

# GO-Surf: Neural Feature Grid Optimization for Fast, High-Fidelity RGB-D Surface Reconstruction

Jingwen Wang\*, Tymoteusz Bleja\*, and Lourdes Agapito

Department of Computer Science, University College London

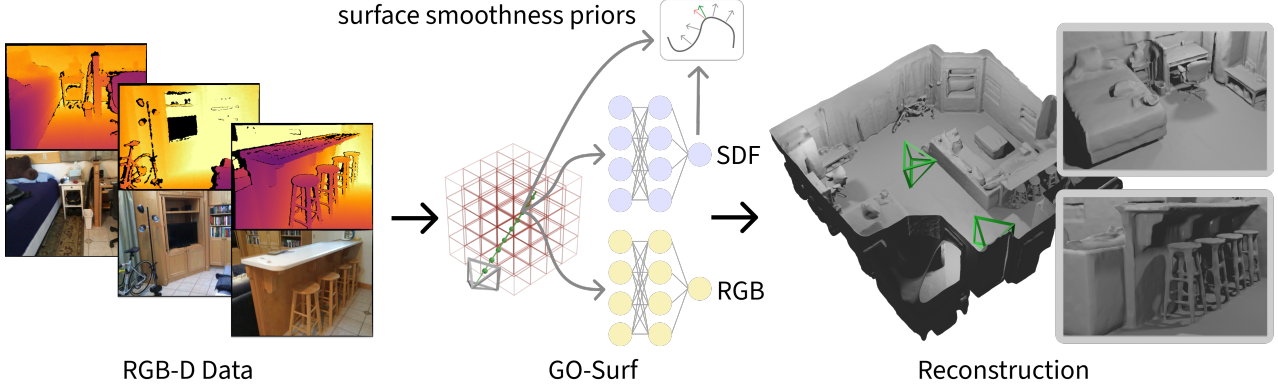


Figure 1: Given an input RGB-D sequence, GO-Surf obtains a high quality 3D surface reconstruction by direct optimization of a multi-resolution feature grid and signed distance value and colour prediction. We formulate a new smoothness prior on the signed distance values that leads to improved hole filling and smoothness properties, while preserving details. Our optimization is  $\times 60$  times faster than MLP-based methods.

## Abstract

We present *GO-Surf*, a direct feature grid optimization method for accurate and fast surface reconstruction from RGB-D sequences. We model the underlying scene with a learned hierarchical feature voxel grid that encapsulates multi-level geometric and appearance local information. Feature vectors are directly optimized such that after being tri-linearly interpolated, decoded by two shallow MLPs into signed distance and radiance values, and rendered via volume rendering, the discrepancy between synthesized and observed RGB/depth values is minimized. Our supervision signals — RGB, depth and approximate SDF — can be obtained directly from input images without any need for fusion or post-processing. We formulate a novel SDF gradient regularization term that encourages surface smoothness and hole filling while maintaining high frequency details. *GO-Surf* can optimize sequences of 1-2K frames in 15-45 minutes, a speedup of  $\times 60$  over *NeuralRGB-D* [1], the most related approach based on an MLP representation, while maintaining on par performance on standard benchmarks. Project page: [https://jingwenwang95.github.io/go\\_surf](https://jingwenwang95.github.io/go_surf).

## 1. Introduction

Recent years have seen impressive progress in learning-based methods for 3D scene modelling from video sequences, with a strong focus on coordinate-based scene representations [21, 11] and the application to highly photo-realistic novel view synthesis [27, 12]. NeRF [12] and its variants represent the scene in the weights of a multi-layer perceptron (MLP) that is trained to predict the volume density and colour of any 3D point location. Combined with classic differentiable volume rendering, NeRF learns to synthesize the input images and can generalize to render nearby unseen views.

Although MLPs have shown an extraordinary ability to represent radiance fields of complex scenes with a low memory footprint, this comes at the cost of very long training and inference times. Direct optimization of trilinearly interpolated features on multi-resolution feature grids has become a powerful alternative representation due to its fast convergence that leads to orders of magnitude faster training and inference times [13, 8, 37, 41, 3, 30]. Their higher memory requirement has also been addressed by using sparse feature grids, such as octrees [37], or multi-resolution hash tables [13].

An added limitation of NeRF [12] is its choice of ge-

\* The first two authors contributed equally.

ometric primitive: the volume density representation does not allow for accurate 3D surface extraction and can lead to blurry renderings or 3D reconstruction artefacts because of the limited capacity of the network to infer the exact surface location. A number of methods have addressed this [33, 16, 36] and proposed differentiable surface rendering pipelines that can be used in conjunction with neural implicit surface representations, thus avoiding representing the scene geometry via density. However, surface representations have mainly been proposed in the context of slow coordinate-based networks, while existing grid-based approaches struggle to achieve the level of smoothness comparable to MLP-based methods [13]. Critically though, most AR/VR and robotics applications require both accurate and fast 3D surface reconstruction, therefore both limitations should be tackled for high fidelity surface inference at interactive runtimes.

GO-Surf combines, for the first time, learnable feature volumes with SDF surface-based representations and rendering to improve both speed and accuracy. To further help overcome these limitations we adopt the advantages of using depth measurements from consumer-level RGB-D cameras, which have become highly accessible, inexpensive, and are now found on many mobile-phones. Other approaches have recently also adopted RGB-D inputs in the context of neural scene representations [29, 41, 20, 1] and use depth images to supervise scene reconstruction. Unlike GO-Surf, iMAP [29] models volume density, while NICE-SLAM [41] uses occupancy, both failing to predict high-fidelity, accurate surface reconstructions. Although NeuralRGB-D [1] adopts an SDF representation, such as GO-Surf, it uses an MLP to encode the scene, leading to very long training times. Our experimental evaluation shows that GO-Surf achieves a speedup of  $\times 60$  through use of learnable feature volumes. Although iSDF [20] also uses an SDF representation, we show that our novel SDF regularization term results in more accurate and higher fidelity reconstructions.

**Contribution:** GO-Surf brings learnable feature grids into the context of SDF reconstruction from RGB-D sequences to achieve both: (i) fast optimization at interactive runtimes, and (ii) highly accurate surface reconstruction. We also apply for the first time Eikonal and smoothness regularisation terms in the context of voxel grids.

## 2. Related Work

**Classic RGB-D Reconstruction** Although this review focuses on learning based approaches to RGB-D reconstruction we cite KinectFusion [15], VoxelHashing [17] and BundleFusion [6] as the best representatives of classic methods, with the last one being one of our baselines.

**Coordinate-based Networks** Pioneered by Scene Representation Networks [27], and popularized by NeRF [12],

coordinate-based methods encode the scene within the weights of a fully connected neural network that take as input a 3D location and predict its geometry and appearance information in the form of density and radiance [12, 29], occupancy [11, 4], signed distance fields [21, 7, 1, 20], colour [19] or semantic labels [40, 9]. NeRF [12] demonstrated photo-realistic novel view synthesis by combining an MLP to encode density and colour with classical volumetric rendering, only requiring a set of posed RGB images for training. Using hybrid scene representations, NVSF [10] and Plenotrees [38] combine neural implicit fields with explicit sparse voxel structures to accelerate rendering.

Closely related to our work are coordinate-based networks trained with additional depth supervision for 3D reconstruction [29, 20, 1]. iMAP [29] is the first to train an MLP in live, real-time operation with RGB-D inputs to learn a scene specific 3D model of occupancy and colour, using rendered depth and RGB as supervision signals, while also tracking camera pose. In iSDF [20] the focus is actual surface reconstruction. Taking as input a live stream of posed depth images, approximate signed distance values are predicted by minimising a loss on both the predicted signed distance and its spatial gradient. Both [29, 20] achieve live operation via active sampling and keyframe selection. Instead, Neural RGB-D [1] operates in batch mode, taking as input a large set of RGB-D frames to train MLPs to map coordinate inputs to approximate SDF values and radiance via differentiable volumetric rendering. While they achieve extremely accurate reconstructions that preserve high frequency details, their training time is on the order of 10 hours for typical scenes. The huge advantage of these MLP-based representations is their compactness, given that they can be sampled continuously without increasing memory footprint, they either suffer from extremely long training times [1] or sacrifice accuracy/details for online operation [29, 20].

**Neural 3D Scene Reconstruction** A recent trend includes architectures trained end-to-end on multiple sequences to aggregate image features and perform fusion over multiple frames using a volumetric representation and decoders to predict signed distance [31, 14, 2] or radiance [39] fields. NeuralRecon [31] reconstructs surfaces sequentially from video fragments as TSDF volumes by performing feature fusion from previous fragments via recurrent units. Atlas [14] aggregates image features over an entire sequence to predict a globally consistent TSDF volume and semantic labels with a 3DNN, while others [2, 28] rely on transformers for feature fusion. While the above methods require direct SDF supervision at training time, inference only requires posed RGB video. NeRFusion [39] uses a 3D CNN for feature fusion but drops the need for 3D supervision by including a volume rendering module and can be trained only from RGB images.

**Neural Implicit Surface Reconstruction** Following a sim-

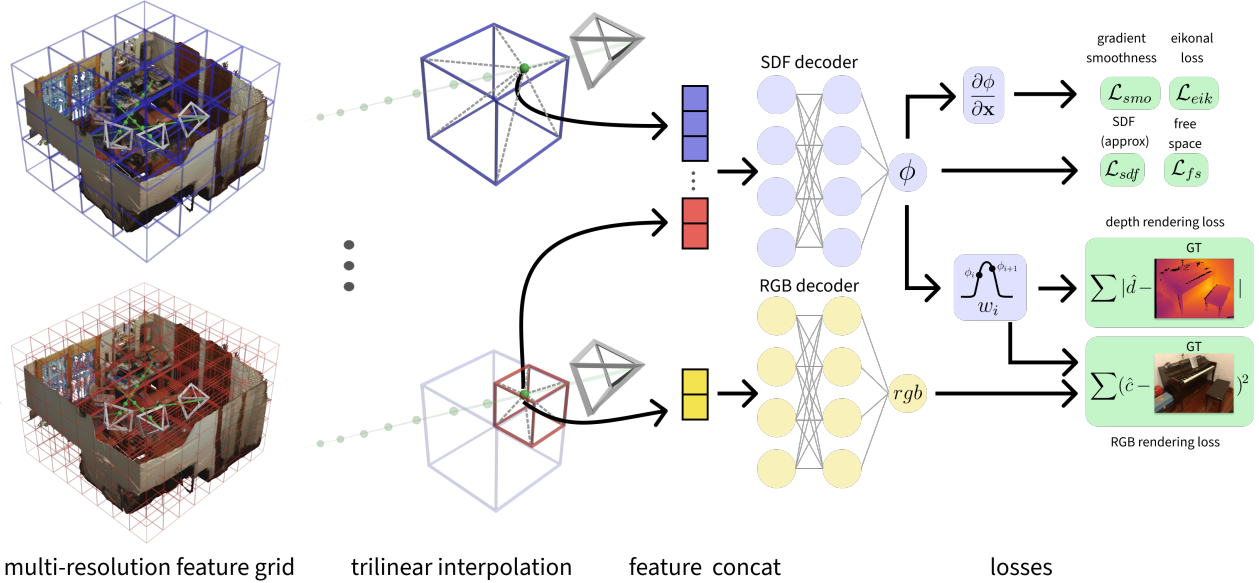


Figure 2: GO-Surf uses multi-level feature grids and two shallow MLP decoders. Given a sample point along a ray, each grid is queried via tri-linear interpolation. Multi-level features are concatenated and decoded into SDF, and used to compute the sample weight. Color is decoded separately from the finest grid. Loss terms are applied to SDF values, and rendered depth and color. The gradient of the SDF is calculated at each query point and used for Eikonal and smoothness regularization.

ilar observation to ours — that the volume density estimated by NeRF [12] does not enable high-quality surface extraction — a recent family of methods focus instead on neural surface representations and formulate compatible differentiable renderers. DVR [16] propose the first differentiable renderer for implicit shape and texture representations needing only multiview RGB images and object masks as supervision. IDR [36] trains an end-to-end architecture with a learned neural renderer that approximates light reflected from the surface and can implicitly model a variety of lighting conditions and materials. NeuS [33] instead formulates a new volume rendering method to train a bias-free neural SDF representation, while VolSDF [35] formulate a novel parameterization for volume density, both leading to more accurate surface reconstruction even without mask supervision. UniSurf [30] unifies neural volume and surface rendering enabling both within the same model.

