

# UNISURF: Unifying Neural Implicit Surfaces and Radiance Fields for Multi-View Reconstruction

Michael Oechsle<sup>1,2,3</sup> Songyou Peng<sup>1,4</sup> Andreas Geiger<sup>1,2</sup>

<sup>1</sup>Max Planck Institute for Intelligent Systems, Tübingen <sup>2</sup>University of Tübingen

<sup>3</sup>ETAS GmbH, Stuttgart <sup>4</sup>ETH Zurich

{firstname.lastname}@tue.mpg.de

## Abstract

Neural implicit 3D representations have emerged as a powerful paradigm for reconstructing surfaces from multi-view images and synthesizing novel views. Unfortunately, existing methods such as DVR or IDR require accurate per-pixel object masks as supervision. At the same time, neural radiance fields have revolutionized novel view synthesis. However, NeRF’s estimated volume density does not admit accurate surface reconstruction. Our key insight is that implicit surface models and radiance fields can be formulated in a unified way, enabling both surface and volume rendering using the same model. This unified perspective enables novel, more efficient sampling procedures and the ability to reconstruct accurate surfaces without input masks. We compare our method on the DTU, BlendedMVS, and a synthetic indoor dataset. Our experiments demonstrate that we outperform NeRF in terms of reconstruction quality while performing on par with IDR without requiring masks.

## 1. Introduction

Capturing the geometry and appearance of 3D scenes from a set of images is one of the cornerstone problems in computer vision. Towards this goal, coordinate-based neural models have emerged as a powerful tool for 3D reconstruction of geometry and appearance within the last years.

Many recent methods employ continuous implicit functions parameterized with neural networks as 3D representations of geometry [3, 8, 12, 32, 33, 37, 41, 43, 47, 57] or appearance [34, 38, 39, 40, 47, 52, 61]. These neural 3D representations have shown impressive performance on geometry reconstruction and novel view synthesis from multi-view images. Besides the choice of the 3D representation (e.g., occupancy field, unsigned or signed distance field), one key element for neural implicit multi-view reconstruction is the rendering technique. While some of these works represent the implicit surface as level set and hence render the appear-

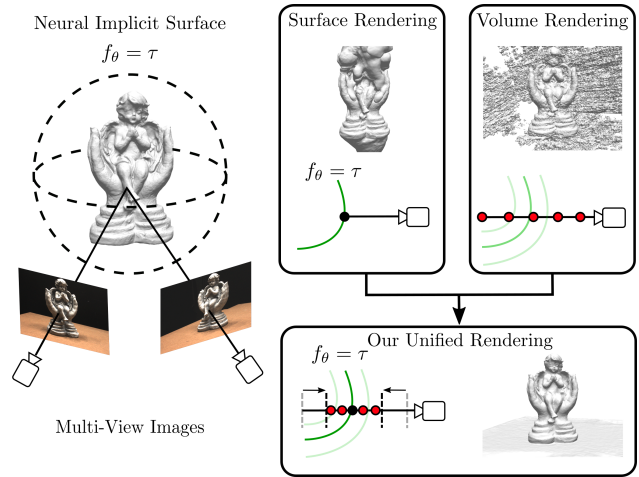


Figure 1: **Illustration.** Implicit models based on surface rendering [38, 61] require input masks and radiance fields [34] do not optimize implicit surfaces directly. UNISURF provides a principled unified formulation, enabling accurate surface reconstruction from images without input masks.

ance from surfaces [38, 52, 61], others integrate densities by drawing samples along the viewing rays [23, 34, 49].

In existing work, surface rendering techniques have shown impressive performance in 3D reconstruction [38, 61]. However, they require per-pixel object masks as input and an appropriate network initialization since surface rendering techniques only provide gradient information locally where a surface intersects with a ray. Intuitively speaking, optimizing wrt. local gradients can be seen as an iterative deformation procedure applied to an initial neural surface which is often initialized as a sphere. Additional constraints such as mask supervision are necessary for converging to a valid surface, see Fig. 2 for an illustration. Due to their reliance on masks, surface rendering methods are limited to object-level reconstruction and do not scale to larger scenes.

On contrary, volume rendering methods like NeRF [34]

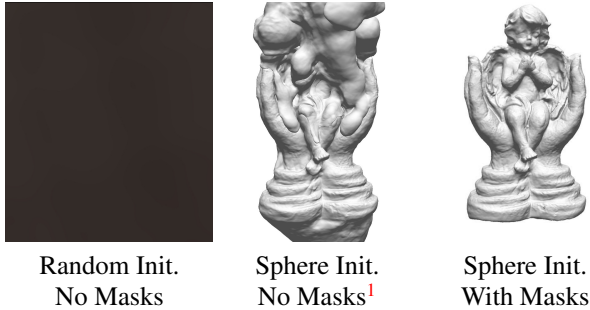


Figure 2: **Surface Rendering (IDR Results)**. State-of-the-art methods like IDR [61] require object masks and careful initialization for capturing accurate geometry.

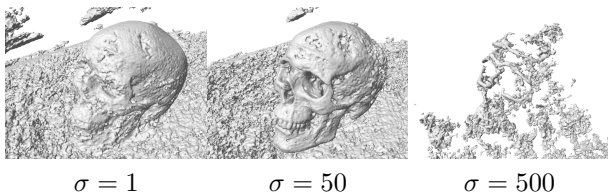


Figure 3: **Volume Rendering (NeRF Results)**. We show level sets of the recovered density volume for a trained NeRF model [34] using different density thresholds  $\sigma$ .

have shown impressive results for novel view synthesis, also for larger scenes. However, surfaces extracted as level sets of the underlying volume density are usually non-smooth and contain artifacts due to the flexibility of the radiance field representation which does not sufficiently constrain the 3D geometry in the presence of ambiguities, see Fig. 3.

**Contributions:** In this paper, we propose *UNISURF (Unified Neural Implicit Surface and Radiance Fields)* a principled unified framework for implicit surfaces and radiance fields, with the goal of reconstructing solid (i.e., non-transparent) objects from a set of RGB images. Our framework combines the benefits of surface rendering with those of volume rendering, enabling the reconstruction of accurate geometry from multi-view images without masks. By recovering implicit surfaces, we are able to gradually decrease the sampling region for volume rendering during optimization. Starting with large sampling regions enables capturing coarse geometry and resolving ambiguities during early iterations. At a later stage, we draw samples closer to the surface which improves reconstruction accuracy. We show that our approach enables capturing accurate geometry without mask supervision on the *DTU MVS* dataset [1], attaining results competitive with state-of-the-art implicit neural reconstruction methods like IDR [61] which use strong mask supervision. Moreover, we also demonstrate

<sup>1</sup>We remark that this result w/o masks from IDR [61] was the only successful case out of several failure cases.

our method on scenes from the *BlendedMVS* dataset [60] as well as synthetic indoor scenes from SceneNet [30].

## 2. Related Work

In this section, we first discuss related work from the domain of 3D reconstruction from multi-view images. Next, we provide an overview of recent works on neural implicit representations as well as differentiable rendering.

**3D Reconstruction from Multi-View Images:** Reconstructing 3D geometry from multiple images has been a longstanding computer vision problem [14]. Before the era of deep learning, classic multi-view stereo (MVS) methods [2, 4, 5, 7, 21, 21, 48, 50, 51] focus on either matching features across views [4, 48] or representing shapes with a voxel grid [2, 5, 7, 21, 28, 42, 50, 54, 55]. The former approaches usually have a complex pipeline requiring additional steps like fusing depth information [9, 31] and meshing [18, 19], while the latter ones are limited to low resolution due to cubic memory requirements. In contrast, neural implicit representations for 3D reconstruction do not suffer from discretization artifacts as they represent surfaces by the level set of a neural network with continuous outputs.

Recent learning-based MVS methods attempt to replace some parts of the classic MVS pipeline. For instance, some works learn to match 2D features [15, 22, 27, 56, 63], fuse depth maps [11, 46], or infer depth maps [16, 58, 59] from multi-view images. Contrary to these learning-based MVS approaches, our method only requires weak 2D supervision during optimization. Moreover, our method yields high-quality 3D geometry *and* synthesizes photorealistic and consistent novel views.

**Neural Implicit Representations:** Recently, neural implicit functions have emerged as an effective representation of 3D geometry [3, 8, 12, 32, 33, 37, 41, 43, 47, 57] and appearance [23, 25, 34, 38, 39, 40, 47, 49, 52] as they represent 3D content continuously and without discretization while simultaneously having a small memory footprint. Most of these methods require 3D supervision. However, several recent works [20, 24, 34, 38, 52, 61] demonstrated differentiable rendering for training directly from images [24, 34, 38, 52, 61]. We divide these methods into two groups: surface rendering and volume rendering.

Surface rendering approaches, including DVR [38], IDR [61] and NLR [20], determine the radiance directly on the surface of an object and provide a differentiable rendering formulation using implicit gradients. This allows for optimizing neural implicit surfaces from multi-view images. Conditioning on the viewing direction allows IDR to capture a high level of detail, even in the presence of non-lambertian surfaces. For real-time photorealistic novel view synthesis, NLR [20] extracts the surface as a mesh and exports the appearance in a Lumigraph representation.

However, DVR, IDR and NLR require pixel-accurate object masks for all views as input. In contrast, our method leads to similar reconstructions without requiring masks.

NeRF [34] and follow-ups [6, 29, 35, 36, 44, 45, 49, 53, 64] use volume rendering by learning alpha-compositing of a radiance field along rays. This method has shown impressive results on novel view synthesis and does not require mask supervision. However, the recovered 3D geometry is far from satisfactory, see Fig. 3. Several follow-up works (Neural Body [44] D-NeRF [45] and NeRD [6]) extract meshes using the volume density from NeRF, but none of them considers optimizing surfaces directly. Unlike these works, we aim at capturing accurate geometry and propose a volume rendering formulation that provably approaches surface rendering in the limit.

