

Splat-SLAM: Globally Optimized RGB-only SLAM with 3D Gaussians

Erik Sandström^{1,2†*} Ganlin Zhang^{1*} Keisuke Tateno² Michael Oechsle² Michael Niemeyer²
 Youmin Zhang⁶ Manthan Patel¹ Luc Van Gool⁵ Martin R. Oswald⁴ Federico Tombari^{2,3}
¹ETH Zürich ²Google ³TU München ⁴University of Amsterdam ⁵INSAIT ⁶Rock Universes
 *Equal contribution †Work done while at internship at Google

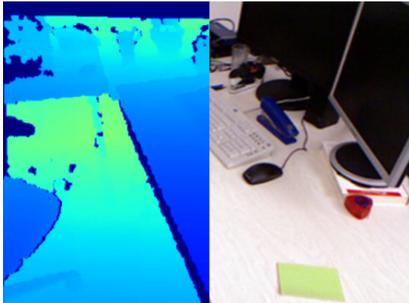
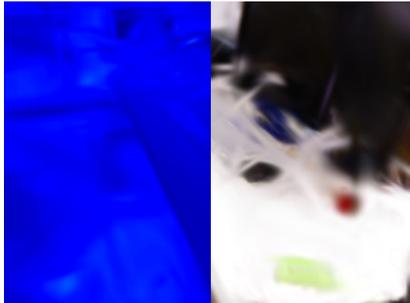
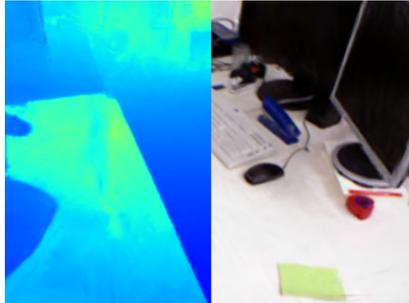
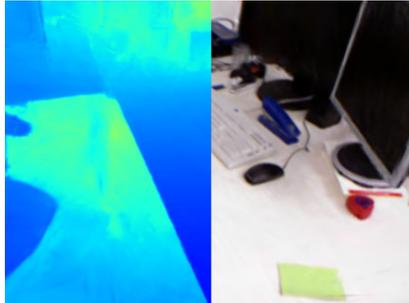
	Ground Truth	MonoGS [38]	Splat-SLAM (Ours)	
				
Depth L1 [cm]↓ PSNR↑		116.71 18.41	15.05	24.06
ATE RMSE [cm]↓		76.56	4.2	

Figure 1. **Splat-SLAM**. Our system yields accurate scene reconstruction (rendering depth L1), rendering (PSNR) and tracking accuracy (ATE RMSE) compared to MonoGS. The results averaged over all keyframes. The scene is from TUM-RGBD [56] fr1 room.

Abstract

3D Gaussian Splatting offers a compact, efficient approach to RGB-only dense SLAM by providing high-quality map rendering with a dense, optimized 3D Gaussian map. Existing methods, however, often underperform in reconstruction quality compared to alternatives like neural point clouds, primarily due to limited map and pose optimization or reliance on monocular depth. We introduce the first RGB-only SLAM system with globally optimized tracking, dynamically adapting the Gaussian map to keyframe pose and depth updates. To address the lack of geometric priors, we incorporate so called Disparity, Scale and Pose Optimization (DSPO) for bundle adjustment, jointly optimizing pose, depth, and monocular depth scale. Our tests on Replica, TUM-RGBD, and ScanNet confirm this approach achieves superior or comparable tracking, mapping, and rendering accuracy with small map sizes and fast runtimes.

1. Introduction

A common factor within the recent trend of dense SLAM is that the majority of works reconstruct a dense map by

optimizing a neural implicit encoding of the scene, either as weights of an MLP [1, 39, 45, 57], as features anchored in dense grids [3, 29, 42, 51, 58, 66, 67, 79, 81], using hierarchical octrees [71], via voxel hashing [8, 40, 49, 76, 77], point clouds [18, 30, 50] or axis-aligned feature planes [33, 47]. We have also seen the introduction of 3D Gaussian Splatting (3DGS) to the dense SLAM field [21, 24, 38, 69, 73].

Out of this 3D representation race there is, however, not yet a clear winner. In the context of dense SLAM, a careful modeling choice needs to be made to achieve accurate surface reconstruction as well as low tracking errors. Some takeaways can be deduced from the literature: neural implicit point cloud representations achieve state-of-the-art reconstruction accuracy [30, 50], especially with RGBD input. At the same time, 3D Gaussian splatting methods yield the highest fidelity renderings [21, 24, 38, 69, 73] and show promise in the RGB-only setting due to their flexibility in optimizing the surface location [21, 38]. However, they are not leveraging any multi-view depth or geometric prior leading to poor geometry in the RGB-only setting. The majority of the aforementioned works *only* deploy so called frame-to-model tracking, and do not implement global trajectory and map optimization, leading to excessive drift, especially in real world conditions. Instead, to this date,

frame-to-frame tracking methods, coupled with loop closure and global bundle adjustment (BA) achieve state-of-the-art tracking accuracy [76, 77]. However, they use hierarchical feature grids [76, 77], not suitable for map deformations at *e.g.* loop closure as they require expensive reintegration strategies.

In this work we propose an RGB-only SLAM system that combines the strengths of frame-to-frame tracking using recurrent dense optical flow [61] with the fidelity of 3D Gaussians as the map representation [38] (see Fig. 1). The 3D Gaussian map enables online map deformations at loop closure and global BA. To enable accurate surface reconstruction, we leverage consistent so called proxy depth that combines multi-view depth estimation with learned monocular depth. Our **contribution** comprises, for the first time, a SLAM pipeline encompassing all the following parts:

- A globally consistent frame-to-frame RGB-only tracker.
- A dense deformable 3D Gaussian map that adapts online to loop closure and global bundle adjustment.
- A novel scheme for joint Disparity, Scale and Pose Optimization (DSPO) that combines pose and geometry estimation. It refines inaccurate parts of the estimated keyframe disparity by tightly coupling a monocular depth prior into the bundle adjustment.
- Improved map sizes and runtimes compared to other dense SLAM approaches.

2. Related Work

Dense Visual SLAM. Curless and Levoy [9] pioneered dense online 3D mapping with truncated signed distance functions, with KinectFusion [42] demonstrating real-time SLAM via depth maps. Enhancements like voxel hashing [11, 22, 40, 43, 44] and octrees [5, 31, 37, 53, 71] improved scalability, while point-based SLAM [4, 6, 22, 25, 30, 50, 52, 68, 74] has also been effective. To address pose drift, globally consistent pose estimation and dense mapping techniques have been developed, often dividing the global map into submaps [2, 4, 7, 11, 15, 17, 22, 23, 30, 34–36, 40, 48, 55, 59, 59]. Loop detection triggers submap deformation via pose graph optimization [4, 7, 13, 14, 16–18, 23, 27, 30, 35, 36, 40, 40, 48, 52, 55, 59, 63, 70]. Sometimes global BA is used for refinement [4, 8, 11, 18, 40, 52, 59, 61, 70, 72]. 3D Gaussian SLAM with RGBD input has also been shown, but these methods do not consider global consistency via *e.g.* loop closure [24, 69, 73]. Other approaches to global consistency minimize reprojection errors directly, with DROID-SLAM [61] refining dense optical flow and camera poses iteratively, and recent enhancements like GO-SLAM [77] and HI-SLAM [76] optimizing factor graphs for accurate tracking. For a recent survey on NeRF-inspired dense SLAM, see [62].

RGB-only Dense Visual SLAM. The majority of NeRF

inspired RGB-only dense SLAM methods do not address the problem of global map consistency or requires expensive reintegration strategies via backpropagation [8, 19, 20, 28, 41, 46, 49, 76, 77, 80]. MonoGS [38] and Photo-SLAM [21] pioneered RGB-only SLAM with 3D Gaussians. However, they lack proxy depth which prevents them from achieving high accuracy mapping. MonoGS [38] also lacks global consistency. MoD-SLAM [78] uses an MLP to parameterize the map via a unique reparameterization.

