on Graphics (TOG)*, 34(4): 1–11, 2015. 2

[54] Jiawei Yang, Marco Pavone, and Yue Wang. Freenerf: Improving few-shot neural rendering with free frequency regularization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8254–8263, 2023. 5

[55] Youtan Yin, Zhoujie Fu, Fan Yang, and Guosheng Lin. Ornerf: Object removing from 3d scenes guided by multiview segmentation with neural radiance fields, 2023. 2, 3, 1

[56] Chubin Zhang, Juncheng Yan, Yi Wei, Jiaxin Li, Li Liu, Yansong Tang, Yueqi Duan, and Jiwen Lu. Occnerf: Self-supervised multi-camera occupancy prediction with neural radiance fields. *arXiv preprint arXiv:2312.09243*, 2023. 1

[57] Xiao Zhao, Bo Chen, Mingyang Sun, Dingkang Yang, Youxing Wang, Xukun Zhang, Mingcheng Li, Dongliang Kou, Xiaoyi Wei, and Lihua Zhang. Hybridocc: Nerf enhanced transformer-based multi-camera 3d occupancy prediction. *IEEE Robotics and Automation Letters*, 2024. 1

[58] Chengxuan Zhu, Renjie Wan, Yunkai Tang, and Boxin Shi. Occlusion-free scene recovery via neural radiance fields. 2023. 1, 2, 4, 6

[59] C Lawrence Zitnick and Sing Bing Kang. Stereo for image-based rendering using image over-segmentation. *International Journal of Computer Vision*, 75(1):49–65, 2007. 1, 3

# DeclutterNeRF: Generative-Free 3D Scene Recovery for Occlusion Removal

## Supplementary Material

## A. Method Details

### A.1. Architecture and Training Details

DeclutterNeRF follows the core architecture and strategy of the original NeRF [23]. Specifically, we build on NeRF-- [48] and apply DeclutterNeRF on top of this structure. Our model is implemented using PyTorch [29] and trained on a single NVIDIA GeForce RTX 4090 GPU. Since our dataset typically requires no more than 10GB of VRAM, and the image size can be adjusted flexibly to control memory usage, GPUs with significantly lower configurations can also be used to train our model. Unlike recent models that employ multiple MLPs and assign distinct names to each, we adhere to the original NeRF approach by using a single MLP for training and rendering. Our MLP consists of 8 fully connected ReLU hidden layers, each with 128 dimensions. Our further camera optimization algorithm mentioned in *Sec. 3.2* and problems encountered based on the logic of NeRF--.

For training settings, we use a scale factor of 4 and a batch size of 4096, with 200K iterations. This aligns with the training methods of current mainstream models. Even with a scale factor of 2 and a batch size of 8192, our GPU memory usage does not exceed 15 GB. We evenly distribute the batch samples across each input image, so the number of samples per image depends on the total number of images in this scene. We train our model using the Adam optimizer [14].

### A.2. Annotation Mapping Details

We directly leverage OR-NeRF's efficient multiview segmentation approach to remove obstacles and construct our dataset [55]. Its multiview segmentation process is both efficient and consistent. When given point prompts on a single view, the system projects these points into 3D space using COLMAP's sparse reconstruction, establishing correspondences between 2D points and the 3D point cloud. These 3D points are then projected back to all 2D images using camera parameters, creating consistent annotations across all views. Once annotations are propagated to all views, the SAM predicts masks for each view at approximately two frames per second, without requiring neural network training for each scene.

### A.3. Evaluation Settings

Due to the irregular occlusion masks in occluded images, we rearrange valid pixels from ground truth and rendered images into rectangular formats suitable for SSIM and LPIPS patch-based evaluation. This rearrangement may introduce slight variations in metrics compared to methods that directly compare original images, as the structural changes can affect SSIM and LPIPS scores. However, these differences are typically minimal and do not impact the overall evaluation results.

Considering the unavoidable occlusions when capturing real-world scenes, we calculate the rendering accuracy only within the valid visible regions using masks. Therefore, we suggest readers interpret the quantitative evaluation metrics reasonably and place more emphasis on the qualitative results, which demonstrate the true rendering performance in scenes with occlusion removal.

## B. Dataset Details

### B.1. Dataset Building Process

For the DeclutterSet, we capture each scene using either a Canon R6 Mark II camera or an iPhone 12 Pro, maintaining consistent exposure and focus settings throughout the capture process. To ensure high-quality multi-view inputs, we record continuous video while moving the camera in a smooth arc trajectory around the scene. From each recording, we extract 30-35 sequential frames at regular intervals, creating a forward-facing dataset similar to the classic NeRF format. We pay attention to select scenes with varying occlusion characteristics - different depths, scales, and geometric complexity. Camera parameters are estimated using COLMAP's structure-from-motion pipeline. For occlusion annotation, we used OR-NeRF's efficient multiview segmentation approach, requiring only point prompts on a single view to generate consistent masks across all views.

