

# Robust SG-NeRF: Robust Scene Graph Aided Neural Surface Reconstruction

Yi Gu    Dongjun Ye    Zhaorui Wang    Jiaxu Wang    Jiahang Cao    Renjing Xu  
HKUST(GZ)

{ygu425, zwang408, jwang457, jcao248, renjingxu}@connect.hkust-gz.edu.cn  
imath@omnispaceai.com

<https://rsg-nerf.github.io/RSG-NeRF/>

## Abstract

Neural surface reconstruction relies heavily on accurate camera poses as input. Despite utilizing advanced pose estimators like COLMAP or ARKit, camera poses can still be noisy. Existing pose-NeRF joint optimization methods handle poses with small noise (inliers) effectively but struggle with large noise (outliers), such as mirrored poses. In this work, we focus on mitigating the impact of outlier poses. Our method integrates an inlier-outlier confidence estimation scheme, leveraging scene graph information gathered during the data preparation phase. Unlike previous works directly using rendering metrics as the reference, we employ a detached color network that omits the viewing direction as input to minimize the impact caused by shape-radiance ambiguities. This enhanced confidence updating strategy effectively differentiates between inlier and outlier poses, allowing us to sample more rays from inlier poses to construct more reliable radiance fields. Additionally, we introduce a re-projection loss based on the current Signed Distance Function (SDF) and pose estimations, strengthening the constraints between matching image pairs. For outlier poses, we adopt a Monte Carlo re-localization method to find better solutions. We also devise a scene graph updating strategy to provide more accurate information throughout the training process. We validate our approach on the SG-NeRF and DTU datasets. Experimental results on various datasets demonstrate that our methods can consistently improve the reconstruction qualities and pose accuracies.

## 1. Introduction

Reconstructing the surfaces of objects from multi-view images is a fundamental challenge in both computer vision and computer graphics. Inspired by Neural Radiance Fields [39] (NeRF), recent strides [29, 40, 59, 64] have marked significant progress in neural surface reconstruction (NSR) area by leveraging implicit scene representations and volume rendering techniques. In NSR, scene geometry is encoded

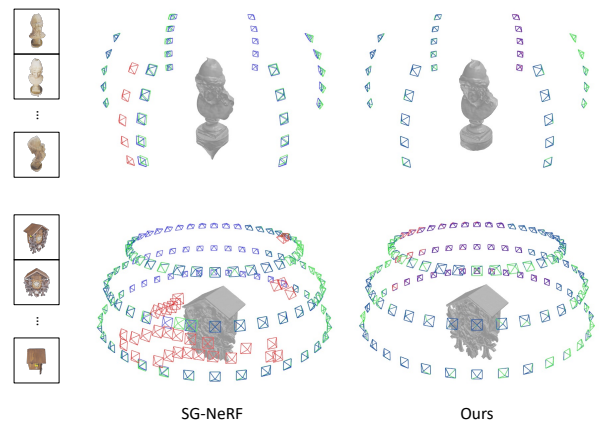


Figure 1. Reconstruction results on the SG-NeRF [6] dataset. Both SG-NeRF [6] and our method take the same initial poses as input, including significant noises. The camera poses are also presented with optimized outlier poses, inlier poses and ground truth poses. More results are illustrated in the supplementary material.

through a signed distance function (SDF), which is learned by a multilayer perceptron (MLP) network trained with an image-based rendering loss. Despite these promising advancements, a key challenge in NSR involves the dependency on accurate camera poses. In practice, NeRF and its variants often rely on COLMAP [50, 51], a widely-used Structure from Motion (SfM) framework, to estimate camera poses prior. Unfortunately, these pose estimations can be significantly erroneous, adversely affecting the reconstruction quality of NeRF. Consequently, recent efforts [2, 3, 5, 8, 9, 26, 30, 57, 62] have aimed to joint optimize scene representations and camera poses to minimize the impact of pose errors. Nevertheless, most of these efforts concentrate on refining relatively small pose errors (referred to as inliers). It is still a challenge to rectify noticeably incorrect camera poses (referred to as outliers). To alleviate the negative effects of outliers, SG-NeRF [6] intro-

duces scene graphs to enhance camera pose optimization for improved geometric reconstruction. The main contribution of SG-NeRF lies in estimating the confidence of each camera pose. By prioritizing ray sampling from images with high confidence poses, SG-NeRF can recover reliable geometry, even in the presence of numerous outliers, as illustrated in the left of Fig. 1. Theoretically, the extreme of the SG-NeRF philosophy is not sampling on outliers. Thus, it is important to recognize the inliers and outliers. However, the heuristic confidence updating strategy in SG-NeRF only depends on the peak signal-to-noise ratio (PSNR) index, which can not well reflect the differences between the inliers and outliers. As shown in Fig. 2, with the wrongly estimated poses, SG-NeRF can still render images with high PSNR. This is a classical shape-radiance ambiguity problem in NeRF series [17, 73, 74]. A potential solution is to add some regularization terms to the loss function, but it may be more complicated coupled with a joint pose-NeRF optimization process.

To address this problem, we explore an improved confidence estimation method to distinguish inliers and outliers. As shown in the last column of Fig. 2, we empirically find that a color network without viewing direction as input can provide more reliable information about pose confidence, albeit possibly at the expense of rendering and geometric quality. It is important to note that our primary objective is to identify inliers and outliers based on estimated confidence. To achieve this, our framework incorporates two color networks: one that aligns with traditional NSR approaches, and another is dedicated to confidence estimation. We detach the latter from the main pose-NeRF optimization graph to maintain the performance of the NSR backbone. Typically, the color network used in NSR is a shallow MLP [38, 59, 73], which means that our method does not substantially increase computational costs. This straightforward design allows us to establish a rule-based threshold to identify inliers and outliers. Subsequently, we can enhance the final results by integrating tailored designs for handling each. We present two strategies: for outliers, we employ a Monte Carlo re-localization method to provide better initialization; for inliers, we enhance constraints with re-projection and Intersection-of-Union(IoU) losses. Additionally, we devise a scene graph updating strategy based on the current SDF to eliminate incorrectly matched pairs. Experiments on the SG-NeRF [6] and the DTU [25] datasets generally show that our method not only yields high-quality 3D reconstructions but also effectively corrects outlier poses, as illustrated in the right of Fig. 1. Our contributions can be summarized as follows:

- We propose a plug-and-play camera pose confidence estimation method that effectively identifies inliers and outliers.
- We introduce strategies such as Monte Carlo re-

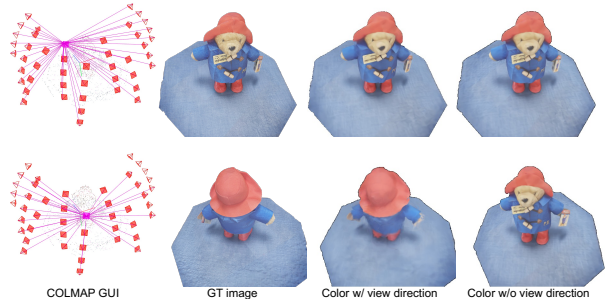


Figure 2. The illustration of the pose ambiguity. The first row is the results from inliers and the second row presents outliers. Images in the first column come from COLMAP [50] GUI, which show that both these two poses are registered in front of the object. However, the ground truth images in the second column show the opposite phenomenon. The third column presents the rendering results of SG-NeRF [6], which use the same color network as NeuS [59] with view direction as input. As shown in the fourth column, our method incorporates an isolated color network, which can well recognize this ambiguity.

localization for outliers and re-projection and IoU losses for inliers to improve geometric constraints.

- Additionally, we implement a scene graph updating strategy to enhance training guidance. To the best of our knowledge, this is the first study to update matching pairs dynamically during the pose-NeRF joint training process.

## 2. Related works

**Neural Surface Reconstruction.** Traditional multi-view stereo (MVS) methods [1, 19, 20, 51] explicitly establish dense correspondences across multiple images to generate depth maps, which are subsequently fused into a global dense point cloud [37, 72]. Surface reconstruction is typically performed as a post-processing step, employing techniques such as screened Poisson surface reconstruction [28]. The processes of searching for correspondences and estimating depth have been significantly enhanced by deep learning-based approaches [66, 67]. Recently, the implicit representation has gained a lot of attention due to its continuity and capability to achieve high spatial resolutions. Building on the pioneering work of Neural Radiance Fields [38] (NeRF), many successors [18, 29, 59, 60, 64, 69, 71] integrate the signed distance function (SDF) into NeRF to enhance geometric modeling. Among these, NeuS [59] is particularly noteworthy for its ability to produce high-quality reconstructions and successfully handle scenes with severe occlusions and complex structures. Thus, in this study, we select NeuS to represent our scenes.

