

# EC-SLAM: Real-time Dense Neural RGB-D SLAM System with Effectively Constrained Global Bundle Adjustment

Guanghao Li<sup>†</sup> Qi Chen<sup>†</sup> Yuxiang Yan Jian Pu<sup>\*</sup>

**Abstract**—We introduce EC-SLAM, a real-time dense RGB-D simultaneous localization and mapping (SLAM) system utilizing Neural Radiance Fields (NeRF). Although recent NeRF-based SLAM systems have demonstrated encouraging outcomes, they have yet to completely leverage NeRF’s capability to constrain pose optimization. By employing an effectively constrained global bundle adjustment (BA) strategy, our system makes use of NeRF’s implicit loop closure correction capability. This improves the tracking accuracy by reinforcing the constraints on the keyframes that are most pertinent to the optimized current frame. In addition, by implementing a feature-based and uniform sampling strategy that minimizes the number of ineffective constraint points for pose optimization, we mitigate the effects of random sampling in NeRF. EC-SLAM utilizes sparse parametric encodings and the truncated signed distance field (TSDF) to represent the map in order to facilitate efficient fusion, resulting in reduced model parameters and accelerated convergence velocity. A comprehensive evaluation conducted on the Replica, ScanNet, and TUM datasets showcases cutting-edge performance, including enhanced reconstruction accuracy resulting from precise pose estimation, 21 Hz run time, and tracking precision improvements of up to 50%. The source code is available at <https://github.com/Lightingooo/EC-SLAM>.

**Index Terms**—Neural Rendering, Neural Radiance Field, Simultaneous Localization and Mapping.

## I. INTRODUCTION

Over the past several decades, the field of real-time visual simultaneous localization and mapping (VSLAM) has rapidly advanced to become a cornerstone of research and development, particularly within the realms of autonomous driving and virtual reality. This groundbreaking technology, which enables the precise mapping of environments while simultaneously determining the observer’s location within that environment, has been instrumental in pushing the boundaries of what’s possible in navigation and immersive experiences. However, despite the impressive achievements of traditional VSLAM methodologies in terms of localization accuracy—allowing devices to pinpoint their position with remarkable precision—they often fall short when it comes to creating maps that are both dense and continuous. Such comprehensive and detailed mapping is crucial for the function-

ality of advanced applications, necessitating ongoing research and innovation to bridge this gap and meet the increasingly complex demands of these cutting-edge technologies.

Learning-based VSLAM systems have increasingly become a pivotal solution within the VSLAM domain, mainly due to their capacity for dense scene representation. Notably, VSLAM systems that utilize NeRF and its variants inputting RGB-D images have excelled in map and pose optimization via differentiable volume rendering. These systems are primarily categorized into two types based on their prediction methodologies: Multilayer Perceptron (MLP) based VSLAM [5], [6] and feature-decoder based VSLAM [1]–[3]. While MLP based VSLAM systems employ one or more MLPs for scene representation, offering smooth depiction, they face a challenge in pose constraints due to slow convergence. This slow convergence is attributed to the necessity of training all parameters in each iteration. Additionally, as the training scenarios increase in size, this integrated neural network struggles to address the issue of catastrophic forgetting.

The feature-decoder based VSLAM systems avoid training all parameters within one iteration and simultaneously ensure the smoothness of the scene through handcrafted loss functions, thus relatively speeding up the training process. However, compared to traditional VSLAM methods, they still operate at a slower speed and have inferior tracking accuracy. Furthermore, due to speed and memory constraints, these systems rely on random sampling methods for joint local/global BA and the continuous tracking of new keyframes. Specifically, Co-SLAM [3] selects a random number and sequence of keyframes and their corresponding pixels during global BA and randomly selects pixels during tracking. ES-SLAM [2] randomly samples pixel positions during local BA and tracking. This randomness, in turn, affects the constraints on the poses that need to be optimized, thereby influencing the accuracy of tracking and reconstruction.

We introduce EC-SLAM, a learning-based dense RGB-D SLAM system to address the above challenges. We employ sparse parametric encodings with TSDF and implement strategies such as effectively constrained global BA, robust sampling, and an accelerating algorithm. This enables us to achieve precise mapping and tracking precision, as demonstrated in Figure 1. Furthermore, we observed that SLAM systems based on NeRF inherently exhibit a specific skill for detecting loop closures. Our system significantly improves the capability of NeRF loop correction. NeRF’s implicit loop detection is more natural and effective compared to the explicit loop detection

<sup>†</sup>Indicates equal contribution.

<sup>\*</sup>Corresponding author.

G. Li, Y. Yan, and J. Pu are with the Institute of Science and Technology for Brain-Inspired Intelligence, Fudan University, Shanghai, 200433, China (E-mails: {ghli22, yxyan22}@m.fudan.edu.cn, jianpu@fudan.edu.cn).

Q. Chen is with Shanghai Key Lab of Intelligent Information Processing and School of Computer Science, Fudan University, Shanghai 200433, China (E-mail: qichen21@m.fudan.edu.cn).

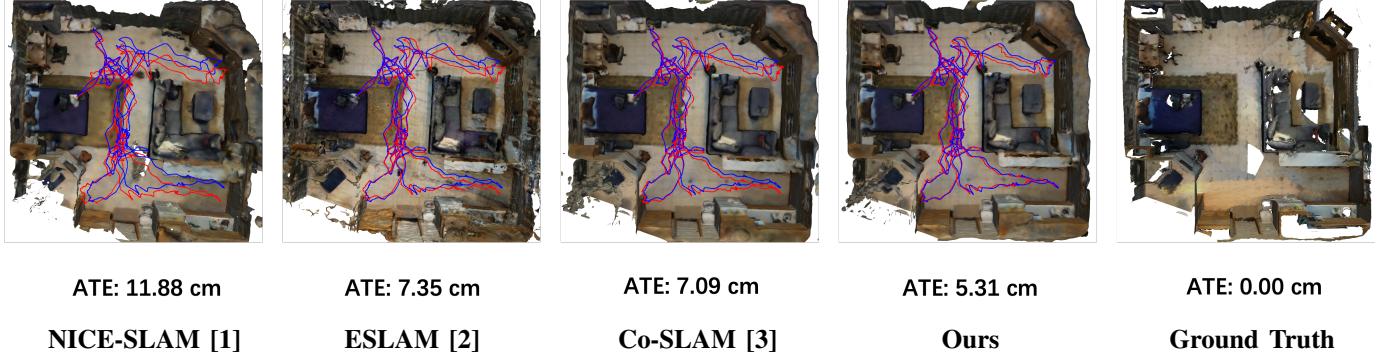


Fig. 1. 3D reconstruction and tracking of the trajectory for scene0000 from the ScanNet dataset [4]. The blue line depicts the ground truth camera trajectory, while the red line shows the estimated camera trajectory. Our method demonstrates superior performance in tracking and quality of scene reconstruction compared to other RGB-D methods such as NICE-SLAM [1], Co-SLAM [3], and ESLAM [2].

methods used in classical SLAMs, such as descriptor and bag-of-words matching. Extensive evaluations and ablation analyses carried out on several datasets (Replica [7], ScanNet [4], Tum [8]) demonstrate that our system provides superior reconstruction and tracking accuracy compared to other cutting-edge systems at a speed of up to 21 Hz.

The contributions of this paper can be summarized as follows. Our first contribution is the development of a novel NeRF-based architecture that effectively enforces constraints within the global bundle adjustment framework. This architecture also takes advantage of the implicit loop closure capabilities of NeRF, which are manifested in the selection of sliding windows, the organization of keyframes, and the optimization of pixel sampling. This strategy greatly improves the accuracy and reliability of tracking and mapping. Furthermore, we integrate sparse parametric encodings with TSDF in order to accurately represent and optimize maps. This method allows for the quick and precise reconstruction of relevant maps. At last, we have created a reliable and immediate RGB-D dense SLAM system that takes advantage of the individual advantages of NeRF and classic SLAM systems for mapping and tracking. Our system has been proven effective through extensive experiments. Our strategy can be used as a robust benchmark for future works.

## II. RELATED WORK

**Traditional Dense VSLAM.** Compared to SLAM systems that perform sparse reconstruction [9]–[25], fewer systems execute dense reconstruction, which are more computationally intensive. However, they provide a detailed and accurate map. DTAM [26], and Mobilefusion [27] were seminal direct methods using photometric error minimization. [28]–[30] were direct methods offering semi-dense reconstruction. Dense mapping SLAM systems also incorporated inertial measurements or depth data instead of visual-only input. KinectFusion [31] was a real-time RGB-D SLAM algorithm for 3D reconstruction and surface mapping, with limitations in handling drift over time. [32]–[34] employed loop closure detection to achieve better tracking and mapping results. Additionally, recent learning-based methods [35]–[48] have emerged. These methods can perform pose and depth estimation via a deep

neural network. However, these emerging systems’ scene representation and pipeline are still based on traditional systems.

**NeRF-based Dense VSLAM.** Although NeRF-based SLAM has lower tracking accuracy compared to traditional VSLAM, it is capable of real-time dense mapping and leverages this capability to utilize its implicit loop closure ability. It can mainly be divided into two broad categories: one integrates NeRF [49] with components from traditional SLAM systems [50]–[52], while the other solely rely on NeRF [1]–[3], [5], [53]–[55]. In the former type, NeRF-SLAM [51] combined the tracking component of DROID-SLAM [36] with its NeRF-based networks, while Orbeez-SLAM [50] integrated with modules from ORB-SLAM3 [18]. GO-SLAM [52] extended the tracking component of DROID-SLAM [36] by several key features. In the latter type, iMAP [5], and NICE-SLAM [1] used a single MLP and a coarse-to-fine feature grid to represent maps, respectively, but they suffered from slow convergence and struggled with scene details. Nicer-SLAM [54] performed mapping and tracking with monocular input instead of RGB-D input. For other RGB-D systems, Co-SLAM [3] and ESLAM [2] employed InstantNGP [56] and Tri-plane [57] as their scene representations, respectively, both utilizing TSDF prediction for quicker convergence. Nevertheless, these two systems overlooked the impact of randomness on the system: Co-SLAM [3] randomly selected the number and order of keyframes for global BA, failing to form effective global constraints for mapping and pose optimization. ESLAM [2] fixed the number of keyframes for BA, which limits the ability of the neural implicit field to perform implicit loop closure correction. Both systems employed random pixel sampling, which could not guarantee the uniformity and constraint of the sampled pixels. The randomness mentioned above results in these methods exhibiting weak robustness across datasets and significant variance in re-identification effectiveness for the same scene.

