

Neural Graph Map: Dense Mapping with Efficient Loop Closure Integration

Leonard Bruns¹ Jun Zhang² Patric Jensfelt¹
¹KTH Royal Institute of Technology, Stockholm, Sweden
²TU Graz, Graz, Austria
 {leonardb, patric}@kth.se, jun.zhang@tugraz.at

Abstract

Neural field-based SLAM methods typically employ a single, monolithic field as their scene representation. This prevents efficient incorporation of loop closure constraints and limits scalability. To address these shortcomings, we propose a novel RGB-D neural mapping framework in which the scene is represented by a collection of lightweight neural fields which are dynamically anchored to the pose graph of a sparse visual SLAM system. Our approach shows the ability to integrate large-scale loop closures, while requiring only minimal reintegration. Furthermore, we verify the scalability of our approach by demonstrating successful building-scale mapping taking multiple loop closures into account during the optimization, and show that our method outperforms existing state-of-the-art approaches on large scenes in terms of quality and runtime. Our code is available open-source at https://github.com/KTH-RPL/neural_graph_mapping.

1. Introduction

Simultaneous localization and mapping (SLAM) using cameras, often referred to as visual SLAM, has been a long standing problem in computer vision [7, 9, 13]. In particular, dense visual SLAM aims to construct a detailed geometric representation of the environment enabling various applications, such as, occlusion handling in augmented reality [26], planning and collision checking in robotics [1], or camera localization [21]. Often volumetric scene representations are employed as they are well-suited for online data integration. Traditional volumetric representations use grid-based structures [25], often accelerated through the use of octrees [10, 18] or voxel hashing [27]. However, incorporating loop closure constraints in volumetric maps is difficult [6, 30, 40, 43] compared to sparse keypoint-based maps, which can easily be deformed.

Recently, neural fields have emerged as a promising vol-

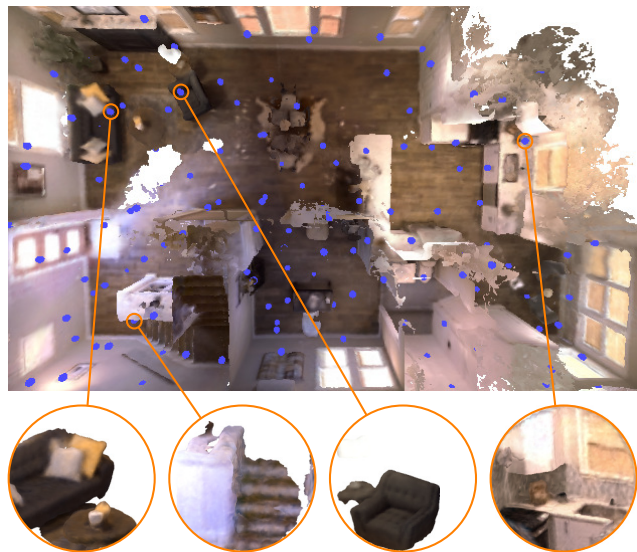


Figure 1. We propose to represent a scene by a set of neural fields (centers indicated by blue spheres) anchored to keyframes in a pose graph with each field capturing the scene within a ball surrounding it. This allows to dynamically extend the scene while also incorporating loop closure deformations into the volumetric scene representation without requiring full reintegration.

umetric scene representation due to their amenability to differentiable rendering-based optimization [37]. Building on the first neural field-based SLAM method iMAP [35], various adaptations of this method have been proposed. These adaptations primarily focused on enhancing optimization speed by altering network architecture, sampling scheme, and rendering formulation. Nevertheless, the majority of existing neural field-based SLAM methods [12, 35, 41, 45, 46] remain constrained by their monolithic data structure (i.e., a single fixed-size architecture). In that regard, neural fields share the same limitation as other volumetric scene representations: after data has been integrated, the volumetric scene cannot be easily deformed to take loop closure constraints into account.

To address this limitation, submap-based approaches

have been proposed that maintain a small number of neural field-based submaps [15, 17, 38]. However, submap-based methods come with their own set of difficulties such as inaccurate submap registration, transition artifacts at submap boundaries, and higher computational cost.

Instead of few, larger submaps, we propose to represent the scene by an extendable set of smaller, lightweight neural fields (see Fig. 1). These fields are dynamically anchored to keyframes in a pose graph, collectively forming the volumetric map while maintaining the ability to deform as the keyframe poses change upon loop closure. Our design allows to combine the benefits of sparse pose graph-based SLAM methods while maintaining a volumetric scene representation that remains consistent with the pose graph without requiring costly reintegration of previous sensor data. Furthermore, our method eliminates the necessity for fixed scene boundaries common among existing neural field-based SLAM methods. It achieves this by allocating additional fields dynamically as new parts of the scene come into view.

To summarize, our contributions are

- a novel RGB-D neural mapping framework that combines the robust, accurate tracking and efficient loop closure handling of sparse visual SLAM with the differentiable-rendering-based dense mapping of neural scene representations,
- a multi-field scene representation in which neural fields are anchored relative to a pose graph allowing for large-scale map deformations while limiting necessary reintegration, and
- a thorough comparison to multiple state-of-the-art methods on scenes of varying scales, including a novel set of sequences for the larger Replica scenes that allow to benchmark robustness and scalability.

2. Related Work

Traditional Volumetric Loop Closure The difficulty of incorporating loop closure constraints in volumetric scene representations is not limited to neural representations. Traditional volumetric representations such as sparse grids and octrees also require solutions to efficiently adapt to pose graph changes. Notable examples include: Kintinuous [43] in which the volumetric map is only used for local fusion and globally a mesh is deformed based on a deformation graph; BundleFusion [6] in which frames whose pose has changed are removed and reintegrated; VoxGraph [30] in which submaps are aligned and resulting constraints are included in the pose graph optimization; and recently LivePose [33] which extends the idea of BundleFusion to an RGB setting. Our work aims to enable efficient loop closure integration specifically for neural field-based representations for which de- and reintegration is difficult.

Neural Field-Based Scene Representations Neural fields are parametrized differentiable functions that map coordinates to quantities at that point in space. Initially, neural fields were trained using 3D supervision to represent shapes as occupancy [4, 19] or signed distance fields (SDFs) [28]. Shortly after, neural radiance fields (NeRFs) were introduced optimized using only 2D images through differentiable volume rendering [20]. Replacing NeRF’s density with an SDF-based representation higher quality surface reconstructions can be achieved [2, 42].

