# NeuralBlox: Real-Time Neural Representation Fusion for Robust Volumetric Mapping

Stefan Lionar*      Lukas Schmid*      Cesar Cadena      Roland Siegwart      Andrei Cramariuc

Autonomous Systems Lab, ETH Zürich, Switzerland

splionar@gmail.com    {schmluk, cesarc, rsiegwart, crandrei}@ethz.ch

## Abstract

*We present a novel 3D mapping method leveraging the recent progress in neural implicit representation for 3D reconstruction. Most existing state-of-the-art neural implicit representation methods are limited to object-level reconstructions and can not incrementally perform updates given new data. In this work, we propose a fusion strategy and training pipeline to incrementally build and update neural implicit representations that enable the reconstruction of large scenes from sequential partial observations. By representing an arbitrarily sized scene as a grid of latent codes and performing updates directly in latent space, we show that incrementally built occupancy maps can be obtained in real-time even on a CPU. Compared to traditional approaches such as Truncated Signed Distance Fields (TS-DFs), our map representation is significantly more robust in yielding a better scene completeness given noisy inputs. We demonstrate the performance of our approach in thorough experimental validation on real-world datasets with varying degrees of added pose noise.*

## 1. Introduction

Mapping the 3D volume of an environment during online operation is a fundamental capability required for various robotic tasks, ranging from autonomous navigation [46] to mobile manipulation [2]. Typically, volumetric maps are built by fusing multiple range measurements into grid-based occupancy or signed distance maps [20, 35]. However, such measurements are often affected adversely by uncertainties arising from sensor noise and pose estimation errors [45]. To ensure safety during robot operation, *e.g.* obstacle avoidance, these uncertainties must be accurately taken into account in the 3D maps.

Classical approaches for robotic volumetric mapping often rely on recursive updates of occupancy probabilities [20] or Signed Distance Functions (SDFs) [10, 35], which both depend heavily on hand-tuned parameters for
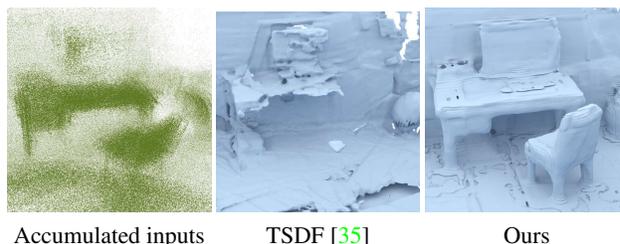
---
*Authors share first authorship.



Figure 1. Reconstructions of sequential observations with high pose uncertainties. Our method leverages shape priors and geometric context for the fusion, resulting in a robust reconstruction given noisy input.

specific sensors, scenes, and uncertainties. These approaches apply their update rules without any geometric or semantic context. Consequently, regular occurrence of noise can heavily afflict the resulting SDFs or occupancy probabilities and break the geometry of reconstructed objects in the map.

Fig. 1 illustrates this undesirable behavior of a TSDF-based method [35]. Given sequential observations with high pose uncertainties, the TSDF-based method fails to reconstruct objects in a scene, increasing the risk of collision of the robot and potentially rendering the map insufficient for manipulation tasks. To overcome this limitation, we propose to use a different map representation to fuse range measurements and provide robust volumetric mapping amidst measurement uncertainties. Motivated by the recent progress in neural implicit representation for 3D reconstruction [37, 5, 38, 25, 1, 52], we consider representing scenes implicitly as a composition of latent codes that are sequentially encoded from the inputs. These latent codes can then be decoded into other representations such as occupancy grids or SDFs during operation. This enables a broad variety of robotic path planning or interaction algorithms to directly interface with and operate in real-time on our new implicit volumetric mapping framework.

In this work, we propose a system that given sparse unoriented point clouds from *e.g.* a LiDAR or RGBD sensor produces occupancy probabilities as the final output. We

represent the map as a dynamically growing grid of large voxels, where each voxel contains a latent code representing the local geometry. Given new measurements, we encode the points into latent implicit representations and perform the map update directly in latent space, instead of updating only the occupancy probabilities. We show that our 3D mapping approach can accurately capture noisy measurements and preserve the overall geometry of the scenes better than classical approaches, as shown in Fig. 1. In particular, we show that our method achieves high recall, *i.e.* it represents the majority of the scene structure in the map even when subject to noise, which is highly desirable for safe robot operation.

In summary, our contributions are as follows:

- We propose a novel method for online volumetric mapping of arbitrary-sized environments based on neural implicit representations and provide a training strategy to incrementally fuse geometric information from sequential point cloud scans directly in latent space.
- We show the flexibility of the presented method to operate on sparse inputs from arbitrary range sensors, generalizing to varying configurations without retraining, and operate in real-time even on just a CPU.
- We demonstrate in thorough experimental evaluations that our method can build maps with significantly better scene completeness in the presence of pose estimation errors.
- We make the code available as open-source for the benefit of the community[1].

## 2. Related Work

**Neural implicit representation:** Recently, many works have proposed methods to learn continuous implicit functions for shape representations. This line of methods has shown superior quality and flexibility compared to explicit representations such as voxels [50, 51, 19], points [14, 47], and meshes [17, 18, 24]. The pioneering implicit representation methods typically encode the geometric input data into a single latent code, which can then be used to condition the decoder to predict the occupancy or SDF values of query locations in continuous space [30, 37]. Due to the use of a single latent code, many of these methods are prone to losing low-level details and do not scale well to large scenes. This has been improved in subsequent works by using more sophisticated structures composed of multiple latent codes, such as planes [38, 25], grids [5, 38] and Gaussians [16, 15]. Several recent works also explore broader applications, such as novel view synthesis [43, 32, 29, 40, 6], representing articulated objects [13, 31], and learning implicit representation only from 2D observations [26, 27, 33].

---

[1] https://github.com/ethz-asl/neuralblox

While most of these works focus on the reconstruction of single objects, Convolutional Occupancy Networks [38] and Neural Distance Fields [7] demonstrate the ability to reconstruct large scenes. However, their methods work offline requiring complete point clouds of the scene beforehand. Concurrently, Jiang *et al.* [23] and Chabra *et al.* [3] proposed methods that represent a scene as coarse latent grid structures which can be decoded to implicit surfaces. However, their methods need to perform optimization for each grid during inference and thus are not suitable for online reconstruction. A common limitation of the aforementioned methods is that they require complete input data and do not consider online updates of the shape representations given sequential partial observations for large-scale mapping.