### 3. Background

The two main ingredients for learning neural implicit 3D representations from multi-view images are the 3D representation and the rendering technique linking the 3D representation and the 2D observations. This section provides the relevant background on implicit surface and volumetric radiance representations which we unify in this paper for the case of solid (non-transparent) objects and scenes.

**Implicit Surface Models:** Occupancy Networks [32, 38] represent surfaces as the decision boundary of a binary occupancy classifier, parameterized by a neural network

$$o_\theta(\mathbf{x}) : \mathbb{R}^3 \rightarrow [0, 1] \quad (1)$$

where  $\mathbf{x} \in \mathbb{R}^3$  is a 3D point and  $\theta$  are the model parameters. The surface is defined as the set of all 3D points where the occupancy probability is one half:  $\mathcal{S} = \{\mathbf{x}_s | o_\theta(\mathbf{x}_s) = 0.5\}$ . To associate a color with every 3D point  $\mathbf{x}_s$  on the surface, a color field  $c_\theta(\mathbf{x}_s)$  can be learned jointly with the occupancy field  $o_\theta(\mathbf{x})$ . The color for a particular pixel/ray  $\mathbf{r}$  is thus predicted as follows

$$\hat{C}(\mathbf{r}) = c_\theta(\mathbf{x}_s) \quad (2)$$

where  $\mathbf{x}_s$  is retrieved by root finding along ray  $\mathbf{r}$ , see [38] for details. The parameters<sup>2</sup>  $\theta$  of the occupancy field  $o_\theta$  and the color field  $c_\theta$  are determined by optimizing a reconstruction loss via gradient descent as described in [25, 38, 61].

While surface rendering allows for accurately estimating geometry and appearance, existing approaches strongly rely on supervision with object masks as surface rendering methods are only able to reason about rays that intersect a surface.

**Volumetric Radiance Models:** In contrast to implicit surface models, NeRF [34] represents scenes as colored volume densities and integrates radiance along rays via alpha

<sup>2</sup>For convenience, we use the same symbol  $\theta$  for all model parameters.

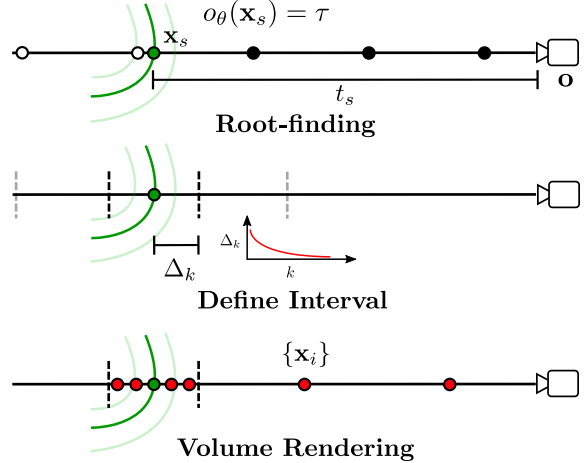


Figure 4: **Concept and Notation.** Our rendering consists of two steps: First, we seek the surface  $\mathbf{x}_s$  (green) in the occupancy field  $o_\theta$ . Second, we define an interval around the surface to sample points  $\{\mathbf{x}_i\}$  (red) for volume rendering.

blending [26, 34]. More specifically, NeRF uses a neural network to map a 3D location  $\mathbf{x} \in \mathbb{R}^3$  and a viewing direction  $\mathbf{d} \in \mathbb{R}^3$  to a volume density  $\sigma_\theta(\mathbf{x}) \in \mathbb{R}^+$  and a color value  $c_\theta(\mathbf{x}, \mathbf{d}) \in \mathbb{R}^3$ . Conditioning on the viewing direction  $\mathbf{d}$  allows for modeling view-dependent effects such as specular reflections [34, 40] and improves reconstruction quality in case the Lambertian assumption is violated [61]. Let  $\mathbf{o}$  denote the location of the camera center. Given  $N$  samples  $\{\mathbf{x}_i\}$  along ray  $\mathbf{r} = \mathbf{o} + t\mathbf{d}$ , NeRF approximates the color of pixel/ray  $\mathbf{r}$  using numerical quadrature:

$$\hat{C}(\mathbf{r}) = \sum_{i=1}^N T_i (1 - \exp(-\sigma_\theta(\mathbf{x}_i) \delta_i)) c_\theta(\mathbf{x}_i, \mathbf{d}) \quad (3)$$

$$T_i = \exp\left(-\sum_{j<i} \sigma_\theta(\mathbf{x}_j) \delta_j\right) \quad (4)$$

Here,  $T_i$  is the accumulated transmittance along the ray and  $\delta_i = |\mathbf{x}_{i+1} - \mathbf{x}_i|$  is the distance between adjacent samples. As Eq. (3) is differentiable, the parameters  $\theta$  of the density field  $\sigma_\theta$  and the color field  $c_\theta$  can be estimated by optimizing a reconstruction loss. We refer to [34] for details.

While NeRF does not require object masks for training due to its volumetric radiance representation, extracting the scene geometry from the volume density requires careful tuning of the density threshold and leads to artifacts due to the ambiguity present in the density field, see Fig. 3.

### 4. Method

We now describe our main contribution. In contrast to NeRF which is also applicable to non-solid scenes (e.g., fog, smoke), we restrict our focus to solid objects that can be

represented by 3D surfaces and view-dependent surface colors. Our method exploits both, the power of volumetric radiance representations to learn coarse scene structure without mask supervision as well as surface rendering which acts as an inductive bias to represent objects by a set of precise 3D surfaces, leading to accurate reconstructions.

#### 4.1. Unifying Surface and Volume Rendering

We start by noting that Eq. (3) can be rewritten as<sup>3</sup>

$$\hat{C}(\mathbf{r}) = \sum_{i=1}^N \alpha_i(\mathbf{x}_i) \prod_{j<i} (1 - \alpha_j(\mathbf{x}_j)) c(\mathbf{x}_i, \mathbf{d}) \quad (5)$$

with alpha values  $\alpha_i(\mathbf{x}) = 1 - \exp(-\sigma(\mathbf{x}) \delta_i)$ . Assuming solid objects,  $\alpha$  becomes a discrete occupancy indicator variable  $o \in \{0, 1\}$  which either takes  $o = 0$  in free space and  $o = 1$  in occupied space as value:

$$\hat{C}(\mathbf{r}) = \sum_{i=1}^N o(\mathbf{x}_i) \prod_{j<i} (1 - o(\mathbf{x}_j)) c(\mathbf{x}_i, \mathbf{d}) \quad (6)$$

We recognize this expression as the image formation model for solid objects [55] where the term  $o(\mathbf{x}_i) \prod_{j<i} (1 - o(\mathbf{x}_j))$  evaluates to 1 for the first occupied sample  $\mathbf{x}_i$  along ray  $\mathbf{r}$  and to 0 for all other samples.  $\prod_{j<i} (1 - o(\mathbf{x}_j))$  is an indicator for visibility which is 1 if there exists no occupied sample  $\mathbf{x}_j$  with  $j < i$  before sample  $\mathbf{x}_i$ . Thus,  $\hat{C}(\mathbf{r})$  takes the color  $c(\mathbf{x}_i, \mathbf{d})$  of the first occupied sample along ray  $\mathbf{r}$ .

To unify implicit surface and volumetric radiance models, we parameterize  $o$  directly by a continuous occupancy field  $o_\theta$  (1) as opposed to predicting volume density  $\sigma$ . Following [61], we condition the color field  $c_\theta$  on the surface normal  $\mathbf{n}$  and a feature vector  $\mathbf{h}$  of the geometry network which empirically induces a useful bias as also observed in [61] for the case of implicit surfaces. Importantly, our unified formulation allows for both volume and surface rendering

$$\hat{C}_v(\mathbf{r}) = \sum_{i=1}^N o_\theta(\mathbf{x}_i) \prod_{j<i} (1 - o_\theta(\mathbf{x}_j)) c_\theta(\mathbf{x}_i, \mathbf{n}_i, \mathbf{h}_i, \mathbf{d}) \quad (7)$$

$$\hat{C}_s(\mathbf{r}) = c_\theta(\mathbf{x}_s, \mathbf{n}_s, \mathbf{h}_s, \mathbf{d}) \quad (8)$$

where  $\mathbf{x}_s$  is retrieved by root-finding along ray  $\mathbf{r}$  and  $\mathbf{n}_s, \mathbf{h}_s$  denote the normal and geometry features at  $\mathbf{x}_s$ , respectively. Note that  $\mathbf{x}_s$  depends on the occupancy field  $o_\theta$ , but we have dropped this dependency here for clarity. For further details, we refer the reader to the supplementary material.

The advantage of this unified formulation is that it allows for both rendering on the surface directly and rendering throughout the entire volume which enables gradually removing ambiguities during optimization. As evidenced

<sup>3</sup>We drop dependencies on the model parameters for clarity.

by our experiments, it is indeed critical to combine both for obtaining accurate reconstructions without mask supervision. Being able to quickly recover the surface  $\mathcal{S}$  via root-finding enables more effective volume rendering, successively focusing on and refining the object surfaces as we will describe in Section 4.3. Furthermore, surface rendering enables faster novel view synthesis as illustrated in Fig. 5.

