# NerfingMVS: Guided Optimization of Neural Radiance Fields for Indoor Multi-view Stereo

Yi Wei[1,2], Shaohui Liu[3], Yongming Rao[1,2], Wang Zhao[4], Jiwen Lu[1,2,*] Jie Zhou[1,2]

[1]Department of Automation, Tsinghua University, China
[2]State Key Lab of Intelligent Technologies and Systems, China
[3]ETH Zurich    [4]Department of Computer Science and Technology, Tsinghua University, China

`y-wei19@mails.tsinghua.edu.cn; b1ueber2y@gmail.com; raoyongming95@gmail.com;`
`zhao-w19@mails.tsinghua.edu.cn; {lujiwen, jzhou}@tsinghua.edu.cn`
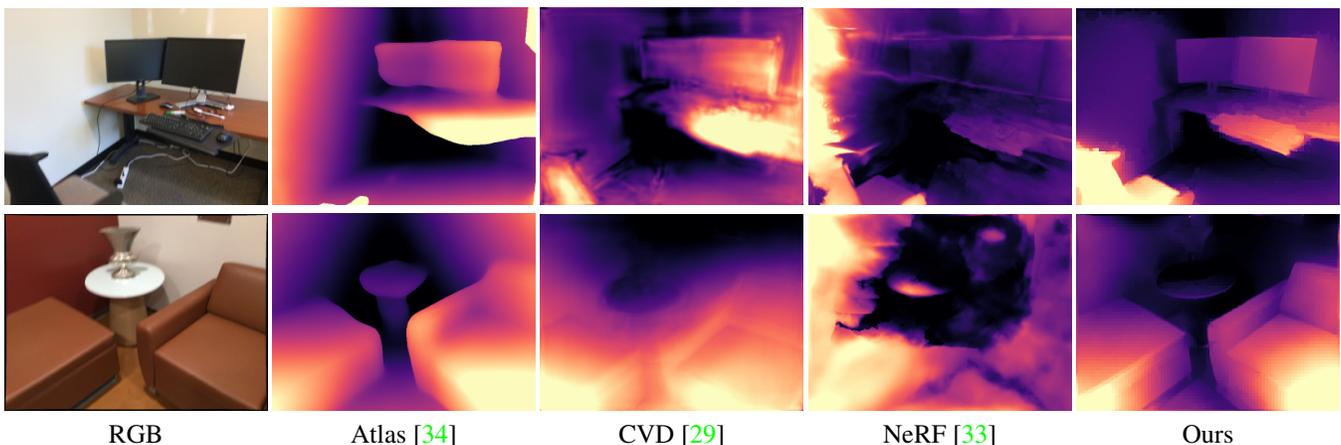
| RGB | Atlas [34] | CVD [29] | NeRF [33] | Ours |

Figure 1: Qualitative results for multi-view depth estimation on ScanNet [4]. Our method clearly surpasses leading multi-view estimation methods [29, 34] by building on top of neural radiance fields [33]. While also using test-time optimization, CVD [29] suffers from inaccurate estimation of flow correspondences. NeRF [33] fails to produce accurate geometry due to the inherent shape-radiance ambiguity [61] (See Figure 3) in indoor scenes. With guided optimization, our method successfully integrates the learning-based depth priors into NeRF, significantly improving the geometry of the radiance fields.

## Abstract

*In this work, we present a new multi-view depth estimation method that utilizes both conventional reconstruction and learning-based priors over the recently proposed neural radiance fields (NeRF). Unlike existing neural network based optimization method that relies on estimated correspondences, our method directly optimizes over implicit volumes, eliminating the challenging step of matching pixels in indoor scenes. The key to our approach is to utilize the learning-based priors to guide the optimization process of NeRF. Our system firstly adapts a monocular depth network over the target scene by finetuning on its sparse SfM+MVS reconstruction from COLMAP. Then, we show that the shape-radiance ambiguity of NeRF still exists in indoor environments and propose to address the issue by employing the adapted depth priors to monitor the sampling process of volume rendering. Finally, a per-pixel confidence map acquired by error computation on the rendered image can be used to further improve the depth quality. Experiments show that our proposed framework significantly outperforms state-of-the-art methods on indoor scenes, with surprising findings presented on the effectiveness of correspondence-based optimization and NeRF-based optimization over the adapted depth priors. In addition, we show that the guided optimization scheme does not sacrifice the original synthesis capability of neural radiance fields, improving the rendering quality on both seen and novel views. Code is available at https://github.com/weiyithu/NerfingMVS.*

## 1. Introduction

Reconstructing 3D scenes from multi-view posed images, also named as multi-view stereo (MVS), has been a

fundamental topic in computer vision over decades. The application varies from robotics, 3D modeling, to virtual reality, etc. Conventional multi-view stereo approaches [2, 9, 13, 60] densely match pixels across views by comparing the similarity of cross-view image patches. While producing impressive results, those methods often suffer from poorly textured regions, thin structures and non-Lambertian surfaces, especially in real-world indoor environments.

Recently, with the success of deep neural networks, several learning-based methods [17, 20, 24, 53] are proposed to tackle the multi-view stereo problem often by employing a cost volume based architecture. Those methods perform a direct neural network inference at test time for multi-view depth estimation and achieve remarkable performance on benchmarks. However, due to the lack of constraints at inference, the predicted depth maps across views are often not consistent and the photometric consistency is often violated. To address this issue, [29] proposed a test-time optimization framework that optimizes over learning-based priors acquired from single-image depth estimation. While being computationally inefficient, the method produces accurate and consistent depth maps that are available for various visual effects. However, the optimization formulation of this method relies heavily on an optical flow network [16] to establish correspondences, which becomes problematic when estimated correspondences are unreliable.

In this paper, we present a new neural network based optimization framework for multi-view depth estimation based on the recently proposed neural radiance fields [33]. Instead of relying on estimated correspondences and cross-view depth reprojection for optimization [29], our method directly optimizes over volumes. However, we show that the shape-radiance ambiguity [61] of NeRF becomes the bottleneck on estimating accurate per-view depths in indoor scenes. To address the issue, we propose a guided optimization scheme to help train NeRF with learning-based depth priors. Specifically, our system firstly adapts a monocular depth network onto the test scene by finetuning on its conventional SfM+MVS reconstruction. Then, we employ the adapted depth priors to guide the sampling process of volume rendering for NeRF. Finally, we acquire a confidence map from the rendered RGB image of NeRF and improve the depth map with a post-filtering step.

Our findings indicate that the scene-specific depth prior adaptation significantly improves the depth quality. However, performing existing correspondence-based optimization on the adapted depth priors will surprisingly degrade the performance. On the contrary, with direct optimization over neural radiance fields, our method consistently improves the depth quality over adapted depth priors. This phenomenon demonstrates the potential of exploiting neural radiance fields for accurate depth estimation.

Experiments show that our proposed framework significantly improves upon state-of-the-art multi-view depth estimation methods on tested indoor scenes. In addition, the guided optimization from learning-based priors can help improve the rendering quality of NeRF on both seen and novel views, achieving comparable or better quality with state-of-the-art novel view synthesis methods. This indicates that conventional non-learning reconstruction method, while demonstrated to be effective on helping image-based view synthesis in [39, 40], can also help improve the synthesis quality on neural implicit representations.