Our system belongs to the type that solely relies on NeRF [49]. Unlike the systems mentioned above, we have meticulously optimized the number and sequence of keyframes participating in BA, as well as the method of pixel sampling, to implement an effectively constrained global BA and utilize NeRF [49] for implicit loop closure correction. Moreover, our

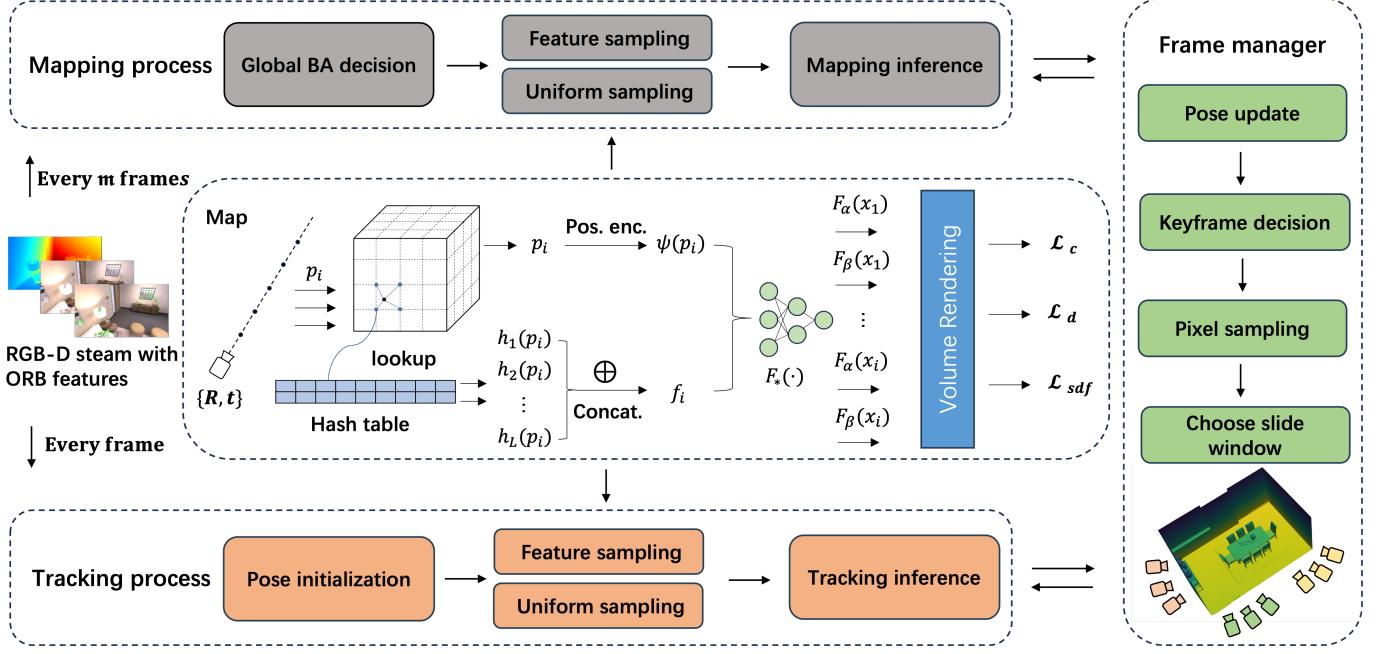


Fig. 2. The overview of EC-SLAM. (1) Mapping process: We jointly optimize the map and the poses of some keyframes in the sliding window. (2) Tracking process: We use a constant velocity motion model to initialize the tracking frame's pose, followed by iterative optimization. (3) Map: We employ a multi-resolution hash grid to store feature values for each point in space and use two MLPs to decode these features, obtaining color and TSDF values. Here  $f_i$ ,  $\psi(\mathbf{p}_i)$  and  $F_*(\cdot)$  denotes the multi-resolution feature, position encoding and decoder function for a certain point, respectively. We compute color, depth, and TSDF losses to optimize the map and pose.  $\mathcal{L}_c$ ,  $\mathcal{L}_d$ , and  $\mathcal{L}_{sdf}$  are the color loss, depth loss, and TSDF loss, respectively.

sparse parametric encodings and specialized loss functions contribute to the effect of real-time, small-parameter, and high-precision in mapping and tracking. This focus on optimization and precision underscores our system's advanced capabilities in the evolving landscape of VSLAM technology.

### III. METHOD

Given a set of sequential RGB-D frames  $\{I_i, D_i\}_{i=1}^M$  with known camera intrinsics  $K \in \mathbb{R}^{3 \times 3}$ , we predict camera poses  $\{R_i | t_i\}_{i=1}^M$  and an implicit TSDF map representation. As shown in Fig. 2, we established two threads to address this problem: the mapping and tracking threads. The mapping thread is responsible for jointly optimizing the map and the poses of the keyframes in the sliding window, while the tracking thread estimates the pose of the new image frame. Consistent with the previous system [3], we optimize the map after a fixed number of frames. The map, represented as a hash map, is a shared variable transmitted between different threads. Each thread performs volume rendering based on the corresponding pose and hash map, then computes the loss against the actual image.

#### A. Robust multi-threaded Pipeline

**System Pipeline.** We introduce our system from two aspects: the mapping thread and the tracking thread. Similar to the majority of NeRF-SLAM systems [1]–[3], [5], the mapping thread is responsible for both map initialization and subsequent mapping tasks. We randomly initialize the network parameters and fix the pose of the first image frame for mapping. In the

subsequent iterations, we perform a mapping operation at fixed intervals of image frames. In each iteration, we first decide on the sliding window containing keyframes, then perform our robust pixel sampling method on all the keyframes in the sliding window. Finally, we jointly optimize our map and the poses of some keyframes in the sliding window, which is called effectively constrained global BA.

During tracking, we first calculate the initial pose of the frame to be tracked based on a constant velocity motion model. Then, we sample the tracking image frame using a sampling strategy similar to the mapping thread. Finally, we fix our map parameters and iteratively optimize the input pose.

**Effectively Constrained Global BA.** In our proposed methodology, we address the limitations observed in certain NeRF-based SLAM systems. Previous systems [1], [2] often employ a process where pixels and points sampled from the current mapping frame are projected across all other keyframes to determine the sliding window. This process, prone to randomness, proves to be excessively time-consuming. Furthermore, these systems typically utilize a sliding window of limited size, which does not adequately support global optimization or facilitate implicit loop closure correction. On the other hand, some systems [3] treat the sliding window as an approximation of all keyframes. Although this considers global optimization, it does not identify the frames most correlated with the current frame's pose. Furthermore, as the number of keyframes grows, the number of sampled pixels per keyframe decreases, and the distribution of sampled pixels among frames becomes uneven, consequently weakening the pixel constraints on optimizing

keyframe poses.

We adopt a fast and effective strategy to determine the number and order of keyframes in the sliding window. To eliminate the impact of randomness, we use only the pose to decide the sliding window. Precisely, we first calculate the optical center distance and parallax angle between all keyframes and the current frame's camera, then filter out keyframes with large parallax angles and sort them based on distance. Regarding the selection of keyframe quantity, to ensure the effectiveness of global BA, we set a maximum number threshold based on the number of sampled pixels. As long as the number of keyframes in the sliding window does not exceed this threshold, we strive to include keyframes that can form loop closure constraints with the current frame, ensuring the presence of historical and the most recent frames. This approach not only forms practical global BA constraints but also fully leverages the inherent implicit loop closure capabilities of NeRF [49].

**Robust Pixel Sampling Method.** A distinctive feature of NeRF-based SLAM [49] is random pixel sampling, which significantly influences the results. We employ two sampling methods that mitigate this issue: feature point sampling and random uniform sampling. We initially extract feature points [58] for each image to serve as a subset of the sampled pixel points. For the remaining pixels to be sampled, we establish a grid over the image and perform uniform sampling within this grid.

Moreover, sampling at each iteration was excessively time-consuming, so we perform a one-time sampling in the camera coordinate system before iteration. During each iteration, we project the relevant pixels to the world coordinate system by using the current optimized pose. Therefore, our robust pixel sampling method significantly reduces time consumption on iterative sampling without substantially increasing memory usage.

### B. Map Representation and Volume Rendering

**Map Representation.** The quantity of grid parameters can significantly influence the convergence speed and reconstruction accuracy when restricting the convergence speed to feature-decoder based methods. To mitigate this issue, we adopt the multi-resolution hash encoding proposed in [56] as our scene representation structure. For a point  $\mathbf{x}_i$  in space, we have:

$$\mathbf{f}_i = \bigoplus_{l=1}^L h_l(\mathbf{x}_i), \quad (1)$$

where  $L$  denotes the number of resolution levels in the hash grid.  $h_l(\cdot)$  is the hash lookup and interpolation function that performs linear interpolation in the corresponding level for a certain point.  $\mathbf{f}_i$  is the final feature obtained after concatenating each level's feature.

In addition to the network architecture, we utilize the One-blob [59] encoding for positional representation as referenced in [3]. We ascertain a given spatial point's One-blob encoding and hash features, which are subsequently fed into the corresponding TSDF and color decoders. Unlike [3], we do not

employ the intermediate embedding generated by the TSDF decoder when decoding color. Thus, the inputs to both the TSDF and color decoders are identical. We observed that such an approach facilitates more rapid network convergence.

**Volume Rendering.** We employ a uniform sampling strategy based on depth and spatial distribution to acquire points from the optical center to the object surface. Here,  $N_u$  and  $N_d$  represent the number of points sampled uniformly in space and based on depth. For depth-based sampling, we sample  $N_d$  points within the interval  $[d - d_t, d + d_t]$ , where  $d$  is the object surface depth and  $d_t$  is a hyperparameter that defines the sampling density.

As described in [2], [3], predicting the TSDF values for spatial points instead of directly predicting their volumetric density can significantly enhance the speed and accuracy of network convergence because it allows constraints to be applied to the TSDF values at each point along a ray rather than constraining the entire ray. As detailed in Equ. 2, we employ the method outlined in [60] to transform the predicted TSDF values into volumetric density:

$$\sigma(s_i) = \beta \cdot \text{Sigmoid}(-\beta \cdot s_i), \quad (2)$$

where  $s_i$  is the TSDF of a certain point, and  $\beta$  is a learnable parameter that controls the sharpness of the surface boundary [2].