**Classical volumetric mapping:** The two common representations for volumetric mapping are occupancy and TSDF maps. In OctoMap [20], a 3D scene is segmented into hierarchical octree nodes, where each node stores an occupancy probability to build memory-efficient occupancy maps. While having a straightforward probabilistic update equation, OctoMap requires hand-tuning of the sensor uncertainties to achieve accurate results.

Similarly, inconvenient hand-tuning of parameters is also required for classical SDF-based approaches that are mostly built upon the TSDF fusion method proposed by Curless and Levoy [10], where uniform grids of SDF values are updated in a weighted averaging fashion. This approach is integrated into many online scene reconstruction frameworks, such as KinectFusion [22], BundleFusion [12], Voxgraph [42], Voxel Hashing [34], or InfiniTAM [39]. However, defining the optimal weighting function for the update rule is a non-trivial task. While Curless and Levoy [10] set it to be a constant of 1, KinectFusion [22] proposes a weighting function that is proportional to the cosine of angle between the ray from the sensor origin and the normal of the surface. In a sensor model-based approach, Voxblox [35] proposes another function using a quadratic weight based on the measured depth.

All of these approaches perform explicit updates without any shape context and therefore are oftentimes not robust to noisy inputs. In contrast, our method leverages shape priors from training data to more robustly deal with noise or partial observations.

**Learning-based volumetric mapping:** More recently, RoutedFusion [48] introduced a learning-based approach for online TSDF updates and weights, resulting in a significant improvement in handling sensor noise. Similar to our approach, NeuralFusion [49] performs depth fusion directly in a learned latent space and extracts an interpretable representation (TSDF) afterward. However, their fusion approach does not scale well to larger scenes or pose noise. Inspired by the recent progress in neural rendering [29],

iMAP [44] represents a scene using a small multi-layer perceptron that is trained in live operation to predict color and volume density, which can then be converted into an occupancy map. However, iMAP has limited scalability as the entire scene is represented in one single code. Since this code takes time to optimize to match the sensor measurements its usability for *e.g.* online planning is not clear. In a concurrent work, Huang *et al.* [21] recently presented DI-Fusion, which similarly to our proposed method uses a grid of neural representations that are temporally fused together. As opposed to previous methods that tightly include RGBD camera models or require dense scans for surface normal estimation, our method does not assume any particular distribution of input points and could potentially be extended to other types of range sensors.

## 3. Method

Fig. 2 presents an overview of our pipeline. Given a sequence of input point clouds $I^{1:t} \in \mathbb{R}^{3 \times N}$, the goal of our method is to produce an occupancy map $o : \mathbb{R}^3 \rightarrow$ {occupied, free}, which contains information about free and occupied space. This map can then be queried at any time step $t$ to enable online robotic tasks such as active path planning or obstacle avoidance.

Since robots oftentimes operate in an open world setting, we aim to generate an occupancy map of any arbitrarily sized scenes. To accurately model the scene, we partition 3D space into large voxels $\mathbb{V} = \{v\}$, similar to [23]. We denote the voxel side length as $d_V$, which is typically within $0.5\,\mathrm{m}$ to $1.0\,\mathrm{m}$. In contrast to [23], we dynamically allocate voxels which are inside the view frustum or contain input points to scale to unknown scene sizes. Each voxel contains a latent code $\mathbf{z}_v$, where the index $v$ is used to denote a specific voxel, to represent the geometry within the fixed volume of each voxel. To prevent discontinuities between voxels, we set an overlapping input dimension $d_I (> d_V)$ that defines the input volume of a particular voxel.

When a new pointcloud $I^t$ arrives, we convert all points that fall into the input volume of each voxel $I_v^t$ into the coordinates of that voxel. Subsequently, we encode $I_v^t$ into a latent code $\mathbf{z}_v^t$ and fuse this shape information into the voxel. This results in fast updates and a compact map representation, where each voxel considers the geometric context within its cell and aggregated shape information from all observations $I_v^{1:t}$.

Afterward, the fused latent codes $\hat{\mathbf{z}}_v^t$ can be decoded via a decoder network into occupancy probabilities for map queries. We set the query dimension $d_q (\geq d_V)$ to construct a query volume, in which occupancy probabilities of points will be queried given the latent code of the corresponding voxel. The larger query volume enables us to smoothen artifacts between neighbouring voxels by interpolating occupancy probabilities in the overlapping regions. Lastly, a

high-resolution output mesh can be generated from the occupancy grids using *e.g.* the marching cubes algorithm [28].

We separate the training pipeline into two steps. First, we train the encoder and decoder in a supervised way using the same subset as Choy *et al.* [8] of ShapeNet [4], a synthetic object-level dataset with the pre-processing steps from Occupancy Networks [30]. Second, using the trained encoder and decoder, we train a fusion network in a self-supervised manner using a real-world scene dataset [36] to perform latent code updates.

### 3.1. Local Geometry Representation

Our encoder and decoder mainly follow the architecture of 3D grid Convolutional Occupancy Networks [38], a neural implicit representation approach that is able to capture local spatial context in their implicit representation. Input point clouds are encoded using PointNet [41] and 3D U-Net [9]. In the decoding phase, occupancy probabilities of query points in continuous space, $\mathbf{p} \in \mathbb{R}^3$, are then predicted using a lightweight decoder network.

**Architecture adaptation:** In the original implementation of Convolutional Occupancy Networks [38], the most compressed representation processed from the input is in the deepest layer of 3D U-Net. However, due to the skip connections in 3D U-Net, this compressed representation cannot be decoded independently from the input. As we are interested in storing the most compressed representation in our map without storing its original input, we modify the 3D U-Net by removing its skip connections. The architecture and training details of our modified encoder and decoder are provided in the supplementary material.