**Structure from motion (SfM) and (re)localization.** NeRF and its variants require accurate camera poses as in-

put [38, 59, 68]. In real-world applications, Structure from Motion (SfM) [12, 31, 34, 44, 50, 53, 55, 63] techniques are commonly employed for data pre-processing. SfM organizes a set of unstructured images by estimating camera poses and triangulating 3D scene points. An essential byproduct of this process is the scene graph, which captures information about matching pairs. However, current advanced SfM frameworks primarily depend on keypoint detection [14, 16, 35, 46, 58] and matching [32, 48, 54] techniques, which can be less effective in textureless or repetitive environments.

The task of (re)localization [4, 15, 47, 49] is also closely related to SfM. Given a database of posed images, the goal of this task is to estimate the camera poses of newly captured images. In the context of NeRF with re-localization, most existing studies [33, 36, 41, 70] concentrate on relocating new images within well-constructed NeRFs. In our approach, we implement the Monte Carlo re-localization method during the training phase to improve the robustness and accuracy of outlier poses.

**Joint NeRF and pose optimization.** NeRFmm [62] and iNeRF [70] demonstrate the potential for jointly learning or refining camera parameters alongside the NeRF framework. Following works [2, 7, 10, 24, 30, 45, 61] also perform different modular modifications. For example, GARF [9] and SiNeRF [65] capitalize on the inherent smoothness of non-traditional activations to mitigate the impact of noisy gradients caused by high-frequency components in positional embeddings. L2G-NeRF [5] and Invertible Neural Warp [11] tackle the camera pose representation with an overparameterization strategy. NoPe-NeRF [3] employs an external monocular depth estimation model to assist in refining camera poses. Some works [2, 6, 26, 57] also incorporate cross-view correspondences to enhance geometry constraints. Commonly, most approaches presume that all images are properly posed initially and focus on local optimizations for pose correction. Notably similar to our method is SG-NeRF [6], which is the first to utilize a scene graph to guide joint optimization. Our work follows this innovative path but with a modified confidence estimation strategy.

### 3. Methods

In this section, we first define the problem setting and provide an overview of the proposed pipeline. Subsequently, we delve into the key technical designs in detail.

**Problem statement.** Our research focuses on the object-level 3D surface reconstruction from a set of unorganized images captured in an inward-facing configuration. Specifically, given a collection of RGB images  $\mathbf{I} =$

$\{I_1, I_2, \dots, I_N\}$ , our objective is to reconstruct the 3D surface  $S$  of the scene. For a specific image  $I_i$ , a key output of our approach is the optimized camera pose  $P_i = (R_i, t_i)$ , where  $R_i$  belongs to  $\mathbf{SO}(3)$  representing the rotation and  $t_i$  is a vector in  $\mathbb{R}^3$  representing the translation. Additionally, each pose is assigned an inlier-outlier confidence score.

**Method overview.** Fig. 3 illustrates the workflow of the proposed pipeline. In the data preparation stage, we first employ a widely-used Structure-from-Motion (SfM) algorithm, specifically COLMAP [50], to obtain the initial camera poses. Given the potential inaccuracies of these poses, proceeding with a direct joint pose-NeRF optimization could be catastrophic. To mitigate this risk, we leverage scene graph information to guide the training process (Sec. 3.1). We update the confidence with our tailored indicator, which can effectively distinguish inlier and outlier poses. For inliers, we introduce additional constraints to enhance the geometric consistency (Sec. 3.2). For outliers, we utilize Monte Carlo re-localization to find better initializations (Sec. 3.1). Additionally, We also devise a scene graph updating strategy to enhance the guidance during training (Sec. 3.4).

#### 3.1. Scene graph guided confidence estimation

A scene graph  $\mathbf{G} = (\mathbf{V}, \mathbf{E})$  in SfM consists of a set of nodes  $\mathbf{V}$  and edges  $\mathbf{E}$ . Each node  $v_i \in \mathbf{V}$  corresponds to an input image  $I_i \in \mathbf{I}$ , and an edge between two nodes contains the matching and co-visibility information about the corresponding images. We annotate all edges as  $\mathbf{M} = \{M_{i,j} | v_i, v_j \in \mathbf{V}, v_i \neq v_j\}$ , where the set  $M_{i,j}$  comprises all matched keypoint pairs between  $I_i$  and  $I_j$ .

The original scene graph tends to be dense and contains many incorrect matches. Following SG-NeRF [6], we set an angular threshold  $\tau$  for the estimated relative rotations and remove any edges exceeding  $\tau$ . Then, each node is assigned a confidence estimate based on this sparsified scene graph.

The confidence score for a node  $v_i$  is defined as the mean number of matching pairs, which can be computed as:

$$CS(v_i) = \frac{\sum_{M_{i,j} \in \mathbf{M}_i} |M_{i,j}|}{|\mathbf{M}_i|}, \quad (1)$$

where  $|\cdot|$  denotes the number of elements in a set, e.g.,  $|M_{i,j}|$  is the total number of matching pairs of  $I_i$  and  $I_j$  and  $|\mathbf{M}_i|$  is the total number of edges of  $v_i$ . A higher score indicates that the image has a better matching quality and a higher likelihood of being an inlier. We normalize this confidence score via  $CS(v_i) = CS(v_i) / \sum_{v \in \mathbf{V}} CS(v)$  to form a probability distribution, which guides the training to sample more rays from poses with higher confidence. All confidence computations involve a normalization step and we omit this step in the following text for brevity.

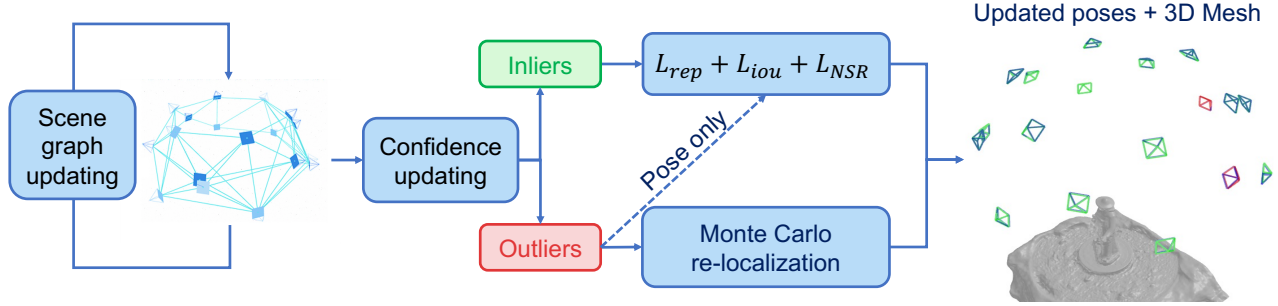


Figure 3. An overview of the proposed pipeline. Given the initial scene graph, we apply a confidence updating strategy based on an indicator from a detached color network, which can identify inlier and outlier poses. For inliers, we utilize re-projection loss and IoU loss to enhance the geometric constraints. For outliers, we utilize Monte Carlo re-localization method to find better initializations. The scene graph is also updated based on current geometry and pose estimations. Eventually, our method can reconstruct the 3D mesh from the trained field and rectify both inlier and outlier poses with high accuracy. The coloration is same as Fig. 1.

These initial scores are derived from keypoint matches, which might lack a comprehensive understanding of the information contained in images. Thus, we adaptively update the confidence scores based on the image rendering quality. Specifically, we estimate the peak signal-to-noise ratio (PSNR) for each image according to current image rendering loss for efficiency. Then, the confidence scores are updated by [6]:

$$CS(v_i) = CS(v_i) + \lambda_c PSNR(v_i). \quad (2)$$

However, as shown in Fig. 2, we empirically found that the PSNR of outliers can be even larger than that of inliers, which means that more outliers will be sampled. The reason behind this phenomenon comes from shape-radiance ambiguity [17, 62, 73, 74].

To solve this problem, we employ a new color network  $C_n$  which does not take viewing direction as input. To mitigate the shape-radiance ambiguity and prevent overfitting, we use the same sampling points and geometry features as the original color network  $C_o$ . We detached all relevant computations from  $C_n$  to ensure that the loss from  $C_n$  does not impact the main networks. Since we only require an indicator that can reflect the relative rendering qualities of the training images, the PSNR estimated by  $C_n$  can serve the same purpose as that by  $C_o$ . As highlighted in NeRF++ [73], most existing works [38, 59] use a shallow MLP for color network, which acts as an implicit regularization. Thus,  $C_n$  will not introduce significant computational overhead. We use the PSNR estimated from  $C_n$  (denoted as  $PSNR_n$ ) to update the confidence score throughout the training process.

It should be noted that we also keep a record of the PSNR with  $C_o$  (denoted as  $PSNR_o$ ), which can be helpful for filtering out outliers. When  $PSNR_o$  and  $PSNR_n$  show a significant discrepancy, it is an indication of anomalous data. Therefore, we recognize poses with  $|PSNR_o -$

$PSNR_n| > \tau_1$  as outliers.