Few works have investigated the use of multiple neural fields. In Block-NeRF [36] multiple fields are combined to represent whole neighborhoods. However, the field positions are predefined. In Nerflets [44], fields of varying size are composed to create an editable scene representation. NeRFuser [8] proposes methods for registering and blending multiple NeRFs. Finally, vMAP [14] uses multiple fields to represent individual objects highlighting the potential of vectorizing neural field evaluations. Our work follows similar motivations: by splitting the scene into multiple independent fields, it becomes both, efficiently deformable upon loop closures, and dynamically extendable as the camera moves.

Neural Field-Based SLAM iMAP [35] was the first SLAM method using a multilayer perceptron (MLP) as the underlying scene representation. During optimization previous keyframes are continuously reintegrated to avoid forgetting. Subsequently, various modifications of this approach have been suggested aiming to alleviate this forgetting issue and improve run-time.

NICE-SLAM [46] addresses the forgetting issue by employing a hierarchical feature grid combined with a fixed, pretrained decoder MLP. By adopting an SDF-based field combined with efficient learnable positional encodings, ESLAM [12] and Co-SLAM [41] demonstrate significantly improved optimization times. Specifically, ESLAM uses a tri-plane [3] encoding, whereas Co-SLAM uses the aforementioned multi-resolution voxel hash encoding [22]. Co-SLAM further complements the hash encoding with one-blob encodings [23] to improve scene completion. PointSLAM [32] replaces NICE-SLAM’s regular feature grid with an irregular neural point cloud allowing to allocate more capacity in complex areas. However, none of the aforementioned approaches support integration of loop closure constraints, preventing their use for larger scenes, in which accumulation of drift is unavoidable.

[17] follows a similar idea to ours, also using multiple fields posed relative to a pose graph. However, they employ space warping in which each field in principle covers the entire scene. Synthesized views from multiple fields are merged through alpha compositing. This approach is well-suited for novel view synthesis; however, it generalizes poorly to mesh extraction and other geometric queries. In-

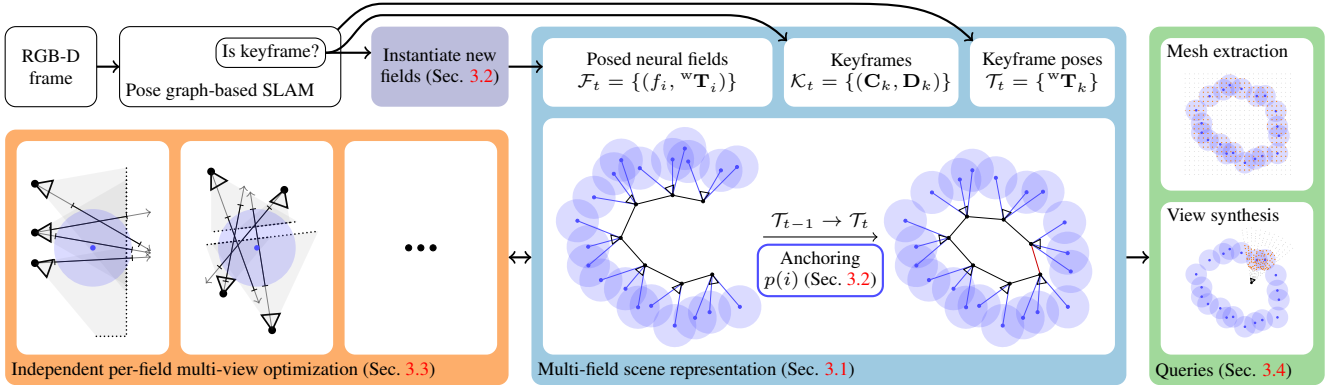


Figure 2. Neural graph map overview. Our framework can be used with any pose graph-based SLAM system and maintains a set of neural fields anchored relative to the keyframes of the pose graph. Each field captures the scene in a sphere surrounding it in local coordinates. This set of fields represents a deformable, volumetric scene representation that can easily adapt to and stay in sync with the pose graph. The fields are optimized independently always using the latest pose graph information for supervision. This design reduces transition artifacts, enables efficient optimization, and allows to average the output of overlapping fields in 3D leading to well-defined queries.

stead, in our work, each field covers a Euclidean ball around it, which allows general geometric queries. In MIPS-Fusion [38] global bundle adjustment at the level of fields is performed. Loopy-SLAM [15] also describes a submap-based system built on top of Point-SLAM [32]. Instead of larger keyframe-centric submaps, our approach represents the scene by many small fields, each covering a sphere surrounding it. By adjusting the field poses relative to the pose graph, while avoiding hard assignments of keyframes to fields during the optimization, our method allows to adapt to local and global pose graph deformations while avoiding transition artifacts common for submap-based methods.

Instead of using submaps, GO-SLAM [45] adapts to loop closures by biasing the optimization to keyframes whose pose changed the most. However, it still uses a single neural field, which requires reoptimization upon loop closure. Instead, we achieve instant map deformation upon loop closure by moving fields relative to the pose graph.

3. Method

We consider the problem of online RGB-D mapping without ground-truth poses. That is, given a stream of RGB-D frames (C_t, D_t) composed of color images $C_t \in \mathbb{R}^{H \times W \times 3}$ and depth maps $D_t \in \mathbb{R}^{H \times W}$, our goal is to build a dense scene representation at each time step t .

Fig. 2 gives an overview of the proposed framework. Similar to other recent neural field-based mapping approaches [14, 17, 45], our approach uses an off-the-shelf keyframe-based SLAM system to provide a set of posed keyframes as well as the pose of the current frame.

3.1. Multi-Field Scene Representation

Let $\mathcal{K}_t = \{(C_k, D_k) \mid k = 1, \dots, K_t\}$ and $\mathcal{T}_t = \{{}^w\mathbf{T}_k \in \text{SE}(3) \mid k = 1, \dots, K_t\}$ denote the set of keyframes and keyframe poses at time step t , respectively.

We propose to represent the scene by an extendable set of posed neural fields

$$\mathcal{F}_t = \{(f_i, {}^w\mathbf{T}_i) \mid i = 1, \dots, F_t\}, \quad (1)$$

where each field

$$f_i : \mathbb{R}^3 \rightarrow \mathbb{R}^3 \times \mathbb{R} \\ \mathbf{x} \mapsto (\mathbf{c}, s) \quad (2)$$

maps a 3D point \mathbf{x} in its local reference frame to a color \mathbf{c} and truncated signed distance s at that point. Each field's reference frame is defined by its pose ${}^w\mathbf{T}_i \in \text{SE}(3)$, which can be decomposed into a position ${}^w\mathbf{t}_i \in \mathbb{R}^3$ and orientation ${}^w\mathbf{R}_i \in \text{SO}(3)$. Each field only captures the scene within a sphere of fixed radius r .