Depth Priors for RGB-only SLAM. NICER-SLAM [80] estimates the scale and shift of a relative mono-depth estimator and supervises all pixels equally. MoD-SLAM [78] combines relative and metric mono-depth estimation, and requires additional finetuning of the metric depth. HI-SLAM [76] proposes a similar technique to ours, but regularizes all available keyframe depth pixels with the mono-depth prior. In our DSPO, we instead split the optimization and use the monocular prior to regularize the high error keyframe depth pixels while the low error keyframe depth is kept fixed to stabilize scale estimation.

3. Method

Splat-SLAM is a monocular SLAM system which tracks the camera pose while reconstructing the dense geometry of the scene in an online manner. This is achieved through the following steps: We first track the camera by performing local BA on selected keyframes by fitting them to dense optical flow estimates. The local BA optimizes the camera pose as well as the dense depth of the keyframe. For global consistency, when loop closure is detected, loop BA is performed on an extended graph including the loop nodes and edges (Sec. 3.1). Interleaved with tracking, mapping is done on a progressively growing 3D Gaussian map which deforms online to the keyframe poses and so called proxy depth maps (Sec. 3.2). For an overview of our method, see Fig. 2.

3.1. Tracking

To predict the motion of the camera during scene exploration, we use a pretrained recurrent optical flow model [60] coupled with our so called Disparity, Scale and Pose Optimization (DSPO) to jointly optimize camera poses and per pixel disparities. In the following, we describe this process in detail.

Optimization is done with the Gauss-Newton algorithm over a factor graph $G(V, E)$, where the nodes V store the keyframe pose and disparity, and edges E store the optical flow between keyframes. Odometry keyframe edges are added to G by computing the optical flow to the last added keyframe. If the mean flow is larger than a threshold $\tau \in \mathbb{R}$, the new keyframe is added to G . Edges for loop closure and global BA are discussed later. Importantly, the same objective is optimized for local BA, loop closure and global

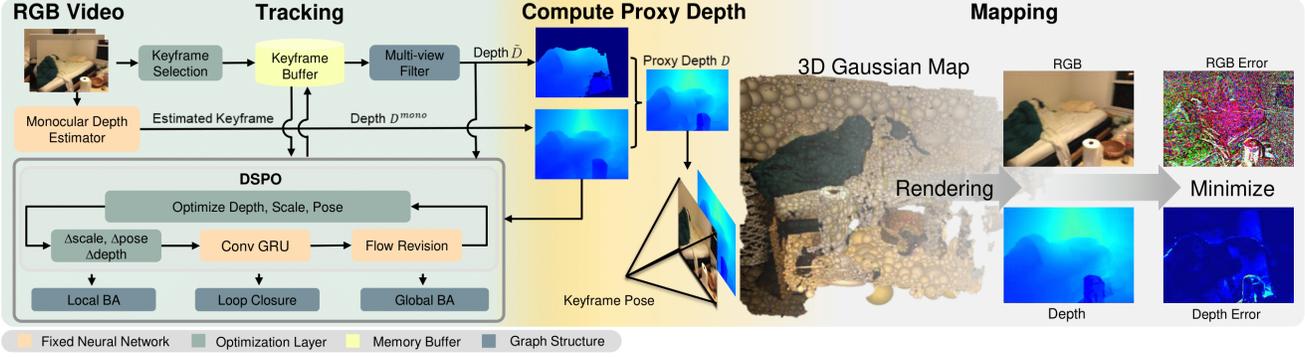


Figure 2. **Splat-SLAM Architecture.** Given an RGB input stream, we track and map each keyframe, initially estimating poses through local bundle adjustment (BA) using DSPO (Disparity, Scale and Pose Optimization). This DSPO integrates pose and depth estimation, enhancing depth with monocular depth. It further refines poses globally via online loop closure and global BA. The proxy depth map merges keyframe depths \tilde{D} from the tracking with monocular depth D^{mono} to fill gaps. Mapping employs a deformable 3D Gaussian map, optimizing its parameters through a re-rendering loss. Notably, the 3D map adjusts for global pose and depth updates before each mapping phase.

BA, but over factor graphs with different structures.

The DSPO consists of two optimization objectives that are optimized alternatively. The first objective, typically termed Dense Bundle Adjustment (DBA) [61] optimizes the pose and disparity of the keyframes jointly, Eq. (1). Specifically, the objective is optimized over a local graph defined within a sliding window over the current frame.

$$\arg \min_{\omega, d} \sum_{(i,j) \in E} \|\tilde{p}_{ij} - K\omega_j^{-1}(\omega_i(1/d_i)K^{-1}[p_i, 1]^T)\|_{\Sigma_{ij}}^2, \quad (1)$$

with $\tilde{p}_{ij} \in \mathbb{R}^{(W \times H \times 2) \times 1}$ being the flattened predicted pixel coordinates when the pixels $p_i \in \mathbb{R}^{(W \times H \times 2) \times 1}$ from keyframe i are projected into keyframe j using optical flow. Further, K is the camera intrinsics, ω_j and ω_i the camera-to-world extrinsics for keyframes j and i , d_i the disparity of pixel p_i and $\|\cdot\|_{\Sigma_{ij}}$ is the Mahalanobis distance with diagonal weighting matrix Σ_{ij} . Each weight denotes the confidence of the optical flow prediction for each pixel in \tilde{p}_{ij} . For clarity of the presentation, we omit homogeneous coordinates.

In the second objective, we introduce monocular depth D^{mono} as two additional data terms, to tackle noisy disparity estimates from the DBA optimization. The monocular depth D^{mono} is predicted at runtime by a pretrained relative depth DPT model [12].

$$\begin{aligned} \arg \min_{d^h, \theta, \gamma} \sum_{(i,j) \in E} & \|\tilde{p}_{ij} - K\omega_j^{-1}(\omega_i(1/d_i^h)K^{-1}[p_i, 1]^T)\|_{\Sigma_{ij}}^2 \\ & + \alpha_1 \sum_{i \in V} \|d_i^h - (\theta(1/D_i^{\text{mono}}) + \gamma_i)\|^2 \\ & + \alpha_2 \sum_{i \in V} \|d_i^l - (\theta(1/D_i^{\text{mono}}) + \gamma_i)\|^2. \end{aligned} \quad (2)$$

Here, the optimizable parameters are the scales $\theta \in \mathbb{R}$, shifts $\gamma \in \mathbb{R}$ and a subset of the disparities d^h classified as being

high error (explained later). This is done since the monocular depth is only deemed useful where the multi-view disparity d_i optimization is inaccurate. Furthermore, $\alpha_1 < \alpha_2$, which is done to ensure that the scales θ and shifts γ are optimized with the preserved low error disparities d^l . The scale θ_i and shift γ_i are initialized using least squares fitting

$$\{\theta_i, \gamma_i\} = \arg \min_{\theta, \gamma} \sum_{(u,v)} \left((\theta(1/D_i^{\text{mono}}) + \gamma) - d_i^l \right)^2. \quad (3)$$

Equation (1) and Eq. (2) are optimized alternatively to avoid the scale ambiguity encountered if d , θ , γ and ω are optimized jointly.

Next, we describe how high and low error disparities are classified. For a given disparity map d_i (separated into low and high error parts $\{d_i^l, d_i^h\}$) for frame i , we denote the corresponding depth $D_i = 1/d_i$. Pixel correspondences (u, v) and (\hat{u}, \hat{v}) between keyframes i and j respectively are established by warping (u, v) into frame j with depth \tilde{D}_i as

$$\begin{aligned} p_j &= \omega_j \tilde{D}_i(u, v) K^{-1}[u, v, 1]^T, \\ [\hat{u}, \hat{v}, 1]^T &\propto K\omega_j^{-1}[p_i, 1]^T. \end{aligned} \quad (4)$$

The corresponding 3D point to (\hat{u}, \hat{v}) is computed from the depth at (\hat{u}, \hat{v}) as

$$p_j = \omega_j \tilde{D}_j(\hat{u}, \hat{v}) K^{-1}[\hat{u}, \hat{v}, 1]^T. \quad (5)$$

If the L2 distance between p_i and p_j is smaller than a threshold, the depth $\tilde{D}_i(u, v)$ is consistent between i and j . By looping over all keyframes except i , the global two-view consistency n_i can be computed for frame i as

$$n_i(u, v) = \sum_{\substack{k \in \text{KFs} \\ k \neq i}} \mathbb{1} \left(\|p_i - p_k\|_2 < \eta \cdot \text{average}(\tilde{D}_i) \right). \quad (6)$$

Here, $\mathbb{1}(\cdot)$ is the indicator function and $\eta \in \mathbb{R}_{\geq 0}$ is a hyper-parameter and n_i is the total two-view consistency for pixel

(u, v) in keyframe i . $\tilde{D}_i(u, v)$ is valid if n_i is larger than a threshold.