## 2. Related Work

**Multi-view Reconstruction:** Recently, 3D vision [14, 48, 54–56, 62] has attracted more and more attention. Early multi-view reconstruction approaches include volumetric optimization [7, 21, 52], which perform global optimization with photo-consistency based assumptions. However, those methods suffer from large computational complexity. Another direction [2, 9] is to estimate per-view depth map. Compared to volumetric approaches, these methods can produce finer geometry. However, they rely on accurately matched pixels by comparing the similarity of cross-view patches at different depth hypotheses, which will be problematic over poorly textured regions in indoor scenes. Recently, a number of learning-based methods are proposed. While some of them predict on voxelized grids [19, 49], they suffer from limited resolution. An exception of this is Atlas [34], which predicts TSDF values via back-projection of the image features. Most learning-based methods [15, 17, 20, 24, 28, 53] follow the spirit of conventional approaches [9] and generate per-view depth map from a cost volume based architecture. Most related to us, [29] performs test-time optimization over per-view depth maps with learning-based priors. While our work also utilizes the learning-based priors, we build on top of the recently proposed neural radiance fields [30] and introduce a new way to accurately estimate multi-view depths by directly optimizing over implicit volumes with the guide of learning-based priors. Our method neither suffers from the resolution problem nor relies on accurately estimated correspondences.

**Neural Implicit Representation:** Recently, several seminal works [3, 31, 36] demonstrate the potential of representing implicit surfaces with a neural network, which enables memory-efficient geometric representation with infinite resolution. Variations include applying neural implicit representations on part hierarchies [11, 18], human reconstruction [41, 42], view synthesis [27, 46], differentiable rendering [26, 35], etc. Neural radiance fields (NeRF) [33] represent scenes as continuous implicit function of positions and orientations for high quality view synthesis, which leads to several follow-up works [1, 38, 61] improving its performance. There are several extensions for NeRF including
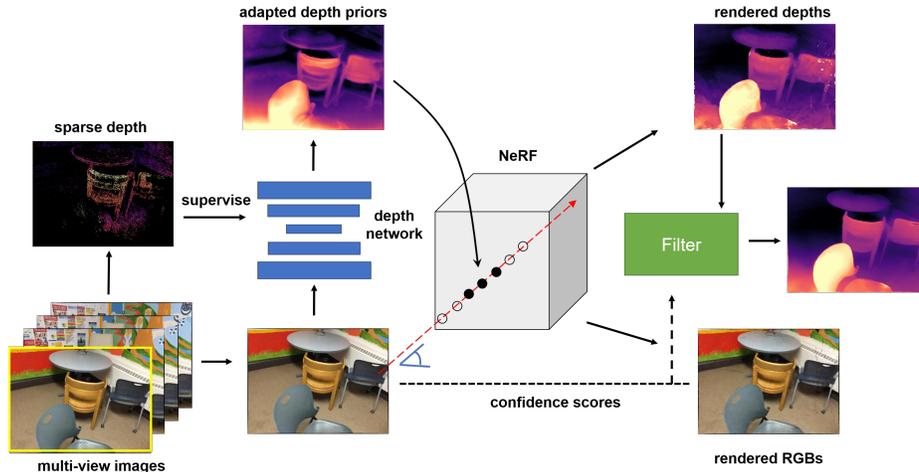
Figure 2: An overview of our method. We first adopt conventional SfM and MVS from COLMAP to get sparse depth (after fusion), which is used to train a monocular depth network to get scene-specific depth priors. Then, we utilize the depth priors to guide volume sampling in the optimization of NeRF [33]. Finally, by computing the errors between the rendered images and the original input images we acquire confidence scores, which enables us to employ a confidence-based filter to improve the rendered depths.

dynamic scenes [37,57], portrait avatars [10], relighting [1], pose estimation [59], etc. In this paper, we propose a guided optimization scheme to enrich NeRF [33] with the ability of accurate depth estimation which surpasses leading multi-view depth estimation approaches.

**View Synthesis:** View synthesis is conventionally often referred as view interpolation [12, 22], where the goal is to interpolate views within the convex hull of the initial camera positions. With the success of deep learning, learning-based methods [8, 32, 47, 64] have been proposed to address the problem and have achieved remarkable improvements. Recently, neural radiance fields [30] demonstrates impressive results of view synthesis by representing scenes as continuous implicit radiance fields. It is further extended to operate on dynamic scenes [37, 57]. [25] employs a sparse voxel octree and achieves great improvement over [33]. [39] employs image-based encoder-decoder architecture to process the proxy generated from the conventional sparse reconstruction, and is later improved by [40]. While view synthesis is not the major focus of this work, we show that our guided optimization scheme consistently improves the synthesis quality of NeRF [33] on both seen and novel views, which shows the potential of using conventional sparse reconstructions to help improve the synthesis quality of NeRF-like methods.

## 3. Approach

### 3.1. Overview

We introduce a multi-view depth estimation method that utilizes conventional sparse reconstruction and learning-

based priors. Our proposed system builds on top of the recently proposed neural radiance fields (NeRF) [33] and performs test-time optimization at inference. Compared to the existing test-time optimization method [29] that relies on estimated correspondences, directly optimizing over volumes eliminates the necessity of accurately matching cross-view pixels. This idea is also exploited by direct methods in the context of simultaneous localization and mapping (SLAM) [6].

The key to our approach is to effectively integrate the additional information from the learning-based priors into the NeRF training pipeline. Figure 2 shows an overview of our proposed system. Section 3.2 shows how we adapt the depth priors to specific scenes at test time. In Section 3.3, we analyze the reason why NeRF fails on producing accurate geometry in indoor scenes and describe our learning-based priors guided optimization scheme. In Section 3.4, we discuss how to infer depth and synthesize views from the neural radiance fields trained with guided optimization.

### 3.2. Scene-specific Adaptation of the Depth Priors

Similar to CVD [29], our method also aims to utilize learning-based depth priors to help optimize the geometry at test time. However, unlike [29] that employs the same monocular depth network for all test scenes, we propose to adapt the network onto each scene to get scene-specific depth priors. Empirically this test-time adaptation method largely improves the quality of the final depth output.

Our proposal on adapting scene-specific depth priors is to finetune a monocular depth network over its conventional sparse reconstruction. Specifically, we run COLMAP [43, 44] on the test scene and acquire per-view sparse depth

maps by projecting the fused 3D point clouds after multi-view stereo. Since geometric consistency check is adopted in the fusion step, the acquired depth map is sparse but robust and can be used as a supervision source for training the scene-specific depth priors.

Due to the scale ambiguity of acquired depth map, we employ the scale-invariant loss [5] to train the depth network, which is written as follows:

$$L(D_p^i, D_{Sparse}^i) = \frac{1}{n} \sum_{j=1}^{n} |\log D_p^i(j) - \log D_{Sparse}^i(j)$$
$$+ \alpha(D_p^i, D_{Sparse}^i)|,$$

(1)

where $D_p^i$ is the predicted depth map and $D_{Sparse}^i$ is the sparse depths acquired from COLMAP [43, 44]. We align the scale of the predicted depth map with the sparse depth supervision by employing the scale factor $\alpha(D_p^i, D_{Sparse}^i)$ in the loss formulation, which can be computed by averaging the difference over all valid pixels:

$$\alpha(D_p^i, D_{Sparse}^i) = \frac{1}{n} \sum_{j} (\log D_p^i(j) - \log D_{Sparse}^i(j)).$$

(2)

The finetuned monocular depth network is a stronger prior that fits the specific target scene. The quality of the adapted priors can be further improved with our guided optimization over NeRF, while Table 2 shows that applying existing correspondence-based neural optimization will surprisingly degrade the quality of the adapted depth priors.