**Direct Grid Optimization** Recent, mostly concurrent, work has also explored the use of directly optimizing neural features on an explicit volumetric grid. NGLoD [32] represents implicit surfaces using a sparse octree feature volume which adaptively fits shapes with multiple levels of detail using 3D supervision. Instant NGP [13] proposes a multi-resolution hash table of trainable feature vectors achieving a speedup of various orders of magnitude, demonstrating instant training on a variety of tasks. DirectVoxGO [30] also improves NeRF’s training time by two orders of magnitude by adopting an explicit density voxel-grid with post-activation interpolation — that enables to

model sharp boundaries — and a feature voxel-grid with a shallow MLP for view-dependent appearance. Plenoxels [37] and ReLU Fields [8] take this idea further dropping the reliance on any neural networks. While Plenoxels [37] relies on a spherical harmonics basis, ReLU Fields [8] apply post-activation interpolation. TensorRF [3] tackles the issue of memory footprint by modelling the volume field as a 4D tensor that is factorized into multiple compact low-rank tensor components. Finally, NiceSLAM [41] is an RGB-D SLAM system that encodes the scene as a multi-level feature grid that is optimized at run-time using pre-trained geometric priors, enabling detailed reconstruction on large indoor scenes.

### 3. Method

We formulate scene reconstruction as a direct optimisation problem, representing scene geometry and color using multi-resolution feature grids, which are decoded into SDF and RGB values with two MLP decoders shared across grids. Given a sequence of RGB-D images  $\{I_t, D_t\}$ , and initial camera poses  $\{\xi_t\} \in \mathbb{SE}(3)$  (provided by SLAM or SfM) we jointly optimize feature volumes, decoders and camera poses.

#### 3.1. Hybrid Scene Representation

Our hybrid scene representation combines multi-level feature grids and shared MLPs for geometry and color prediction. We encode the scene geometry into a four-level feature grid  $\mathcal{V}_\theta = \{V^l\}$ , where  $l \in \{0, 1, 2, 3\}$  encode multi-

level local detail from coarse to fine. Finer features capture high-frequency detailed geometry while the coarser features encapsulate larger structures and are crucial for hole-filling.

Given a scene point  $\mathbf{x} \in \mathbb{R}^3$ , its geometry feature is obtained by concatenating the trilinearly interpolated features at each level. The feature is decoded into an SDF value  $\phi(\mathbf{x})$  via the geometry MLP  $f_\omega(\cdot)$ :

$$\mathcal{V}_\theta(\mathbf{x}) = [V^0(\mathbf{x}), V^1(\mathbf{x}), V^2(\mathbf{x}), V^3(\mathbf{x})] \quad (1)$$

$$\phi(\mathbf{x}) = f_\omega(\mathcal{V}_\theta(\mathbf{x})) \quad (2)$$

We encode color information only at the finest level, as [41], using a separate feature grid  $\mathcal{W}_\beta$  and decoder  $g_\gamma(\cdot)$ :

$$\mathbf{c} = g_\gamma(\mathcal{W}_\beta(\mathbf{x}), \mathbf{d}) \quad (3)$$

where  $\mathbf{d}$  is the viewing direction. Here  $\theta, \beta, \omega, \gamma$  represent the optimizable parameters, i.e. geometry and color local features and decoder parameters.

**Architecture Details.** The multi-resolution feature grids, use geometry features of dimension  $h_g = 4$  per-level, giving a final geometry feature vector of dimension 16. Color features have dimension  $h_c = 6$ . For both geometry and color decoders we use light-weight MLPs with 2 hidden layers, 32 neurons each.

### 3.2. Depth and Color Rendering

Inspired by the recent success of volume rendering [12], we render depth and color by integrating predicted colors and sampled depth along camera rays. Specifically, for each back-projected ray parameterised by camera centre  $\mathbf{o}$  and ray direction  $\mathbf{r}$  we sample  $N$  points  $\mathbf{x}_i = \mathbf{o} + d_i \mathbf{r}, i \in \{1, \dots, N\}$  and take the expected values of predicted colors  $\mathbf{c}_i$  and sampled depth  $d_i$ :

$$\hat{\mathbf{c}} = \sum_{i=1}^N w_i \mathbf{c}_i, \quad \hat{d} = \sum_{i=1}^N w_i d_i \quad (4)$$

where  $\{w_i\}$  are unbiased and occlusion-aware weights [33] given by  $w_i = T_i \alpha_i$ , where  $T_i = \prod_{j=1}^{i-1} (1 - \alpha_j)$  represents the *accumulated transmittance* at point  $\mathbf{x}_i$ , and  $\alpha_i$  is the *opacity value* computed by:

$$\alpha_i = \max \left( \frac{\sigma_s(\phi(\mathbf{x}_i)) - \sigma_s(\phi(\mathbf{x}_{i+1}))}{\sigma_s(\phi(\mathbf{x}_i))}, 0 \right) \quad (5)$$

where  $\sigma_s(x) = (1 + e^{-sx})^{-1}$  is a Sigmoid function modulated by a learnable parameter  $s$  which controls the smoothness of the transition at the surface. The RGB and depth per-pixel rendering losses are:

$$\ell_{rgb} = \|\mathbf{I}[u, v] - \hat{\mathbf{c}}\|, \quad \ell_d = |D[u, v] - \hat{d}| \quad (6)$$

To obtain sampling points along the ray we perform 3 steps of importance sampling, as [33]: starting from  $N_c = 96$  coarse samples we iteratively add 12 samples each time based on weights computed with previously sampled points.

### 3.3. Approximate SDF Supervision

Similar to [1, 20], for each sampled point along the ray we also approximate ground-truth SDF supervision based on the distance to observed depth values along the ray direction:  $b(\mathbf{x}) = D[u, v] - d$ . With this bound we have  $|\phi(\mathbf{x})| \leq |b(\mathbf{x})|, \forall \mathbf{x}$ , which is expected to be tighter near the surface, so for near-surface points ( $|D[u, v] - d| \leq t$ ) we apply the following SDF loss:

$$\ell_{sdf}(\mathbf{x}) = |\phi(\mathbf{x}) - b(\mathbf{x})| \quad (7)$$

The truncation threshold  $t$  is a hyper-parameter (we set  $t = 16\text{cm}$ ). Unlike [1] which models TSDF for points far from the surface ( $|D[u, v] - d| > t$ ), encouraging the MLP to predict a fixed truncation value; we apply a relaxed loss to encourage free space prediction as [20]:

$$\ell_{fs}(\mathbf{x}) = \max \left( 0, e^{-\alpha \phi(\mathbf{x})} - 1, \phi(\mathbf{x}) - b(\mathbf{x}) \right) \quad (8)$$

We apply no penalty if the SDF prediction is positive and smaller than the bound, a linear loss if it is larger, and exponential penalty ( $\alpha = 5$ ) if it is negative. Despite being approximated SDF values, these terms provide more direct supervision than the rendering terms in Sec. 3.2. Empirically we observed that they add robustness to the optimisation.

### 3.4. Grid-based SDF Regularisation

Directly optimizing the feature volume from RGB-D observations is not sufficiently constrained and can lead to artifacts and invalid predictions. We employ two regularization terms on the SDF values: Eikonal term  $\ell_{eik}$  and smoothness term  $\ell_{smooth}$ . Eikonal regularisation has been used in [7, 33, 20] to encourage valid SDF values, especially in areas without direct supervision. Specifically, at any query point, the gradient of the SDF value w.r.t. 3D query coordinates should have unit length:

$$\ell_{eik}(\mathbf{x}) = (1 - \|\nabla \phi(\mathbf{x})\|)^2 \quad (9)$$

Intuitively, this term encourages the absolute value of SDF to increase uniformly as we move away from the surface. Unlike [1, 33, 29, 41] that use an MLP scene representation with built-in smoothness priors, we must explicitly add a smoothness term. Our loss minimises the difference between the gradients of nearby points.

$$\ell_{smooth}(\mathbf{x}) = \|\nabla \phi(\mathbf{x}) - \nabla \phi(\mathbf{x} + \epsilon)\|^2 \quad (10)$$

Here  $\mathbf{x}$  is taken from near-surface points sampled randomly across the whole model,  $\epsilon$  is a small perturbation of length  $\delta_s$  and random direction (we set  $\delta_s$  between 4mm and 1mm, see Fig. 3).

**Customised Grid Sampler.** Adding the above terms to the final loss requires second order derivatives of the SDF to be



Figure 3: The distance between sampled points for gradient regularisation ( $\delta_s$ ) controls surface smoothness. Left:  $\delta_s = 1$  mm, Middle:  $\delta_s = 4$  mm, Right:  $\delta_s = 2$  cm.

computed. The SDF value is a function of 3D coordinates  $\mathbf{x}$ , feature vectors  $\theta$ , and geometry network parameters  $\omega$ :

$$\phi(\mathbf{x}, \theta, \omega) = f_\omega(\mathcal{V}_\theta(\mathbf{x})) \quad (11)$$

The gradient of the SDF w.r.t. 3D query coordinates is:

$$\frac{\partial \phi}{\partial \mathbf{x}} = \frac{\partial f_\omega(\mathcal{V}_\theta(\mathbf{x}))}{\partial \mathcal{V}_\theta(\mathbf{x})} \frac{\partial \mathcal{V}_\theta(\mathbf{x})}{\partial \mathbf{x}} \quad (12)$$

To regularise this gradient we need to compute the first three rows (since  $\mathbf{x}$  is 3-dimensional) of the Hessian of the SDF (Eq. 11). In practice, Pytorch’s automatic differentiation package supports double back-propagation through an MLP in the computation graph, but we still need to compute second order derivatives of the trilinear interpolation w.r.t. features  $\theta$  and query coordinates  $\mathbf{x}$ . We implement a drop-in replacement of Pytorch’s `grid_sampler` that supports double back-propagation. Specifically, we reuse Pytorch’s implementation for forward and backward passes, and add a new CUDA kernel for the second backward pass. Our implementation will be open-sourced. For derivation of the second order derivatives of trilinear interpolation, please refer to the supplementary material.

