# Neural 3D Scene Reconstruction from Multiple 2D Images without 3D Supervision

Yi Guo[1], Che Sun[1], Yunde Jia[2,1], and Yuwei Wu[1,2]

[1]Beijing Key Laboratory of Intellegent Information Technology,
School of Computer Science & Technology, Beijing Institute of Technology, China
[2]Guangdong Laboratory of Machine Perception and Intelligent Computing,
Shenzhen MSU-BIT University, China

{guoyi,sunche,jiayunde,,wuyuwei}@bit.edu.cn

## Abstract

*Neural 3D scene reconstruction methods have achieved impressive performance when reconstructing complex geometry and low-textured regions in indoor scenes. However, these methods heavily rely on 3D data which is costly and time-consuming to obtain in real world. In this paper, we propose a novel neural reconstruction method that reconstructs scenes using sparse depth under the plane constraints without 3D supervision. We introduce a signed distance function field, a color field, and a probability field to represent a scene. We optimize these fields to reconstruct the scene by using differentiable ray marching with accessible 2D images as supervision. We improve the reconstruction quality of complex geometry scene regions with sparse depth obtained by using the geometric constraints. The geometric constraints project 3D points on the surface to similar-looking regions with similar features in different 2D images. We impose the plane constraints to make large planes parallel or vertical to the indoor floor. Both two constraints help reconstruct accurate and smooth geometry structures of the scene. Without 3D supervision, our method achieves competitive performance compared with existing methods that use 3D supervision on the ScanNet dataset.*

## 1. Introduction

Learning-based 3D scene reconstruction methods have achieved good performance from multi-view images. Some works [41, 2, 34] use deep networks to perform some processes in conventional formulation, such as extracting features, matching features, estimating depth, and fusing depth. Others [20, 33] use deep networks to recover 3D representations directly from a whole sequence of images. However, these learning-based methods heavily rely on 3D supervised data to optimize the large scale parameters in
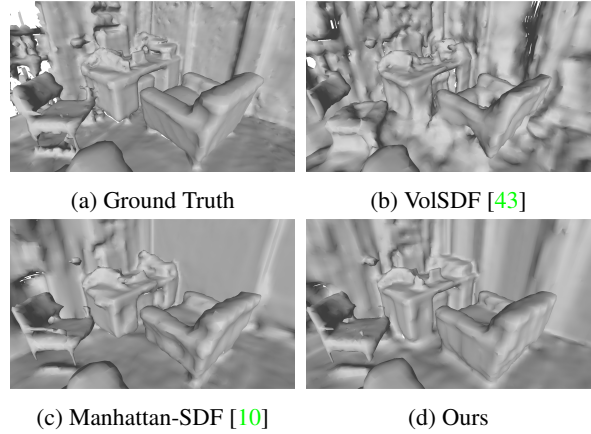


(a) Ground Truth      (b) VolSDF [43]

(c) Manhattan-SDF [10]      (d) Ours

Figure 1. **Comparisons of neural reconstructions.** (a) The ground truth of the scene. (b) The reconstruction result of VolSDF using 2D color supervision. (c) The reconstruction result of Manhattan-SDF with dense depth predictions and the Manhattan assumption. (d) Our result with 2D color supervision by using sparse depth and constraints.

deep networks. And obtaining 3D supervised data of real scenes is costly and time-consuming. In contrast, acquiring 2D image data is simple and inexpensive, which inspires researchers to use 2D images to reconstruct 3D scenes.

Several works [39, 43, 44, 22] reconstruct the 3D geometry structures of scenes from 2D images based on NeRF [19]. These methods represent scenes as neural implicit fields and combine the surface representation with volume rendering, and achieve impressive results. However, they do not work well in reconstructing some challenging indoor scenes. For example, VoISDF [43] fails to reconstruct the regions with complex geometric structures and large low-textured planes, as shown in Figure 1b.

In order for handling this issue, Guo *et al.* [10] use dense

depth and the Manhattan-world assumption to reconstruct indoor scenes with 2D images, and achieve impressive performance. However, dense depth from 2D images is inevitably noisy, incurring poor reconstruction quality of the planes, such as tables and chairs shown in Figure 1c, because these planes do not fit the Manhattan-world assumption. To improve the work [10], Wang *et al.* [38] and Yu *et al.* [46] use pre-trained networks to obtain extra 3D data, such as dense depth and normal maps for improving the reconstruction quality of plane regions. However, the pre-trained networks have to acquire extra 3D supervised data for fine-tuning. Obtaining sufficient 3D supervised data is often costly and time-consuming. Our work, without using any 3D supervised data, focuses on improving the reconstruction quality of indoor scenes.

In this paper, we propose a novel method that reconstructs indoor scenes using spare depth maps under the plane constraints from 2D images. The sparse depth from 2D images is easily obtained and has less noise compared with dense depth. The plane constraints would help reconstruct low-textured plane regions in the scene, not only walls and floors in the Manhattan-world assumption but also planes of tables and chairs, as shown in Figure 1d. That is to say, our plane constraints are more general, and have a wider scope of applications and more relaxed assumptions, compared with the Manhattan constraints. Specifically, our method represents scenes as a signed distance function field, a color field, and a plane probability field, and optimizes these fields by volume rendering [43] to reconstruct the scenes. We obtain sparse depth by using the geometric constraints. The geometric constraints project 3D points on the surface to similar-looking regions with similar features in different views, which ensures that corresponding points obtained by image feature matching represent the same 3D point on the surface. Besides, we utilize plane constraints to make large planes in the scene parallel or vertical to the wall or floor. We estimate large planes in images and impose the plane constraints by making the normal of planes parallel or orthogonal to the normal of the ground. The plane constraints ensure smooth reconstruction for low-textured regions.

Our method is evaluated on the ScanNet [4]. Our method with sparse depth achieves comparable results with Manhattan-SDF with dense depth. Our method with dense depth outperforms the Manhattan-SDF with dense depth. Our method with dense depth achieves comparable results with existing methods that use 3D supervision.

In summary, our contributions are as follows:

• We propose a novel neural reconstruction method that uses sparse depth to achieve high quality reconstruction of indoor scenes without 3D supervision.

• We introduce the plane constraints to improve the reconstruction quality of planes that do not fit the Manhattan-

world assumption. Our plane constraints are able to apply to more general planes.

