

OG-Mapping: Octree-based Structured 3D Gaussians for Online Dense Mapping

Meng Wang^{1,2}, Junyi Wang³, Changqun Xia², Chen Wang⁴, Yue Qi¹

¹State Key Laboratory of Virtual Reality Technology and Systems, Beihang University.

²PengCheng Laboratory. ³School of Computer Science and Technology, Shandong University.

⁴Beijing Technology and Business University.

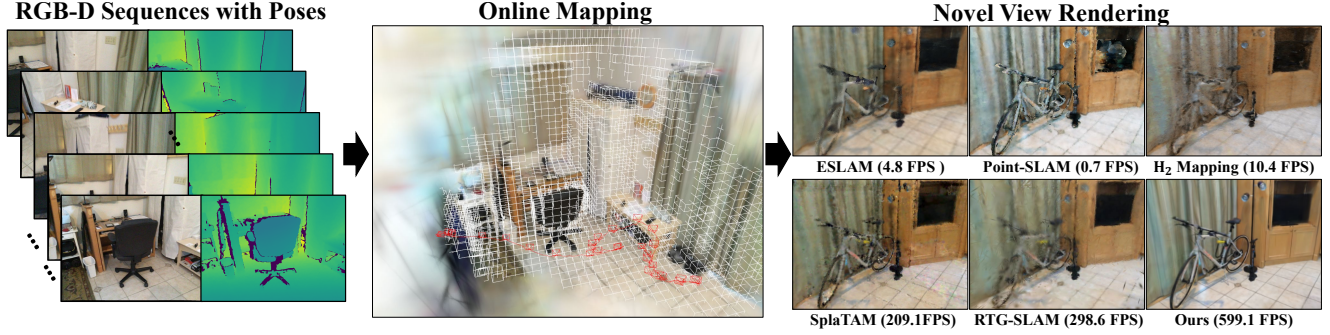


Figure 1. In this work, we introduce OG-Mapping, a novel online dense mapping framework with an octree-based structured 3D Gaussian representation. By combining our proposed anchor-based progressive map refinement strategy and dynamic keyframe window, OG-Mapping achieves fast, high-fidelity online reconstruction with efficient memory usage, and demonstrates superior realism in novel view synthesis compared to other existing RGB-D online mapping methods. The rendering FPS is indicated to the right of each method.

Abstract

3D Gaussian splatting (3DGS) has recently demonstrated promising advancements in RGB-D online dense mapping. Nevertheless, existing methods excessively rely on per-pixel depth cues to perform map densification, which leads to significant redundancy and increased sensitivity to depth noise. Additionally, explicitly storing 3D Gaussian parameters of room-scale scene poses a significant storage challenge. In this paper, we introduce OG-Mapping, which leverages the robust scene structural representation capability of sparse octrees, combined with structured 3D Gaussian representations, to achieve efficient and robust online dense mapping. Moreover, OG-Mapping employs an anchor-based progressive map refinement strategy to recover the scene structures at multiple levels of detail. Instead of maintaining a small number of active keyframes with a fixed keyframe window as previous approaches do, a dynamic keyframe window is employed to allow OG-Mapping to better tackle false local minima and forgetting issues. Experimental results demonstrate that OG-Mapping delivers more robust and superior realism mapping results than existing Gaussian-based RGB-D online mapping methods with a compact model, and no additional post-processing is required.

1. Introduction

Constructing highly detailed dense maps in real-time holds significant importance in AR/VR, robotics, and digi-

tal twins applications. For the past several decades, research on mapping has extensively centered around the scene representation, resulting in various representations such as occupancy grids [9], TSDF [3, 44], surfels [1, 37] and point clouds [7]. Although systems utilizing these map representations have reached a production-ready standard over the past years, there remain notable deficiencies that demand attention and resolution. These methods exhibit deficiencies in rendering realism of reconstructed results, and fine-detailed maps based on the above representations require massive amounts of storage space.

In recent years, implicit representations [2, 26, 35] have demonstrated promising outcomes in various domains following the advent of Neural Radiance Fields (NeRF) [25]. Numerous studies [6, 12, 14, 30, 34, 46] employ these implicit representations to enhance mapping methodologies and exhibit strengths in generating high-quality dense maps with low memory consumption. Nevertheless, volumetric sampling significantly limits these methods' efficiency. Consequently, they opt for optimization over a sparse set of pixels instead of dense per-pixel photometric error, resulting in the reconstruction dense maps lack the richness and intricacy of texture details.

More recently, several works [23, 41] based on the 3D Gaussian representation [16] and tile-based splatting technique have shown great superiority in the efficiency of differentiable rendering. However, these methods are meticu-

lously designed for offline reconstruction scenarios. Concurrent works [11, 15, 24, 27, 38, 42] attempt to incorporate 3D Gaussians representation into RGB-D online dense mapping to deliver high-fidelity rendering, overcoming the limitation of Nerf-based methods. While promising results are demonstrated, these methods neglect the scene structure and rely excessively on pixel-level depth information to expand 3D Gaussians, resulting in significant redundancy and being sensitive to depth noise effects.

To overcome these challenges, in this paper, we introduce an innovative framework named OG-Mapping to perform efficient and highly detailed online dense mapping. Our solution comprises three main components. First, we incrementally build a sparse voxel octree with Morton coding [33] for fast allocation and retrieval of anchors to dynamically expand the map. Each anchor tethers a set of 3D Gaussians with learnable offsets. The attributes (color, opacity, quaternion, and scale) of these Gaussians are predicted based on the viewing direction and anchor feature encodings. By leveraging this structured representation, we can effectively mitigate pixel-level depth noise effects and avoid the substantial memory consumption associated with explicitly storing large amounts of 3D Gaussian attributes. Second, as the anchors organized by the sparse octree can only provide a rough description of the scene structure, we further design an anchor-based progressive map refinement strategy to recover the scene structures at different levels of detail by adaptively adding finer-level anchors. Unlike the previous method [23] utilizes an error-based policy, when an area is marked as under-optimized, we grow anchors according to the level hierarchy of anchors already assigned to that area, resulting in a more compact map. Finally, we develop a dynamic keyframe window to alleviate the false local minima and forgetting issue to improve mapping quality. Through extensive experiments, we empirically demonstrate that our method achieves superior mapping results with approximately a 5dB enhancement in PSNR on real-world scenes [4, 27], with a more compact model, while maintaining a more compact model and operating at high processing speed.

To summarize, our contributions are as follows:

- We introduce a novel framework that naturally integrates the sparse octree and structured 3D Gaussian representation to perform efficient and detailed online dense mapping.
- We design an anchor-based progressive map refinement strategy for better scene coverage and mapping quality.
- We develop a dynamic keyframe window to mitigate the false local minima and forgetting issue.