#### 4.2. Loss Function

We optimize the following regularized loss function

$$\mathcal{L} = \mathcal{L}_{rec} + \lambda \mathcal{L}_{reg} \quad (9)$$

with  $\ell_1$  reconstruction loss and  $\ell_2$  surface regularization which encourages the normal of a surface point  $\mathbf{x}_s$  and a point sampled in its neighborhood to be similar:

$$\mathcal{L}_{rec} = \sum_{\mathbf{r} \in \mathcal{R}} \|\hat{C}_v(\mathbf{r}) - C(\mathbf{r})\|_1 \quad (10)$$

$$\mathcal{L}_{reg} = \sum_{\mathbf{x}_s \in \mathcal{S}} \|\mathbf{n}(\mathbf{x}_s) - \mathbf{n}(\mathbf{x}_s + \epsilon)\|_2 \quad (11)$$

Here,  $\mathcal{R}$  denotes the set of all pixels/rays in the minibatch,  $\mathcal{S}$  is the set of corresponding surface points,  $C(\mathbf{r})$  is the observed color for pixel/ray  $\mathbf{r}$  and  $\epsilon$  is a small random uniform 3D perturbation. The normal at  $\mathbf{x}_s$  is given by

$$\mathbf{n}(\mathbf{x}_s) = \frac{\nabla_{\mathbf{x}_s} o_\theta(\mathbf{x}_s)}{\|\nabla_{\mathbf{x}_s} o_\theta(\mathbf{x}_s)\|_2} \quad (12)$$

which can be computed using double backpropagation [38].

#### 4.3. Optimization

The key hypothesis of implicit surface models [38,61] is that only the region at the first intersection point with the surface contributes to the rendering equation. However, this assumption is not true during early iterations where the surface is not well defined. Consequently, existing methods [38,61] require strong mask supervision. Conversely, during later iterations, knowledge of the approximate surface is valuable for drawing informative samples when evaluating the volume rendering equation in Eq. (7). Therefore, we utilize a training schedule with a monotonically decreasing sampling interval for drawing samples during volume rendering, as visualized in Fig. 4. In other words, during early iterations, the samples  $\{\mathbf{x}_i\}$  cover the entire optimization volume, effectively bootstrapping the reconstruction process using volume rendering. During later iterations, the samples  $\{\mathbf{x}_i\}$  are drawn closer around the estimated surface. As the surface can be estimated directly from the occupancy field  $o_\theta$  via root-finding [38], this eliminates the need for hierarchical two-stage sampling as in NeRF. Our experiments demonstrate that this procedure is particularly effective for estimating accurate geometry, while it allows for resolving ambiguities during early iterations.

More formally, let  $\mathbf{x}_s = \mathbf{o} + t_s \mathbf{d}$ . We obtain samples  $\mathbf{x}_i = \mathbf{o} + t_i \mathbf{d}$  by drawing  $N$  depth values  $t_i$  using stratified sampling within the interval  $[t_s - \Delta, t_s + \Delta]$  centered at  $t_s$ :

$$t_i \sim \mathcal{U} \left[ t_s + \left( \frac{2i-2}{N} - 1 \right) \Delta, t_s + \left( \frac{2i}{N} + 1 \right) \Delta \right] \quad (13)$$

During training, we start with a large interval  $\Delta_{\max}$  and gradually decrease  $\Delta$  for more accurate sampling and optimization of the surface using the following decay schedule

$$\Delta_k = \max(\Delta_{\max} \exp(-k\beta), \Delta_{\min}) \quad (14)$$

where  $k$  denotes the iteration number and  $\beta$  is a hyperparameter. In fact, it can be shown that for  $\Delta \rightarrow 0$  and  $N \rightarrow \infty$ , volume rendering (7) indeed approaches surface rendering (8):  $\hat{C}_v(\mathbf{r}) \rightarrow \hat{C}_s(\mathbf{r})$ . A formal proof of this limit is provided in the supplementary material.

As evidenced by our experiments, the decay schedule in (14) is critical for capturing detailed geometry as it combines volume rendering of large and uncertain volumes in the beginning of training with surface rendering towards the end of training. To reduce free space artifacts, we combine these samples with points sampled randomly between the camera and the surface. For rays without surface intersection, we use stratified sampling on the entire ray.

#### 4.4. Implementation Details

**Architecture:** Similar to Yariv et al. [61], we use an 8-layer MLP with a Softplus activation function and a hidden dimension of 256 for the occupancy field  $o_\theta$ . We initialize the network such that the decision boundary is a sphere [13]. In contrast, the radiance field  $c_\theta$  is parameterized as a 4-layer ReLU MLP. We encode the 3D location  $\mathbf{x}$  and the viewing direction  $\mathbf{d}$  using Fourier features [34] at  $k$  octaves. We empirically found  $k = 6$  for the 3D location  $\mathbf{x}$  and  $k = 4$  for the viewing direction  $\mathbf{d}$  to work best.

**Optimization:** In all experiments, we fit our model to multi-view images of a single scene. During optimization of the model parameters, we first randomly sample a view and then  $M$  pixels/rays  $\mathcal{R}$  from this view based on the camera intrinsics and extrinsics. Next, we render all rays to compute the loss function in Eq. (9). For root-finding, we use 256 uniform sampled points and apply the secant method with 8 steps [32]. For our rendering procedure, we use  $N = 64$  query points inside the interval and 32 in the free space between the camera and the lower bound of the interval. The interval decay parameters are  $\beta = 1.5e - 5$ ,  $\Delta_{\min} = 0.05$  and  $\Delta_{\max} = 1.0$ . We use Adam with a learning rate of 0.0001 and optimize for  $M = 1024$  pixels per iteration with two decay steps after 200k and 400k iterations. In total, we train our models for 450k iterations.

<sup>4</sup>We used a single NVIDIA V100 GPU to measure the inference time.

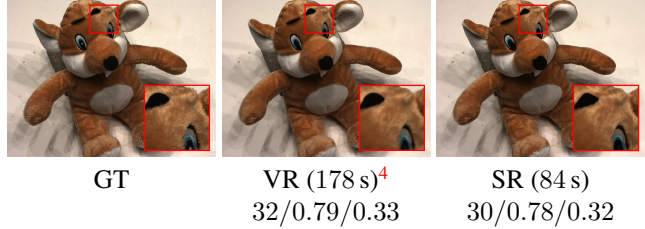


Figure 5: **Volume vs. Surface Rendering.** We compare images generated using volume rendering (VR,  $\Delta = \Delta_{\min}$ ) and surface rendering (SR), reporting three image metrics (PSNR $\uparrow$ /SSIM $\uparrow$ /LPIPS $\downarrow$ ). Both approaches yield similar results while surface rendering is twice as fast.

**Inference:** Our method allows for inferring 3D shapes as well as for synthesizing novel view images. For synthesizing images, we can render our representation in two different ways, we can either use volume rendering or surface rendering. In Fig. 5, we show that both rendering approaches lead to similar results. However, we observe that surface rendering is faster than volume rendering.

To extract meshes, we apply the Multiresolution IsoSurface Extraction (MISE) algorithm from [32]. We use  $64^3$  as the initial resolution and up-sample the mesh in 3 steps without gradient-based refinement.

## 5. Experimental Evaluation

We conduct experiments on multi-view 3D reconstruction to assess our method. First, we provide a qualitative and quantitative comparison of our approach to existing methods (IDR [61], NeRF [34], Colmap [48]) on the widely used *DTU MVS* dataset [17]. Second, we show a qualitative comparison for samples from the *BlendedMVS* dataset [60] and synthetic renderings of scenes from the *SceneNet* dataset [30]. Third, we analyze our rendering procedure and loss function in an ablation study.

### 5.1. Baselines

To validate the effectiveness of our method, we compare it to three different baselines.

**Colmap [48]:** We consider Colmap [48] as a classical MVS baseline, as it shows strong performance on multi-view reconstruction and is widely used in related works [38, 61]. We reconstruct a mesh from the output of Colmap using screened Poisson Surface reconstruction (sPSR) [19]. Following [38, 61], we show results for the quantitatively best trim parameter 7 and for the setting that results in watertight meshes (trim parameter 0).

**NeRF [34]:** Although NeRF targets novel view synthesis, its volume density admits the extraction of geometry. To extract meshes from NeRF, we define a density threshold of 50. We validate this choice in the supplementary.

**IDR [61]:** IDR is the state-of-the-art multi-view reconstruction method for neural implicit surfaces. IDR reconstructs surfaces with an impressively high level of detail and handles specular surfaces, but requires input masks. We do not compare to DVR [38] as IDR’s view-dependent modeling has been demonstrated to be superior to DVR, see [61].

## 5.2. Datasets

For our investigations, we use three different datasets.

**DTU MVS Dataset [17]:** The *DTU MVS* dataset contains 49 to 64 images at a resolution of  $1200 \times 1600$  as well as extrinsic and intrinsic camera parameters for all views. The dataset consists of objects with different shapes and appearances. Non-lambertian appearance effects make some of the objects particularly challenging. For each scan, ground truth 3D shapes, as well as the official evaluation procedure, are available<sup>5</sup>. Like previous works [38, 61], we use the “Surface” method of the evaluation script, and evaluate all methods on meshes cleaned with the respective masks. The official evaluation procedure calculates the Chamfer distance between sampled points of the predicted shapes and ground truth shapes provided in the dataset. For evaluating IDR [61], we use the pixel-accurate masks for all images which are provided by the authors of IDR.

**BlendedMVS Dataset [60]:** The *BlendedMVS* dataset is a large-scale dataset containing multi-view images with respective camera extrinsics and intrinsics. We use examples from the *BlendedMVS* low-res set with an image resolution of  $768 \times 576$ . These examples contain 24 to 64 different views of unmasked images. We define the scene as such that the object is in the center and the closest camera lies near the unit sphere.