### 3.2. Joint optimization

We build up our framework based on NeuS [59]. The neural surface reconstruction loss function is defined as follows:

$$\mathcal{L}_{NSR} = \mathcal{L}_{color}(C_o) + \mathcal{L}_{color}(C_n) + \lambda \mathcal{L}_{reg}. \quad (3)$$

The  $\mathcal{L}_{color}(C_o)$  represents calculating  $\mathcal{L}_{color}$  by  $C_o$ . The  $\mathcal{L}_{color}(C_n)$  is calculated by  $C_n$ , with gradients only back-propagated to  $C_n$ . The  $\mathcal{L}_{color}$  is a photometric loss:

$$\mathcal{L}_{color} = \left\| \hat{\mathbf{C}} - \mathbf{C} \right\|_1, \quad (4)$$

where  $\hat{\mathbf{C}}$  is obtained by volume rendering equation [27, 38] and  $\mathbf{C}$  is the ground truth color.

The  $\mathcal{L}_{reg}$  incorporates the Eikonal term [21] applied to the sampled points to regularize the learned SDF, which can be expressed as:

$$\mathcal{L}_{reg} = \frac{1}{k} \sum_{i=1}^k (\|\nabla f(p_i)\|_2 - 1)^2, \quad (5)$$

where  $f(p_i)$  represents the distance estimate for each sampled 3D location along the ray.

We also utilize the Intersection-of-Union (IoU) loss  $\mathcal{L}_{iou}$  and re-projection loss  $\mathcal{L}_{rep}$  [23] to further improve the pose accuracy. The  $\mathcal{L}_{iou}$  loss, firstly proposed by SG-NeRF [6], not only enhances geometry consistency but also accelerates convergence. Another related constraint is the epipolar loss proposed by PoRF [2]. However, we find that both epipolar and IoU losses do not handle outliers effectively. In fact, the re-projection loss can fulfill the same role as the epipolar loss [23] but is more reliable to the scene geometry. The epipolar loss does not require depth for back-projection



but is invariant to the scale of the translation part. Considering the aforementioned analysis, we opt for the IoU loss and the re-projection loss in our framework.

Given a pair of matched keypoints  $kp_i$  from image  $I_i$  and  $kp_j$  from image  $I_j$ , we define the Intersection Volume as:

$$I = MoG(kp_i) \cdot MoG(kp_j), \quad (6)$$

and Union Volume as:

$$U = MoG(kp_i) + MoG(kp_j) - I, \quad (7)$$

where  $MoG(\cdot)$  is a mixture of Gaussians for sampling points along a ray. The IoU loss can be computed as:

$$\mathcal{L}_{iou} = 1 - \frac{I}{U}. \quad (8)$$

With a set of points sampled from the ray corresponding to  $kp_i$ , we approximate the depth  $d_i$  of  $kp_i$  by selecting the point with maximal weights. Thus, the re-projection loss can be achieved by:

$$\mathcal{L}_{rep} = L_\delta(kp'_i, kp_j), \quad (9)$$

where  $kp'_i$  is the re-projected point of  $kp_i$  in image  $I_j$  and  $L_\delta$  represents the Huber loss.

We jointly optimize inlier-inlier pairs, while bypassing outlier-outlier pairs. For inlier-outlier pairs, we only optimize the poses of outliers and keep inliers and NeuS backbone fixed. Our overall loss is defined as:

$$\mathcal{L} = \mathcal{L}_{NSR} + \alpha \mathcal{L}_{iou} + \beta \mathcal{L}_{rep}. \quad (10)$$

### 3.3. Monte Carlo re-localization

Geometry constraints in Sec. 3.2 can still struggle with certain extreme cases. One intractable case comes from the mirror-symmetry ambiguity, which has been extensively studied in SfM [42, 43, 52]. In the context of pose-NeRF joint optimization, NeRFmm [62] and LU-NeRF [8] also mentioned the same problem. LU-NeRF solves this problem by training two NeRF models, one of which uses reflected poses, requiring significant time to find the mirror poses.

Leveraging our confidence scheme, we can easily detect outliers, particularly those mirrored outliers. To maximize the use of training images, we propose to utilize Monte Carlo re-localization techniques [13, 36] to assist outliers in finding better initializations. Specifically, as we focus on inward-facing scenes, we first estimate a coarse main axis of the scene using inlier poses. The rotation around this main axis is defined as  $R_{axis}(\theta)$ , where  $\theta$  is the angle of rotation. We then distribute the initial particles uniformly around this axis. Given an outlier pose  $R_o, t_o$ , the poses of these particles can be obtained by:

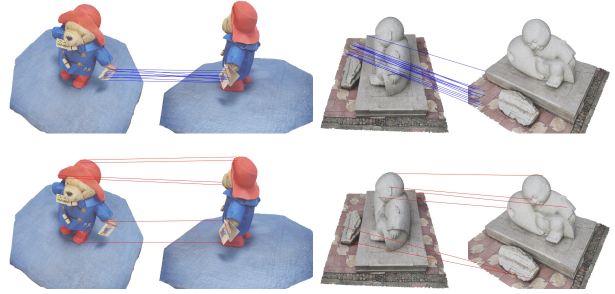


Figure 4. The illustration of scene graph updating. We filter out the wrong matched keypoint pairs (colored in red lines) and keep the correct pairs (colored in blue lines). We select image pairs with relatively few matching pairs for clearer visualization. Additional results are detailed in the supplementary materials.

$$R_{pi} = R_{axis}\left(\frac{i \cdot 2\pi}{N_p}\right) \cdot R_o, \quad i \in \{1, 2, \dots, N_p\}, \quad (11)$$

and

$$t_{pi} = R_{axis}\left(\frac{i \cdot 2\pi}{N_p}\right) \cdot t_o, \quad i \in \{1, 2, \dots, N_p\}, \quad (12)$$

where  $(R_{pi}, t_{pi})$  is the pose of  $i$ -th particle and  $N_p$  is the number of particles. We fix all network components and only optimize the poses of particles. Initially, each particle is sampled equally for optimization. Subsequently, the sampling probability is adjusted based on the estimated  $PSNR_n$ . If the maximum  $PSNR_n$  of the particles exceeds that of the current outlier at the end of the re-localization stage, we replace  $(R_o, t_o)$  with the pose of this particle. Additional details about our Monte Carlo re-localization are provided in the supplementary materials.

### 3.4. Scene graph updating

The initial confidence, based on results from SfM, may be sub-optimal. Therefore, we periodically update the scene graph according to current geometry and pose estimations. Similar to Sec. 3.1, we use the same angular threshold  $\tau$  to remove edges from the raw graph. For remaining keypoint matching pairs, we remove those with a re-projection loss surpassing the threshold  $\tau_{rep}$ , which is gradually reduced throughout the training. As illustrated in Fig. 4, our method effectively eliminates wrong matches, providing more reliable information for subsequent training iterations.

## 4. Experiments

### 4.1. Experiment setup

Following SG-NeRF [6], we conduct our experiments on 8 cases from the SG-NeRF dataset and 5 cases from the

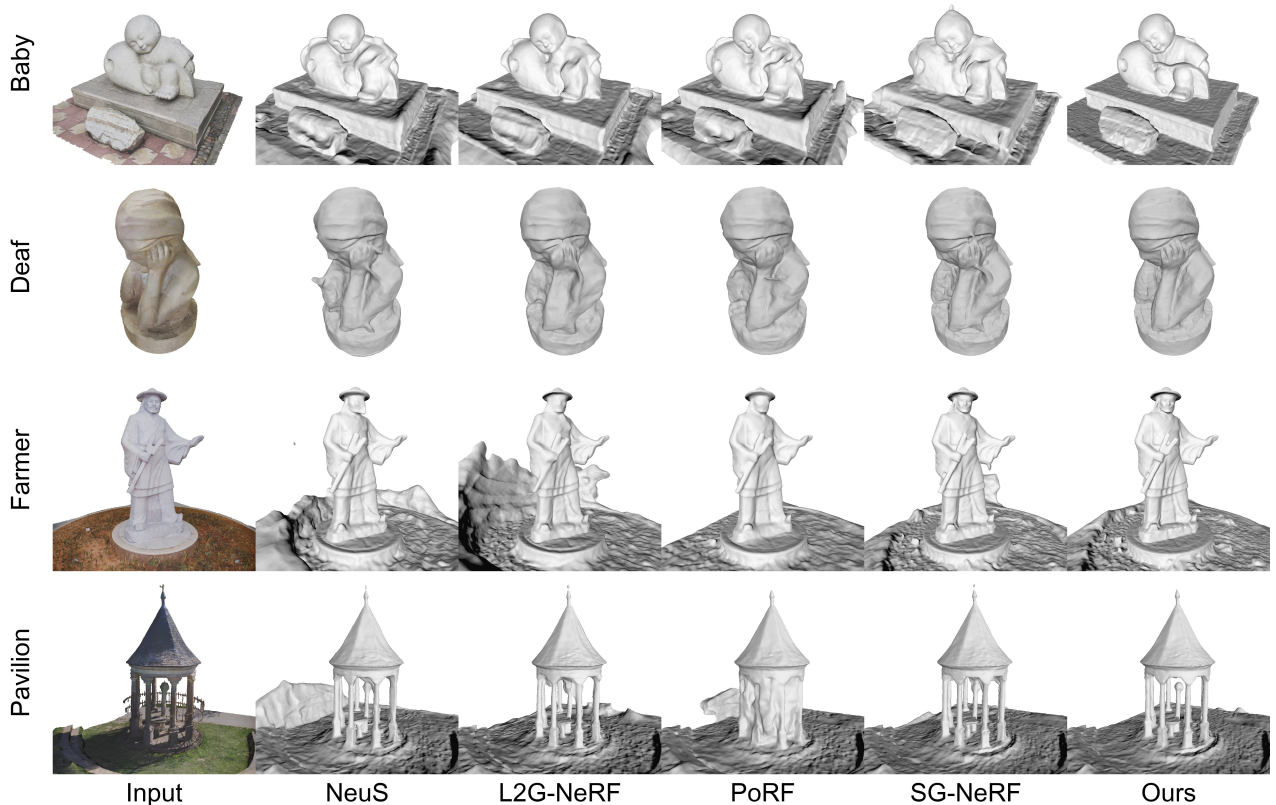


Figure 5. Qualitative comparisons on the SG-NeRF [6] dataset. Our method can generally recover high-fidelity geometry with only one-stage training. More visual comparisons are provided in supplementary materials.

