

Flow4DGS-SLAM: Optical Flow-Guided 4D Gaussian Splatting SLAM

Yunsong Wang Gim Hee Lee

School of Computing, National University of Singapore

yunsong@comp.nus.edu.sg gimhee.lee@nus.edu.sg

<https://github.com/wangys16/Flow4DGS-SLAM>

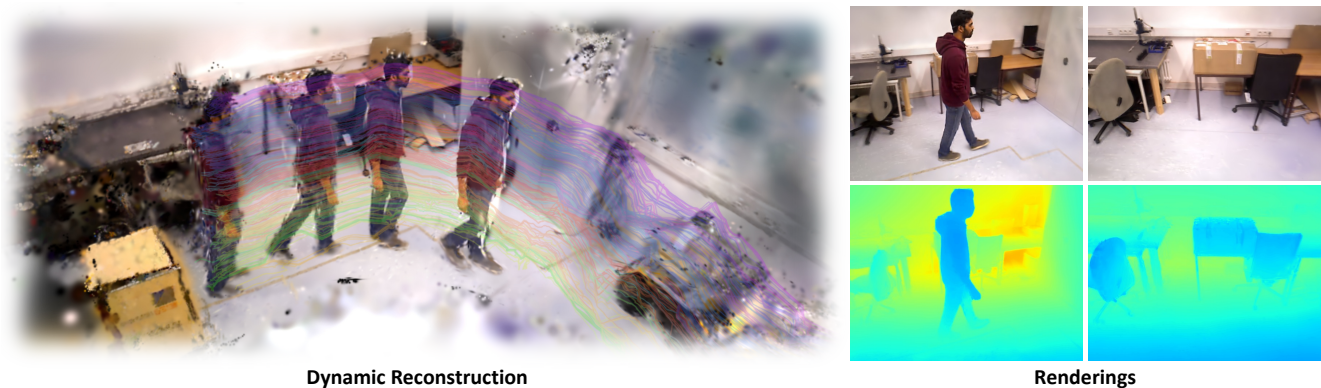


Figure 1. **Overview of our results.** Our method achieves high-quality renderings with spatially and temporally coherent Gaussian motion.

Abstract

Handling the dynamic environments is a significant research challenge in Visual Simultaneous Localization and Mapping (SLAM). Recent research combines 3D Gaussian Splatting (3DGS) with SLAM to achieve both robust camera pose estimation and photorealistic renderings. However, using SLAM to efficiently reconstruct both static and dynamic regions remains challenging. In this work, we propose an efficient framework for dynamic 3DGS SLAM guided by optical flow. Using the input depth and prior optical flow, we first propose a category-agnostic motion mask generation strategy by fitting a camera ego-motion model to decompose the optical flow. This module separates dynamic and static Gaussians and simultaneously provides flow-guided camera pose initialization. We boost the training speed of dynamic 3DGS by explicitly modeling their temporal centers at keyframes. These centers are propagated using 3D scene flow priors and are dynamically initialized with an adaptive insertion strategy. Alongside this, we model the temporal opacity and rotation using a Gaussian Mixture Model (GMM) to adaptively learn the complex dynamics. The empirical results demonstrate our state-of-the-art performance in tracking, dynamic reconstruction, and training efficiency.

1. Introduction

Visual Simultaneous Localization and Mapping (SLAM) remains a fundamental yet challenging research task. It jointly estimates camera motion and builds a 3D scene map to enable applications such as robotics, AR/VR, and autonomous systems. However, most existing systems still assume a static world [2, 17, 26, 33, 39, 48, 57]. In real-world dynamic environments, this assumption often causes drift or tracking failure. Recently, 3D Gaussian Splatting (3DGS) [19] has been integrated into SLAM to achieve real-time, photorealistic rendering and explicit 3D mapping [18, 25, 47, 53]. These methods improve convergence in pose optimization and support online mapping, but typically treat dynamic objects as outliers, removing them to stabilize tracking [20, 46, 55]. Therefore, the final maps capture only the static backgrounds.

Dynamic reconstruction studies have extended 3DGS to the temporal domain by learning MLP-based deformation field [13, 15, 23, 50], parameterizing Gaussian motion [22, 40], or directly modeling temporal offsets [11, 31, 49]. However, these approaches typically rely on precomputed multi-view poses and lengthy offline training. Therefore, combining dynamic 3DGS with SLAM offers great potential for joint camera tracking and efficient dynamic reconstruction.