### 3.5. Optimization

To optimize our scene representation described in Sec 3.1, at each iteration we sample a batch of  $R$  rays from all pixels across all images (including pixels that have missing depth values). Our global objective function  $\mathcal{L}(\mathcal{P})$  is:

$$\mathcal{L}(\mathcal{P}) = \lambda_{rgb} \mathcal{L}_{rgb} + \lambda_d \mathcal{L}_d + \lambda_{sdf} \mathcal{L}_{sdf} + \lambda_{fs} \mathcal{L}_{fs} + \lambda_{eik} \mathcal{L}_{eik} + \lambda_{smooth} \mathcal{L}_{smooth} \quad (13)$$

where  $\mathcal{P} = \{\theta, \omega, \beta, \gamma, \{\xi_t\}\}$  is the list of parameters being optimised, including features, decoders and camera poses.  $\mathcal{L}_{rgb}$  and  $\mathcal{L}_d$  measure the difference between observed pixel color/depth and rendered values.  $\mathcal{L}_{rgb}$  is computed over all the sampled rays, while for  $\mathcal{L}_d$  only rays with valid depth values  $R_d$  are considered.

$$\mathcal{L}_{rgb} = \frac{1}{M} \sum_{m=1}^M \ell_{rgb}^m, \quad \mathcal{L}_d = \frac{1}{|R_d|} \sum_{r \in R_d} \ell_d^m \quad (14)$$

SDF and free space losses,  $\mathcal{L}_{sdf}$  and  $\mathcal{L}_{fs}$  are applied on two disjoint subsets  $S_{tr}$  and  $S_{fs}$  of sampled ray points:

$$\mathcal{L}_{sdf} = \frac{1}{M} \sum_{m=1}^M \frac{1}{|S_{tr}|} \sum_{s \in S_{tr}} \ell_{sdf}(\mathbf{x}_s) \quad (15)$$

$$\mathcal{L}_{fs} = \frac{1}{M} \sum_{m=1}^M \frac{1}{|S_{fs}|} \sum_{s \in S_{fs}} \ell_{fs}(\mathbf{x}_s) \quad (16)$$

$\mathcal{L}_{eik}$  encourages points far from the surface to have valid SDF predictions and is applied on  $S_{fs}$ :

$$\mathcal{L}_{eik} = \frac{1}{M} \sum_{m=1}^M \frac{1}{|S_{fs}|} \sum_{s \in S_{fs}} \ell_{eik}(\mathbf{x}_s) \quad (17)$$

$\mathcal{L}_{smooth}$  encourages surface smoothness and is applied to randomly sampled near-surface points  $S_g$  over the whole voxel grid:

$$\mathcal{L}_{smooth} = \frac{1}{|S_g|} \sum_{s \in S_g} \ell_{smooth}(\mathbf{x}_s) \quad (18)$$

**Geometric Initialisation.** To improve convergence, we initialize our feature grid and geometry decoder such that the initial surface is a sphere [7] centred at volume origin and with a radius half the size of the smallest dimension.

## 4. Results

### 4.1. Experimental Setup

**Datasets.** We consider two main datasets for evaluation. We qualitatively evaluate our reconstruction quality on real sequences from ScanNet [5] showing our approach is able to achieve complete and smooth reconstructions. We also quantitatively evaluate our method on 10 synthetic datasets following the same evaluation protocol as in [1].

**Metrics.** We measure accuracy, completion, chamfer  $\ell_1$  distance, normal consistency and F-score for reconstruction quality. As for pose accuracy we measure translation and rotation error. The metrics are computed between point clouds sampled at a density of 1 point per  $\text{cm}^2$ . F-score is computed using athreshold of 5cm.

**Baselines.** We consider NeuralRGB-D Surface Reconstruction [1] (NeuralRGB-D) as our main baseline, and also compare quantitatively against some other state-of-the-art traditional as well as learning-based RGB-D reconstruction methods: BundleFusion [6], COLMAP with Poisson surface reconstruction [25, 24, 23], RoutedFusion [34], Convolutional Occupancy Networks [22] and SIREN [26].

**Implementation Details.** We run our reconstruction method on a desktop PC with a 3.60GHz Intel i7-9700K CPU and an NVIDIA RTX-2080Ti GPU. For all our experiments, we sample  $M = 6144$  rays per batch, and  $N_c = 96$

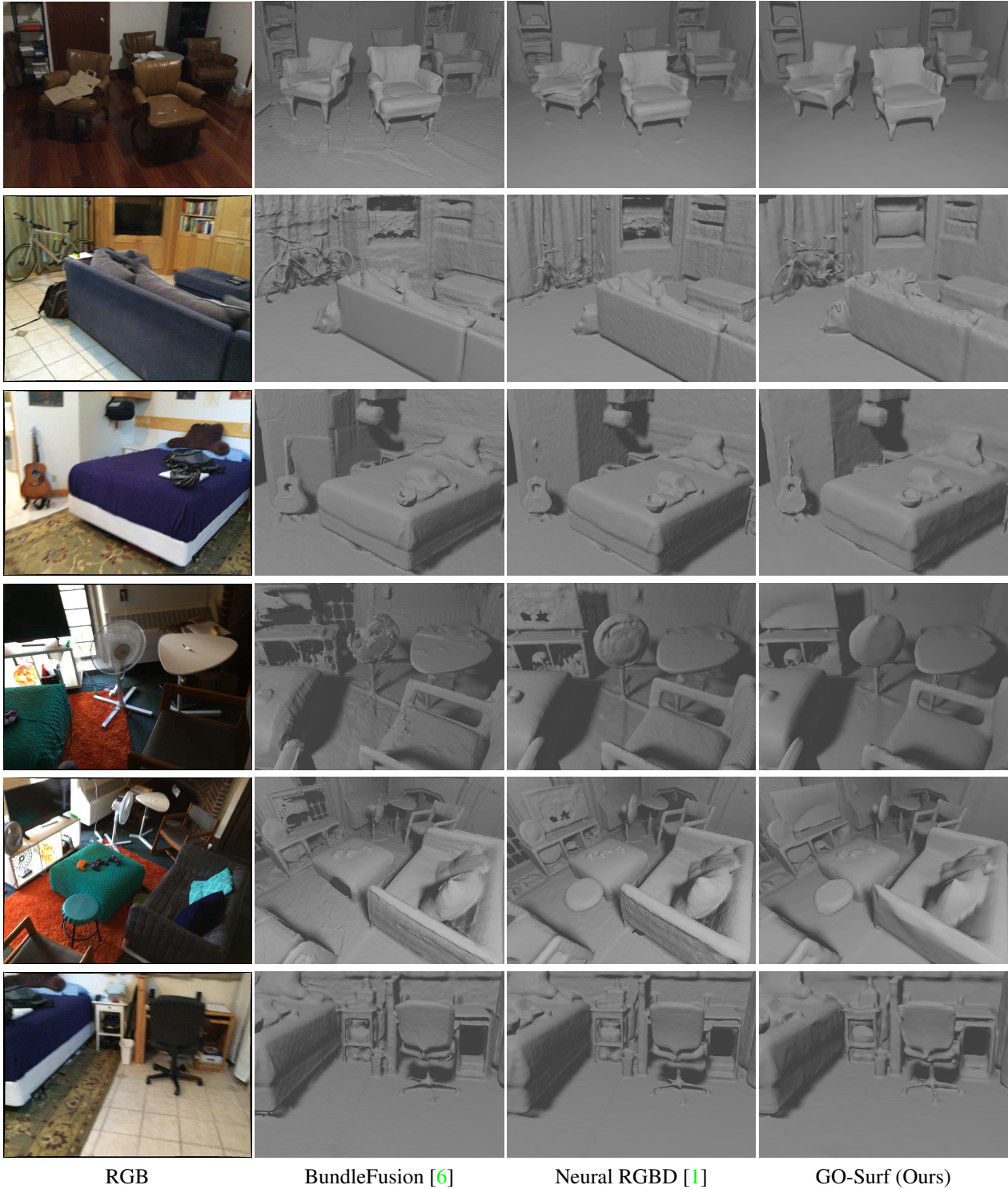


Figure 4: We compare our method to Neural RGB-D Surface Reconstruction [1] and BundleFusion [6] on scenes 0, 2, and 12 of the ScanNet dataset. Qualitative results show that our method displays higher completion and smoothness than the baselines. Our method is also able to produce high-quality reconstructions for missing depth regions and thin structures.

coarse samples and  $N_f = 36$  fine samples along each ray. Our method is implemented in Pytorch using the ADAM

optimizer with learning rate of  $1 \times 10^{-3}$ ,  $1 \times 10^{-2}$  and  $5 \times 10^{-4}$  for MLP decoders, feature grids and camera poses.

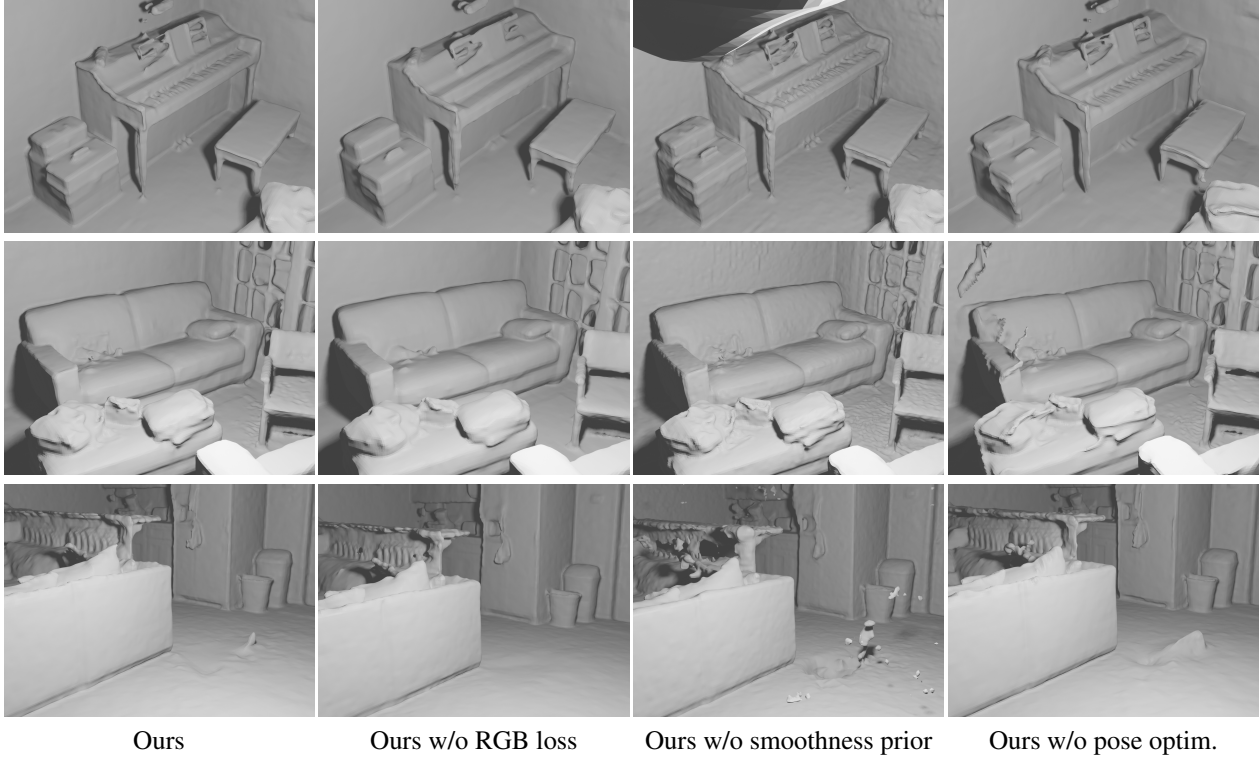


Figure 5: We evaluate our method without different loss components. We observe better reconstruction of fine details with RGB loss which, however, without smoothness prior results in noisy surfaces. Similarly to [1], we observe improved surface consistency with pose refinement.