## 2. Related Work

Reconstructing 3D geometry structures of a scene from a sequence of images has been a longstanding computer vision problem. Traditional multi-view 3D scene reconstruction methods [29, 30, 31] usually focus on estimating depth from a sequence of images, then fusing the depth maps and reconstructing the surface by screened Poisson surface reconstruction [12]. Some works use deep neural networks to learn extracting features [37, 47, 32], matching features [14, 17, 36], estimating depth maps [11, 41, 42, 45] or fusing depth maps [7, 26] from a sequence of images. These deep methods improve performance compared with traditional 3D scene reconstruction methods. However, the fused reconstruction results are prone to be either layered or scattered, due to estimating depth of each key frame individually and estimation errors. To address the problem, neural scene reconstruction methods have been proposed. For example, Atlas [20] extracts the image features and regresses the TSDF of the scene. NeuralRecon [33] regresses TSDF directly using a coarse-to-fine framework which achieves real-time indoor scene reconstruction. The TSDF fails to perform high-resolution reconstruction very well, because the TSDF representation is based on volume voxel. To this end, some works use the neural implicit function to represent the scene, such as occupancy [18, 25, 3, 35] and SDF [23, 28]. The neural implicit representation is naturally continuous, which means free of the limitation of finite resolution. These learning-based methods require 3D supervision for training, but obtaining it is often time-consuming and costly in real scenes. In contrast, our method reconstructs large-scale scenes without 3D supervision.

Several recent methods [21, 19, 40, 22] have demonstrated differentiable rendering for successfully reconstructing 3D scenes directly from 2D images. These methods can be divided into two groups based on the rendering technique: surface rendering and volume rendering. Surface rendering based methods, such as DVR [21] and IDR [44], achieve good performance in most cases, but these methods require pixel-accurate object masks for all images as input. Therefore, these methods achieve unsatisfied performance for complex objects and scenes. Volume rendering based methods, such as NeRF [19] and its variants [40, 5, 48], render an image by learning alpha-compositing of a radiance field along rays. These methods have shown good performance on novel view synthesis. They encode a scene as continuous radiance fields of color and volume density, and map the position and view direction to an image by using differentiable ray marching. Using volume density to represent the scene fails to extract high-quality surfaces well since the volume density representation does not have suf-
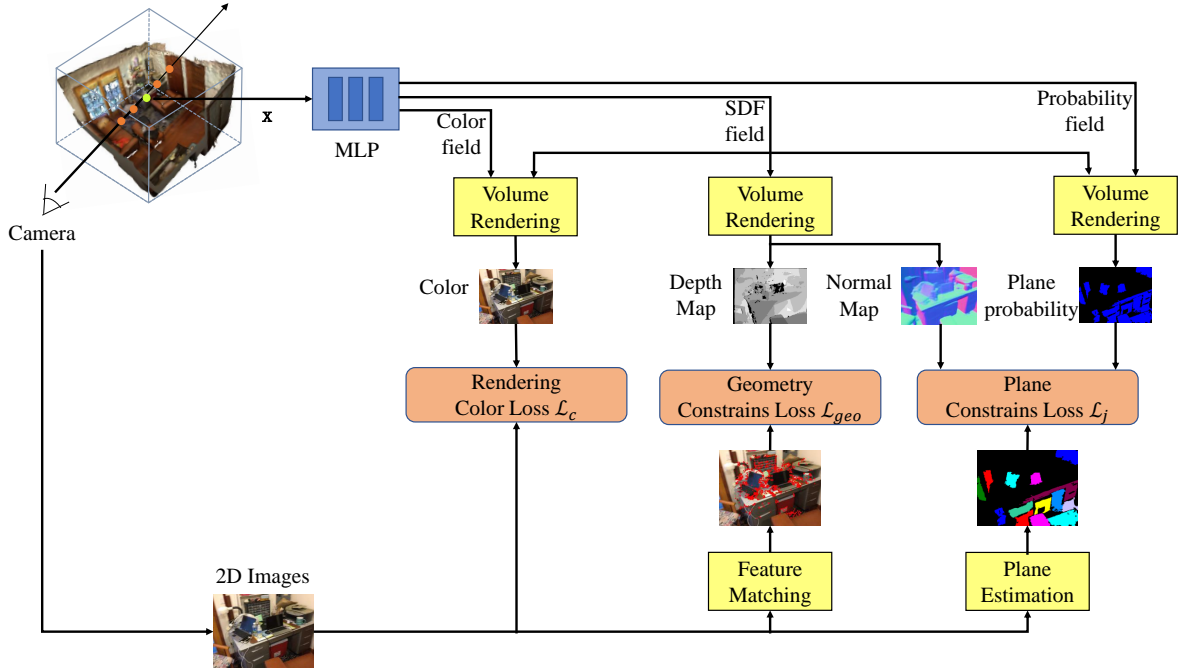
Figure 2. **The overview of our method.**

ficient constraints on geometry. UNISURF [22], NeuS [39] and VolSDF [43] combine the surface representation such as occupancy values and signed distance function with the volume rendering to reconstruct the scene well and achieved good results. However, these methods may fail to handle complex geometry and large low-textured regions in large scenes. Differently, our method introduces geometry constraints and plane constraints to handle these regions in the large-scale indoor scene without 3D supervision.

Recently, some works focus on large-scale indoor scene reconstruction by using volume rendering, and achieve good results. Manhattan-SDF [10] uses dense depth maps predicted by COLMAP [29] as supervision and uses Manhattan assumption to handle the walls and floor regions. However, obtaining dense depth maps is time-consuming by traditional methods without 3D supervision and the depth maps are noisy. The plane reconstruction is poor in regions where the Manhattan World assumption does not hold, due to the noise in depth estimation. To address these problems, NeuRIS [38] utilizes dense normal maps predicted by a monocular method [6] as supervision. MonoSDF [46] utilizes both depth and normal maps as supervision, and uses multi-resolution feature grids to help reconstruct the indoor scene. However, these two methods heavily rely on pre-trained networks supervised by 3D data. In contrast to these methods, we perform differentiable volume rendering for scene reconstruction by using sparse depth under plane constraints without 3D supervision, to improve the plane

reconstruction quality of scenes.

