

Gaussian Splatting SLAM

Hideobu Matsuki^{1*}

Riku Murai^{2*}

Paul H. J. Kelly²

Andrew J. Davison¹

¹Dyson Robotics Laboratory, Imperial College London

²Software Performance Optimisation Group, Imperial College London

{h.matsuki20, riku.murai15, p.kelly, a.davison}@imperial.ac.uk

Website: <https://rmurai.co.uk/projects/GaussianSplattingSLAM/>

Video: https://youtu.be/x604ghp9R_Q/

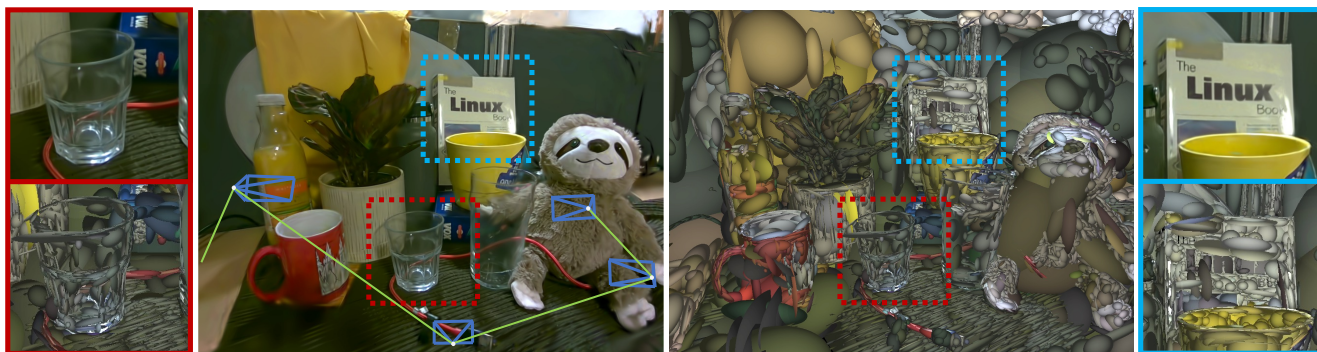


Figure 1. From a single monocular camera, we reconstruct a high fidelity 3D scene live at 3fps. For every incoming RGB frame, 3D Gaussians are incrementally formed and optimised together with the camera poses. We show both the rasterised Gaussians (left) and Gaussians shaded to highlight the geometry (right). Notice the details and the complex material properties (e.g. transparency) captured. Thin structures such as wires are accurately represented by numerous small, elongated Gaussians, and transparent objects are effectively represented by placing the Gaussians along the rim. Our system significantly advances the fidelity a live monocular SLAM system can capture.

Abstract

We present the first application of 3D Gaussian Splatting in monocular SLAM, the most fundamental but the hardest setup for Visual SLAM. Our method, which runs live at 3fps, utilises Gaussians as the only 3D representation, unifying the required representation for accurate, efficient tracking, mapping, and high-quality rendering. Designed for challenging monocular settings, our approach is seamlessly extendable to RGB-D SLAM when an external depth sensor is available. Several innovations are required to continuously reconstruct 3D scenes with high fidelity from a live camera. First, to move beyond the original 3DGS algorithm, which requires accurate poses from an offline Structure from Motion (SfM) system, we formulate camera tracking for 3DGS using direct optimisation against the 3D Gaussians, and

show that this enables fast and robust tracking with a wide basin of convergence. Second, by utilising the explicit nature of the Gaussians, we introduce geometric verification and regularisation to handle the ambiguities occurring in incremental 3D dense reconstruction. Finally, we introduce a full SLAM system which not only achieves state-of-the-art results in novel view synthesis and trajectory estimation but also reconstruction of tiny and even transparent objects.

1. Introduction

A long-term goal of online reconstruction with a single moving camera is near-photorealistic fidelity, which will surely allow new levels of performance in many areas of Spatial AI and robotics as well as opening up a whole range of new applications. While we increasingly see the benefit of applying powerful pre-trained priors to 3D recon-

*Authors contributed equally to this work.

struction, a key avenue for progress is still the invention and development of core 3D representations with advantageous properties. Many “layered” SLAM methods exist which tackle the SLAM problem by integrating multiple different 3D representations or existing SLAM components; however, the most interesting advances are when a new unified dense representation can be used for all aspects of a system’s operation: local representation of detail, large-scale geometric mapping and also camera tracking by direct alignment.

In this paper, we present the first online visual SLAM system based solely on the 3D Gaussian Splatting (3DGS) representation [11] recently making a big impact in offline scene reconstruction. In 3DGS a scene is represented by a large number of Gaussian blobs with orientation, elongation, colour and opacity. Other previous world/map-centric scene representations used for visual SLAM include occupancy or Signed Distance Function (SDF) voxel grids [24]; meshes [30]; point or surfel clouds [10, 31]; and recently neural fields [35]. Each of these has disadvantages: grids use significant memory and have bounded resolution, and even if octrees or hashing allow more efficiency they cannot be flexibly warped for large corrections [26, 39]; meshes require difficult, irregular topology to fuse new information; surfel clouds are discontinuous and difficult to fuse and optimise; and neural fields require expensive per-pixel raycasting to render. We show that 3DGS has none of these weaknesses. As a SLAM representation, it is most similar to point and surfel clouds, and inherits their efficiency, locality and ability to be easily warped or modified. However, it also represents geometry in a smooth, continuously differentiable way: a dense cloud of Gaussians merge together and jointly define a continuous volumetric function. And crucially, the design of modern graphics cards means that a large number of Gaussians can be efficiently rendered via “splatting” rasterisation, up to 200fps at 1080p. This rapid, differentiable rendering is integral to the tracking and map optimisation loops in our system.

The 3DGS representation has up until now only been used in offline systems for 3D reconstruction with known camera poses, and we present several innovations to enable online SLAM. We first derive the analytic Jacobian on Lie group of camera pose with respect to a 3D Gaussians map, and show that this can be seamlessly integrated into the existing differentiable rasterisation pipeline to enable camera poses to be optimised alongside scene geometry. Second, we introduce a novel Gaussian isotropic shape regularisation to ensure geometric consistency, which we have found is important for incremental reconstruction. Third, we propose a novel Gaussian resource allocation and pruning method to keep the geometry clean and enable accurate camera tracking. Our experimental results demonstrate photorealistic online local scene reconstruction, as well as

state-of-the-art camera trajectory estimation and mapping for larger scenes compared to other rendering-based SLAM methods. We further show the uniqueness of the Gaussian-based SLAM method such as an extremely large camera pose convergence basin, which can also be useful for map-based camera localisation. Our method works with only monocular input, one of the most challenging scenarios in SLAM. To highlight the intrinsic capability of 3D Gaussian for camera localisation, our method does not use any pre-trained monocular depth predictor or other existing tracking modules, but relies solely on RGB image inputs in line with the original 3DGS. Since this is one of the most challenging SLAM scenario, we also show our method can easily be extended to RGB-D SLAM when depth measurements are available.

In summary, our contributions are as follows:

- The first near real-time SLAM system which works with a 3DGS as the only underlying scene representation, which can handle monocular only inputs.
- Novel techniques within the SLAM framework, including the analytic Jacobian on Lie group for direct camera pose estimation, isotropic regularisation of the Gaussian shape, and geometric verification.
- Extensive evaluations on a variety of datasets both for monocular and RGB-D settings, demonstrating competitive performance, particularly in real-world scenarios.

2. Related Work

Dense SLAM: Dense visual SLAM focuses on reconstructing detailed 3D maps, unlike sparse SLAM methods which excel in pose estimation [5, 6, 22] but typically yield maps useful mainly for localisation. In contrast, dense SLAM creates interactive maps beneficial for broader applications, including AR and robotics. Dense SLAM methods are generally divided into two primary categories: Frame-centric and Map-centric. **Frame-centric SLAM** minimises photometric error across consecutive frames, jointly estimating per-frame depth and frame-to-frame camera motion. Frame-centric approaches [2, 38] are efficient, as individual frames host local rather than global geometry (e.g. depth maps), and are attractive for long-session SLAM, but if a dense global map is needed, it must be constructed on demand by assembling all of these parts which are not necessarily fully consistent. In contrast, **Map-centric SLAM** uses a unified 3D representation across the SLAM pipeline, enabling a compact and streamlined system. Compared to purely local frame-to-frame tracking, a map-centric approach leverages global information by tracking against the reconstructed 3D consistent map. Classical map-centric approaches often use voxel grids [3, 24, 27, 42] or points [10, 31, 43] as the underlying 3D representation. While voxels enable a fast look-up of features in 3D, the representation is expensive, and the fixed voxel resolution and distribution

are problematic when the spatial characteristics of the environment are not known in advance. On the other hand, a point-based map representation, such as surfel clouds, enables adaptive changes in resolution and spatial distribution by dynamic allocation of point primitives in the 3D space. Such flexibility benefits online applications such as SLAM with deformation-based loop closure [31, 43]. However, optimising the representation to capture high fidelity is challenging due to the lack of correlation among the primitives. Recently, in addition to classical graphic primitives, neural network-based map representations are a promising alternative. iMAP [35] demonstrated the interesting properties of neural representation, such as sensible hole filling of unobserved geometry. Many recent approaches combine the classical and neural representations to capture finer details [9, 29, 48, 49]; however, the large amount of computation required for neural rendering makes the live operation of such systems challenging.

Differentiable Rendering: The classical method for creating a 3D representation was to unproject 2D observations into 3D space and to fuse them via weighted averaging [17, 24]. Such an averaging scheme suffers from over-smooth representation and lacks the expressiveness to capture high-quality details. To capture a scene with photo-realistic quality, differentiable volumetric rendering [25] has recently been popularised with Neural Radiance Fields (NeRF) [18]. Using a single Multi-Layer Perceptron (MLP) as a scene representation, NeRF performs volume rendering by marching along pixel rays, querying the MLP for opacity and colour. Since volume rendering is naturally differentiable, the MLP representation is optimised to minimise the rendering loss using multiview information to achieve high-quality novel view synthesis. The main weakness of NeRF is its training speed. Recent developments have introduced explicit volume structures such as multi-resolution voxel grids [7, 15, 36] or hash functions [20] to improve performance. Interestingly, these projects demonstrate that the main contributor to high-quality novel view synthesis is not the neural network but rather differentiable volumetric rendering, and that it is possible to avoid the use of an MLP and yet achieve comparable rendering quality to NeRF [7]. However, even in these systems, per-pixel ray marching remains a significant bottleneck for rendering speed. This issue is particularly critical in SLAM, where immediate interaction with the map is essential for tracking. In contrast to NeRF, 3DGS performs differentiable rasterisation. Similar to regular graphics rasterisations, by iterating over the primitives to be rasterised rather than marching along rays, 3DGS leverages the natural sparsity of a 3D scene and achieves a representation which is expressive to capture high-fidelity 3D scenes while offering significantly faster rendering. Several works have applied 3D Gaussians and differentiable rendering to static scene cap-

ture [12, 40], and in particular more recent works utilise 3DGS and demonstrate superior results in vision tasks such as dynamic scene capture [16, 44, 46] and 3D generation [37, 47]. Our method adopts a Map-centric approach, utilising 3D Gaussians as the only SLAM representation. Similar to surfel-based SLAM, we dynamically allocate the 3D Gaussians, enabling us to model an arbitrary spatial distribution in the scene. Unlike other methods such as ElasticFusion [43] and PointFusion [10], however, by using differentiable rasterisation, our SLAM system can capture high-fidelity scene details and represent challenging object properties by direct optimisation against information from every pixel.