Method	Acc ↓	Com ↓	C- $\ell_1$ ↓	NC ↑	F-score ↑
BundleFusion	0.0191	0.0581	0.0386	0.9027	0.8439
RoutedFusion	0.0223	0.0364	0.0293	0.8765	0.8736
COLMAP	0.0271	0.0322	0.0296	0.9134	0.8744
ConvOccNets	0.0498	0.0524	0.0511	0.8607	0.6822
SIREN	0.0229	0.0412	0.0320	0.9049	0.8515
Neural RGBD	<b>0.0151</b>	0.0197	<b>0.0174</b>	0.9316	<b>0.9635</b>
Ours	0.0158	<b>0.0195</b>	0.0177	<b>0.9317</b>	0.9591

Table 1: Quantitative evaluation of the reconstruction quality on a dataset of 10 synthetic scenes using the protocol established in [1]. We compare with a number of baselines and with NeuralRGB-D [1]. GO-Surf is on par in terms of performance but is significantly faster.

We set the loss weights to  $\lambda_{rgb} = 10.0$ ,  $\lambda_{depth} = 1.0$ ,  $\lambda_{sdf} = 10.0$ ,  $\lambda_{fs} = 1.0$ ,  $\lambda_{eik} = 1.0$  and  $\lambda_{smooth} = 1.0$ . We use 3cm, 6cm, 24cm and 96cm voxel sizes for each level. Camera pose is parameterised as translation vector  $\mathbf{t} \in \mathbb{R}^3$  and Lie Algebra  $\nu \in \mathfrak{so}(3)$ . We run the optimisation for 10k for all scenes. The optimisation process takes 15 to 45 minutes depending on the scene size.

## 4.2. Reconstruction Quality

**Qualitative Evaluation on ScanNet [5].** We test our method on real-world sequences from the ScanNet RGB-

Method	Trans. error (meters) ↓	Rot. error (degrees) ↓
BundleFusion	0.033	0.571
COLMAP	0.038	0.692
Neural RGBD	0.023	0.146
Ours	<b>0.014</b>	<b>0.143</b>

Table 2: We evaluate the average translation and rotation errors of the estimated camera poses. Our method further improves camera poses from BundleFusion initialization and outperforms NeuralRGB-D which adopts similar pose refinement scheme.

D dataset. Depth image measurements are often noisy and miss valid depth values in areas such as glass and thin structures like chair legs. We show that our RGB loss and SDF regularisation can mitigate these effects significantly. We also jointly refine the camera poses provided in ScanNet.

Fig. 4 shows comparisons of our method with reconstructions obtained via BundleFusion [6], and NeuralRGB-D [1]. Our method produces more complete and smoother reconstruction results. For missing depth on areas such as TV screens (2nd, 4th and 5th row), our method demonstrates strong hole-filling ability and produces high-quality reconstruction. Our method also captures thin structures (bike in the 2nd row and chair leg in the 5th row) better than

	scene 0000	scene 0012	morning apartment
dimension	$8.8 \times 9.1 \times 3.4$	$5.8 \times 5.7 \times 2.9$	$3.5 \times 4.0 \times 2.3$
voxel dim	$321 \times 321 \times 129$	$225 \times 225 \times 129$	$129 \times 161 \times 97$
runtime	44 min	31 min	19 min
model size	268 MB	132 MB	41 MB
num params	66.5 M	32.7 M	10.1 M

Table 3: Performance Analysis on 3 scenes from ScanNet [5] and synthetic dataset. Our method requires much less runtime than NeuralRGB-D [1] at the cost of larger model size.

our baselines. More reconstruction results can be found in the supplementary material.

**Quantitative Evaluation on Synthetic Sequences.** We perform a quantitative evaluation of our method on a synthetic dataset of 10 sequences [1]. Noise and artifacts are added to the rendered ground-truth depth images to simulate a real depth sensor. The camera poses obtained with BundleFusion [6] are used as initialization for reconstruction, and are optimized along with the feature grids and decoders. We run marching cubes at the resolution of 1cm to extract meshes. For fair comparison, we also run BundleFusion [6] and RoutedFusion at the same resolution. Regions that are not observed in any camera views are culled before evaluation.

As shown in Tab. 1, our method performs on par with NeuralRGB-D, and outperforms all the other baselines by a large margin. Specifically, we achieve better Completion and Normal Consistency which shows the benefit of our proposed smoothness prior in terms of hole-filling and smoothness. In Tab. 2 we evaluate the pose estimation accuracy of our method. We take BundleFusion as initial camera poses and refine them along with other optimisable parameters. Results show that our pose-refinement significantly improved the initial poses and also outperforms NeuralRGB-D in terms of both translation and rotation error.

### 4.3. Performance Analysis

**Runtime.** GO-Surf achieves high-quality reconstructions in 15 to 45 minutes, while NeuralRGB-D [1] takes 15 to 25 hours on the same hardware. Tab. 3 shows our method’s runtime on 3 example scenes with various scales from ScanNet and synthetic datasets.

**Memory Footprint.** We report memory usage for storing our model in Tab. 3. Our model requires hundreds of MB to store the full feature grid and scales up rapidly as the scene size increases. This is one of our major limitations and we leave it as future work to optimize memory footprint. In contrast, our method runs much faster than the low-memory-footprint implicit representations while achieving on par reconstruction quality and still requires fewer parameters than traditional methods that use explicit voxel grid.

### 4.4. Ablation Studies

We conduct ablation studies to show the effect of individual building blocks and justify our design choices.

**Effect of RGB Term.** We compare our full model against our model without the RGB loss term. Results show that the RGB loss term enables reconstruction of high frequency details that could not be recovered using solely the depth supervision while still preserving overall smoothness (Fig. 5).

**Effect of Smoothness Term** We also verify the effectiveness of the smoothness prior term, which has a twofold role. Firstly, it completes the unobserved parts of the scene, for example holes on the floor. Secondly, it regularises the RGB and depth losses that would otherwise overfit and produce noisy surfaces (Fig. 5).

**Effect of Pose Refinement** Similarly to [1], we observe that joint optimisation of the model and camera poses corrects the imperfect trajectory estimate from SLAM system and leads to more consistent surface reconstruction (Fig. 5).

## 5. Conclusion

We presented GO-Surf, a novel approach to surface reconstruction from a sequence of RGB-D images. We achieved a level of smoothness and hole filling on-par with MLP-based approaches while reducing the training time by an order of magnitude. Our system produces accurate and complete meshes of indoor scenes.

**Limitations and Future Work.** The biggest limitation of our system is the memory consumption which at the moment scales cubically with the scene dimensions. This is a well-known problem of voxel-like architectures and various solutions have been proposed to mitigate it. Voxel hashing [18] or octree-based sparsification would significantly reduce the memory footprint of our system and we intend to explore this direction in future work. We also recognize the high complexity of our loss term which makes it difficult to fully explore the hyperparameter space. Specifically, we would like to better understand the relation between SDF, free-space and depth loss and preferably merge them into one loss term. Finally, our method currently overfits to a single scene. We are interested in learning the geometry priors on large datasets to use them at inference time.

### Acknowledgements

Research presented here has been supported by the UCL Centre for Doctoral Training in Foundational AI under UKRI grant number EP/S021566/1. TB was supported from a sponsored research award by Cisco Research. We also thank Dejan Azinovic for providing additional details and results of NeuralRGBD.

## References

- [1] Dejan Azinović, Ricardo Martin-Brualla, Dan B Goldman, Matthias Nießner, and Justus Thies. Neural rgb-d surface reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022. 1, 2, 4, 5, 6, 7, 8
- [2] Aljaz Bozic, Pablo Palafox, Justus Thies, Angela Dai, and Matthias Nießner. Transformerfusion: Monocular rgb scene reconstruction using transformers. *Advances in Neural Information Processing Systems*, 34, 2021. 2
- [3] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. *arXiv preprint arXiv:2203.09517*, 2022. 1, 3
- [4] Julian Chibane, Thiemo Alldieck, and Gerard Pons-Moll. Implicit functions in feature space for 3d shape reconstruction and completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6970–6981, 2020. 2
- [5] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, IEEE, 2017. 5, 7, 8
- [6] Angela Dai, Matthias Nießner, Michael Zollhöfer, Shahram Izadi, and Christian Theobalt. Bundlefusion: Real-time globally consistent 3d reconstruction using on-the-fly surface reintegration. *ACM Transactions on Graphics (TOG)*, 36(4):76a, 2017. 2, 5, 6, 7, 8
- [7] Amos Gropp, Lior Yariv, Niv Haim, Matan Atzmon, and Yaron Lipman. Implicit geometric regularization for learning shapes. In *International Conference on Machine Learning*, pages 3789–3799. PMLR, 2020. 2, 4, 5
- [8] Animesh Karnear, Tobias Ritschel, Oliver Wang, and Niloy J Mitra. Relu fields: The little non-linearity that could. *arXiv preprint arXiv:2205.10824*, 2022. 1, 3
- [9] Abhijit Kundu, Kyle Genova, Xiaoqi Yin, Alireza Fathi, Caroline Pantofaru, Leonidas Guibas, Andrea Tagliasacchi, Frank Dellaert, and Thomas Funkhouser. Panoptic neural fields: A semantic object-aware neural scene representation. *arXiv preprint arXiv:2205.04334*, 2022. 2
- [10] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. *Advances in Neural Information Processing Systems*, 33:15651–15663, 2020. 2
- [11] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019. 1, 2
- [12] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020. 1, 2, 3, 4
- [13] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multi-resolution hash encoding. *arXiv preprint arXiv:2201.05989*, 2022. 1, 2, 3
- [14] 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 *European Conference on Computer Vision*, pages 414–431. Springer, 2020. 2
- [15] Richard A Newcombe, Shahram Izadi, Otmar Hilliges, David Molyneaux, David Kim, Andrew J Davison, Pushmeet Kohi, Jamie Shotton, Steve Hodges, and Andrew Fitzgibbon. Kinectfusion: Real-time dense surface mapping and tracking. In *2011 10th IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*, pages 127–136. IEEE, 2011. 2
- [16] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3504–3515, 2020. 2, 3
- [17] Matthias Nießner, Michael Zollhöfer, Shahram Izadi, and Marc Stamminger. Real-time 3d reconstruction at scale using voxel hashing. *ACM Transactions on Graphics (TOG)*, 32(6):1–11, 2013. 2
- [18] M. Nießner, M. Zollhöfer, S. Izadi, and M. Stamminger. Real-time 3d reconstruction at scale using voxel hashing. *ACM Transactions on Graphics (TOG)*, 2013. 8
- [19] Michael Oechsle, Lars Mescheder, Michael Niemeyer, Thilo Strauss, and Andreas Geiger. Texture fields: Learning texture representations in function space. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4531–4540, 2019. 2
- [20] Joseph Ortiz, Alexander Clegg, Jing Dong, Edgar Sucar, David Novotny, Michael Zollhoefer, and Mustafa Mukadam. isdf: Real-time neural signed distance fields for robot perception. *arXiv preprint arXiv:2204.02296*, 2022. 2, 4
- [21] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 165–174, 2019. 1, 2
- [22] Songyou Peng, Michael Niemeyer, Lars Mescheder, Marc Pollefeys, and Andreas Geiger. Convolutional occupancy networks. In *European Conference on Computer Vision (ECCV)*, Cham, Aug. 2020. Springer International Publishing. 5
- [23] Johannes L. Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016. 5
- [24] Johannes Lutz Schönberger, True Price, Torsten Sattler, Jan-Michael Frahm, and Marc Pollefeys. A vote-and-verify strategy for fast spatial verification in image retrieval. In *Asian Conference on Computer Vision (ACCV)*, 2016. 5
- [25] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise view selection for unstructured multi-view stereo. In *European Conference on Computer Vision (ECCV)*, 2016. 5
- [26] Vincent Sitzmann, Julien N.P. Martel, Alexander W. Bergman, David B. Lindell, and Gordon Wetzstein. Implicit