Loop Closure. To mitigate scale and pose drift, we incorporate loop closure along with online global bundle adjustment (BA) in addition to local window frame tracking. Loop detection is achieved by calculating the mean optical flow magnitude between the current active keyframes (within the local window) and all previous keyframes. Two criteria are evaluated for each keyframe pair: First, the optical flow must be below a specified threshold τ_{loop} , ensuring sufficient co-visibility between the views. Second, the time interval between the frames must exceed a predefined threshold τ_t to prevent the introduction of redundant edges into the graph. When both criteria are met, a unidirectional edge is added to the graph. During the loop closure optimization process, only the active keyframes and their connected loop nodes are optimized to keep the computational load manageable.

Global BA. For the online global BA, a separate graph that includes all keyframes up to the present is constructed. Edges are introduced based on the temporal and spatial relationships between the keyframes, as outlined in [77]. We execute online global BA every 20 keyframes. To maintain numerical stability, the scales of the disparities and poses are normalized prior to each global BA optimization. This normalization involves calculating the average disparity \bar{d} across all keyframes and then adjusting the disparity to $d_{\text{norm}} = d/\bar{d}$ and the pose translation to $t_{\text{norm}} = dt$.

3.2. Deformable 3D Gaussian Scene Representation

We adopt a 3D Gaussian Splatting representation [26] which deforms under DSPO or loop closure optimizations to achieve global consistency. Thus, the scene is represented by a set $\mathcal{G} = \{g_i\}_{i=1}^N$ of 3D Gaussians. Each Gaussian primitive g_i , is parameterized by a covariance matrix $\Sigma_i \in \mathbb{R}^{3 \times 3}$, a mean $\boldsymbol{\mu}_i \in \mathbb{R}^3$, opacity $o_i \in [0, 1]$, and color $\mathbf{c}_i \in \mathbb{R}^3$. All attributes of each Gaussian are optimized through back-propagation. The density function of a single Gaussian is described as

$$g_i(\mathbf{x}) = \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^\top \Sigma_i^{-1}(\mathbf{x} - \boldsymbol{\mu}_i)\right). \quad (7)$$

Here, the spatial covariance Σ_i defines an ellipsoid and is decomposed as $\Sigma_i = R_i S_i S_i^\top R_i^\top$, where $S_i = \text{diag}(s_i) \in \mathbb{R}^{3 \times 3}$ is the spatial scale and $R_i \in \mathbb{R}^{3 \times 3}$ represents the rotation.

Rendering. Rendering color and depth from \mathcal{G} , given a camera pose, involves first projecting (known as ‘‘splatting’’) 3D Gaussians onto the 2D image plane. This is done by projecting the covariance matrix Σ and mean $\boldsymbol{\mu}$ as $\Sigma' = JR\Sigma R^\top J^\top$ and $\boldsymbol{\mu}' = K\omega^{-1}\boldsymbol{\mu}$, where R is the rotation component of world-to-camera extrinsics ω^{-1} and J is the Jacobian of the affine approximation of the projective

transformation [82]. The final pixel color C and depth D^r at pixel \mathbf{x}' is computed by blending 3D Gaussian splats that overlap at a given pixel, sorted by their depth as

$$C = \sum_{i \in \mathcal{N}} \mathbf{c}_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j)$$

$$D^r = \sum_{i \in \mathcal{N}} \hat{d}_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j), \quad (8)$$

where \hat{d}_i is the z-axis depth of the center of the i -th 3D Gaussian and the final opacity α_i is the product of the opacity o_i and the 2D Gaussian density as

$$\alpha_i = o_i \exp\left(-\frac{1}{2}(\mathbf{x}' - \boldsymbol{\mu}'_i)^\top \Sigma_i'^{-1}(\mathbf{x}' - \boldsymbol{\mu}'_i)\right). \quad (9)$$

Map Initialization. For every new keyframe, we adopt the RGBD strategy of MonoGS [38] for adding new Gaussians to the unexplored scene space. As we do not have access to a depth sensor, we construct a proxy depth map D by combining the inlier multi-view depth \tilde{D} and the monocular depth D^{mono} as

$$D(u, v) = \begin{cases} \tilde{D}(u, v) & \text{if } \tilde{D}(u, v) \text{ is valid} \\ \theta D^{\text{mono}}(u, v) + \gamma & \text{otherwise} \end{cases} \quad (10)$$

Here, θ and γ are computed as in Eq. (3) but using depth instead of disparity.

Keyframe Selection and Optimization. Apart from the keyframe selection based on a mean optical flow threshold τ , we additionally adopt the keyframe selection strategy from [38] to avoid mapping redundant frames.

To optimize the 3D Gaussian parameters, we batch the parameter updates to a local window similar to [38] and apply a photometric and geometric loss to the proxy depth as well as a scale regularizer to avoid artifacts from elongated Gaussians. Inspired by [38], we further use exposure compensation by optimizing an affine transformation for each keyframe. The final loss is

$$\min_{\mathcal{G}, \mathbf{a}, \mathbf{b}} \sum_{k \in \text{KFs}} \frac{\lambda}{N_k} |(a_k C_k + b_k) - C_k^{\text{gt}}|_1$$

$$+ \frac{1 - \lambda}{N_k} |D_k^r - D_k|_1 + \frac{\lambda_{\text{reg}}}{|\mathcal{G}|} \sum_i |s_i - \tilde{s}_i|_1, \quad (11)$$

where KFs contains the set of keyframes in the local window, N_k is the number of pixels per keyframe, λ and λ_{reg} are hyperparameters, $\mathbf{a} = \{a_1, \dots, a_k, \dots\}$ and $\mathbf{b} = \{b_1, \dots, b_k, \dots\}$ are the parameters for the exposure compensation and \tilde{s} is the mean scaling, repeated over the three dimensions.

Metric	GO-SLAM [77]	NICER-SLAM [80]	MoD-SLAM [28]	Photo-SLAM [21]	Mono-GS [38]	Q-SLAM [46]	Ours
PSNR \uparrow	22.13	25.41	27.31	33.30	31.22	32.49	36.45
SSIM \uparrow	0.73	0.83	0.85	0.93	0.91	0.89	0.95
LPIPS \downarrow	-	0.19	-	-	0.21	0.17	0.06
ATE RMSE \downarrow	0.39	1.88	0.35	1.09	14.54	-	0.35

Table 1. **Rendering and Tracking Results on Replica [54] for RGB-Methods.** Our method outperforms all methods on rendering and performs on par for tracking accuracy. Results are from [62] except ours (average over 8 scenes). Best results are highlighted as **first**, **second**, **third**.

Metrics	NeRF-SLAM [62]	DIM-SLAM [28]	GO-SLAM [77]	NICER-SLAM [80]	HI-SLAM [76]	MoD-SLAM [78]	Mono-GS [38]	Q-SLAM [46]	Ours
Render Depth L1 \downarrow	4.49	-	-	-	-	-	27.24	2.76	2.41
Accuracy \downarrow	-	4.03	3.81	3.65	3.62	2.48	30.61	-	2.43
Completion \downarrow	-	4.20	4.79	4.16	4.59	-	12.19	-	3.64
Comp. Rat. \uparrow	-	79.60	78.00	79.37	80.60	-	40.53	-	84.69

Table 2. **Reconstruction Results on Replica [54] for RGB-Methods.** Our method outperforms existing works on all metrics. Results are averaged over 8 scenes.