## 3. Method

Given multi-view 2D images with camera poses, our goal is to reconstruct high-quality scenes without 3D supervised data. We represent a scene as a signed distance function (SDF) field, a color field, and a plane probability field, and optimize them by volume rendering. We use geometry constraints to obtain sparse depth for improving the reconstruction quality of the regions that have complex geometry structures. We estimate the plane of the scene and use plane constraints to improve the reconstruction quality of the large low-textured regions. The overview of our method is illustrated in Figure 2.

### 3.1. Implicit Scene Representation

We utilize two fields of SDF and color to model the scene geometry and scene appearance, respectively. We represent the scene geometry as an SDF field. The SDF takes a 3D point $\mathbf{x} \in \mathbb{R}^3$ as the input, and generates its distance $s(\mathbf{x})$ to the closest surface by using a multi-layer perceptron (MLP) neural network $f_g$, given by

$$\big(s(\mathbf{x}), \mathbf{z}(\mathbf{x})\big) = f_g(\mathbf{x}; \theta_g), \tag{1}$$

where $f_g$ is implemented as the MLP network with parameters $\theta_g$, and $\mathbf{z}(\mathbf{x})$ is the geometry feature calculated in [44]. The surface is defined as the zero level set of the SDF, that

is,

$$S = \{\mathbf{x}|s(\mathbf{x}) = 0\}. \tag{2}$$

We represent the appearance of the scene as a color field. We define a network $f_c$ to predict the RGB color $\mathbf{c}(\mathbf{x})$ for a 3D point $\mathbf{x}$ and a viewing direction $\mathbf{v}$, given by

$$\mathbf{c}(\mathbf{x}) = f_c(\mathbf{x}, \mathbf{v}, \mathbf{n}(\mathbf{x}), \mathbf{z}(\mathbf{x}); \theta_c), \tag{3}$$

where $\mathbf{n}(\mathbf{x})$ is the unit normal obtained by computing the gradient of our SDF function $f_g$, $\mathbf{z}(\mathbf{x})$ is the geometry feature of the output of the MLP, and $\theta_c$ is the parameters of the network.

## 3.2. Volume Rendering of Implicit Surfaces

Following [43, 39], we apply differentiable volume rendering to optimize the scene representation from images. Specifically, we cast a ray $\mathbf{r} = \mathbf{o} + t\mathbf{v}$ with the origin in the camera center $\mathbf{o}$ and the viewing direction $\mathbf{v}$ to render a pixel of the image. We sample $N$ points $\{\mathbf{x}_i | \mathbf{x}_i = \mathbf{o} + t_i\mathbf{v}, i = 1, 2, 3..., N\}$ along the camera ray and predict the $s(\mathbf{x}_i)$ and $\mathbf{c}(\mathbf{x}_i)$ for each point. For convenience, here we use the $s_i$ and $\mathbf{c}_i$ to represent the $s(\mathbf{x}_i)$ and $\mathbf{c}(\mathbf{x}_i)$. We transform the SDF $s_i$ to the volume density $\sigma_i$ by

$$\sigma(s) = \begin{cases} \frac{1}{2\beta}\exp\left(\frac{s}{\beta}\right) & \text{if } s \leq 0, \\ \frac{1}{\beta}\left(1 - \frac{1}{2}\exp\left(-\frac{s}{\beta}\right)\right) & \text{if } s > 0, \end{cases} \tag{4}$$

where $\beta$ is a learnable parameter. We accumulate the color $\hat{C}(\mathbf{r})$ along the ray $\mathbf{r}$ via numerical quadrature [19]:

$$\hat{C}(\mathbf{r}) = \sum_{i=1}^{N} T_i(1 - \exp(-\sigma_i\delta_i))\mathbf{c}_i,$$

$$\text{where} \quad T_i = \exp\left(-\sum_{j=1}^{i-1}\sigma_j\delta_j\right), \tag{5}$$

and $\delta_i$ is the distance between adjacent samples.

## 3.3. Geometry Constraints

We utilize geometry constraints to obtain sparse depth for helping reconstruct the regions with complex geometry structures of the indoor scene. We observe that (1) the 3D point on the surface is projected to similar-looking regions with similar features in different views, as shown in Figure 3; (2) the regions with complex geometry structures always have sharp features. These two observations inspire us to use geometric constraints to generate sparse depth for improving the reconstruction quality of the regions with complex geometry.

As illustrated in Figure 3, given multi-view images, we first extract the feature points (e.g., ORB feature points [27]
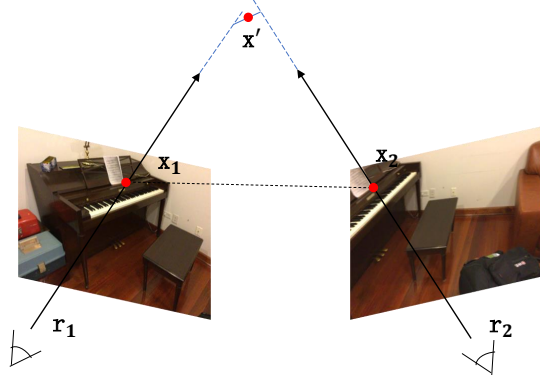


Figure 3. **The illustration of geometry constraints.** Given two images, we match the feature points $\mathbf{x}_1$ and $\mathbf{x}_1$. We cast rays $\mathbf{r}_1$ and $\mathbf{r}_2$ from the camera center through the matched points. The midpoint $\mathbf{x}'$ of the common perpendiculars of rays is on the scene surface.

or SIFT feature points [16]) of the images, and then we obtain the point correspondences by matching the feature points in adjacent images. We cast a ray from the camera center through the pixel of each matching point. For one correspondence of two points $\mathbf{x}_1$ and $\mathbf{x}_2$, two rays $\mathbf{r}_1$ and $\mathbf{r}_2$ will intersect at a point on the surface of the scene, theoretically, but they usually do not intersect due to the existence of errors. Therefore, we calculate the point closest to the two lines as an approximate intersection point. For the two rays, we calculate their common perpendiculars, and the approximate intersection point is the midpoint $\mathbf{x}'$ of the line between the common perpendiculars and the actual intersection point of the two rays. The approximate depth $D_{app}$ of the two rays can be obtained by projecting the approximate intersection points on the two rays respectively. If the distance between the two rays is larger than a threshold, we consider that this correspondence is wrong and discard it.