- neural representations with periodic activation functions. In *arXiv*, 2020. 5
- [27] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. In *Advances in Neural Information Processing Systems*, 2019. 1, 2
- [28] Noah Stier, Alexander Rich, Pradeep Sen, and Tobias Höllerer. VoRTX: Volumetric 3d reconstruction with transformers for voxelwise view selection and fusion. In *2021 International Conference on 3D Vision (3DV)*, pages 320–330. IEEE, 2021. 2
- [29] Edgar Sucar, Shikun Liu, Joseph Ortiz, and Andrew J Davison. imap: Implicit mapping and positioning in real-time. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6229–6238, 2021. 2, 4
- [30] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. *arXiv preprint arXiv:2111.11215*, 2021. 1, 3
- [31] Jiaming Sun, Yiming Xie, Linghao Chen, Xiaowei Zhou, and Hujun Bao. Neuralrecon: Real-time coherent 3d reconstruction from monocular video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15598–15607, 2021. 2
- [32] Towaki Takikawa, Joey Litalien, Kangxue Yin, Karsten Kreis, Charles Loop, Derek Nowrouzezahrai, Alec Jacobson, Morgan McGuire, and Sanja Fidler. Neural geometric level of detail: Real-time rendering with implicit 3D shapes. 2021. 3
- [33] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *NeurIPS*, 2021. 2, 3, 4
- [34] Silvan Weder, Johannes L. Schönberger, Marc Pollefeys, and Martin R. Oswald. Routedfusion: Learning real-time depth map fusion. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 5
- [35] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces. In *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021. 3
- [36] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Basri Ronen, and Yaron Lipman. Multiview neural surface reconstruction by disentangling geometry and appearance. *Advances in Neural Information Processing Systems*, 33:2492–2502, 2020. 2, 3
- [37] Alex Yu, Sara Fridovich-Keil, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. *arXiv preprint arXiv:2112.05131*, 2021. 1, 3
- [38] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenotrees for real-time rendering of neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5752–5761, 2021. 2
- [39] Xiaoshuai Zhang, Sai Bi, Kalyan Sunkavalli, Hao Su, and Zexiang Xu. Nerfusion: Fusing radiance fields for large-scale scene reconstruction. *CVPR*, 2022. 2
- [40] Shuaifeng Zhi, Tristan Laidlow, Stefan Leutenegger, and Andrew J Davison. In-place scene labelling and understanding with implicit scene representation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15838–15847, 2021. 2
- [41] Zihan Zhu, Songyou Peng, Viktor Larsson, Weiwei Xu, Hujun Bao, Zhaopeng Cui, Martin R. Oswald, and Marc Pollefeys. Nice-slam: Neural implicit scalable encoding for slam. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 1, 2, 3, 4

# Supplementary Material

## GO-Surf: Neural Feature Grid Optimization for Fast, High-Fidelity RGB-D Surface Reconstruction

Jingwen Wang\*, Tymoteusz Bleja\*, and Lourdes Agapito

Department of Computer Science, University College London

### 1. Second-order Grid Sampler

In this section we provide additional details on our implementation of grid sampler that allows double back-propagation in PyTorch. As described in our main paper, the predicted SDF value  $\phi$  is a function of 3D coordinates of the query point  $\mathbf{x}$ , feature vectors in feature grid  $\theta$ , and geometry network parameters  $\omega$ :

$$\phi(\mathbf{x}, \theta, \omega) = f_{\omega}(\mathbf{z}(\theta, \mathbf{x})) \quad (1)$$

where  $\mathbf{z} = \mathcal{V}_{\theta}(\mathbf{x}) \in \mathbb{R}^C$  is the tri-linearly interpolated feature vector at the query point  $\mathbf{x}$ . The first-order gradient of the SDF w.r.t. 3D query coordinates is:

$$\begin{aligned} \frac{\partial \phi}{\partial \mathbf{x}} &= \frac{\partial f_{\omega}(\mathbf{z})}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{x}} \\ &= \sum_{i=1}^C \frac{\partial f_{\omega}(\mathbf{z})}{\partial z_i} \frac{\partial z_i}{\partial \mathbf{x}} \end{aligned} \quad (2)$$

Here, to avoid the usage of higher-order tensors in derivation of second-order derivatives we consider the derivatives for each individual feature dimension  $z_i$ , as in Eq. 2. To regularise this gradient term we need to obtain the gradient Eq. 2 w.r.t.  $\mathbf{x}$ ,  $\theta$  and  $\omega$  respectively, which corresponds to the first three rows of the Hessian of SDF:

$$\frac{\partial^2 \phi}{\partial \mathbf{x}^2} = \sum_{i=1}^C \frac{\partial^2 f_{\omega}(\mathbf{z})}{\partial z_i^2} \left( \frac{\partial z_i}{\partial \mathbf{x}} \right)^2 + \frac{\partial f_{\omega}(\mathbf{z})}{\partial z_i} \frac{\partial^2 z_i}{\partial \mathbf{x}^2} \quad (3)$$

$$\frac{\partial^2 \phi}{\partial \mathbf{x} \partial \theta} = \sum_{i=1}^C \frac{\partial^2 f_{\omega}(\mathbf{z})}{\partial z_i^2} \frac{\partial z_i}{\partial \theta} \frac{\partial z_i}{\partial \mathbf{x}} + \frac{\partial f_{\omega}(\mathbf{z})}{\partial z_i} \frac{\partial^2 z_i}{\partial \mathbf{x} \partial \theta} \quad (4)$$

$$\frac{\partial^2 \phi}{\partial \mathbf{x} \partial \omega} = \sum_{i=1}^C \frac{\partial^2 f_{\omega}(\mathbf{z})}{\partial z_i \partial \omega} \frac{\partial z_i}{\partial \mathbf{x}} \quad (5)$$

We don't need to implement Eq. 5 as Pytorch's automatic differentiation package supports double back-propagation

\* The first two authors contributed equally.

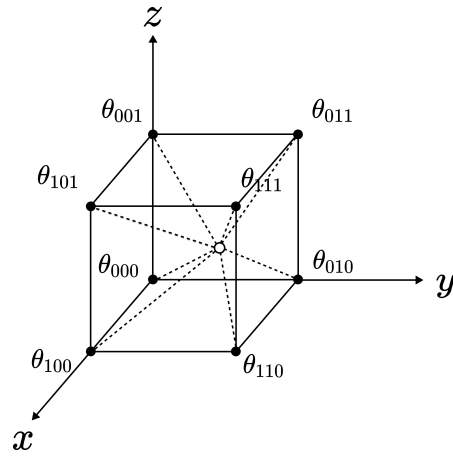


Figure 1: Example of trilinear interpolation in a 3D voxel.

through an MLP and Pytorch's `grid_sampler` also has first-order gradient implementation. However Eq. 3 and 4 need to be implemented as they require double back-propagation through the `grid_sampler` which is not implemented in PyTorch.

**Derivative Derivation.** Now we will derive the second-order derivatives of the tri-linear interpolation. More specifically, we only need two Hessian blocks:  $\frac{\partial^2 z_i}{\partial \mathbf{x}^2}$  and  $\frac{\partial^2 z_i}{\partial \mathbf{x} \partial \theta}$ . For simplicity we derive with 1-D features, but it generalises to any feature dimensions as the derivation is equivalent to all dimensions. Assume we have a query point  $\mathbf{x} = [x, y, z]$  and 8 feature vectors (points) stacking in a column  $\theta = [\theta_{000}, \theta_{001}, \theta_{010}, \theta_{011}, \theta_{100}, \theta_{101}, \theta_{110}, \theta_{111}]^T \in \mathbb{R}^8$  that correspond to the eight vertices of the voxel that encloses the point  $\mathbf{x}$ . Then, the tri-linearly interpolated feature  $f$  at point  $\mathbf{x}$  is given by:

$$f = \theta^T \mathbf{w}(\mathbf{x}) \quad (6)$$

where  $\mathbf{w}(\mathbf{x}) \in \mathbb{R}^8$  is the coefficient vector of the 8 vertices:

$$\mathbf{w}(\mathbf{x}) = \begin{bmatrix} (1-x)(1-y)(1-z) \\ (1-x)(1-y)z \\ (1-x)y(1-z) \\ (1-x)yz \\ x(1-y)(1-z) \\ x(1-y)z \\ xy(1-z) \\ xyz \end{bmatrix} \quad (7)$$

Then the Jacobian of  $f$  w.r.t.  $\mathbf{x}$  is given by:

$$\frac{\partial f}{\partial \mathbf{x}} = \theta^T \frac{\partial \mathbf{w}(\mathbf{x})}{\partial \mathbf{x}} \quad (8)$$

where the Jacobian of the coefficient vector  $\mathbf{w}$  w.r.t. query point  $\mathbf{x}$  is a  $8 \times 3$  matrix and is given by:

$$\frac{\partial \mathbf{w}}{\partial \mathbf{x}} = [\mathbf{J}_1 \quad \mathbf{J}_2 \quad \mathbf{J}_3] \quad (9)$$

$$\mathbf{J}_1 = \begin{bmatrix} -(1-y)(1-z) \\ -(1-y)z \\ -y(1-z) \\ -yz \\ (1-y)(1-z) \\ (1-y)z \\ y(1-z) \\ yz \end{bmatrix} \quad (10)$$

$$\mathbf{J}_2 = \begin{bmatrix} -(1-x)(1-z) \\ -(1-x)z \\ (1-x)(1-z) \\ (1-x)z \\ -x(1-z) \\ -xz \\ x(1-z) \\ xz \end{bmatrix} \quad (11)$$

$$\mathbf{J}_3 = \begin{bmatrix} -(1-x)(1-y) \\ (1-x)(1-y) \\ -(1-x)y \\ (1-x)y \\ -x(1-y) \\ x(1-y) \\ -xy \\ xy \end{bmatrix} \quad (12)$$