3. Method

3.1. Gaussian Splatting

Our SLAM representation is 3DGS, mapping the scene with a set of anisotropic Gaussians \mathcal{G} . Each Gaussian \mathcal{G}^i contains optical properties: colour c^i and opacity α^i . For continuous 3D representation, the mean μ_W^i and covariance Σ_W^i , defined in the world coordinate, represent the Gaussian’s position and its ellipsoidal shape. We omit the spherical harmonics (SHs) representing view-dependent radiance for simplicity but report the ablation with SHs in the supplementary. Since 3DGS uses volume rendering, explicit extraction of the surface is not required. Instead, by splatting and blending \mathcal{N} Gaussians, a pixel colour C_p is synthesised:

$$C_p = \sum_{i \in \mathcal{N}} c_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j). \quad (1)$$

3DGS performs rasterisation, iterating over the Gaussians rather than marching along the camera rays, and hence, free spaces are ignored during rendering. During rasterisation, the contributions of α are decayed via a Gaussian function, based on the 2D Gaussian formed by splatting a 3D Gaussian. The 3D Gaussians $\mathcal{N}(\mu_W, \Sigma_W)$ in world coordinates are related to the 2D Gaussians $\mathcal{N}(\mu_I, \Sigma_I)$ on the image plane through a projective transformation:

$$\mu_I = \pi(T_{CW} \cdot \mu_W), \Sigma_I = \mathbf{J} \mathbf{W} \Sigma_W \mathbf{W}^T \mathbf{J}^T, \quad (2)$$

where π is the projection operation and $T_{CW} \in \mathbf{SE}(3)$ is the camera pose of the viewpoint. \mathbf{J} is the Jacobian of the linear approximation of the projective transformation and \mathbf{W} is the rotational component of T_{CW} . This formulation enables the 3D Gaussians to be differentiable and the blending operation provides gradient flow to the Gaussians. Using first-order gradient descent [13], Gaussians gradually refine both their optic and geometric parameters to represent the captured scene with high fidelity.

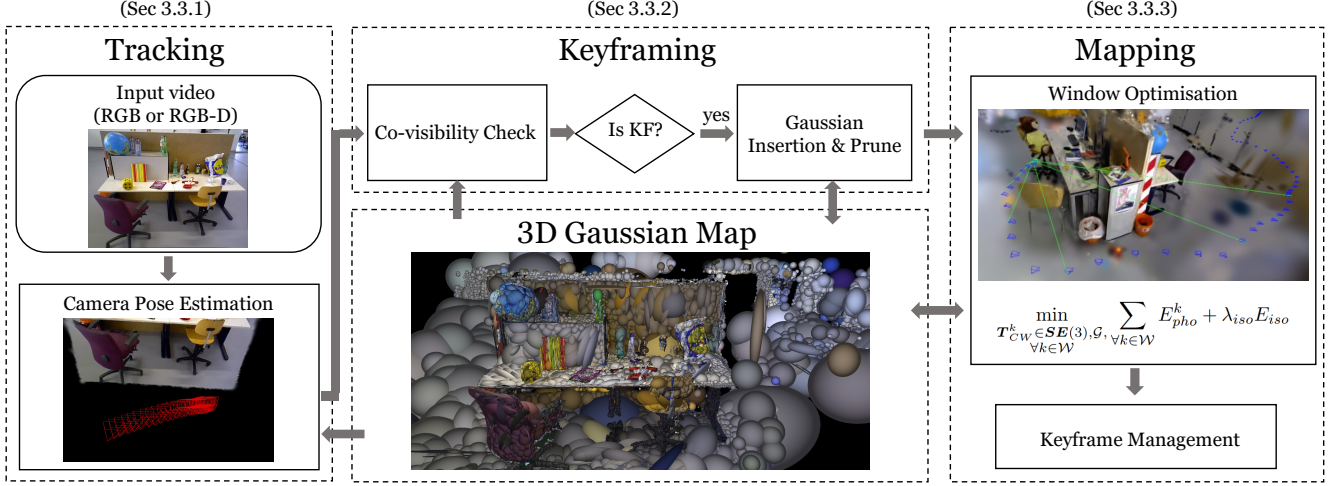


Figure 2. **SLAM System Overview:** Our SLAM system uses 3D Gaussians as the only representation, unifying all components of SLAM, including tracking, mapping, keyframe management, and novel view synthesis.

3.2. Camera Pose Optimisation

To achieve accurate tracking, we typically require at least 50 iterations of gradient descent per frame. This requirement emphasises the necessity of a representation with computationally efficient view synthesis and gradient computation, making the choice of 3D representation a crucial part of designing a SLAM system.

In order to avoid the overhead of automatic differentiation, 3DGS implements rasterisation with CUDA with derivatives for all parameters calculated explicitly. Since rasterisation is performance critical, we similarly derive the camera Jacobians explicitly.

To the best of our knowledge, we provide the first analytical Jacobian of $SE(3)$ camera pose with respect to the 3D Gaussians used in EWA splatting [50] and 3DGS. This opens up new applications of 3DGS beyond SLAM.

We use Lie algebra to derive the minimal Jacobians, ensuring that the dimensionality of the Jacobians matches the degrees of freedom, eliminating any redundant computations. The terms of Eq. (2) are differentiable with respect to the camera pose T_{CW} ; using the chain rule:

$$\frac{\partial \mu_I}{\partial T_{CW}} = \frac{\partial \mu_I}{\partial \mu_C} \frac{\mathcal{D} \mu_C}{\mathcal{D} T_{CW}}, \quad (3)$$

$$\frac{\partial \Sigma_I}{\partial T_{CW}} = \frac{\partial \Sigma_I}{\partial \Sigma} \frac{\partial \mathbf{J}}{\partial \mu_C} \frac{\mathcal{D} \mu_C}{\mathcal{D} T_{CW}} + \frac{\partial \Sigma_I}{\partial \mathbf{W}} \frac{\mathcal{D} \mathbf{W}}{\mathcal{D} T_{CW}}. \quad (4)$$

where T_{CW} represents the 3D position of Gaussian in the camera coordinate. We take the derivatives on the manifold to derive minimal parameterisation. Borrowing the notation from [32], let $T \in SE(3)$ and $\tau \in \mathfrak{se}(3)$. We define the partial derivative on the manifold as:

$$\frac{\mathcal{D} f(T)}{\mathcal{D} T} \triangleq \lim_{\tau \rightarrow 0} \frac{\text{Log}(f(\text{Exp}(\tau) \circ T) \circ f(T)^{-1})}{\tau}, \quad (5)$$

where \circ is a group composition, and Exp , Log are the exponential and logarithmic mappings between Lie algebra and Lie Group. With this, we derive the following:

$$\frac{\mathcal{D} \mu_C}{\mathcal{D} T_{CW}} = [\mathbf{I} \quad -\mu_C^\times], \quad \frac{\mathcal{D} \mathbf{W}}{\mathcal{D} T_{CW}} = \begin{bmatrix} \mathbf{0} & -\mathbf{W}_{:,1}^\times \\ \mathbf{0} & -\mathbf{W}_{:,2}^\times \\ \mathbf{0} & -\mathbf{W}_{:,3}^\times \end{bmatrix}, \quad (6)$$

where $^\times$ denotes the skew symmetric matrix of a 3D vector, and $\mathbf{W}_{:,i}$ refers to the i th column of the matrix.

3.3. SLAM

In this section, we present details of full SLAM framework. The overview of the system is summarised in Fig. 2. Please refer to the supplementary material for the further parameter details.

3.3.1 Tracking

In tracking only the current camera pose is optimised, without updates to the map representation. In the monocular case, we minimise the following photometric residual:

$$E_{pho} = \|I(\mathcal{G}, T_{CW}) - \bar{I}\|_1, \quad (7)$$

where $I(\mathcal{G}, T_{CW})$ renders the Gaussians \mathcal{G} from T_{CW} , and \bar{I} is an observed image.

We further optimise affine brightness parameters for varying exposure and penalise non-edge or low-opacity pixels. When depth observations are available, we define the geometric residual as:

$$E_{geo} = \|D(\mathcal{G}, T_{CW}) - \bar{D}\|_1, \quad (8)$$

where $D(\mathcal{G}, \mathbf{T}_{CW})$ is depth rasterisation and \bar{D} is the observed depth. Rather than simply using the depth measurements to initialise the Gaussians, we minimise both photometric and geometric residuals: $\lambda_{pho}E_{pho} + (1 - \lambda_{pho})E_{geo}$, where λ_{pho} is a hyperparameter.

As in Eq. (1), per-pixel depth is rasterised by alpha-blending:

$$\mathcal{D}_p = \sum_{i \in \mathcal{N}} z_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j), \quad (9)$$

where z_i is the distance to the mean μ_W of Gaussian i along the camera ray. We derive analytical Jacobians for the camera pose optimisation in a similar manner to Eq. (3), (4).

3.3.2 Keyframing

Since using all the images from a video stream to jointly optimise the Gaussians and camera poses online is infeasible, we maintain a small window \mathcal{W}_k consisting of carefully selected keyframes based on inter-frame covisibility. Ideal keyframe management will select non-redundant keyframes observing the same area, spanning a wide baseline to provide better multiview constraints. The parameters are detailed in the supplementary.

Selection and Management Every tracked frame is checked for keyframe registration based on our simple yet effective criteria. We measure the covisibility by measuring the intersection over the union of the observed Gaussians between the current frame i and the last keyframe j . If the covisibility drops below a threshold, or if the relative translation t_{ij} is large with respect to the median depth, frame i is registered as a keyframe. For efficiency, we maintain only a small number of keyframes in the current window \mathcal{W}_k following the keyframe management heuristics of DSO [5]. The main difference is that a keyframe is removed from the current window if the overlap coefficient with the latest keyframe drops below a threshold.

Gaussian Covisibility An accurate estimate of covisibility simplifies keyframe selection and management. 3DGS respects visibility ordering since the 3D Gaussians are sorted along the camera ray. This property is desirable for covisibility estimation as occlusions are handled by design. A Gaussian is marked to be visible from a view if used in the rasterisation and if the ray’s accumulated α has not yet reached 0.5. This enables our estimated covisibility to handle occlusions without requiring additional heuristics.

Gaussian Insertion and Pruning At every keyframe, new Gaussians are inserted into the scene to capture newly visible scene elements and to refine the fine details. When depth measurements are available, Gaussian means μ_W are

initialised by back-projecting the depth. In the monocular case, we render the depth at the current frame. For pixels with depth estimates, μ_W are initialised around those depths with low variance; for pixels without the depth estimates, we initialise μ_W around the median depth of the rendered image with high variance.

In the monocular case, the positions of many newly inserted Gaussians are incorrect. While the majority will quickly vanish during optimisation as they violate multi-view consistency, we further prune the excess Gaussians by checking the visibility amongst the current window \mathcal{W}_k . If the Gaussians inserted within the last 3 keyframes are unobserved by at least 3 other frames, we prune them out as they are geometrically unstable.