After determining the volumetric density of all points along a ray, we follow the standard volume rendering procedure in [49] to render the weights for all points on the ray, thereby obtaining the final color and depth corresponding to the pixel associated with that ray (Equ. 3). Here,  $N = N_u + N_d$  represents the total number of points sampled along the ray,  $w_n$  is the weight for space point  $n$ ,  $z_n$  is the sampled depth along each ray, and  $\hat{\mathbf{c}}$  and  $\hat{d}$  are the estimated color and depth.

$$w_n = \exp \left( - \sum_{i=1}^{n-1} \sigma(s_i) \right) (1 - \exp(-\sigma(s_n))), \\ \hat{\mathbf{c}} = \sum_{n=1}^N w_n \mathbf{c}_n \quad \text{and} \quad \hat{d} = \sum_{n=1}^N w_n z_n. \quad (3)$$

### C. Objective Functions

The objective function for measuring network performance consists of three components: color loss, depth loss, and TSDF loss. We employ the L2 loss to quantify the discrepancy between rendered and actual values for color and depth:

$$\mathcal{L}_d = \frac{\lambda_{ud}}{|R_u|} \sum_{r \in R_u} (\hat{d}_r - d_r)^2 + \frac{\lambda_{fd}}{|R_f|} \sum_{r \in R_f} (\hat{d}_r - d_r)^2, \\ \mathcal{L}_c = \frac{\lambda_{uc}}{|R_u|} \sum_{r \in R_u} (\hat{\mathbf{c}}_r - \mathbf{c}_r)^2 + \frac{\lambda_{fc}}{|R_f|} \sum_{r \in R_f} (\hat{\mathbf{c}}_r - \mathbf{c}_r)^2. \quad (4)$$

$R_f$  and  $R_u$  are the feature-based and uniform sampled rays, respectively.  $R = R_f + R_u$  represents all the sampled rays.  $\lambda_*$  represents the weight for each loss.

To further reinforce constraints on each sampled spatial point along each ray, we calculate the corresponding TSDF

loss for each point based on its distance from the scene surface, following the approach in [2]. For points beyond the truncation region  $T$ , we assess the difference between the predicted TSDF value and the actual value of 1:

$$\mathcal{L}_{fs} = \frac{1}{|R|} \sum_{r \in R} \frac{1}{|P_r^{fs}|} \sum_{p \in P_r^{fs}} (s_p - 1)^2, \quad (5)$$

where  $P_r^{fs}$  denotes the points along the ray beyond the truncation region  $T$ . For points within the truncation region  $T$ , the closer a point is to the surface, the more its TSDF value approaches zero, and TSDF values are negative for points inside the surface:

$$\begin{aligned} \mathcal{L}_t &= \frac{1}{|R|} \sum_{r \in R} \frac{1}{|P_r^t|} \sum_{p \in P_r^t} (z_p + s_p \cdot T - d_r)^2, \\ \mathcal{L}_m &= \frac{1}{|R|} \sum_{r \in R} \frac{1}{|P_r^m|} \sum_{p \in P_r^m} (z_p + s_p \cdot T - d_r)^2, \end{aligned} \quad (6)$$

here,  $P_r^m$  denotes the set of points that  $|z_p - d_r| < \gamma T$  for a ray  $r$ , and  $P_r^t$  is the other points in the truncation region.  $\gamma$  is the parameter that controls the number of spatial points within  $P_r^m$ .

The total loss is the weighted sum of the losses above, where the weights act as hyperparameters reflecting each loss component's relative importance:

$$\mathcal{L} = \lambda_{fs} \mathcal{L}_{fs} + \lambda_m \mathcal{L}_m + \lambda_t \mathcal{L}_t + \mathcal{L}_d + \mathcal{L}_c. \quad (7)$$

We apply smooth loss from [3] at the last image frame to make the scene more consistent.

#### IV. EXPERIMENTS

This section demonstrates that our system outperforms current NeRF based systems on three widely-used datasets while maintaining a rapid processing speed. The ablation studies further confirm the efficacy of our approach.

##### A. Experimental Setup

**Baselines.** We benchmark our method against four existing state-of-the-art NeRF based RGB-D dense VSLAM methods: iMAP [5], NICE-SLAM [1], ESLAM [2] and Co-SLAM [3] and three classic methods: Kintuino [31], BAD-SLAM [61], ORB-SLAM2 [17].

**Datasets.** We evaluate our method on three standard 3D benchmarks. Following [1]–[3], we quantitatively evaluate the reconstruction and tracking quality on eight synthetic scenes (room0, room1, room2, office0, office1, office2, office3, office4) from Replica [7]. We also evaluate the tracking results on six scenes (scene0000, scene0059, scene0106, scene0169, scene0181, scene0207) from ScanNet [4] and three scenes (fr1/desk, fr2/xyz, fr3/office) from TUM RGB-D [8] datasets.

**Metrics.** Our assessment of reconstruction quality includes metrics such as Depth L1 (in cm), Accuracy (in cm), Completion (in cm), and Completion Ratio (in %) with a threshold set at 5 cm. In line with [1], [3], we exclude areas not observed by any camera frustum and apply mesh culling to remove extraneous points within the camera's view but outside the intended scene. We utilize the ATE RMSE metric [8] (in cm

for Replica [7], m for others) for camera tracking evaluation. In addition, we calculate the Frames Per Second (FPS) following the way in Co-SLAM [3].

**Implementation Details** We run our system on a desktop PC with an Intel Core i7-12700 CPU and NVIDIA RTX 3090 GPU. The resolution to generate the scene mesh is set to 0.02 m for all systems to compare the reconstruction metric fairly. The decoders used for predicting SDF and color consist of two two-layer fully connected networks whose hidden feature size is 32. Following the configuration of [3], we sample 5% of the total pixels for each keyframe using our sampling strategy and store them in the keyframe set. The number of iterations during mapping and tracking and the number of points sampled on the rays corresponding to pixels varies across different datasets. We use two configurations for the Replica dataset [7], which differs in the iteration times in mapping and tracking.

In the mapping thread, we perform mapping every five frames. For the Replica dataset [7], we employ two parameter configurations. In our lite version, the mapping thread iterates 10 times, sampling 2000 pixels per iteration. In our full version, it iterates 20 times, sampling 4000 pixels each time. The lite and full versions operate at speeds of 21 Hz and 12 Hz, respectively. For the ScanNet dataset [4], the mapping thread iterates 10 times with speed of 8 Hz, sampling 4000 pixels per iteration. For the TUM dataset [8], it iterates 20 times with speed of 7 Hz, sampling 2000 pixels each time.

In the tracking thread, we track every frame and use a constant velocity assumption to compute the initial pose of the tracking frame. For the two versions of the Replica dataset [7], the lite version iterates 8 times, sampling 1000 pixels per iteration, while the full version iterates 20 times, sampling 2000 pixels each time. In the ScanNet dataset [4], we iterate 10 times, sampling 2000 pixels per iteration. In the TUM dataset [8], we iterate 20 times, sampling 2000 pixels each time.

For the number of samples per ray, following the approach of Co-SLAM [3], we employ both spatially uniform sampling and depth-guided sampling along a ray. We set  $N_u$  to 32 and  $N_d$  to 11 for the Replica dataset [7]. For ScanNet dataset [4], we set  $N_u$  to 96 and  $N_d$  to 21. On the TUM dataset [8],  $N_u$  is set to 64 and  $N_d$  to 21.

Please note that we have implemented a series of acceleration optimizations, such as pre-sampling pixels in the camera coordinate system and transforming them into the world coordinate system, thereby significantly enhancing the efficiency of our system. These acceleration optimizations enable our system to operate up to 50% faster than other systems under the same parameter configurations across various datasets and scenarios.

##### B. Experimental Results

**Evaluation on Replica [7].** We evaluated the reconstruction and tracking performance of our system in eight scenes of the Replica dataset [7], with results detailed in Tab. I. Each metric consists of the mean and variance obtained by randomly running the system five times on each scene of the dataset. It can be observed that other methods exhibit significant fluctuations in accuracy across different scenes, indicating