We further define a function $p(i)$ that anchors a field i to a parent keyframe k . At each iteration, given a field i and its parent keyframe $k_i = p(i)$, we compute the field's pose according to ${}^w\mathbf{T}_i = {}^w\mathbf{T}_{k_i} {}^{k_i}\tilde{\mathbf{T}}_w \tilde{{}^w\mathbf{T}}_i$, where $\tilde{\mathbf{T}}$ denotes transforms from the previous time step. That is, fields move as if they are rigidly connected to their parent keyframe. The parent keyframe can however change over time depending on the anchoring strategy defined by $p(i)$.

In the following sections, the proposed strategies to instantiate and dynamically anchor fields to the pose graph (Sec. 3.2), optimize the fields (Sec. 3.3), and query (e.g., for mesh extraction and view synthesis) the full scene representation (Sec. 3.4) are described. Finally, further architecture details are provided (Sec. 3.5).

3.2. Field Instantiation and Anchoring

Whenever a new keyframe k is added to the pose graph, new fields are instantiated such that all observed 3D points ${}^w\mathcal{X}_k$ in the keyframe are covered by at least one field of radius r . An approximate two-stage algorithm is used (see

Fig. 3). First, the uncovered 3D points $\mathcal{X}_{\text{unc}} = \{\mathbf{x} \in {}^w\mathcal{X}_k \mid \|\mathbf{x} - {}^w\mathbf{t}_i\| > r \forall i \in \mathcal{N}_f\}$ are found. Second, the space is divided into voxel cells with a side length of $g = 2r/\sqrt{3}$, such that a cell is fully covered by a field of radius r when the center of the field is at the center of the cell. New fields are instantiated in the center of all cells that contain a point from \mathcal{X}_{unc} and no field center. This scheme ensures that a minimum distance of $g/2$ to existing fields is maintained; however, points in \mathcal{X}_{unc} might remain uncovered when fields move away from the cell centers through pose graph deformations. To alleviate the chance of uncovered points over time, the voxel grid is randomly shifted for every added keyframe.

New fields added this way are initially anchored to the keyframe that triggered the creation. Over time fields are supervised by all keyframes that observe it (see Sec. 3.3), some of which might be significantly closer than the initial keyframe. Therefore, we define the parenting function $p(i)$ as the closest keyframe (Euclidean distance) with valid depth observation used to supervise the field.

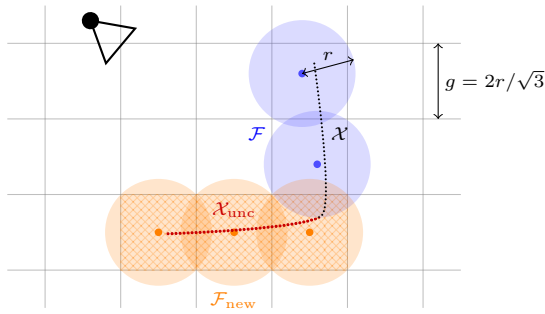


Figure 3. Grid-based instantiation scheme. A new keyframe (top-left) observes world points \mathcal{X} (black & red dots). The points \mathcal{X}_{unc} (red dots) that are not covered by any of the existing fields \mathcal{F} (blue circles) are determined. Cells that contain uncovered points and no existing field center (hatched orange) are used to instantiate new fields \mathcal{F}_{new} (orange circles). The new fields are positioned in the center of the uncovered cells and anchored to the keyframe.

3.3. TSDF-Based Optimization

At each time step t , a fixed number of optimization steps N_{it} are performed on up to N_f fields. During the optimization the fields f_i are trained independently, allowing for efficient parallelization [14]. Only the fields’ parameters are used as optimization variables. The fields’ poses are updated solely according to the pose graph as described in Sec. 3.1.

Due to its well-defined isosurface for mesh extraction, we model the geometry as a truncated signed distance field (TSDF) adopting the TSDF-based rendering model from [2] and modify it to be occlusion-aware. Specifically, instead of directly converting signed distances to sample weights as done in [2, 41], we first convert them to occupancy prob-

abilities. That is, given N_s samples along a ray $\mathbf{x}_i = \mathbf{o} + l_i\mathbf{d}, i = 1, \dots, N_s$ with origin \mathbf{o} and direction \mathbf{d} , the colors and signed distances $(\mathbf{c}_i, s_i) = f(\mathbf{x}_i)$ along the ray are computed. The signed distances are converted to occupancy probabilities according to

$$o_i = 4\sigma\left(\frac{\eta s_i}{\tau}\right)\sigma\left(-\frac{\eta s_i}{\tau}\right), \quad (3)$$

where σ denotes the sigmoid function, τ is the truncation distance, and η is a parameter that determines how sharply occupancy probability decays around the surface. Note that this definition ensures $o_i = 1.0$ for $s_i = 0.0$. The ray’s rendered color and depth are then computed via the weight $w_i = o_i \prod_{j=1}^{i-1} (1 - o_j)$ as

$$\mathbf{c} = \sum_{i=1}^{N_s} w_i \mathbf{c}_i \quad d = \sum_{i=1}^{N_s} w_i l_i \quad (4)$$

assuming the ray’s direction \mathbf{d} was scaled appropriately.

Sampling Strategy Each optimization iteration follows a three-stage sampling procedure. Specifically, we sample: (1) N_f fields that will be optimized in the next iteration; (2) N_r ray segments (with associated observed color \mathbf{c} and distance d) for each sampled field; (3) N_s query points along each sampled ray segment. Importantly, no assignments of keyframes to fields are made. Instead, which keyframes observe which field are recomputed in each iteration based on the latest pose graph information. This avoids transition effects stemming from different optimization targets of neighboring fields. The sampling strategy is visualized in Fig. 4 and details are provided in the supplementary material.

Loss Our loss is a weighted sum of four terms: color, depth, TSDF, and free-space. The weighing has been found experimentally to achieve a good trade-off between appearance and geometry quality.

The color loss is computed as the mean L1-norm of the difference between observed and estimated color

$$l_{\text{color}} = \frac{1}{|\mathcal{C}|} \sum_{(\mathbf{c}, \tilde{\mathbf{c}}) \in \mathcal{C}} \|\mathbf{c} - \tilde{\mathbf{c}}\|_1 \quad (5)$$

and the depth loss is computed as the mean Huber loss [11] between observed and estimated depth with $\delta = 5$ cm

$$l_{\text{depth}} = \frac{1}{|\mathcal{D}|} \sum_{(d, \tilde{d}) \in \mathcal{D}} \text{Huber}_\delta(d, \tilde{d}), \quad (6)$$

where \mathcal{C} and \mathcal{D} are sets containing tuples of the observed and estimated color and depth, respectively.