### 3.3. Guided Optimization of NeRF

Neural radiance fields were initially proposed in [33], which achieves impressive results on view synthesis. Our system exploits its potential for accurate depth estimation. By integrating the aforementioned adapted depth priors, we directly optimize on implicit volumes. The key to the success of NeRF is to employ a fully connected network parameterized by $\theta$ to represent implicit radiance fields with $F_\theta : (\mathbf{x}, \mathbf{d}) \rightarrow (\mathbf{c}, \sigma)$, where $\mathbf{x}$ and $\mathbf{d}$ denotes the location and direction, $\mathbf{c}$ and $\sigma$ denotes the color and density as the network outputs. View synthesis can be easily achieved over NeRF with volume rendering, which enables NeRF to train itself directly over multi-view RGB images. During volume rendering, NeRF adopts the near bound $t_n$ and the far bound $t_f$ computed from the sparse 3D reconstruction to monitor the sampling space along each ray. Specifically, it partitions $[t_n, t_f]$ into $M$ bins and one query point is randomly sampled for each bin with a uniform distribution:

$$t_i \sim \mathcal{U}\left[t_n + \frac{i-1}{M}(t_f - t_n), \ t_n + \frac{i}{M}(t_f - t_n)\right]. \quad (3)$$

The rendered RGB value $C(\mathbf{r})$ for each ray can be calculated from the finite samples with volume rendering. More-



(a) rendered RGB      (b) sampled points

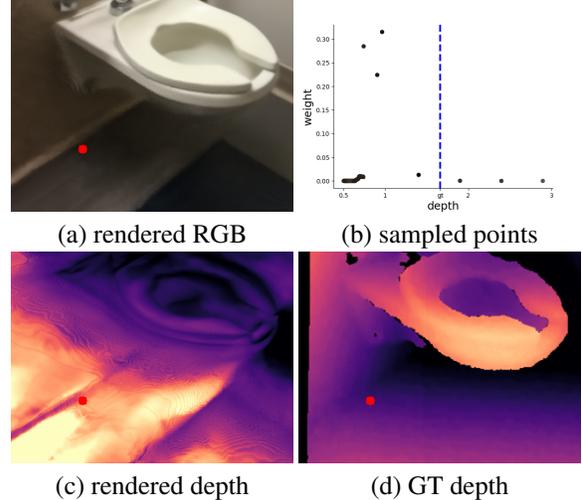(c) rendered depth      (d) GT depth

Figure 3: The inherent shape-radiance ambiguity [61] becomes a bottleneck in indoor scenes. **Top row:** (a) rendered RGB of NeRF [33]. (b) visualization of the sampled points along the camera ray at the position colored in red. The blue line indicates the groundtruth depth value. **Bottom row:** (c) the rendered depth map of NeRF [33]. (d) the groundtruth depth map. While NeRF produces high quality rendered image (PSNR: 31.53), the rendered depth largely deviates from the groundtruth.

over, per-view depth $D(\mathbf{r})$ can also be approximated by calculating the expectation of the samples along the ray:

$$C(\mathbf{r}) = \sum_{i=1}^{M} T_i(1 - \exp(-\sigma_i \delta_i)) c_i$$
$$D(\mathbf{r}) = \sum_{i=1}^{M} T_i(1 - \exp(-\sigma_i \delta_i)) t_i$$

(4)

where $T_i = \exp\left(-\sum_{j=1}^{i-1} \sigma_j \delta_j\right)$ indicates the accumulated transmittance from $t_n$ to $t_i$ and $\delta_i = t_{i+1} - t_i$ is the distance between adjacent samples.

While simply satisfying the radiance field over the input image does not guarantee a correct geometry, the shape-radiance ambiguity between the 3D geometry and radiance has been studied in [61]. It is believed in the paper that because incorrect geometry leads to high intrinsic complexity, the correct shape, with smoother surface light field, is more favored by the learned neural radiance fields with limited network capacity. This assumption generally holds for rich textured outdoor scenes. However, we empirically observe that NeRF struggles on poorly textured areas (e.g. walls), which are common in indoor environments. Figure 3 shows one failure case of NeRF that suffers from shape-radiance ambiguity in texture-less areas, where NeRF perfectly synthesizes the input image with a geometry largely deviated from the groundtruth. The failure comes from the fact that

while extremely implausible shapes are ignored with the favor of smoothed surface light field [61], there still exists a family of smoothed radiance fields that perfectly explains the training images. Further, the blurred images and large-motion real-world indoor scenes will reduce the capacity of NeRF and aggravate the shape-radiance ambiguity issue. We find that this is a common issue in all tested indoor scenes.

In Figure 3(b), we show that all the sampled points along the camera ray that corresponds to a poorly textured pixel predict roughly the same RGB values, with the confidence distribution concentrated only in a limited range. Motivated by this observation, we consider guiding the NeRF sampling process with our adapted depth priors from the monocular depth network. By explicitly limiting the sampling range to be distributed around the depth priors, we avoid most degenerate cases for NeRF in indoor scenes. This enables accurate depth estimation by directly optimizing over RGB images.

Specifically, we first acquire error maps of the adapted depth priors with a geometric consistency check. Denote the adapted depth priors as $\{D^i\}_{i=1}^N$ for the N input views. We project the depth map of each view to all the other views:

$$p^{i \to j}, D^{i \to j} = proj(K, T^{i \to j}, D^i)$$
$$D^{j'} = D^j(p^{i \to j}), \tag{5}$$

where $K$ is the camera intrinsics, $T^{i \to j}$ is the relative pose. $p^{s \to t}$ and $D^{i \to j}$ are the 2D coordinates and depth of the projection in $j$th view. Then we calculate the depth reprojection error using the relative error between $D^{j'}$ and $D^{i \to j}$. Note that there are pixels that do not overlap across some view pairs. Thus, we define the error map of the depth priors for each view $e_i$ as the average value of the top $K$ minimum cross-view depth projection error.

We use the error maps $\{e^i\}_{i=1}^N$ to calculate adaptive sample ranges $[t_n, t_f]$ for each camera ray:

$$t_n = D(1 - clamp(e, \alpha_l, \alpha_h))$$
$$t_f = D(1 + clamp(e, \alpha_l, \alpha_h)) \tag{6}$$

where $\alpha_l$ and $\alpha_h$ defines the relative lower and higher bounds of the ranges. With the adaptive ranges we achieve a balance between diversity and precision of the confidence distribution along camera rays. As illustrated in Figure 4, the sampling over pixels with relatively low error is more concentrated around the adapted depth priors, while the sampling over pixels with large error is close to the original NeRF formulation.