### B.2. Considerations

While OCC-NeRF [58] provides some occlusion datasets, community feedback (as evidenced by multiple issues raised in its repository) has identified several issues with their data. These include blurry images, missing parameters, and even mismatches in ground truth for testing. Even the authors' model and code failed to reproduce their reported results.

To address these shortcomings, we constructed *DeclutterSet*, which includes a variety of occlusion types, varying occlusion sizes and camera motions, and different occluder distances. As stated in the main text, it combines reliable data from existing references and is augmented with newly captured scenes, offering a new and robust benchmark for the community.
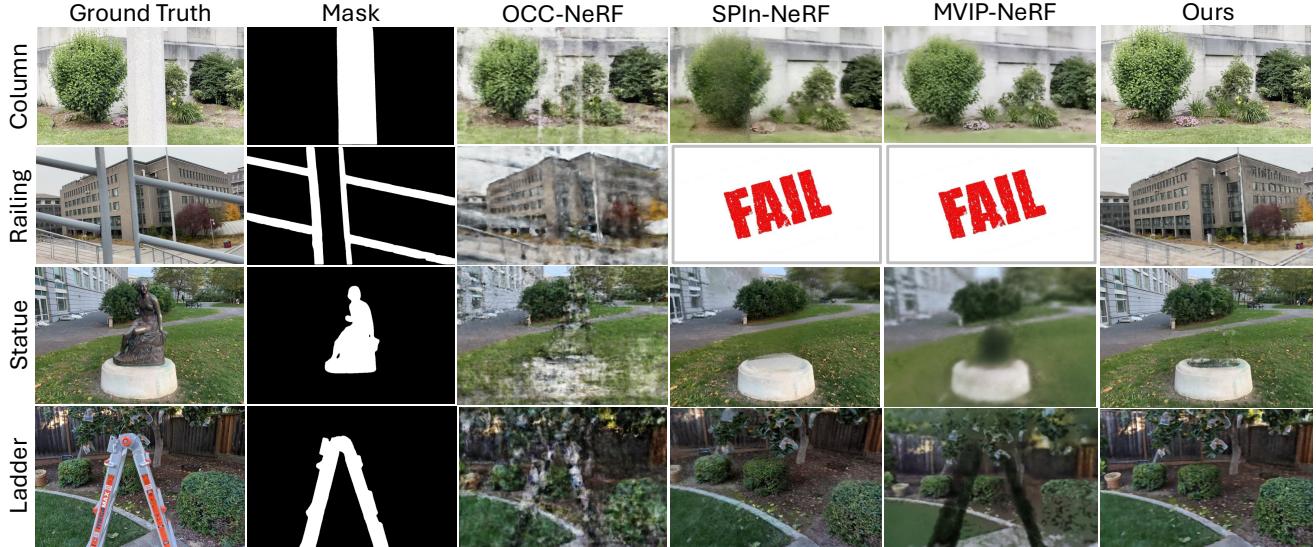
Figure 8. **Additional Qualitative Comparisons With Baselines.** Our method consistently produces desirable results, while generative models still suffer from artifacts and floaters during rendering. Notably, DeclutterNeRF maintains geometric fidelity and cross-view consistency in challenging occlusion scenarios with complex depth relationships. A detailed analysis of failure cases is provided in Sec. C.

## B.3. Samples Exhibition

Figure 9 and Fig. 10 show more samples from our DeclutterSet. We select image frames that are evenly distributed to characterize our dataset: (i) wider distance distribution, (ii) larger occluded regions, (iii) greater relative motion between viewpoints and occluders, and (iv) more uncertain occluder shapes and mask layouts.

## C. Additional Qualitative Results

Figure 8 shows additional visual results on our collected dataset. Beyond normal results, our method demonstrates remarkable robustness by producing high-quality renderings even when faces with incorrect camera parameters from OCC-NeRF data. This issue originates from the OCC-NeRF dataset itself. Specifically, while incorporating existing scenes to complement DeclutterSet, we observed that the *Railing* scene in OCC-NeRF suffers from camera calibration inconsistencies. Although we attempted to re-estimate the camera poses using COLMAP, the anomalies persisted. Nonetheless, we retained this scene in our dataset to reflect the realistic challenges posed by imperfect calibration—an inherent difficulty in occlusion removal tasks. As shown, baseline methods without camera parameter optimization fail to generate converged results and coherent reconstructions. OCC-NeRF produces only blurred representations, while our method successfully recovers a clear scene despite the adverse calibration conditions.
**Failure Cases.** The label "FAIL" in qualitative results is used to denote two distinct failure cases. (i) For SPIn-NeRF, it indicates that reconstruction was not accessible even be-

fore rendering, due to the lack of reliable depth information provided by COLMAP. (ii) For MVIP-NeRF, it refers to a failure that occurred during rendering, where the training process did not converge, resulting in extremely blurred and semantically meaningless images.

To balance reconstruction quality and memory usage when using SPIn-NeRF with COLMAP, we uniformly apply a downsampling factor of 4.