DTU [25] dataset to validate our method. Following literature [29, 59], we assess the mesh quality with Chamfer distance (CD) and F-score metrics. The baseline methods for comparison include BARF [30], SCNeRF [26], GARF [10], L2G-NeRF [5], Joint-TensorRF [7], PoRF [2] and SG-NeRF [6]. Results with \* are achieved in a two-stage manner, including official implementations and NeuS [59] with optimized poses. The initial camera poses of SG-NeRF are obtained by using Superpoint [14] and SuperGlue [48], with COLMAP [47, 63] backend optimization. As presented in SG-NeRF [6], this combination consistently outperforms the standard COLMAP but still results in a proportion of significant incorrect poses, ranging from 1/9 to 1/3. The initial poses for DTU are obtained by conventional COLMAP first. To simulate outlier poses, SG-NeRF randomly selects 1/7 to 1/4 of the images for each scene and injects random noises to their poses. For a fair comparison, all methods, including ours, use the same initial poses as input.

## 4.2. Implementation details

We implement our method based on NeuS [59]. The camera poses are implemented by Lictorch [56], which can perform backpropagation on  $SE(3)$  Groups. Following SG-

NeRF [6], the angular threshold  $\tau$  for scene graph sparsification is set as 70 for SG-NeRF dataset [6] and 45 for DTU [25] dataset, respectively. The inlier-outlier threshold is set as  $\tau_1 = 9$ , which is an extremely large performance gap for  $PSNR_o$  and  $PSNR_n$ . The weights of loss are set as  $\lambda = 0.1$ ,  $\alpha = 0.2$ , and  $\beta = 0.001$  respectively. The particle number  $N_p$  is set as 24 for efficiency. Other configurations are kept the same as NeuS. All experiments are conducted on NVIDIA RTX 3090 GPUs. Our method runs an average of 13 hours for 150k iterations on the SG-NeRF dataset, and 22 hours for 300k iterations on the DTU dataset.

## 4.3. Comparisons

**Results on SG-NeRF.** The quantitative results are reported in Table 1. Both NeuS [59] and Neuralangelo [29] degenerate severely due to significantly noisy camera poses. PoRF [2] and SCNeRF [26] demonstrate commendable results in certain cases, highlighting the importance of incorporating cross-view correspondences. Among the competitors, SG-NeRF [6] achieves the best overall performance, underscoring the effectiveness of scene graph guidance. Our method consistently outperforms other approaches by a

Table 1. Quantitative results on SG-NeRF [6]. The **red** and **blue** numbers indicate the first and second performer for each scene. † denotes that only valid values are used for the average. Methods with \* are trained in a two-stage manner.

		Baby	Bear	Bell	Clock	Deaf	Farmer	Pavilion	Sculpture	Mean
Chamfer distance ↓	NeuS [59]	0.69	0.31	3.33	1.16	0.55	2.49	0.29	0.66	1.18
	Neuralangelo [29]	0.70	0.65	-	0.38	0.59	4.89	1.95	0.31	1.35 <sup>†</sup>
	BARF [30]*	1.08	0.28	3.31	0.19	0.46	2.13	0.38	0.57	1.05
	SCNeRF [26]*	1.19	0.27	3.74	1.33	0.46	1.45	0.23	0.81	1.19
	GARF [10]*	2.04	2.25	3.08	2.01	0.59	1.58	0.96	0.57	1.64
	L2G-NeRF [5]*	1.15	0.29	1.26	0.24	0.40	2.18	-	4.36	1.41 <sup>†</sup>
	Joint-TensoRF [7]*	3.11	-	2.49	0.36	0.88	2.51	1.35	0.70	1.63 <sup>†</sup>
	PoRF [2]	<b>0.31</b>	0.49	-	-	<b>0.30</b>	3.80	2.20	-	1.42 <sup>†</sup>
	SG-NeRF [6]	0.56	<b>0.25</b>	<b>0.98</b>	<b>0.15</b>	0.45	<b>0.87</b>	<b>0.20</b>	<b>0.22</b>	<b>0.46</b>
	Ours	<b>0.07</b>	<b>0.09</b>	<b>1.22</b>	<b>0.15</b>	<b>0.13</b>	<b>0.62</b>	<b>0.17</b>	<b>0.09</b>	<b>0.32</b>
F-score ↑	NeuS [59]	0.65	0.93	0.48	0.72	0.84	0.54	0.93	0.70	0.74
	Neuralangelo [29]	0.57	0.80	-	0.85	0.66	0.14	0.47	0.89	0.63 <sup>†</sup>
	BARF [30]*	0.58	0.91	0.49	0.95	0.86	0.51	0.86	0.87	0.75
	SCNeRF [26]*	0.56	0.93	0.49	0.69	0.86	0.59	<b>0.95</b>	0.73	0.72
	GARF [10]*	0.18	0.21	0.50	0.27	0.78	0.57	0.41	0.83	0.47
	L2G-NeRF [5]*	0.58	0.92	0.65	0.92	0.89	0.49	-	0.21	0.67 <sup>†</sup>
	Joint-TensoRF [7]*	0.20	-	0.38	0.84	0.60	0.24	0.34	0.63	0.46 <sup>†</sup>
	PoRF [2]	<b>0.92</b>	0.78	-	-	<b>0.92</b>	0.39	0.35	-	0.67 <sup>†</sup>
	SG-NeRF [6]	0.74	<b>0.93</b>	<b>0.71</b>	<b>0.96</b>	0.87	<b>0.76</b>	0.94	<b>0.92</b>	<b>0.85</b>
	Ours	<b>0.99</b>	<b>0.99</b>	<b>0.65</b>	<b>0.96</b>	<b>0.99</b>	<b>0.79</b>	<b>0.94</b>	<b>0.99</b>	<b>0.91</b>
APE ↓	SG-NeRF [6]	2.16	1.84	2.15	1.80	1.17	0.72	0.6	1.10	1.44
	Ours	0.004	0.004	0.37	0.003	0.016	0.018	0.003	0.006	<b>0.053</b>
RPE ↓	SG-NeRF [6]	2.26	1.38	2.29	0.51	2.10	1.12	1.04	2.02	1.59
	Ours	0.005	0.004	0.70	0.002	0.021	0.022	0.004	0.008	<b>0.096</b>

considerable margin, which shows the effectiveness of our framework. The visual comparison is provided in Fig. 5,

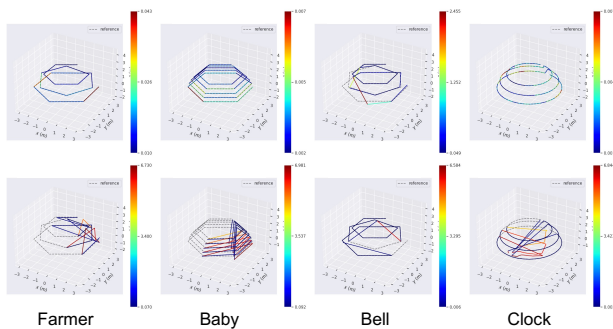


Figure 6. Visualization of pose accuracy. The first row presents our results and the second is SG-NeRF [6]. Dashed lines represent ground truth poses and solid lines are optimized poses. The poses rectified by our method are well aligned with the ground truth poses. In the hardest Bell case, our method can still refine most outliers, while SG-NeRF failed in all cases.

where our method distinctly excels in capturing finer geometric details. However, we empirically observed that all methods, including ours, struggle with the Bell scene, likely due to the sparsity of training images.