Then the second-order derivative  $\frac{\partial^2 f}{\partial \mathbf{x} \partial \theta}$  is given by:

$$\begin{aligned} \frac{\partial^2 f}{\partial \mathbf{x} \partial \theta} &= \frac{\partial}{\partial \theta} \left( \frac{\partial f}{\partial \mathbf{x}^T} \right) \\ &= \left( \frac{\partial \mathbf{w}}{\partial \mathbf{x}} \right)^T \end{aligned} \quad (13)$$

which is simply the transpose of Eq. 9. And to obtain the

other second-order derivative  $\frac{\partial^2 f}{\partial \mathbf{x}^2}$ , we just need to differentiate Eq. 8 further w.r.t. the query points:

$$\begin{aligned} \frac{\partial^2 f}{\partial \mathbf{x}^2} &= \frac{\partial}{\partial \mathbf{x}} \left( \frac{\partial f}{\partial \mathbf{x}^T} \right) \\ &= \frac{\partial}{\partial \mathbf{x}} \begin{bmatrix} \mathbf{J}_1^T \theta \\ \mathbf{J}_2^T \theta \\ \mathbf{J}_3^T \theta \end{bmatrix} \end{aligned} \quad (14)$$

$$= \begin{bmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{33} \end{bmatrix} \quad (15)$$

where  $\{\mathbf{J}_i^T \theta\}_{i=1}^3$  are inner-products between each column of Eq. 9 and  $\theta$ , so the result is a  $3 \times 3$  matrix with each column being the derivative w.r.t.  $x$ ,  $y$  and  $z$ . It is trivial to show that all diagonal elements are zero as each  $\mathbf{J}_i$  term only contains the variables of other two dimensions, as in Eq. 10, 11 and 12. The off-diagonal elements can be easily computed as:

$$\begin{aligned} h_{12} = h_{21} &= (1-z)(\theta_{000} + \theta_{110} - \theta_{010} - \theta_{100}) \\ &\quad + z(\theta_{001} + \theta_{111} - \theta_{011} - \theta_{101}) \end{aligned} \quad (16)$$

$$\begin{aligned} h_{13} = h_{31} &= (1-y)(\theta_{000} + \theta_{101} - \theta_{001} - \theta_{100}) \\ &\quad + y(\theta_{010} + \theta_{111} - \theta_{011} - \theta_{110}) \end{aligned} \quad (17)$$

$$\begin{aligned} h_{23} = h_{32} &= (1-x)(\theta_{000} + \theta_{011} - \theta_{001} - \theta_{010}) \\ &\quad + x(\theta_{100} + \theta_{111} - \theta_{101} - \theta_{110}) \end{aligned} \quad (18)$$

## 2. Per-scene Quantitative Evaluation

In this section we provide additional per-scene breakdown of our quantitative evaluation on the 10 synthetic sequences.

**Evaluation Protocol.** For all the methods we run marching cubes at 1cm resolution to extract the meshes for evaluation. We measure accuracy (Acc), completion (Comp), chamfer- $\ell_1$  ( $C-\ell_1$ ), normal consistency (NC) and F-score for evaluation of reconstruction quality. Specifically, all the metrics are computed between point clouds sampled on ground-truth and predicted mesh. Instead of sampling a fixed number of points we sample point cloud at density of 1 point per  $\text{cm}^2$  to take into account the scene scale. The threshold for computing F-score is set to 5 cm.

**Mesh Culling.** To prevent the evaluation from falsely penalizing the scene completion ability of our method, surfaces that are not observed in RGB-D images are culled. Following [1] we subdivide the meshes such that all the faces have maximum edge length of below 1.5cm. A face will be removed if 1. it is not inside any camera frusta, or 2. it is occluded by other geometry, or 3. it has no valid depth measurements. Note for thin geometry sequence, we only apply the first two criteria.






Scene	Method	Acc ↓	Comp ↓	C- $\ell_1$ ↓	NC ↑	F-score ↑	Trans. ↓	Rot. ↓
<b>Complete kitchen</b> 	BundleFusion	0.0303	0.1475	0.0889	0.8570	0.6943	0.050	0.566
	RoutedFusion	0.0270	0.0854	0.0562	0.8484	0.7939	-	-
	COLMAP	0.0365	0.0354	0.0360	0.9245	0.8248	0.009	0.210
	ConvOccNets	0.0502	0.0527	0.0514	0.8667	0.6610	-	-
	SIREN	0.0319	0.0700	0.0509	0.9031	0.7415	-	-
	NeuralRGBD	<b>0.0224</b>	0.0394	0.0309	0.9098	0.8962	0.083	0.450
	Ours	<b>0.0224</b>	<b>0.0258</b>	<b>0.0241</b>	<b>0.9413</b>	<b>0.8998</b>	<b>0.017</b>	<b>0.137</b>
	BundleFusion	0.0253	0.0578	0.0416	0.9112	0.7967	0.038	0.327
	RoutedFusion	0.0281	0.0362	0.0322	0.8553	0.8484	-	-
	COLMAP	0.0228	0.0282	0.0255	0.9332	0.9170	0.103	0.641
<b>Kitchen</b> 	ConvOccNets	0.0420	0.049	0.0455	0.8752	0.6253	-	-
	SIREN	0.0327	0.0575	0.0451	0.8996	0.7071	-	-
	NeuralRGBD	0.0218	0.0297	0.0257	0.9296	0.9005	0.030	<b>0.114</b>
	Ours	<b>0.0214</b>	<b>0.0271</b>	<b>0.0243</b>	<b>0.9316</b>	<b>0.9379</b>	<b>0.026</b>	0.145
<b>Breakfast room</b> 	BundleFusion	0.0129	0.0235	0.0182	0.9582	0.9606	0.037	0.697
	RoutedFusion	0.0181	0.0202	0.0191	0.9341	0.9758	-	-
	COLMAP	0.0191	0.0194	0.0192	0.9522	0.9533	0.009	0.210
	ConvOccNets	0.0311	0.0329	0.0320	0.8925	0.9602	-	-
	SIREN	0.0150	0.0454	0.0302	0.9371	0.9230	-	-
	NeuralRGBD	0.0145	0.0148	0.0146	<b>0.9657</b>	<b>0.9898</b>	<b>0.007</b>	<b>0.135</b>
	Ours	<b>0.0144</b>	<b>0.0136</b>	<b>0.0139</b>	0.9629	0.9829	0.009	0.137
<b>Morning apartment</b> 	BundleFusion	<b>0.0079</b>	0.0146	0.0112	0.8891	0.9740	0.008	0.165
	RoutedFusion	0.0100	0.0143	0.0121	0.8754	0.9795	-	-
	COLMAP	0.0133	0.0183	0.0158	0.8810	0.9666	0.017	0.380
	ConvOccNets	0.0408	0.0482	0.0445	0.8105	0.7912	-	-
	SIREN	0.0105	0.0146	0.0125	0.8765	0.9718	-	-
	NeuralRGBD	0.0087	<b>0.0121</b>	<b>0.0104</b>	<b>0.8918</b>	<b>0.9866</b>	<b>0.005</b>	<b>0.093</b>
	Ours	0.0095	0.0129	0.0112	0.8874	0.9778	<b>0.005</b>	0.101
<b>Grey white room</b> 	BundleFusion	0.0297	0.0456	0.0377	0.8612	0.7537	0.056	1.891
	RoutedFusion	0.0303	0.0347	0.0325	0.8531	0.7908	-	-
	COLMAP	0.0287	0.0293	0.0290	0.9013	0.9036	0.029	0.296
	ConvOccNets	0.0470	0.0488	0.0479	0.8434	0.6057	-	-
	SIREN	0.0323	0.0335	0.0329	0.8697	0.8142	-	-
	NeuralRGBD	<b>0.0132</b>	<b>0.0151</b>	<b>0.0142</b>	<b>0.9318</b>	<b>0.9923</b>	0.014	<b>0.146</b>
	Ours	0.0140	0.0158	0.0149	0.9261	0.9895	<b>0.013</b>	0.205

Table 1: We compare the quality of our reconstruction on several synthetic scenes for which ground truth data is available. We measure accuracy, completion, chamfer  $\ell_1$  distance, normal consistency and F-score for reconstruction quality. The numbers are computed between point clouds sampled with a density of 1 point per  $\text{cm}^2$ .




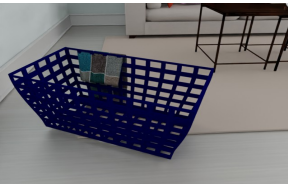
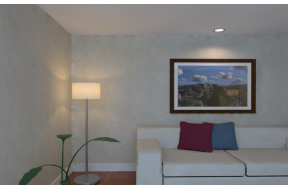
Scene	Method	Acc ↓	Comp ↓	C- $\ell_1$ ↓	NC ↑	F-score ↑	Trans. ↓	Rot. ↓
<b>White room</b> 	BundleFusion	0.0276	0.0918	0.0597	0.8788	0.7286	0.045	0.375
	RoutedFusion	0.0289	0.0430	0.0360	0.8280	0.8222	-	-
	COLMAP	0.0309	0.0342	0.0325	0.9188	0.8259	<b>0.018</b>	0.167
	ConvOccNets	0.0537	0.0583	0.0560	0.8653	0.5012	-	-
	SIREN	0.0276	0.0588	0.0432	0.8992	0.7788	-	-
	NeuralRGBD	<b>0.0204</b>	<b>0.0256</b>	<b>0.0230</b>	<b>0.9297</b>	<b>0.9551</b>	0.028	<b>0.146</b>
	Ours	0.0210	0.0325	0.0268	0.9281	0.9233	<b>0.024</b>	0.157
<b>Green room</b> 	BundleFusion	0.0118	0.0339	0.0228	0.9254	0.9314	0.027	0.546
	RoutedFusion	0.0156	0.0193	0.0174	0.9095	0.9735	-	-
	COLMAP	0.0159	0.0194	0.0177	0.9270	0.9712	0.014	0.227
	ConvOccNets	0.0548	0.0493	0.0521	0.8600	0.7434	-	-
	SIREN	0.0183	0.0253	0.0218	0.9143	0.9448	-	-
	NeuralRGBD	<b>0.0106</b>	<b>0.0142</b>	<b>0.0124</b>	<b>0.9348</b>	<b>0.9913</b>	<b>0.012</b>	0.104
	Ours	0.0138	0.0169	0.0153	0.9256	0.9838	0.014	<b>0.085</b>
<b>Staircase</b> 	BundleFusion	0.0257	0.1146	0.0701	0.8792	0.7108	0.039	0.643
	RoutedFusion	0.0411	0.0512	0.0461	0.8909	0.6896	-	-
	COLMAP	0.0454	0.058	0.0517	0.9253	0.6875	0.043	0.305
	ConvOccNets	0.0618	0.0562	0.059	0.8601	0.5646	-	-
	SIREN	0.0355	0.0514	0.0434	0.9117	0.7487	-	-
	NeuralRGBD	<b>0.0216</b>	<b>0.0254</b>	<b>0.0235</b>	0.9471	<b>0.9333</b>	0.016	0.123
	Ours	0.0221	0.0257	0.024	<b>0.9496</b>	0.9235	<b>0.015</b>	0.144
<b>Thin geometry</b> 	BundleFusion	0.0072	0.0305	0.0188	0.9063	0.9199	<b>0.009</b>	0.126
	RoutedFusion	<b>0.0070</b>	0.0396	0.0233	0.8243	0.8785	-	-
	COLMAP	0.0372	0.0558	0.0465	0.8181	0.7209	0.079	2.4
	ConvOccNets	0.0115	0.0329	0.0222	0.8800	0.9072	-	-
	SIREN	0.0086	0.0335	0.0210	0.8823	0.9115	-	-
	NeuralRGBD	0.0079	<b>0.0092</b>	<b>0.0086</b>	<b>0.9077</b>	<b>0.9956</b>	<b>0.010</b>	<b>0.037</b>
	Ours	0.0093	0.0121	0.0107	0.8986	0.9817	0.011	0.146
<b>ICL living room</b> 	BundleFusion	0.0129	0.0214	0.0172	0.9606	0.9694	0.022	0.382
	RoutedFusion	0.0168	0.0201	0.0185	0.9456	0.9841	-	-
	COLMAP	0.0209	0.0238	0.0224	0.9528	0.9730	0.029	0.836
	ConvOccNets	0.1049	0.0956	0.1003	0.8535	0.5619	-	-
	SIREN	0.0167	0.0219	0.0193	0.9555	0.9734	-	-
	NeuralRGBD	<b>0.0095</b>	<b>0.0115</b>	<b>0.0105</b>	<b>0.9689</b>	<b>0.9944</b>	<b>0.007</b>	<b>0.109</b>
	Ours	0.0105	0.0127	0.0117	0.9661	0.9909	<b>0.007</b>	0.167