**SceneNet Dataset [30]:** For testing our model on complex indoor scenes, we take two scenes from [30] for evaluation. *BlenderProc* [10] is applied for rendering images of a part of the scenes containing multiple objects. The first scene is a bedroom scene with a bed, a lamp and a bedside table. The other scene is a living room scene with a sofa, a curtain and a round table. We use 83 and 40 images, respectively.

## 5.3. Comparison on DTU

In Table 1, we quantitatively compare our method to the baselines on the *DTU MVS* dataset. While the Colmap baseline with trim parameter  $\zeta = 7$  shows the best performance on the Chamfer distance, it produces non-watertight meshes with incomplete surfaces. Our method performs nearly on par with the state-of-the-art neural implicit model IDR while not relying on strong mask supervision. NeRF and Colmap ( $\zeta = 0$ ) also do not use input masks but exhibit worse performance in terms of Chamfer distance.

<sup>5</sup><https://roboimagedata.compute.dtu.dk>

	Colmap	IDR	Colmap	NeRF	Ours
masks	<b>X</b>	<b>✓</b>	<b>X</b>	<b>X</b>	<b>X</b>
trim	7	-	0	-	-
scan24	0.45	1.63	<b>0.81</b>	1.90	1.32
scan37	0.91	1.87	2.05	1.60	<b>1.36</b>
scan40	0.37	0.63	<b>0.73</b>	1.85	1.72
scan55	0.37	0.48	1.22	0.58	<b>0.44</b>
scan63	0.90	1.04	1.79	2.28	<b>1.35</b>
scan65	1.00	0.79	1.58	1.27	<b>0.79</b>
scan69	0.54	0.77	1.02	1.47	<b>0.80</b>
scan83	1.22	1.33	3.05	1.67	<b>1.49</b>
scan97	1.08	1.16	1.40	2.05	<b>1.37</b>
scan105	0.64	0.76	2.05	1.07	<b>0.89</b>
scan106	0.48	0.67	1.00	0.88	<b>0.59</b>
scan110	0.59	0.90	<b>1.32</b>	2.53	1.47
scan114	0.32	0.42	0.49	1.06	<b>0.46</b>
scan118	0.45	0.51	0.78	1.15	<b>0.59</b>
scan122	0.43	0.53	1.17	0.96	<b>0.62</b>
mean	0.65	0.90	1.36	1.49	<b>1.02</b>

Table 1: **Quantitative Comparison on DTU [17].** We show a quantitative comparison against the baselines for the reconstructed geometry on 15 scans from the DTU dataset. Our method performs on par with IDR, even though IDR uses masks as input. Bold numbers are only considering methods without mask supervision and trimming.

In Fig. 6, we show qualitative results for our method and the baselines. While Colmap provides detailed reconstructions, it leads to incomplete geometry due to trimming. For NeRF, holes and noise artifacts can be observed in the reconstructions. In contrast, our approach and IDR (with input masks) produce accurate surfaces with high-quality details. We remark that our model captures the overall spatial arrangement of the scene accurately while being able to also capture geometric details, e.g., the teeth of the skull and other surface details.

## 5.4. Comparison on BlendedMVS and SceneNet

To show our model’s capabilities on more diverse scenes, we use samples of the *BlendedMVS* dataset and indoor scenes from *SceneNet*. As there exist no object masks for these scenes, we only consider Colmap and NeRF as baseline methods. While we have tried running IDR, none of the scenes converged, resulting in degenerate outputs. For scenes with complex backgrounds, we use a background model that learns to capture appearance information outside of our region of interest, we refer the reader to the supplementary material for more details.

Our qualitative results in Fig. 8 provide evidence that, unlike existing implicit surface models, our method is able

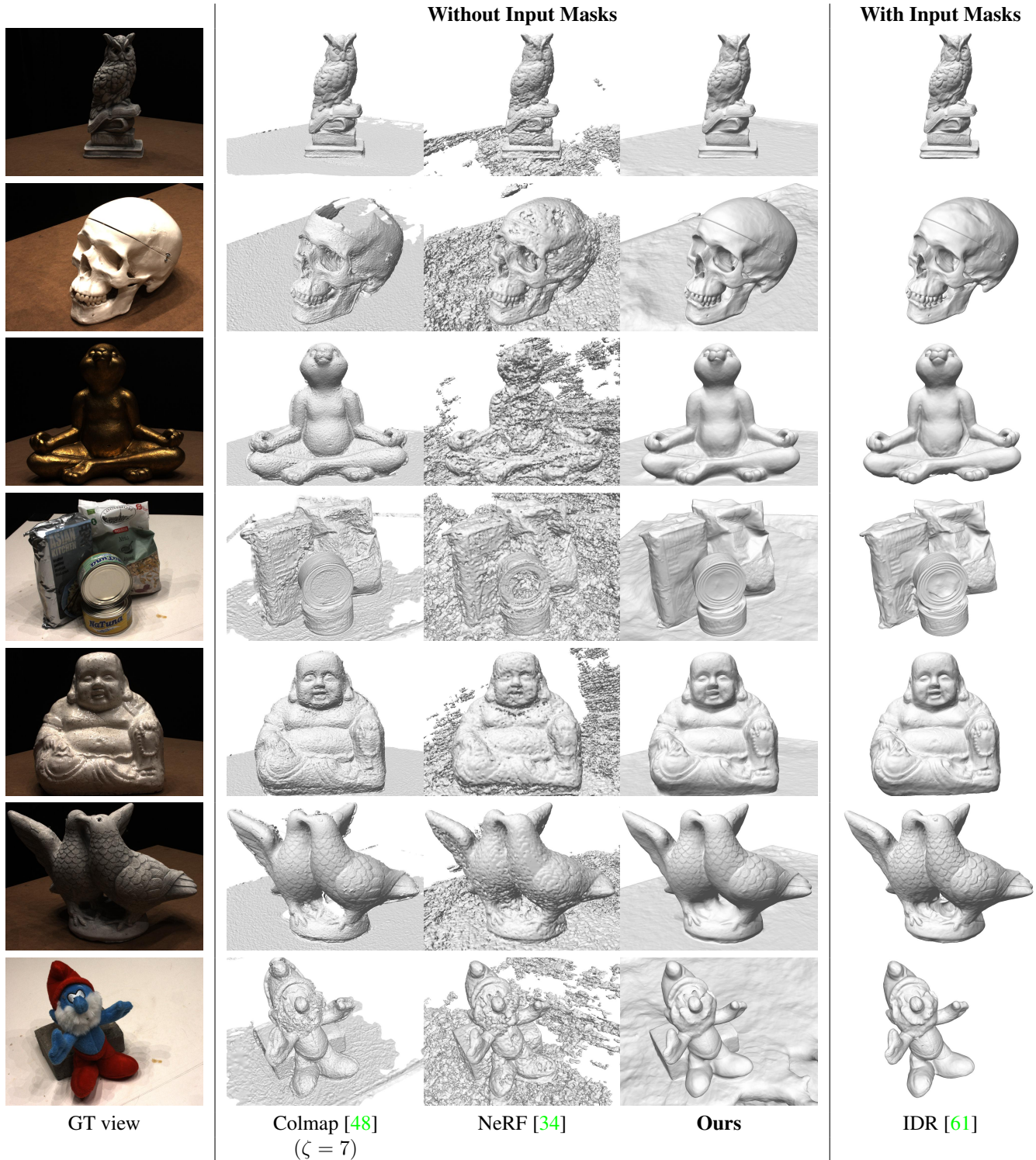


Figure 6: **Reconstructed 3D Shapes for DTU.** We show a qualitative comparison of reconstructed surfaces. We use Colmap [48], NeRF [34] and IDR [61] as baselines. While IDR requires pixel-accurate object masks for all images, the other methods reason about the full scene extent. Our method performs visually on par with IDR while not requiring masks.

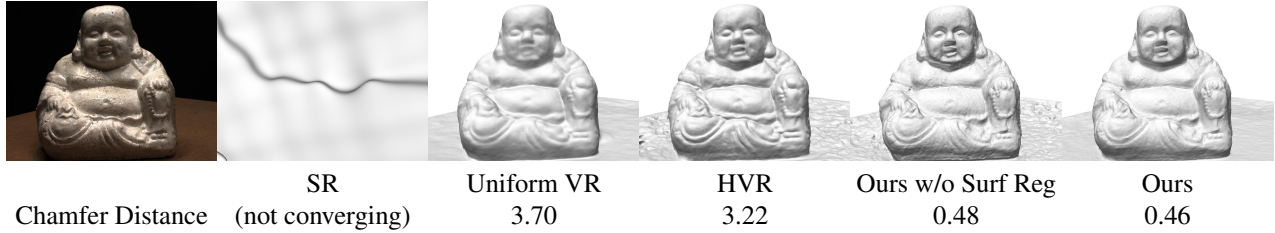


Figure 7: **Ablation Study.** We compare our method to different rendering procedures: Surface rendering in isolation (SR), volume rendering with uniformly sampled points (Uniform VR) and hierarchical volume rendering (HVR). This analysis on the Buddha scene (DTU) shows that our full method performs best, both visually and quantitatively. Moreover, our model without surface regularizer (Ours w/o Surf Reg) leads to non-smooth regions, especially at the textureless flat table.



Figure 8: **Reconstructing Indoor Scenes [30] and BlendedMVS [60].** We show a qualitative comparison of reconstructed surfaces on indoor scenes (1st and 2nd row) and *BlendedMVS* (3rd, 4th and 5th row).

to reconstruct plausible geometry for complex scenes with multiple objects and backgrounds. While Colmap works well on indoor scenes, it shows artifacts on uniformly colored regions (e.g., round table in the second row). As for the *BlendedMVS* experiment, NeRF can reason about the overall spatial structure but shows less accurate surfaces with significantly higher levels of noise compared to UNISURF. More results can be found in the supplementary material.