2. Related Work

Classical RGB-D online dense mapping. For online 3D reconstruction of scenes, various explicit representations

have been employed to store scene information, including point clouds [7], surface elements (surfels) [1, 37] and truncated signed distance functions (TSDF) [3, 44]. Several works [17, 31] leverage deep learning to improve the accuracy and robustness of above mentioned representations, and even achieve dense reconstruction with monocular input [32, 39]. These methods are renowned for their rapid processing speed, which is attributable to the inherent physical properties of explicit representations. However, they demand significant memory resources to manage high-detail mapping [45], and are incapable of realistically rendering from novel views.

NeRF-based RGB-D online dense mapping. Following the significant success of neural radiance fields (NeRF) [25], several studies leveraged latent features and neural networks as implicit representations to integrated NeRF with RGB-D dense mapping. iMAP [30] presents the first NeRF-style online dense mapping, using Multi-Layer Perceptrons (MLP) as the scene representation. NICE-SLAM [46] represents scenes as hierarchical feature grids, utilizing pre-trained MLPs for decoding. To enhance mapping speed and expand representational capacity, various representations have been investigated, including multi-resolution hash grids [12, 13, 34], factored grids [14], sparse octree grids [40] and neural point clouds [28]. Nevertheless, the aforementioned methods struggle to achieve fine-detailed rendering results and maintain fast rendering speed. These limitations arise from their reliance on time-consuming volumetric rendering techniques. In contrast, our approach leverages fast rasterization, enabling complete use of per-pixel dense photometric errors.

Gaussian-based RGB-D online dense mapping. The high-fidelity and rapid rasterization capabilities of 3D Gaussian Splatting (3DGS) [16] facilitate superior quality and efficiency in scene reconstruction. Recently, many works [11, 15, 24, 42] have attempted to apply 3DGS in online dense mapping. SplatAM [15] adopts an explicit volumetric approach using isotropic Gaussians, enabling precise map densification. However, this method necessitates projecting and densifying each pixel in the depth image to perform gaussian densification, resulting in substantial map storage usage. MonoGS [24] randomly select a subset of pixels for projection, which necessitates more optimization time for map densification. Concurrent to our work, CG-SLAM [10] introduces a depth uncertainty model to select more valuable Gaussian primitives during optimization. RTG-SLAM [27] treats each opaque Gaussian as an ellipsoid disc on the dominant plane of Gaussian to maintain a compact representation. NGM-SLAM [18] and Gaussian-SLAM [42] focus on building submap for 3D Gaussian representation. More recently MG-SLAM [21] leverages Manhattan World hypothesis and additional semantic informations to refine and complete scene geome-

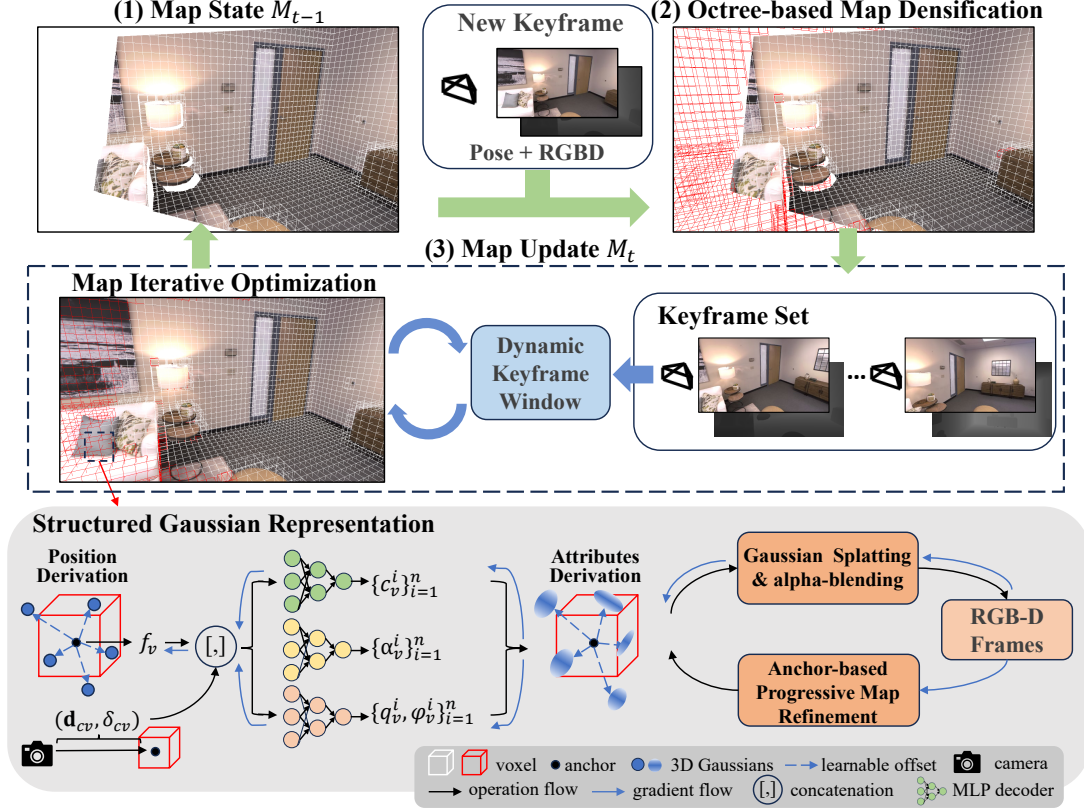


Figure 2. **Overview of OG-Mapping.** Given a set of sequential RGB-D frames and camera poses, we utilize an octree-based structured 3D Gaussians as the scene representation to perform efficient online dense mapping. When a new keyframe is detected, we employ a sparse octree to swiftly capture the rough structure of the new observed region to guide anchor densification (Sec. 3.1). During the map update process, we perform anchor-based progressive map refinement to enhance the geometry and appearance quality (Sec. 3.2), and construct a dynamic keyframe window to effectively mitigate false local minima and forgetting issues (Sec. 3.3).

tries. Different from these methods, our approach is based on the structured 3D Gaussian representation [23] and utilizes the structured information of the octree to guide the distribution of Gaussian kernels, resulting in better robustness to depth noise and smaller map size.

3. Method

The overview of OG-Mapping is shown in Fig. 2. Taking RGB-D images from sensors and poses from other tracking modules, we utilize a structured 3D Gaussians representation [23] managed by a sparse octree (Sec. 3.1) to depict the scene geometry and appearance. During mapping process, we employ an anchor-based progressive map refinement strategy to enhance map reconstruction quality (Sec. 3.2). In Sec. 3.3, we introduce how to build our dynamic keyframe window to mitigate the false local minima and forgetting problem. Sec. 3.4 elaborate the online optimization details.

3.1. Structured Gaussian Representation for Online Mapping

Scene Representation. We represent the underlying map of the scene as a set of anchors [23]. Specifically, we first vox-

elize the scene using the provided camera pose and depth image. For each voxel, the center v is treated as an anchor point, equipped with a level mark $l_v \in \mathbb{N}_0$, a scaling factor $s_v \in \mathbb{R}^3$, n learnable offsets $\{o_v^i\}_{i=1}^n \in \mathbb{R}^{n \times 3}$ and a feature vector $f_v = \text{encoding}(v)$ (In our dense version, the encoding function is the multi-hash encoding [26], while in the sparse version, it represents local context encoding. Implementation details are in supplementary).