3.3.3 Mapping

The purpose of mapping is to maintain a coherent 3D structure and to optimise the newly inserted Gaussians. During mapping, the keyframes in \mathcal{W}_k are used to reconstruct currently visible regions. Additionally, two random past keyframes \mathcal{W}_r are selected per iteration to avoid forgetting the global map. Rasterisation of 3DGS imposes no constraint on the Gaussians along the viewing ray direction, even with a depth observation. This is not a problem when sufficient carefully selected viewpoints are provided (e.g. in the novel view synthesis case); however, in continuous SLAM this causes many artefacts, making tracking challenging. We therefore introduce an isotropic regularisation:

$$E_{iso} = \sum_{i=1}^{|\mathcal{G}|} \|\mathbf{s}_i - \tilde{\mathbf{s}}_i \cdot \mathbf{1}\|_1 \quad (10)$$

to penalise the scaling parameters \mathbf{s}_i (i.e. stretch of the ellipsoid) by its difference to the mean $\tilde{\mathbf{s}}_i$. As shown in Fig 3, this encourages sphericity, and avoids the problem of Gaussians which are highly elongated along the viewing direction creating artefacts. Let the union of the keyframes in the current window and the randomly selected one be $\mathcal{W} = \mathcal{W}_k \cup \mathcal{W}_r$. For mapping, we solve the following problem:

$$\min_{\mathbf{T}_{CW}^k \in \mathbf{SE}(3), \mathcal{G}, \forall k \in \mathcal{W}} \sum E_{pho}^k + \lambda_{iso} E_{iso}. \quad (11)$$

If depth observations are available, as in tracking, geometric residuals Eq. (8) are added to the optimisation problem.

4. Evaluation

We conduct a comprehensive evaluation of our system across a range of both real and synthetic datasets. Additionally, we perform an ablation study to justify our design

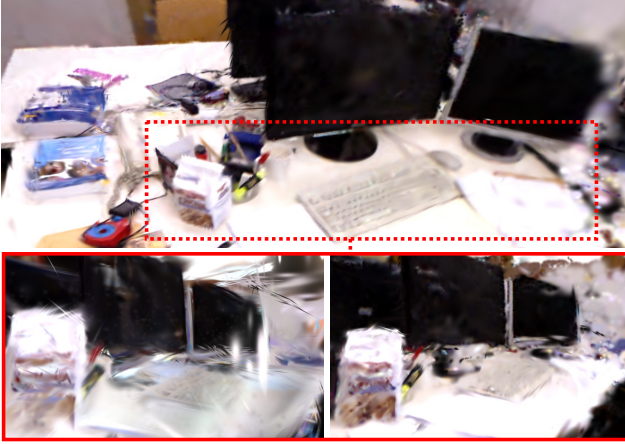


Figure 3. **Effect of isotropic regularisation:** **Top:** Rendering close to a training view (looking at the keyboard). **Bottom:** Rendering 3D Gaussians far from the training views (view from a side of the keyboard) without (left) and with (right) the isotropic loss. When the photometric constraints are insufficient, the Gaussians tend to elongate along the viewing direction, creating artefacts in the novel views, and affecting the camera tracking.

choices. Finally, we present qualitative results of our system operating live using a monocular camera, illustrating its practicality and high fidelity reconstruction.

4.1. Experimental Setup

Datasets For our quantitative analysis, we evaluate our method on the TUM RGB-D dataset [34] (3 sequences) and the Replica dataset [33] (8 sequences), following the evaluation in [35]. For qualitative results, we use self-captured real-world sequences recorded by Intel Realsense d455. Since the Replica dataset is designed for RGB-D SLAM evaluation, it contains challenging purely rotational camera motions. We hence use the Replica dataset for RGB-D evaluation only. The TUM RGB-D dataset is used for both monocular and RGB-D evaluation.

Implementation Details We run our SLAM on a desktop with Intel Core i9 12900K 3.50GHz and a single NVIDIA GeForce RTX 4090. We present results from our multi-process implementation aimed at real-time applications. For a fair comparison with other methods on Replica, we additionally report result for single-process implementation which performs more mapping iterations. As with 3DGS, time-critical rasterisation and gradient computation are implemented using CUDA. The rest of the SLAM pipeline is developed with PyTorch. Details of hyperparameters are provided in the supplementary material.

Metrics For camera tracking accuracy, we report the Root Mean Square Error (RMSE) of the Absolute Trajectory Error (ATE) of the keyframes. To evaluate map quality, we report standard photometric rendering quality metrics (PSNR,

SSIM and LPIPS) following the evaluation protocol used in [29]. To evaluate the map quality, on every fifth frame, rendering metrics are computed. We exclude the keyframes (training views). We report the average across three runs for all our evaluations. In the tables, the best result is in bold, and the second best is underlined.

Baseline Methods We primarily benchmark our SLAM method against other approaches that, like ours, do not have explicit loop closure. In monocular settings, we compare with state-of-the-art classical and learning-based direct visual odometry (VO) methods. Specifically, we compare DSO [5], DepthCov [4], and DROID-SLAM [38] in VO configurations. These methods are selected based on their public reporting of results on the benchmark (TUM dataset) or the availability of their source code for getting the benchmark result. Since one of our focuses is the online scale estimation under monocular scale ambiguity, the method which uses ground truth poses for the system initialisation such as [14] is not considered for the comparison. In the RGB-D case, we compare against neural-implicit SLAM methods [8, 9, 29, 35, 41, 45, 48] which are also map-centric, rendering-based and do not perform loop closure.

4.2. Quantitative Evaluation

Camera Tracking Accuracy Table 1 shows the tracking results on the TUM RGB-D dataset. In the monocular setting, our method surpasses other baselines without requiring any deep priors. Furthermore, our performance is comparable to systems which perform explicit loop closure. This clearly highlights that there still remains potential for enhancing the tracking of monocular SLAM by exploring fundamental SLAM representations.

Our RGB-D method shows better performance than any other baseline method. Notably, our system surpasses ORB-SLAM in the fr1 sequences, narrowing the gap between Map-centric SLAM and the state-of-the-art sparse frame-centric methods. Table 2 reports results on the synthetic Replica dataset. Our single-process implementation shows competitive performance and achieves the best result in 6 out of 8 sequences. Our multi-process implementation which performs fewer mapping iterations still performs comparably. In contrast to other methods, our system demonstrates higher performance on real-world data (TUM RGB-D), by optimising the Gaussian positions to compensate for the sensor noise.

Novel View Rendering Table 5 summarises the novel view rendering performance of our method with RGB-D input. We consistently show the best performance across most sequences and is least second best. Our rendering FPS is hundreds of times faster than other methods, offering a significant advantage for applications which require

Input	Loop-closure	Method	fr1/desk	fr2/xyz	fr3/office	Avg.
Monocular	w/o	DSO [5]	22.4	1.10	9.50	11.0
		DROID-VO [38]	<u>5.20</u>	10.7	<u>7.30</u>	<u>7.73</u>
		DepthCov-VO [4]	5.60	<u>1.20</u>	68.8	25.2
		Ours	3.78	4.60	3.50	3.96
	w/	DROID-SLAM [38]	1.80	0.50	2.80	1.70
RGB-D	w/o	ORB-SLAM2 [21]	1.90	0.60	2.40	1.60
		iMAP [35]	4.90	2.00	5.80	4.23
		NICE-SLAM [48]	4.26	6.19	3.87	4.77
		DI-Fusion [8]	4.40	2.00	5.80	4.07
		Vox-Fusion [45]	3.52	1.49	26.01	10.34
		ESLAM [9]	2.47	1.11	2.42	<u>2.00</u>
		Co-SLAM [41]	<u>2.40</u>	1.70	<u>2.40</u>	2.17
		Point-SLAM [29]	4.34	<u>1.31</u>	3.48	3.04
		Ours	1.50	1.44	1.49	1.47
	w/	BAD-SLAM [31]	1.70	1.10	1.70	1.50
		Kintinous [42]	3.70	2.90	3.00	3.20
		ORB-SLAM2 [21]	1.60	0.40	1.00	1.00

Table 1. **Camera tracking result on TUM for monocular and RGB-D.** ATE RMSE in cm is reported. In both monocular and RGB-D cases, we achieve state-of-the-art performance. In particular, in the monocular case, not only do we outperform systems which use deep prior, but we achieve comparable performance with many of the RGB-D systems.

Method	r0	r1	r2	o0	o1	o2	o3	o4	Avg.
iMAP [35]	3.12	2.54	2.31	1.69	1.03	3.99	4.05	1.93	2.58
NICE-SLAM	0.97	1.31	1.07	0.88	1.00	1.06	1.10	1.13	1.07
Vox-Fusion [45]	1.37	4.70	1.47	8.48	2.04	2.58	1.11	2.94	3.09
ESLAM [9]	0.71	0.70	0.52	0.57	0.55	0.58	0.72	0.63	0.63
Point-SLAM [29]	0.61	0.41	0.37	<u>0.38</u>	<u>0.48</u>	0.54	0.69	<u>0.72</u>	<u>0.53</u>
Ours	<u>0.44</u>	<u>0.32</u>	<u>0.31</u>	0.44	0.52	0.23	<u>0.17</u>	2.25	0.58
Ours (sp)	0.33	0.22	0.29	0.36	0.19	<u>0.25</u>	0.12	0.81	0.32

Table 2. **Camera tracking result on Replica for RGB-D SLAM.** ATE RMSE in cm is reported. We achieve best performance across most sequences. Here, Ours is our multi-process implementation and Ours (sp) is the single-process implementation which ensures a certain amount of mapping iteration similar to other works.

Input	Method	fr1/desk	fr2/xyz	fr3/office	Avg.
Mono	w/o E_{iso}	4.16	4.66	5.73	4.83
	w/o kf selection	13.2	4.36	8.65	8.73
	Ours	3.78	4.60	3.50	3.96
RGB-D	w/o E_{geo}	2.39	0.62	4.98	2.66
	w/o kf selection	1.64	1.49	2.60	1.90
	Ours	1.50	1.44	1.49	1.47

Table 3. **Ablation Study on TUM RGB-D dataset.** We analyse the usefulness of isotropic regularisation, geometric residual, and keyframe selection to our SLAM system. Further isotropic regularisation ablation is available in supplementary.

Memory Usage [MB]				
iMAP [35]	NICE-SLAM [48]	Co-SLAM [41]	Ours (Mono)	Ours (RGB-D)
0.8MB	40.3.4MB	6.4MB	<u>2.6MB</u>	3.97MB

Table 4. **Memory Analysis on TUM RGB-D dataset.** The baseline numbers are computed from the parameter numbers in [41]

real-time map interaction. While Point-SLAM is competitive, that method focuses on view synthesis rather than novel-view synthesis. Their view synthesis is conditional on the availability of depth due to the depth-guided ray-sampling, making novel-view synthesis challenging. On the other hand, our rasterisation-based approach does not require depth guidance and achieves efficient, high-quality,

Method	PSNR[db]↑	SSIM↑	LPIPS↓	Rendering FPS
NICE-SLAM[48]	24.42	0.809	0.233	0.54
Vox-Fusion[45]	24.41	0.801	0.236	<u>2.17</u>
Point-SLAM [29]	<u>35.17</u>	0.975	0.124	1.33
ours	38.94	<u>0.968</u>	0.070	769

Table 5. **Average rendering performance on Replica (RGB-D).** Our method outperforms most of the rendering metrics compared to existing methods. Note that Point-SLAM uses ground-truth depth to guide sampling along rays. The full detail is available in supplementary.