For the scene representation, we calculate the depth $\hat{D}(\mathbf{r})$ from a viewpoint by using the volume rendering:

$$\hat{D}(\mathbf{r}) = \sum_{i=1}^{N} T_i(1 - \exp(-\sigma_i\delta_i))t_i. \tag{6}$$

We use a geometry loss function $\mathcal{L}_{geo}$ of matching points to assist the learning of the textured regions,

$$\mathcal{L}_{geo} = \sum_{\mathbf{r}\in\mathcal{M}} |\hat{D}(\mathbf{r}) - D_{app}(\mathbf{r})|, \tag{7}$$

where $\mathcal{M}$ is the set of the camera rays along the matching points and $|\cdot|$ is the absolute value. This geometry loss improves the reconstruction quality of the regions which have sharp features.
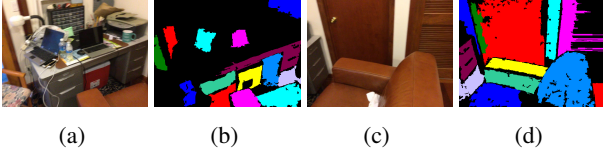
Figure 4. **Two examples of plane estimation.** (a) and (c) are the input images, (b) and (d) are the results of plane estimation. The colored regions represent the large planes in the scene that need to be imposed by plane constraints. The black regions do not contain large planes. We can see that not only walls and floors, but also the surfaces of tables have been successfully estimated.

## 3.4. Plane Constraints

We utilize the plane constraints to help reconstruct the low-textured plane regions. We observe that the large low-textured plane regions are located not only on floors and walls, but also lie on tables and beds, as shown in Figure 4. These regions are often parallel or vertical to the floor. For example, the tabletop and bed in Figure 4 are both parallel to the floor. The wall, the side of the table, the bookshelf, and other planes are all perpendicular to the floor.

According to the observation, we apply the plane constraints to help reconstruct the large low-textured regions. We utilize the Felzenswalb superpixel segmentation algorithm [8] to obtain the plane regions. The algorithm follows a greedy approach and segments areas with low gradients, and produces more plane regions. The plane segmentation results contain many small planes which are discarded in the plane constraints. We only keep the planes which are larger than a certain proportion of the image to impose the constraints. We calculate the normal $\hat{N}(\mathbf{r})$ from a viewpoint by using the volume rendering:

$$\hat{N}(\mathbf{r}) = \sum_{i=1}^{N} T_i(1 - \exp(-\sigma_i \delta_i))\mathbf{n}_i,$$

$$\text{where} \quad \mathbf{n}_i = \nabla f_g(\mathbf{x}_i; \theta_g), \tag{8}$$

and $\nabla f_g(\mathbf{x}_i; \theta_g)$ is the spatial gradient of SDF. We assume the floors are vertical to the z-axis. We design a plane loss function $\mathcal{L}_{pla}$ that enforces the normal of large regions to be parallel or orthogonal to the upper unit vector as

$$\mathcal{L}_{pla} = \sum_{\mathbf{r} \in \mathcal{P}} \min_{i \in \{-1,0,1\}} |i - \hat{N}(\mathbf{r}) \cdot \mathbf{n}_f|, \tag{9}$$

where $\mathcal{P}$ is the set of camera rays of images pixels that are segmented as the large plane regions, and $\mathbf{n}_f = (0, 0, 1)$ is the upper unit vector that denotes the normal of floors.

We use an MLP network $f_p$ to denote the probability of a point on the large plane. The probability logits are defined as

$$p(\mathbf{x}) = f_p(\mathbf{x}; \theta_p), \tag{10}$$

where $\theta_p$ is the parameters of the MLP. For each image, we render the probability logits similar to image rendering as

$$\hat{P}(\mathbf{r}) = \sum_{i=1}^{N} T_i(1 - \exp(-\sigma_i \delta_i))p_i. \tag{11}$$

We use a joint optimization loss to jointly optimize the scene representation and plane region estimation results, that is

$$\mathcal{L}_j = \sum_{\mathbf{r} \in \mathcal{P}} \hat{P}(\mathbf{r})\mathcal{L}_{pla}(\mathbf{r}) + \mathcal{L}_p,$$

$$\text{where} \quad \mathcal{L}_p = -\sum_{\mathbf{r} \in \mathcal{R}} P(\mathbf{r}) \log \hat{P}(\mathbf{r}), \tag{12}$$

and $\hat{P}(\mathbf{r})$ is the rendered probability, $P(\mathbf{r})$ is the plane obtained by the Felzenswalb. $\mathcal{L}_p$ is the cross entropy loss to avoid the $\hat{P}(\mathbf{r})$ converging to zero.

## 3.5. Training

During the training stage, we sample a batch of pixels and minimize the loss functions of the color and the constraints. The overall loss is defined as

$$\mathcal{L} = \lambda_c \mathcal{L}_c + \lambda_{geo} \mathcal{L}_{geo} + \lambda_j \mathcal{L}_j + \lambda_{eik} \mathcal{L}_{eik}. \tag{13}$$

The color loss $\mathcal{L}_c$ is defined as

$$\mathcal{L}_c = \sum_{\mathbf{r} \in \mathcal{R}} ||\hat{C}(\mathbf{r}) - C(\mathbf{r})||, \tag{14}$$

where $\mathcal{R}$ is the set of sample pixel, $C(\mathbf{r})$ is the ground truth pixel color, and $|| \cdot ||$ is the 1-norm.

The Eikonal loss [9] is introduced to regularize SDF values in 3D space:

$$\mathcal{L}_{eik} = \sum_{\mathbf{x} \in \mathcal{X}} (||\nabla f_g(\mathbf{x}; \theta_g)||_2 - 1)^2, \tag{15}$$

where $\mathcal{X}$ are a set of uniform sampling points and near surface points, and $|| \cdot ||_2$ is the 2-norm.