**Encoder:** We denote the encoder with weights $\theta_e$ as $f_{\theta_e}$ and input point clouds as $I_v^t$, where no particular structure or density is imposed on $I_v^t$. The latent code $\mathbf{z}_v^t$ is obtained as:

$$\mathbf{z}_v^t = f_{\theta_e}(I_v^t) \tag{1}$$

**Decoder:** Given the latent code, the goal of the decoder is to predict the occupancy probability of any point $\mathbf{p} \in \mathbb{R}^3$ within the query volume of voxel $v$. The decoder $g_{\theta_d}$, with weights $\theta_d$ processes a latent code $\mathbf{z}_v^t$ and generates the occupancy probability at the query point $\mathbf{p}$:

$$P_{occ}(\mathbf{p}|\mathbf{z}_v^t) = g_{\theta_d}(\mathbf{p}, \mathbf{z}_v^t), \quad P_{occ} \in [0, 1] \tag{2}$$

### 3.2. Latent Code Fusion

The goal of the latent code fusion is to update the current map estimate given a new input, *i.e.*, the updated latent code will contain more complete geometrical information from sequential partial observations. Given the trained encoder $f_{\theta_e}$ and decoder $g_{\theta_d}$, we train the fusion network $h_{\theta_f}$
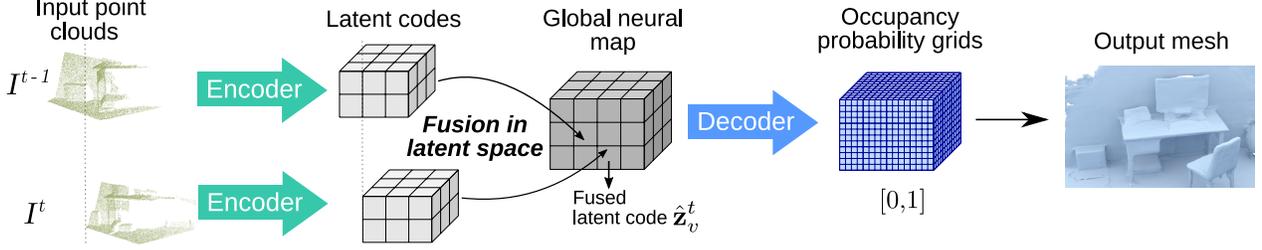
Figure 2. Overview of our method. We encode a stream of input points clouds into latent codes and fuse them in latent space to incrementally build a global neural map. This neural map can be decoded at anytime into an occupancy map.

in a self-supervised manner on a real-world scene dataset. During training, we set input dimension $d_I = 1.2\,\mathrm{m}$, voxel size $d_V = 1.0\,\mathrm{m}$ and query dimension $d_q = d_I$ and do not perform boundary interpolation. With the frozen encoder and decoder, we train the fusion network to predict a fused latent code at any time step $t$, denoted as $\hat{\mathbf{z}}^t$. An objective function is designed such that $\hat{\mathbf{z}}^t$ has similar properties to the latent code encoded from the accumulated input point clouds up to time step $t$, $\mathbf{z}^{*t}$. The latent code of voxel $v$ based on the accumulated input point clouds that have fallen into its input volume until time step $t$ is defined as:

$$\mathbf{z}_v^{*t} = f_{\theta_e}(I_v^{1:t}) \qquad (3)$$

Since we only consider measurements that affect a voxel $v$, we denote the set of relevant time steps as:

$$\mathbb{T}_v^t = \{\tau \mid |I_v^\tau| > 0\}_{\tau=0,\dots,t} \qquad (4)$$

To generate the fused latent code prediction, we keep track of previous latent codes in an averaging fashion. This enables a constant memory footprint for each voxel, which stores the summation of latent codes from input point clouds that fall into it:

$$\bar{\mathbf{z}}_v^t = \sum_{\tau \in \mathbb{T}_v^t} f_{\theta_e}(I_v^\tau) \qquad (5)$$

and the number of measurements $N_v = |\mathbb{T}_v^t|$. The fused latent code of voxel is generated by first dividing the summation by $N_v$, then feeding it into the fusion network $h_{\theta_f}$:

$$\hat{\mathbf{z}}_v^t = h_{\theta_f}\left(\frac{\bar{\mathbf{z}}_v^t}{N_v}\right) \qquad (6)$$

**Feature alignment loss:** We employ a feature alignment loss to minimize the distance between $\hat{\mathbf{z}}$ and $\mathbf{z}^*$ in latent space:

$$\mathcal{L}_{fea} = \frac{1}{|\mathbb{V}|} \sum_{v \in \mathbb{V}} \left\| \mathbf{z}_v^{*t} - \hat{\mathbf{z}}_v^t \right\|_1 \qquad (7)$$

This feature alignment loss allows the fusion network to correct for errors introduced by the averaging operation if the latent space is not perfectly linear and encourages the fused codes $\hat{\mathbf{z}}_v$ to stay in the domain of the pre-trained encoder and decoder.

**Reconstruction loss:** Additionally, we add a reconstruction loss that minimizes the logit differences of query points $\mathbf{p}$ decoded using $\hat{\mathbf{z}}$ and $\mathbf{z}^*$, where $\mathbb{P}$ are randomly distributed sample points:

$$\mathcal{L}_{rec} = \frac{1}{|\mathbb{V}||\mathbb{P}|} \sum_{\mathbf{p} \in \mathbb{P}} \sum_{v \in \mathbb{V}} \left\| g_{\theta_d}(\mathbf{p}, \mathbf{z}_v^{*t}) - g_{\theta_d}(\mathbf{p}, \hat{\mathbf{z}}_v^t) \right\|_1 \qquad (8)$$

**Backward pass:** During training, we always take 8 sequences of input point clouds and then perform the backward pass ($t = 8$). Therefore, the number of updates of a voxel ($N_v$) can range between 1 and 8. However, we found that this training procedure generalizes well to longer sequences. The total loss function is obtained as the summation of the feature alignment and reconstruction loss:

$$\mathcal{L} = \mathcal{L}_{fea} + \mathcal{L}_{rec} \qquad (9)$$

**Fusion network:** The fusion network $h_{\theta_f}$ is composed of two 3D convolution layers with ReLU activations. The architecture details are in supplementary materials. Due to its shallow structure, the fusion process does not add any noticeable runtime to our pipeline.

### 3.3. Occupancy Map Generation

A threshold $\tau_{occ}$ is set to define whether queried points are considered free or occupied space, where points with occupancy probability less than $\tau_{occ}$ are considered free. We uniformly sample query points in a structure of dense 3D grid for each voxel. If a voxel is not yet allocated, it is considered unknown. Otherwise, if a voxel is allocated but has not yet received any point cloud ($N_v = 0$), it directly outputs free space when queried. In our scene reconstruction experiments discussed in Sec. 4, we store the summation of latent codes and $N_v$ for every voxel and only generate the occupancy after the last input scan. In practice, the occupancy can be decoded at any time step by dividing the summed latent codes by $N_v$, predict the fused latent code,

and then feed query points as well as the predicted latent code into our decoder.