Table 2: We compare the quality of our reconstruction on several synthetic scenes for which ground truth data is available. We measure accuracy, completion, chamfer  $\ell_1$  distance, normal consistency and F-score for reconstruction quality. The numbers are computed between point clouds sampled with a density of 1 point per  $\text{cm}^2$ .

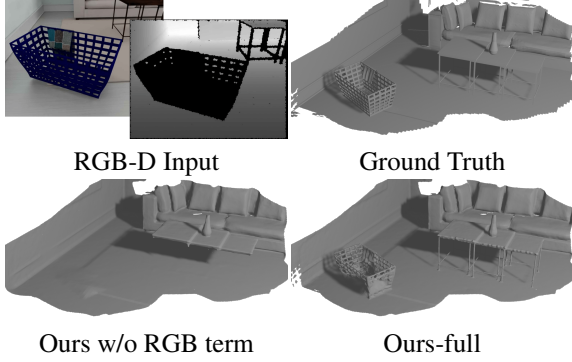


Figure 2: Qualitative reconstruction results on synthetic sequence with thin geometries and missing depth measurements.

Method	Acc.) ↓	Com. ↓	C- $\ell_1$ ↓	NC ↑	F-score ↑
Ours (no rgb)	<b>0.0087</b>	0.0291	0.0189	0.8967	0.9273
Ours (full)	0.0093	<b>0.0121</b>	<b>0.0107</b>	<b>0.8986</b>	<b>0.9817</b>

Table 3: We evaluate the reconstruction results of our model with out RGB loss term and our full model on a synthetic sequence with thin geometries.

**Synthetic Dataset.** We evaluate on the same synthetic dataset as in [1] which consists of 10 scenes published under the CC-BY or CC-0 license. For more details please refer to their paper [1].

### 3. More Ablation Studies

In this section, we show additional ablation studies on the effect of RGB loss term, regularisation terms .

#### 3.1. Effect of RGB Loss Term

Similar to [1], in the main paper we showed in Fig. 5 that the RGB loss term is able to capture better high-frequency details and recover missing depth regions. In this section we show quantitative evaluation on a synthetic scene with thin geometries that have no depth measurements, and also show more qualitative ablation results on real-world ScanNet scenes.

For the experiment on the synthetic scene, we simulate missing depth by removing the depth measurements from the baskets and table legs (top left corner in Fig. 2). Comparison in Fig. 2 demonstrates that RGB loss term is able to recover thin structures with missing depth and produce complete reconstruction. Tab. 3 also shows our full model achieves significantly better reconstruction quality, especially completeness.

#### 3.2. Effect of SDF Regularisation Terms

In this section we provide additional ablation studies on the two SDF regularisation terms.

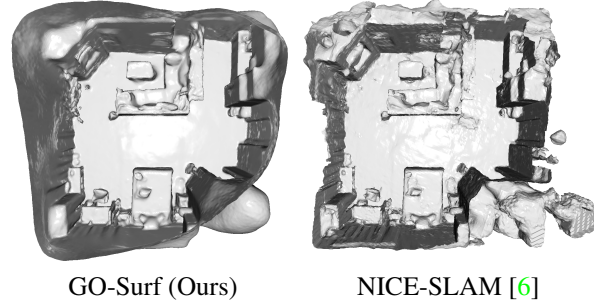


Figure 3: Reconstructed meshes (unculled) of sequential mapping on ScanNet scene0000. Note that for fair comparison NICE-SLAM is running in mapping mode with ground truth camera poses. GO-Surf runs over  $\times 30$  faster while achieving much smoother high-quality reconstruction.

**Eikonal Term** The Eikonal term encourages the model prediction to be a valid signed distance field. We observe that without the Eikonal term the SDF values are being reconstructed only within the truncation region (i.e. very close to the surface). The Eikonal term helps correct SDF values to be propagated outside of the truncation region, although some local artefacts still remain (see Fig. 4) which we suspect are due to the local nature of gradients in the feature grid.

**Normal Smoothness Term** We found that the normal smoothness term fills the holes in unobserved regions which is particularly useful for fixing discontinuities in large planar structures like walls or floor. It also encourages surface smoothness in all other areas, usually at the cost of some fine details. We investigate different settings of normal smoothness radius (see Fig. 5).

### 4. More Results

In this section, we show more qualitative reconstruction results on synthetic scenes. Note that in order to also showcase the scene completion ability of different methods, the results shown here are from unculled meshes. Fig. 6 shows the comparison of our method to NeuralRGB-D and other learning-based methods. Overall, our method produces smoother, and more complete reconstruction without losing tiny details.

From the first two columns it can be seen that both NeuralRGB-D and our method have the ability to fill in unobserved regions (windows and the top ceiling). However, NeuralRGB-D tends to produce noisier scene completion results with many artefacts whereas ours is much smoother and looks more natural. In Fig 7 we provide more results to further showcase our method’s advantage over NeuralRGB-D in terms of scene completion.

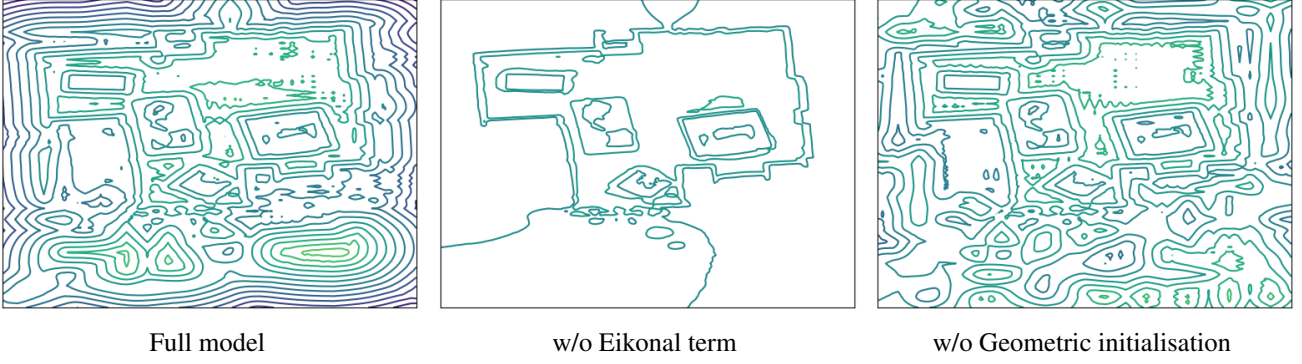


Figure 4: We found that applying the Eikonal term to some extent propagates the SDF gradient to empty spaces. Without the Eikonal term SDF is reconstructed only within the truncation region. Geometric initialisation reduces the noise in the unobserved regions (e.g. behind walls).