We also utilize evo [22] to evaluate the pose accuracy of our method and SG-NeRF [6]. Due to the original SG-NeRF dataset does not provide inlier-outlier information, we utilize our indicator to filter out outliers. We align inliers to ground truth poses to get a global  $SIM(3)$  transformation, which is then applied to all poses. The results of absolute pose error (APE) and relative pose error (RPE) w.r.t. full transformation (including both rotation and translation parts) are reported in Table 1. Our pose accuracy surpasses that of SG-NeRF by more than two orders of magnitude on both RPE and APE. Fig. 6 shows the visual comparison of the camera pose accuracy.

**Results on DTU.** The quantitative results are shown in Table 2. We report a new result of SG-NeRF [6] on Scan 37 with a better performance (originally reported as 2.39),

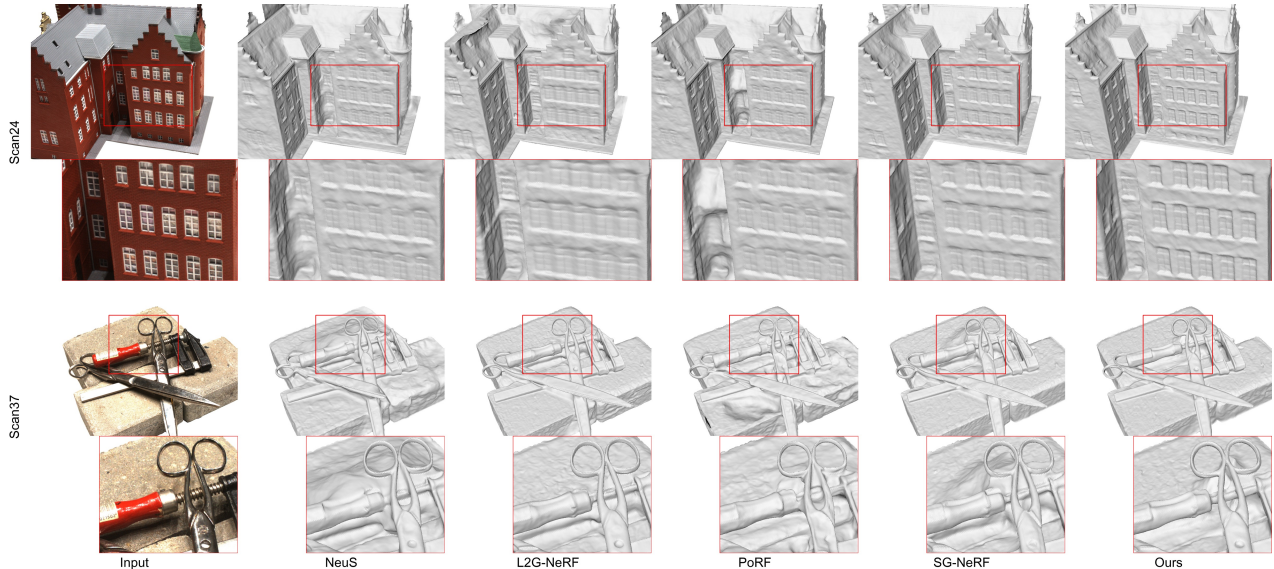


Figure 7. Qualitative comparison on the DTU [25] dataset. L2G-NeRF [5] is trained in a two-stage manner and others are trained in one stage with the same iterations. All methods take the same initial poses as input.

Table 2. Quantitative results on the DTU [25] dataset with noisy camera poses as input.

Chamfer distance ↓	24	37	40	55	63	Mean
NeuS [59]	1.07	2.80	1.52	1.30	3.20	1.98
Neuralangelo [29]	1.06	2.96	1.22	<b>0.42</b>	1.23	1.38
BARF [30]*	1.46	<b>1.40</b>	5.16	1.78	1.80	2.32
SCNeRF [26]*	1.45	2.84	2.60	0.78	1.83	1.90
GARF [10]*	1.18	2.00	2.61	2.37	8.74	3.38
L2G-NeRF [5]*	1.08	1.60	3.27	1.79	6.97	2.94
Joint-TensorRF [7]*	1.00	2.60	-	-	7.71	3.77 <sup>†</sup>
PoRF [2]	1.15	2.33	0.97	0.76	1.30	1.30
SG-NeRF [6]	<b>0.87</b>	1.83	<b>0.88</b>	<b>0.38</b>	<b>1.13</b>	<b>1.01</b>
Ours	<b>0.80</b>	<b>1.30</b>	<b>0.61</b>	0.44	<b>1.09</b>	<b>0.85</b>

due to the fact of our experiment. In DTU [6] dataset, our method performs slightly better than SG-NeRF. We empirically find that our Monte Carlo re-localization has not been triggered. Thus, the experiment on the DTU [25] dataset can be viewed as an improved version of SG-NeRF. Our method outperforms the competitors on four scans and achieves a similar performance with SG-NeRF on Scan 55. A qualitative comparison can be found in Figure 7.

#### 4.4. Ablation studies

To validate the effectiveness of each component of our method, including the re-projection loss, Monte Carlo re-localization, and the scene updating strategy, we select 4 cases from both SG-NeRF [6] and DTU [25] datasets. We evaluate the re-projection loss using the DTU dataset, while the other components are assessed with the SG-NeRF

dataset. All experiments are conducted using our proposed confidence updating strategy. The results are reported in Table 3. Monte Carlo re-localization significantly enhances the results, and the scene updating strategy also contributes to a slight improvement in the final outcomes, confirming the effectiveness of these two components. Although the experiments without the re-projection loss include a cross-view constraint loss (IoU loss), our complete model still demonstrates further improvements.

## 5. Conclusions

This paper addresses neural surface reconstruction from image sets characterized by significant outlier poses. By leveraging the scene graph to guide training, we introduce a novel confidence updating strategy that effectively recognizes inliers and outliers. We enhance geometric constraints through the integration of Intersection-of-Union (IoU) loss and re-projection loss, while employing Monte Carlo re-localization techniques to accurately reposition outliers. These methods, combined with our scene graph updating

Table 3. Ablation studies on SG-NeRF [6] and DTU [25] datasets. We individually remove scene updating (S.U.) and Monte Carlo re-localization (M.C.) on the SG-NeRF dataset. On the DTU dataset, we validate the effectiveness of the re-projection (Rep.) loss.

scene	SG-NeRF [6]			DTU [25]		
	w/o S.U.	w/o M.C.	full	scene	w/o Rep.	full
baby	0.11	0.38	0.07	scan24	0.91	0.80
bear	0.17	0.28	0.09	scan40	0.78	0.61



strategy, enable our framework to achieve state-of-the-art performance on the challenging SG-NeRF dataset. One limitation of our approach is its dependency on a substantial number of inlier poses. As a promising direction for future research, incorporating prior models could make our framework more robust, especially in sparsely captured scenes.