**Boundary interpolation:** We set $d_q = d_V + (d_I - d_V) \times 0.5$, giving overlapping regions between adjacent voxels. When a voxel is decoded, it looks up to its neighboring voxels and performs linear interpolations of probabilities on the overlapping regions. A comparison with and without the boundary interpolation is in the supplementary material.

## 4. Experiments

We provide the experimental details comparing the performance of our model and a standard TSDF method implemented in Voxblox [35]. We mainly evaluate the robustness of every model subject to inputs with varying levels of pose noise. In Sec. 4.3, we also provide thorough details of the runtime of our models in different configurations.

### 4.1. Model Comparison

Without any retraining, we set up two configurations with different voxel size $d_V$ and input dimension $d_I$ during occupancy map generation, as follows:

- Ours (0.5 m): $d_V = 0.5$ m, $d_I = 0.7$ m
- Ours (1 m): $d_V = 1$ m, $d_I = 1.2$ m

A smaller voxel size results in a finer reconstruction at the expense of storing more latent codes and thus slower inference speed. In both configurations, we employ query points structured in a grid with a resolution of $100^3$ points/m$^3$. By including the query points in overlapping regions, this results in $60^3$ and $110^3$ points/voxel for the 0.5 m and 1 m configurations, respectively.

We set the initial 3D grid resolution to 24 voxels, hidden layer size to 32, and 3D U-Net depth to 2. Thus, for each voxel, input point clouds are processed into a latent code $\mathbf{z}_v$ with a dimension of $6 \times 6 \times 6 \times 128$. By corresponding to the same number of bytes used per m$^3$ between our method and TSDF, we benchmark our model with $d_V = 0.5$ m to a TSDF with a voxel resolution of $0.02$ m and ours with $d_V = 1$ m to a TSDF with a voxel resolution of $0.04$ m. For fast runtime, our method can operate accurately using sparse input points, which is further discussed in Sec. 4.3. Thus in our experiments, we only use 10% of the points per scan as the input to our models. With this input percentage, we empirically set the threshold $\tau_{occ} = 0.05$. Additionally, we do not update a latent code if its corresponding input volume receives less than 1% of the input points. Conversely, we feed 100% of the point cloud for TSDF.

### 4.2. Metrics

We use the following metrics for our evaluation:
- **Accuracy** (lower better): mean absolute distance (MAD) of points sampled densely from the predicted mesh to the closest surface on the ground truth mesh.

- **Completeness** (lower better): MAD of points sampled densely from the ground truth mesh to the closest surface on the predicted mesh.
- **Recall** (higher better): the percentage of points sampled from the ground truth mesh that have closest absolute distances to the predicted mesh within a threshold $\tau_r$. We use $\tau_r = 0.05$ m.

### 4.3. Redwood Indoor RGBD Dataset

**Training:** We leverage the Redwood dataset [36] to train the fusion network. It contains $640 \times 480$ resolution RGBD sequences of five real-world indoor scenes captured with an Asus Xtion Live camera. We use four scenes, namely *bedroom, boardroom, lobby,* and *loft* as our training data. Nonetheless, we find that using only one scene (*bedroom*) already produces a model that generalizes well. During training, we take 8 depth sequences with alternating separations of 30 and 90 frames in each iteration. Then, we convert them into point clouds with a truncated depth of $3$ m. We sample 25,000 points each scan, which corresponds to about 10% of total points. Further training details are in the supplementary material.

**Evaluation:** We use the *apartment* scene for evaluation, which consists of approximately 31,910 RGBD frames. We utilize the high-quality mesh, reconstructed offline using the method presented in [36], as our ground truth. Evaluations are carried out on the scene thoroughly covered by the first 13,000 frames where we found the ground truth to be most reliable. We take every 30$^{\text{th}}$ frame of the RGBD scans and convert them to point clouds with a truncation distance of 3.0 m as input to the evaluated models. The ground truth mesh is slightly trimmed so that it covers roughly the same area as the input point clouds, as shown in Fig. 3.

To evaluate the robustness of our system given state estimation uncertainties, we add artificial pose errors to the provided extrinsic camera parameters. Since most robots are usually equipped with an Inertial Measurement Unit (IMU), the direction of gravity can typically be well estimated. We thus add noise only to the $x$ and $y$ directions (assuming $z$ is the direction of gravity) of the camera pose. Given the extrinsic camera parameters of a frame:

$$C_{ex} = \left( \begin{array}{cc} R & T \\ 0 & 1 \end{array} \right)$$

where $R \in SE(3)$ is the rotation matrix and $T \in \mathbb{R}^3$ is the translation vector, we sample two values from a normal distribution with zero mean and standard deviation $\sigma_T$ and add them to the translation components $T_x$ and $T_y$.

To better illustrate the influence of different levels of $\sigma_T$, we show the sampled accumulated inputs in Fig. 3. Quantitative results with varying values of $\sigma_T$ are presented in Tab. 1 and qualitative visualizations of the reconstructions
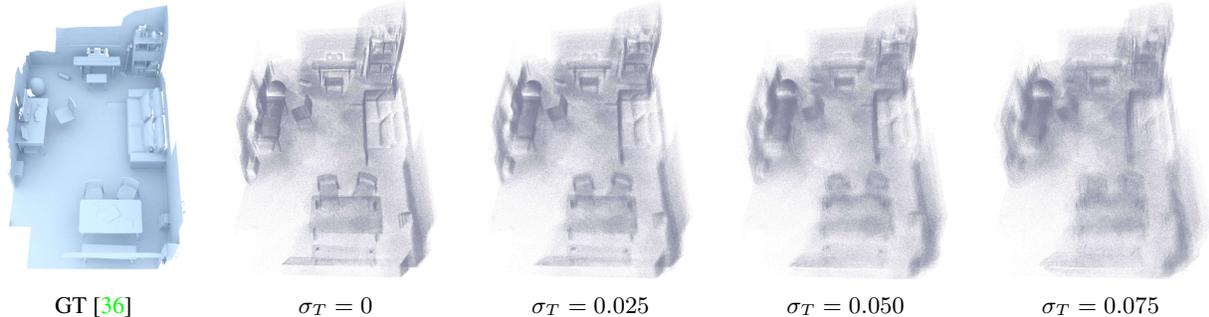
| GT [36] | $\sigma_T = 0$ | $\sigma_T = 0.025$ | $\sigma_T = 0.050$ | $\sigma_T = 0.075$ |