Map Deformation. Since our tracking framework is globally consistent, changes in the keyframe poses and proxy depth maps need to be accounted for in the 3D Gaussian map by a non-rigid deformation. Though the Gaussian means are directly optimized, one could in theory let the optimizer deform the map as refined poses and proxy depth maps are provided. We find, however, that in particular rendering is aided by actively deforming the 3D Gaussian map. We apply the deformation to all Gaussians which receive updated poses and depths before mapping.

Each Gaussian g_i is associated with a keyframe that anchored it to the map \mathcal{G} . Assume that a keyframe with camera-to-world pose ω and proxy depth D is updated such that $\omega \rightarrow \omega'$ and $D \rightarrow D'$. We update the mean, scale and rotation of all Gaussians g_i associated with the keyframe. Association is determined by what keyframe added the Gaussian to the scene. The mean μ_i is projected into ω to find the pixel correspondence (u, v) . Since the Gaussians are not necessarily anchored on the surface, instead of re-anchoring the mean at D' , we opt to shift the mean by $D'(u, v) - D(u, v)$ along the optical axis. We update R_i and s_i accordingly as

$$\mu'_i = \left(1 + \frac{D'(u, v) - D(u, v)}{(\omega^{-1}\mu_i)_z}\right)\omega'\omega^{-1}\mu_i, \quad (12)$$

$$R'_i = R'R^{-1}R_i, \quad s'_i = \left(1 + \frac{D'(u, v) - D(u, v)}{(\omega^{-1}\mu_i)_z}\right)s_i.$$

Here, $(\cdot)_z$ denotes the z-axis depth. For Gaussians which project into pixels with missing depth or outside the viewing frustum, we *only* rigidly deform them. After the final global BA optimization, we additionally deform the Gaussian map and perform a set of final refinements (see suppl. material).

4. Experiments

We first describe our experimental setup and then evaluate our method against state-of-the-art dense RGB and RGBD SLAM methods on Replica [54] as well as the real world TUM-RGBD [56] and the ScanNet [10] datasets. For more experiments and details, we refer to the supplementary material.

Implementation Details. For the proxy depth, we use $\eta = 0.01$ to filter points and use the condition $n_c \geq 2$ to ensure multi-view consistency. For the mapping loss function, we use $\lambda = 0.8$, $\lambda_{reg} = 10.0$. We use 60 iterations during mapping. For tracking, we use $\alpha_1 = 0.01$ and $\alpha_2 = 0.1$ as weights for the DSPO. We use the flow threshold $\tau = 4.0$ on ScanNet, $\tau = 3.0$ on TUM-RGBD and $\tau = 2.25$ on Replica. The threshold for loop detection is $\tau_{loop} = 25.0$. The time interval threshold is $\tau_t = 20$. We conducted the experiments on a cluster with an NVIDIA A100 GPU.

Evaluation Metrics. For rendering we report PSNR, SSIM [65] and LPIPS [75] on the rendered keyframe images against the sensor images. For reconstruction, we first extract the meshes with marching cubes [32] as in [50] and evaluate the meshes using accuracy [cm], completion [cm] and completion ratio [%] (threshold 5 cm) against the ground truth meshes. We also report the re-rendering depth L1 [cm] metric to the ground truth sensor depth as in [49]. We use ATE RMSE [cm] [56] to evaluate the estimated trajectory.

Datasets. We use the RGBD trajectories from [57] captured from the synthetic Replica dataset [54]. We also test on real-world data using the TUM-RGBD [56] and the ScanNet [10] datasets.

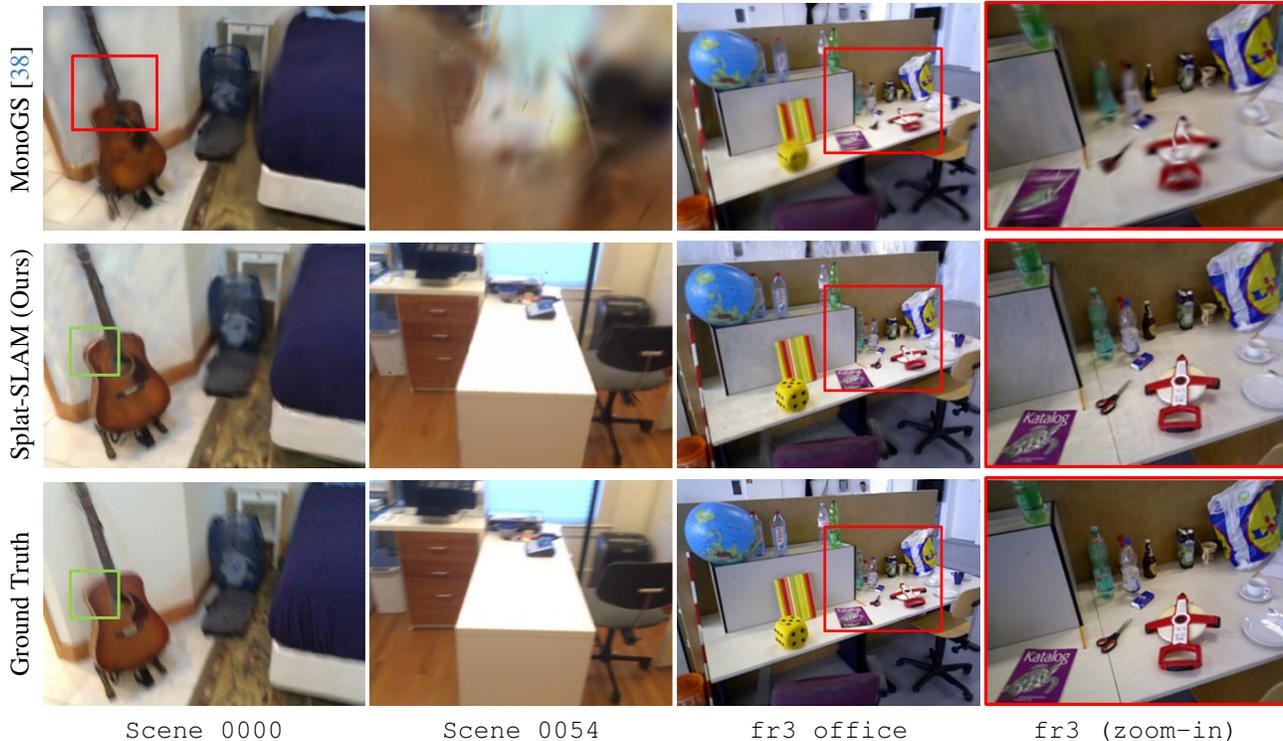


Figure 3. **Rendering Results on ScanNet [10] and TUM-RGBD [56].** Our method yields better rendering quality MonoGS. First column: The red box shows a rendering distortion, likely from the large trajectory error. The green boxes show that our method fuses information from multiple views to avoid motion blur, present in the input. Fourth column: The rendering is from the pose of the red box in the third column.

Method	Metric	0000	0059	0106	0169	0181	0207	Avg.
<i>RGB-D Input</i>								
SplaTaM [24]	PSNR \uparrow	19.33	19.27	17.73	21.97	16.76	19.80	19.14
	SSIM \uparrow	0.66	0.79	0.69	0.78	0.68	0.70	0.72
	LPIPS \downarrow	0.44	0.29	0.38	0.28	0.42	0.34	0.36
MonoGS [38]	PSNR \uparrow	18.70	20.91	19.84	22.16	22.01	18.90	20.42
	SSIM \uparrow	0.71	0.79	0.81	0.78	0.82	0.75	0.78
	LPIPS \downarrow	0.48	0.32	0.32	0.34	0.42	0.41	0.38
Gaussian-SLAM [73]	PSNR \uparrow	28.54	26.21	26.26	28.60	27.79	28.63	27.67
	SSIM \uparrow	0.93	0.93	0.93	0.92	0.92	0.91	0.92
	LPIPS \downarrow	0.27	0.21	0.22	0.23	0.28	0.29	0.25
<i>RGB Input</i>								
GO-SLAM [77]	PSNR \uparrow	15.74	13.15	14.58	14.49	15.72	15.37	14.84
	SSIM \uparrow	0.42	0.32	0.46	0.42	0.53	0.39	0.42
	LPIPS \downarrow	0.61	0.60	0.59	0.57	0.62	0.60	0.60
MonoGS [38]	PSNR \uparrow	16.91	19.15	18.57	20.21	19.51	18.37	18.79
	SSIM \uparrow	0.62	0.69	0.74	0.74	0.75	0.70	0.71
	LPIPS \downarrow	0.70	0.51	0.55	0.54	0.63	0.58	0.59
Ours	PSNR \uparrow	28.68	27.69	27.70	31.14	31.15	30.49	29.48
	SSIM \uparrow	0.83	0.87	0.86	0.87	0.84	0.84	0.85
	LPIPS \downarrow	0.19	0.15	0.18	0.15	0.23	0.19	0.18

Table 3. **Rendering Performance on ScanNet [10].** Our method performs even better or on par with all RGB-D methods. We take the numbers for SplaTaM and Gaussian-SLAM from [73].