For each visible anchor within the viewing frustum, n 3D Gaussians are generated. The positions $\{\mu_v^i\}_{i=1}^n$ of these 3D Gaussians are calculated as:

$$\{\mu_v^i\}_{i=1}^n = p_v + s_v \cdot \{o_v^i\}_{i=1}^n \quad (1)$$

The other attributes (color $c_v^i \in \mathbb{R}^3$, opacity $\alpha_v^i \in \mathbb{R}_{>0}$, quaternions $q_v^i \in \mathbb{R}^4$, and scale $\varphi_v^i \in \mathbb{R}^3$) of n 3D Gaussians are decode from the relative distance δ_{cv} , view direction \mathbf{d}_{cv} , and the anchor feature f_v using individual multilayer perceptron (MLP) decoders, denoted as F_{color} , $F_{opacity}$ and F_{cov} :

$$\begin{aligned} \{c_v^i\}_{i=1}^n &= \mathbf{F}_{color}(\delta_{cv}, \mathbf{d}_{cv}, f_v), \\ \{\alpha_v^i\}_{i=1}^n &= \mathbf{F}_{opacity}(\delta_{cv}, \mathbf{d}_{cv}, f_v), \\ \{q_v^i, \varphi_v^i\}_{i=1}^n &= \mathbf{F}_{cov}(\delta_{cv}, \mathbf{d}_{cv}, f_v), \end{aligned} \quad (2)$$

where the relative distance δ_{cv} and the viewing direction \mathbf{d}_{cv} between camera position p_c and anchor point position p_v are calculated as follows :

$$\delta_{cv} = \|p_c - p_v\|_2, \quad \mathbf{d}_{cv} = \frac{p_c - p_v}{\delta_{cv}}. \quad (3)$$

Then, these generated \mathcal{N} 3D Gaussians are used for fast rasterization rendering to produce color and depth maps. Given a viewpoint, the rendered color of each pixel \mathbf{p} can be written as:

$$C(\mathbf{p}) = \sum_{i \in \mathcal{N}} c^i \sigma^i \prod_{j=1}^{i-1} (1 - \sigma^j), \sigma^i = \alpha^i G_{2d}^i(\mathbf{p}), \quad (4)$$

where the 2D Gaussians $G_{2d}(\mathbf{p})$ are transformed from 3D Gaussian $G(\mathbf{p})$ introduced by [16]. Similarly, per-pixel depth is rendered via alpha-blending:

$$D(\mathbf{p}) = \sum_{i \in \mathcal{N}} z^i \sigma^i \prod_{j=1}^{i-1} (1 - \sigma^j), \quad (5)$$

where z^i is the distance of the i^{th} 3D Gaussian's position u^i along the camera ray.

Map Densification. To achieve a comprehensive representation of the environment, we need to add new anchors to the scene during online scanning to cover newly observed regions. The adaptive map densification approach in [15, 27], which is based on pixel-level projection error, is sensitive to depth noise and is unreliable when handling edges. Upon receiving a new keyframe, we dynamically allocate new voxels using the provided pose and depth image, incrementally updating a sparse octree to roughly encompass all visible regions.

3.2. Anchor-based Progressive Map Refinement

A progressive optimization strategy can better shape the loss landscape, reducing the risk of the algorithm becoming trapped in misleading local minima. Such a strategy has been successfully applied in various computer vision applications like registration [20] and surface reconstruction [19]. OG-Mapping also utilizes a coarse-to-fine optimization scheme to reconstruct the scene with progressive levels of detail, by adaptively adding finer-level anchors to under-optimized regions.

To ensure system efficiency, the sparse octree only maintains the coarsest granularity of anchors to quickly reveal scene structure information. The 3D Gaussians inferred from these newly added anchors have a large scale, which

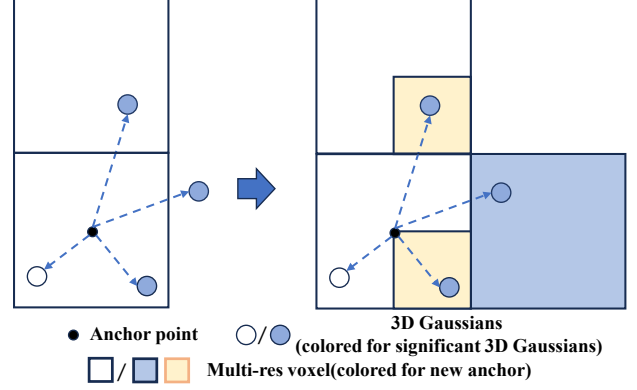


Figure 3. **Anchor-based Progressive Map Refinement.** We grow new anchors in under-optimized regions based on the level of 3D Gaussians and their gradients. If the regions already contain anchors at the same level, the granularity of these new anchors will be increased.

makes it difficult to fit areas with high-frequency texture changes. Therefore, anchors with finer granularity need to be added to these under-optimized areas. Scaffold-GS [23] employs an error-based anchor growth method, which allows the 3D Gaussians of fine-level anchors to participate in the growth of coarser ones. Nevertheless, during online reconstruction, it fails to provide sufficient optimization iterations to stabilize the gradients of such fine-level 3D Gaussians, resulting in the generation of numerous unnecessary anchors.

To address the aforementioned issues, we develop an anchor-based progressive map refinement strategy as illustrated in Fig. 3. Specifically, for each level l , we predefine the corresponding voxel size γ_l and significance threshold τ_l . To determine under-optimized areas, the gradients of 3D Gaussians belonging to anchor v are denoted as $\{\nabla_v^i\}_{i=1}^n$. Subsequently, the 3D Gaussian satisfying $\nabla_v^i \geq \tau_l$ marks its respective region as candidate region. A new anchor v' with level $l_{v'} = l_v$ is deployed at the position $u_{v'}^i$. If the candidate area already has anchor with level $l_{v'}$, indicating that the current level of detail is still insufficient to meet the optimization requirements, the granularity of anchor v' will be increased. Specifically, the 3D Gaussians that have already participated in growth will not be used for further anchor growth. In practice, the variations of voxel size and significance threshold between different levels are defined as:

$$\gamma_{l+1} = \gamma_l / 4, \quad \tau_{l+1} = \tau_l * 2. \quad (6)$$

This coarse-to-fine strategy allows us to effectively suppress the expansion of anchors caused by gradient instability.

3.3. Dynamic Keyframe Window Construction

Keyframe Insert Selection. The keyframe set is constructed based on the overlap between the visible coarsest anchors of the current frame and the last keyframe. A

new input RGB-D frame is added to the keyframe set if the ratio $N_{overlap}/N_{union}$ smaller than a threshold, where $N_{overlap}$ is the number of coarsest granularity anchors visible in both the current frame and the last keyframe, and N_{union} represents the total number of visible coarsest anchors of them. By this insertion strategy, we ensure the views in the keyframe set have relatively little overlap.