Figure 3. Ground truth and sampled accumulated inputs. Different levels of translation noise $\sigma_T$ (m) are applied on $x$ and $y$ axes.

are presented in Fig. 4. In the supplementary material, we also present quantitative and qualitative results of the cases where translation noise is applied to the $x, y$, and $z$ axes as well as when orientation errors around $z$ (yaw) are present.

| $\sigma_T$ | *Left: TSDF* (0.02 m) [35]. ***Right:*** Ours (0.5 m). | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| (m) | *Accuracy* (m) | | *Completeness* (m) | | *Recall* | | *Recall\** | |
| 0 | **0.0175** | 0.0225 | 0.0137 | **0.0135** | 0.988 | **0.990** | 0.984 | **0.986** |
| 0.025 | **0.0213** | 0.0264 | 0.0202 | **0.0159** | 0.926 | **0.989** | 0.888 | **0.985** |
| 0.050 | 0.0377 | **0.0372** | 0.0352 | **0.0216** | 0.756 | **0.932** | 0.625 | **0.898** |
| 0.075 | 0.0535 | **0.0503** | 0.0486 | **0.0249** | 0.635 | **0.876** | 0.438 | **0.809** |

| $\sigma_T$ | *Left: TSDF* (0.04 m) [35]. ***Right:*** Ours (1 m). | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| (m) | *Accuracy* (m) | | *Completeness* (m) | | *Recall* | | *Recall\** | |
| 0 | **0.0206** | 0.0292 | **0.0160** | 0.0187 | 0.964 | **0.970** | 0.950 | **0.957** |
| 0.025 | **0.0218** | 0.0309 | **0.0191** | 0.0201 | 0.936 | **0.964** | 0.906 | **0.947** |
| 0.050 | **0.0290** | 0.0367 | 0.0299 | **0.0245** | 0.817 | **0.923** | 0.722 | **0.885** |
| 0.075 | **0.0446** | 0.0456 | 0.0404 | **0.0286** | 0.705 | **0.849** | 0.545 | **0.768** |

\*The floor is excluded.

Table 1. Performance of our method on our evaluation set subject to translation noise $\sigma_T$ on the $x$ and $y$ axes. Notably, our method shows high robustness to state estimation noise and is always able to faithfully capture the underlying geometry, as shown by the high recall. This is particularly pronounced for objects of interest, *e.g.* when the floor is excluded.

From the quantitative and qualitative results, we can see that our models are significantly better at preserving the completeness of the reconstruction from noisy input. While accuracy and completeness are comparable for low levels of noise, our method shows significantly increased robustness in the presence of pose errors. Most importantly for use in robotic planning, our method is always able to faithfully capture the underlying geometry, shown in the high recall. This is particularly pronounced for objects of interest, *e.g.* when the floor is excluded. Note that we train our networks with a 1 m voxel configuration and generate the 0.5 m voxel resolution maps without retraining, thus demonstrating the flexibility of our method.

We also present the reconstructed meshes using a static fusion method, *i.e.* encoding and decoding is performed in a sliding window manner after all input point clouds are accumulated, as in [38]. The performance of the static approach using our encoder and decoder with $\sigma_T = 0.075$ m is shown in Tab. 2. Without our fusion strategy, all of the noise is encoded into the latent codes, resulting in inaccu-

rate meshes with overly thick structures, visible in Fig. 5. In contrast, our method captures the estimated statistics of input data and results in more accurate reconstructions.

| Voxel | *Left:* Static fusion [38]. ***Right:*** Ours. | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| size | *Accuracy* (m) | | *Completeness* (m) | | *Recall* | | *Recall\** | |
| 0.5 m | 0.0669 | **0.0503** | 0.0416 | **0.0249** | 0.690 | **0.876** | 0.539 | **0.809** |
| 1 m | 0.0789 | **0.0456** | 0.0573 | **0.0286** | 0.579 | **0.849** | 0.399 | **0.768** |

Table 2. Comparison of static fusion [38] using our encoder and decoder and our sequential method. $\sigma_T = 0.075$ m.

**Pointcloud sparsity and runtime:** Next, we evaluate the performance of our models given different levels of input point cloud sparsity by feeding sub-sampled point clouds per scan into our model. We also report the encoding speed of our PyTorch implementation evaluated on a machine with an NVIDIA GeForce GTX 1650 GPU (4GB) and an Intel i7-9750H@2600GHz CPU. Tab. 3 shows the results with different levels of input sparsity.

| Input % | *Left:* 0.5 m voxel. ***Right:*** 1.0 m voxel. | | | | | | |
|---|---|---|---|---|---|---|---|
| | *Encoding speed (FPS)* | | *Accuracy* (m) | | *Completeness* (m) | | *Recall* | |
| 5 | 5.2 | 11.7 | **0.0208** | 0.0281 | **0.0129** | 0.0181 | 0.988 | 0.969 |
| 10 | 4.7 | 11.4 | 0.0225 | 0.0292 | 0.0135 | 0.0187 | 0.990 | **0.970** |
| 50 | 2.9 | 6.3 | 0.0255 | 0.0326 | 0.0151 | 0.0206 | 0.992 | 0.965 |
| 100 | 1.8 | 3.1 | 0.0263 | 0.0335 | 0.0157 | 0.0213 | **0.993** | 0.961 |

Table 3. Performance with varying levels of input point cloud sparsity. Our system achieves similar reconstruction performance even for sparse inputs, enabling it to run in real-time and with low resolution sensors. The reconstruction quality appears slightly lower for high input percentages as we tune the threshold $\tau_{occ}$ for 10% of inputs. This can be improved by adjusting $\tau_{occ}$.