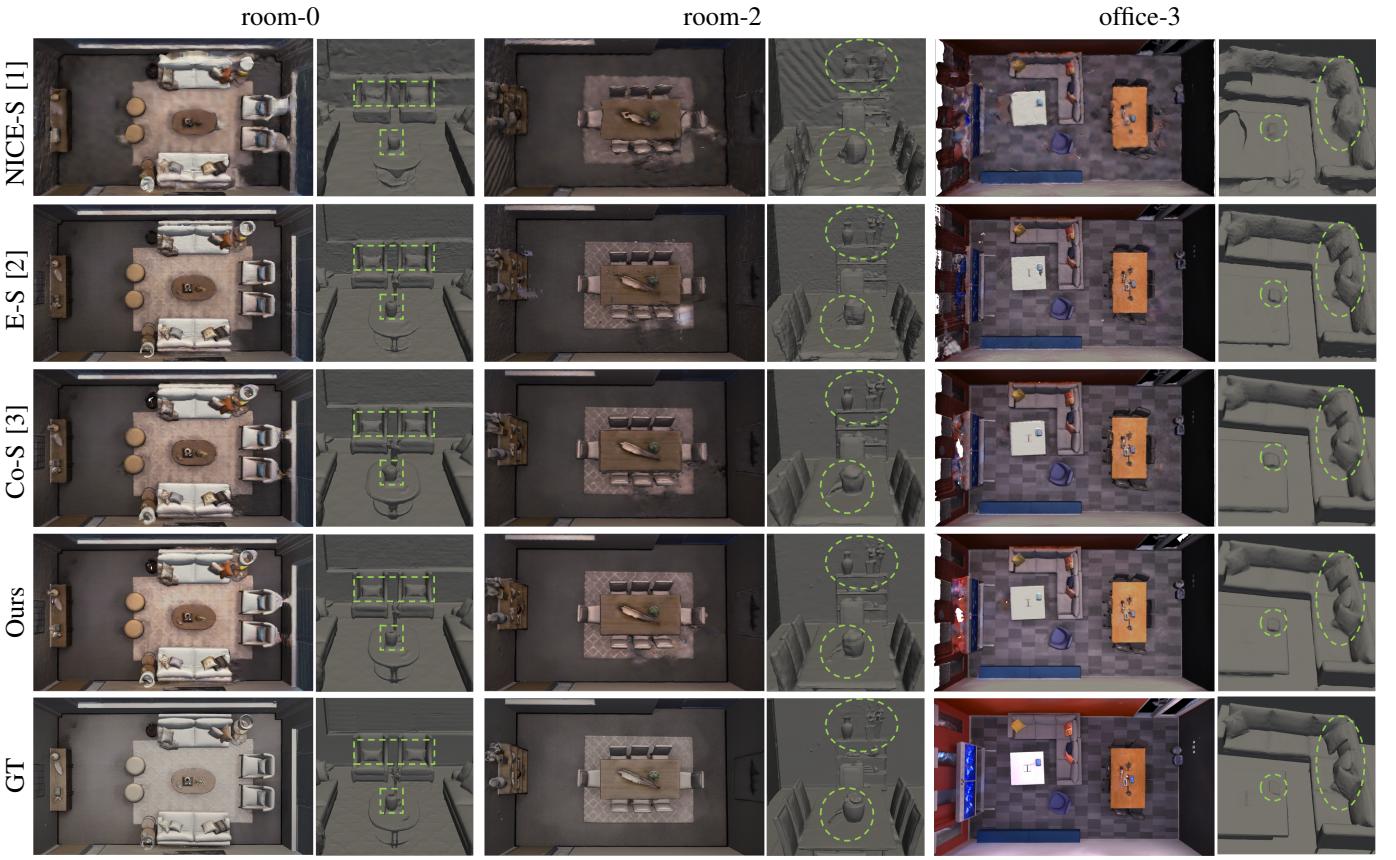


Fig. 3. The reconstruction results of our system with other SOTA NeRF-based RGB-D dense visual SLAM systems on the Replica dataset [7]. Our system reconstructed more accurate scenes compared with other systems, and the reconstruction details encircled by the green dashed line further demonstrate our system’s better trade-off between sharpness and smoothness.

their poor robustness. This is attributed to the randomness of sampling in NeRF. Both configurations of our system outperform other methods in tracking metrics with speed of 12-21 Hz, achieving up to a 50% improvement in accuracy. For the Lite configuration, we achieved a precision similar to ESLAM [2] (the previous state-of-the-art) with an FPS of 21Hz. For the Full configuration, we improved the precision by 50% relative to ESLAM [2] , at a similar FPS of approximately 12Hz.

Regarding mapping metrics, our system also reached SOTA levels, and our model has a minimal number of parameters (Tab. IV). The qualitative reconstruction results of the Replica dataset [7] are shown in Fig. 3, where our system constructs more accurate scenes than other systems. The mapping results of other systems are either overly smooth, failing to adequately display high-frequency information, or excessively sharp, unable to represent the natural continuity of the scene. The reconstructed scenes of our system exhibit a good tradeoff between sharpness and smoothness, which is particularly evident in the reconstruction of small objects within the scenes, benefiting from our accurate pose estimation and sampling strategy.

**Evaluation on ScanNet [4].** Although the Replica dataset [7] provides accurate reconstructed scenes for evaluating reconstruction metrics, as a synthetic indoor dataset, it does not fully represent the real world. Therefore, we also tested our system on the ScanNet dataset [7], a larger-scale real indoor

dataset, which can exhibits the robustness of our system. The tracking results of our system on ScanNet [4] are detailed in Tab. II. Each metric consists of the mean and variance obtained by randomly running the system five times on each scene. Our system outperforms other systems in most scenes with a speed of 8 Hz, demonstrating its robustness to real-world environments. Our system also demonstrates higher stability in multiple runs, with much lower standard deviations.

Fig. 4 shows our system’s reconstruction depth maps and tracking trajectory compared with other baselines on some ScanNet [4] scenes. The tracking results show that our estimated trajectory closely matches the ground truth trajectory. The mapping results indicate that the scenes reconstructed by our system more closely resemble the original scenes, especially in terms of overall structure and shape. For instance, regarding the upper left overall structure of scene 0000, the results computed by our system align more closely with the actual structure. In contrast, the upper left overall structure of the scene generated by other methods is slightly more slanted and larger compared to the Ground Truth . Fig. 5 presents a comparison of scene detail depth maps under the same camera pose, where our method yields more precise and smoother scene details. For instance, in the scene000 scenario, the positions of various components of the scene displayed by our method under the same camera pose are most similar to those in the ground truth scene.

TABLE I  
MAPPING AND TRACKING RESULTS ON THE REPLICA DATASET [7]