### 3.4. Inference and View Synthesis

For inference, we can directly predict the depth map for each input view by resampling within the sampling range
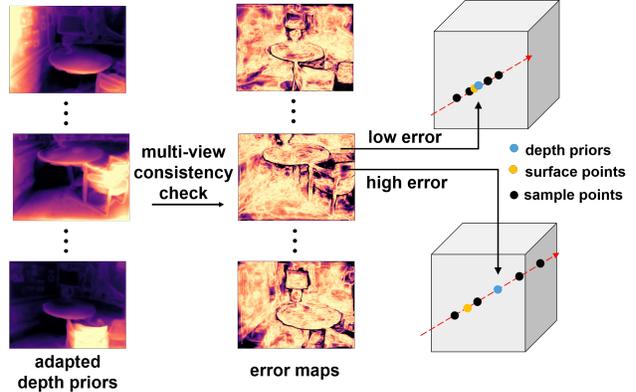


Figure 4: Guided optimization of NeRF [33]. We adopt multi-view consistency check on adapted depth priors to get error maps, which help calculate adaptive depth ranges for each camera ray to sample points for NeRF optimization.

defined in Eq. (6) and applying Eq. (3) to compute the expectation. This gives an accurate output depth for the NeRF equipped with our proposed guided optimization scheme.

To further improve depth quality, we exploit the potential of using the view synthesis results of NeRF to compute per-pixel confidence for the predicted geometry. If the rendered RGB at a specific pixel does not match the input training image well, we attach a relatively low confidence for the depth prediction of this pixel. The confidence $S^i_j$ for the $j$th pixel in the $i$th view is specifically defined as:

$$S^i_j = 1 - \frac{1}{3}||C^i_{gt}(j) - C^i_{render}(j)||_1, \tag{7}$$

where $C^i_{gt}$ and $C^i_{render}$ are the groundtruth images and rendered images for each seen view with all the values divided by 255. The absolute difference is employed. This confidence map can be further used to refine the predicted depth map with off-the-shelf post-filtering techniques. We employ plane bilateral filtering introduced in [51] over the depth to get the final output, which improves depth quality especially for the regions where rendered RGB images are not accurate.

While the proposed guided optimization strategy needs the adapted depth priors as input to guide point sampling along the camera ray, we can still perform novel view synthesis by directly using the adapted depth priors from the nearest seen view. Empirically this is sufficient to produce accurate depth maps and significantly outperforms the original NeRF in terms of view synthesis quality (See Table 5).

## 4. Experiments

### 4.1. Experimental Setup

**Dataset:** We conducted experiments on ScanNet [4] dataset. Following the experimental setup in NeRF [33],

| Method | Abs Rel | Sq Rel | RMSE | RMSE log | $\delta < 1.25$ | $\delta < 1.25^2$ | $\delta < 1.25^3$ |
|---|---|---|---|---|---|---|---|
| COLMAP [43,44] | 0.4619 | 0.6308 | 1.0125 | 1.7345 | 0.4811 | 0.5139 | 0.5333 |
| ACMP [58] | 0.1945 | 0.1710 | 0.4551 | 0.3056 | 0.7309 | 0.8810 | 0.9419 |
| DELTAS [45] | 0.1001 | 0.0319 | 0.2070 | 0.1284 | 0.8618 | 0.9920 | 0.9991 |
| Atlas [34] | 0.0776 | 0.0631 | 0.2441 | 0.2693 | 0.9289 | 0.9536 | 0.9594 |
| DeepV2D [50] | 0.0818 | 0.0226 | 0.1714 | 0.1095 | 0.9414 | 0.9908 | 0.9979 |
| NeRF [33] | 0.3929 | 1.4849 | 1.0901 | 0.5210 | 0.4886 | 0.7318 | 0.8285 |
| Mannequin [23] | 0.1554 | 0.0636 | 0.2969 | 0.1806 | 0.7859 | 0.9735 | 0.9953 |
| CVD [29] | 0.0995 | 0.0304 | 0.1945 | 0.1269 | 0.9008 | 0.9879 | 0.9971 |
| Ours w/o filter | 0.0635 | 0.0145 | 0.1455 | 0.0936 | 0.9541 | 0.9910 | 0.9989 |
| Ours | **0.0614** | **0.0126** | **0.1345** | **0.0861** | **0.9601** | **0.9955** | **0.9996** |

Table 1: Quantitative comparisons for multi-view depth estimation. Scores are averaged over 8 scenes from ScanNet.

| Method | Abs Rel | Sq Rel | $\delta < 1.25$ |
|---|---|---|---|
| Mannequin [23] | 0.1554 | 0.0636 | 0.7859 |
| adapted depth priors | 0.0844 | 0.0223 | 0.9410 |
| CVD optimization | 0.0886 | 0.0251 | 0.9128 |
| Our optimization | **0.0635** | **0.0145** | **0.9541** |

Table 2: Comparison of the effectiveness of test-time optimization between CVD [29] and our method. Both methods perform optimization over the adapted depth priors, which is acquired by training the Mannequin Challenge depth network [23] with sparse supervision from COLMAP [43,44]. Scores are averaged over 8 scenes.

we randomly selected 8 scenes to evaluate our method. For each scene, we picked 40 images covering a local region and held out $1/8$ of these as the test set for novel view synthesis. All images are resized as $484 \times 648$ resolution. Due to the scale ambiguity issue, we adopted the median groundtruth scaling strategy [63] for depth evaluation.

**Implementation Details:** For the adapted depth priors, following CVD [29], we used the network architecture introduced in Mannequin Challenge [23] with its pretrained weights as our monocular depth network. 15 finetuning epochs were used in the scene-specific adaptation. We set $K = 4$ for multi-view consistency check and $\alpha_l = 0.05, \alpha_h = 0.15$ as the bounds of sample ranges. Please refer to our supplementary material for more details.

### 4.2. Results on Multi-view Depth Estimation

Table 1 shows the results for depth estimation task on ScanNet [4]. For all methods, we used their released implementation in the experiments. We also report results without applying the filtering step. Our method outperforms state-of-the-art depth estimation methods in all metrics. Note that DeepV2D [50], DELTAS [45] and Atlas [34] are all trained on ScanNet with groundtruth depth supervision. With the proposed guided optimization scheme, our method mitigates the problem of the shape-radiance ambi-



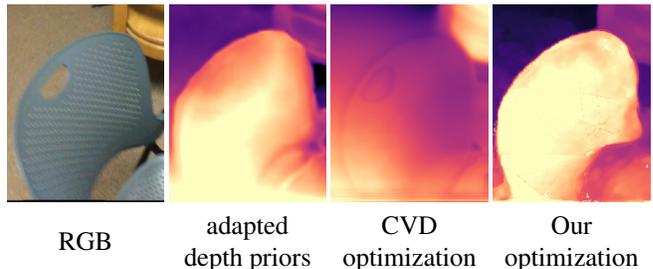| RGB | adapted depth priors | CVD optimization | Our optimization |
|---|---|---|---|

Figure 5: The optimization of CVD [29] surprisingly degrades the quality of the depth priors due to unreliable flow correspondences, while our method achieves improvement with guided optimization of NeRF [33]

| NeRF | depth priors | filter | Abs Rel | Sq Rel | $\delta < 1.25$ |
|---|---|---|---|---|---|
| ✓ | | | 0.302 | 0.210 | 0.518 |
| | ✓ | | 0.067 | 0.010 | 0.960 |
| ✓ | | ✓ | 0.287 | 0.167 | 0.546 |
| | ✓ | ✓ | 0.065 | 0.009 | 0.966 |
| ✓ | ✓ | | 0.053 | 0.006 | 0.979 |
| ✓ | ✓ | ✓ | **0.051** | **0.005** | **0.987** |

Table 3: Ablation studies on each component of our system. For the experiments 'NeRF + filter' and 'depth priors + filter', we compute the confidence scores by using the relative errors between the prediction depths and the groundtruth. The experiment was conducted on scene0521.

guity and demonstrates the potential of exploiting NeRF for accurate depth estimation. Figure 6 shows some qualitative results. While the original NeRF [33] fails to predict reasonable geometry, our method generates visually appealing depth maps. The confidence-based filter can further refine the predicted depth by smoothing the per-pixel estimation of NeRF [33].