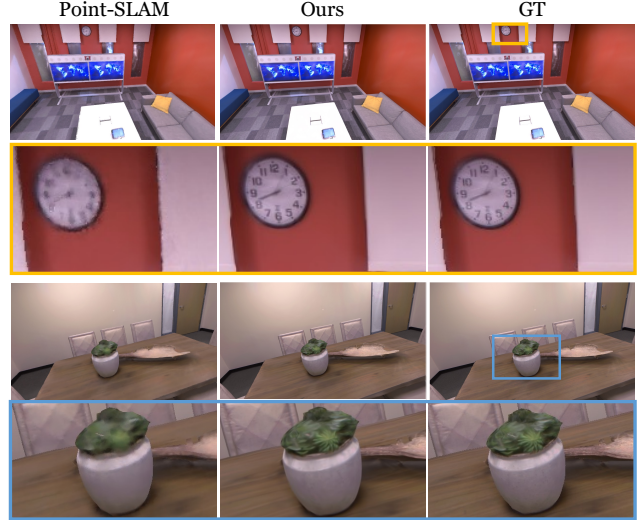


Figure 4. **Rendering examples on Replica.** Point-SLAM struggle with rendering fine details due to the stochastic ray sampling.

novel view synthesis. Fig. 4 provides a qualitative comparison of the rendering of ours and Point-SLAM (with depth guidance).

Ablative Analysis In Table 3, we perform ablation to confirm our design choices. Isotropic regularisation and geometric residual improve the tracking of monocular and RGB-D SLAM respectively, as they aid in constraining the geometry when photometric signals are weak. For both cases, keyframe selection significantly improves systems performance, as it automatically chooses suitable keyframes based on our occlusion-aware keyframe selection and management. We further compare the memory usage of different 3D representations in Table 4. MLP-based iMAP is clearly more memory efficient, but it struggles to express high-fidelity 3D scenes due to the limited capacity of small MLP. Compared with a voxel grid of features used in NICE-SLAM, our method uses significantly less memory.

Convergence Basin Analysis In our SLAM experiments, we discovered that 3D Gaussian maps have a notably large convergence basin for camera localisation. To investigate further, we conducted a convergence funnel analysis, an evaluation methodology proposed in [19] and used in [23]. Here, we train a 3D representation (e.g. 3DGS) using 9

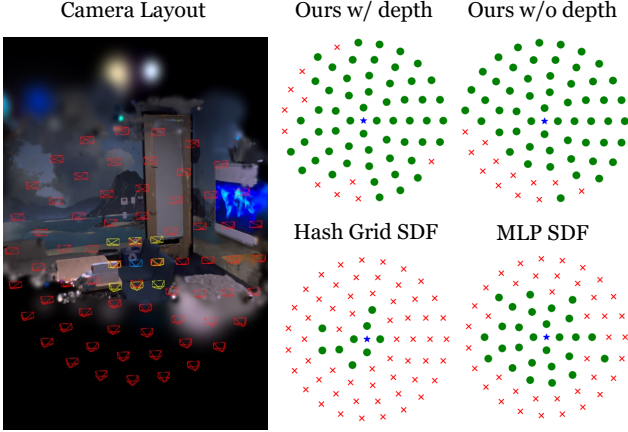


Figure 5. **Convergence basin analysis:** **Left:** 3D Gaussian map from training views (Yellow) and visualisation of the test poses (Red) and target pose (Blue). **Right:** Convergence basin of our method. The green marks success, and the red marks failure.

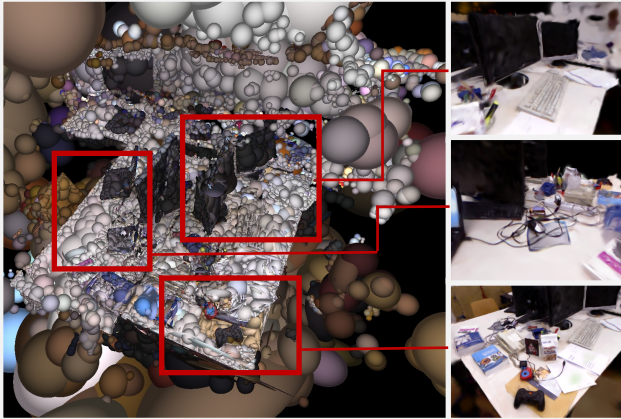


Figure 6. **Monocular SLAM result on fr1/desk sequence:** We show the reconstructed 3D Gaussian maps (Left) and novel view synthesis result (Right).

Method	seq1	seq2	seq3	Avg.
Neural SDF (Hash Grid)	0.13	0.15	0.16	0.14
Neural SDF (MLP)	0.40	0.38	0.22	0.33
Ours w/o depth	<u>0.82</u>	<u>0.91</u>	0.65	<u>0.79</u>
Ours w/ depth	0.83	1.0	0.65	0.82

Table 6. **Camera convergence analysis.** We report the ratio of successful camera convergence for the different sequences, across different differentiable 3D representations.

fixed views arranged in a square. We set the viewpoint in the middle of the square to be the target view. As shown in Fig 5, we uniformly sample a position, creating a funnel. From the sampled position, given the RGB image of the target view, we perform camera pose optimisation for 1000 iterations. The optimisation is successful if it converges to within 1cm of the target view within the fixed iterations. We compare our Gaussian approach with Co-SLAM [41]’s network (Hash Grid SDF) and iMAP’s [35] network with Co-SLAM’s SDF loss for further geometric accuracy (MLP



Figure 7. **Self-captured Scenes:** Challenging scenes and objects, for example, transparent glasses and crinkled texture of salad are captured by our monocular SLAM running live.

Neural SDF). We render the training views using a synthetic Replica dataset and create three sequences for testing (seq1, seq2 and seq3). The width of the square formed by the training view is 0.5m, and the test cameras are distributed with radii ranging from 0.2m to 1.2m, covering a larger area than the training view. When training the map, the three methods—Ours w/depth, Hash Grid SDF, and MLP SDF—use RGB-D images, whereas Ours w/o depth utilises only colour images. Fig. 5 shows the qualitative results and Table 6 reports the success rate. For both with and without depth for training, our method shows better convergence. Unlike hashing and positional encoding which can lead to signal conflict, anisotropic Gaussians form a smooth gradient in 3D space, increasing the convergence basin. Further experimental details are available in the supplementary.

4.3. Qualitative Results

We report both the 3D reconstruction of the SLAM dataset and self-captured sequences. In Fig. 6, we visualise the monocular SLAM reconstruction of fr1/desk. The placements of the Gaussians are geometrically sensible and are 3D coherent, and our rendering from the different viewpoints highlights the quality of our systems’ novel view synthesis. In Fig. 7, we self-capture challenging scenes for monocular SLAM. By not explicitly modelling a surface, our system naturally handles transparent objects which is challenging for many other SLAM systems.

5. Conclusion

We have proposed the first SLAM method using 3D Gaussians as a SLAM representation. Via efficient volume rendering, our system significantly advances the fidelity and diversity of object materials a live SLAM system can capture. Our system achieves state-of-the-art performance across benchmarks for both monocular and RGB-D cases. Inter-

esting directions for future research are the integration of loop closure for handling large-scale scenes and extraction of geometry such as surface normal as Gaussians do not explicitly represent the surface.

6. Acknowledgement

Research presented in this paper has been supported by Dyson Technology Ltd. We are very grateful to Eric Dexheimer, Kirill Mazur, Xin Kong, Marwan Taher, Ignacio Alzugaray, Gwangbin Bae, Aalok Patwardhan, and members of the Dyson Robotics Lab for their advice and insightful discussions.

Supplementary Material

7. Implementation Details

7.1. System Details and Hyperparameters

7.1.1 Tracking and Mapping (Sec. 3.3.1 and 3.3.3)

Learning Rates We use the Adam optimiser for both camera poses and Gaussian parameters optimisation. For camera poses, we used 0.003 for rotation and 0.001 for translation. For 3D Gaussians, we used the default learning parameters of the original Gaussian Splatting implementation [11], apart from in monocular setting where we increase the learning rate of the positions of the Gaussians μ_W by a factor of 10.

Iteration numbers 100 tracking iterations are performed per frame for across all experiments. However, we terminate the iterations early if the magnitude of the pose update becomes less than 10^{-4} . For mapping, 150 iterations are used for the single-process implementation.

Loss Weights Given a depth observation, for tracking we minimise both photometric Eq. (7) and geometric residual Eq. (8) as:

$$\min_{T_{CW} \in \mathbf{SE}(3)} \lambda_{pho} E_{pho} + (1 - \lambda_{pho}) E_{geo}, \quad (12)$$

and similarly, for mapping we modify Eq. (11) to:

$$\min_{\substack{T_{CW} \in \mathbf{SE}(3), \mathcal{G}, \\ \forall k \in \mathcal{W}}} \sum (\lambda_{pho} E_{pho}^k + (1 - \lambda_{pho}) E_{geo}^k) + \lambda_{iso} E_{iso}. \quad (13)$$

We set $\lambda_{pho} = 0.9$ for all RGB-D experiments, and $\lambda_{iso} = 10$ for both monocular and RGB-D experiments.

7.1.2 Keyframing (Sec. 3.3.2)

Gaussian Covisibility Check (Sec. 3.3.2) As described in Sec. 3.3.2, keyframe selection is based on the covisibility of the Gaussians. Between two keyframes i, j , we define the covisibility using the Intersection of Union (IOU) and Overlap Coefficient (OC):

$$IOU_{cov}(i, j) = \frac{|\mathcal{G}_i^v \cap \mathcal{G}_j^v|}{|\mathcal{G}_i^v \cup \mathcal{G}_j^v|}, \quad (14)$$

$$OC_{cov}(i, j) = \frac{|\mathcal{G}_i^v \cap \mathcal{G}_j^v|}{\min(|\mathcal{G}_i^v|, |\mathcal{G}_j^v|)}, \quad (15)$$

where \mathcal{G}_i^v is the Gaussians visible in keyframe i , based on visibility check described in Section 3.3.2, Gaussian Covisibility. A keyframe i is added to the keyframe window \mathcal{W}_k if

given last keyframe j , $IOU_{cov}(i, j) < kf_{cov}$ or if the relative translation $t_{ij} > kf_m \hat{D}_i$, where \hat{D}_i is the median depth of frame i . For Replica $kf_{cov} = 0.95, kf_m = 0.04$ and for TUM $kf_{cov} = 0.90, kf_m = 0.08$. We remove the registered keyframe j in \mathcal{W}_k if the $OC_{cov}(i, j) < kf_c$, where keyframe i is the latest added keyframe. For both Replica and TUM, we set the cutoff to $kf_c = 0.3$. We set the size of the keyframe window to be for Replica, $|\mathcal{W}_k| = 10$, and for TUM, $|\mathcal{W}_k| = 8$.