Baseline Methods. We compare our method to numerous works on dense RGB and RGBD SLAM. The main baseline is MonoGS [38].

Method	Method	f1/desk	f2/xyz	f3/off	f1/desk2	f1/room	Avg.
<i>RGB-D Input</i>							
SplaTaM [24]	PSNR \uparrow	22.00	24.50	21.90	-	-	-
	SSIM \uparrow	0.86	0.95	0.88	-	-	-
	LPIPS \downarrow	0.23	0.10	0.20	-	-	-
Gaussian-SLAM [73]	PSNR \uparrow	24.01	25.02	26.13	23.15	22.98	24.26
	SSIM \uparrow	0.92	0.92	0.94	0.91	0.89	0.92
	LPIPS \downarrow	0.18	0.19	0.14	0.20	0.24	0.19
<i>RGB Input</i>							
Photo-SLAM [21]	PSNR \uparrow	20.97	21.07	19.59	-	-	-
	SSIM \uparrow	0.74	0.73	0.69	-	-	-
	LPIPS \downarrow	0.23	0.17	0.24	-	-	-
MonoGS [38]	PSNR \uparrow	19.67	16.17	20.63	19.16	18.41	18.81
	SSIM \uparrow	0.73	0.72	0.77	0.66	0.64	0.70
	LPIPS \downarrow	0.33	0.31	0.34	0.48	0.51	0.39
Ours	PSNR \uparrow	25.61	29.53	26.05	23.98	24.06	25.85
	SSIM \uparrow	0.84	0.90	0.84	0.81	0.80	0.84
	LPIPS \downarrow	0.18	0.08	0.20	0.23	0.24	0.19

Table 4. **Rendering Performance on TUM-RGBD [56].** Our method performs competitively or better than RGB-D methods. For all RGB-D methods, we take the numbers from [73].

Rendering. In Tab. 1, we evaluate the rendering performance on Replica [54] and find that our method performs superior among all baseline RGB-methods. Table 3 and Table 4 show the rendering accuracy on the ScanNet [10] and TUM-RGBD [56] datasets. In particular, we outperform existing RGB-only works with a clear margin, while even beating

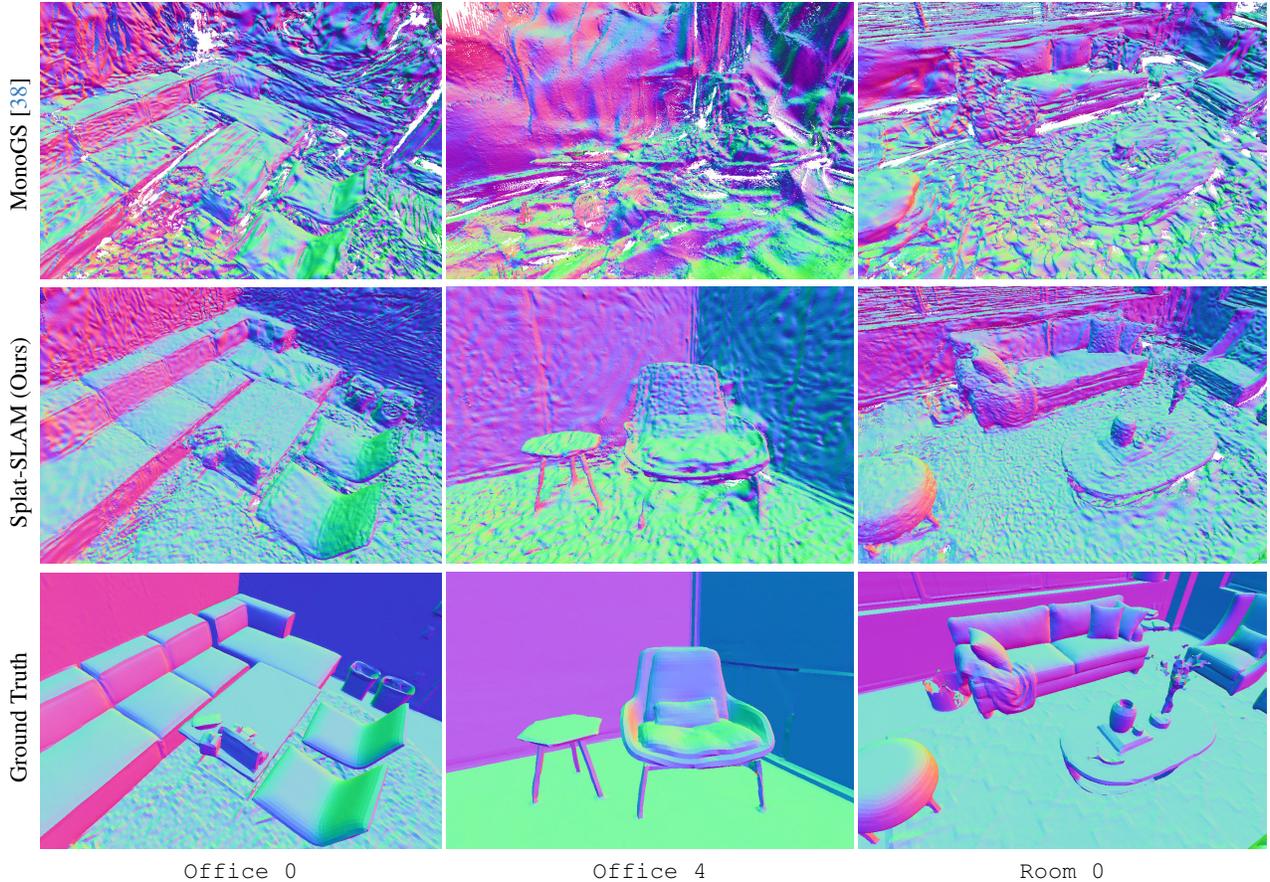


Figure 4. **Reconstruction Results on Replica [54] on Normal Shaded Meshes.** Our method achieves higher geometric accuracy compared to existing works. MonoGS suffers significantly from a lack of proxy depth, despite multiview optimization.

Method	00	59	106	169	181	207	Avg.-6	54	233	Avg.-8
<i>RGB-D Input</i>										
NICE-SLAM [79]	12.0	14.0	7.9	10.9	13.4	6.2	10.7	20.9	9.0	11.8
Co-SLAM [64]	7.1	11.1	9.4	5.9	11.8	7.1	8.7	-	-	-
ESLAM [33]	7.3	8.5	7.5	6.5	9.0	5.7	7.4	36.3	4.3	10.6
MonoGS[38]	16.1	6.4	8.1	8.7	26.4	9.2	12.5	20.6	13.1	13.6
<i>RGB Input</i>										
MonoGS[38]	149.2	96.8	155.5	140.3	92.6	101.9	122.7	206.4	89.1	129.0
GO-SLAM [77]	5.9	8.3	8.1	8.4	8.3	6.9	7.7	13.3	5.3	8.1
HI-SLAM[76]	6.4	7.2	6.5	8.5	7.6	8.4	7.4	-	-	-
Q-SLAM[46]	5.8	8.5	8.4	8.7	8.8	-	-	12.6	5.3	-
Ours	5.5	9.1	7.0	8.2	8.3	7.5	7.6	9.4	5.1	7.5