### 5.5. Ablation study

We finally investigate the impact of rendering design choices and show ablations of the loss function.

**Rendering Procedure:** We argue for our choices of the rendering procedure by comparing different variations in Fig. 7: First, we consider a baseline which only uses surface rendering during optimization (SR). We use an  $\ell_1$  reconstruction loss on the surface color and backpropagate through implicit differentiation following [38]. A second baseline uses uniform volume rendering with 96 query points (Uniform VR). Third, we use NeRF’s hierarchical volume sampling (HVR) [34] with 64+64 sample points for querying the occupancy field. While surface rendering does not converge, Uniform VR and HVR result in overly smooth and bloated shapes with missing details. This provides evidence that the proposed unified model leads to more accurate reconstructions compared to the baselines.

**Losses:** In Fig. 7, we also show an ablation study of our surface regularization term in (9). Without this regularizer surfaces become less smooth in flat and ambiguous areas, e.g., at the table. This regularization term is particularly useful for regions that are observed less frequently, as it incorporates an inductive bias towards smooth surfaces.

### 6. Conclusion and Future Work

This work presents UNISURF, a unified formulation of implicit surfaces and radiance fields for capturing high-quality implicit surface geometry from multi-view images without input masks. We believe that neural implicit surfaces and advanced differentiable rendering procedures play a key role in future 3D reconstruction methods. Our unified formulation shows a path towards optimizing implicit surfaces in a more general setting than possible before.

As true for MVS methods in general, rarely visible and texture-less regions are hard to reconstruct. Thus, a prior becomes necessary to resolve ambiguities. While we incorporate an explicit smoothness prior during optimization,



learning a probabilistic neural surface model which captures regularities and uncertainty across objects would help to resolve ambiguities, leading to more accurate reconstructions.

**Acknowledgement** This work was supported by an NVIDIA research gift. Andreas Geiger was supported by the ERC Starting Grant LEGO-3D (850533) and the DFG EXC number 2064/1 - project number 390727645. The authors thank the Max Planck ETH Center for Learning Systems (CLS) for supporting Songyou Peng.

## References

- [1] Henrik Aanæs, Rasmus Ramsbøl Jensen, George Vogiatzis, Engin Tola, and Anders Bjarholm Dahl. Large-scale data for multiple-view stereopsis. *International Journal of Computer Vision (IJCV)*, 2016. 2
- [2] Motilal Agrawal and Larry S Davis. A probabilistic framework for surface reconstruction from multiple images. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2001. 2
- [3] Matan Atzmon, Niv Haim, Lior Yariv, Ofer Israelov, Haggai Maron, and Yaron Lipman. Controlling neural level sets. In *Advances in Neural Information Processing Systems (NIPS)*, 2019. 1, 2
- [4] Michael Bleyer, Christoph Rhemann, and Carsten Rother. Patchmatch stereo - stereo matching with slanted support windows. In *Proc. of the British Machine Vision Conf. (BMVC)*, 2011. 2
- [5] JS De Bonet and Paul Viola. Poxels: Probabilistic voxelized volume reconstruction. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 1999. 2
- [6] Mark Boss, Raphael Braun, Varun Jampani, Jonathan T Barron, Ce Liu, and Hendrik Lensch. Nerf: Neural reflectance decomposition from image collections. *arXiv.org*, 2020. 3
- [7] A Broadhurst, T W Drummond, and R Cipolla. A probabilistic framework for space carving. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2001. 2
- [8] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2018. 1, 2
- [9] Brian Curless and Marc Levoy. A volumetric method for building complex models from range images. In *ACM Trans. on Graphics*, 1996. 2
- [10] Maximilian Denninger, Martin Sundermeyer, Dominik Winkelbauer, Youssef Zidan, Dmitry Olefir, Mohamad Elbadrawy, Ahsan Lodhi, and Harinandan Katam. Blenderproc. *arXiv.org*, 2019. 6
- [11] Simon Donne and Andreas Geiger. Learning non-volumetric depth fusion using successive rejections. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [12] Kyle Genova, Forrester Cole, Daniel Vlasic, Aaron Sarna, William T Freeman, and Thomas Funkhouser. Learning shape templates with structured implicit functions. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2019. 1, 2
- [13] Amos Gropp, Lior Yariv, Niv Haim, Matan Atzmon, and Yaron Lipman. Implicit geometric regularization for learning shapes. In *Proc. of the International Conf. on Machine Learning (ICML)*, 2020. 5
- [14] Richard Hartley and Andrew Zisserman. *Multiple View Geometry in Computer Vision*. Cambridge university press, 2003. 2
- [15] Wilfried Hartmann, Silvano Galliani, Michal Havlena, Luc Van Gool, and Konrad Schindler. Learned multi-patch similarity. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2017. 2
- [16] Po-Han Huang, Kevin Matzen, Johannes Kopf, Narendra Ahuja, and Jia-Bin Huang. Deepmvs: Learning multi-view stereopsis. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2
- [17] Rasmus Ramsbøl Jensen, Anders Lindbjerg Dahl, George Vogiatzis, Engin Tola, and Henrik Aanæs. Large scale multi-view stereopsis evaluation. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2014. 5, 6
- [18] Michael M. Kazhdan, Matthew Bolitho, and Hugues Hoppe. Poisson surface reconstruction. In *Proceedings of the Fourth Eurographics Symposium on Geometry Processing*, 2006. 2
- [19] Michael M. Kazhdan and Hugues Hoppe. Screened poisson surface reconstruction. *ACM Trans. on Graphics*, 2013. 2, 5
- [20] Petr Kellnhofer, Lars Jebe, Andrew Jones, Ryan Spicer, Kari Pulli, and Gordon Wetzstein. Neural lumigraph rendering. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [21] Kiriakos N. Kutulakos and Steven M. Seitz. A theory of shape by space carving. *International Journal of Computer Vision (IJCV)*, 2000. 2
- [22] Vincent Leroy, Jean-Sébastien Franco, and Edmond Boyer. Shape reconstruction using volume sweeping and learned photoconsistency. In *Proc. of the European Conf. on Computer Vision (ECCV)*, 2018. 2
- [23] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 1, 2
- [24] Shichen Liu, Shunsuke Saito, Weikai Chen, and Hao Li. Learning to infer implicit surfaces without 3d supervision. In *Advances in Neural Information Processing Systems (NIPS)*, 2019. 2
- [25] 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 *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2, 3
- [26] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural volumes: Learning dynamic renderable volumes from images. In *ACM Trans. on Graphics*, 2019. 3
- [27] W. Luo, A. Schwing, and R. Urtasun. Efficient deep learning for stereo matching. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2

- [28] David Marr and Tomaso Poggio. Cooperative computation of stereo disparity. *Science*, 1976. 2
- [29] Ricardo Martin-Brualla, Noha Radwan, Mehdi S. M. Sajjadi, Jonathan T. Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021. 3
- [30] John McCormac, Ankur Handa, Stefan Leutenegger, and Andrew J. Davison. Scenenet RGB-D: can 5m synthetic images beat generic imagenet pre-training on indoor segmentation? In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2017. 2, 5, 6, 8
- [31] 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 *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2007. 2
- [32] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019. 1, 2, 3, 5
- [33] Mateusz Michalkiewicz, Jhony K Pontes, Dominic Jack, Mahsa Baktashmotlagh, and Anders Eriksson. Implicit surface representations as layers in neural networks. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2019. 1, 2
- [34] 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 *Proc. of the European Conf. on Computer Vision (ECCV)*, 2020. 1, 2, 3, 5, 7, 8, 17, 19
- [35] Thomas Neff, Pascal Stadlbauer, Mathias Parger, Andreas Kurz, Chakravarty R. Alla Chaitanya, Anton Kaplanyan, and Markus Steinberger. Donerf: Towards real-time rendering of neural radiance fields using depth oracle networks. *arXiv.org*, 2021. 3
- [36] Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021. 3
- [37] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Occupancy flow: 4d reconstruction by learning particle dynamics. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2019. 1, 2
- [38] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020. 1, 2, 3, 4, 5, 6, 8, 14
- [39] Michael Oechsle, Lars Mescheder, Michael Niemeyer, Thilo Strauss, and Andreas Geiger. Texture fields: Learning texture representations in function space. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2019. 1, 2
- [40] Michael Oechsle, Michael Niemeyer, Christian Reiser, Lars Mescheder, Thilo Strauss, and Andreas Geiger. Learning implicit surface light fields. In *Proc. of the International Conf. on 3D Vision (3DV)*, 2020. 1, 2, 3
- [41] Jeong Joon Park, Peter Florence, Julian Straub, Richard A. Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019. 1, 2
- [42] Despoina Paschalidou, Ali Osman Ulusoy, Carolin Schmitt, Luc van Gool, and Andreas Geiger. Raynet: Learning volumetric 3d reconstruction with ray potentials. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2
- [43] Songyou Peng, Michael Niemeyer, Lars Mescheder, Marc Pollefeys, and Andreas Geiger. Convolutional occupancy networks. In *Proc. of the European Conf. on Computer Vision (ECCV)*, 2020. 1, 2
- [44] Sida Peng, Yuanqing Zhang, Yinghao Xu, Qianqian Wang, Qing Shuai, Hujun Bao, and Xiaowei Zhou. Neural body: Implicit neural representations with structured latent codes for novel view synthesis of dynamic humans. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020. 3
- [45] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. *arXiv.org*, 2020. 3
- [46] Gernot Riegler, Ali Osman Ulusoy, Horst Bischof, and Andreas Geiger. OctNetFusion: Learning depth fusion from data. In *Proc. of the International Conf. on 3D Vision (3DV)*, 2017. 2
- [47] 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 *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2019. 1, 2
- [48] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise view selection for unstructured multi-view stereo. In *Proc. of the European Conf. on Computer Vision (ECCV)*, 2016. 2, 5, 7, 8, 14, 17, 19
- [49] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. Graf: Generative radiance fields for 3d-aware image synthesis. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 1, 2, 3
- [50] S.M. Seitz and C.R. Dyer. Photorealistic scene reconstruction by voxel coloring. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 1997. 2
- [51] Steven M. Seitz, Brian Curless, James Diebel, Daniel Scharstein, and Richard Szeliski. A comparison and evaluation of multi-view stereo reconstruction algorithms. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2006. 2
- [52] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. In *Advances in Neural Information Processing Systems (NIPS)*, 2019. 1, 2
- [53] Pratul Srinivasan, Boyang Deng, Xiuming Zhang, Matthew Tancik, Ben Mildenhall, and Jonathan T. Barron. NeRV: Neural reflectance and visibility fields for relighting and view synthesis. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021. 3