Gaussian Insertion and Pruning (Sec. 3.3.2) As we optimise the positions of Gaussians and prune geometrically unstable Gaussians, we do not require any strong prior such as depth observation for Gaussian initialisation. When **inserting** new Gaussians in a monocular setting, we randomly sample the Gaussians position μ_W using rendered depth D . Since the estimated depth may sometimes be incorrect, we account for this by initialising the Gaussians with some variance. For a pixel p where the rendered depth \mathcal{D}_p exists, we sample the depth from $\mathcal{N}(\mathcal{D}_p, 0.2\sigma_D)$. Otherwise, for unobserved regions, we initialise the Gaussians by sampling from $\mathcal{N}(\hat{D}, 0.5\sigma_D)$, where \hat{D} is the median of D . For **pruning**, as described in Section 3.3.2, we perform visibility-based pruning, where if new Gaussians inserted within the last 3 keyframes are not observed by at least 3 other frames, they are pruned. We only perform visibility-based pruning once the keyframe window \mathcal{W}_k is full. Additionally, we prune all Gaussians with opacity of less than 0.7.

8. Evaluation details

8.1. Camera Tracking Accuracy (Table 1 and Table 2)

8.1.1 Evaluation Metric

We measured the keyframe absolute trajectory error (ATE) RMSE. For monocular evaluation, we perform scale alignment between the estimated scale-free and ground-truth trajectories. For RGB-D evaluation, we only align the estimated trajectory and ground truth without scale adjustment.

8.1.2 Baseline Results

Table 1 Numbers for monocular DROID-SLAM [38] and ORB-SLAM [21] is taken from [14]. We have locally run DSO [5], DepthCov [4] and DROID-VO [38] – which is DROID-SLAM without loop closure and global bundle adjustment. For the RGB-D case, numbers for NICE-SLAM [48], DI-Fusion [8], Vox-Fusion [45], Point-SLAM [29] are taken from Point-SLAM [29], and numbers for iMAP [35], BAD-SLAM [31], Kintinous [42], ORB-SLAM [21] are from iMAP [35], and all the other base-

lines: ESLAM [9], Co-SLAM [41] are from each individual papers.

Table 2 and 5 We took the numbers from Point-SLAM [29] paper.

Table 4 The numbers are from Co-SLAM [41] paper.

8.2. Rendering Performance (Table 5)

We provide the full detail of the rendering performance evaluation in Table 7.

In Table 5, we reported the photometric quality metrics (PSNR, SSIM and LPIPS) and rendering fps of our methods. We demonstrated that our rendering fps (769) is much higher than other existing methods (VoxFusion is the second best with 2.17fps). Here we describe the detail of how we measured the fps. The rendering time refers to the duration necessary for full-resolution rendering (1200×680 for the Replica sequence). For each method, we perform 100 renderings and report the average time taken per rendering. The reported rendering fps is found by taking 1 and dividing it by the average rendering time. We summarise the numbers in Table 8. Note that the “rendering fps” means the fps just for the forward rendering, which differs from the end-to-end system fps reported in Table 9 and 10.

8.3. The convergence basin analysis (Table 6 and Fig 5)

8.3.1 The detail of the benchmark Dataset

For convergence basin analysis, we create three datasets by rendering the synthetic Replica dataset. In addition to the qualitative visualisation in Figure 5, we report more detailed camera pose distributions in Figure 8. Figure 8 shows the camera view frustums of the test (red), training (yellow) and target (blue) views. As we mentioned in the main paper, we set the training view in the shape of a square with a width of 0.5m and test views are distributed with radii ranging from 0.2m to 1.2m, covering a larger area than the training views. We only apply displacements to the camera translation but not to the rotation. For each sequence, we use a total of 67 test views.

8.3.2 Training setup

For each method, the 3D representation is trained for 30000 iterations using the training views. Here, we detail the training setup of each of the methods:

Ours We evaluated our method under two settings: “w/ depth” and “w/o depth”, where we train the initial 3D Gaussian map \mathcal{G}_{init} with and without depth supervision. In the

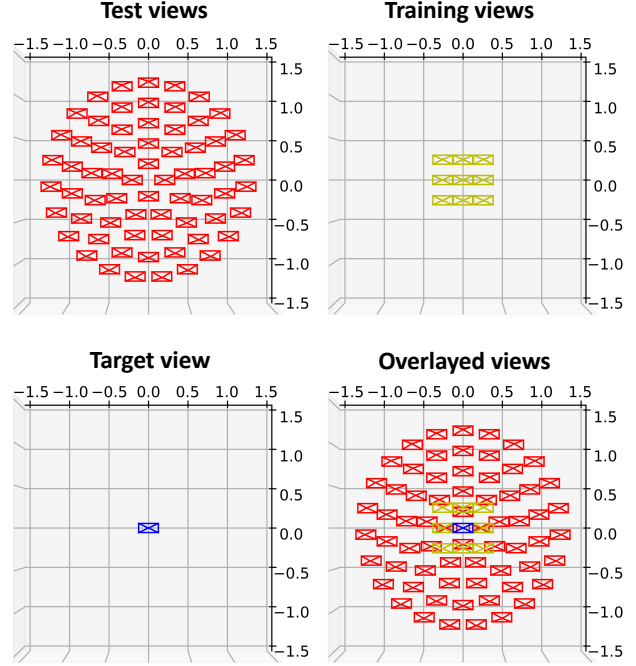


Figure 8. **2D Visualisation of the camera pose distributions used for convergence basin analysis in Figure 5.**

“w/o depth” setting, the 3D Gaussians’ positions are randomly initialised, and we minimise the monocular mapping cost Eq. (11) for the 3D Gaussian training, but keeping the camera poses fixed. Specifically, let $k \in \mathbb{N}$ be a number of training views and 3D Gaussians \mathcal{G} , we find \mathcal{G}_{init} by:

$$\mathcal{G}_{init} = \arg \min_{\mathcal{G}} \sum_{\forall k \in \mathcal{W}} E_{pho}^k + \lambda_{iso} E_{iso}. \quad (16)$$

Note that training views’ camera poses T_{CW}^k are fixed during the optimisation.

In the “w/ depth” setting, we train the Gaussian map by minimising the same cost function as our RGB-D SLAM system:

$$\mathcal{G}_{init} = \arg \min_{\mathcal{G}} \sum_{\forall k \in \mathcal{W}} (\lambda_{pho} E_{pho}^k + (1 - \lambda_{pho}) E_{geo}^k) + \lambda_{iso} E_{iso}, \quad (17)$$

where we use $\lambda_{pho} = 0.9$ and $\lambda_{iso} = 10$ for all the experiments

Baseline Methods For Hash Grid SDF, we trained the same network architecture as Co-SLAM [41]. For MLP SDF, we trained the network of iMAP [35]. For both baselines, we supervised networks with the same loss functions as Co-SLAM, which are colour rendering loss L_{rgb} , depth rendering loss L_{depth} , SDF loss L_{fs} , free-space loss L_{fs} , and smoothness loss L_{smooth} . Please refer to the original

Method	Metric	room0	room1	room2	office0	office1	office2	office3	office4	Avg.	Rendering FPS
NICE-SLAM [48]	PSNR[dB] ↑	22.12	22.47	24.52	29.07	30.34	19.66	22.23	24.94	24.42	0.54
	SSIM ↑	0.689	0.757	0.814	0.874	0.886	0.797	0.801	0.856	0.809	
	LPIPS↓	0.33	0.271	0.208	0.229	0.181	0.235	0.209	0.198	0.233	
Vox-Fusion [45]	PSNR[dB] ↑	22.39	22.36	23.92	27.79	29.83	20.33	23.47	25.21	24.41	<u>2.17</u>
	SSIM ↑	0.683	0.751	0.798	0.857	0.876	0.794	0.803	0.847	0.801	
	LPIPS↓	0.303	0.269	0.234	0.241	0.184	0.243	0.213	0.199	0.236	
Point-SLAM [29]	PSNR[dB] ↑	<u>32.40</u>	<u>34.08</u>	<u>35.5</u>	<u>38.26</u>	<u>39.16</u>	<u>33.99</u>	<u>33.48</u>	<u>33.49</u>	<u>35.17</u>	1.33
	SSIM ↑	0.974	0.977	0.982	0.983	0.986	<u>0.96</u>	<u>0.960</u>	0.979	0.975	
	LPIPS↓	<u>0.113</u>	<u>0.116</u>	<u>0.111</u>	<u>0.1</u>	<u>0.118</u>	<u>0.156</u>	<u>0.132</u>	<u>0.142</u>	<u>0.124</u>	
Ours	PSNR[dB] ↑	34.83	36.43	37.49	39.95	42.09	36.24	36.7	36.07	37.50	769
	SSIM ↑	<u>0.954</u>	<u>0.959</u>	<u>0.965</u>	<u>0.971</u>	<u>0.977</u>	0.964	0.963	<u>0.957</u>	<u>0.960</u>	
	LPIPS↓	0.068	0.076	0.075	0.072	0.055	0.078	0.065	0.099	0.070	

Table 7. **Rendering performance comparison of RGB-D SLAM methods on Replica.** Our method outperforms most of the rendering metrics compared to existing methods. Note that Point-SLAM uses sensor depth (ground-truth depth in Replica) to guide sampling along rays, which limits the rendering performance to existing views. The numbers for the baselines are taken from [29].

Method	Rendering FPS ↑	Rendering time per image [s] ↓
NICE-SLAM [48]	0.54	1.85
Vox-Fusion [45]	<u>2.17</u>	<u>0.46</u>
Point-SLAM [29]	1.33	0.75
Ours	769	0.0013

Table 8. **Further detail of Rendering FPS and Rendering Time comparison based on Table 5.**

Co-SLAM paper for the exact formulation (equation (6) - (9)). All the training hyperparameters (e.g. learning rate of the network, number of sampling points, loss weight) are the same as Co-SLAM’s default configuration of the Replica dataset. While Co-SLAM stores training view information by downsampling the colour and depth images, we store the full pixel information because the number of training views is small.

8.3.3 Testing Setup

For testing, we localise the camera pose by minimising only the photometric error against the ground-truth colour image of the target view.

Ours Let the camera pose $T_{CW} \in SE(3)$ and initial 3D Gaussians \mathcal{G}_{init} , the localised camera pose T_{CW}^{est} is found by:

$$T_{CW}^{est} = \arg \min_{T_{CW}} \|I(\mathcal{G}_{init}, T_{CW}) - \bar{I}_{target}\|_1. \quad (18)$$

Note that \mathcal{G}_{init} is fixed during the optimisation. We initialise T_{CW} at one of the test view’s positions, and optimisation is performed for 1000 iterations. We perform this localisation process for all the test views and measure the success rate. Camera localisation is successful if the estimated pose converges to within 1cm of the target view within the 1000 iterations.

Method	Total Time [s]	FPS
Monocular	798.9	3.2
RGB-D	986.7	2.5

Table 9. **Performance Analysis using fr3/office.** Both monocular and RGB-D implementations use multiprocessing. We report **the total execution time of our system**, FPS computed by dividing the total number of processed frames by the total time.

Method	Total Time [s]	FPS
RGB-D	1111.1	1.8
RGB-D (sp)	1904.7	1.1

Table 10. **Performance Analysis using replica/office1.** RGB-D uses a multi-process implementation and RGB-D-sp is the single-process implementation. We report **the total execution time of our system**, FPS computed by dividing the total number of processed frames by the total time.