To further study the advantages of optimizing over implicit volumes, we also applied the optimization of CVD [29] on our adapted depth priors. Results are shown in Table 2 and one example is exhibited in Figure 5. We surpris-

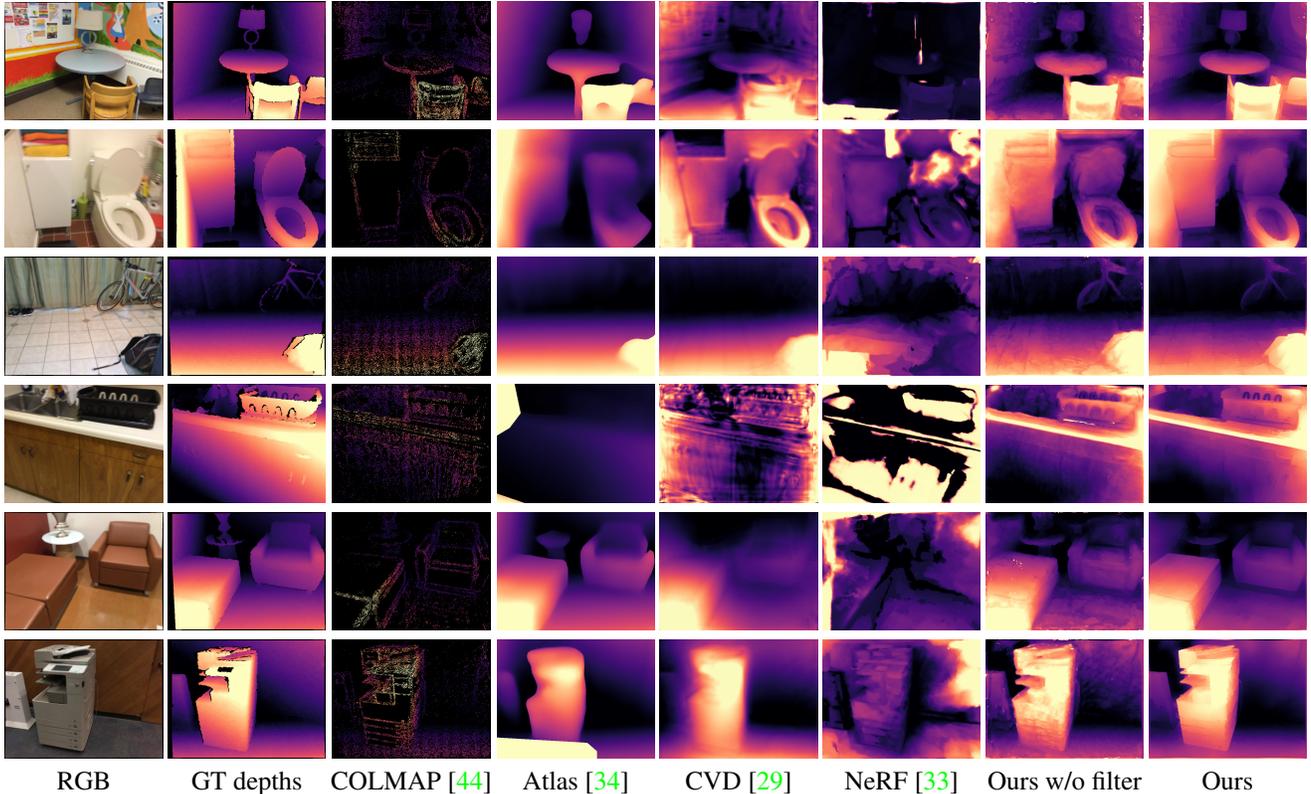| RGB | GT depths | COLMAP [44] | Atlas [34] | CVD [29] | NeRF [33] | Ours w/o filter | Ours |

Figure 6: Qualitative comparisons on ScanNet [4] dataset. Our method, without the post-filtering step, outperforms all compared methods in terms of depth quality. The filter further smooths the per-pixel estimated depth maps. **Better viewed when zoomed in.**

ingly find that the optimization of CVD degrades the depth quality of the initial depth priors. This is mainly due to wrong estimated correspondences from the employed flow network in [29]. The flow estimation is particularly challenging over poorly textured regions, which is ubiquitous in indoor scenes. The proposed guided optimization enables us to integrate depth priors on top of NeRF [33], which directly optimizes on raw RGB images, avoiding the challenging step of correspondence estimation in indoor scenes.

### 4.3. Ablation Studies

To better understand the working mechanism of our method, we performed ablation studies over each component of the proposed system. Results in Table 3 show that each component is beneficial to the final depth quality. This verifies the advantages of integrating depth priors into the optimization of NeRF [33].

We further study the design of adaptive ranges used in the guided optimization. It is shown that both the adaptive strategy and the use of bounds contribute to the performance gain. With the computed error maps, $\alpha_l$ and $\alpha_h$ avoid the samples being over-concentrated or overly random respectively, which enables the sampling to reach a balance between diversity and precision of the sampled points.

| adaptive | bound | Abs Rel | Sq Rel | $\delta < 1.25$ |
|:---:|:---:|:---:|:---:|:---:|
| | | 0.065 | 0.009 | 0.971 |
| ✓ | | 0.056 | 0.009 | 0.978 |
| ✓ | ✓ | **0.051** | **0.005** | **0.987** |

Table 4: Ablation studies on the design of the proposed guided optimization with adaptive ranges. 'bound' denotes the use of $\alpha_l$ and $\alpha_h$ in Eq. (6). For the experiment without using adaptive depth ranges for each camera ray, we set a fixed relative depth range to $[0.9, 1.1]$. The experiment was conducted on scene0521.