	Methods	Reconstruction(cm)			Localization(cm)		FPS	
		Depth L1↓	Acc.↓	Comp.↓	Comp. Ratio(%)↑	ATE RMSE↓		
Room0	iMAP [5]	5.17 ± 0.33	4.14 ± 0.25	5.99 ± 0.22	77.84 ± 1.12	5.43 ± 1.28	3.20 ± 0.75	0.4
	NICE-SLAM [1]	1.83 ± 0.13	2.48 ± 0.06	2.68 ± 0.08	91.66 ± 0.13	1.64 ± 0.11	1.43 ± 0.09	0.9
	Co-SLAM [3]	1.05 ± 0.03	2.03 ± 0.05	2.07 ± 0.04	95.16 ± 0.04	0.70 ± 0.03	0.63 ± 0.02	17.1
	ESLAM [2]	0.86 ± 0.03	2.52 ± 0.04	<b>1.98 ± 0.01</b>	<b>96.04 ± 0.05</b>	0.70 ± 0.13	0.60 ± 0.06	12.1
	Ours (Lite)	0.66 ± 0.01	1.90 ± 0.04	2.09 ± 0.04	95.36 ± 0.06	0.77 ± 0.03	0.66 ± 0.03	<b>21.4</b>
	Ours (Full)	<b>0.59 ± 0.01</b>	<b>1.89 ± 0.02</b>	2.07 ± 0.03	95.79 ± 0.04	<b>0.27 ± 0.01</b>	<b>0.24 ± 0.01</b>	11.9
Room1	iMAP [5]	3.56 ± 0.68	3.16 ± 0.36	4.57 ± 0.40	85.39 ± 1.05	3.22 ± 0.36	2.66 ± 0.26	0.4
	NICE-SLAM [1]	1.41 ± 0.16	2.14 ± 0.06	2.23 ± 0.09	93.42 ± 0.19	2.08 ± 0.08	1.70 ± 0.29	0.9
	Co-SLAM [3]	0.85 ± 0.02	1.60 ± 0.04	1.86 ± 0.05	95.19 ± 0.03	1.09 ± 0.34	0.96 ± 0.35	17.1
	ESLAM [2]	<b>0.88 ± 0.04</b>	2.51 ± 0.05	1.78 ± 0.01	95.09 ± 0.06	0.70 ± 0.02	0.55 ± 0.02	12.1
	Ours (Lite)	0.86 ± 0.01	1.67 ± 0.04	1.77 ± 0.05	95.49 ± 0.07	1.07 ± 0.16	0.95 ± 0.16	<b>21.4</b>
	Ours (Full)	0.89 ± 0.01	<b>1.56 ± 0.02</b>	<b>1.66 ± 0.05</b>	<b>96.51 ± 0.05</b>	<b>0.28 ± 0.02</b>	<b>0.24 ± 0.02</b>	11.9
Room2	iMAP [5]	5.87 ± 0.78	3.96 ± 0.39	5.23 ± 0.27	79.02 ± 1.63	2.85 ± 0.16	2.41 ± 0.20	0.4
	NICE-SLAM [1]	2.22 ± 0.20	2.21 ± 0.09	2.85 ± 0.08	91.32 ± 0.29	1.80 ± 0.28	1.41 ± 0.24	0.9
	Co-SLAM [3]	2.37 ± 0.02	1.95 ± 0.04	1.99 ± 0.05	93.48 ± 0.04	1.21 ± 0.14	0.82 ± 0.02	17.1
	ESLAM [2]	1.18 ± 0.03	1.76 ± 0.06	1.78 ± 0.01	95.88 ± 0.06	0.51 ± 0.01	0.42 ± 0.01	12.1
	Ours (Lite)	1.10 ± 0.01	1.72 ± 0.03	1.67 ± 0.05	95.95 ± 0.06	0.64 ± 0.01	0.55 ± 0.02	<b>21.4</b>
	Ours (Full)	<b>1.11 ± 0.01</b>	<b>1.69 ± 0.01</b>	<b>1.59 ± 0.04</b>	<b>97.30 ± 0.04</b>	<b>0.32 ± 0.03</b>	<b>0.29 ± 0.02</b>	11.9
Office0	iMAP [5]	3.99 ± 0.34	3.37 ± 0.32	3.91 ± 0.37	83.01 ± 1.60	2.60 ± 0.82	1.74 ± 0.90	0.4
	NICE-SLAM [1]	1.45 ± 0.14	1.87 ± 0.06	1.92 ± 0.09	94.80 ± 0.39	1.23 ± 0.24	1.12 ± 0.22	0.9
	Co-SLAM [3]	1.24 ± 0.02	1.47 ± 0.05	1.63 ± 0.05	96.09 ± 0.05	0.56 ± 0.01	0.49 ± 0.01	17.1
	ESLAM [2]	0.77 ± 0.02	1.61 ± 0.07	1.50 ± 0.01	97.32 ± 0.05	0.56 ± 0.04	0.41 ± 0.03	12.1
	Ours (Lite)	0.71 ± 0.01	1.42 ± 0.04	1.37 ± 0.04	97.70 ± 0.07	0.58 ± 0.03	0.49 ± 0.02	<b>21.4</b>
	Ours (Full)	<b>0.72 ± 0.01</b>	<b>1.39 ± 0.03</b>	<b>1.33 ± 0.04</b>	<b>97.79 ± 0.05</b>	<b>0.24 ± 0.01</b>	<b>0.21 ± 0.00</b>	11.9
Office1	iMAP [5]	3.81 ± 0.39	2.13 ± 0.25	3.97 ± 0.24	88.05 ± 1.75	1.30 ± 0.23	1.18 ± 0.13	0.4
	NICE-SLAM [1]	1.64 ± 0.11	1.62 ± 0.13	1.86 ± 0.10	93.94 ± 0.21	0.79 ± 0.17	0.74 ± 0.19	0.9
	Co-SLAM [3]	1.48 ± 0.02	1.27 ± 0.04	1.64 ± 0.04	94.55 ± 0.03	0.60 ± 0.06	0.53 ± 0.06	17.1
	ESLAM [2]	1.22 ± 0.02	1.98 ± 0.08	<b>1.41 ± 0.01</b>	96.66 ± 0.08	0.54 ± 0.04	0.45 ± 0.05	12.1
	Ours (Lite)	1.25 ± 0.01	1.30 ± 0.03	1.53 ± 0.05	95.73 ± 0.07	0.70 ± 0.05	0.63 ± 0.05	<b>21.4</b>
	Ours (Full)	<b>1.15 ± 0.01</b>	<b>1.22 ± 0.02</b>	1.46 ± 0.05	<b>96.84 ± 0.05</b>	<b>0.25 ± 0.01</b>	<b>0.23 ± 0.01</b>	11.9
Office2	iMAP [5]	4.02 ± 0.53	4.18 ± 0.41	4.89 ± 0.37	79.17 ± 1.16	5.94 ± 1.60	4.08 ± 0.78	0.4
	NICE-SLAM [1]	2.71 ± 0.13	3.32 ± 0.13	3.23 ± 0.08	88.12 ± 0.28	1.69 ± 0.14	1.42 ± 0.10	0.9
	Co-SLAM [3]	1.86 ± 0.07	2.74 ± 0.08	2.48 ± 0.07	91.63 ± 0.06	2.08 ± 0.02	1.84 ± 0.02	17.1
	ESLAM [2]	<b>1.06 ± 0.02</b>	2.87 ± 0.09	2.05 ± 0.02	94.38 ± 0.02	0.57 ± 0.09	0.46 ± 0.03	12.1
	Ours (Lite)	1.44 ± 0.01	3.13 ± 0.04	1.99 ± 0.06	94.20 ± 0.06	0.69 ± 0.03	0.60 ± 0.03	<b>21.4</b>
	Ours (Full)	1.36 ± 0.01	<b>2.39 ± 0.02</b>	<b>2.00 ± 0.05</b>	<b>94.72 ± 0.05</b>	<b>0.31 ± 0.02</b>	<b>0.29 ± 0.02</b>	11.9
Office3	iMAP [5]	5.71 ± 0.95	4.28 ± 0.33	5.65 ± 0.20	73.42 ± 0.65	5.18 ± 1.24	4.16 ± 0.78	0.4
	NICE-SLAM [1]	2.17 ± 0.40	3.05 ± 0.20	3.27 ± 0.09	87.53 ± 0.45	3.90 ± 1.34	2.31 ± 0.31	0.9
	Co-SLAM [3]	1.66 ± 0.03	3.01 ± 0.05	2.77 ± 0.04	90.62 ± 0.03	1.58 ± 0.04	1.50 ± 0.05	17.1
	ESLAM [2]	1.02 ± 0.04	2.53 ± 0.07	<b>2.27 ± 0.01</b>	94.01 ± 0.03	0.71 ± 0.02	0.60 ± 0.03	12.1
	Ours (Lite)	0.84 ± 0.01	2.60 ± 0.03	2.36 ± 0.05	93.36 ± 0.05	0.88 ± 0.01	0.74 ± 0.01	<b>21.4</b>
	Ours (Full)	<b>0.79 ± 0.01</b>	<b>2.49 ± 0.02</b>	2.33 ± 0.04	<b>94.13 ± 0.04</b>	<b>0.38 ± 0.03</b>	<b>0.35 ± 0.03</b>	11.9
Office4	iMAP [5]	5.81 ± 1.06	4.53 ± 0.24	6.81 ± 0.20	74.29 ± 0.93	2.42 ± 0.23	2.04 ± 0.16	0.4
	NICE-SLAM [1]	2.10 ± 0.41	2.58 ± 0.23	3.72 ± 0.24	87.08 ± 1.35	2.77 ± 0.60	2.22 ± 0.39	0.9
	Co-SLAM [3]	1.54 ± 0.02	2.41 ± 0.06	2.50 ± 0.05	90.32 ± 0.04	0.71 ± 0.01	0.62 ± 0.02	17.1
	ESLAM [2]	1.10 ± 0.04	2.14 ± 0.10	<b>2.28 ± 0.01</b>	93.10 ± 0.14	0.62 ± 0.03	0.51 ± 0.02	12.1
	Ours (Lite)	0.80 ± 0.01	2.15 ± 0.03	2.47 ± 0.06	91.95 ± 0.05	0.89 ± 0.07	0.78 ± 0.07	<b>21.4</b>
	Ours (Full)	<b>0.79 ± 0.01</b>	<b>1.96 ± 0.02</b>	2.36 ± 0.05	<b>93.24 ± 0.04</b>	<b>0.30 ± 0.03</b>	<b>0.28 ± 0.03</b>	11.9
Average	iMAP [5]	4.74 ± 0.63	3.71 ± 0.31	5.12 ± 0.29	80.02 ± 1.23	3.61 ± 0.74	2.68 ± 0.49	0.4
	NICE-SLAM [1]	1.94 ± 0.21	2.41 ± 0.12	2.72 ± 0.11	90.98 ± 0.41	1.98 ± 0.37	1.54 ± 0.23	0.9
	Co-SLAM [3]	1.50 ± 0.03	2.06 ± 0.05	2.13 ± 0.05	93.38 ± 0.04	1.07 ± 0.08	0.93 ± 0.07	17.1
	ESLAM [2]	1.01 ± 0.03	2.24 ± 0.07	1.88 ± 0.08	95.31 ± 0.06	0.62 ± 0.05	0.51 ± 0.03	12.1
	Ours (Lite)	0.96 ± 0.01	1.99 ± 0.04	1.90 ± 0.05	94.96 ± 0.07	0.77 ± 0.05	0.67 ± 0.05	<b>21.4</b>
	Ours (Full)	<b>0.93 ± 0.01</b>	<b>1.83 ± 0.02</b>	<b>1.86 ± 0.05</b>	<b>95.79 ± 0.04</b>	<b>0.29 ± 0.02</b>	<b>0.27 ± 0.02</b>	11.9

Quantitative comparison of our system with other state-of-the-art NeRF-based RGB-D dense VSLAM systems on the scenes of the Replica dataset [7]. Numbers in **bold black font** represent the best results in each column. Our tracking results outperform those of other systems, and our mapping results reach state-of-the-art levels with the minor model parameters and highest FPS.

**Evaluation on TUM RGB-D [8].** Consistent with previous literature [1]–[3], we further tested our system on three scenes of the TUM dataset. Each result is the average of five random runs of the system. As shown in Tab. III, our system achieved state-of-the-art results in most scenes, further demonstrating our system’s cross-dataset robustness. Meanwhile, our speed is the fastest in NeRF based SLAM. Although SLAM systems based on NeRF have lower tracking accuracy compared to traditional SLAM systems, our system significantly narrows this gap with accurate and dense map, providing a robust

baseline for future work.

### C. Implicit loop closure capabilities of NeRF

Our observations reveal that NeRF intrinsically supports implicit loop closure correction. This capability is illustrated in the upper figure of Fig. 6(a), using the Replica dataset’s *Room0* scene, in the absence of feature-based and uniform sampling (FUS) or effective constraint through global bundle adjustment (EBA). Notably, at frames 270 and 1625, a significant reduction in error is observed, corresponding with the

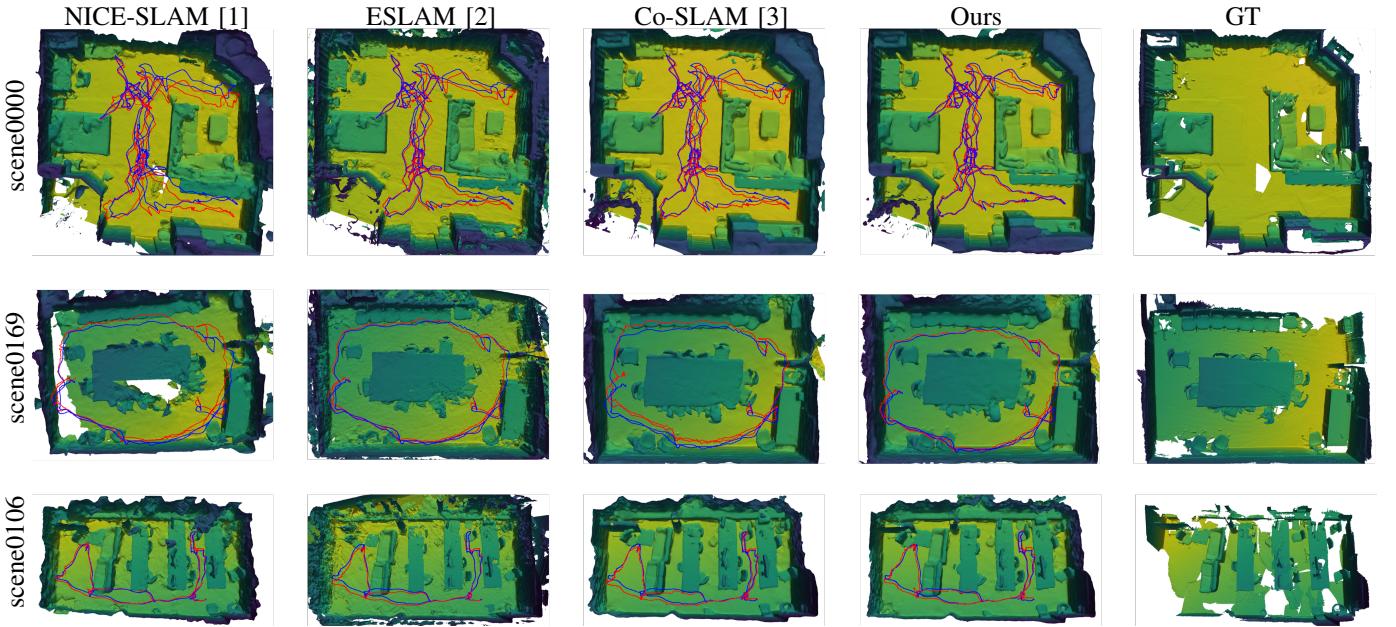


Fig. 4. Comparison of our system with other state-of-the-art NeRF-based RGB-D dense visual SLAM systems conducted on the ScanNet dataset [4], focusing on the reconstructed scene depth maps and tracking metrics. Blue lines represent the GT camera trajectory, and red lines indicate the estimated trajectory. Under the same camera perspectives, our system reconstructed overall scene structures more consistent with the GT, along with more accurate camera trajectories.