One prior work, 4DGS-SLAM [21], represents an initial

step in this direction. It leverages sparse control points from SC-GS [15] to reconstruct dynamic elements, and estimates camera motion using the static 3DGS. However, its deformation field training is computationally expensive, hindering real-time performance. Moreover, its category-based segmentation model struggles in more general dynamic scenes, and its reconstruction capability is limited in handling complex dynamics (*e.g.*, people leaving and re-entering the view). To address these challenges, we propose **Flow4DGS-SLAM**, a dynamic SLAM framework guided by optical flow for efficient tracking and reconstruction in complex dynamic environments.

To handle more general dynamic scenes, we first design a Camera-Induced Motion Decomposition module to perform category-agnostic motion segmentation. This module leverages the input depth map and prior optical flow to efficiently solve for the camera ego-motion, estimating the rigid flow induced by camera motion and filter out the outlier dynamic pixels. We further represent the dynamic Gaussians in a hybrid form, with explicitly learned Gaussian centers at keyframes and parametric Gaussian-Mixture Model (GMM)-based temporal opacity and rotation. With the explicit Gaussian centers, we take inspiration from GFlow [41] to leverage the off-the-shelf optical flow models for explicit propagation to efficiently model dynamics of the Gaussians, which is combined with a KNN-based smoothing strategy to ensure local rigidity. Furthermore, we leverage optical flow to back-track newly appearing motion regions, in order to adaptively insert new dynamic Gaussians. Additionally, the GMM-based temporal opacity and rotation improves the capability to reconstruct complex dynamic scenes without significantly enlarging the model size.

Our **contributions** can be summarized as follows:

1. **Camera-Induced Motion Decomposition.** We design a parametric model that efficiently generates a category-agnostic motion mask to handle more general dynamic scenes. The model fits a 6-DoF camera motion to the depth map and optical flow, providing a better camera initialization and filtering out the outlier dynamic pixels.
2. **Hybrid 4DGS guided by optical flow.** We propose a hybrid 4D Gaussian representation with explicit temporal positions and GMM-based opacity and rotation. The explicit position modeling is guided by a scene-flow Gaussian propagation module and an adaptive gaussian insertion module to accelerate training. The GMM-based opacity and rotation improves the representation capability of the model to reconstruct complex dynamic scenes.
3. **Superior Performance.** The empirical results demonstrate our state-of-the-art performance in tracking, photorealistic rendering, and computational efficiency.

2. Related Work

Static SLAM. Classic visual SLAM decomposes tracking and mapping, relying on sparse keypoints to estimate camera motion [2, 26]. Dense SLAM extends this by reconstructing per-frame geometry from depth maps anchored to keyframes [28, 34, 37, 56]. The frame-centric methods [9, 36] efficiently estimate frame-to-frame camera motion, but may face inconsistencies when assembling into a global map. In contrast, map-centric methods use a unified 3D map representation, such as surfels [30, 44] or voxel-based TSDF grids [8, 10, 27]. More recent works leverage NeRF or 3DGS [18, 25, 33, 39, 57] via differentiable rendering for photorealistic tracking and mapping. Although effective in static environments, these systems often fail in dynamic scenes where moving objects significantly hinder tracking.

Gaussian Splatting. 3D Gaussian Splatting (3DGS) [19] achieves real-time rendering by optimizing anisotropic Gaussians, showing strong performance in surface reconstruction [5, 6, 12], feed-forward reconstruction [3, 7, 42, 43], and SLAM [18, 25]. To handle dynamic scenes, some subsequent works learn temporal deformation fields through MLPs [13, 15, 23, 50], leveraging the low-frequency bias of MLPs for smooth temporal and spatial deformation. For instance, SC-GS [15] trains sparse control points and an MLP-based deformation field where Gaussian motion is computed based on weighted sum of its neighboring control points’ motion, ensuring spatial rigidity and temporal coherence. However, such approaches require hours of optimization due to costly deformation MLPs training. Others optimize Gaussian trajectories parametrically [22, 40] or directly learn temporal Gaussian parameters [24, 49]. Since dynamic SLAM requires efficient reconstruction from sparse keyframes, we propose a hybrid dynamic Gaussian representation that combines the efficiency of explicit modeling with the temporal coherence of parametric approaches.