Baseline Methods For the baseline methods, the camera localisation is performed by minimising colour volume rendering loss L_{rgb} , while all the other trainable network parameters are fixed. The learning rates of the pose optimiser are also the same as Co-SLAM’s default configuration of Replica dataset.

9. Further Ablation Analysis (Table 3)

9.1. Pruning Ablation (Monocular input)

In Table 9.1, we report the ablation study of our proposed Gaussian pruning, which prunes randomly initialised 3D Gaussians effectively in a monocular SLAM setting. As the result shows, Gaussian pruning plays a significant role in enhancing camera tracking performance. This improvement is primarily because, without pruning, randomly initialised Gaussians persist in the 3D space, potentially leading to incorrect initial geometry for other views.

Input	Method	fr1/desk	fr2/xyz	fr3/office	Avg.
Mono	w/o pruning	78.2	4.5	57.0	46.6
	Ours	3.78	4.60	3.50	3.96

Table 11. **Pruning Ablation Study on TUM RGB-D dataset (Monocular Input)**. Numbers are camera tracking error (ATE RMSE) in cm.

Input	Method	fr1/desk	fr2/xyz	fr3/office	Avg.
RGB-D	w/o E_{iso}	1.60	1.42	1.32	1.43
	Ours	1.50	1.44	1.49	1.47

Table 12. **Isotropic Loss Ablation Study on TUM RGB-D dataset (RGB-D input)**. Numbers are camera tracking error (ATE RMSE) in cm.

Method	r0	r1	r2	o0	o1	o2	o3	o4	Avg.
w/o E_{iso}	0.44	0.86	0.28	0.75	0.99	0.36	0.28	2.6	0.82
Ours	0.44	0.32	0.31	0.44	0.52	0.23	0.17	2.25	0.58

Table 13. **Isotropic Loss Ablation Study on Replica dataset (RGB-D input)**. Numbers are camera tracking error (ATE RMSE) in cm.

9.2. Isotropic Loss Ablation (RGB-D input)

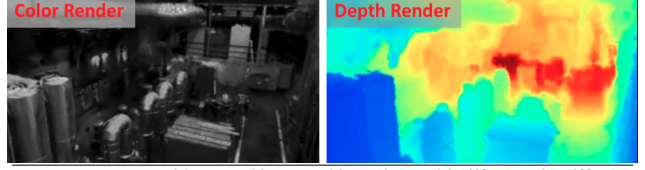
Table 12 and 13 report the ablation study of the effect of isotropic loss E_{iso} for RGB-D input. In TUM, as Table 12 shows, isotropic regularisation does not improve the performance but only shows a marginal difference. However, for Replica, as summarised in Table 13, isotropic loss significantly improves camera tracking performance. Even with the depth measurement, since rasterisation does not consider the elongation along the viewing axis. Isotropic regularisation is required to prevent the Gaussians from over-stretching, especially for textureless regions, which are common in Replica.

9.3. Effect of Spherical Harmonics (SH)

While we disabled SHs in the main paper for simplicity, here we report the ablation study of the effect of SHs. The 3DGS paper [11] shows that addition of SH leads to small improvements in rendering metrics, and we have found similar improvement with SH enabled in our system (Tab.15a). We did not observe a significant change in runtime with SH enabled, but it notably increases Gaussian map size and hence GPU memory usage. Though an analytical Jacobian propagates the gradients from SH to camera poses, ATE marginally gets worse when SH is enabled (Tab. 16), as SH may incorrectly explain non-view directional effects caused by the camera motion, degrading the trajectory estimate.

9.4. Mapping Performance with ORB-SLAM

One of the most straightforward approaches for real-time operation is to combine an existing tracking system and 3DGS. In particular, frame-based SLAM methods have been well-studied for years and is capable of providing reliable tracking. In this section, we compare our unified 3DGS-based method to the combined approach. We have



	01-easy	02-easy	03-medium	04-difficult	05-difficult
Point-SLAM [29]	-	-	-	-	-
Ours	<u>0.121</u>	<u>0.141</u>	2.197	4.515	3.190
Vins-Fusion [28]	0.540	0.460	<u>0.330</u>	<u>0.780</u>	<u>0.500</u>
SVO [6]	<u>0.040</u>	<u>0.070</u>	<u>0.270</u>	<u>0.170</u>	<u>0.120</u>
ORB-SLAM3 [1]	0.029	0.019	0.024	0.085	0.052

Table 14. ATE RMSE (meter) on EuRoC Machine Hall with Stereo Depth. Baseline numbers of classical methods are from [1]. The third best result is highlighted with a dash line.

		TUM			Replica		
	Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
	Ours (w/o SH)	21.89	0.733	0.327	38.94	0.968	0.0703
(a)	Ours (w. SH)	24.37	0.804	0.225	-	-	-
	Point-SLAM	21.39	0.727	0.463	24.37	0.840	0.185
(b)	ORB+GS (w/o SH)	25.12	0.837	0.161	37.11	0.964	0.040
	ORB+GS (w.SH)	25.44	0.842	0.146	-	-	-

Table 15. Mean Rendering metrics for TUM and Replica (RGBD).

Memory Usage for RGB-D SLAM					ATE RMSE	
Ours (w/o SH)	Ours (w. SH)	Point-SLAM	ORB+GS (w/o SH)	ORB+GS (w. SH)	Ours (w/o SH)	Ours (w. SH)
3.97MB	11.47MB	38.0MB	45.97MB	186.5MB	1.47cm	1.56cm

Table 16. Mean Memory and ATE metrics for TUM (RGBD).

run RGB-D ORB-SLAM to recover the poses and train 3DGS with the poses and sensor depth of the keyframes, equivalent to performing offline splatting. Though ORB-SLAM is best in terms of ATE (Tab.1 main), we find no significant difference across the rendering metrics (Tab.15b). SH is omitted in the synthetic Replica dataset as it contains no view-directional effects. While using an off-the-shelf tracker with a 3DGS mapper is possible, this work has focused on the value of the 3DGS throughout the entire algorithms, which is unexplored and therefore provides new insights. Further performance improvement of the unified approach will be an interesting future work.

9.5. Large-scale Scenes with Stereo Inputs:

This work focuses on pioneering 3DGS-based SLAM for live operation in small-scale scenes. However, we tested our method on the large-scale EuRoC Machine Hall dataset with depth from stereo (Tab.14). Fig.1 is a qualitative reconstruction result from our system. Our method is competitive in “easy” sequences, although performance drops for more difficult, longer sequences. Note that Point-SLAM [29] fails on all sequences in this dataset. In future work, we expect to improve our method by incorporating loop closure. In principle, loop closure will be easier to incorporate compared to other representations such as voxel grids (where feature allocations are fixed), via a method similar to surfel-based approaches like ElasticFusion [43].

9.6. Memory Consumption and Frame Rate (Table. 4)

9.6.1 Memory Analysis

In memory consumption analysis, for Table. 4, we measure the final size of the created Gaussians. The memory footprint of our system is lower than the original Gaussian Splatting, which uses approximately 300-700MB for the standard novel view synthesis benchmark dataset [11], as we only maintain well-constrained Gaussians via pruning and do not store the spherical harmonics.

9.6.2 Timing Analysis

To analyse the processing time of our monocular/RGB-D SLAM system, we measure the total time required to process all frames in the TUM-RGBD fr3/office dataset. This approach assesses the performance of our system as a whole, rather than isolating individual components. By adopting this approach, we gain a more realistic understanding of the system's true performance which better reflects the real-world operating conditions, as it avoids the assumption of an idealised, sequential interleaving of the tracking and mapping processes. As shown in Table 10, our system operates at 3.2 FPS with monocular and 2.5 FPS with depth. The FPS is found by dividing the number of processed frames by the total time. We conducted a similar analysis with the Replica dataset office2. Here, we compare the RGB-D method with and without multiprocessing. As expected, single-process implementation takes longer as it performs more mapping iterations.

10. Camera Pose Jacobian

Use of 3D Gaussian as a primitive and performing camera pose optimisation is discussed in [12]; however, the method assumes a smaller number of Gaussians and is based on ray-intersection not splatting; hence, is not applicable to 3DGS. While many applications of 3DGS exist, for example, dynamic tracking and 4D scene representation [16, 44], they assume offline application and require accurate camera position. In contrast, we perform camera pose optimisation by deriving the minimal analytical Jacobians on Lie group, and for completeness, we provide the derivation of the Ja-

cobians presented in Eq. (6).

$$\frac{\mathcal{D}\mu_C}{\mathcal{D}T_{CW}} = \lim_{\tau \rightarrow 0} \frac{\text{Exp}(\tau) \cdot \mu_C - \mu_C}{\tau} \quad (19)$$

$$= \lim_{\tau \rightarrow 0} \frac{(\mathbf{I} + \tau^\wedge) \cdot \mu_C - \mu_C}{\tau} \quad (20)$$

$$= \lim_{\tau \rightarrow 0} \frac{\tau^\wedge \cdot \mu_C}{\tau} \quad (21)$$

$$= \lim_{\tau \rightarrow 0} \frac{\theta^\times \mu_C + \rho}{\tau} \quad (22)$$

$$= \lim_{\tau \rightarrow 0} \frac{-\mu_C^\times \theta + \rho}{\tau} \quad (23)$$

$$= [\mathbf{I} \quad -\mu_C^\times] \quad (24)$$

where $T \cdot \mathbf{x}$ is the group action of $T \in SE(3)$ on $\mathbf{x} \in \mathbb{R}^3$.

Similarly, we compute the Jacobian with respect to \mathbf{W} . Since the translational component is not involved, we only consider the rotational part R_{CW} of T_{CW} .

$$\frac{\mathcal{D}\mathbf{W}}{\mathcal{D}R_{CW}} = \lim_{\theta \rightarrow 0} \frac{\text{Exp}(\theta) \circ \mathbf{W} - \mathbf{W}}{\theta} \quad (25)$$

$$= \lim_{\theta \rightarrow 0} \frac{(\mathbf{I} + \theta^\wedge) \circ \mathbf{W} - \mathbf{W}}{\theta} \quad (26)$$

$$= \lim_{\theta \rightarrow 0} \frac{\theta^\wedge}{\theta} \circ \mathbf{W} \quad (27)$$

$$= \lim_{\theta \rightarrow 0} \frac{\theta^\times}{\theta} \circ \mathbf{W} \quad (28)$$

Since skew symmetric matrix is:

$$\theta^\times = \begin{bmatrix} 0 & -\theta_z & \theta_y \\ \theta_z & 0 & -\theta_x \\ -\theta_y & \theta_x & 0 \end{bmatrix} \quad (29)$$

The partial derivative of one of the component (e.g. θ_x) is:

$$\frac{\partial \theta^\times}{\partial \theta_x} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & -1 \\ 0 & 1 & 0 \end{bmatrix} = \mathbf{e}_1^\times \quad (30)$$

where $\mathbf{e}_1 = [1, 0, 0]^\top$, $\mathbf{e}_2 = [0, 1, 0]^\top$, $\mathbf{e}_3 = [0, 0, 1]^\top$.

$$\frac{\partial \mathbf{W}}{\partial \theta_x} = \mathbf{e}_1^\times \mathbf{W} = \begin{bmatrix} \mathbf{0}_{1 \times 3} \\ -\mathbf{W}_{3,:} \\ \mathbf{W}_{2,:} \end{bmatrix} \quad (31)$$

$$\frac{\partial \mathbf{W}}{\partial \theta_y} = \mathbf{e}_2^\times \mathbf{W} = \begin{bmatrix} \mathbf{W}_{3,:} \\ \mathbf{0}_{1 \times 3} \\ -\mathbf{W}_{1,:} \end{bmatrix} \quad (32)$$

$$\frac{\partial \mathbf{W}}{\partial \theta_z} = \mathbf{e}_3^\times \mathbf{W} = \begin{bmatrix} -\mathbf{W}_{2,:} \\ \mathbf{W}_{1,:} \\ \mathbf{0}_{1 \times 3} \end{bmatrix} \quad (33)$$

where $\mathbf{W}_{i,:}$ refers to the i th row of the matrix. After column-wise vectorisation of Eq. (31), (32), (33), and stacking horizontally we get:

$$\frac{\mathcal{D}\mathbf{W}}{\mathcal{D}\mathbf{R}_{CW}} = \begin{bmatrix} -\mathbf{W}_{:,1}^\times \\ -\mathbf{W}_{:,2}^\times \\ -\mathbf{W}_{:,3}^\times \end{bmatrix}, \quad (34)$$

where $\mathbf{W}_{:,i}$ refers to the i th column of the matrix. Since the translational part is all zeros, with this we get Eq. (6).