- [54] Shubham Tulsiani, Tinghui Zhou, Alexei A Efros, and Jitendra Malik. Multi-view supervision for single-view reconstruction via differentiable ray consistency. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2
- [55] Ali Osman Ulusoy, Andreas Geiger, and Michael J. Black. Towards probabilistic volumetric reconstruction using ray potentials. In *Proc. of the International Conf. on 3D Vision (3DV)*, 2015. 2, 4
- [56] 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 *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2
- [57] Weiyue Wang, Xu Qiangeng, Duygu Ceylan, Radomir Mech, and Ulrich Neumann. DISN: Deep implicit surface network for high-quality single-view 3d reconstruction. *Advances in Neural Information Processing Systems (NIPS)*, 2019. 1, 2
- [58] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. *Proc. of the European Conf. on Computer Vision (ECCV)*, 2018. 2
- [59] Yao Yao, Zixin Luo, Shiwei Li, Tianwei Shen, Tian Fang, and Long Quan. Recurrent mvsnet for high-resolution multi-view stereo depth inference. *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [60] Yao Yao, Zixin Luo, Shiwei Li, Jingyang Zhang, Yufan Ren, Lei Zhou, Tian Fang, and Long Quan. Blendedmvs: A large-scale dataset for generalized multi-view stereo networks. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2, 5, 6, 8
- [61] 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 (NeurIPS)*, 2020. 1, 2, 3, 4, 5, 6, 7, 14, 17
- [62] Lin Yen-Chen. PyTorchNeRF: a PyTorch implementation of NeRF, 2020. 15
- [63] Sergey Zagoruyko and Nikos Komodakis. Learning to compare image patches via convolutional neural networks. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2015. 2
- [64] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv.org*, 2020. 3, 14

# Supplementary Materials

## Abstract

In this **supplementary document**, we first formally show that volume rendering converges to surface rendering in the limit of small intervals and a large number of query points (Section I). Next, we provide more information about the experimental setups for all datasets (Section II) and provide more details about the baselines (Section III). Next, we discuss the limiting factors of our method and possible solutions (Section IV). Finally, we present additional results on the DTU dataset, indoor scenes from SceneNet and the BlendenMVS dataset (Section V).

## I. Convergence Proof

We consider the volume and surface rendering equations

$$\hat{C}_v(\mathbf{r}) = \sum_{i=1}^N o_\theta(\mathbf{x}_i) \prod_{j<i} (1 - o_\theta(\mathbf{x}_j)) c_\theta(\mathbf{x}_i, \mathbf{n}_i, \mathbf{h}_i, \mathbf{d}) \quad (15)$$

$$\hat{C}_s(\mathbf{r}) = c_\theta(\mathbf{x}_s, \mathbf{n}_s, \mathbf{h}_s, \mathbf{d}) \quad (16)$$

with occupancy field  $o_\theta$  and color field  $c_\theta$ . The samples  $\mathbf{x}_i = \mathbf{o} + t_i \mathbf{d}$  are obtained by drawing  $N$  depth values  $t_i$  randomly from the interval  $[t_s - \Delta, t_s + \Delta]$  centered at  $t_s$ . For notation clarity, we will drop the model parameters  $\theta$  in the following. Without loss of generality, we will consider  $c$  as function of  $\mathbf{x}$  only as  $\mathbf{n}$  and  $\mathbf{h}$  are functions of the sample point  $\mathbf{x}$  and a single color channel, i.e.,  $c(\mathbf{x}) \in \mathbb{R}$ :

$$\hat{C}_v(\mathbf{r}) = \sum_{i=1}^N o(\mathbf{x}_i) \prod_{j<i} (1 - o(\mathbf{x}_j)) c(\mathbf{x}_i) \quad (17)$$

$$\hat{C}_s(\mathbf{r}) = c(\mathbf{x}_s) \quad (18)$$

We will now show that  $\hat{C}_v(\mathbf{r})$  approaches  $\hat{C}_s(\mathbf{r})$  for  $\Delta \rightarrow 0$  and  $N \rightarrow \infty$ .

**Theorem 1.** *Assuming Multi-layer Perceptrons for the occupancy and color fields with Softplus and ReLU activation functions, respectively, as well as Fourier features at  $k$  octaves, we have*

$$\lim_{\substack{\Delta \rightarrow 0 \\ N \rightarrow \infty}} \hat{C}_v(\mathbf{r}) = \hat{C}_s(\mathbf{r}) \quad (19)$$

Thus, volume and surface rendering become equivalent when reducing the interval and increasing the number of samples.

*Proof.* Linear layers, Softplus, Sigmoid and ReLU activation functions as well as Fourier features are Lipschitz continuous. Compositions of Lipschitz continuous functions are Lipschitz continuous. Thus,  $o(\mathbf{x})$  and  $c(\mathbf{x})$  are Lipschitz continuous wrt.  $\mathbf{x}$ . Let  $k_o$  and  $k_c$  denote the respective Lipschitz constants. As  $o(\mathbf{x}_s) = 0.5$ , we have

$$\begin{aligned} |o(\mathbf{x}_i) - 0.5| &\leq k_o \|\mathbf{x}_i - \mathbf{x}_s\|_2 = \xi \\ |c(\mathbf{x}_i) - c(\mathbf{x}_s)| &\leq k_c \|\mathbf{x}_i - \mathbf{x}_s\|_2 = \gamma \end{aligned}$$

with  $\xi \rightarrow 0$  and  $\gamma \rightarrow 0$  for  $\Delta \rightarrow 0$ . Further,  $o(\mathbf{x}) \in (0, 1)$ , so that  $\xi < 0.5$ . Furthermore:

$$\lim_{\substack{\xi \rightarrow 0 \\ N \rightarrow \infty}} \sum_{i=1}^N (0.5 + \xi)^i = 1 \quad (20)$$

as the limits are interchangeable:

$$\begin{aligned}\lim_{\xi \rightarrow 0} \lim_{N \rightarrow \infty} \sum_{i=1}^N (0.5 + \xi)^i &= 1 \\ \lim_{\xi \rightarrow 0} \sum_{i=1}^{\infty} (0.5 + \xi)^i &= 1 \\ \lim_{\xi \rightarrow 0} \left( 1 + \frac{2\xi}{0.5 + \xi} \right) &= 1\end{aligned}$$

as well as

$$\begin{aligned}\lim_{N \rightarrow \infty} \lim_{\xi \rightarrow 0} \sum_{i=1}^N (0.5 + \xi)^i &= 1 \\ \lim_{N \rightarrow \infty} \sum_{i=1}^N 0.5^i &= 1\end{aligned}$$

Therefore, for any  $\epsilon > 0$ , there exists  $\delta > 0$  and  $N > 0$  such that:

$$\begin{aligned}|\xi| < \delta &\Rightarrow \left| \sum_{i=1}^N (0.5 + \xi)^i - 1 \right| < \epsilon \\ &\Rightarrow \left| \sum_{i=1}^N (0.5 + \xi)^i c(\mathbf{x}_s) - c(\mathbf{x}_s) \right| < \epsilon c(\mathbf{x}_s) \\ &\Rightarrow \left| \sum_{i=1}^N (0.5 + \xi)^i (c(\mathbf{x}_s) + \gamma) - c(\mathbf{x}_s) \right| < \epsilon c(\mathbf{x}_s) + \left| \gamma \sum_{i=1}^N (0.5 + \xi)^i \right| \\ &\Rightarrow \left| \sum_{i=1}^N o(\mathbf{x}_i) \prod_{j < i} (1 - o(\mathbf{x}_j)) c(\mathbf{x}_i) - c(\mathbf{x}_s) \right| < \epsilon c(\mathbf{x}_s) + \left| \gamma \sum_{i=1}^N (0.5 + \xi)^i \right|\end{aligned}$$

Since  $\gamma \rightarrow 0$  for  $\Delta \rightarrow 0$ , the right-hand term can become arbitrarily small, completing the proof. □

**Remark:** While  $N \rightarrow \infty$  might appear as a strong assumption, we remark that even for moderate values of  $N$  as those used in our experiments  $\sum_{i=1}^N 0.5^i$  becomes very close to 1 very quickly.

## II. Experimental Setup

In this section, we provide more details about the experimental setup of all experiments. We focus here specifically on the spatial definition of the scenes.

### II.1. DTU MVS dataset

Previous works [38, 61] consider a unit cube or unit sphere as the region of interest for 3D reconstruction of DTU objects. These methods reconstruct objects using object masks, where the respective visual hull lies inside the unit cube/sphere. As we do not consider masks in our work, we need to model the entire scene within the field of view e.g. including the table. Hence, we define a larger region of interest. We also consider a sphere as the region of interest but with a four-times larger radius. Our ray evaluations take place only inside the area of interest. Note that, our model is not sensitive to this choice, but an overly large area of interest requires adaptation of the sampling accuracy for the root-finding, and a too-small region can not represent the whole scene. As we assume to cover the entire scene inside this region of interest, we consider the background as black.