Dynamic Keyframe Window. When a new keyframe is added, it typically indicates that new unexplored areas require optimization. Relying solely on the data from this single keyframe for map optimization can lead to severe forgetting and overfitting issues, resulting in poor final reconstruction quality. Existing online mapping methods usually select a specific number of keyframes from the keyframe set, along with the newly added keyframe(some works additionally add the most recent keyframe), to form a fixed keyframe window. Only the keyframes within this fixed window are used for map optimization. However, these methods face

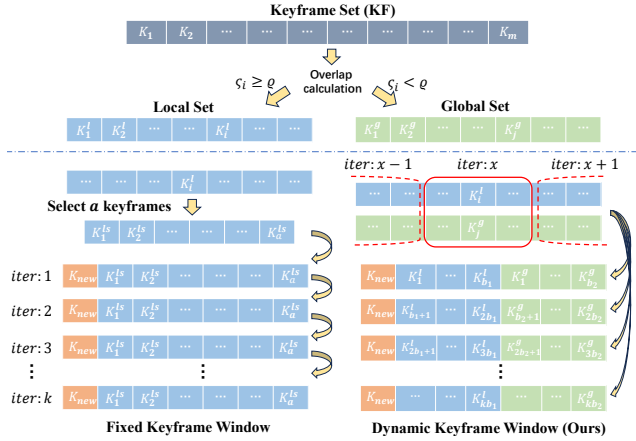


Figure 4. **Illustration of the differences between Fixed Keyframe Window and our Dynamic Keyframe Window.** Left: Fixed Keyframe Window maintains a static window content across all optimization iterations. Right: our dynamic Keyframe Window updates the keyframe window with new data each optimization iteration.

challenges in balance optimization iterations with mapping quality. Suppose the new keyframe K_{new} observes a new region with distribution ζ . Since the information for ζ is only observed by K_{new} , it is necessary to optimize K_{new} with k times to accurately fit the new region. Let a represent the number of keyframes in current window, the total number of inference operations required is $k \times (a+1)$. A diminutive b is prone to causing significant overfitting and forgetting problem, which hampers the model’s performance on previously visited area. Conversely, optimize K_{new} with a larger b can lead to an excessively high number of inference operations, potentially causing the model to become computationally intensive and time-consuming.

Further considering this issue, it becomes apparent that only the K_{new} necessitates multiple optimization iterations.

The role of the other keyframes is to alleviate local overfitting and forgetting issues. Thus, having diverse sources for these keyframes is advantageous. Based on the above analysis, instead of using a fixed keyframe window, we develop a simple yet effective dynamic keyframe window. Specifically, we first calculate the overlap weight ς between the existing keyframes in the keyframe set \mathbf{KF} and the K_{new} . \mathbf{KF} is subsequently divided into two sets with threshold ϱ : the local set ($\{K_i | K_i \in \mathbf{KF}, \varsigma_i \geq \varrho\}$) and global set ($\{K_i | K_i \in \mathbf{KF}, \varsigma_i < \varrho\}$). During each new optimization iteration, we clear the keyframe window of all keyframes except K_{new} . Then, we select b_1 keyframes from the local set and b_2 keyframes from the global set without replacement and add them to the keyframe window ($a = b_1 + b_2$). By employing the aforementioned method, we ensure that the contents of the keyframe window are always dynamically changing. This approach enables effective adaptation to newly observed regions while incorporating past experiences, thereby addressing the limitations of the fixed keyframe window method. A more intuitive comparison between the two methods can be found in Fig. 4.

3.4. Losses Design

The learnable parameters and MLPs are optimized with respect to the $L1$ loss over rendered pixel colors, denoted as \mathcal{L}_c , and depths, denoted as \mathcal{L}_d . The loss function is further augmented by SSIM term [8] \mathcal{L}_{SSIM} and scale term [22] \mathcal{L}_s :

$$\mathcal{L} = \lambda_c \mathcal{L}_c + \lambda_{SSIM} \mathcal{L}_{SSIM} + \lambda_d \mathcal{L}_d + \lambda_s \mathcal{L}_s. \quad (7)$$

where the Gaussians scale term \mathcal{L}_s is:

$$\mathcal{L}_s = \sum_{i=1}^{\mathcal{N}} Prod(\varphi^i). \quad (8)$$

The $Prod(\cdot)$ is the product computation function and \mathcal{N} is the the number of all 3D Gaussians.

4. Experiment

4.1. Experimental Setup

Dataset. To evaluate the performance of our proposed method, we compare its mapping accuracy and time consumption with other RGB-D mapping systems currently open-source on both the synthetic Replica dataset [29] and the real-world ScanNet dataset [4] in addition to a self-scanned scene provided by [27].

Metrics. Following the evaluation protocol of mapping results used in SplatAM [15], for measuring RGB rendering performance, we report standard photometric rendering quality metrics: PSNR [36], SSIM [8] and LPIPS(AlexNet) [43]. Depth rendering performance is measured by Depth L1 loss(cm). All rendering metrics are computed on each

frame with valid pose to evaluate the map quality. We report the average across five runs for all our evaluations.

Baseline Methods. We select several advanced Gaussian-based dense RGB-D SLAM methods currently open-source, MonoGS [24], SplaTAM [15] and RTG-SLAM [27] for comparison. We also benchmark our method against other RGB-D based approaches [12–14, 28] that, like ours, do not have explicit loop closure and submap.

Implementation Details. All methods are benchmarked on a desktop computer with an intel i7-14700K CPU and an Nvidia RTX 4090 GPU. As we exclusively focus on incremental mapping, we omit the tracking component of [14, 15, 24, 28] and instead employ the ground truth pose. All methods are evaluated using their single-threaded versions. Specifically, the scene refinement module used after online mapping process in MonoGS [24] was removed to ensure a fair comparison. An anchor is pruned if all the opacity of its 3D Gaussians is less than $\rho = 0.01$. To provide a more comprehensive comparison with other methods, in addition to the full version (marked as ‘Ours’), we also offer the ‘Ours-sparse’ version (using a larger ρ to prune more anchors) and the ‘Ours-sparse²’ version (based on the last version but using only half the number of iterations for mapping). More details of hyperparameters are provided in the supplementary material.