11. Additional Qualitative Results

We urge readers to view our supplementary video for convincing qualitative results. In Fig. 9 - Fig. 16, we further show additional qualitative results. We visually compare other state-of-the-art SLAM methods using differentiable rendering (Point-SLAM [29] and ESLAM [9]).

12. Limitation of this work

Although our novel Gaussian Splatting SLAM shows competitive performance on experimental results, the method also has several limitations.

- Currently, the proposed method is tested only on small room-scale scenes. For larger real-world scenes, the trajectory drift is inevitable. This could be addressed by integrating a loop closure module into our existing pipeline.
- Although we achieve interactive live operation, hard real-time operation on the benchmark dataset (30 fps on TUM sequences) is not achieved in this work. To improve speed, exploring a second-order optimiser would be an interesting direction.

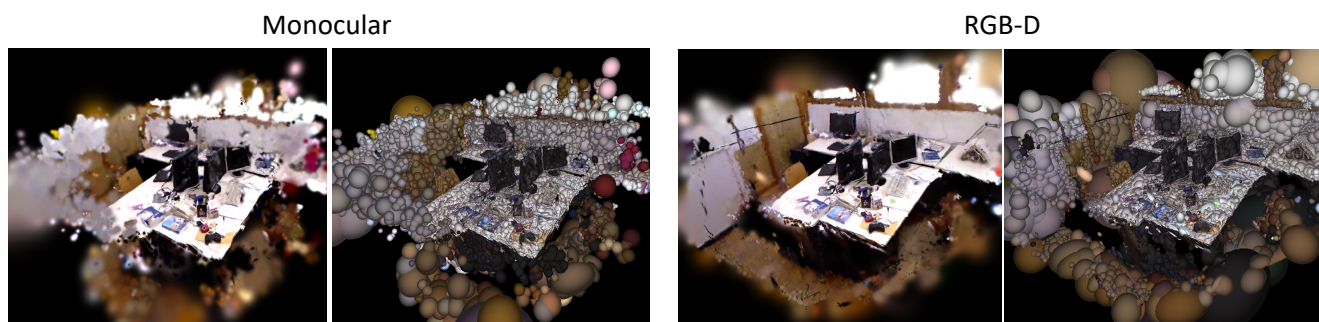


Figure 9. Novel view rendering and Gaussian visualizations on TUM fr1/desk



Figure 10. Rendering comparison on TUM fr1/desk

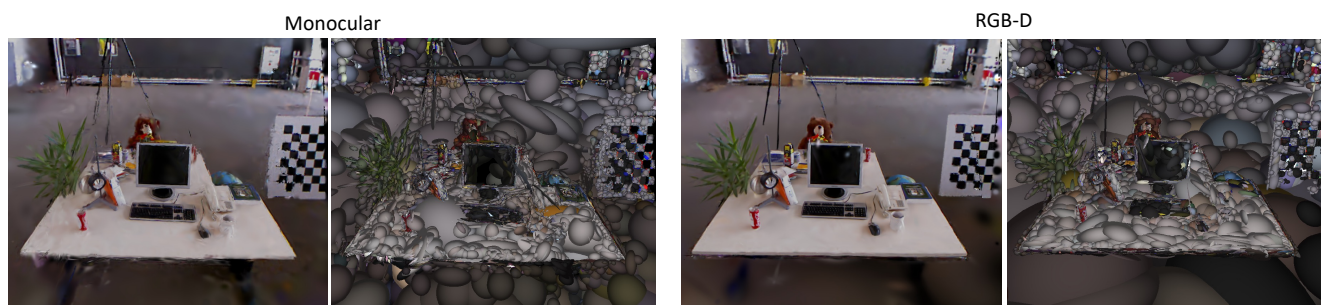


Figure 11. Novel view rendering and Gaussian visualizations on TUM fr2/xyz

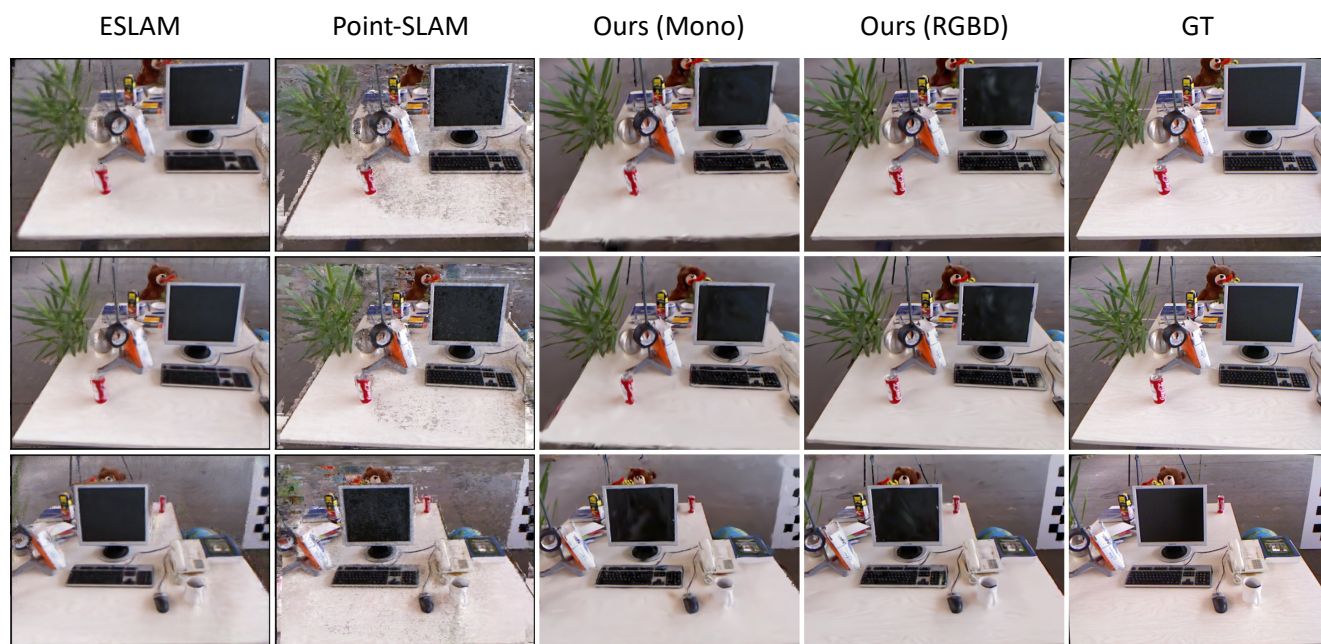


Figure 12. Rendering comparison on TUM fr2/xyz

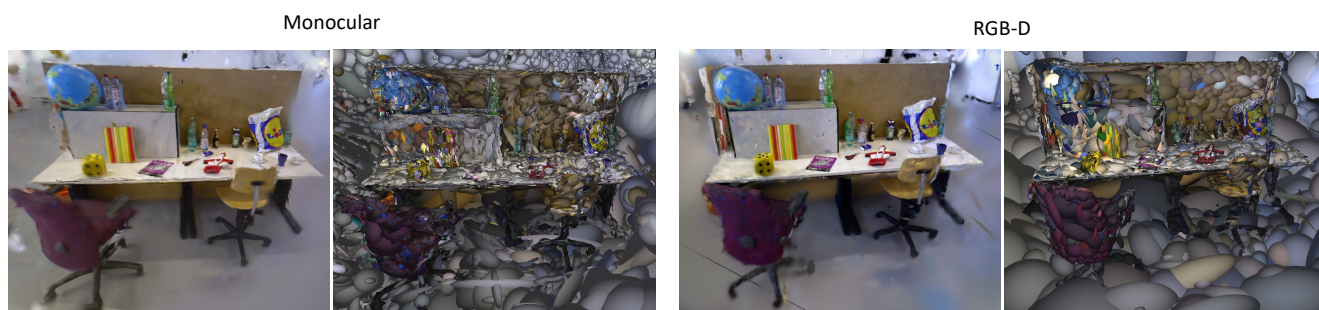


Figure 13. Novel view rendering and Gaussian visualizations on TUM fr3/office



Figure 14. Rendering comparison on TUM fr3/office



Figure 15. Novel view rendering and Gaussian visualizations on Replica

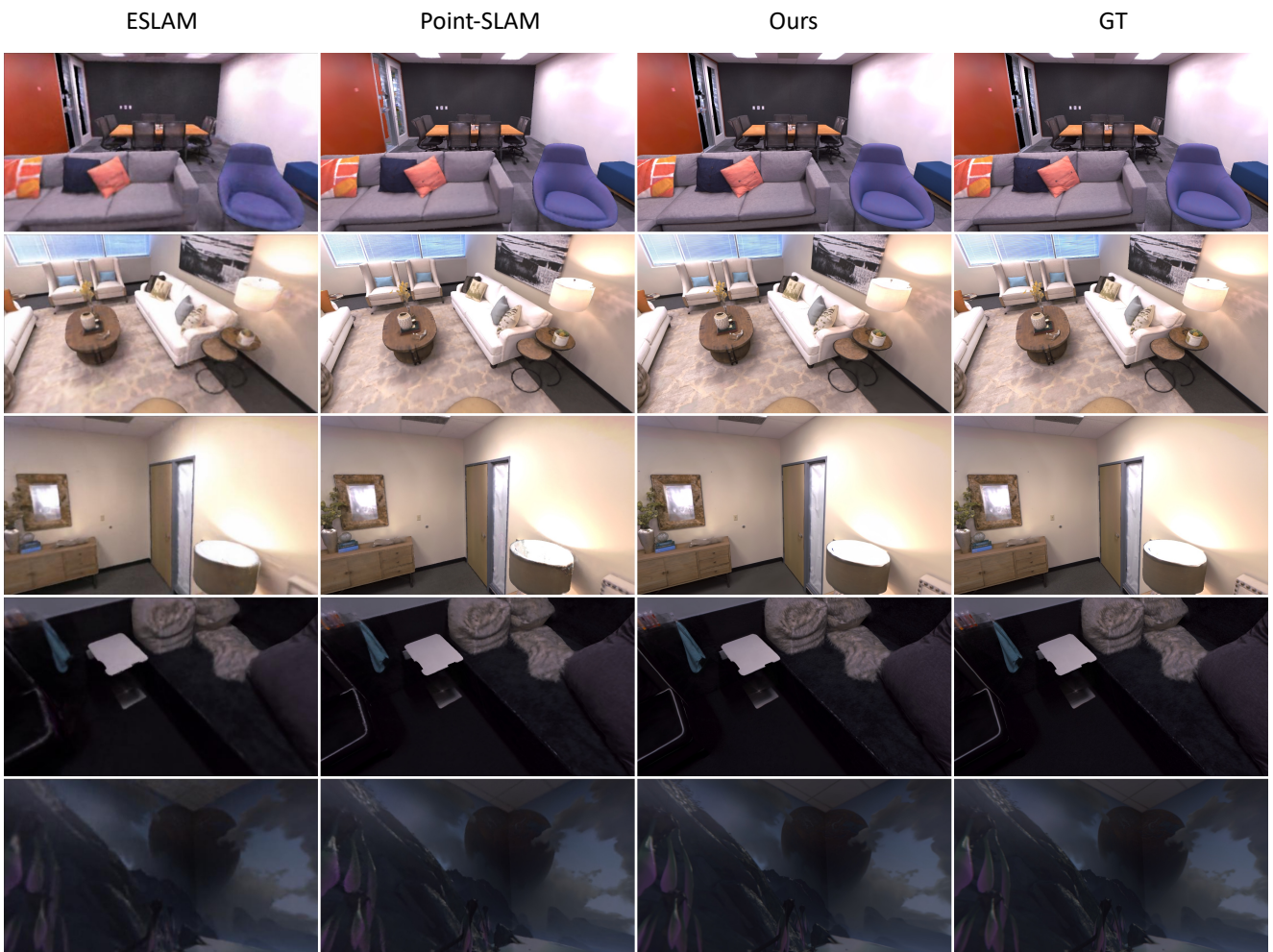


Figure 16. Rendering comparison on Replica

References

- [1] Carlos Campos, Richard Elvira, Juan J. Gómez, José M. M. Montiel, and Juan D. Tardós. ORB-SLAM3: An accurate open-source library for visual, visual-inertial and multi-map SLAM. *IEEE Transactions on Robotics (T-RO)*, 37(6):1874–1890, 2021.
- [2] J. Czarnowski, T. Laidlow, R. Clark, and A. J. Davison. Deepfactors: Real-time probabilistic dense monocular SLAM. *IEEE Robotics and Automation Letters (RAL)*, 5(2): 721–728, 2020.
- [3] Angela Dai, Matthias Nießner, Michael Zollhöfer, Shahram Izadi, and Christian Theobalt. BundleFusion: Real-time Globally Consistent 3D Reconstruction using On-the-fly Surface Re-integration. *ACM Transactions on Graphics (TOG)*, 36(3):24:1–24:18, 2017.
- [4] Eric Dexheimer and Andrew J. Davison. Learning a Depth Covariance Function. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.
- [5] J. Engel, V. Koltun, and D. Cremers. Direct sparse odometry. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 2017.
- [6] C. Forster, M. Pizzoli, and D. Scaramuzza. SVO: Fast Semi-Direct Monocular Visual Odometry. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, 2014.
- [7] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [8] Jiahui Huang, Shi-Sheng Huang, Haoxuan Song, and Shi-Min Hu. Di-fusion: Online implicit 3d reconstruction with deep priors. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [9] M. M. Johari, C. Carta, and F. Fleuret. ESLAM: Efficient dense slam system based on hybrid representation of signed distance fields. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.
- [10] M. Keller, D. Lefloch, M. Lambers, S. Izadi, T. Weyrich, and A. Kolb. Real-time 3D Reconstruction in Dynamic Scenes using Point-based Fusion. In *Proc. of Joint 3DIM/3DPVT Conference (3DV)*, 2013.
- [11] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3D gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics (TOG)*, 2023.
- [12] Leonid Keselman and Martial Hebert. Approximate differentiable rendering with algebraic surfaces. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2022.
- [13] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2015.
- [14] Heng Li, Xiaodong Gu, Weihao Yuan, Luwei Yang, Zilong Dong, and Ping Tan. Dense rgb slam with neural implicit maps. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2023.
- [15] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. *NeurIPS*, 2020.
- [16] Jonathon Luiten, Georgios Kopanas, Bastian Leibe, and Deva Ramanan. Dynamic 3d gaussians: Tracking by persistent dynamic view synthesis. *3DV*, 2024.
- [17] J. McCormac, A. Handa, A. J. Davison, and S. Leutenegger. SemanticFusion: Dense 3D semantic mapping with convolutional neural networks. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, 2017.
- [18] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020.
- [19] N. J. Mitra, N. Gelfand, H. Pottmann, and L. J. Guibas. Registration of Point Cloud Data from a Geometric Optimization Perspective. In *Proceedings of the Symposium on Geometry Processing*, 2004.
- [20] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multi-resolution hash encoding. *ACM Transactions on Graphics (TOG)*, 2022.
- [21] R. Mur-Artal and J. D. Tardós. ORB-SLAM2: An Open-Source SLAM System for Monocular, Stereo, and RGB-D Cameras. *IEEE Transactions on Robotics (T-RO)*, 33(5): 1255–1262, 2017.
- [22] R. Mur-Artal, J. M. M. Montiel, and J. D. Tardós. ORB-SLAM: a Versatile and Accurate Monocular SLAM System. *IEEE Transactions on Robotics (T-RO)*, 31(5):1147–1163, 2015.