## 4. Implementation Details

We implement our method in PyTorch [24] and use Adam optimizer [13] with a learning rate of 5e-4. We sample 1024 rays from pixel points for each batch to train the network. The network is trained for 50k iterations on one NVIDIA RTX3090 GPU. We use sphere initialization [1] to initialize the network parameters. In the early stage of training, we aim to reconstruct the structure of the scene first, so we tend to select matching pixel points for training and set a larger geometry weight $\lambda_{geo}$. As the training process goes on, we reduce the weight of the geometric loss and random sample pixel points for training. We use Marching Cubes algorithm [15] to extract surface mesh from the learned signed distance function.

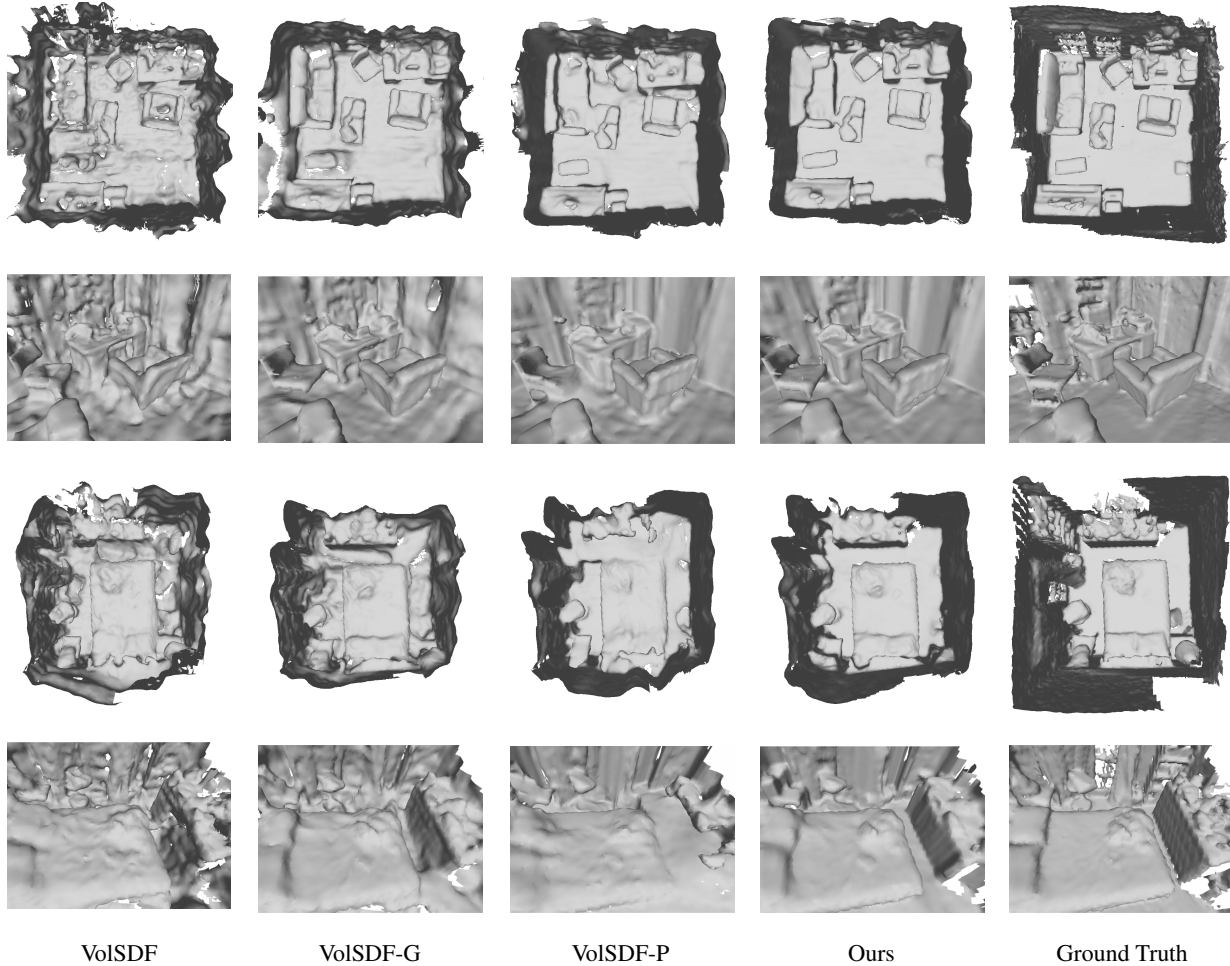| | | | | |
|:---:|:---:|:---:|:---:|:---:|
| VolSDF | VolSDF-G | VolSDF-P | Ours | Ground Truth |

Figure 5. **Qualitative results of ablation studies on ScanNet.** Our method produces much more accurate and smooth reconstruction results compared with our baselines, i.e., VolSDF, VolSDF-G, and VolSDF-P. VolSDF use 2D images to reconstruct the scenes without constraints.VolSDF-G uses the geometry constraints to reconstruct more accurate results than the VolSDF. VolSDF-P uses the plane constraints to reconstruct smoother and more complete planes compared with VolSDF. Our method reconstructs more accurate and smoother results than VolSDF-G and VolSDF-P, by imposing the two constraints.

## 5. Experiments

**Dataset.** We perform the experiments on the ScanNet(V2) [4], which contains 1613 indoor scenes. It provides ground truth camera poses of 2D images, surface reconstructions and instance-level semantic segmentations of the scenes. For each scene, it contains 1K-5K RGB-D images, and we uniform sample one tenth images for reconstruction. The experiment setup is the same as the work of Guo *et al.* [10].

**Metrics.** We evaluate the 3D surface geometry results using 5 standard metrics defined in [20]: accuracy, completeness, precision, recall, and F-score. Among these metrics, F-score is usually considered as the most suitable metric to measure 3D reconstruction quality according to the work of Sun *et al.* [33]. The definitions of these metrics are detailed

| Method | Acc↓ | Comp↓ | Pre↑ | Recall↑ | **F-score↑** |
|:---:|:---:|:---:|:---:|:---:|:---:|
| VolSDF | 0.414 | 0.120 | 0.321 | 0.394 | 0.346 |
| VolSDF-G | 0.106 | 0.101 | 0.507 | 0.469 | 0.487 |
| VolSDF-P | 0.100 | 0.118 | 0.493 | 0.427 | 0.459 |
| Ours | **0.068** | **0.079** | **0.626** | **0.551** | **0.586** |

Table 1. **Quantitative results of ablation studies on ScanNet.** Our method significantly improves the accuracy and completeness compared of our baselines of VolSDF, VolSDF-G and VolSDF-P.

in the *supplementary material*.