TABLE II  
TRACKING RESULTS ON THE SCANNET DATASET [4]

Methods	ATE	Sc. 0000	Sc. 0059	Sc. 0106	Sc. 0169	Sc. 0181	Sc. 0207	Avg.	FPS
iMAP [5]	Mean	$34.2 \pm 12.8$	$13.0 \pm 2.4$	$12.9 \pm 1.7$	$33.6 \pm 15.3$	$20.8 \pm 3.8$	$18.6 \pm 6.0$	$22.2 \pm 7.0$	
	RMSE	$42.7 \pm 16.6$	$17.8 \pm 7.4$	$15.0 \pm 1.7$	$39.1 \pm 18.2$	$24.7 \pm 5.8$	$20.1 \pm 6.8$	$26.6 \pm 9.4$	0.1
NICE-S [1]	Mean	$9.9 \pm 0.4$	$11.9 \pm 1.8$	$7.0 \pm 0.2$	$9.2 \pm 1.0$	$12.2 \pm 0.3$	$5.5 \pm 0.3$	$9.3 \pm 0.7$	
	RMSE	$12.0 \pm 0.5$	$14.0 \pm 1.8$	$7.9 \pm 0.2$	$10.9 \pm 1.1$	$13.4 \pm 0.3$	$6.2 \pm 0.4$	$10.7 \pm 0.7$	0.7
Co-S [3]	Mean	$6.0 \pm 0.1$	$9.4 \pm 0.8$	$8.4 \pm 0.1$	$5.1 \pm 0.1$	$10.3 \pm 0.3$	$6.6 \pm 0.4$	$7.6 \pm 0.3$	
	RMSE	$7.1 \pm 0.2$	$11.2 \pm 0.7$	$9.3 \pm 0.2$	$5.8 \pm 0.2$	$11.6 \pm 0.6$	$7.1 \pm 0.4$	$8.7 \pm 0.4$	6.4
E-S [2]	Mean	$6.5 \pm 0.1$	$6.4 \pm 0.4$	$6.7 \pm 0.1$	$5.9 \pm 0.1$	$8.3 \pm 0.2$	<b><math>5.4 \pm 0.1</math></b>	$6.5 \pm 0.2$	
	RMSE	$7.3 \pm 0.2$	$8.5 \pm 0.5$	$7.5 \pm 0.1$	$6.5 \pm 0.1$	$9.0 \pm 0.2$	<b><math>5.7 \pm 0.1</math></b>	$7.4 \pm 0.2$	4.1
Ours	Mean	<b><math>4.7 \pm 0.1</math></b>	<b><math>5.8 \pm 0.2</math></b>	<b><math>6.4 \pm 0.1</math></b>	<b><math>4.6 \pm 0.1</math></b>	<b><math>8.1 \pm 0.1</math></b>	$6.6 \pm 0.1$	<b><math>6.0 \pm 0.1</math></b>	
	RMSE	<b><math>5.3 \pm 0.2</math></b>	<b><math>8.2 \pm 0.4</math></b>	<b><math>7.3 \pm 0.1</math></b>	<b><math>5.1 \pm 0.1</math></b>	<b><math>8.9 \pm 0.2</math></b>	$7.0 \pm 0.1$	<b><math>6.9 \pm 0.2</math></b>	8.6

Quantitative comparison of our system against other SOTA NeRF-based RGB-D dense visual SLAM systems, using the ScanNet dataset [4]. We get the results of ESLAM [2], Nice-SLAM [1] and iMap [5] from ESLAM [2]. Numbers in **bold black font** represent the best results in each column. Our tracking results outperformed those of other systems in most scenes, with the fastest FPS.

system revisiting previously encountered locations (Fig. 6(b)). Our observations of implicit loop closure capabilities in NeRF are not isolated to a single scene or dataset. We consistently observe similar phenomena across multiple scenes in various datasets, which further corroborates NeRF’s proficiency in implicit loop closure. Upon integrating FUS along with EBA, there is a notable reduction in the overall pose error of the system. Additionally, this integration facilitated the identification of an extra loop closure event at frame 875, as depicted in the lower figure of Figure 6(a). This event serves as compelling evidence of NeRF’s implicit loop closure capabilities.

The observed capability for implicit loop closure in NeRF stems from its continuous and dense spatial representation, while our FUS sampling further enhances this capability. The optimized implicit mapping ability, to some extent, resembles

the biological perception pattern of the environment: the brain creates a dense map of the space while focusing on recording some meaningful objects (such as tables and chairs). When the eyes perceive a scene, the brain continuously compares it with existing scenes in the mind, relying on certain feature objects (attention-based) for comprehensive localization. Compared to the explicit recognition of loop closures in classic SLAM, NeRF’s implicit loop closure appears more natural because it can comprehensively consider both the continuity and specificity of space. It demonstrates the powerful spatial representation capabilities of NeRF.

#### D. Runtime Analysis

We evaluated the model size and FPS produced by all systems on the Replica [7], and ScanNet [4] datasets. For the

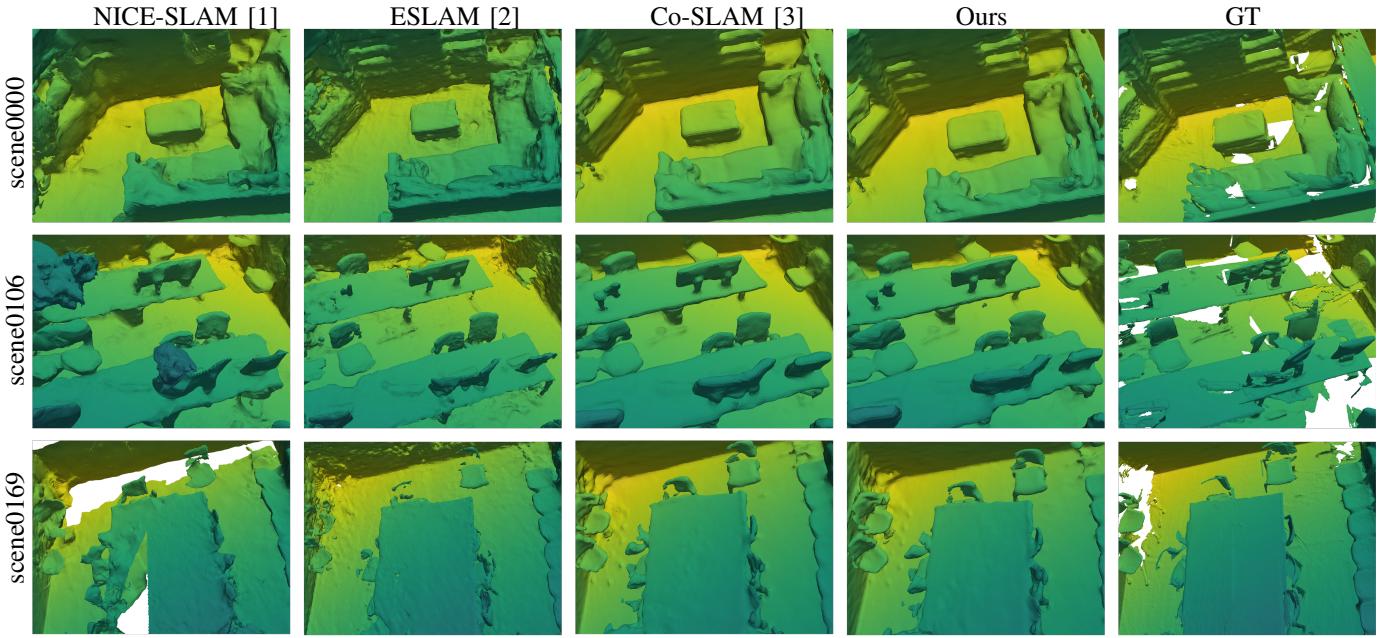


Fig. 5. Comparison of our system with other state-of-the-art NeRF-based RGB-D dense visual SLAM systems conducted on the ScanNet dataset [4] (scene0000, scene0106, scene0169), focusing on the reconstructed scene depth maps. Under the same camera perspectives, our system reconstructed scene details more consistent with the GT.

TABLE III  
TRACKING RESULTS ON THE TUM DATASET [8]

	Method	fr1/desk	fr2/xyz	fr3/office	FPS
Classic	Kintinous [31]	3.70	2.90	3.00	<b>40</b>
	BAD-SLAM [61]	1.70	1.10	1.70	2
	ORB-SLAM2 [17]	<b>1.60</b>	<b>0.40</b>	<b>1.00</b>	30
NeRF	iMAP [5]	4.90	2.00	5.80	0.1
	NICE-SLAM [1]	2.85	1.84	2.95	0.1
	Co-SLAM [3]	2.44	1.71	2.46	6.9
	E-SLAM [2]	2.54	<b>1.09</b>	2.47	0.2
	Ours	<b>2.32</b>	1.90	<b>2.41</b>	<b>7.1</b>

Numbers in **bold black font** represent the best results in each column. Our method surpasses the accuracy of other NeRF-based methods and narrows the gap with traditional approaches.

TABLE IV  
TIME AND MEMORY ANALYSIS.

	Method	FPS	Memory Param.
Replica	iMAP [5]	0.4	<b>0.20 M</b>
	NICE-SLAM [1]	0.9	17.40 M
	Co-SLAM [3]	17.1	0.26 M
	ESLAM [2]	12.1	9.29 M
	Ours (Lite)	<b>21.4</b>	0.26 M
	Ours (Full)	11.9	0.26 M
ScanNet	iMAP [5]	0.1	<b>0.2 M</b>
	NICE-SLAM [1]	0.7	10.3 M
	Co-SLAM [3]	6.4	0.8 M
	ESLAM [2]	4.1	10.5 M
	Ours	<b>8.6</b>	0.8 M

Comparison of the model size and runtime for maps generated by different systems on the Replica [7] and ScanNet [4] datasets. Numbers in **bold black font** represent the best results in each column.