Methods	Metrics	R0	R1	R2	O0	O1	O2	O3	O4	Avg
H ₂ -Mapping [12] (RAL23)	PSNR ↑	29.67	32.34	32.23	37.87	38.93	31.02	30.64	33.07	33.22
	SSIM ↑	0.927	0.956	0.963	0.982	0.984	0.958	0.95	0.968	0.932
	LPIPS ↓	0.217	0.168	0.151	0.09	0.0945	0.203	0.221	0.166	0.164
ESLAM [14] (CVPR23)	PSNR ↑	26.40	28.34	30.25	35.10	34.76	29.07	28.83	31.15	30.49
	SSIM ↑	0.763	0.874	0.923	0.928	0.921	0.879	0.876	0.908	0.871
	LPIPS ↓	0.300	0.292	0.232	0.180	0.205	0.235	0.191	0.199	0.230
Point-SLAM [28] (ICCV23)	PSNR ↑	34.38	35.05	36.80	39.35	40.29	34.95	34.54	35.66	36.38
	SSIM ↑	0.937	0.942	0.955	0.962	0.960	0.916	0.917	0.942	0.941
	LPIPS ↓	0.093	0.106	0.102	0.088	0.103	0.147	0.116	0.129	0.111
H ₃ -Mapping [13] (arXiv24)	PSNR ↑	33.16	34.99	35.24	39.85	40.12	33.89	34.10	35.99	35.92
	SSIM ↑	0.921	0.939	0.948	0.969	0.966	0.936	0.933	0.950	0.945
	LPIPS ↓	-	-	-	-	-	-	-	-	-
MonoGS [24] (CVPR24)	PSNR ↑	33.16	34.99	35.24	39.85	40.12	34.59	34.32	36.50	35.68
	SSIM ↑	0.928	0.928	0.942	0.972	0.968	0.941	0.939	0.952	0.946
	LPIPS ↓	0.128	0.153	0.129	0.085	0.084	0.129	0.107	0.129	0.118
SplaTAM [15] (CVPR24)	PSNR ↑	33.57	34.40	36.36	40.37	40.51	33.38	32.56	35.07	35.78
	SSIM ↑	0.978	0.973	0.984	0.985	0.980	0.974	0.968	0.970	0.976
	LPIPS ↓	0.052	0.064	0.053	0.049	0.067	0.073	0.081	0.106	0.068
RTG-SLAM [27] (SIGGRAPH24)	PSNR ↑	31.30	33.95	34.95	39.28	39.63	32.90	32.86	36.12	35.12
	SSIM ↑	0.967	0.979	0.983	0.988	0.990	0.981	0.982	0.985	0.982
	LPIPS ↓	0.144	0.111	0.116	0.080	0.097	0.140	0.133	0.118	0.117
Ours	PSNR ↑	36.24	38.13	38.73	42.07	42.26	36.15	36.34	38.58	38.56
	SSIM ↑	0.971	0.974	0.976	0.986	0.981	0.973	0.971	0.975	0.976
	LPIPS ↓	0.044	0.049	0.055	0.036	0.049	0.054	0.047	0.050	0.048

Table 1. Quantitative comparison of our method against baselines for rendering results on the Replica [29] dataset.

4.2. Experiments Results.

Evaluation on Replica [29]. In Tab. 1, we evaluate our method’s rendering quality results on Replica dataset [29]. Our approach achieves the best PSNR and LPIPS results, while maintaining a highly competitive SSIM result. Fig. 5 provides a qualitative comparison of the rendering of ours and baseline method. Thanks to our hierarchical representation optimization method, our approach is able to better

Methods	Metrics	0000	0059	0106	0169	0181	0207	Avg
H ₂ -Mapping [12] (RAL23)	PSNR ↑	21.02	17.60	15.59	20.65	18.83	21.04	19.12
	SSIM ↑	0.809	0.765	0.742	0.819	0.844	0.822	0.800
	LPIPS ↓	0.475	0.445	0.522	0.419	0.492	0.472	0.471
ESLAM [14] (CVPR23)	PSNR ↑	19.10	17.84	16.92	20.42	17.60	19.03	18.64
	SSIM ↑	0.636	0.630	0.616	0.690	0.708	0.664	0.657
	LPIPS ↓	0.560	0.520	0.590	0.542	0.566	0.591	0.561
Point-SLAM [28] (ICCV23)	PSNR ↑	24.09	22.14	21.33	22.74	22.29	24.36	22.83
	SSIM ↑	0.715	0.683	0.619	0.630	0.746	0.717	0.685
	LPIPS ↓	0.471	0.480	0.554	0.560	0.520	0.512	0.516
MonoGS [24] (CVPR24)	PSNR ↑	17.20	15.44	16.92	18.79	13.76	18.10	16.72
	SSIM ↑	0.636	0.574	0.657	0.704	0.594	0.705	0.645
	LPIPS ↓	0.560	0.628	0.571	0.553	0.723	0.542	0.596
SplaTAM [15] (CVPR24)	PSNR ↑	19.70	19.65	19.11	23.23	18.58	20.64	20.15
	SSIM ↑	0.654	0.807	0.747	0.787	0.749	0.749	0.749
	LPIPS ↓	0.422	0.240	0.315	0.276	0.351	0.278	0.313
RTG-SLAM [27] (SIGGRAPH24)	PSNR ↑	18.98	17.44	16.87	19.39	17.51	18.97	18.19
	SSIM ↑	0.804	0.705	0.684	0.778	0.716	0.757	0.741
	LPIPS ↓	0.488	0.534	0.579	0.499	0.613	0.538	0.542
Ours	PSNR ↑	25.87	24.04	24.86	26.84	25.91	26.13	25.61
	SSIM ↑	0.807	0.836	0.864	0.838	0.873	0.819	0.840
	LPIPS ↓	0.345	0.270	0.271	0.285	0.322	0.340	0.306

Table 2. Quantitative comparison of our method against baselines for rendering results on the ScanNet [4] dataset.

Methods	PSNR↑	SSIM↑	LPIPS↓	Depth L1(cm)↓	Rendering FPS↑	Model Size(MB)↓
MonoGS [24]	25.31	0.808	0.543	3.94	862.0	5.6
SplaTAM [15]	22.99	0.723	0.378	1.60	54.9	540.3
RTG-SLAM [27]	24.40	0.823	0.415	2.72	316.5	270.8
Ours-sparse	28.91	0.846	0.423	1.10	647.8	5.4
Ours	30.32	0.874	0.300	0.89	559.7	30.5

Table 3. Quantitative comparison in terms of rendering, Depth L1, and memory performance on the real scanned scene provided by [27].

capture details. Tab. 4 provides the performance analysis of MonoGS [24], SplaTAM [15], RTG-SLAM [27] and our method. Since we don’t need to explicitly store all 3D Gaussian parameters, our method is more storage-efficient. Despite using fewer iterations and a more sparse distribution of 3D Gaussians, our approach still achieves superior mapping results.