The query points \mathbf{x}_i are filtered into two sets $\mathcal{X}_{\text{tsdf}}$ and \mathcal{X}_{fs} based on whether a query point’s projected depth l_i is within the truncation distance τ of the corresponding observed depth $l_{\text{obs},i}$ or is more than τ in front of it,

respectively. A query point associated estimated signed distance \tilde{s}_i and its signed distance to the observed depth $s_i = l_{\text{obs},i} - l_i$ are then used to compute the TSDF loss as

$$l_{\text{tsdf}} = \frac{1}{|\mathcal{X}_{\text{tsdf}}|} \sum_{\mathbf{x}_i \in \mathcal{X}_{\text{tsdf}}} (\tilde{s}_i - s_i)^2 \quad (7)$$

and the free-space loss as

$$l_{\text{fs}} = \frac{1}{|\mathcal{X}_{\text{fs}}|} \sum_{\mathbf{x}_i \in \mathcal{X}_{\text{fs}}} (\tilde{s}_i - \tau)^2. \quad (8)$$

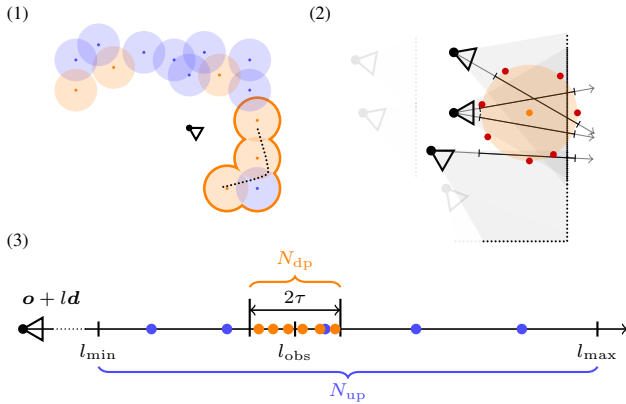


Figure 4. Multi-view sampling procedure. (1) A subset of fields (orange circles) to optimize in the next iteration is sampled biased to the currently observed fields (orange outline). (2) For each field, the observing keyframes (opaque frustums) are approximated based on whether samples on the field boundary (red dots) fall into the observed keyframe region. Ray segments to supervise the field are then sampled from all observing keyframes. (3) Each sampled ray segment $[l_{\min}, l_{\max}]$ is approximated by N_{up} uniformly sampled query points that cover the whole field and N_{dp} depth-guided query points distributed within the truncation distance τ around the observed distance l_{obs} .

3.4. Novel View Synthesis and Mesh Extraction

While the independent optimization ensures that overlapping regions will be supervised with the same data, fields might be at different stages of optimization and unobserved regions might converge to different solutions due to epistemic uncertainty. To reduce the resulting transition artifacts at the boundary of fields, it is possible to query the k nearest fields and average the outputs (see supplementary material for details). For query points outside of all fields’ spheres empty space is assumed, that is, $s = 1$. For novel view synthesis, since no depth observation is available at inference time, a fixed far distance l_{far} is used. Query points are uniformly distributed within the ray interval $[0, l_{\text{far}}]$ at the same density as the depth-guided points, that is, $N_{\text{p,inf}} = l_{\text{far}} N_{\text{dp}} / (2\tau)$ samples are used. To extract a mesh the marching cubes algorithm [16] with an unbiased isosurface level of 0 is used.

3.5. Architecture Details

ORB-SLAM2 is used as the keyframe-based SLAM system [24]. However, our proposed method is agnostic to the exact SLAM framework as long as a pose graph as described in Sec. 3.1 is available.

Each field is parameterized by a permutohedral hash encoding [31] followed by a linear layer with ReLU activation and another linear layer mapping to the color and distance. We also experimented with other positional encodings, such as Fourier encoding [37], triplane encoding [3, 12], and voxel hash encoding [22]; however, the best results given the same time budget were achieved with the permutohedral hash encoding closely followed by the voxel hash encoding.

To efficiently evaluate and optimize multiple neural fields in parallel, optimization and evaluation is vectorized following [14]. One limitation imposed by PyTorch’s [29] vectorization framework is that the batch size has to be the same for all evaluated networks. The design choices described in Sec. 3.3 are tailored to this limitation sampling the same number of rays per supervised field.

4. Experiments

4.1. Experimental Setup

Datasets We evaluate our method quantitatively on three synthetic datasets of varying difficulty: the *Replica* [34] sequences from iMAP [35], the *NRGBD* dataset by [2], and three novel sequences on the larger *Replica-Big*. The latter contains challenging trajectories with larger loops spanning multiple rooms. We also include results on two real-world datasets: ScanNet [5] and Kintinuous [43]. The latter contains a large-scale loop and the tested ScanNet sequences typically include one or two smaller loops. Quantitative evaluation of the map is not possible for the real-world datasets, because no accurate ground-truth is available in these datasets: ScanNet provides meshes from BundleFusion [6], Kintinuous provides none.

Metrics To assess the reconstruction quality of the methods we report the accuracy ratio, completion ratio, and F1-score [39] at a 5 cm threshold. In the supplementary material we further report the non-robust metrics, accuracy and completion. We do not evaluate tracking metrics, since we use ORB-SLAM2 as the tracking backend and our focus is on dense mapping, while efficiently supporting map deformations due to loop closures.

Baselines We compare our method to five state-of-the-art neural SLAM methods: NICE-SLAM [46] and Co-SLAM [41], which do not support loop closure; GO-SLAM [45], MIPS-Fusion [38], and Loopy-SLAM [15], which support loop closure. Of these, GO-SLAM uses a monolithic scene representation; MIPS-Fusion few, large submaps; and

Table 1. Comparison of mesh reconstruction quality (best ●, second best ●, third best ●).