Despite being trained on 10% input data, due to the local pooling operations that aggregate the features from $N$ point clouds into a 3D grid with dimension $24 \times 24 \times 24$, our method shows strong robustness w.r.t sparse input data and obtains similar reconstruction results given sparser inputs. It further generalizes well to more input points with only a slight drop in performance. Since the density of occupancy prediction is not directly affected by input sparsity, our models suffer less from reduced surface density and smoothness that can appear in TSDFs given sparse input. This enables our method to achieve good reconstruction re-
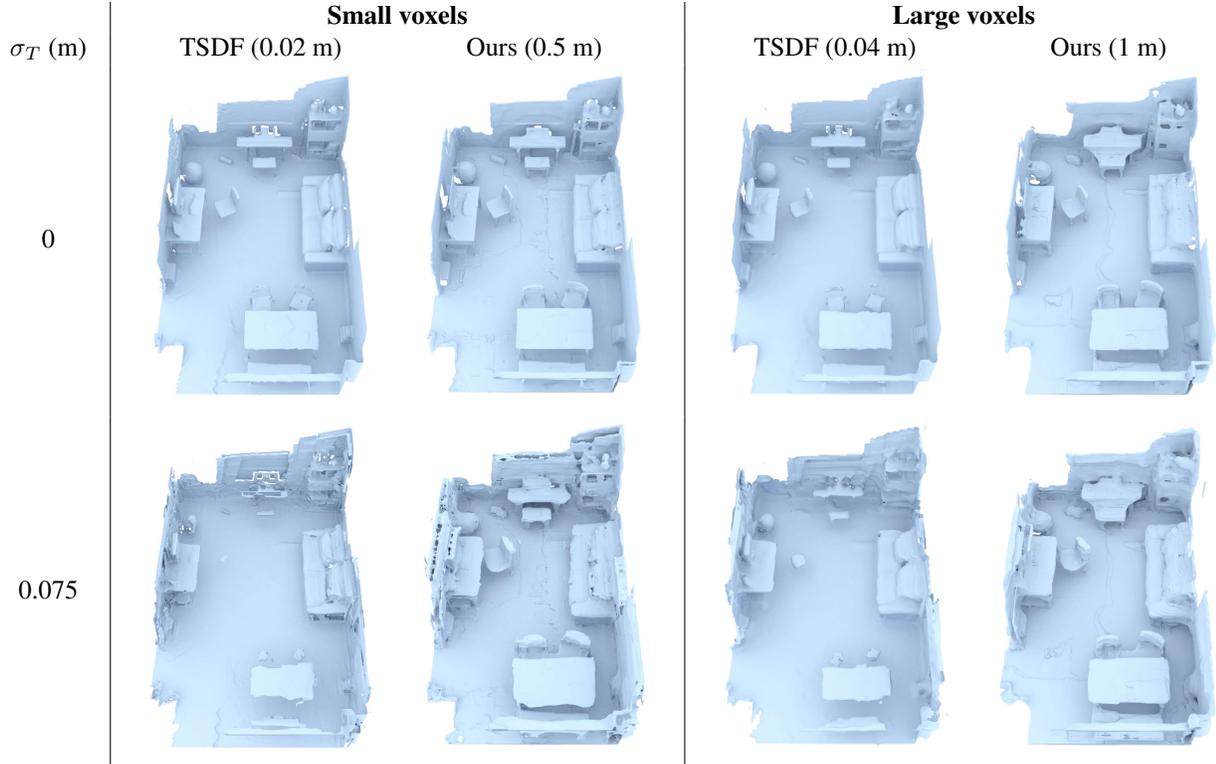
|  | **Small voxels** | | **Large voxels** | |
| $\sigma_T$ (m) | TSDF (0.02 m) | Ours (0.5 m) | TSDF (0.04 m) | Ours (1 m) |
| 0 | | | | |
| 0.075 | | | | |

Figure 4. Qualitative results on our evaluation set. Translation noise $\sigma_T$ is applied on the $x$ and $y$ axes.
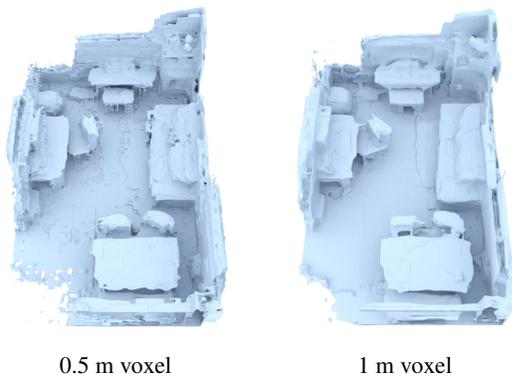


| 0.5 m voxel | 1 m voxel |

Figure 5. Reconstruction results using the sliding window method (as in [38]) applied to accumulated noisy inputs. Simply processing accumulated noisy inputs results in inaccurate meshes with overly thick structures.

sults in real-time and makes it amenable to cost-efficient low resolution sensors. Reconstruction comparisons with extremely sparse input are provided in the supplementary.

We report the speed of decoding latent codes in series, as shown in Tab. 4. In our implementation, we directly output free space when a voxel has not received any input point cloud. Thus, the decoding voxel / s can be faster than the reported speed in Tab. 4 depending on the layout of the scene. It can be accelerated further by processing voxels in a batch. The encoding and decoding speeds can give an overall oper-

ational frame rate of $2.4 - 2.6$ FPS when querying $25^3$ and $50^3$ points/voxel for $0.5\,\mathrm{m}$ and $1.0\,\mathrm{m}$ voxels, respectively[2].

| | | **Query points / voxel** | | | | |
| | | $25^3$ | $50^3$ | $60^3$ | $100^3$ | $110^3$ |
|---|---|---|---|---|---|---|
| **Decoding speed** | *Voxel / s* | 136.0 | 35.7 | 27.0 | 4.2 | 3.4 |
| | *FPS (0.5 m voxel)[3]* | 5.0 | 1.3 | 1.0 | - | - |
| | *FPS (1.0 m voxel)[3]* | 12.6 | 3.3 | 2.5 | 0.4 | 0.3 |

Table 4. Decoding speed given different query point densities.

**CPU-only configuration:** To allow our system to be employed on a compute-constrained mobile device, such as a Micro Aerial Vehicle (MAV), we demonstrate that it can run in real-time using only a CPU. We use the 1 m configuration, query dimension $d_q = d_V$, $25^3$ query points per voxel, and $5\%$ of input sparsity. Additionally, we do not update a latent code if its corresponding input volume receives less than $5\%$ of the input points. The reconstruction metrics and inference speed evaluated on an Intel i7-9750H@2600GHz CPU shown in Tab. 5 highlight that our method can be readily employed in mobile robotics.

| **Encoding** | **Decoding** | | **Accuracy** | **Completeness** | **Recall** |
|---|---|---|---|---|---|
| 4.7 FPS | 15.6 voxel/s | 2.1 FPS[3] | 0.0257 m | 0.0192 m | 0.953 |

Table 5. CPU-only configuration runtime and evaluation metrics.

---

[2]Decoding is performed every frame, thus the minimum overall FPS.
[3]Average speed when decoding all updated voxels every frame.