- [23] R. A. Newcombe. *Dense Visual SLAM*. PhD thesis, Imperial College London, 2012.
- [24] 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 *Proceedings of the International Symposium on Mixed and Augmented Reality (ISMAR)*, 2011.
- [25] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [26] M. Nießner, M. Zollhöfer, S. Izadi, and M. Stamminger. Real-time 3D Reconstruction at Scale using Voxel Hashing. In *Proceedings of SIGGRAPH*, 2013.
- [27] Victor Adrian Prisacariu, Olaf Kähler, Ming-Ming Cheng, Carl Yuheng Ren, Julien P. C. Valentin, Philip H. S. Torr, Ian D. Reid, and David W. Murray. A framework for the volumetric integration of depth images. *CoRR*, abs/1410.0925, 2014.
- [28] Tong Qin, Jie Pan, Shaozu Cao, and Shaojie Shen. A general optimization-based framework for local odometry estimation with multiple sensors, 2019.

- [29] Erik Sandström, Yue Li, Luc Van Gool, and Martin R. Oswald. Point-slam: Dense neural point cloud-based slam. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2023.
- [30] Thomas Schöps, Torsten Sattler, and Marc Pollefeys. Surfelmeshing: Online surfel-based mesh reconstruction. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 2020.
- [31] Thomas Schöps, Torsten Sattler, and Marc Pollefeys. Bad slam: Bundle adjusted direct rgb-d slam. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [32] J. Solà, J. Deray, and D. Atchuthan. A micro Lie theory for state estimation in robotics. *arXiv:1812.01537*, 2018.
- [33] Julian Straub, Thomas Whelan, Lingni Ma, Yufan Chen, Erik Wijmans, Simon Green, Jakob J. Engel, Raul Mur-Artal, Carl Ren, Shobhit Verma, Anton Clarkson, Mingfei Yan, Brian Budge, Yajie Yan, Xiaqing Pan, June Yon, Yuyang Zou, Kimberly Leon, Nigel Carter, Jesus Briales, Tyler Gillingham, Elias Mueggler, Luis Pesqueira, Manolis Savva, Dhruv Batra, Hauke M. Strasdat, Renzo De Nardi, Michael Goesele, Steven Lovegrove, and Richard Newcombe. The Replica dataset: A digital replica of indoor spaces. *arXiv preprint arXiv:1906.05797*, 2019.
- [34] J. Sturm, N. Engelhard, F. Endres, W. Burgard, and D. Cremers. A Benchmark for the Evaluation of RGB-D SLAM Systems. In *Proceedings of the IEEE/RSJ Conference on Intelligent Robots and Systems (IROS)*, 2012.
- [35] E. Sucar, S. Liu, J. Ortiz, and A. J. Davison. iMAP: Implicit mapping and positioning in real-time. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2021.
- [36] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [37] Jiaxiang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang Zeng. Dreamgaussian: Generative gaussian splatting for efficient 3d content creation. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2024.
- [38] Zachary Teed and Jia Deng. DROID-SLAM: Deep Visual SLAM for Monocular, Stereo, and RGB-D Cameras. In *Neural Information Processing Systems (NIPS)*, 2021.
- [39] 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 (RAL)*, 2018.
- [40] Angtian Wang, Peng Wang, Jian Sun, Adam Kortylewski, and Alan Yuille. Voge: a differentiable volume renderer using gaussian ellipsoids for analysis-by-synthesis. 2022.
- [41] Hengyi Wang, Jingwen Wang, and Lourdes Agapito. Co-slam: Joint coordinate and sparse parametric encodings for neural real-time slam. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.
- [42] T. Whelan, M. Kaess, H. Johannsson, M. F. Fallon, J. J. Leonard, and J. B. McDonald. Real-time large scale dense RGB-D SLAM with volumetric fusion. *International Journal of Robotics Research (IJRR)*, 34(4-5):598–626, 2015.
- [43] T. Whelan, S. Leutenegger, R. F. Salas-Moreno, B. Glocker, and A. J. Davison. ElasticFusion: Dense SLAM without a pose graph. In *Proceedings of Robotics: Science and Systems (RSS)*, 2015.
- [44] Guanjun Wu, Taoran Yi, Jiemin Fang, Lingxi Xie, Xiaopeng Zhang, Wei Wei, Wenyu Liu, Qi Tian, and Xinggang Wang. 4d gaussian splatting for real-time dynamic scene rendering. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024.
- [45] 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 *Proceedings of the International Symposium on Mixed and Augmented Reality (ISMAR)*, 2022.
- [46] Zeyu Yang, Hongye Yang, Zijie Pan, Xiatian Zhu, and Li Zhang. Real-time photorealistic dynamic scene representation and rendering with 4d gaussian splatting. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2024.
- [47] Taoran Yi, Jiemin Fang, Guanjun Wu, Lingxi Xie, Xiaopeng Zhang, Wenyu Liu, Qi Tian, and Xinggang Wang. Gaussian-dreamer: Fast generation from text to 3d gaussian splatting with point cloud priors. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024.
- [48] 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 Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [49] Zihan Zhu, Songyou Peng, Viktor Larsson, Zhaopeng Cui, Martin R. Oswald, Andreas Geiger, and Marc Pollefeys. Nicer-slam: Neural implicit scene encoding for rgb slam. *International Conference on 3D Vision (3DV)*, 2024.
- [50] M. Zwicker, H. Pfister, J. van Baar, and M. Gross. Ewa splatting. *IEEE Transactions on Visualization and Computer Graphics*, 8(3):223–238, 2002.