		Replica	NRGBD	Replica-Big		
		Avg.*	Avg.*	apt0	apt1	apt2
NICE-S. [46]	Acc. R. (%)	93.17 ●	78.10	30.67	41.48	28.20
	Comp. R. (%)	90.11	70.58	70.58	42.21	26.80
	F1-Score (%)	91.60	74.07	74.07	41.84	27.49
Co-S. [41]	Acc. R. (%)	87.02	87.23	18.85	62.11	34.71
	Comp. R. (%)	86.39	82.86 ●	21.93	66.37	36.75
	F1-Score (%)	86.70	84.94 ●	20.27	64.17	35.70
GO-S. [45]	Acc. R. (%)	89.31	84.10	57.93 ●	74.38	80.22 ●
	Comp. R. (%)	74.20	66.19	59.21 ●	73.90 ●	81.06 ●
	F1-Score (%)	81.03	73.98	58.56 ●	74.14 ●	80.64 ●
MIPS-F. [38]	Acc. R. (%)	92.63	88.29 ●	28.00	75.97 ●	33.16
	Comp. R. (%)	91.38 ●	75.85	28.11	72.14	34.76
	F1-Score (%)	92.00 ●	81.40	28.06	74.01	33.94
Loopy-S.† [15]	Acc. R. (%)	99.99	80.52	×	31.34	59.87
	Comp. R. (%)	89.28	72.66	×	21.71	79.04
	F1-Score (%)	94.32	76.29	×	25.65	68.14
Ours-SF	Acc. R. (%)	93.11 ●	93.11 ●	54.69 ●	93.14 ●	76.86 ●
	Comp. R. (%)	91.28 ●	85.45 ●	56.07 ●	95.42 ●	79.57 ●
	F1-Score (%)	92.19 ●	89.02 ●	55.37 ●	94.27 ●	78.19 ●
Ours	Acc. R. (%)	92.76 ●	93.09 ●	56.25 ●	94.14 ●	79.28 ●
	Comp. R. (%)	90.87 ●	84.70 ●	56.89 ●	95.25 ●	81.81 ●
	F1-Score (%)	91.80 ●	88.58 ●	56.57 ●	94.69 ●	80.52 ●

* Full results are available in the supplementary material.

† Excluded from ranking because it uses ground-truth depth during evaluation.

Loopy-SLAM uses neural point-based submaps. In addition we compare to a monolithic, single-field ablation of our method (Ours-SF), which is discussed in Sec. 4.4.

Loopy-SLAM is the only non-volumetric method representing only the volume within centimeters of the surfaces. When querying (for rendering, which in their case is used for subsequent mesh extraction), Loopy-SLAM relies on ground-truth depth maps. Because none of the other methods do so, it leads to an unfair advantage especially on the noise-free synthetic scenes. Therefore, we exclude Loopy-SLAM from the ranking in the tables.

Implementation Details For all experiments we use the same parameters, most notably, $r = 1$ m, $N_{it} = 5$, $N_f = 32$, $N_r = 512$, $N_{up} = 8$, $N_{dp} = 16$ and $k = 2$. We adjust the truncation distance τ from 0.1 m for synthetic data to 0.2 m for real data to account for the increased depth noise. For ORB-SLAM2 we use the default parameters in the RGB-D setup increasing the number of extracted features in challenging scenes with low texture. To facilitate reproducibility of our results, the full evolution of the pose graph was precomputed and played back as if ORB-SLAM is running in parallel. We will publish these recordings in conjunction with our code. All remaining parameters are provided in the supplementary material and as part of the released code.

4.2. Reconstruction Quality Evaluation

Tab. 1 reports quantitative results on all datasets. Our method achieves state-of-the-art results on the smaller scenes (Replica, NRGBD) and outperforms existing methods on larger scenes (Replica-Big). The qualitative results on Replica-Big shown in Fig. 5 further illustrate the differences. NICE-SLAM and Co-SLAM drift off and cannot correct for the drift upon returning to previous frames leading to poor results. GO-SLAM also benefiting from loop closures successfully maps large parts; however, fails in the last section of apt0 and generally exhibits worse reconstruction quality. While GO-SLAM achieves quantitatively slightly better results as our method on apt0, we note that qualitatively our method exhibits better scene completion and higher fidelity overall¹. The submap-based methods MIPS-Fusion and Loopy-SLAM fail to achieve globally consistent result on most scenes. Particularly Loopy-SLAM failed on all scenes of Replica-Big². Our method and GO-SLAM benefit from their underlying SLAM systems ORB-SLAM2 and DROID-SLAM, respectively.

Fig. 5 shows results on three scenes of the ScanNet dataset (the “ground-truth” was obtained using BundleFusion in this case [6]). The benefit of neural scene representations over classic representations in terms of scene completion can easily be observed. While less drift accumulates on these scenes, our method still achieves on par results to the best performing method Co-SLAM. The point-based Loopy-SLAM appears less robust to noisy measurements and creates more noisy, less complete reconstructions.

Fig. 6 shows additional results on the Kintinuous sequences before and after loop closure highlighting the efficient loop closure integration of our method. Because none of the other methods successfully closed the loop, we compare with the monolithic ablation of our approach (see Sec. 4.4 for further discussion).

4.3. Run-Time and Model Size Analysis

Tab. 2 reports run-time³ and model size on three representative scenes of different sizes. Our approach compares favorably in terms of run-time despite the overhead of parallelized training of multiple fields. GO-SLAM and Loopy-SLAM slow down significantly for longer sequences, whereas our method slows down the least of all evaluated methods (5% difference between the smallest and largest scene). Model size is competitive with the submap-

¹The upper floor has a slight consistent orientation error even after loop closure, which makes many walls fall outside the 5 cm threshold explaining the mismatch between perceived quality and evaluation result.

²We could not achieve a successful run on apt0 without exceeding 24 GB of GPU memory due to too many created submaps.

³We report time per frame by benchmarking the execution time for the whole sequence and dividing it by the number of frames. Mesh extraction has been deactivated for all methods for fair comparison. All evaluations were run with an Intel Core i9-13900KF and NVIDIA GeForce RTX 4090.