## References

- [1] Connelly Barnes, Eli Shechtman, Adam Finkelstein, and Dan B Goldman. Patchmatch: A randomized correspondence algorithm for structural image editing. *ACM Trans. Graph.*, 28(3):24, 2009. 2
- [2] Jia-Wang Bian, Wenjing Bian, Victor Adrian Prisacariu, and Philip Torr. Porf: Pose residual field for accurate neural surface reconstruction. In *ICLR*, 2024. 1, 3, 4, 6, 7, 8
- [3] Wenjing Bian, Zirui Wang, Kejie Li, Jia-Wang Bian, and Victor Adrian Prisacariu. Nope-nerf: Optimising neural radiance field with no pose prior. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4160–4169, 2023. 1, 3
- [4] Eric Brachmann and Carsten Rother. Visual camera relocalization from rgb and rgb-d images using dsac. *IEEE transactions on pattern analysis and machine intelligence*, 44(9):5847–5865, 2021. 3
- [5] Yue Chen, Xingyu Chen, Xuan Wang, Qi Zhang, Yu Guo, Ying Shan, and Fei Wang. Local-to-global registration for bundle-adjusting neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8264–8273, 2023. 1, 3, 6, 7, 8
- [6] Yiyang Chen, Siyan Dong, Xulong Wang, Lulu Cai, Youyi Zheng, and Yanchao Yang. Sg-nerf: Neural surface reconstruction with scene graph optimization. In *European Conference on Computer Vision (ECCV)*, 2024. 1, 2, 3, 4, 5, 6, 7, 8
- [7] Bo-Yu Cheng, Wei-Chen Chiu, and Yu-Lun Liu. Improving robustness for joint optimization of camera poses and decomposed low-rank tensorial radiance fields. *arXiv preprint arXiv:2402.13252*, 2024. 3, 6, 7, 8
- [8] Zezhou Cheng, Carlos Esteves, Varun Jampani, Abhishek Kar, Subhransu Maji, and Ameesh Makadia. Lu-nerf: Scene and pose estimation by synchronizing local unposed nerfs. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 18312–18321, 2023. 1, 5
- [9] Shin-Fang Chng, Sameera Ramasinghe, Jamie Sherrah, and Simon Lucey. Gaussian activated neural radiance fields for high fidelity reconstruction and pose estimation. In *European Conference on Computer Vision*, pages 264–280. Springer, 2022. 1, 3
- [10] Shin-Fang Chng, Sameera Ramasinghe, Jamie Sherrah, and Simon Lucey. Gaussian activated neural radiance fields for high fidelity reconstruction and pose estimation. In *European Conference on Computer Vision*, pages 264–280. Springer, 2022. 3, 6, 7, 8
- [11] Shin-Fang Chng, Ravi Garg, Hemanth Saratchandran, and Simon Lucey. Invertible neural warp for nerf. In *European Conference on Computer Vision (ECCV)*, 2024. 3
- [12] Zhaopeng Cui and Ping Tan. Global structure-from-motion by similarity averaging. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 864–872, 2015. 3
- [13] Frank Dellaert, Dieter Fox, Wolfram Burgard, and Sebastian Thrun. Monte carlo localization for mobile robots. In *Proceedings 1999 IEEE international conference on robotics and automation (Cat. No. 99CH36288C)*, pages 1322–1328. IEEE, 1999. 5, 1
- [14] Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superpoint: Self-supervised interest point detection and description. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 224–236, 2018. 3, 6
- [15] Siyan Dong, Qingnan Fan, He Wang, Ji Shi, Li Yi, Thomas Funkhouser, Baoquan Chen, and Leonidas J Guibas. Robust neural routing through space partitions for camera relocalization in dynamic indoor environments. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8544–8554, 2021. 3
- [16] Mihai Dusmanu, Ignacio Rocco, Tomás Pajdla, Marc Pollefeys, Josef Sivic, Akihiko Torii, and Torsten Sattler. D2-net: A trainable CNN for joint description and detection of local features. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 8092–8101. Computer Vision Foundation / IEEE, 2019. 3
- [17] Qihang Fang, Yafei Song, Keqiang Li, and Liefeng Bo. Reducing shape-radiance ambiguity in radiance fields with a closed-form color estimation method. *Advances in Neural Information Processing Systems*, 36, 2024. 2, 4
- [18] Qiancheng Fu, Qingshan Xu, Yew Soon Ong, and Wenbing Tao. Geo-neus: Geometry-consistent neural implicit surfaces learning for multi-view reconstruction. *Advances in Neural Information Processing Systems*, 35:3403–3416, 2022. 2
- [19] Yasutaka Furukawa and Jean Ponce. Accurate, dense, and robust multiview stereopsis. *IEEE transactions on pattern analysis and machine intelligence*, 32(8):1362–1376, 2009. 2
- [20] Silvano Galliani, Katrin Lasinger, and Konrad Schindler. Gipuma: Massively parallel multi-view stereo reconstruction. *Publikationen der Deutschen Gesellschaft für Photogrammetrie, Fernerkundung und Geoinformation e. V.*, 25 (361-369):2, 2016. 2
- [21] Amos Gropp, Lior Yariv, Niv Haim, Matan Atzmon, and Yaron Lipman. Implicit geometric regularization for learning shapes. *arXiv preprint arXiv:2002.10099*, 2020. 4
- [22] Michael Grupp. evo: Python package for the evaluation of odometry and slam. <https://github.com/MichaelGrupp/evo>, 2017. 7
- [23] Richard Hartley and Andrew Zisserman. *Multiple view geometry in computer vision*. Cambridge university press, 2003. 4
- [24] Hwan Heo, Taekyung Kim, Jiyoung Lee, Jaewon Lee, Soohyun Kim, Hyunwoo J Kim, and Jin-Hwa Kim. Robust camera pose refinement for multi-resolution hash encoding. *arXiv preprint arXiv:2302.01571*, 2023. 3

- [25] Rasmus Jensen, Anders Dahl, George Vogiatzis, Engin Tola, and Henrik Aanæs. Large scale multi-view stereopsis evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 406–413, 2014. 2, 6, 8
- [26] Yoonwoo Jeong, Seokjun Ahn, Christopher Choy, Anima Anandkumar, Minsu Cho, and Jaesik Park. Self-calibrating neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5846–5854, 2021. 1, 3, 6, 7, 8
- [27] James T Kajiya and Brian P Von Herzen. Ray tracing volume densities. *ACM SIGGRAPH computer graphics*, 18(3):165–174, 1984. 4
- [28] Michael Kazhdan and Hugues Hoppe. Screened poisson surface reconstruction. *ACM Transactions on Graphics (ToG)*, 32(3):1–13, 2013. 2
- [29] Zhaoshuo Li, Thomas Müller, Alex Evans, Russell H Taylor, Mathias Unberath, Ming-Yu Liu, and Chen-Hsuan Lin. Neuralangelo: High-fidelity neural surface reconstruction. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 1, 2, 6, 7, 8
- [30] Chen-Hsuan Lin, Wei-Chiu Ma, Antonio Torralba, and Simon Lucey. Barf: Bundle-adjusting neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5741–5751, 2021. 1, 3, 6, 7, 8
- [31] Philipp Lindenberger, Paul-Edouard Sarlin, Viktor Larsson, and Marc Pollefeys. Pixel-Perfect Structure-from-Motion with Featuremetric Refinement. In *ICCV*, 2021. 3
- [32] Philipp Lindenberger, Paul-Edouard Sarlin, and Marc Pollefeys. Lightglue: Local feature matching at light speed. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 17627–17638, 2023. 3
- [33] Jianlin Liu, Qiang Nie, Yong Liu, and Chengjie Wang. Nerfloc: Visual localization with conditional neural radiance field. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9385–9392. IEEE, 2023. 3
- [34] Shaohui Liu, Yifan Yu, Rémi Pautrat, Marc Pollefeys, and Viktor Larsson. 3d line mapping revisited. In *Computer Vision and Pattern Recognition (CVPR)*, 2023. 3
- [35] David G Lowe. Distinctive image features from scale-invariant keypoints. *International journal of computer vision*, 60:91–110, 2004. 3
- [36] Dominic Maggio, Marcus Abate, Jingnan Shi, Courtney Mario, and Luca Carlone. Loc-nerf: Monte carlo localization using neural radiance fields. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 4018–4025. IEEE, 2023. 3, 5, 1
- [37] Paul Merrell, Amir Akbarzadeh, Liang Wang, Philippos Mordohai, Jan-Michael Frahm, Ruigang Yang, David Nistér, and Marc Pollefeys. Real-time visibility-based fusion of depth maps. In *2007 IEEE 11th International Conference on Computer Vision*, pages 1–8. Ieee, 2007. 2
- [38] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020. 2, 3, 4
- [39] 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. *Communications of the ACM*, 65(1):99–106, 2021. 1
- [40] Bailey Miller, Hanyu Chen, Alice Lai, and Ioannis Gkioulekas. Objects as volumes: A stochastic geometry view of opaque solids. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 87–97, 2024. 1
- [41] Arthur Moreau, Nathan Piasco, Moussab Bennehar, Dzmitry Tsishkou, Bogdan Stanculescu, and Arnaud de La Fortelle. Crossfire: Camera relocalization on self-supervised features from an implicit representation. *arXiv preprint arXiv:2303.04869*, 2023. 3
- [42] Denis Oberkamp, Daniel F DeMenthon, and Larry S Davis. Iterative pose estimation using coplanar feature points. *Computer Vision and Image Understanding*, 63(3):495–511, 1996. 5
- [43] Kemal E Ozden, Konrad Schindler, and Luc Van Gool. Multibody structure-from-motion in practice. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(6):1134–1141, 2010. 5
- [44] Linfei Pan, Daniel Barath, Marc Pollefeys, and Johannes Lutz Schönberger. Global Structure-from-Motion Revisited. In *European Conference on Computer Vision (ECCV)*, 2024. 3
- [45] Keunhong Park, Philipp Henzler, Ben Mildenhall, Jonathan T Barron, and Ricardo Martin-Brualla. Camp: Camera preconditioning for neural radiance fields. *ACM Transactions on Graphics (TOG)*, 42(6):1–11, 2023. 3
- [46] Jerome Revaud, Cesar De Souza, Martin Humenberger, and Philippe Weinzaepfel. R2d2: Reliable and repeatable detector and descriptor. *Advances in neural information processing systems*, 32, 2019. 3
- [47] Paul-Edouard Sarlin, Cesar Cadena, Roland Siegwart, and Marcin Dymczyk. From coarse to fine: Robust hierarchical localization at large scale. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12716–12725, 2019. 3, 6
- [48] Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. SuperGlue: Learning feature matching with graph neural networks. In *CVPR*, 2020. 3, 6
- [49] Torsten Sattler, Bastian Leibe, and Leif Kobbelt. Efficient & effective prioritized matching for large-scale image-based localization. *IEEE transactions on pattern analysis and machine intelligence*, 39(9), 2016. 3
- [50] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 1, 2, 3
- [51] 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. 1, 2
- [52] Gerald Schweighofer and Axel Pinz. Robust pose estimation from a planar target. *IEEE transactions on pattern analysis and machine intelligence*, 28(12):2024–2030, 2006. 5