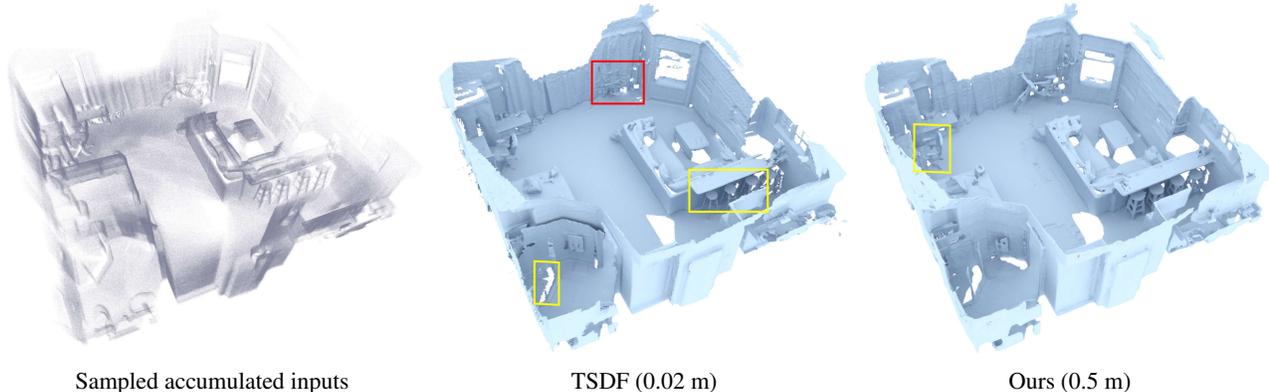|                           |                 |              |
|---------------------------|-----------------|--------------|
| Sampled accumulated inputs | TSDF (0.02 m)  | Ours (0.5 m) |

Figure 6. Qualitative results on ScanNet dataset. *Red:* Failed reconstruction. *Yellow:* Incomplete reconstruction.

## 4.4. ScanNet

To investigate the ability of our method to generalize to out-of-domain data, we evaluate it on ScanNet [11], an RGBD dataset of real-world indoor scenes captured by a Structure sensor mounted on an iPad. Since no accurate ground-truth reconstruction is provided, only qualitative comparisons are presented in Fig. 6. Similar to the in-domain evaluations, we observe high reconstruction quality and better object perseverance than TSDF, showing the capability of our model to generalize to scenes with different layouts and sensor characteristics. More results are provided in the supplementary materials.

## 5. Discussion

**Reconstruction:** From the experiments in Sec. 4, we can see a contrasting behavior between our method and TSDF. Given inputs subject to high state estimation errors, TSDF tends to produce a reconstruction thinner than the ground truth and fails to reconstruct objects with thin geometries. In contrast, our method captures the overall statistics of the input in a large spatial context and can better preserve object existence. High recall shows that our method is able to faithfully capture the underlying structure, albeit with occasional thickening. However, this is preferable to vanishing surfaces for safe navigation.

**Fusion strategy:** Through our experiments, we show that our fusion strategy (Eq. 6) can combine multiple latent codes of incomplete observations into a latent code of a more complete observation. Our conjecture to this ability is that the sequence of 3D convolutions in our encoder converts the input into a specific high-dimensional coordinate in latent space, as frequently is seen in autoencoder networks. The transformation $h_{\theta_f}$ in Eq. 6 then corrects the coordinate given from simply averaging latent codes into the desired coordinate of fused latent code. This ensures that our fused data can be both efficiently compressed via aver-

aging and remain in the support domain of the pre-trained shape representation networks.

**Limitations:** From Eq. 6, it is evident that the represented shape of a voxel can be dominated by frequent repeated measurements. To tackle this, the structure of our approach is amenable to many approaches such as key-framing or moving average integration, similar to [35], which gives more weight to recent observations. However, we leave the exploration of such methods for future work.

**Reproducibility:** Upon re-training our encoder, decoder, and fusion network from scratch, we observe that the threshold $\tau_{occ}$ that defines whether a query point is free has to be re-adjusted. Also, we find that a different total loss on convergence is observed during fusion training. This can result from different latent space landscapes being reached each time the encoder and decoder are trained from scratch.

## 6. Conclusion

We introduced a novel method for online volumetric mapping based on neural implicit representations that scales to arbitrary size environments and can incrementally update its geometric information. We formulate a method to integrate multiple neural implicit representations from sequential partial observations into a more complete representation. By fusing data in a formation of coarse voxels directly in latent space, our method leverages shape priors and captures local geometric context over temporally independent frames. We show in thorough experimental validation that this approach enables robust mapping in the presence of state estimation errors. Our system generalizes to a large variety of configurations and low-resolution range sensors and can operate in real-time even in CPU-only configurations. We make our system available as open source.

# References

[1] Matan Atzmon and Yaron Lipman. Sal: Sign agnostic learning of shapes from raw data. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020. 1

[2] Jonathan Bohren, Radu Bogdan Rusu, E Gil Jones, Eitan Marder-Eppstein, Caroline Pantofaru, Melonee Wise, Lorenz Mösenlechner, Wim Meeussen, and Stefan Holzer. Towards autonomous robotic butlers: Lessons learned with the pr2. In *Proc. IEEE International Conf. on Robotics and Automation (ICRA)*, 2011. 1

[3] Rohan Chabra, Jan Eric Lenssen, Eddy Ilg, Tanner Schmidt, Julian Straub, Steven Lovegrove, and Richard Newcombe. Deep local shapes: Learning local sdf priors for detailed 3d reconstruction. In *Proc. of the European Conf. on Computer Vision (ECCV)*, 2020. 2

[4] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015. 3

[5] Julian Chibane, Thiemo Alldieck, and Gerard Pons-Moll. Implicit functions in feature space for 3d shape reconstruction and completion. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020. 1, 2

[6] Julian Chibane, Aayush Bansal, Verica Lazova, and Gerard Pons-Moll. Stereo radiance fields (srf): Learning view synthesis for sparse views of novel scenes. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2

[7] Julian Chibane, Aymen Mir, and Gerard Pons-Moll. Neural unsigned distance fields for implicit function learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 2

[8] Christopher B Choy, Danfei Xu, JunYoung Gwak, Kevin Chen, and Silvio Savarese. 3d-r2n2: A unified approach for single and multi-view 3d object reconstruction. In *Proc. of the European Conf. on Computer Vision (ECCV)*, 2016. 3