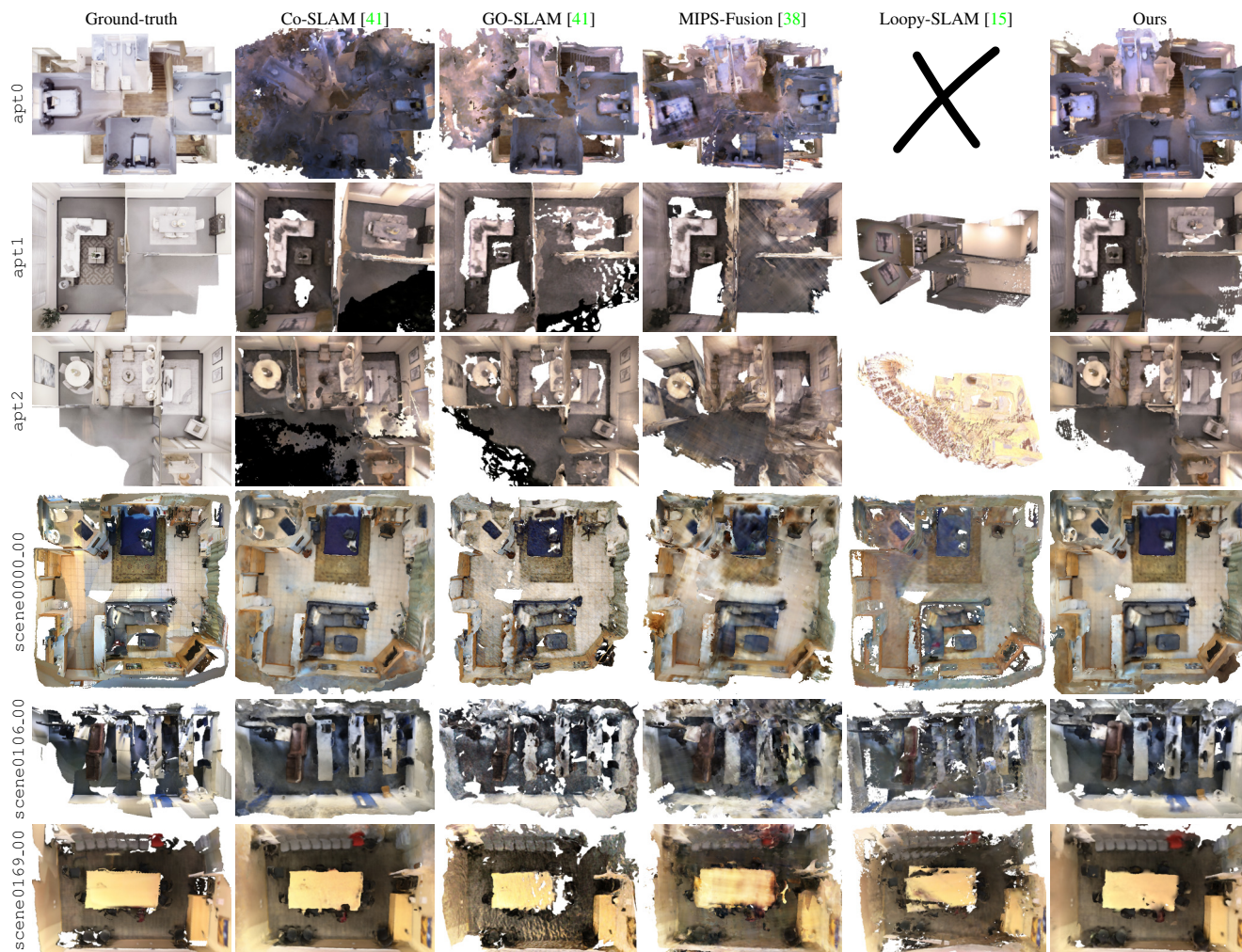


Figure 5. Qualitative comparison of final reconstruction on Replica-Big and ScanNet. Alternative view is shown for significant failure cases. The submap-based methods MIPS-Fusion and Loopy-SLAM exhibit artifacts at submap borders and erroneous submap registrations. Our method allows submap-like map deformation at much finer granularity, while significantly reducing transition artifacts.

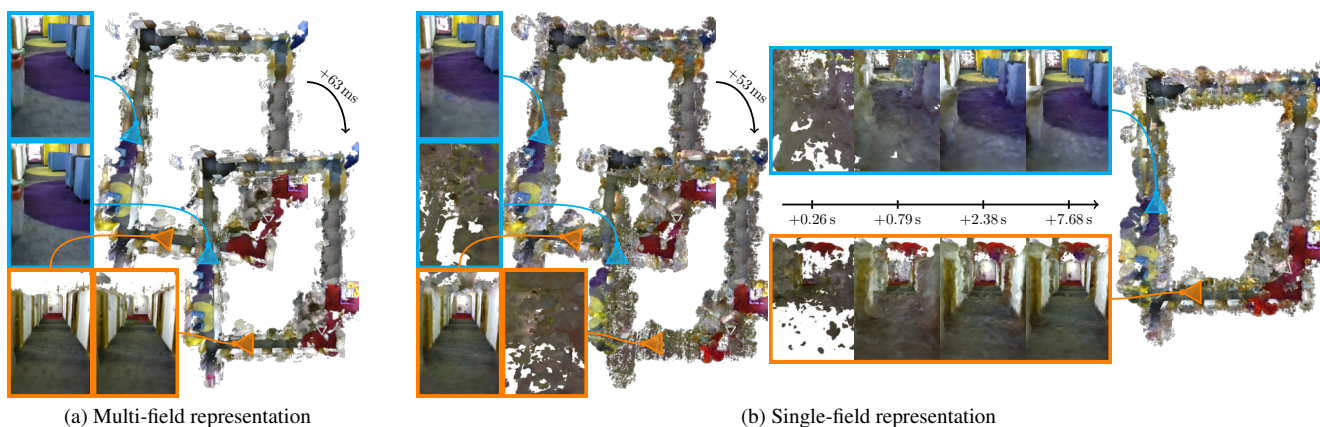


Figure 6. Qualitative results on the Kintuous sequence for our multi-field representation and the ablated, monolithic variant. We show a top-down and two additional views of the mesh right before and after specific durations. Our multi-field representation adapts within just one frame with minimal loss in quality, whereas the monolithic variant does not reach the same quality even after 7 s of optimization.

Table 2. Run-time and model size comparison of all methods. Our method compares favorably in terms of processing time compared to all baselines. The model size is larger than for the monolithic baselines, but competitive with the submap-based methods. Loopy-SLAM failed after 16 hours (0.1 Hz) on apt0.