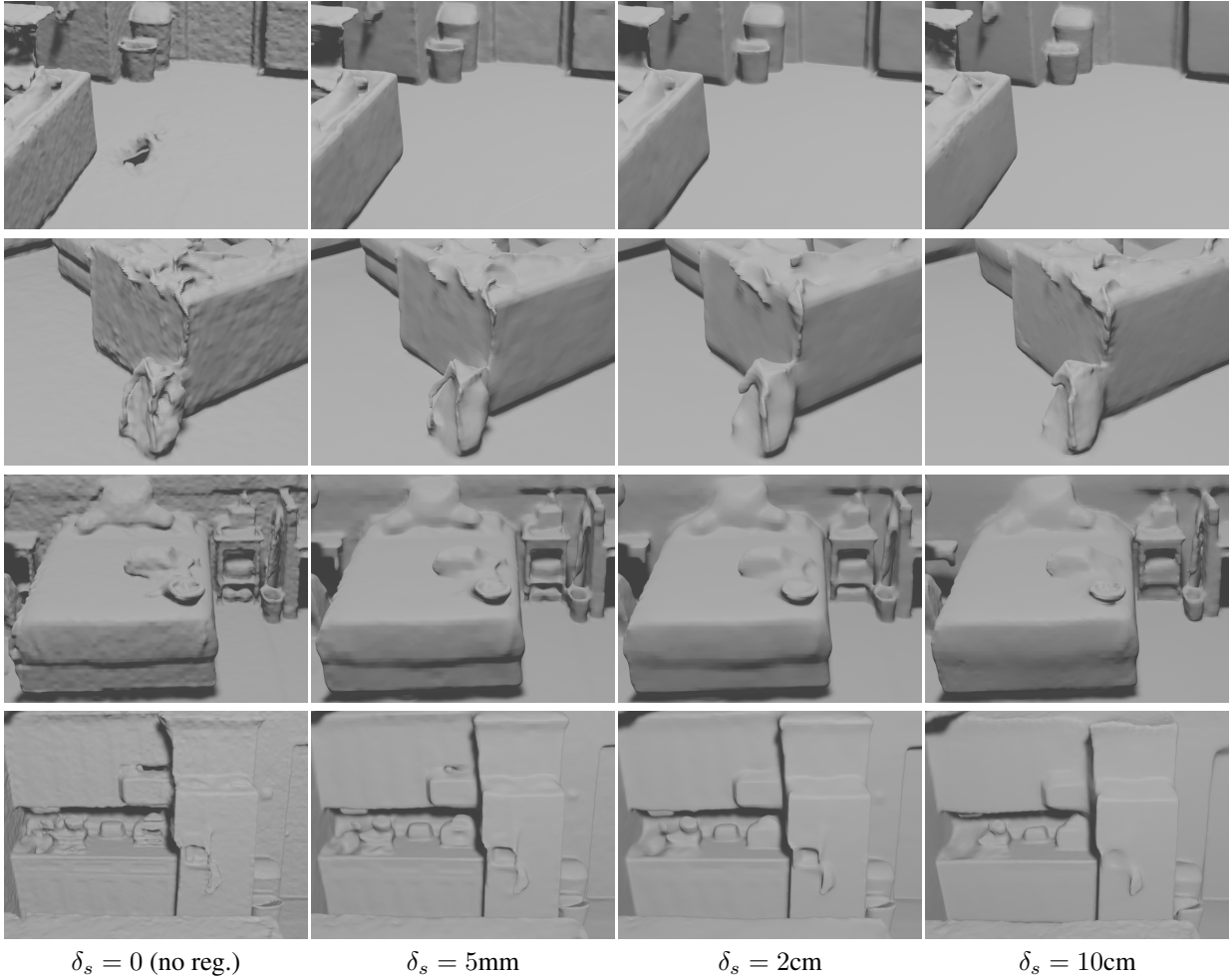


Figure 5: Ablation study on effect of smoothness prior. We compare the reconstruction results with different values of  $\delta_s$  for normal regularisation term. With  $\delta_s = 0$  the results capture more details like the backpack strap and tiny objects on the kitchen table, but also have more high frequency noises. Also the hole shows up near the sofa. Increasing  $\delta_s$  to 10cm results in over-smoothed reconstruction. We empirically found 2 – 5mm is a good trade-off between smoothness and details.



Figure 6: Qualitative comparison on synthetic scenes. GO-Surf produces smoother, and more complete reconstruction without losing tiny details. Both NeuralRGBD and our method have the scene completion ability. However, NeuralRGBD produces much noisier completion results with many artefacts whereas ours is much smoother and looks more natural.

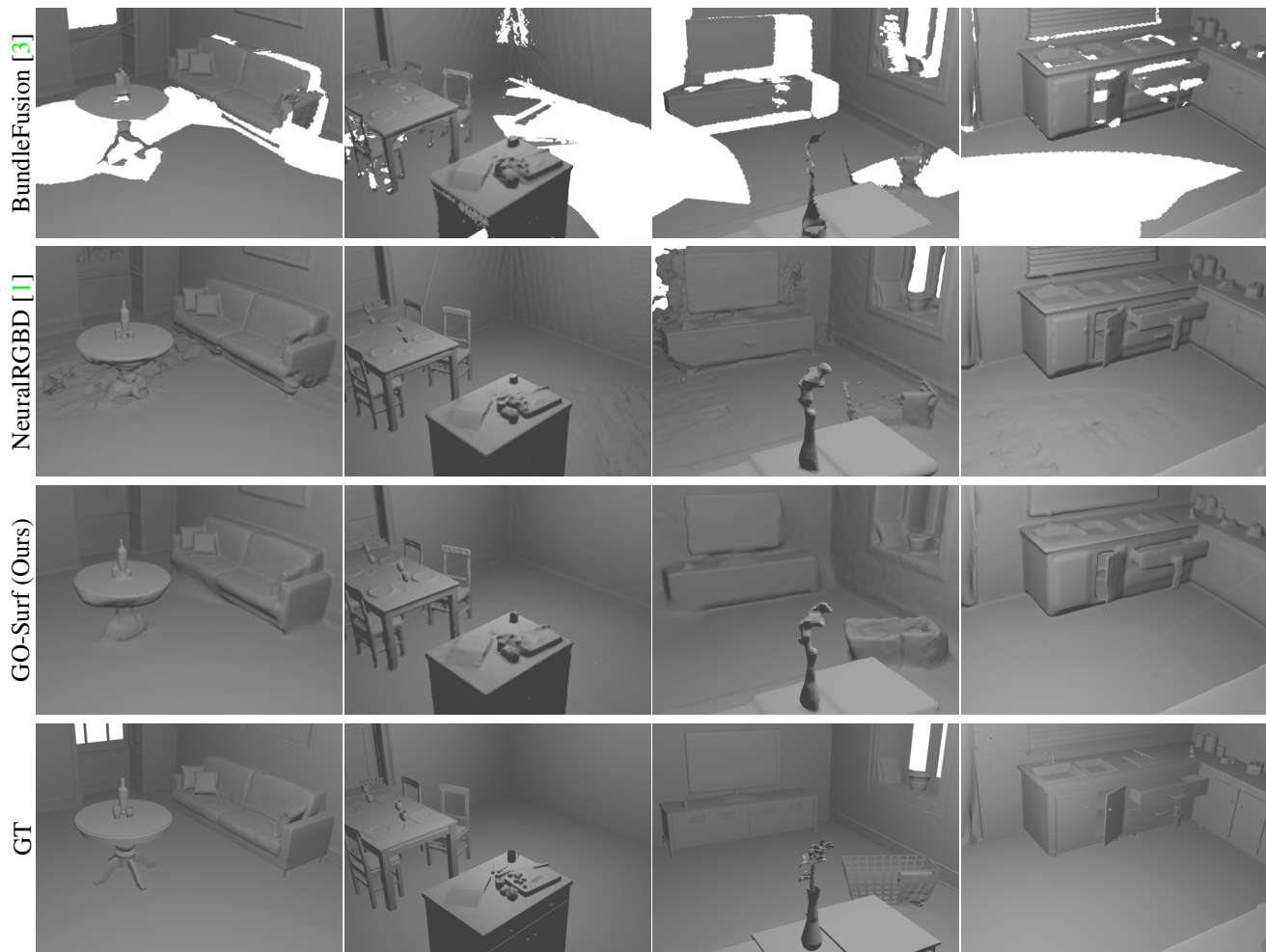


Figure 7: Scene completion comparison. NeuralRGBD produces much noisier completion results with many artefacts whereas ours is much smoother and looks more natural.

scene	scene dim.	voxel dim.	runtime	num params.	model size
complete kitchen	$9.3 \times 10.0 \times 3.3$	$321 \times 353 \times 129$	45 min	73.1 M	294 MB
kitchen	$7.0 \times 8.7 \times 3.4$	$257 \times 321 \times 129$	40 min	53.2 M	214 MB
breakfast room	$4.1 \times 4.7 \times 3.3$	$161 \times 193 \times 129$	25 min	20.1 M	81 MB
morning apartment	$3.5 \times 4.0 \times 2.3$	$129 \times 161 \times 97$	19 min	10.1 M	41 MB
grey white room	$5.9 \times 4.4 \times 3.1$	$225 \times 161 \times 129$	24 min	23.4 M	94 MB
white room	$5.6 \times 7.9 \times 3.8$	$193 \times 289 \times 161$	39 min	44.9 M	181 MB
staircase	$6.8 \times 6.5 \times 3.7$	$257 \times 225 \times 161$	39 min	46.6 M	188 MB
green room	$8.0 \times 4.7 \times 3.1$	$289 \times 193 \times 129$	30 min	36.0 M	145 MB
thin geometry	$3.4 \times 3.6 \times 1.2$	$129 \times 129 \times 65$	15 min	5.4 M	28 MB
ICL living room	$5.3 \times 2.9 \times 5.4$	$193 \times 193 \times 129$	24 min	24.1 M	97 MB
scene 0000	$8.8 \times 9.1 \times 3.4$	$321 \times 321 \times 129$	44 min	66.5 M	268 MB
scene 0002	$4.4 \times 6.0 \times 3.5$	$161 \times 225 \times 129$	28 min	23.4 M	113 MB
scene 0005	$5.8 \times 5.3 \times 2.9$	$225 \times 193 \times 125$	33 min	27.1 M	113 MB
scene 0012	$5.8 \times 5.7 \times 2.9$	$225 \times 225 \times 129$	31 min	32.7 M	132 MB
scene 0024	$7.6 \times 8.4 \times 2.9$	$289 \times 289 \times 129$	40 min	53.9 M	217 MB
scene 0050	$5.9 \times 4.5 \times 3.1$	$225 \times 161 \times 129$	26 min	23.4 M	94 MB

Table 4: Runtime and memory requirement on synthetic and ScanNet [2] scenes. We list scene dimension ( $\text{m}^3$ ), voxel dimension, runtime (min), number of model parameters and model size (MB). Our method requires much less runtime than NeuralRGBD [1] (15 – 45 min vs 15 – 25 hours) at the cost of larger model size (tens to hundreds of MB vs several MB).

## 5. Runtime and Memory Analysis

In Tab. 4 we provide detailed breakdown of the runtime and memory usage of our method on all of the scenes that we had run experiments on. Our method requires 15 to 45 minutes to converge which is approximately  $\times 60$  faster than NeuralRGB-D. However our model size is much larger (tens to hundreds million of parameters vs several million parameters) and scales rapidly as the scene gets larger. This is an inherent problem of voxel-like architectures as the number of voxels scales cubically with the scene dimension. Voxel hashing or octree-based sparsification could be adopted to significantly reduce the memory footprint of our system and we intend to explore this direction in future work.

## 6. Sequential Mapping Experiments

GO-Surf’s fast runtime also enables sequential/online mapping at interactive framerate with slightly reduced resolution. In this section we run GO-Surf in a sequential fashion using 3-level feature grid with voxel sizes of 6cm, 24cm and 96cm on ScanNet scene0000 and compare against NICE-SLAM [6] running on the same scene with ground truth camera poses.

With increased voxel size GO-Surf runs in near real-time at 15Hz while NICE-SLAM runs much slower below 0.5Hz. Fig 3 shows the comparison of reconstruction results. Ours is much smoother and has less artifacts.

## References

- [1] Dejan Azinović, Ricardo Martin-Brualla, Dan B Goldman, Matthias Nießner, and Justus Thies. Neural rgb-d surface reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022. 2, 5, 7, 8, 9
- [2] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, IEEE, 2017. 9
- [3] Angela Dai, Matthias Nießner, Michael Zollhöfer, Shahram Izadi, and Christian Theobalt. Bundlefusion: Real-time globally consistent 3d reconstruction using on-the-fly surface reintegration. *ACM Transactions on Graphics (TOG)*, 36(4):76a, 2017. 7, 8
- [4] Songyou Peng, Michael Niemeyer, Lars Mescheder, Marc Pollefeys, and Andreas Geiger. Convolutional occupancy networks. In *European Conference on Computer Vision (ECCV)*, Cham, Aug. 2020. Springer International Publishing. 7
- [5] Vincent Sitzmann, Julien N.P. Martel, Alexander W. Bergman, David B. Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. In *arXiv*, 2020. 7
- [6] Zihan Zhu, Songyou Peng, Viktor Larsson, Weiwei Xu, Hujun Bao, Zhaopeng Cui, Martin R. Oswald, and Marc Pollefeys. Nice-slam: Neural implicit scalable encoding for slam. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 5, 9