[9] Özgün Çiçek, Ahmed Abdulkadir, Soeren S Lienkamp, Thomas Brox, and Olaf Ronneberger. 3D U-Net: learning dense volumetric segmentation from sparse annotation. In *Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, 2016. 3

[10] Brian Curless and Marc Levoy. A volumetric method for building complex models from range images. In *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, 1996. 1, 2

[11] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2017. 8

[12] 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)*, 2017. 2

[13] Boyang Deng, John P Lewis, Timothy Jeruzalski, Gerard Pons-Moll, Geoffrey Hinton, Mohammad Norouzi, and Andrea Tagliasacchi. NASA neural articulated shape approximation. In *Proc. of the European Conf. on Computer Vision (ECCV)*, 2020. 2

[14] Haoqiang Fan, Hao Su, and Leonidas J Guibas. A point set generation network for 3d object reconstruction from a single image. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2

[15] Kyle Genova, Forrester Cole, Avneesh Sud, Aaron Sarna, and Thomas Funkhouser. Local deep implicit functions for 3d shape. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2

[16] 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. 2

[17] Georgia Gkioxari, Jitendra Malik, and Justin Johnson. Mesh r-cnn. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2019. 2

[18] Thibault Groueix, Matthew Fisher, Vladimir G Kim, Bryan C Russell, and Mathieu Aubry. A papier-mâché approach to learning 3d surface generation. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2

[19] Christian Häne, Shubham Tulsiani, and Jitendra Malik. Hierarchical surface prediction for 3d object reconstruction. In *Proc. of the International Conf. on 3D Vision (3DV)*, 2017. 2

[20] Armin Hornung, Kai M Wurm, Maren Bennewitz, Cyrill Stachniss, and Wolfram Burgard. OctoMap: An efficient probabilistic 3D mapping framework based on octrees. *Autonomous robots*, 2013. 1, 2

[21] Jiahui Huang, Shi-Sheng Huang, Haoxuan Song, and Shi-Min Hu. DI-Fusion: Online implicit 3d reconstruction with deep priors. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021. 3

[22] Shahram Izadi, David Kim, Otmar Hilliges, David Molyneaux, Richard Newcombe, Pushmeet Kohli, Jamie Shotton, Steve Hodges, Dustin Freeman, Andrew Davison, et al. KinectFusion: real-time 3d reconstruction and interaction using a moving depth camera. In *Proc. of the 24th annual ACM symposium on User interface software and technology*, 2011. 2

[23] Chiyu Jiang, Avneesh Sud, Ameesh Makadia, Jingwei Huang, Matthias Nießner, and Thomas Funkhouser. Local implicit grid representations for 3d scenes. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2, 3

[24] Angjoo Kanazawa, Michael J Black, David W Jacobs, and Jitendra Malik. End-to-end recovery of human shape and pose. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2

[25] Stefan Lionar, Daniil Emtsev, Dusan Svilarkovic, and Songyou Peng. Dynamic plane convolutional occupancy networks. In *Proc. IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2021. 1, 2

[26] Shichen Liu, Shunsuke Saito, Weikai Chen, and Hao Li. Learning to infer implicit surfaces without 3d supervision. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. 2

[27] 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

[28] William E Lorensen and Harvey E Cline. Marching cubes: A high resolution 3d surface construction algorithm. *ACM SIGGRAPH Computer Graphics*, 1987. 3

[29] Ricardo Martin-Brualla, Noha Radwan, Mehdi SM Sajjadi, Jonathan T Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2

[30] 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. 2, 3

[31] Marko Mihajlovic, Yan Zhang, Michael J Black, and Siyu Tang. Leap: Learning articulated occupancy of people. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2

[32] 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. 2

[33] 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. 2

[34] Matthias Nießner, Michael Zollhöfer, Shahram Izadi, and Marc Stamminger. Real-time 3d reconstruction at scale using voxel hashing. *ACM Transactions on Graphics (ToG)*, 2013. 2

[35] Helen Oleynikova, Zachary Taylor, Marius Fehr, Roland Siegwart, and Juan Nieto. Voxblox: Incremental 3d euclidean signed distance fields for on-board MAV planning. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2017. 1, 2, 5, 6, 8

[36] Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Colored point cloud registration revisited. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2017. 3, 5, 6

[37] Jeong Joon Park, Peter Florence, Julian Straub, Richard 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

[38] 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, 3, 6, 7

[39] Victor Adrian Prisacariu, Olaf Kähler, Stuart Golodetz, Michael Sapienza, Tommaso Cavallari, Philip HS Torr, and David W Murray. Infinitam v3: A framework for large-scale 3d reconstruction with loop closure. *arXiv preprint arXiv:1708.00783*, 2017. 2

[40] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2

[41] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2017. 3

[42] Victor Reijgwart, Alexander Millane, Helen Oleynikova, Roland Siegwart, Cesar Cadena, and Juan Nieto. Voxgraph: Globally consistent, volumetric mapping using signed distance function submaps. *IEEE Robotics and Automation Letters*, 2019. 2

[43] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. 2

[44] Edgar Sucar, Shikun Liu, Joseph Ortiz, and Andrew J Davison. iMAP: Implicit mapping and positioning in real-time. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2021. 3

[45] Sebastian Thrun et al. Robotic mapping: A survey. 2002. 1

[46] Emmanouil G Tsardoulias, A Iliakopoulou, Andreas Kargakos, and Loukas Petrou. A review of global path planning methods for occupancy grid maps regardless of obstacle density. *Journal of Intelligent & Robotic Systems*, 2016. 1

[47] Xiaogang Wang, Marcelo H Ang Jr, and Gim Hee Lee. Point cloud completion by learning shape priors. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2020. 2

[48] Silvan Weder, Johannes Schonberger, Marc Pollefeys, and Martin R Oswald. Routedfusion: Learning real-time depth map fusion. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2

[49] Silvan Weder, Johannes L Schonberger, Marc Pollefeys, and Martin R Oswald. NeuralFusion: Online depth fusion in latent space. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2

[50] Jiajun Wu, Chengkai Zhang, Tianfan Xue, Bill Freeman, and Josh Tenenbaum. Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2016. 2

[51] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shapes. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2015. 2

[52] Qiangeng Xu, Weiyue Wang, Duygu Ceylan, Radomir Mech, and Ulrich Neumann. DISN: Deep implicit surface network for high-quality single-view 3d reconstruc-

tion. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. 1