	FPS (Hz)			Model Size (MB)		
	br	ck	apt0	br	ck	apt0
NICE-SLAM [46]	2.74	2.70	2.25	5.3	21.6	57.1
Co-SLAM [41]	12.70	12.84	10.75	7.0	7.0	7.0
GO-SLAM [45]	17.35	12.06	6.36	66.5	66.5	66.5
MIPS-Fusion [38]	4.02	3.64	3.31	36.2	72.4	398.3
Loopy-SLAM [15]	0.38	0.36	X	21.9	48.1	X
Ours-SF	19.78	18.92	18.84	8.4	8.4	8.4
Ours	16.67	15.93	15.68	28.1	57.7	229.2

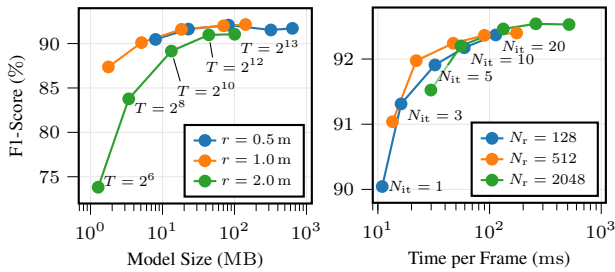


Figure 7. Left: reconstruction quality and model size for varying field radii r and hash table sizes T . Right: reconstruction quality and run-time for varying number of iterations per frame N_{it} and number of rays per field N_r .

based baselines, but worse in comparison to monolithic methods. This might be due to the more constrained hash tables, which can be overparameterized for simple areas.

Run-time results such as these are difficult to interpret due to varying implementation quality and large number of parameters. Especially NICE-SLAM and Loopy-SLAM are an order-of-magnitude slower than the other methods. Whether this is due to inefficient implementation or inherent complexity is unclear. Also, the baselines perform tracking and mapping, whereas our method uses precomputed ORB-SLAM results. However, ORB-SLAM is a CPU-only method and could therefore run in parallel without significantly affecting the run-time of our approach.

Field Radius and Hash Map Size In Fig. 7 (left) the memory requirement and F1-score for different field radii r and hash table sizes T is reported. In general, larger fields require larger hash maps to maintain the same quality. Our approach works across a wide range of r , but the ability to adapt to loop closures will be limited for larger r .

Batch Size and Iterations In Fig. 7 (right) the processing time per frame and F1-score for varying batch sizes and number of iterations is reported. Increasing the number of iterations per frame or the batch size leads to improved reconstruction quality. Notably, our approach maintains high reconstruction quality even at frame rates of over 50 Hz.



Figure 8. Single-view supervision leads to artifacts when the opposite side of a wall is observed. Multi-view supervision significantly reduces this local forgetting effect.

4.4. Ablation Study

Multi-Field Representation In Tab. 1 and Tab. 2 a comparison to a single-field variant of our method to isolate the effect of the multi-field representation. This variant uses the same rendering, poses, and optimization scheme as the multi-field variant, but represents the scene by a global monolithic field with $T = 2^{16}$. On the scenes with significant loop closures (apt0 and apt2) the multi-field representation outperforms the monolithic variant.

These results are at the end of the sequence and hence the monolithic variant has time to reintegrate with updated poses. Fig. 6 shows the map quality right after a large-scale loop closures. Our proposed multi-field representation adapts immediately to the pose adjustments keeping the map intact, while the single-field variant degrades significantly and only slowly reoptimizes without reaching the previous quality after more than 7 s.

Multi-View Supervision Fig. 8 compares optimization with targets sampled from all keyframes as described in Sec. 3.3 to targets sampled from a single view. In the latter case, optimization alternates between a random previous keyframe and the latest keyframe. Single-view supervision shows local forgetting effects, particularly notable when a previously seen wall is observed from the other side. Multi-view optimization avoids this by combining previous and current observations in each optimization step.

5. Conclusion

We have presented a neural mapping framework that anchors lightweight neural fields to the pose graph of a sparse visual SLAM system. This allows to incorporate loop closures into the volumetric map at near-zero cost achieving global map consistency over time. Our approach deforms the map reducing the need for reintegration while allowing geometric queries without requiring any space warping or image-space fusion. Compared to submap-based approaches, it deforms the map at a more granular level and reduces commonly observed artifacts at submap boundaries by avoiding hard assignments of frames to submaps.

Acknowledgments This work was partially supported by the Wallenberg AI, Autonomous Systems and Software Program (WASP) funded by the Knut and Alice Wallenberg Foundation.