Table 5. **Tracking Accuracy on ScanNet [10]** Our method performs on average competitively with HI-SLAM and better than all other methods. Results for the RGB-D methods are from [30].

the currently best RGBD method, Gaussian-SLAM [73] on most metrics, despite the fact that we do not implement view-dependent rendering in the form of spherical harmonics. We attribute this to our deformable 3D Gaussian map, optimized with strong proxy depth along a globally consistent tracking backend. In Fig. 3 and Fig. 1 we show renderings on the real-world ScanNet [10] and TUM-RGBD [56] datasets. Due to high tracking errors, MonoGS [38] performs poorly

Method	f1/dsk	f2/xyz	f3/off	Avg.-3	f1/dsk2	f1/rm	Avg.-5
<i>RGB-D Input</i>							
SplaTAM [24]	3.4	1.2	5.2	3.3	6.5	11.1	5.5
GS-SLAM [69]	1.5	1.6	1.7	1.6	-	-	-
GO-SLAM [77]	1.5	0.6	1.3	1.1	-	4.7	-
MonoGS [38]	1.4	1.4	1.5	1.5	5.1	6.3	3.1
<i>RGB Input</i>							
MonoGS [38]	3.8	5.2	2.9	4.0	75.7	76.6	32.8
Photo-SLAM [21]	1.5	1.0	1.3	1.3	-	-	-
DIM-SLAM [28]	2.0	0.6	2.3	1.6	-	-	-
GO-SLAM [77]	1.6	0.6	1.5	1.2	2.8	5.2	2.3
MoD-SLAM [78]	1.5	0.7	1.1	1.1	-	-	-
Q-SLAM [46]	1.3	0.9	-	-	2.3	4.9	-
Ours	1.6	0.2	1.4	1.1	2.8	4.2	2.1

Table 6. **Tracking Accuracy on TUM-RGBD [56].** Our method performs even better than RGB-D methods.

on some scenes, yet fails to achieve the same fidelity as our method when the tracking error is low, as a result of the weak geometric constraints during optimization.

Reconstruction. We show quantitative and qualitative results on the Replica [54] dataset in Tab. 2 and Fig. 4 respectively. Our method achieves the best performance on all metrics. Qualitatively, we show normal shaded meshes



Figure 5. **Comparison of Estimated Depth.** We show the depth output \tilde{D} from the tracker. The pixels which are invalid (high error) are colored dark blue. DBA is the method that Droid-SLAM [61] uses. The DBA+mono prior strategy is used in HI-SLAM [76], *i.e.* the mono prior supervises all pixels directly. It is clear that our formulation (DSPO) provides the most consistent keyframe depth.

Mono Depth	Multiview Depth	Multiview Filtering	PSNR [dB] \uparrow	Acc. [cm] \downarrow	Comp. [cm] \downarrow	Comp. Ratio [cm] \uparrow
✓	✓	✗	36.02	3.62	4.08	81.16
✗	✓	✓	36.17	2.64	4.73	80.12
✗	✓	✗	36.21	18.71	4.06	80.29
✓	✓	✓	36.45	2.43	3.64	84.69

Table 7. **Ablation Study on Replica [54].** We show that the combination of filtered multiview depth completed with monocular depth yields the best performance on all metrics. Mono Depth refers to D^{mono} , Multiview Depth refers to \tilde{D} and Multiview Filtering means enabling Eq. (6). All results are averaged over 8 scenes.

from different viewpoints. Our method can reconstruct finer details than existing works, especially around thin structures (*e.g.* second row), where our strong proxy depth coupled with the 3D Gaussian map representation yields superior depth rendering, which directly influences the mesh quality. MonoGS [38] suffers significantly from the lack of proxy depth, visible in all scenes. Figure 1 shows depth rendering on the real-world TUM-RGBD [56] room scene. We compute the average depth L1 error over all keyframes, achieving 15.05 cm, beating existing works.

Tracking. In Tab. 1, Tab. 5 and Tab. 6, we report the tracking accuracy of the estimated trajectory on Replica [54], ScanNet [10] and TUM-RGBD [56]. On all datasets, our method shows competitive results in every single scene and gives the best average value among the RGB and RGB-D methods.

Ablation Study. In Tab. 7, we conduct a set of ablation studies, by enabling and disabling certain parts. We find that the combination of filtered multiview depth completed with monocular depth yields the best performance in terms of rendering and reconstruction metrics.

In Fig. 5, we show the benefit of the DSPO on the the valid estimated depth maps \tilde{D} , yielding more consistent depth estimation.

Memory and Runtime. In Tab. 8, we evaluate the peak GPU memory usage, map size and runtime of our method. We achieve a comparable GPU memory usage with GO-SLAM [77] and SplatAM [24]. Our map size is similar to MonoGS [38]. Regarding runtime, we are faster than SplatAM and comparable to MonoGS. GO-SLAM has the

	GO-SLAM [77]	SplatAM [24]	MonoGS [38]	Ours
GPU Usage [GiB]	18.50	18.54	14.62	17.57
Map Size [MB]	-	-	6.8	6.5
Avg. FPS	8.36	0.14	0.32	1.24

Table 8. **Memory and Running Time Evaluation on Replica [54] room0.** Our peak memory usage and runtime are comparable to existing works. We take the numbers from [62] except for ours and MonoGS and we add the Map Size, which denotes the size of the final 3D representation. GPU Usage denotes the peak usage during runtime. All methods are evaluated on an NVIDIA RTX 3090 GPU using single threading for fairness.

fastest runtime, but as shown in Tab. 1 and Tab. 2, it sacrifices rendering and reconstruction quality for speed.

Limitations. We currently do not model the appearance with spherical harmonics, since it only yields a marginal gains in rendering accuracy, while requiring more memory. It is straightforward to add. We only make use of globally optimized frame-to-frame tracking, which fails to leverage frame-to-model queues from the 3D Gaussian map. Another limitation is that our construction of the final proxy depth D is quite simple and does not fuse the monocular and keyframe depths in an informed manner, *e.g.* using normal consistency. Finally, as future work, it is interesting to study how surface regularization can be enforced via *e.g.* quadric surface elements as in [46].

5. Conclusion

We proposed Splat-SLAM, a dense RGB-only SLAM system which uses a deformable 3D Gaussian map for mapping and globally optimized frame-to-frame tracking via optical flow. Importantly, the inclusion of monocular depth into the tracking loop, to refine the scale and to correct the erroneous keyframe depth predictions, leads to better rendering and mapping. By using the monocular depth for completion, mapping is further improved. Our experiments demonstrate that Splat-SLAM outperforms existing solutions regarding reconstruction and rendering accuracy while being on par or better with respect to tracking as well as runtime and memory usage.