### 4.4. Results on View Synthesis

We also observe that the proposed guided optimization scheme is beneficial to the view synthesis quality of NeRF. Figure 7 illustrates some visualizations. Table 5 shows results on novel view synthesis, where our method consistently improves NeRF on all 8 scenes. Although view synthesis is not the main focus of our work, we achieve comparable or even better results compared to state-of-the-art novel view synthesis methods [1, 40]. Note that SVS [40] employs image-based novel view synthesis methods over the information extracted from sparse reconstruction. Our

| Method | scene 0616 | | scene 0521 | | scene 0000 | | scene 0158 | |
|--------|------|------|------|------|------|------|------|------|
| | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ |
| NSVF [25] | 15.71 | 0.704 | 27.73 | 0.892 | **23.36** | 0.823 | **31.98** | **0.951** |
| SVS [40] | **21.38** | **0.899** | **27.97** | **0.924** | 21.39 | **0.914** | 29.43 | **0.953** |
| NeRF [33] | 15.76 | 0.699 | 24.41 | 0.871 | 18.75 | 0.751 | 29.19 | 0.928 |
| Ours | 18.07 | 0.748 | **28.07** | 0.901 | 22.10 | 0.880 | 30.55 | **0.948** |
| Method | scene 0316 | | scene 0553 | | scene 0653 | | scene 0079 | |
| | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ |
| NSVF [25] | **22.29** | 0.917 | 31.15 | 0.947 | 28.95 | 0.929 | 26.88 | 0.887 |
| SVS [40] | 20.63 | **0.941** | 30.95 | **0.968** | 27.91 | **0.965** | 25.18 | **0.923** |
| NeRF [33] | 17.09 | 0.828 | 30.76 | 0.950 | 30.89 | 0.953 | 25.48 | 0.896 |
| Ours | 20.88 | 0.899 | **32.56** | **0.965** | **31.43** | **0.964** | **27.27** | **0.916** |

Table 5: Quantitative comparisons for novel view synthesis. Numbers in bold are within 1% of the best.

| Method | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
|--------|--------|--------|---------|
| NeRF [33] | 28.62 | 0.909 | 0.319 |
| Ours | **31.55** | **0.942** | **0.200** |

Table 6: Comparison between NeRF [33] and our method on seen views. Results are averaged over 8 scenes.

method, with the guided optimization scheme, opens a new way to employ the robust conventional sparse reconstruction to improve the synthesis quality directly over implicit 3D volumes. In addition, results in Table 6 show that our method can improve the view synthesis quality of NeRF on seen views. The guided optimization helps NeRF to focus on more informative regions and improves its capacity for rendering RGB images.

## 5. Conclusion and Future Work

In this work, we present a multi-view depth estimation method that integrates learning-based depth priors into the optimization of NeRF. Contrary to existing studies, we show that the shape-radiance ambiguity of NeRF becomes a bottleneck for NeRF-based depth estimation in indoor scenes. To address the issue, we propose a guided optimization framework to regularize the sampling process of NeRF during volume rendering with the adapted depth priors. Our proposed system demonstrates the significant improvement over prior works for indoor multi-view depth estimation, with a surprising finding that correspondence-based optimization can degrade the quality of depth priors in indoor scenes due to wrongly estimated flow correspondence. In addition, we also observe that the guided optimization improves the view synthesis quality of NeRF.

While our optimization is 3x faster than NeRF due to the advantages of guided optimization, the current method is still not efficient and thus hard to be scaled up to large datasets. Nonetheless, our work demonstrates the potential of using neural radiance fields for accurate depth estima-
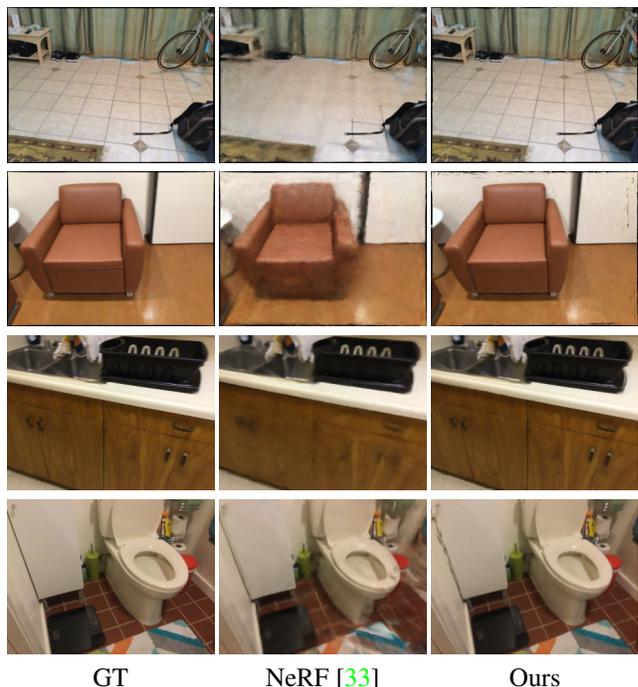


|  GT | NeRF [33] | Ours |

Figure 7: Results on view synthesis. The top two rows are rendering results on seen (training) views while the bottom two rows are on the novel views. With the adapted depth priors, our method improves the rendering quality for both seen and novel views. **Better viewed when zoomed in.**

tion. Future work includes efficient optimization, non-rigid reconstruction and visual effects based on the improved geometric structure in the learned neural radiance fields.

# References

[1] Sai Bi, Zexiang Xu, Pratul Srinivasan, Ben Mildenhall, Kalyan Sunkavalli, Miloš Hašan, Yannick Hold-Geoffroy, David Kriegman, and Ravi Ramamoorthi. Neural reflectance fields for appearance acquisition. *arXiv preprint arXiv:2008.03824*, 2020. 2, 3, 7

[2] Michael Bleyer, Christoph Rhemann, and Carsten Rother. PatchMatch Stereo-Stereo Matching with Slanted Support Windows., 2011. 2, 11

[3] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *CVPR*, pages 5939–5948, 2019. 2

[4] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *CVPR*, pages 5828–5839, 2017. 1, 5, 6, 7

[5] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. In *NeurIPS*, 2014. 4

[6] Jakob Engel, Thomas Schöps, and Daniel Cremers. Lsd-slam: Large-scale direct monocular slam. In *ECCV*, pages 834–849. Springer, 2014. 3

[7] Olivier Faugeras and Renaud Keriven. *Variational principles, surface evolution, PDE's, level set methods and the stereo problem*. IEEE, 2002. 2

[8] John Flynn, Michael Broxton, Paul Debevec, Matthew DuVall, Graham Fyffe, Ryan Overbeck, Noah Snavely, and Richard Tucker. Deepview: View synthesis with learned gradient descent. In *CVPR*, pages 2367–2376, 2019. 3

[9] David Gallup, Jan-Michael Frahm, Philippos Mordohai, Qingxiong Yang, and Marc Pollefeys. Real-time plane-sweeping stereo with multiple sweeping directions. In *CVPR*, pages 1–8. IEEE, 2007. 2

[10] Chen Gao, Yichang Shih, Wei-Sheng Lai, Chia-Kai Liang, and Jia-Bin Huang. Portrait neural radiance fields from a single image. *arXiv preprint arXiv:2012.05903*, 2020. 3

[11] Kyle Genova, Forrester Cole, Avneesh Sud, Aaron Sarna, and Thomas Funkhouser. Local deep implicit functions for 3d shape. In *CVPR*, pages 4857–4866, 2020. 2

[12] Steven J Gortler, Radek Grzeszczuk, Richard Szeliski, and Michael F Cohen. The lumigraph. In *SIGGRAPH*, pages 43–54, 1996. 3