**Competing methods.** The competing methods include: (1) Classical MVS method: COLMAP [29], we use ground truth camera poses to reconstruct the point cloud, and use

| 3D | Method | dense | Acc↓ | Comp↓ | Pre↑ | Recall↑ | **F-score↑** |
|---|---|---|---|---|---|---|---|
| | NeRF [19] | ✗ | 0.735 | 0.177 | 0.131 | 0.290 | 0.176 |
| | UNISURF [22] | ✗ | 0.554 | 0.164 | 0.212 | 0.362 | 0.267 |
| | NeuS [39] | ✗ | 0.179 | 0.208 | 0.313 | 0.275 | 0.291 |
| ✗ | VolSDF [43] | ✗ | 0.414 | 0.120 | 0.321 | 0.394 | 0.346 |
| | Manhattan-SDF-s [10] | ✗ | 0.076 | 0.079 | 0.577 | 0.506 | 0.541 |
| | **Ours** | ✗ | **0.068** | **0.079** | **0.626** | **0.551** | **0.586** |
| | COLMAP [29] | ✓ | **0.047** | 0.235 | 0.711 | 0.441 | 0.537 |
| | Manhattan-SDF [10] | ✓ | 0.072 | 0.068 | 0.621 | 0.586 | 0.602 |
| | **Ours-d** | ✓ | 0.057 | **0.057** | **0.708** | **0.677** | **0.692** |
| ✓ | Atlas [20] | ✓ | 0.124 | 0.074 | 0.413 | **0.711** | 0.520 |
| | NeuRIS [38] | ✓ | **0.050** | **0.049** | **0.717** | 0.669 | **0.692** |

Table 2. **Quantitative results of our method and state-of-the-art mthods on ScanNet.** The first column means whether the method uses 3D data or pre-trained networks supervised by 3D data. The third column means whether the method uses dense depth or normal maps. Our method with sparse depth outperforms the existing method using sparse depth.

Screened Poisson Surface Reconstruction [12] to reconstruct mesh from a point cloud. (2) TSDF based method: Atlas [20]. Atlas directly regresses a TSDF from a set of posed RGB images. (3) Neural volume rendering based methods: NeRF [19], UNISURF [22], NeuS [39] and VolSDF [43]. These methods do not use 3D data. (4) State-of-the-art neural scene reconstruction methods: Manhattan-SDF [10], NeuRIS [38]. For neural volume rendering based methods and neural scene reconstruction methods, we extract mesh by using Marching Cubes algorithm[15].

## 5.1. Ablation Studies

We conduct ablation studies on ScanNet to verify the effectiveness of each component in our method. We conduct experiments in four different settings: (1) Baseline setting of VolSDF [43], we train networks only with 2D image supervision. (2) VolSDF-G, we add the sparse depth obtained by geometry constraints to the VolSDF. (3) VolSDF-P, we add the plane constraints to the VolSDF. (4) Ours, we learn the scene with both sparse depth and the plane constraints. The quantitative results are shown in Table 1, and the qualitative results are shown in Figure 5.

The comparison between VolSDF and VolSDF-G shows that the sparse depth obtained by geometry constraints in our method improves 0.141 precision in terms of F-score, as shown in Table 1. The sparse depth obtained by geometry constraints helps reconstruct the regions that have complex geometry structures, but the reconstruction quality is still poor for the large low-textured regions, as shown in VolSDF and VolSDF-G in Figure 5.

The comparison between VolSDF and VolSDF-P shows that the plane constraints in our method improve 0.113 precision in terms of F-score, as shown in Table 1. The plane constraints help reconstruct the large low-textured regions, such as the wall, floor and table. Nevertheless, the reconstruction quality is still poor for the regions with sharp fea-

tures, as shown in VolSDF and VolSDF-P in Figure 5.

Our method improves 0.240 precision in terms of F-score compared with VolSDF. Our method combines the sparse depth and plane constraints, and reconstructs both sharp features regions and large low-textured regions well.

## 5.2. Comparisons

We compare our method with state-of-the-art reconstruction methods on ScanNet. The averaged quantitative results of the geometry quality are shown in Table 2. The qualitative results are shown in the Figure 6. The individual scene results are in the *supplementary material*.