References

- [1] Dejan Azinović, Ricardo Martin-Brualla, Dan B Goldman, Matthias Nießner, and Justus Thies. Neural rgb-d surface reconstruction. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6290–6301, 2022. 1
- [2] Michael Bosse, Paul Newman, John Leonard, Martin Soika, Wendelin Feiten, and Seth Teller. An atlas framework for scalable mapping. In *2003 IEEE International Conference on Robotics and Automation (Cat. No. 03CH37422)*, pages 1899–1906. IEEE, 2003. 2
- [3] Aljaž Božič, Pablo Palafox, Justus Thies, Angela Dai, and Matthias Nießner. Transformerfusion: Monocular rgb scene reconstruction using transformers. *arXiv preprint arXiv:2107.02191*, 2021. 1
- [4] 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. 2
- [5] Jiawen Chen, Dennis Bautembach, and Shahram Izadi. Scalable real-time volumetric surface reconstruction. *ACM Transactions on Graphics (ToG)*, 32(4):1–16, 2013. 2
- [6] Hae Min Cho, HyungGi Jo, and Euntai Kim. Sp-slam: Surfel-point simultaneous localization and mapping. *IEEE/ASME Transactions on Mechatronics*, 27(5):2568–2579, 2021. 2
- [7] Sungjoon Choi, Qian-Yi Zhou, and Vladlen Koltun. Robust reconstruction of indoor scenes. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 5556–5565, 2015. 2
- [8] Chi-Ming Chung, Yang-Che Tseng, Ya-Ching Hsu, Xiang-Qian Shi, Yun-Hung Hua, Jia-Fong Yeh, Wen-Chin Chen, Yi-Ting Chen, and Winston H Hsu. Orbeez-slam: A real-time monocular visual slam with orb features and nerf-realized mapping. *arXiv preprint arXiv:2209.13274*, 2022. 1, 2
- [9] Brian Curless and Marc Levoy. Volumetric method for building complex models from range images. In *SIGGRAPH Conference on Computer Graphics*. ACM, 1996. 2
- [10] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. ScanNet: Richly-annotated 3D reconstructions of indoor scenes. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE/CVF, 2017. 5, 6, 7, 8
- [11] 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. 2
- [12] Ainaz Eftekhari, Alexander Sax, Jitendra Malik, and Amir Zamir. Omnidata: A scalable pipeline for making multi-task mid-level vision datasets from 3d scans. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10786–10796, 2021. 3
- [13] Felix Endres, Jürgen Hess, Nikolas Engelhard, Jürgen Sturm, Daniel Cremers, and Wolfram Burgard. An evaluation of the rgb-d slam system. In *2012 IEEE international conference on robotics and automation*, pages 1691–1696. IEEE, 2012. 2
- [14] Jakob Engel, Thomas Schöps, and Daniel Cremers. Lsd-slam: Large-scale direct monocular slam. In *European conference on computer vision*, pages 834–849. Springer, 2014. 2
- [15] Nicola Fioraio, Jonathan Taylor, Andrew Fitzgibbon, Luigi Di Stefano, and Shahram Izadi. Large-scale and drift-free surface reconstruction using online subvolume registration. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4475–4483, 2015. 2
- [16] Peter Henry, Michael Krainin, Evan Herbst, Xiaofeng Ren, and Dieter Fox. Rgb-d mapping: Using kinect-style depth cameras for dense 3d modeling of indoor environments. *The international journal of Robotics Research*, 31(5):647–663, 2012. 2
- [17] Peter Henry, Dieter Fox, Achintya Bhowmik, and Rajiv Mongia. Patch volumes: Segmentation-based consistent mapping with rgb-d cameras. In *2013 International Conference on 3D Vision-3DV 2013*, pages 398–405. IEEE, 2013. 2
- [18] Jiarui Hu, Mao Mao, Hujun Bao, Guofeng Zhang, and Zhaopeng Cui. CP-SLAM: Collaborative neural point-based SLAM system. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. 1, 2
- [19] Tongyan Hua, Haotian Bai, Zidong Cao, and Lin Wang. Fmapping: Factorized efficient neural field mapping for real-time dense rgb slam. *arXiv preprint arXiv:2306.00579*, 2023. 2
- [20] Tongyan Hua, Haotian Bai, Zidong Cao, Ming Liu, Dacheng Tao, and Lin Wang. Hi-map: Hierarchical factorized radiance field for high-fidelity monocular dense mapping. *arXiv preprint arXiv:2401.03203*, 2024. 2
- [21] 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. 1, 2, 5, 6, 7
- [22] Olaf Kähler, Victor Adrian Prisacariu, Carl Yuheng Ren, Xin Sun, Philip H. S. Torr, and David William Murray. Very high frame rate volumetric integration of depth images on mobile devices. *IEEE Trans. Vis. Comput. Graph.*, 21(11): 1241–1250, 2015. 2
- [23] Olaf Kähler, Victor A Prisacariu, and David W Murray. Real-time large-scale dense 3d reconstruction with loop closure. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part VIII 14*, pages 500–516. Springer, 2016. 2
- [24] Nikhil Keetha, Jay Karhade, Krishna Murthy Jatavallabhula, Gengshan Yang, Sebastian Scherer, Deva Ramanan, and Jonathon Luiten. Splatam: Splat, track and map 3d gaussians for dense rgb-d slam. *arXiv preprint*, 2023. 1, 2, 6, 7, 8
- [25] Maik Keller, Damien Lefloch, Martin Lambers, Shahram Izadi, Tim Weyrich, and Andreas Kolb. Real-time 3d reconstruction in dynamic scenes using point-based fusion. In *International Conference on 3D Vision (3DV)*, pages 1–8. IEEE, 2013. 2
- [26] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4), 2023. 4
- [27] Christian Kerl, Jürgen Sturm, and Daniel Cremers. Dense visual slam for rgb-d cameras. In *2013 IEEE/RSJ Interna-*

- tional Conference on Intelligent Robots and Systems*, pages 2100–2106. IEEE, 2013. 2
- [28] 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*, 2023. 2, 5, 7
- [29] Kejie Li, Yansong Tang, Victor Adrian Prisacariu, and Philip HS Torr. Bnv-fusion: Dense 3d reconstruction using bi-level neural volume fusion. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6166–6175, 2022. 1
- [30] Lorenzo Liso, Erik Sandström, Vladimir Yugay, Luc Van Gool, and Martin R Oswald. Loopy-slam: Dense neural slam with loop closures. *arXiv preprint arXiv:2402.09944*, 2024. 1, 2, 7
- [31] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. *Advances in Neural Information Processing Systems*, 33:15651–15663, 2020. 2
- [32] William E Lorensen and Harvey E Cline. Marching cubes: A high resolution 3d surface construction algorithm. *ACM siggraph computer graphics*, 21(4):163–169, 1987. 5
- [33] Mohammad Mahdi Johari, Camilla Carta, and François Fleuret. Eslam: Efficient dense slam system based on hybrid representation of signed distance fields. *arXiv e-prints*, pages arXiv–2211, 2022. 1, 7
- [34] Robert Maier, Jürgen Sturm, and Daniel Cremers. Submap-based bundle adjustment for 3d reconstruction from rgb-d data. In *Pattern Recognition: 36th German Conference, GCPR 2014, Münster, Germany, September 2-5, 2014, Proceedings 36*, pages 54–65. Springer, 2014. 2
- [35] R Maier, R Schaller, and D Cremers. Efficient online surface correction for real-time large-scale 3d reconstruction. *arXiv 2017. arXiv preprint arXiv:1709.03763*, 2017. 2
- [36] Yunxuan Mao, Xuan Yu, Kai Wang, Yue Wang, Rong Xiong, and Yiyi Liao. Ngel-slam: Neural implicit representation-based global consistent low-latency slam system. *arXiv preprint arXiv:2311.09525*, 2023. 2
- [37] Nico Marniok, Ole Johannsen, and Bastian Goldluecke. An efficient octree design for local variational range image fusion. In *German Conference on Pattern Recognition (GCPR)*, pages 401–412. Springer, 2017. 2
- [38] Hidenobu Matsuki, Riku Murai, Paul HJ Kelly, and Andrew J Davison. Gaussian splatting slam. *arXiv preprint arXiv:2312.06741*, 2023. 1, 2, 4, 5, 6, 7, 8
- [39] Hidenobu Matsuki, Edgar Sucar, Tristan Laidow, Kentaro Wada, Raluca Scona, and Andrew J Davison. imode: Real-time incremental monocular dense mapping using neural field. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 4171–4177. IEEE, 2023. 1
- [40] Hidenobu Matsuki, Keisuke Tateno, Michael Niemeyer, and Federic Tombari. Newton: Neural view-centric mapping for on-the-fly large-scale slam. *arXiv preprint arXiv:2303.13654*, 2023. 1, 2
- [41] Jens Naumann, Binbin Xu, Stefan Leutenegger, and Xingxing Zuo. Nerf-vo: Real-time sparse visual odometry with neural radiance fields. *arXiv preprint arXiv:2312.13471*, 2023. 2
- [42] Richard A Newcombe, Shahram Izadi, Otmar Hilliges, David Molyneaux, David Kim, Andrew J Davison, Pushmeet Kohli, Jamie Shotton, Steve Hodges, and Andrew W Fitzgibbon. Kinectfusion: Real-time dense surface mapping and tracking. In *ISMAR*, pages 127–136, 2011. 1, 2
- [43] Matthias Nießner, Michael Zollhöfer, Shahram Izadi, and Marc Stamminger. Real-time 3d reconstruction at scale using voxel hashing. *ACM Transactions on Graphics (TOG)*, 32, 2013. 2
- [44] Helen Oleynikova, Zachary Taylor, Marius Fehr, Roland Siegwart, and Juan I Nieto. Voxblox: Incremental 3d euclidean signed distance fields for on-board MAV planning. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS 2017, Vancouver, BC, Canada, September 24-28, 2017*, pages 1366–1373. IEEE, 2017. 2
- [45] Joseph Ortiz, Alexander Clegg, Jing Dong, Edgar Sucar, David Novotny, Michael Zollhoefer, and Mustafa Mukadam. isdf: Real-time neural signed distance fields for robot perception. *arXiv preprint arXiv:2204.02296*, 2022. 1
- [46] Chensheng Peng, Chenfeng Xu, Yue Wang, Mingyu Ding, Heng Yang, Masayoshi Tomizuka, Kurt Keutzer, Marco Pavone, and Wei Zhan. Q-slam: Quadric representations for monocular slam. *arXiv preprint arXiv:2403.08125*, 2024. 2, 5, 7, 8
- [47] Songyou Peng, Michael Niemeyer, Lars Mescheder, Marc Pollefeys, and Andreas Geiger. Convolutional Occupancy Networks. In *European Conference Computer Vision (ECCV)*. CVF, 2020. 1
- [48] Victor Reijgwart, Alexander Millane, Helen Oleynikova, Roland Siegwart, Cesar Cadena, and Juan Nieto. Voxgraph: Globally consistent, volumetric mapping using signed distance function submaps. *IEEE Robotics and Automation Letters*, 5(1):227–234, 2019. 2
- [49] Antoni Rosinol, John J. Leonard, and Luca Carlone. NeRF-SLAM: Real-Time Dense Monocular SLAM with Neural Radiance Fields. *arXiv*, 2022. 1, 2, 5
- [50] Erik Sandström, Yue Li, Luc Van Gool, and Martin R Oswald. Point-slam: Dense neural point cloud-based slam. In *International Conference on Computer Vision (ICCV)*. IEEE/CVF, 2023. 1, 2, 5
- [51] Erik Sandström, Kevin Ta, Luc Van Gool, and Martin R. Oswald. Uncle-SLAM: Uncertainty learning for dense neural SLAM. In *International Conference on Computer Vision Workshops (ICCVW)*, 2023. 1
- [52] Thomas Schops, Torsten Sattler, and Marc Pollefeys. BAD SLAM: Bundle adjusted direct RGB-D SLAM. In *CVF/IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [53] Frank Steinbrucker, Christian Kerl, and Daniel Cremers. Large-scale multi-resolution surface reconstruction from rgb-d sequences. In *IEEE International Conference on Computer Vision*, pages 3264–3271, 2013. 2
- [54] 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