Evaluation on Real-world Scene. Since the camera trajectory provided by the ScanNet [4] originates from BundleFusion [5] rather than ground truth, and the depth images contain significant noise, conducting online mapping on ScanNet is extremely challenging. Tab. 2 presents the rendering evaluation results on six scenes of ScanNet. Due to the lack of structured organization, previous Gaussian-based methods [15, 24, 27] are highly sensitive to depth noise and camera trajectory inaccuracies. In contrast, our method effectively filters out these noises, resulting in more realistic rendering outcomes. In Tab. 5, we compare the performance of OG-Mapping with other Gaussian-based methods. The results demonstrate that our method can construct more compact maps at a faster processing speed. Fig. 1 and Fig. 6 also illustrate that our method can achieve better visual quality with more reliable geometry and texture details. The rendering viewpoints of these results are selected outside of the training views, with random perturbations added to the poses. Since the depth information provided by ScanNet

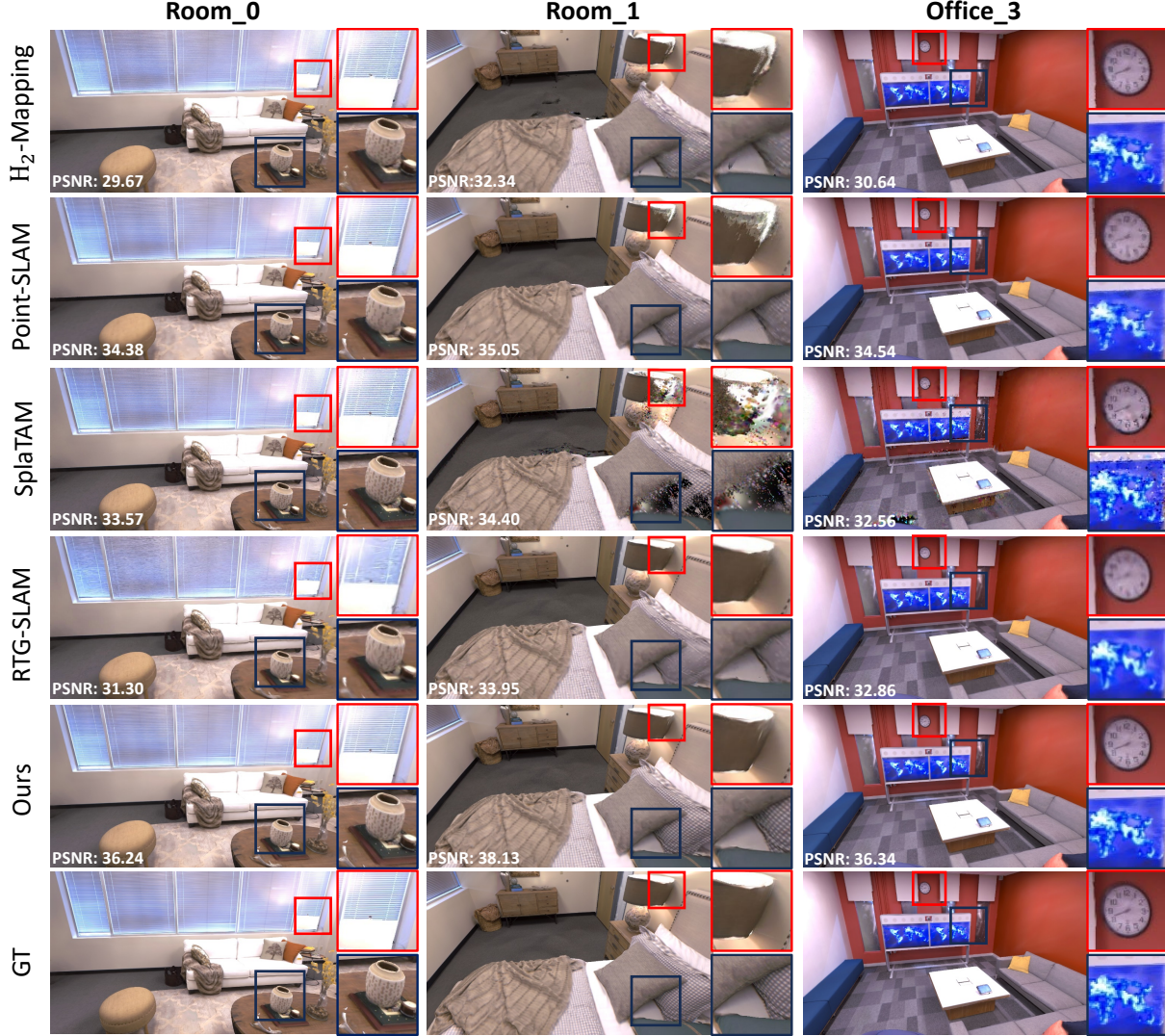


Figure 5. Qualitative comparison of rendering results across three scenes from the Replica dataset. Key details are highlighted by colored boxes. The average PSNR metric for each scene is indicated in the lower-left corner.

[4] contains excessive noise, the Depth-L1 method used in SplatTAM [15] is not suitable for measuring geometric accuracy. In Tab. 3, we provide additional experimental results on the high-quality large indoor scene [27] to further validate the effectiveness of our method in real-world scenarios.

Methods	PSNR \uparrow	Depth L1 [cm] \downarrow	Processing FPS \uparrow	Model Size[MB] \downarrow
MonoGS [24]	35.68	2.12	1.2	9.7
SplatTAM [15]	35.78	0.61	0.4	275.1
RTG-SLAM [27]	35.12	1.69	15.6	71.5
Ours-sparse	36.75	0.72	16.1	8.9
Ours-sparse	37.00	0.69	8.5	8.8
Ours	38.56	0.46	5.6	34.6

Table 4. Comparison of key performance metrics on the Replica [29] dataset between ours and baseline Gaussian-based methods [15, 24, 27].

Methods	PSNR \uparrow	Processing FPS \uparrow	Model Size[MB] \downarrow
MonoGS [24]	16.72	2.9	5.6
SplatTAM [15]	19.99	1.8	160.7
RTG-SLAM [27]	18.19	11.4	153.2
Ours-sparse	23.87	12.5	3.1
Ours	25.61	10.1	36.1

Table 5. Comparison of key performance metrics on ScanNet [4] between ours and baseline Gaussian-based methods [15, 24, 27].

Methods	PSNR \uparrow	Processing FPS \uparrow	Model Size(MB) \downarrow
Projection-based [15]	36.18	2.95	71.8
Naive Unique	36.24	3.37	39.9
Octree-based(ours)	36.24	5.16	39.9

Table 6. Ablation study of map densification methods quantitative results on room0 [29].

Dataset	Methods	PSNR \uparrow	Model Size(MB) \downarrow
Replica	w/o Growing	36.85	29.4
	Error-based [23]	38.54	40.5
	Anchor-based(ours)	38.56	34.6
ScanNet	w/o Growing	24.89	19.9
	Error-based [23]	25.27	78.4
	Anchor-based(ours)	25.61	36.1

Table 7. Ablation study of our progressive map refinement strategy on Replica [29] dataset and ScanNet [4] dataset.

Type	Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
Fixed Window	Random [14]	37.14	0.969	0.057
	Overlap [15, 24]	35.85	0.963	0.063
	Coverage-Maximizing [12]	37.75	0.972	0.050
Dynamic Window	Global(ours)	38.56	0.976	0.048

Table 8. Ablation study of our dynamic keyframe window methods on Replica [29] dataset.