We categorize existing methods into two types according to whether they use 3D supervision, as shown in Table 2. We first compare our method with existing methods that do not use 3D supervision and do not need dense depth. The performance of NeRF [19] is poor, because its volume density representation of the scene has geometry ambiguity. UNISURF [22], Neus [39] and VolSDF [43] represent the scene as occupancy and SDF to improve the reconstruction of the scene surface. However, they do not perform well in indoor scenes. Manhattan-SDF [10] with sparse depth, *i.e.*, Manhattan-SDF-s, has achieved performance improvement, but compared with Manhattan-SDF [10], the reconstruction quality of Manhattan-SDF-s is significantly reduced because of the lack of sufficient constraints on the planes. Our method with sparse depth outperforms these methods in indoor scenes, which proves the superiority of our method. Then, we compare our method with existing methods using dense depth without 3D supervision. COLMAP [29] achieves high accuracy and precision, because it filters out inconsistent reconstruction points between multiple views in the fusion phase. Its completeness and recall are low. Manhattan-SDF [10] performs well, but it is limited by the Manhattan-world assumption. Our method with dense depth outperforms Manhattan-SDF [10],
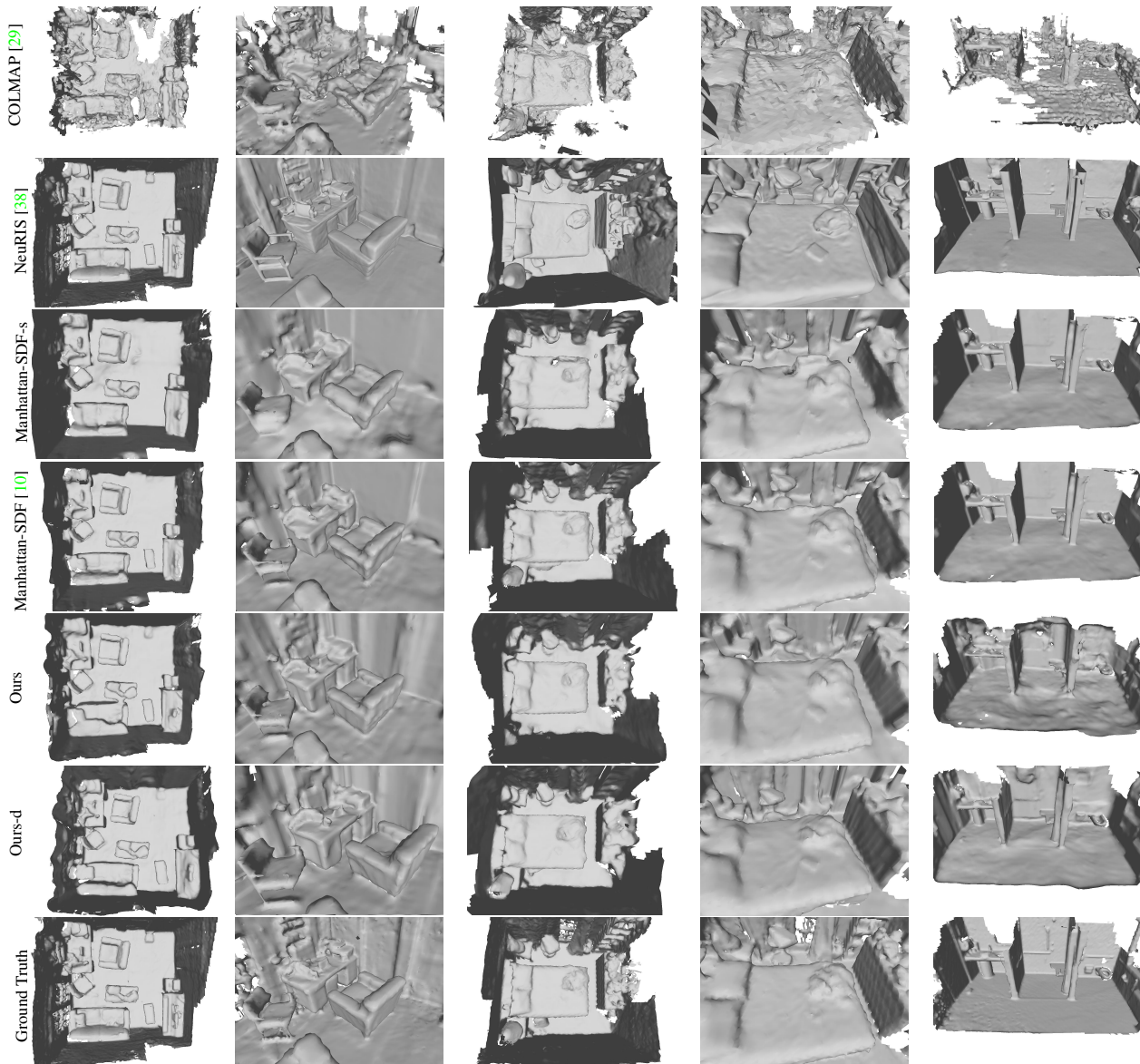
Figure 6. **Qualitative results of our method and sate-of-the-art methods on ScanNet.** Our method outperforms COLMAP and Manhattan-SDF methods, and performs comparable results with NeuRIS with pre-trained networks. Our method produces accurate results in the regions of complex structures and coherent results in the regions of large planes.

verifying the effectiveness of the plane constraints (compared with the Manhattan-world assumption).

We finally compare our method with two methods using dense depth **with 3D supervision**. Atlas [20] achieves high completeness and recall due to its TSDF completion capability, but the accuracy and precision are low. NeuRIS [38] performs much well in indoor scenes, but it needs pretrained networks finetuned by 3D supervision to provide normal maps to supervise the networks. Our method using dense depth without any 3D supervision achieves comparable results with NeuRIS [38] that utilizes pre-training

networks.

## 6. Conclusion

In this paper, we present a novel 3D scene reconstruction method without 3D supervision. Our method can reconstruct indoor scenes with accessible 2D images as supervision by using sparse depth under the plane constraints. The sparse depth obtained by using geometry constraints can guarantee the good reconstruction quality of the scene regions with complex geometry structures. The plane constraints can guarantees the good reconstruction quality of

the scene regions with large low-textured planes by making large planes keep parallel or vertical to the wall or floor. Experiments can show that our method reconstructs the scene completely and accurately, and performs competitive reconstruction quality to the methods that require 3D data or pretrained networks finetuned by 3D data.

# References

[1] Matan Atzmon and Yaron Lipman. SAL: Sign agnostic learning of shapes from raw data. In *CVPR*, pages 2565–2574, 2020. 5

[2] Shuo Cheng, Zexiang Xu, Shilin Zhu, Zhuwen Li, Li Erran Li, Ravi Ramamoorthi, and Hao Su. Deep stereo using adaptive thin volume representation with uncertainty awareness. In *CVPR*, pages 2524–2534, 2020. 1

[3] Christopher B Choy, Danfei Xu, JunYoung Gwak, Kevin Chen, and Silvio Savarese. 3d-r2n2: A unified approach for single and multi-view 3d object reconstruction. In *ECCV*, pages 628–644, 2016. 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. 2, 6

[5] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised nerf: Fewer views and faster training for free. In *CVPR*, pages 12882–12891, 2022. 2

[6] Tien Do, Khiem Vuong, Stergios I Roumeliotis, and Hyun Soo Park. Surface normal estimation of tilted images via spatial rectifier. In *ECCV*, pages 265–280, 2020. 3

[7] Simon Donne and Andreas Geiger. Learning non-volumetric depth fusion using successive reprojections. In *CVPR*, pages 7634–7643, 2019. 2

[8] Pedro F Felzenszwalb and Daniel P Huttenlocher. Efficient graph-based image segmentation. *IJCV*, 59(2):167–181, 2004. 5

[9] Amos Gropp, Lior Yariv, Niv Haim, Matan Atzmon, and Yaron Lipman. Implicit geometric regularization for learning shapes. *T-PAMI*, pages 3569–3579, 2020. 5