[13] Asmaa Hosni, Christoph Rhemann, Michael Bleyer, Carsten Rother, and Margrit Gelautz. Fast cost-volume filtering for visual correspondence and beyond. *TPAMI*, 35(2):504–511, 2012. 2

[14] Ji Hou, Angela Dai, and Matthias Nießner. 3d-sis: 3d semantic instance segmentation of rgb-d scans. In *CVPR*, pages 4421–4430, 2019. 2

[15] Yuxin Hou, Juho Kannala, and Arno Solin. Multi-view stereo by temporal nonparametric fusion. In *ICCV*, pages 2651–2660, 2019. 2

[16] Eddy Ilg, Nikolaus Mayer, Tonmoy Saikia, Margret Keuper, Alexey Dosovitskiy, and Thomas Brox. Flownet 2.0: Evolution of optical flow estimation with deep networks. In *CVPR*, pages 2462–2470, 2017. 2

[17] Sunghoon Im, Hae-Gon Jeon, Stephen Lin, and In So Kweon. Dpsnet: End-to-end deep plane sweep stereo. *arXiv preprint arXiv:1905.00538*, 2019. 2

[18] Chiyu Jiang, Avneesh Sud, Ameesh Makadia, Jingwei Huang, Matthias Nießner, Thomas Funkhouser, et al. Local implicit grid representations for 3d scenes. In *CVPR*, pages 6001–6010, 2020. 2

[19] Abhishek Kar, Christian Häne, and Jitendra Malik. Learning a multi-view stereo machine. *arXiv preprint arXiv:1708.05375*, 2017. 2

[20] Uday Kusupati, Shuo Cheng, Rui Chen, and Hao Su. Normal assisted stereo depth estimation. In *CVPR*, pages 2189–2199, 2020. 2, 11

[21] Kiriakos N Kutulakos and Steven M Seitz. A theory of shape by space carving. *IJCV*, 38(3):199–218, 2000. 2

[22] Marc Levoy and Pat Hanrahan. Light field rendering. In *SIGGRAPH*, pages 31–42, 1996. 3

[23] Zhengqi Li, Tali Dekel, Forrester Cole, Richard Tucker, Noah Snavely, Ce Liu, and William T Freeman. Learning the depths of moving people by watching frozen people. In *CVPR*, pages 4521–4530, 2019. 6, 11

[24] Chao Liu, Jinwei Gu, Kihwan Kim, Srinivasa G Narasimhan, and Jan Kautz. Neural rgb (r) d sensing: Depth and uncertainty from a video camera. In *CVPR*, pages 10986–10995, 2019. 2

[25] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. In *NeurIPS*, 2020. 3, 8

[26] Shaohui Liu, Yinda Zhang, Songyou Peng, Boxin Shi, Marc Pollefeys, and Zhaopeng Cui. Dist: Rendering deep implicit signed distance function with differentiable sphere tracing. In *CVPR*, pages 2019–2028, 2020. 2

[27] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural volumes: Learning dynamic renderable volumes from images. *arXiv preprint arXiv:1906.07751*, 2019. 2

[28] Xiaoxiao Long, Lingjie Liu, Christian Theobalt, and Wenping Wang. Occlusion-aware depth estimation with adaptive normal constraints. In *ECCV*, pages 640–657. Springer, 2020. 2

[29] Xuan Luo, Jia-Bin Huang, Richard Szeliski, Kevin Matzen, and Johannes Kopf. Consistent video depth estimation. *TOG*, 39(4):71–1, 2020. 1, 2, 3, 6, 7, 11

[30] Ricardo Martin-Brualla, Noha Radwan, Mehdi SM Sajjadi, Jonathan T Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *CVPR*, 2021. 2, 3

[31] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *CVPR*, pages 4460–4470, 2019. 2

[32] 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. *TOG*, 38(4):1–14, 2019. 3

[33] 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*, pages 405–421. Springer, 2020. 1, 2, 3, 4, 5, 6, 7, 8, 11

[34] Zak Murez, Tarrence van As, James Bartolozzi, Ayan Sinha, Vijay Badrinarayanan, and Andrew Rabinovich. Atlas: End-to-end 3d scene reconstruction from posed images. In *ECCV*, 2020. 1, 2, 6, 7, 11

[35] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. In *CVPR*, pages 3504–3515, 2020. 2

[36] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *CVPR*, pages 165–174, 2019. 2

[37] Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M Seitz, and Ricardo-Martin Brualla. Deformable Neural Radiance Fields. *arXiv preprint arXiv:2011.12948*, 2020. 3

[38] Daniel Rebain, Wei Jiang, Soroosh Yazdani, Ke Li, Kwang Moo Yi, and Andrea Tagliasacchi. DeRF: Decomposed Radiance Fields. *arXiv preprint arXiv:2011.12490*, 2020. 2

[39] Gernot Riegler and Vladlen Koltun. Free view synthesis. In *ECCV*, pages 623–640. Springer, 2020. 2, 3

[40] Gernot Riegler and Vladlen Koltun. Stable View Synthesis. In *CVPR*, 2021. 2, 3, 7, 8

[41] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In *ICCV*, pages 2304–2314, 2019. 2

[42] Shunsuke Saito, Tomas Simon, Jason Saragih, and Hanbyul Joo. Pifuhd: Multi-level pixel-aligned implicit function for high-resolution 3d human digitization. In *CVPR*, pages 84–93, 2020. 2

[43] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *CVPR*, pages 4104–4113, 2016. 3, 4, 6, 11

[44] Johannes L Schönberger, Enliang Zheng, Jan-Michael Frahm, and Marc Pollefeys. Pixelwise view selection for unstructured multi-view stereo. In *ECCV*, pages 501–518. Springer, 2016. 3, 4, 6, 7, 11

[45] Ayan Sinha, Zak Murez, James Bartolozzi, Vijay Badrinarayanan, and Andrew Rabinovich. Deltas: Depth estimation by learning triangulation and densification of sparse points. In *ECCV*, 2020. 6, 11

[46] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. *arXiv preprint arXiv:1906.01618*, 2019. 2

[47] Pratul P Srinivasan, Richard Tucker, Jonathan T Barron, Ravi Ramamoorthi, Ren Ng, and Noah Snavely. Pushing the boundaries of view extrapolation with multiplane images. In *CVPR*, pages 175–184, 2019. 3

[48] An Tao, Yueqi Duan, Yi Wei, Jiwen Lu, and Jie Zhou. Seggroup: Seg-level supervision for 3d instance and semantic segmentation. *arXiv preprint arXiv:2012.10217*, 2020. 2

[49] Maxim Tatarchenko, Alexey Dosovitskiy, and Thomas Brox. Octree generating networks: Efficient convolutional architectures for high-resolution 3d outputs. In *ICCV*, pages 2088–2096, 2017. 2

[50] Zachary Teed and Jia Deng. Deepv2d: Video to depth with differentiable structure from motion. In *ICLR*, 2020. 6, 11

[51] Julien Valentin, Adarsh Kowdle, Jonathan T Barron, Neal Wadhwa, Max Dzitsiuk, Michael Schoenberg, Vivek Verma, Ambrus Csaszar, Eric Turner, Ivan Dryanovski, et al. Depth from motion for smartphone ar. *ACM Transactions on Graphics (ToG)*, 37(6):1–19, 2018. 5