## D. Statement

### D.1. Ethics Statement

Due to concerns about the misuse of generative models and image processing techniques, both 2D and 3D generation have to face these issues. Our DeclutterNeRF, which does not employ any generative priors, mitigates these concerns to a certain extent. This approach helps to avoid potential ethical issues associated with generative models while still achieving effective results in our specific domain.

### D.2. Open Source Statement

Through extensive experimentation with numerous baseline methods, we have identified some opportunities for improvement in the field. Many technical repositories lack proper maintenance and guidance. We recognize that to achieve occlusion removal in NeRF, 3DGS and similar fields, it is first necessary to remove the barriers that exist in the dissemination and communication of these technologies. To this end, all code and data will be open-sourced under the MIT license for community use, fostering transparency and collaborative advancement in the field.
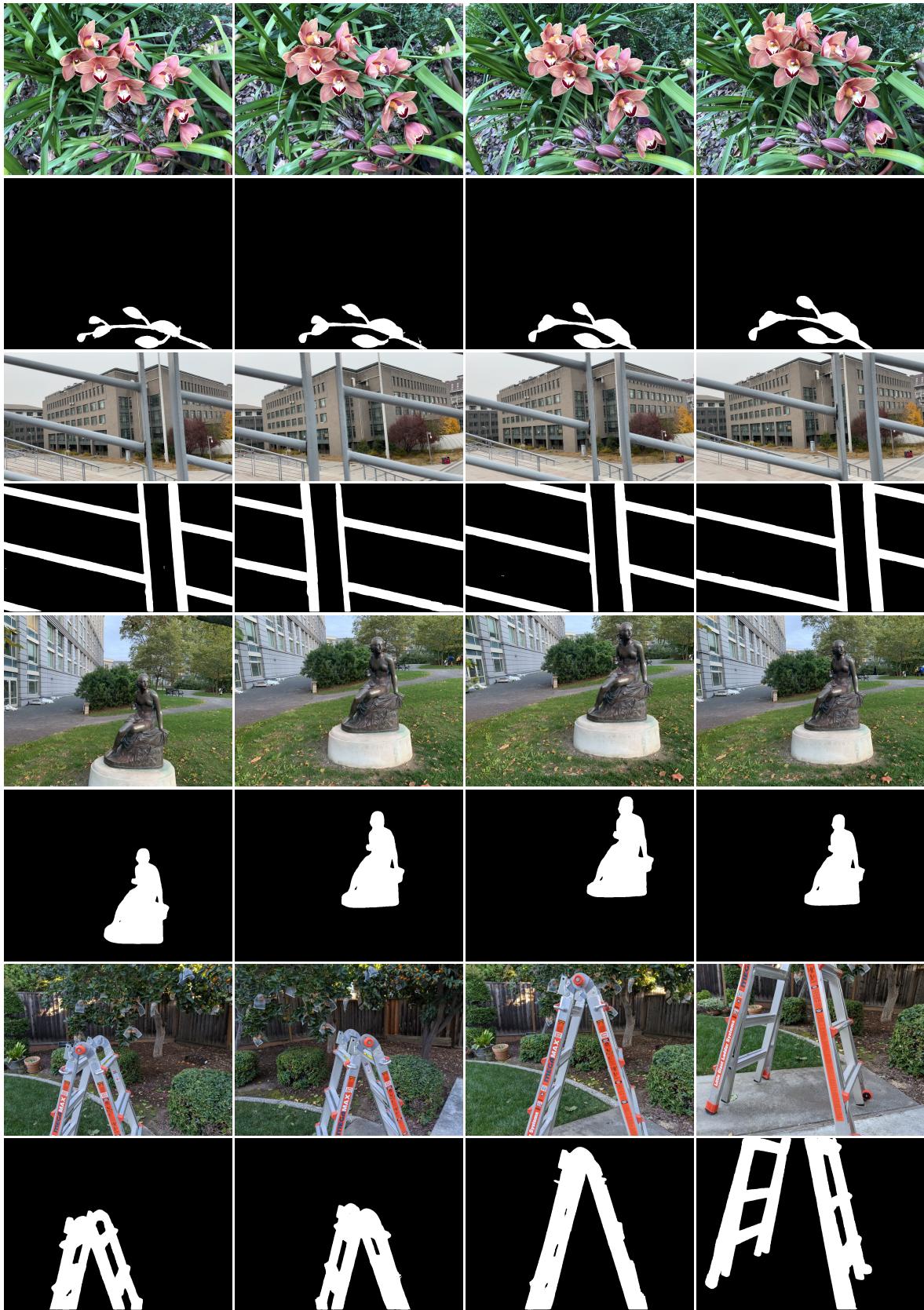
Figure 9. **DeclutterSet Illustration (Part I).** From the top to the bottom: (a) Orchids, (b) Railing, (c) Statue, (d) Ladder.

Figure 10. **DeclutterSet Illustration (Part II).** (e) Stone Column, (f) Lamp Post, (g) Chain Fence, (h) Chair Back.

4