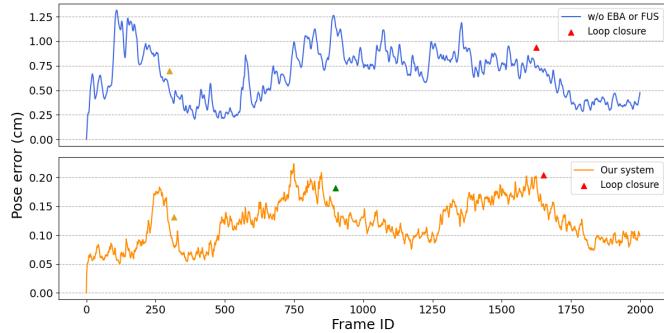
Replica dataset [7], we tested the average parameter count and FPS for eight scenes, and for the ScanNet dataset [4], we assessed the parameter count and runtime for scene0000. The details are provided in Tab. IV. All results were obtained on the same computer configuration. The results indicate that our system operates faster than others, with a small model size as Co-SLAM [3], attributed to using a hash sparse grid representation of the map and our one-time sampling iteration optimization strategy.

### E. Ablations

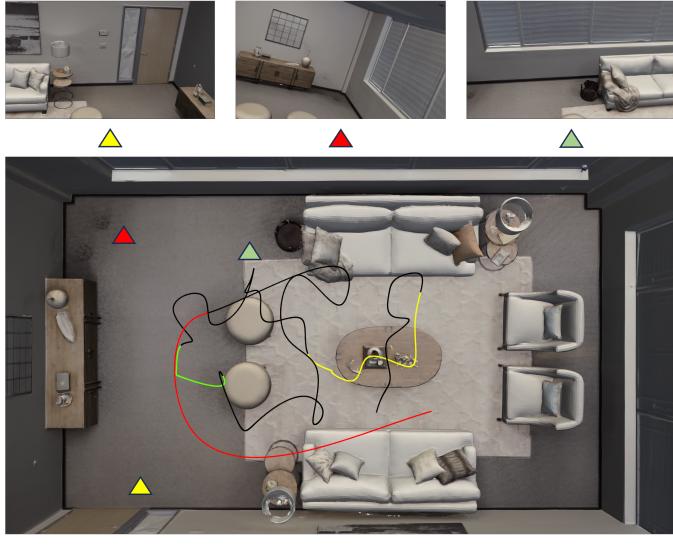
We conducted ablation experiments on our effectively constrained global BA strategy and feature-based and uniform sampling, detailed in Tab. V. There are four comparative experiments: one without the effectively constrained global BA strategy, one without feature-based and uniform sampling, one with neither, and one with both. Without the effectively constrained global BA strategy, we randomly selected the number and order of keyframes in the sliding window and the number of sampled pixels per keyframe. Without feature-based and uniform sampling, we randomly sampled pixels. The rest of the parameters remained the same in all experiments except for the listed variable differences. The results show that both optimization strategies significantly reduce our tracking error and improve our mapping accuracy. It is noted that FUS contributes to a more significant improvement in system accuracy compared to EBA, as it markedly reduces randomness and fully leverages the implicit loop closure capabilities of NeRF.

## V. CONCLUSION

We present EC-SLAM, a real-time, NeRF-based RGB-D SLAM system that delivers exceptional accuracy and efficiency. Our innovations include an effectively constrained



(a) The frame-wise quaternion discrepancies.



(b) Detailed loop closure location

Fig. 6. The implicit loop closure capabilities of NeRF on the Replica dataset: (a) The upper figure shows the results of frame-by-frame pose quaternion differences without FUS and EBA, while the lower figure presents the results with FUS and EBA. Different colored solid triangles represent the locations of implicit loop closures. (b) Shows the actual areas of loop closures corresponding to figure (a), with colored solid triangles matching the three loop closure colors in figure (a). The black trajectory corresponds to the actual trajectory, with other colored trajectories indicating the locations of the three loop closures identified in figure (a).

TABLE V  
ABLATION STUDY.

Name	FUS	EBA	RMSE.(cm)	Mean.(cm)
w/o both			$6.9 \pm 0.2$	$6.1 \pm 0.1$
w/o EBA	✓		$6.4 \pm 0.2$	$5.5 \pm 0.2$
w/o FUS		✓	$5.7 \pm 0.2$	$5.1 \pm 0.2$
Full model	✓	✓	<b><math>5.3 \pm 0.2</math></b>	<b><math>4.7 \pm 0.1</math></b>

Ablation studies on our system on scene0000 of the ScanNet dataset [4]: EBA represents our effectively constrained global Bundle Adjustment strategy; FUS represents our feature-based and uniform sampling. Numbers in **bold black font** represent the best results in each column.

global bundle adjustment strategy, feature-based and uniform sampling, sparse positional encoding, and a TSDF-based loss function. These enhancements fully exploit NeRF's implicit loop closure capabilities, resulting in precise camera trajectory tracking and map reconstruction with minimal model parameters. Extensive experiments on multiple datasets demonstrate

EC-SLAM's state-of-the-art performance and impressive 21 Hz runtime.

While EC-SLAM offers significant advancements, opportunities for improvement remain. Our future work concentrates on extending the capabilities of EC-SLAM. We will explore multi-sensor integration to boost robustness in diverse environments, investigate combining classic SLAM algorithms with NeRF-based methods for enhanced performance, and work towards adapting EC-SLAM to tackle the complexities of large-scale outdoor exploration. These advancements will drive the development of even more powerful and versatile SLAM systems.

## REFERENCES

- [1] Z. Zhu, S. Peng, V. Larsson, W. Xu, H. Bao, Z. Cui, M. R. Oswald, and M. Pollefeys, "Nice-slam: Neural implicit scalable encoding for slam," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 12 786–12 796.
- [2] M. M. Johari, C. Carta, and F. Fleuret, "Eslam: Efficient dense slam system based on hybrid representation of signed distance fields," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023, pp. 17 408–17 419.
- [3] H. Wang, J. Wang, and L. Agapito, "Co-slam: Joint coordinate and sparse parametric encodings for neural real-time slam," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023, pp. 13 293–13 302.
- [4] A. Dai, A. X. Chang, M. Savva, M. Halber, T. Funkhouser, and M. Nießner, "Scannet: Richly-annotated 3d reconstructions of indoor scenes," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 5828–5839.
- [5] E. Sucar, S. Liu, J. Ortiz, and A. J. Davison, "imap: Implicit mapping and positioning in real-time," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 6229–6238.
- [6] X. Kong, S. Liu, M. Taher, and A. J. Davison, "vmap: Vectorised object mapping for neural field slam," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023, pp. 952–961.
- [7] J. Straub, T. Whelan, L. Ma, Y. Chen, E. Wijmans, S. Green, J. J. Engel, R. Mur-Artal, C. Ren, S. Verma *et al.*, "The replica dataset: A digital replica of indoor spaces," *ArXiv Preprint ArXiv:1906.05797*, 2019.
- [8] J. Sturm, N. Engelhard, F. Endres, W. Burgard, and D. Cremers, "A benchmark for the evaluation of rgbd slam systems," in *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2012, pp. 573–580.
- [9] A. J. Davison, I. D. Reid, N. D. Molton, and O. Stasse, "Monoslam: Real-time single camera slam," *IEEE transactions on pattern analysis and machine intelligence*, vol. 29, no. 6, pp. 1052–1067, 2007.
- [10] B. Vincke, A. Elouardi, and A. Lambert, "Design and evaluation of an embedded system based slam applications," in *IEEE/SICE International Symposium on System Integration*, 2010, pp. 224–229.
- [11] B. Vincke, A. Elouardi, A. Lambert, and A. Merigot, "Efficient implementation of ekf-slam on a multi-core embedded system," in *IECON 2012-38th Annual Conference on IEEE Industrial Electronics Society*, 2012, pp. 3049–3054.
- [12] G. Klein and D. Murray, "Parallel tracking and mapping for small ar workspaces," in *2007 6th IEEE and ACM International Symposium on Mixed and Augmented Reality*, 2007, pp. 225–234.
- [13] A. A. Jiménez Serrata, S. Yang, and R. Li, "An intelligible implementation of fastslam2. 0 on a low-power embedded architecture," *EURASIP Journal on Embedded Systems*, vol. 2017, pp. 1–11, 2017.
- [14] C. Forster, M. Pizzoli, and D. Scaramuzza, "Svo: Fast semi-direct monocular visual odometry," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2014, pp. 15–22.
- [15] C. Forster, Z. Zhang, M. Gassner, M. Werlberger, and D. Scaramuzza, "Svo: Semidirect visual odometry for monocular and multicamera systems," *IEEE Transactions on Robotics*, vol. 33, no. 2, pp. 249–265, 2016.
- [16] R. Mur-Artal, J. M. M. Montiel, and J. D. Tardos, "Orb-slam: a versatile and accurate monocular slam system," *IEEE Transactions on Robotics*, vol. 31, no. 5, pp. 1147–1163, 2015.
- [17] R. Mur-Artal and J. D. Tardós, "Orb-slam2: An open-source slam system for monocular, stereo, and rgbd cameras," *IEEE Transactions on Robotics*, vol. 33, no. 5, pp. 1255–1262, 2017.