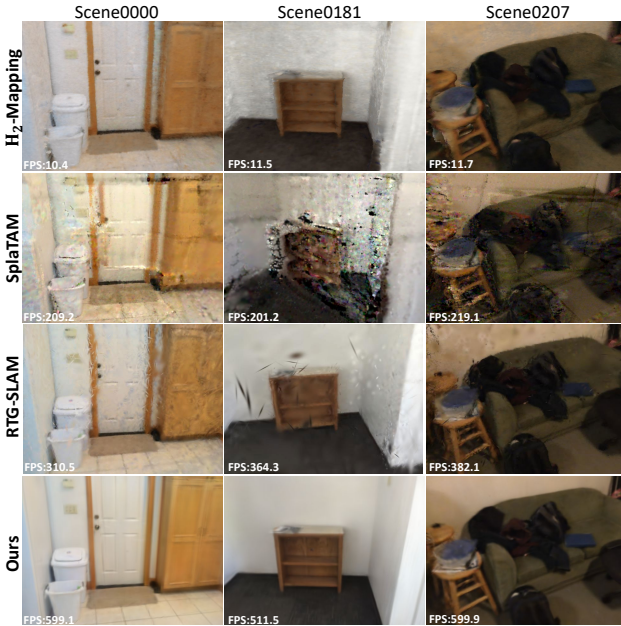


Figure 6. Qualitative comparison of rendering results across three scenes from ScanNet. The average rendering FPS are indicated at the bottom of the images.

Map Densification Ablations. To prove the effectiveness of our octree-based structured map densification method, we compare the mapping performance of Depth error-based method [15], Naive Unique method (Using the unique function to determine whether the anchors are newly observed), and our octree-based method on the room0 of Replica. As shown in Tab. 6, The projection-based method inevitably adds extra anchors at object edges incorrectly, resulting in

larger map size and worse rendering results. While the Naive Unique method can handle anchor growth correctly like ours, the inefficient unique computation results in an intolerable time overhead. Our octree-based method leverages the efficiency of sparse octrees, combining both processing accuracy and high efficiency.

Progressive Map Refinement Ablations. We evaluated our progressive map refinement described in Sec. 3.2. Tab. 7 shows the results of disabling growing operation and employing the error-based method [23] on the Replica and the ScanNet datasets. The results show the growing operation is crucial for accurately reconstructing details. Compared to the error-based method [23], our anchor-based method makes better use of structured information, resulting in smaller map occupancy while maintaining comparable rendering accuracy.

Keyframe Window Ablations. In Tab. 8, we compared our dynamic keyframe window method with other existing keyframe window construction approaches. The overlap method, which only uses keyframes that overlap with the newly added keyframe, is not suitable for addressing the problem of forgetting. Consequently, it results in the poorest rendering quality. The Coverage-Maximizing method [12] uses the most recently unoptimized keyframes to construct the keyframe window. However, the reconstruction process tends to overfit these areas as this window is fixed. In contrast, our dynamic keyframe window overcomes these issues, achieving optimal results.

5. Conclusion

In this work, we introduce OG-Mapping, a novel framework for effective online dense mapping. By utilizing an octree-based structured 3D Gaussians representation, OG-Mapping achieves efficient map densification and compaction. Additionally, we propose an anchor-based progressive map refinement strategy to enhance the capture of finer details. Furthermore, we develop a dynamic keyframe window to mitigate the issues of local overfitting and forgetting problems encountered during the reconstruction process. Experiment results demonstrate that this approach, leveraging a more compact map, outperforms existing algorithms. The advantages of our structural representation and dynamic keyframe window are particularly evident in challenging real scenes where existing Gaussian-based online mapping methods typically falter.

Limitation. OG-Mapping relies on inputs from an RGB-D sensor to build a sparse octree. Future research is anticipated to explore the use of monocular image input alone. Our method employs MLPs as feature decoders, successfully constructing compact maps. However, in particularly large scenes, it may still encounter forgetting issues. Further research on submap construction [18, 42] could further reduce frame processing time.

References

- [1] Yan-Pei Cao, Leif Kobbelt, and Shi-Min Hu. Real-time high-accuracy three-dimensional reconstruction with consumer rgb-d cameras. *ACM Transactions on Graphics (TOG)*, 37(5):1–16, 2018. 1, 2
- [2] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16123–16133, 2022. 1
- [3] Jiawen Chen, Dennis Bautembach, and Shahram Izadi. Scalable real-time volumetric surface reconstruction. *ACM Trans. Graph.*, 32(4):113–1, 2013. 1, 2
- [4] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5828–5839, 2017. 2, 5, 6, 7, 8
- [5] Angela Dai, Matthias Nießner, Michael Zollhöfer, Shahram Izadi, and Christian Theobalt. Bundlefusion: Real-time globally consistent 3d reconstruction using on-the-fly surface reintegration. *ACM Transactions on Graphics (ToG)*, 36(4):1, 2017. 6
- [6] Tianchen Deng, Guole Shen, Tong Qin, Jianyu Wang, Wentao Zhao, Jingchuan Wang, Danwei Wang, and Weidong Chen. Plgslam: Progressive neural scene representation with local to global bundle adjustment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19657–19666, 2024. 1
- [7] Hao Du, Peter Henry, Xiaofeng Ren, Marvin Cheng, Dan B Goldman, Steven M Seitz, and Dieter Fox. Interactive 3d modeling of indoor environments with a consumer depth camera. In *Proceedings of the 13th international conference on Ubiquitous computing*, pages 75–84, 2011. 1, 2
- [8] Alain Hore and Djemel Ziou. Image quality metrics: Psnr vs. ssim. In *2010 20th international conference on pattern recognition*, pages 2366–2369. IEEE, 2010. 5
- [9] 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
- [10] Jiarui Hu, Xianhao Chen, Boyin Feng, Guanglin Li, Liangjing Yang, Hujun Bao, Guofeng Zhang, and Zhaopeng Cui. Cg-slam: Efficient dense rgb-d slam in a consistent uncertainty-aware 3d gaussian field. *arXiv preprint arXiv:2403.16095*, 2024. 2
- [11] Huajian Huang, Longwei Li, Hui Cheng, and Sai-Kit Yeung. Photo-slam: Real-time simultaneous localization and photo-realistic mapping for monocular, stereo, and rgb-d cameras. *arXiv preprint arXiv:2311.16728*, 2023. 2
- [12] Chenxing Jiang, Hanwen Zhang, Peize Liu, Zehuan Yu, Hui Cheng, Boyu Zhou, and Shaojie Shen. H₂-mapping: Real-time dense mapping using hierarchical hybrid representation. *IEEE Robotics and Automation Letters*, 2023. 1, 2, 6, 8
- [13] Chenxing Jiang, Yiming Luo, Boyu Zhou, and Shaojie Shen. H3-mapping: Quasi-heterogeneous feature grids for real-time dense mapping using hierarchical hybrid representation. *arXiv preprint arXiv:2403.10821*, 2024. 2, 6
- [14] Mohammad Mahdi Johari, Camilla Carta, and François Fleuret. Eslam: Efficient dense slam system based on hybrid representation of signed distance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17408–17419, 2023. 1, 2, 6, 8
- [15] Nikhil Keetha, Jay Karhade, Krishna Murthy Jatavallabhula, Gengshan Yang, Sebastian Scherer, Deva Ramanan, and Jonathan Luiten. Splatam: Splat, track & map 3d gaussians for dense rgb-d slam. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. 2, 4, 5, 6, 7, 8
- [16] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4):1–14, 2023. 1, 2, 4
- [17] Lukas Koestler, Nan Yang, Niclas Zeller, and Daniel Cremers. Tandem: Tracking and dense mapping in real-time using deep multi-view stereo. In *Conference on Robot Learning*, pages 34–45. PMLR, 2022. 2
- [18] Mingrui Li, Jingwei Huang, Lei Sun, Aaron Xuxiang Tian, Tianchen Deng, and Hongyu Wang. Ngm-slam: Gaussian splatting slam with radiance field submap. *arXiv preprint arXiv:2405.05702*, 2024. 2, 8
- [19] Zhaoshuo Li, Thomas Müller, Alex Evans, Russell H Taylor, Mathias Unberath, Ming-Yu Liu, and Chen-Hsuan Lin. Neuralangelo: High-fidelity neural surface reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8456–8465, 2023. 4
- [20] Chen-Hsuan Lin, Wei-Chiu Ma, Antonio Torralba, and Simon Lucey. Barf: Bundle-adjusting neural radiance fields. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5741–5751, 2021. 4
- [21] Shuhong Liu, Heng Zhou, Liuzhuozheng Li, Yun Liu, Tianchen Deng, Yiming Zhou, and Mingrui Li. Structure gaussian slam with manhattan world hypothesis. *arXiv preprint arXiv:2405.20031*, 2024. 2
- [22] Stephen Lombardi, Tomas Simon, Gabriel Schwartz, Michael Zollhoefer, Yaser Sheikh, and Jason Saragih. Mixture of volumetric primitives for efficient neural rendering. *ACM Transactions on Graphics (ToG)*, 40(4):1–13, 2021. 5
- [23] Tao Lu, Mulin Yu, Linning Xu, Yuanbo Xiangli, Limin Wang, Dahua Lin, and Bo Dai. Scaffold-gs: Structured 3d gaussians for view-adaptive rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20654–20664, 2024. 1, 2, 3, 4, 8
- [24] Hidenobu Matsuki, Riku Murai, Paul H. J. Kelly, and Andrew J. Davison. Gaussian Splatting SLAM. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. 2, 6, 7, 8
- [25] 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. *Communications of the ACM*, 65(1):99–106, 2021. 1, 2