[52] George Vogiatzis, Philip HS Torr, and Roberto Cipolla. Multi-view stereo via volumetric graph-cuts. In *CVPR*, volume 2, pages 391–398. IEEE, 2005. 2

[53] Kaixuan Wang and Shaojie Shen. Mvdepthnet: Real-time multiview depth estimation neural network. In *3DV*, pages 248–257. IEEE, 2018. 2

[54] Yi Wei, Shaohui Liu, Wang Zhao, and Jiwen Lu. Conditional single-view shape generation for multi-view stereo reconstruction. In *CVPR*, pages 9651–9660, 2019. 2

[55] Yi Wei, Shang Su, Jiwen Lu, and Jie Zhou. FGR: Frustum-Aware Geometric Reasoning for Weakly Supervised 3D Vehicle Detection. In *ICRA*, 2021. 2

[56] Yi Wei, Ziyi Wang, Yongming Rao, Jiwen Lu, and Jie Zhou. PV-RAFT: Point-Voxel Correlation Fields for Scene Flow Estimation of Point Clouds. In *CVPR*, pages 6954–6963, 2021. 2

[57] Wenqi Xian, Jia-Bin Huang, Johannes Kopf, and Changil Kim. Space-time Neural Irradiance Fields for Free-Viewpoint Video. *arXiv preprint arXiv:2011.12950*, 2020. 3

[58] Qingshan Xu and Wenbing Tao. Planar prior assisted patchmatch multi-view stereo. In *AAAI*, volume 34, pages 12516–12523, 2020. 6, 11

[59] Lin Yen-Chen, Pete Florence, Jonathan T Barron, Alberto Rodriguez, Phillip Isola, and Tsung-Yi Lin. iNeRF: Inverting Neural Radiance Fields for Pose Estimation. *arXiv preprint arXiv:2012.05877*, 2020. 3

[60] Kuk-Jin Yoon and In So Kweon. Adaptive support-weight approach for correspondence search. *TPAMI*, 28(4):650–656, 2006. 2

[61] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*, 2020. 1, 2, 4, 5

[62] Wang Zhao, Shaohui Liu, Yezhi Shu, and Yong-Jin Liu. Towards better generalization: Joint depth-pose learning without posenet. In *CVPR*, pages 9151–9161, 2020. 2

[63] Tinghui Zhou, Matthew Brown, Noah Snavely, and David G Lowe. Unsupervised learning of depth and ego-motion from video. In *CVPR*, pages 1851–1858, 2017. 6, 11

[64] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: Learning view synthesis using multiplane images. In *SIGGRAPH*, 2018. 3

| $K$ | $\alpha_l$ | $\alpha_h$ | Abs Rel | Sq Rel | RMSE | RMSE log | $\delta < 1.25$ | $\delta < 1.25^2$ | $\delta < 1.25^3$ |
|-----|-----------|-----------|---------|--------|------|----------|-----------------|-------------------|-------------------|
| 2 | 0.05 | 0.15 | 0.055 | 0.006 | 0.083 | 0.075 | 0.977 | 0.998 | **1.000** |
| 8 | 0.05 | 0.15 | 0.054 | 0.006 | 0.084 | 0.074 | 0.979 | **0.999** | **1.000** |
| 4 | 0.01 | 0.3 | 0.054 | 0.007 | 0.087 | 0.080 | 0.971 | 0.997 | **1.000** |
| 4 | 0.05 | 0.3 | 0.055 | 0.007 | 0.087 | 0.079 | 0.976 | 0.998 | **1.000** |
| 4 | 0.01 | 0.15 | 0.053 | 0.006 | 0.083 | 0.075 | 0.980 | 0.998 | **1.000** |
| 4 | 0.05 | 0.15 | **0.051** | **0.005** | **0.076** | **0.069** | **0.987** | 0.998 | **1.000** |

Table 7: Hyperparameter analysis. The experiment was conducted on scene0521.

# Appendix

## A. Implementation Details

To train the proposed system, we mostly followed NeRF [33]. Specifically, we sampled 64 points in each ray and used a batch of 1024 rays. Since we did not adopt coarse-to-fine strategy in the sampling process, we only need one network (the architecture is same with [33]) to optimize the neural radiance fields. We added random Gaussian noise with zero mean and unit variance to the density $\sigma$ to regularize the network. In addition, following [33], positional encoding was also employed. Adam was adopted as our optimizer with the initial learning rate as $5 \times 10^{-4}$ and decayed exponentially to $5 \times 10^{-5}$. We utilized PyTorch in our implementation. Each scene was trained with 200K iterations on a single RTX 2080 Ti.

**Error metrics.** We follow the metrics in [20, 29, 34, 45, 50, 63] to evaluate depth estimation results:

- Abs Rel: $\frac{1}{|T|} \sum_{y \in T} |y - y^*|/y^*$

- Sq Rel: $\frac{1}{|T|} \sum_{y \in T} ||y - y^*||^2/y^*$

- RMSE: $\sqrt{\frac{1}{|T|} \sum_{y \in T} ||y - y^*||^2}$

- RMSE log: $\sqrt{\frac{1}{|T|} \sum_{y \in T} ||\log y - \log y^*||^2}$

- $\delta < t$: % of $y$ s.t. $\max(\frac{y}{y^*}, \frac{y^*}{y}) = \delta < t$

where $y$ and $y^*$ indicate predicted and groundtruth depths respectively, and T indicates all pixels on the depth image.

## B. Baseline Method Details

We compared our results with several state-of-the-art depth estimation method, which can be roughly classified as four categories:

**Conventional multi-view stereo:** COLMAP [43, 44], ACMP [58]. COLMAP is a non-learning MVS method for 3D reconstruction building upon PatchMatch stereo [2]. Based on COLMAP, ACMP introduces planar models to solve low-textured areas in complex indoor environments.

**Learning-based multi-view stereo:** DELTAS [45], Atlas [34]. These two methods are trained on ScanNet with groundtruth depth supervision. For DELTAS, we used two neighboring frames as the reference frames.

**Monocular depth estimation:** Mannequin Challenge [23]. Mannequin Challenge is a state-of-the-art monocular depth estimation method. We directly used their pretrained weight for evaluation.

**Video-based depth estimation:** CVD [29], DeepV2D [50]. For video-based methods, we sorted images in a scene according to the timeline. DeepV2D is trained on ScanNet with groundtruth depth supervision.

## C. Hyperparameter Analysis

To further demonstrate the effectiveness of our method, we did hyperparameter analysis for the number of used minimum errors $K$, and the bounds $\alpha_l$, $\alpha_h$ used in the guided sampling process. The experiments were conducted on scene0521. Table 7 shows experimental results. We find that using a $K$ that is too small or too large will degrade the performance. On the one hand, it is possible to satisfy the multi-view consistency check although the depths are not correct. Small $K$ will increase the probability of this phenomenon. On the other hand, there are pixels that do not overlap across some view pairs. Thus, the projection errors on some views are invalid and a large $K$ may cover these invalid views. In addition, a large upper bound $\alpha_h$ or a small lower bound $\alpha_l$ for sampling range will lead to worse results, which indicates the necessity to set bounds in sampling process.