[10] Haoyu Guo, Sida Peng, Haotong Lin, Qianqian Wang, Guofeng Zhang, Hujun Bao, and Xiaowei Zhou. Neural 3d scene reconstruction with the manhattan-world assumption. In *CVPR*, pages 5511–5520, 2022. 1, 2, 3, 6, 7, 8

[11] Po-Han Huang, Kevin Matzen, Johannes Kopf, Narendra Ahuja, and Jia-Bin Huang. Deepmvs: Learning multi-view stereopsis. In *CVPR*, pages 2821–2830, 2018. 2

[12] Michael Kazhdan and Hugues Hoppe. Screened poisson surface reconstruction. *ACM TOG*, pages 1–13, 2013. 2, 7

[13] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 5

[14] Vincent Leroy, Jean-Sébastien Franco, and Edmond Boyer. Shape reconstruction using volume sweeping and learned photoconsistency. In *ECCV*, pages 781–796, 2018. 2

[15] William E Lorensen and Harvey E Cline. Marching cubes: A high resolution 3D surface construction algorithm. *SIGGRAPH*, pages 163–169, 1987. 5, 7

[16] David G Lowe. Distinctive image features from scale-invariant keypoints. *IJCV*, 60(2):91–110, 2004. 4

[17] Wenjie Luo, Alexander G Schwing, and Raquel Urtasun. Efficient deep learning for stereo matching. In *CVPR*, pages 5695–5703, 2016. 2

[18] 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

[19] 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, 2020. 1, 2, 4, 7

[20] 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*, pages 414–431, 2020. 1, 2, 6, 7, 8

[21] 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

[22] Michael Oechsle, Songyou Peng, and Andreas Geiger. UNISURF: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction. In *ICCV*, 2021. 1, 2, 3, 7

[23] 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

[24] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca

Antiga, et al. PyTorch: An imperative style, high-performance deep learning library. In *NeurIPS*, pages 8026–8037, 2019. 5

[25] Songyou Peng, Michael Niemeyer, Lars Mescheder, Marc Pollefeys, and Andreas Geiger. Convolutional occupancy networks. In *ECCV*, pages 523–540, 2020. 2

[26] Gernot Riegler, Ali Osman Ulusoy, Horst Bischof, and Andreas Geiger. Octnetfusion: Learning depth fusion from data. In *3DV*, pages 57–66. IEEE, 2017. 2

[27] Ethan Rublee, Vincent Rabaud, Kurt Konolige, and Gary Bradski. Orb: An efficient alternative to sift or surf. In *ICCV*, pages 2564–2571. IEEE, 2011. 4

[28] 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

[29] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *CVPR*, pages 4104–4113, 2016. 2, 3, 6, 7, 8

[30] Johannes L Schönberger, Enliang Zheng, Jan-Michael Frahm, and Marc Pollefeys. Pixelwise view selection for unstructured multi-view stereo. In *ECCV*, pages 501–518, 2016. 2

[31] Steven M Seitz, Brian Curless, James Diebel, Daniel Scharstein, and Richard Szeliski. A comparison and evaluation of multi-view stereo reconstruction algorithms. In *CVPR*, volume 1, pages 519–528. IEEE, 2006. 2

[32] Iago Suárez, Ghesn Sfeir, José M Buenaposada, and Luis Baumela. Beblid: Boosted efficient binary local image descriptor. *PR*, 133:366–372, 2020. 2

[33] Jiaming Sun, Yiming Xie, Linghao Chen, Xiaowei Zhou, and Hujun Bao. NeuralRecon: Real-time coherent 3D reconstruction from monocular video. In *CVPR*, pages 15598–15607, 2021. 1, 2, 6

[34] Chengzhou Tang and Ping Tan. Ba-net: Dense bundle adjustment network. *arXiv preprint arXiv:1806.04807*, 2018. 1

[35] 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

[36] Benjamin Ummenhofer, Huizhong Zhou, Jonas Uhrig, Nikolaus Mayer, Eddy Ilg, Alexey Dosovitskiy, and Thomas Brox. Demon: Depth and motion network for learning monocular stereo. In *CVPR*, pages 5038–5047, 2017. 2

[37] Fangjinhua Wang, Silvano Galliani, Christoph Vogel, Pablo Speciale, and Marc Pollefeys. PatchmatchNet: Learned multi-view patchmatch stereo. In *CVPR*, pages 14194–14203, 2021. 2

[38] Jiepeng Wang, Peng Wang, Xiaoxiao Long, Christian Theobalt, Taku Komura, Lingjie Liu, and Wenping Wang. Neuris: Neural reconstruction of indoor scenes using normal priors. *arXiv preprint arXiv:2206.13597*, 2022. 2, 3, 7, 8

[39] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. NeuS: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. In *NeurIPS*, 2021. 1, 3, 4, 7

[40] Yi Wei, Shaohui Liu, Yongming Rao, Wang Zhao, Jiwen Lu, and Jie Zhou. NerfingMVS: Guided optimization of neural radiance fields for indoor multi-view stereo. In *ICCV*, pages 5610–5619, 2021. 2

[41] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. MVSNet: Depth inference for unstructured multi-view stereo. In *ECCV*, pages 767–783, 2018. 1, 2

[42] Yao Yao, Zixin Luo, Shiwei Li, Tianwei Shen, Tian Fang, and Long Quan. Recurrent mvsnet for high-resolution multi-view stereo depth inference. In *CVPR*, pages 5525–5534, 2019. 2

[43] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces. In *NeurIPS*, 2021. 1, 2, 3, 4, 7

[44] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Ronen Basri, and Yaron Lipman. Multiview neural surface reconstruction by disentangling geometry and appearance. In *NeurIPS*, 2020. 1, 2, 3

[45] Zehao Yu and Shenghua Gao. Fast-mvsnet: Sparse-to-dense multi-view stereo with learned propagation and gauss-newton refinement. In *CVPR*, pages 1949–1958, 2020. 2

[46] Zehao Yu, Songyou Peng, Michael Niemeyer, Torsten Sattler, and Andreas Geiger. Monosdf: Exploring monocular geometric cues for neural implicit surface reconstruction. *arXiv preprint arXiv:2206.00665*, 2022. 2, 3

[47] Sergey Zagoruyko and Nikos Komodakis. Learning to compare image patches via convolutional neural networks. In *CVPR*, pages 4353–4361, 2015. 2

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