- [55] Jörg Stückler and Sven Behnke. Multi-resolution surfel maps for efficient dense 3d modeling and tracking. *Journal of Visual Communication and Image Representation*, 25(1):137–147, 2014. 2
- [56] Jürgen Sturm, Nikolas Engelhard, Felix Endres, Wolfram Burgard, and Daniel Cremers. A benchmark for the evaluation of RGB-D SLAM systems. In *International Conference on Intelligent Robots and Systems (IROS)*. IEEE/RSJ, 2012. 1, 5, 6, 7, 8
- [57] Edgar Sucar, Shikun Liu, Joseph Ortiz, and Andrew J. Davison. iMAP: Implicit Mapping and Positioning in Real-Time. In *International Conference on Computer Vision (ICCV)*. IEEE/CVF, 2021. 1, 5
- [58] Jiaming Sun, Yiming Xie, Linghao Chen, Xiaowei Zhou, and Hujun Bao. Neuralrecon: Real-time coherent 3d reconstruction from monocular video. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15598–15607, 2021. 1
- [59] Yijie Tang, Jiazhao Zhang, Zhinan Yu, He Wang, and Kai Xu. Mips-fusion: Multi-implicit-submaps for scalable and robust online neural rgb-d reconstruction. *arXiv preprint arXiv:2308.08741*, 2023. 2
- [60] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*, pages 402–419. Springer, 2020. 2
- [61] 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, 3, 8
- [62] Fabio Tosi, Youmin Zhang, Ziren Gong, Erik Sandström, Stefano Mattoccia, Martin R. Oswald, and Matteo Poggi. How nerfs and 3d gaussian splatting are reshaping slam: a survey, 2024. 2, 5, 8
- [63] Hao Wang, Jun Wang, and Wang Liang. Online reconstruction of indoor scenes from rgb-d streams. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3271–3279, 2016. 2
- [64] 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/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13293–13302, 2023. 7
- [65] 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
- [66] Silvan Weder, Johannes Schonberger, Marc Pollefeys, and Martin R Oswald. Routedfusion: Learning real-time depth map fusion. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4887–4897, 2020. 1
- [67] Silvan Weder, Johannes L Schonberger, Marc Pollefeys, and Martin R Oswald. Neurralfusion: Online depth fusion in latent space. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3162–3172, 2021. 1
- [68] Thomas Whelan, Stefan Leutenegger, Renato Salas-Moreno, Ben Glocker, and Andrew Davison. Elasticfusion: Dense slam without a pose graph. In *Robotics: Science and Systems (RSS)*, 2015. 2
- [69] Chi Yan, Delin Qu, Dong Wang, Dan Xu, Zhigang Wang, Bin Zhao, and Xuelong Li. Gs-slam: Dense visual slam with 3d gaussian splatting. *arXiv preprint arXiv:2311.11700*, 2023. 1, 2, 7
- [70] Zhixin Yan, Mao Ye, and Liu Ren. Dense visual slam with probabilistic surfel map. *IEEE transactions on visualization and computer graphics*, 23(11):2389–2398, 2017. 2
- [71] 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 *IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*, pages 499–507. IEEE, 2022. 1, 2
- [72] Xingrui Yang, Yuhang Ming, Zhaopeng Cui, and Andrew Calway. Fd-slam: 3-d reconstruction using features and dense matching. In *2022 International Conference on Robotics and Automation (ICRA)*, pages 8040–8046. IEEE, 2022. 2
- [73] Vladimir Yugay, Yue Li, Theo Gevers, and Martin R. Oswald. Gaussian-slam: Photo-realistic dense slam with gaussian splatting, 2023. 1, 2, 6, 7
- [74] Heng Zhang, Guodong Chen, Zheng Wang, Zhenhua Wang, and Lining Sun. Dense 3d mapping for indoor environment based on feature-point slam method. In *2020 the 4th International Conference on Innovation in Artificial Intelligence*, pages 42–46, 2020. 2
- [75] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018. 5
- [76] Wei Zhang, Tiecheng Sun, Sen Wang, Qing Cheng, and Norbert Haala. Hi-slam: Monocular real-time dense mapping with hybrid implicit fields. *IEEE Robotics and Automation Letters*, 2023. 1, 2, 5, 7, 8
- [77] Youmin Zhang, Fabio Tosi, Stefano Mattoccia, and Matteo Poggi. Go-slam: Global optimization for consistent 3d instant reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3727–3737, 2023. 1, 2, 4, 5, 6, 7, 8
- [78] Heng Zhou, Zhetao Guo, Shuhong Liu, Lechen Zhang, Qihao Wang, Yuxiang Ren, and Mingrui Li. Mod-slam: Monocular dense mapping for unbounded 3d scene reconstruction, 2024. 2, 5, 7
- [79] 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 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12786–12796, 2022. 1, 7
- [80] 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. *arXiv preprint arXiv:2302.03594*, 2023. 2, 5
- [81] Zi-Xin Zou, Shi-Sheng Huang, Yan-Pei Cao, Tai-Jiang Mu, Ying Shan, and Hongbo Fu. Mononeurralfusion: Online monocular neural 3d reconstruction with geometric priors. *arXiv preprint arXiv:2209.15153*, 2022. 1
- [82] Matthias Zwicker, Hanspeter Pfister, Jeroen Van Baar, and Markus Gross. Surface splatting. In *Proceedings of the 28th*

annual conference on Computer graphics and interactive techniques, pages 371–378, 2001. [4](#)