- [53] Noah Snavely, Steven M Seitz, and Richard Szeliski. Modeling the world from internet photo collections. *International journal of computer vision*, 80:189–210, 2008. 3
- [54] Jiaming Sun, Zehong Shen, Yuang Wang, Hujun Bao, and Xiaowei Zhou. Loftr: Detector-free local feature matching with transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8922–8931, 2021. 3
- [55] Chris Sweeney. Theia multiview geometry library: Tutorial & reference. <http://theia-sfm.org>. 3
- [56] Zachary Teed and Jia Deng. Tangent space backpropagation for 3d transformation groups. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10338–10347, 2021. 6
- [57] Prune Truong, Marie-Julie Rakotosaona, Fabian Manhardt, and Federico Tombari. Sparf: Neural radiance fields from sparse and noisy poses. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4190–4200, 2023. 1, 3
- [58] Michał Tyszkiewicz, Pascal Fua, and Eduard Trulls. Disk: Learning local features with policy gradient. *Advances in Neural Information Processing Systems*, 33:14254–14265, 2020. 3
- [59] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *NeurIPS*, 2021. 1, 2, 3, 4, 6, 7, 8
- [60] Yiming Wang, Qin Han, Marc Habermann, Kostas Daniilidis, Christian Theobalt, and Lingjie Liu. Neus2: Fast learning of neural implicit surfaces for multi-view reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3295–3306, 2023. 2
- [61] 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. 3
- [62] 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. 1, 3, 4, 5
- [63] Changchang Wu. Visualsfm: A visual structure from motion system. <http://www.cs.washington.edu/homes/ccwu/vsfm>, 2011. 3, 6
- [64] Tong Wu, Jiaqi Wang, Xingang Pan, Xudong Xu, Christian Theobalt, Ziwei Liu, and Dahua Lin. Voxurf: Voxel-based efficient and accurate neural surface reconstruction. *arXiv preprint arXiv:2208.12697*, 2022. 1, 2
- [65] Yitong Xia, Hao Tang, Radu Timofte, and Luc Van Gool. Sinerf: Sinusoidal neural radiance fields for joint pose estimation and scene reconstruction. *arXiv preprint arXiv:2210.04553*, 2022. 3
- [66] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. In *Proceedings of the European conference on computer vision (ECCV)*, pages 767–783, 2018. 2
- [67] Yao Yao, Zixin Luo, Shiwei Li, Tianwei Shen, Tian Fang, and Long Quan. Recurrent mvsnet for high-resolution multi-view stereo depth inference. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5525–5534, 2019. 2
- [68] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Basri Ronen, and Yaron Lipman. Multiview neural surface reconstruction by disentangling geometry and appearance. *Advances in Neural Information Processing Systems*, 33, 2020. 3
- [69] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces. *Advances in Neural Information Processing Systems*, 34:4805–4815, 2021. 2
- [70] Lin Yen-Chen, Pete Florence, Jonathan T Barron, Alberto Rodriguez, Phillip Isola, and Tsung-Yi Lin. inerf: Inverting neural radiance fields for pose estimation. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1323–1330. IEEE, 2021. 3
- [71] Zehao Yu, Songyou Peng, Michael Niemeyer, Torsten Sattler, and Andreas Geiger. Monosdf: Exploring monocular geometric cues for neural implicit surface reconstruction. *Advances in neural information processing systems*, 35:25018–25032, 2022. 2
- [72] Christopher Zach, Thomas Pock, and Horst Bischof. A globally optimal algorithm for robust tv-l1 range image integration. In *2007 IEEE 11th International Conference on Computer Vision*, pages 1–8. IEEE, 2007. 2
- [73] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv:2010.07492*, 2020. 2, 4
- [74] Bingfan Zhu, Yanchao Yang, Xulong Wang, Youyi Zheng, and Leonidas Guibas. Vdn-nerf: Resolving shape-radiance ambiguity via view-dependence normalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 35–45, 2023. 2, 4

# Robust SG-NeRF: Robust Scene Graph Aided Neural Surface Reconstruction

## Supplementary Material

### 5.1. More visual results

The interactive visual comparisons are presented in <https://rsg-nerf.github.io/RSG-NeRF/>.

Fig. 8 and 9 illustrate more results of the scene graph updating. We present 4 cases in Bear and Baby scenes, including all rejected, all accepted, more inlier matching and more outlier matching in each.

Fig. 10 illustrate the comparisons of SG-NeRF [6] and our method in 8 scenes of SG-NeRF dataset, including meshes and poses.

Fig. 11 to 14 present the detailed mesh comparisons.

### 5.2. Details about Monte Carlo re-localization

We follow [13, 36] to implement our Monte Carlo re-localization. We fix the NeuS [59] backbone first and utilize  $\mathcal{L}_{color}(C_n)$  of optimize pose parameters. To make the gradients backpropagate to pose parameters, we remove the detach flag in the re-localization process. In the first stage, each particle is sampled with the same probability. We maintain a  $PSNR_n$  list for each particle and approximate the current  $PSNR_n$  by the mean value of last 10 elements in the list.

Subsequently, we compute the distribution of the particles according to the  $PSNR_n$  in current state, abbreviated as  $P_n$ . For a specific particle  $p_i$ , we define the sampling probability of  $p_i$  by:

$$\mathcal{P}(p_i) = \frac{e^{(P_n(p_i) - \min P_n)}}{\sum_{p \in \{p_1, p_2, \dots, p_{N_p}\}} e^{(P_n(p) - \min P_n)}}. \quad (13)$$

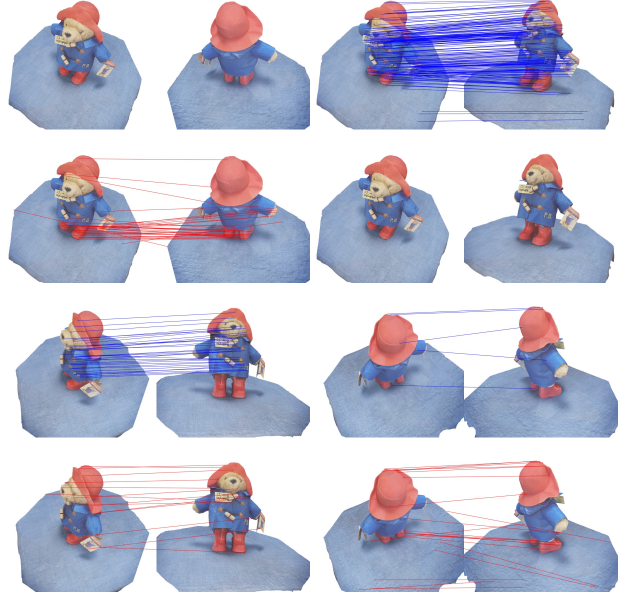


Figure 8. Scene graph updating on Bear.

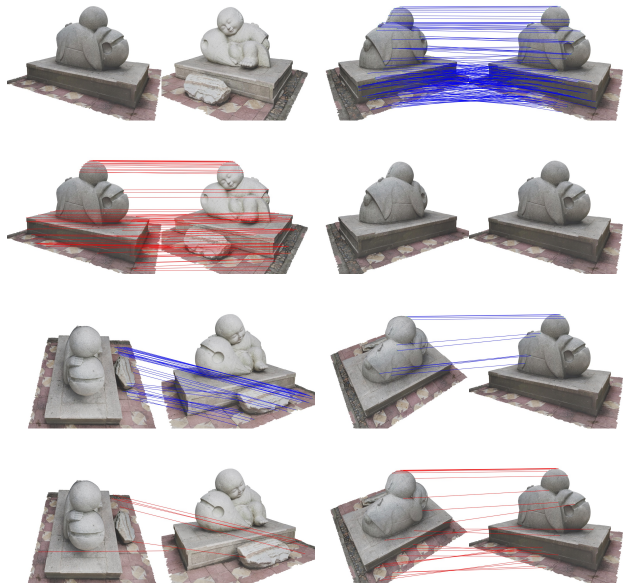


Figure 9. Scene graph updating on Baby.