- [26] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multi-resolution hash encoding. *ACM Transactions on Graphics (ToG)*, 41(4):1–15, 2022. 1, 3
- [27] Zhexi Peng, Tianjia Shao, Yong Liu, Jingke Zhou, Yin Yang, Jingdong Wang, and Kun Zhou. Rtg-slam: Real-time 3d reconstruction at scale using gaussian splatting. *arXiv preprint arXiv:2404.19706*, 2024. 2, 4, 5, 6, 7
- [28] Erik Sandström, Yue Li, Luc Van Gool, and Martin R Oswald. Point-slam: Dense neural point cloud-based slam. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 18433–18444, 2023. 2, 6
- [29] Julian Straub, Thomas Whelan, Lingni Ma, Yufan Chen, Erik Wijmans, Simon Green, Jakob J Engel, Raul Mur-Artal, Carl Ren, Shobhit Verma, et al. The replica dataset: A digital replica of indoor spaces. *arXiv preprint arXiv:1906.05797*, 2019. 5, 6, 7, 8
- [30] Edgar Sucar, Shikun Liu, Joseph Ortiz, and Andrew J Davison. imap: Implicit mapping and positioning in real-time. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6229–6238, 2021. 1, 2
- [31] Keisuke Tateno, Federico Tombari, Iro Laina, and Nassir Navab. Cnn-slam: Real-time dense monocular slam with learned depth prediction. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6243–6252, 2017. 2
- [32] Zachary Teed and Jia Deng. Droid-slam: Deep visual slam for monocular, stereo, and rgb-d cameras. *Advances in neural information processing systems*, 34:16558–16569, 2021. 2
- [33] Emanuele Vespa, Nikolay Nikolov, Marius Grimm, Luigi Nardi, Paul HJ Kelly, and Stefan Leutenegger. Efficient octree-based volumetric slam supporting signed-distance and occupancy mapping. *IEEE Robotics and Automation Letters*, 3(2):1144–1151, 2018. 2
- [34] Hengyi Wang, Jingwen Wang, and Lourdes Agapito. Co-slam: Joint coordinate and sparse parametric encodings for neural real-time slam. In *Proceedings of the IEEE international conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 1, 2
- [35] 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. *arXiv preprint arXiv:2106.10689*, 2021. 1
- [36] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. 5
- [37] Thomas Whelan, Stefan Leutenegger, Renato F Salas-Moreno, Ben Glocker, and Andrew J Davison. Elasticfusion: Dense slam without a pose graph. In *Robotics: science and systems*, page 3. Rome, Italy, 2015. 1, 2
- [38] Chi Yan, Delin Qu, Dong Wang, Dan Xu, Zhigang Wang, Bin Zhao, and Xuelong Li. Gs-slam: Dense visual slam with 3d gaussian splatting. 2024. 2
- [39] Nan Yang, Lukas von Stumberg, Rui Wang, and Daniel Cremers. D3vo: Deep depth, deep pose and deep uncertainty for monocular visual odometry. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1281–1292, 2020. 2
- [40] Xingrui Yang, Hai Li, Hongjia Zhai, Yuhang Ming, Yuqian Liu, and Guofeng Zhang. Vox-fusion: Dense tracking and mapping with voxel-based neural implicit representation. In *2022 IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*, pages 499–507. IEEE, 2022. 2
- [41] Zehao Yu, Anpei Chen, Binbin Huang, Torsten Sattler, and Andreas Geiger. Mip-splatting: Alias-free 3d gaussian splatting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19447–19456, 2024. 1
- [42] Vladimir Yugay, Yue Li, Theo Gevers, and Martin R. Oswald. Gaussian-slam: Photo-realistic dense slam with gaussian splatting, 2023. 2, 8
- [43] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018. 5
- [44] Yizhong Zhang, Weiwei Xu, Yiyong Tong, and Kun Zhou. Online structure analysis for real-time indoor scene reconstruction. *ACM Transactions on Graphics (TOG)*, 34(5):1–13, 2015. 1, 2
- [45] Xingguang Zhong, Yue Pan, Jens Behley, and Cyrill Stachniss. Shine-mapping: Large-scale 3d mapping using sparse hierarchical implicit neural representations. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 8371–8377. IEEE, 2023. 2
- [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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12786–12796, 2022. 1, 2