While we do not consider masks during optimization, we evaluate all methods only inside the mask areas (inside the visual hull). The reason of doing so is that the visual hull is the only region that is represented by all methods including IDR. Hence, we evaluate all methods inside this area. With this procedure, we can guarantee a fair comparison among all baselines and our method. More qualitative results for the *DTU MVS* dataset are shown in Fig. 11.

### II.2. Indoor Scene from SceneNet

For the indoor scene, we use Colmap [48] to obtain the camera extrinsics and intrinsics. We then define the area of interest such that all camera locations are inside a sphere with the sphere’s center approximately at the scene’s center. As before, we assume a black background. In addition to the main paper, we show the result of one more scene in Fig. 13.

### II.3. BlendedMVS

Besides the datasets mentioned previously where we either have a black background (DTU) or a closed scene definition (indoor scene), we also consider the BlendedMVS dataset that has scenes with more complex backgrounds. Here, the multi-view images contain objects as well as complex backgrounds that can be located further away or appear blurred. Since the BlendedMVS dataset consists of a large variety of scene layouts, we must model not only the foreground but also the background of the scenes.

To this end, we found that the setup of NeRF++ [64] is useful for extending our model to complex backgrounds. The main idea is to model foreground and background using two separate models by spatially separating the representations. Similarly, we define the area of interest for reconstruction as a sphere that covers all cameras centered at an approximate scene center. Within this sphere, we use our model for representing the scene. Everything that is located outside, we represent with a NeRF model that has an *inverted sphere parameterization*. We refer the reader to [64] for more details. The inverted sphere parameterization allows for both, representing far-away background as well as background elements that are closer to the area of interest.

For rendering during the optimization process, we apply our rendering procedure for each ray with 64 samples inside the interval and 32 samples in the free space between the camera and interval bound. Furthermore, we uniformly sample 32 points outside of the sphere to roughly capture the background. We use both sets of sampled points for volume rendering. As we aim for 3D reconstruction inside the area of interest, we use significantly fewer samples for capturing the background as used in [64]. This results in less computational effort, while the model is still able to separate foreground and background properly. Additional results are shown in Fig. 13.

## III. Baselines

This section provides additional details about the baselines we compare against.

**IDR:** For the DTU evaluations, we use the official code<sup>6</sup> and the provided pre-trained models. To test IDR without mask supervision, we set the weight of the mask loss to zero and consider the RGB loss for all rays intersecting surfaces.

**Colmap:** We use the same Colmap procedure as reported in previous works [38, 61]. Therefore, the original Colmap [48] implementation is used to output meshes with the following steps: a) `exhaustive_matcher`, b) `point_triangulator`, c) `patch_match_stereo`, d) `stereo_fusion` and e) `poisson_mesher`.

<sup>6</sup><https://github.com/lioryariv/idr>

threshold $\sigma$	1	5	10	<b>50</b>	100	500
scan65	2.29	1.53	1.26	1.27	1.80	3.15
scan105	3.46	2.27	1.85	1.07	1.32	5.99
scan114	2.88	1.74	1.37	1.06	1.21	2.86

Table 2: **Volume Density Thresholds of NeRF.** We show the Chamfer distance for meshes extracted from a trained NeRF model considering different density thresholds  $\sigma$ . We applied this analysis for three models of the DTU dataset and found that a threshold  $\sigma = 50$  leads to the best overall performance.

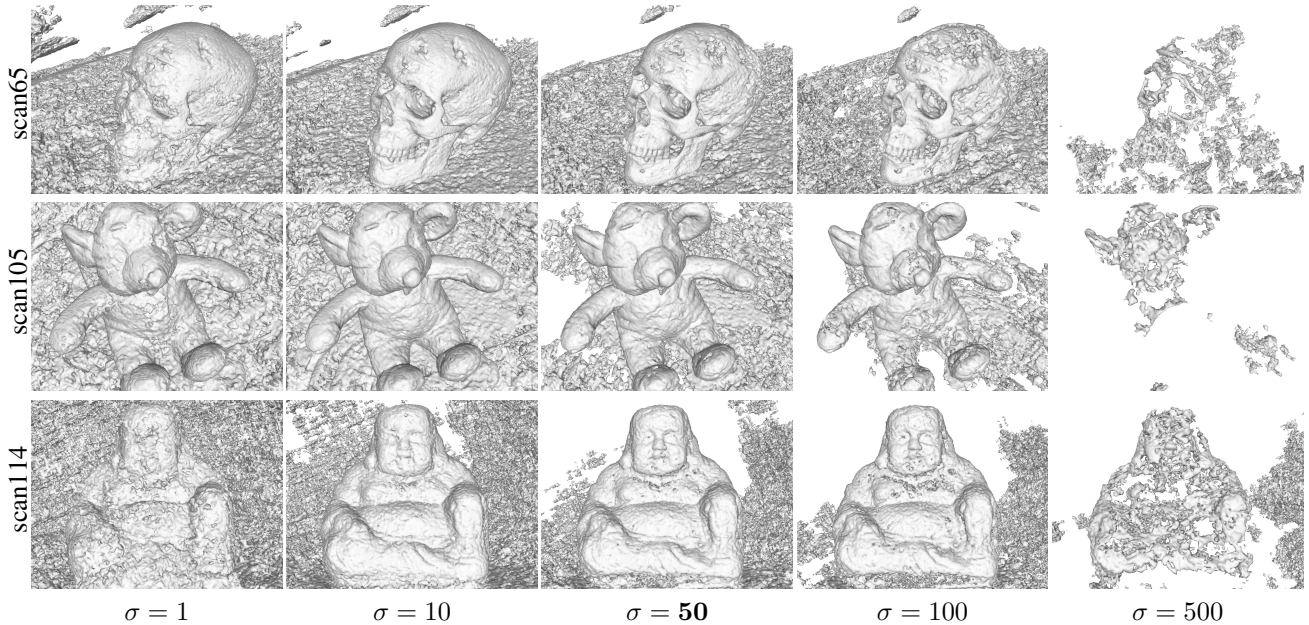


Figure 9: **Volume Density Thresholds of NeRF.** We show a qualitative comparison of meshes extracted from NeRF for different thresholds of the volume density. We see that the extracted meshes significantly differ for various threshold parameter and identify a threshold of 50 as superior, qualitatively and quantitatively

**NeRF:** For the NeRF baseline, we adapt the Pytorch reimplementation [62] to our framework. We apply the NeRF model to the same scene setups as it is used for our method. For extracting meshes we need to choose a threshold for the volume density. Therefore we evaluate NeRF for three DTU objects and five different threshold parameters  $\sigma$ .

In Table 2, we provide a quantitative comparison that shows a superior performance at  $\sigma = 50$ . In Fig. 9, we depict qualitative results for different threshold parameters. We see that small thresholds, e.g. 1 and 10, lead to bloated reconstructions, while the surface for  $\sigma > 50$  shows significantly more missing regions. The qualitative comparison verifies our findings from our quantitative evaluation, and hence, we choose 50 as the threshold for all evaluations in the paper.

#### IV. Limiting Factors

**Overexposed regions:** Challenging for nearly all multi-view reconstruction methods are overexposed and textureless image regions. In the DTU dataset, the table appears completely white without any texture in some of the scenes, e.g., second row in Fig. 11. While this problematic region is masked during training of the IDR baseline, our method and the other baselines do not utilize pixel-accurate masks and hence struggle predicting accurate surfaces in these areas. Surfaces extracted by NeRF are typically non-smooth with numerous geometric artifacts. Our method also struggles to output smooth surfaces even though we use a surface regularizer during optimization. We attribute this to the inherent ambiguity of the texture-less overexposed regions that can not be entirely resolved by such simple prior assumptions, encouraging the development of more advanced priors and implicit models that are able to expose their uncertainty in the future.

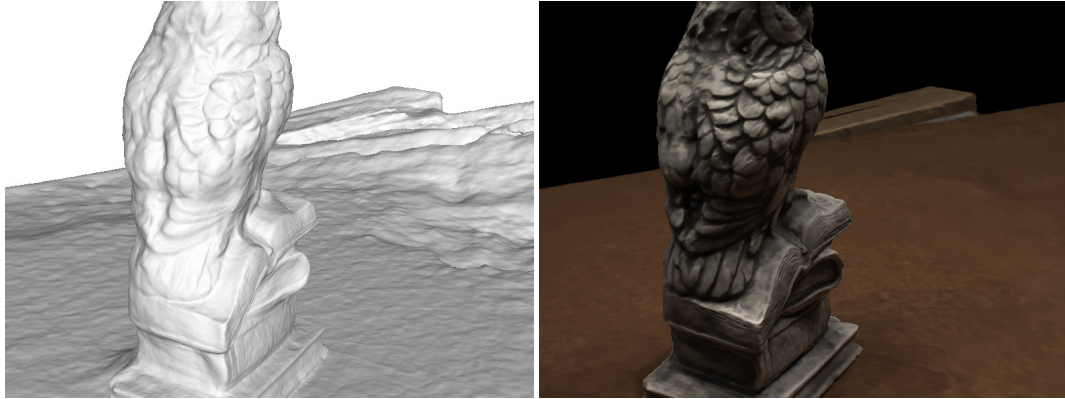


Figure 10: **Rarely visible regions.** In this figure, we show an extreme view from the input images. The table area in the back is reconstructed poorly due to missing observations of this surface.