- [18] C. Campos, R. Elvira, J. J. G. Rodríguez, J. M. Montiel, and J. D. Tardós, “Orb-slam3: An accurate open-source library for visual, visual-inertial, and multimap slam,” *IEEE Transactions on Robotics*, vol. 37, no. 6, pp. 1874–1890, 2021.
- [19] J. Engel, V. Koltun, and D. Cremers, “Direct sparse odometry,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 3, pp. 611–625, 2017.
- [20] A. I. Mourikis and S. I. Roumeliotis, “A multi-state constraint kalman filter for vision-aided inertial navigation,” in *Proceedings 2007 IEEE International Conference on Robotics and Automation*, 2007, pp. 3565–3572.
- [21] S. Leutenegger, S. Lynen, M. Bosse, R. Siegwart, and P. Furgale, “Keyframe-based visual-inertial odometry using nonlinear optimization,” *The International Journal of Robotics Research*, vol. 34, no. 3, pp. 314–334, 2015.
- [22] M. Bloesch, S. Omari, M. Hutter, and R. Siegwart, “Robust visual inertial odometry using a direct ekf-based approach,” in *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2015, pp. 298–304.
- [23] R. Mur-Artal and J. D. Tardós, “Visual-inertial monocular slam with map reuse,” *IEEE Robotics and Automation Letters*, vol. 2, no. 2, pp. 796–803, 2017.
- [24] T. Qin, P. Li, and S. Shen, “Vins-mono: A robust and versatile monocular visual-inertial state estimator,” *IEEE Transactions on Robotics*, vol. 34, no. 4, pp. 1004–1020, 2018.
- [25] L. Von Stumberg, V. Usenko, and D. Cremers, “Direct sparse visual-inertial odometry using dynamic marginalization,” in *2018 IEEE International Conference on Robotics and Automation (ICRA)*, 2018, pp. 2510–2517.
- [26] R. A. Newcombe, S. J. Lovegrove, and A. J. Davison, “Dtam: Dense tracking and mapping in real-time,” in *2011 International Conference On Computer Vision*. IEEE, 2011, pp. 2320–2327.
- [27] P. Ondruška, P. Kohli, and S. Izadi, “Mobilefusion: Real-time volumetric surface reconstruction and dense tracking on mobile phones,” *IEEE Transactions On Visualization and Computer Graphics*, vol. 21, no. 11, pp. 1251–1258, 2015.
- [28] J. Engel, T. Schöps, and D. Cremers, “Lsd-slam: Large-scale direct monocular slam,” in *European Conference on Computer Vision*. Springer, 2014, pp. 834–849.
- [29] K. Boikos and C.-S. Bouganis, “Semi-dense slam on an fpga soc,” in *2016 26th International Conference on Field Programmable Logic and Applications (FPL)*. IEEE, 2016, pp. 1–4.
- [30] ———, “A high-performance system-on-chip architecture for direct tracking for slam,” in *2017 27th International Conference on Field Programmable Logic and Applications (FPL)*, 2017, pp. 1–7.
- [31] R. A. Newcombe, S. Izadi, O. Hilliges, D. Molyneaux, D. Kim, A. J. Davison, P. Kohli, J. Shotton, S. Hodges, and A. Fitzgibbon, “Kinectfusion: Real-time dense surface mapping and tracking,” in *2011 10th IEEE International Symposium on Mixed and Augmented Reality*, 2011, pp. 127–136.
- [32] R. F. Salas-Moreno, R. A. Newcombe, H. Strasdat, P. H. Kelly, and A. J. Davison, “Slam++: Simultaneous localisation and mapping at the level of objects,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2013, pp. 1352–1359.
- [33] C. Kerl, J. Sturm, and D. Cremers, “Dense visual slam for rgb-d cameras,” in *2013 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2013, pp. 2100–2106.
- [34] F. Endres, J. Hess, J. Sturm, D. Cremers, and W. Burgard, “3-d mapping with an rgb-d camera,” *IEEE Transactions on Robotics*, vol. 30, no. 1, pp. 177–187, 2013.
- [35] K. Tateno, F. Tombari, I. Laina, and N. Navab, “Cnn-slam: Real-time dense monocular slam with learned depth prediction,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 6243–6252.
- [36] Z. Teed and J. Deng, “Droid-slam: Deep visual slam for monocular, stereo, and rgb-d cameras,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 16 558–16 569, 2021.
- [37] N. Yang, L. v. Stumberg, R. Wang, and D. Cremers, “D3vo: Deep depth, deep pose and deep uncertainty for monocular visual odometry,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 1281–1292.
- [38] R. Li, S. Wang, Z. Long, and D. Gu, “Undeepvo: Monocular visual odometry through unsupervised deep learning,” in *IEEE International Conference on Robotics and Automation (ICRA)*, 2018, pp. 7286–7291.
- [39] R. Li, S. Wang, and D. Gu, “Deepslam: A robust monocular slam system with unsupervised deep learning,” *IEEE Transactions on Industrial Electronics*, vol. 68, no. 4, pp. 3577–3587, 2020.
- [40] R. Kang, J. Shi, X. Li, Y. Liu, and X. Liu, “Df-slam: A deep-learning enhanced visual slam system based on deep local features,” *ArXiv Preprint ArXiv:1901.07223*, 2019.
- [41] B. Ummenhofer, H. Zhou, J. Uhrig, N. Mayer, E. Ilg, A. Dosovitskiy, and T. Brox, “Demon: Depth and motion network for learning monocular stereo,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 5038–5047.
- [42] H. Zhou, B. Ummenhofer, and T. Brox, “Deeptam: Deep tracking and mapping,” in *European Conference on Computer Vision*, 2018, pp. 822–838.
- [43] M. Bloesch, J. Czarnowski, R. Clark, S. Leutenegger, and A. J. Davison, “Codeslam—learning a compact, optimisable representation for dense visual slam,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 2560–2568.
- [44] S. Zhi, M. Bloesch, S. Leutenegger, and A. J. Davison, “Scenecode: Monocular dense semantic reconstruction using learned encoded scene representations,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 11 776–11 785.
- [45] E. Sucar, K. Wada, and A. Davison, “Nodeslam: Neural object descriptors for multi-view shape reconstruction,” in *International Conference on 3D Vision (3DV)*, 2020, pp. 949–958.
- [46] Z. Teed and J. Deng, “Deepv2d: Video to depth with differentiable structure from motion,” in *8th International Conference on Learning Representations*, 2020.
- [47] C. Tang and P. Tan, “Ba-net: Dense bundle adjustment network,” *ArXiv Preprint ArXiv:1806.04807*, 2018.
- [48] J. Czarnowski, T. Laidlow, R. Clark, and A. J. Davison, “Deepfactors: Real-time probabilistic dense monocular slam,” *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 721–728, 2020.
- [49] B. Mildenhall, P. Srinivasan, M. Tancik, J. Barron, R. Ramamoorthi, and R. Ng, “Nerf: Representing scenes as neural radiance fields for view synthesis,” in *European conference on computer vision*, 2020.
- [50] C.-M. Chung, Y.-C. Tseng, Y.-C. Hsu, X.-Q. Shi, Y.-H. Hua, J.-F. Yeh, W.-C. Chen, Y.-T. Chen, and W. H. Hsu, “Orbeez-slam: A real-time monocular visual slam with orb features and nerf-realized mapping,” in *2023 IEEE International Conference on Robotics and Automation (ICRA)*, 2023, pp. 9400–9406.
- [51] A. Rosinol, J. J. Leonard, and L. Carlone, “Nerf-slam: Real-time dense monocular slam with neural radiance fields,” in *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2023, pp. 3437–3444.
- [52] Y. Zhang, F. Tosi, S. Mattoccia, and M. Poggi, “Go-slam: Global optimization for consistent 3d instant reconstruction,” in *IEEE/CVF International Conference on Computer Vision*, 2023, pp. 3727–3737.
- [53] X. Yang, H. Li, H. Zhai, Y. Ming, Y. Liu, and G. Zhang, “Vox-fusion: Dense tracking and mapping with voxel-based neural implicit representation,” in *IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*, 2022, pp. 499–507.
- [54] Z. Zhu, S. Peng, V. Larsson, Z. Cui, M. R. Oswald, A. Geiger, and M. Pollefeys, “Nicer-slam: Neural implicit scene encoding for rgb slam,” *ArXiv Preprint ArXiv:2302.03594*, 2023.
- [55] H. Li, X. Gu, W. Yuan, L. Yang, Z. Dong, and P. Tan, “Dense rgb slam with neural implicit maps,” *ArXiv Preprint ArXiv:2301.08930*, 2023.
- [56] T. Müller, A. Evans, C. Schied, and A. Keller, “Instant neural graphics primitives with a multiresolution hash encoding,” *ACM Transactions on Graphics (ToG)*, vol. 41, no. 4, pp. 1–15, 2022.
- [57] E. R. Chan, C. Z. Lin, M. A. Chan, K. Nagano, B. Pan, S. De Mello, O. Gallo, L. J. Guibas, J. Tremblay, S. Khamis *et al.*, “Efficient geometry-aware 3d generative adversarial networks,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 16 123–16 133.
- [58] J. Shi *et al.*, “Good features to track,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 1994, pp. 593–600.
- [59] T. Müller, B. McWilliams, F. Rousselle, M. Gross, and J. Novák, “Neural importance sampling,” *ACM Transactions on Graphics (ToG)*, vol. 38, no. 5, pp. 1–19, 2019.
- [60] R. Or-El, X. Luo, M. Shan, E. Shechtman, J. J. Park, and I. Kemelmacher-Shlizerman, “Stylesdf: High-resolution 3d-consistent image and geometry generation,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 13 503–13 513.
- [61] T. Schops, T. Sattler, and M. Pollefeys, “Bad slam: Bundle adjusted direct rgb-d slam,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 134–144.



**Guanghao Li** received the B.S. degree from Xiamen University. He is currently pursuing his Ph.D. at the Institute of Science and Technology for Brain-Inspired Intelligence (ISTBI), Fudan University. His research areas are focused on Visual SLAM and Neural Radiance Fields (NeRF).



**Qi Chen** received a B.S. degree from the Beijing Institute of Technology in 2016 and an M.S. degree from New York University in 2018. From 2019 to 2021, he was an algorithm engineer at the Institute of Acoustics of the Chinese Academy of Sciences. He is currently a Ph.D. student in computer science at Fudan University, Shanghai, China. His research interests include computer vision, digital signal processing, and machine learning.



**Yuxiang Yan** received a B.S. degree from East China University of Science and Technology. He is currently pursuing his Ph.D. at the Institute of Science and Technology for Brain-Inspired Intelligence, Fudan University, Shanghai, China. His current research areas are focused on computer vision and machine learning for autonomous driving.



**Jian Pu** received a Ph.D. degree from Fudan University, Shanghai, China, in 2014. Currently, he is a Young Principal Investigator at the Institute of Science and Technology for Brain-Inspired Intelligence (ISTBI), Fudan University. He was an associate professor at the School of Computer Science and Software Engineering, East China Normal University from 2016 to 2019, and a postdoctoral researcher of the Institute of Neuroscience, Chinese Academy of Sciences in China from 2014 to 2016. His current research interest is to develop machine learning and computer vision methods for autonomous driving.