References

- [1] Michal Adamkiewicz, Timothy Chen, Adam Caccavale, Rachel Gardner, Preston Culbertson, Jeannette Bohg, and Mac Schwager. Vision-only robot navigation in a neural radiance world. *RAL*, 7(2):4606–4613, 2022. [1](#)
- [2] Dejan Azinović, Ricardo Martin-Brualla, Dan B Goldman, Matthias Nießner, and Justus Thies. Neural RGB-D surface reconstruction. In *CVPR*, pages 6290–6301, 2022. [2](#), [4](#), [5](#)
- [3] Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J Guibas, Jonathan Tremblay, Sameh Khamis, et al. Efficient geometry-aware 3D generative adversarial networks. In *CVPR*, pages 16123–16133, 2022. [2](#), [5](#)
- [4] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *CVPR*, pages 5939–5948, 2019. [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 *CVPR*, pages 5828–5839, 2017. [5](#)
- [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 TOG*, 36(4):1, 2017. [1](#), [2](#), [5](#), [6](#)
- [7] Andrew J Davison, Ian D Reid, Nicholas D Molton, and Olivier Stasse. MonoSLAM: Real-time single camera SLAM. *IEEE TPAMI*, 29(6):1052–1067, 2007. [1](#)
- [8] Jiading Fang, Shengjie Lin, Igor Vasiljevic, Vitor Guizilini, Rares Ambrus, Adrien Gaidon, Gregory Shakhnarovich, and Matthew R Walter. NeRFuser: Large-scale scene representation by NeRF fusion. *arXiv preprint arXiv:2305.13307*, 2023. [2](#)
- [9] Christopher G Harris and JM Pike. 3D positional integration from image sequences. *Image and Vision Computing*, 6(2):87–90, 1988. [1](#)
- [10] Armin Hornung, Kai M Wurm, Maren Bennewitz, Cyrill Stachniss, and Wolfram Burgard. Octomap: An efficient probabilistic 3d mapping framework based on octrees. *Autonomous Robots*, 34:189–206, 2013. [1](#)
- [11] Peter J Huber. Robust estimation of a location parameter. *Ann. Math. Statist.*, 35:73–101, 1964. [4](#)
- [12] Mohammad Mahdi Johari, Camilla Carta, and François Fleuret. ESLAM: Efficient dense SLAM system based on hybrid representation of signed distance fields. In *CVPR*, pages 17408–17419, 2023. [1](#), [2](#), [5](#)
- [13] Georg Klein and David Murray. Parallel tracking and mapping for small ar workspaces. In *ISMAR*, pages 225–234, 2007. [1](#)
- [14] Xin Kong, Shikun Liu, Marwan Taher, and Andrew J Davison. vMAP: Vectorised object mapping for neural field SLAM. In *CVPR*, pages 952–961, 2023. [2](#), [3](#), [4](#), [5](#)
- [15] Lorenzo Liso, Erik Sandström, Vladimir Yugay, Luc Van Gool, and Martin R Oswald. Loopy-SLAM: Dense neural SLAM with loop closures. In *CVPR*, pages 20363–20373, 2024. [2](#), [3](#), [5](#), [6](#), [7](#), [8](#)
- [16] William E Lorensen and Harvey E Cline. Marching cubes: A high resolution 3d surface construction algorithm. In *ACM SIGGRAPH*, pages 163–169, 1987. [5](#)
- [17] Hidenobu Matsuki, Keisuke Tateno, Michael Niemeyer, and Federico Tombari. NEWTON: Neural view-centric mapping for on-the-fly large-scale SLAM. *RAL*, 2024. [2](#), [3](#)
- [18] Donald Meagher. Geometric modeling using octree encoding. *Computer Graphics and Image Processing*, 19(2):129–147, 1982. [1](#)
- [19] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3D reconstruction in function space. In *CVPR*, pages 4460–4470, 2019. [2](#)
- [20] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. NeRF: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, pages 405–421. Springer, 2020. [2](#)
- [21] Arthur Moreau, Nathan Piasco, Dzmitry Tsishkou, Bogdan Stanculescu, and Arnaud de La Fortelle. LENS: Localization enhanced by NeRF synthesis. In *CoRL*, pages 1347–1356, 2022. [1](#)
- [22] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM TOG*, 41(4), 2022. [2](#), [5](#)
- [23] Thomas Müller, Brian McWilliams, Fabrice Rousselle, Markus Gross, and Jan Novák. Neural importance sampling. *ACM TOG*, 38(5):1–19, 2019. [2](#)
- [24] Raul Mur-Artal and Juan D Tardós. ORB-SLAM2: an open-source SLAM system for monocular, stereo and RGB-D cameras. *TRO*, 33(5):1255–1262, 2017. [5](#)
- [25] 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 *ISMAR*, pages 127–136, 2011. [1](#)
- [26] Richard A Newcombe, Steven J Lovegrove, and Andrew J Davison. DTAM: Dense tracking and mapping in real-time. In *ICCV*, pages 2320–2327, 2011. [1](#)
- [27] Matthias Nießner, Michael Zollhöfer, Shahram Izadi, and Marc Stamminger. Real-time 3D reconstruction at scale using voxel hashing. *ACM TOG*, 32(6):1–11, 2013. [1](#)
- [28] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. In *CVPR*, pages 165–174, 2019. [2](#)
- [29] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An imperative style, high-performance deep learning library. In *NeurIPS*, volume 32, pages 8026–8037, 2019. [5](#)
- [30] Victor Reijgwart, Alexander Millane, Helen Oleynikova, Roland Siegwart, Cesar Cadena, and Juan Nieto. Voxgraph: Globally consistent, volumetric mapping using signed distance function submaps. *RAL*, 5(1):227–234, 2019. [1](#), [2](#)

- [31] Radu Alexandru Rosu and Sven Behnke. PermutoSDF: Fast multi-view reconstruction with implicit surfaces using permutohedral lattices. In *CVPR*, pages 8466–8475, 2023. 5
- [32] Erik Sandström, Yue Li, Luc Van Gool, and Martin R Oswald. Point-SLAM: Dense neural point cloud-based SLAM. In *ICCV*, pages 18433–18444, 2023. 2, 3
- [33] Noah Stier, Baptiste Angles, Liang Yang, Yajie Yan, Alex Colburn, and Ming Chuang. LivePose: Online 3D reconstruction from monocular video with dynamic camera poses. In *ICCV*, pages 7921–7930, 2023. 2
- [34] Julian Straub et al. The Replica dataset: A digital replica of indoor spaces. *arXiv preprint arXiv:1906.05797*, 2019. 5
- [35] Edgar Sucar, Shikun Liu, Joseph Ortiz, and Andrew J Davison. iMAP: Implicit mapping and positioning in real-time. In *ICCV*, pages 6229–6238, 2021. 1, 2, 5
- [36] Matthew Tancik, Vincent Casser, Xincheng Yan, Sabeek Pradhan, Ben Mildenhall, Pratul P Srinivasan, Jonathan T Barron, and Henrik Kretzschmar. Block-NeRF: Scalable large scene neural view synthesis. In *CVPR*, pages 8248–8258, 2022. 2
- [37] Matthew Tancik, Pratul Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. *NeurIPS*, 33:7537–7547, 2020. 1, 5
- [38] Yijie Tang, Jiazhao Zhang, Zhinan Yu, He Wang, and Kai Xu. MIPS-FUSION: Multi-implicit-submaps for scalable and robust online neural RGB-D reconstruction. *ACM TOG*, 42(6):1–16, 2023. 2, 3, 5, 6, 7, 8
- [39] Maxim Tatarchenko, Stephan R Richter, René Ranftl, Zhuwen Li, Vladlen Koltun, and Thomas Brox. What do single-view 3d reconstruction networks learn? In *CVPR*, pages 3405–3414, 2019. 5
- [40] Diego Thomas and Akihiro Sugimoto. Modeling large-scale indoor scenes with rigid fragments using RGB-D cameras. *Computer Vision and Image Understanding*, 157:103–116, 2017. 1
- [41] Hengyi Wang, Jingwen Wang, and Lourdes Agapito. Co-SLAM: Joint coordinate and sparse parametric encodings for neural real-time SLAM. In *CVPR*, pages 13293–13302, 2023. 1, 2, 4, 5, 6, 7, 8
- [42] 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*, 34, 2021. 2
- [43] Thomas Whelan, Michael Kaess, Hordur Johannsson, Maurice Fallon, John J Leonard, and John McDonald. Real-time large-scale dense RGB-D SLAM with volumetric fusion. *IJRR*, 34(4-5):598–626, 2014. 1, 2, 5
- [44] Xiaoshuai Zhang, Abhijit Kundu, Thomas Funkhouser, Leonidas Guibas, Hao Su, and Kyle Genova. Nerflets: Local radiance fields for efficient structure-aware 3D scene representation from 2D supervision. In *CVPR*, pages 8274–8284, 2023. 2
- [45] Youmin Zhang, Fabio Tosi, Stefano Mattoccia, and Matteo Poggi. GO-SLAM: Global optimization for consistent 3D instant reconstruction. In *ICCV*, pages 3727–3737, 2023. 1, 3, 5, 6, 8
- [46] 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 *CVPR*, 2022. 1, 2, 5, 6, 8