**Rarely visible regions:** Similar difficulties arise from rarely visible regions in the input views. In Fig. 10, we show an extreme view from the dataset, where we see a rarely visible region in the back. In this area, the reconstructed surface of the table appears to be less accurate compared to the rest of the table which is a consequence of missing observations.

**Shape-Appearance Ambiguity:** For reconstructing neural surfaces from multi-view images, it is necessary to optimize neural implicit surfaces and the appearance representations at the same time. The network architecture and the optimization procedure yield an inductive bias whether predicted views are explained by adapting the shape or the view-dependent appearance. In the bottom row in Fig. 11, we show an example where our method is not able to correctly model the geometry and our model explains the inner part with the view-dependent appearance instead, hence not reconstructing the inner part correctly. NeRF does not show this behavior as its capacity for modeling view-dependent effects is much smaller compared to our model. Hence, for objects with weakly view-dependent appearance, it NeRF less prone to this behavior. However, this also limits NeRF’s capabilities in modeling plausible geometry and strong view-dependent effects. IDR can resolve this particular example due to its strong mask supervision.

**Possible Solution:** To circumvent these limiting cases, we hypothesize that a learning-based prior should be applied to resolve the underlying ambiguities. In future work, we, therefore, consider learning a probabilistic neural surface model which captures regularities and uncertainty across objects. We believe that such priors will help to resolve the ambiguities in texture-less areas, rarely visible regions as well as the aforementioned shape-appearance ambiguity.

## V. Additional Results

We show additional qualitative results for the *DTU MVS* dataset in Fig. 11. Fig. 12 shows the geometry and appearance from different views of the objects. In Fig. 13, we show reconstructions on more scenes from *SceneNet* and *BlendedMVS*.



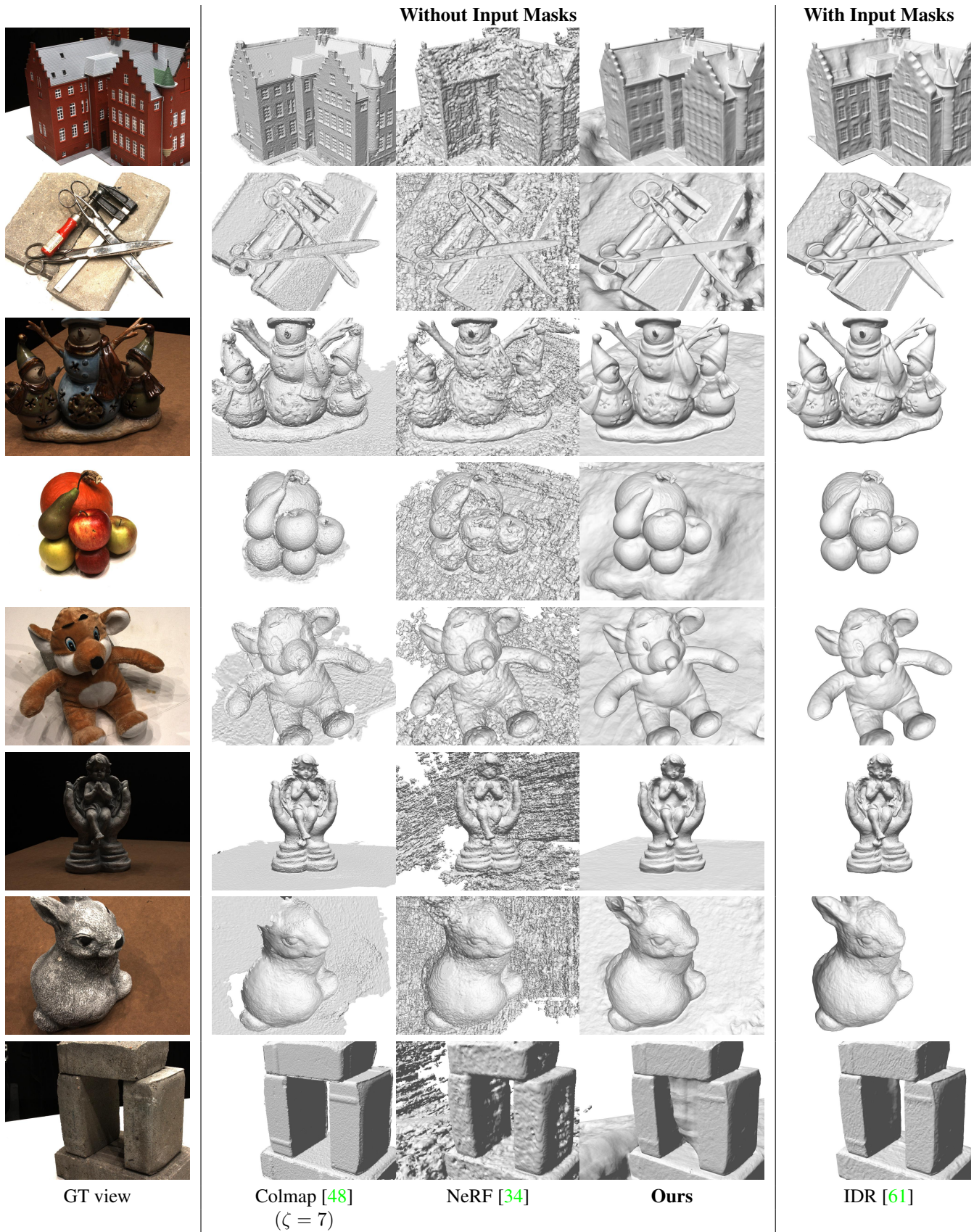


Figure 11: DTU MVS dataset. Additional results for objects from the DTU MVS dataset.



Figure 12: **Rotations.** In this figure, we show different views from DTU objects and the respective surfaces.

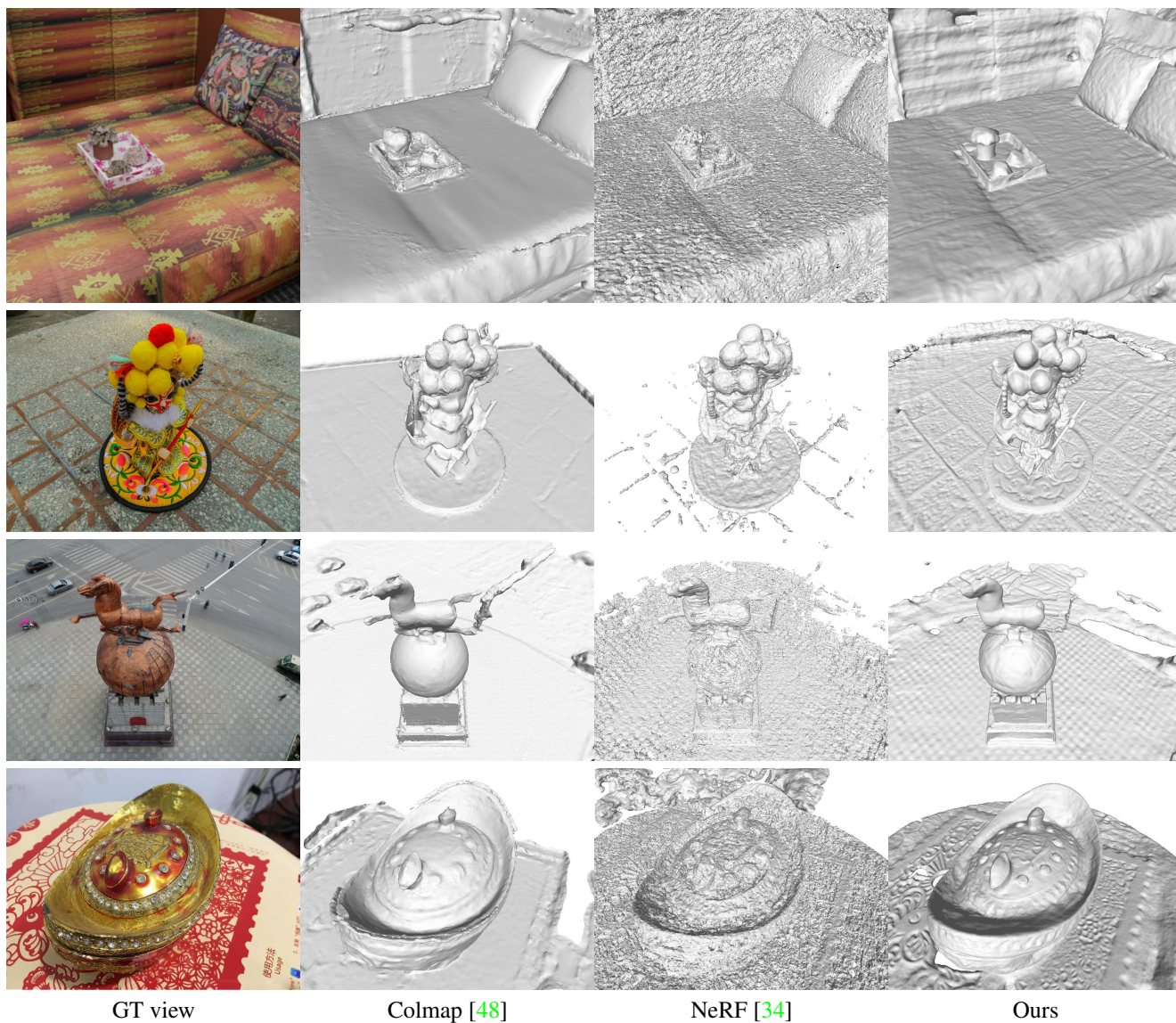


Figure 13: **Indoor Scene and BlendedMVS Scenes.** We show additional results for one indoor scene (first row) and objects from the *Blended MVS* dataset.