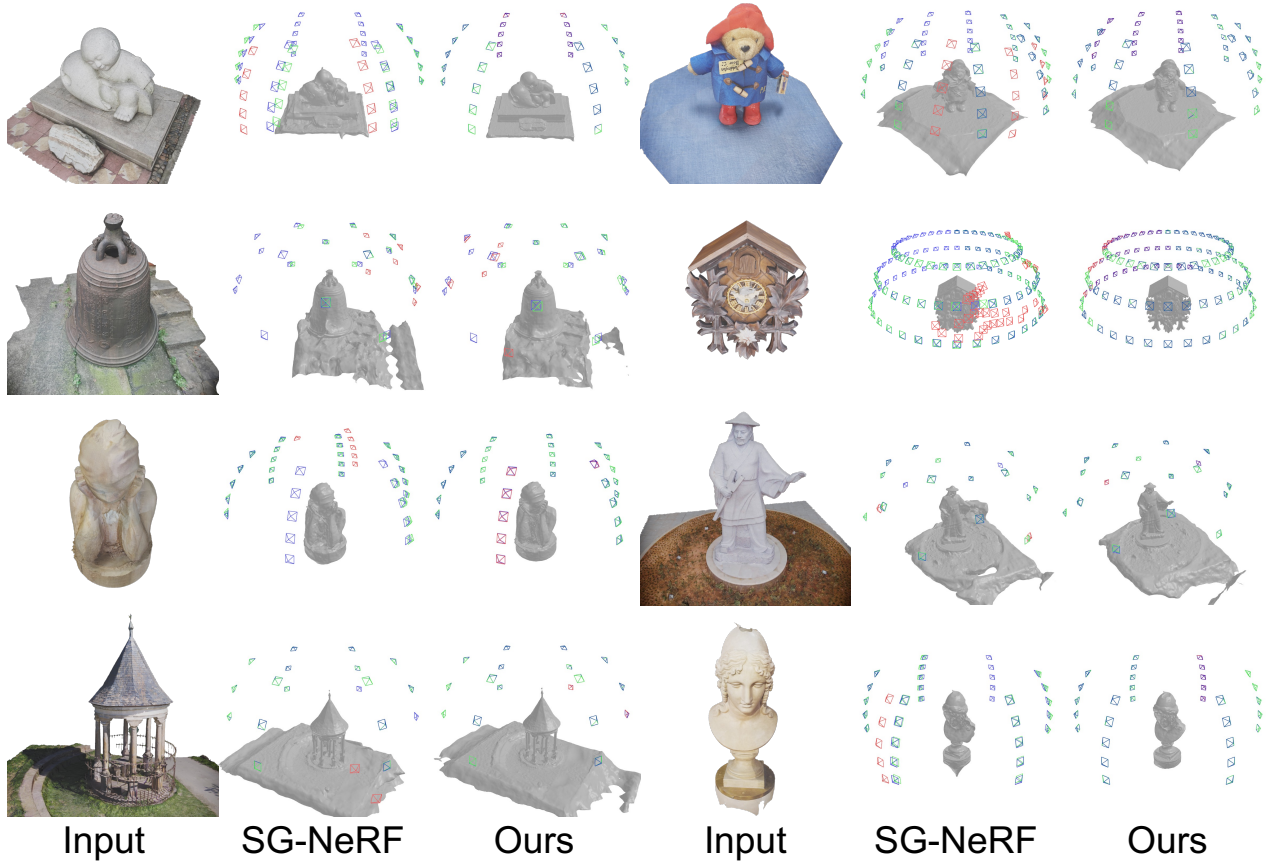


Figure 10. Reconstruction results on the SG-NeRF [6] dataset. Both SG-NeRF [6] and our method take the same initial poses as input, including significant noises. The camera poses are also presented with optimized outlier poses, inlier poses and ground truth poses.

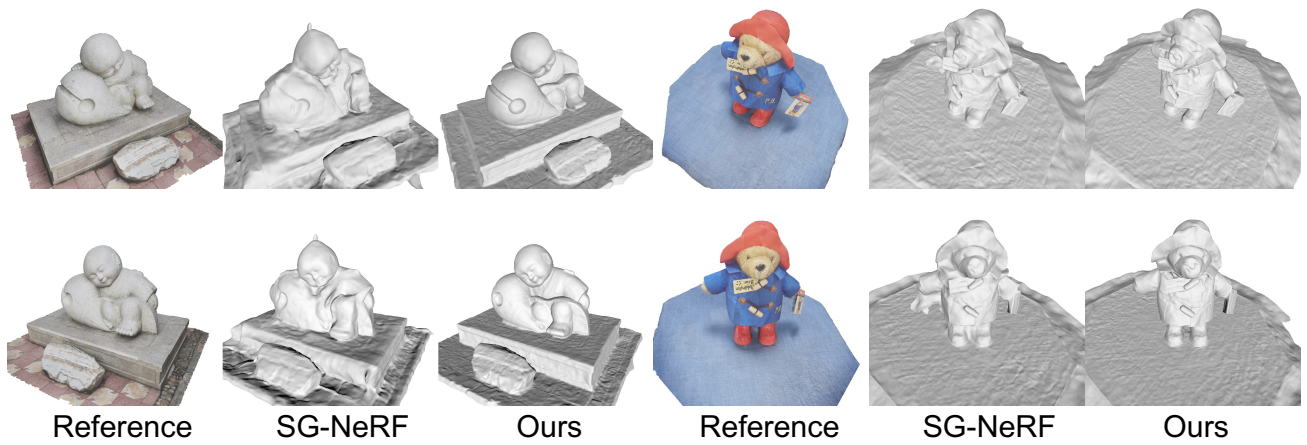


Figure 11. Reconstruction results on the SG-NeRF [6] dataset (Baby, Bear).

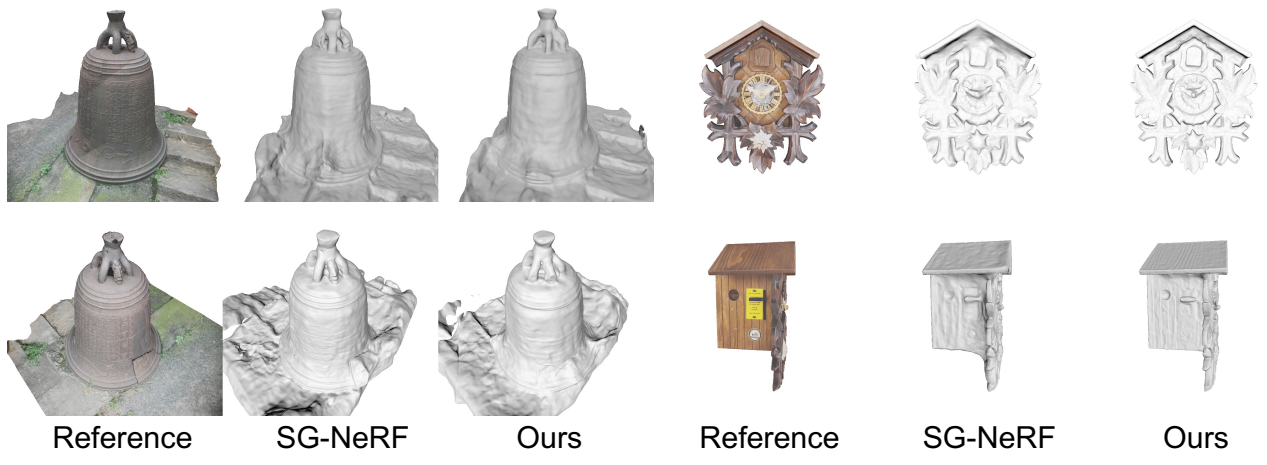


Figure 12. Reconstruction results on the SG-NeRF [6] dataset (Bell, Clock).



Figure 13. Reconstruction results on the SG-NeRF [6] dataset (Deaf, Farmer).

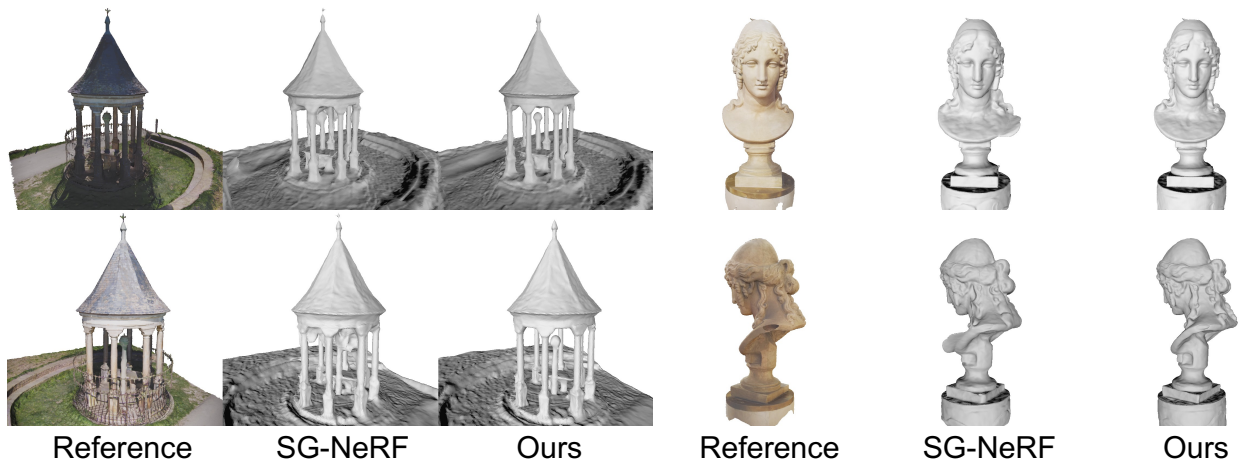


Figure 14. Reconstruction results on the SG-NeRF [6] dataset (Pavilion, Sculpture).