SLAM in Dynamic Scenes. Typical SLAM systems handle dynamics by detecting and excluding moving regions using geometric residuals or semantic priors [1, 46, 51, 54, 55]. For example, WildGS-SLAM [55] employs DINO [52] features and an uncertainty MLP to distinguish dynamic from static regions, stabilizing tracking while yielding only static background maps. ADD-SLAM [45] reconstructs both the static and dynamic regions on keyframes, but lacks deformation of the dynamic Gaussians for temporal view interpolation. 4DGS-SLAM [21] advances this by reconstructing dynamic Gaussians online via a sparse-control-point deformation field, but faces intensive computational cost and becomes less effective in more general and complex dynamic scenes. In contrast, our method explicitly models Gaussian motion over time with scene-flow priors, enabling fast and robust reconstruction of both static and dynamic elements within the SLAM pipeline.

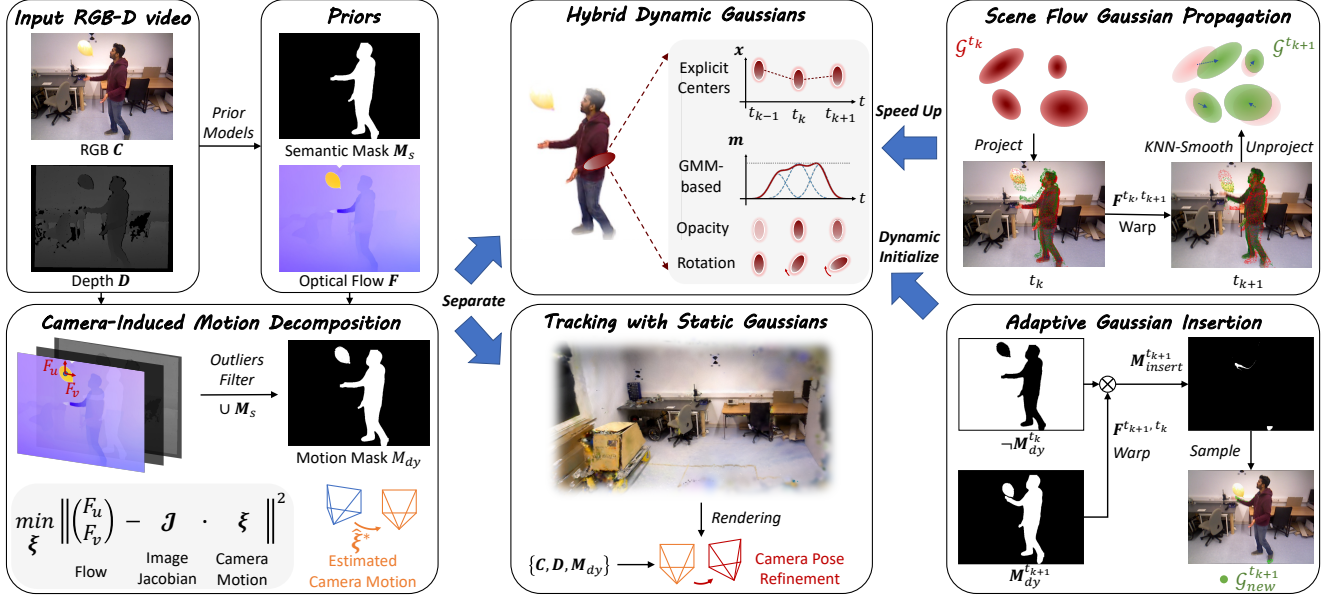


Figure 2. **Overall framework of Flow4DGS-SLAM.** Given input RGB-D video, we first extract the prior semantic mask and optical flow, and feed them into a camera-induced motion decomposition module to filter out category-agnostic motion mask and solve an optical-flow guided camera initialization. The static Gaussians help refine the camera pose during tracking, and the dynamic Gaussians are represented in a hybrid form, combined with a scene flow Gaussian propagation module and an adaptive Gaussian insertion module to accelerate training.

3. Our Methodology

Overview. As shown in Figure 2, our framework performs online tracking of camera motion and mapping of a 4D scene map. For tracking, we propose a Camera-Induced Motion Decomposition approach to generate a category-agnostic motion mask, which is combined with the semantic mask generated by YOLOv9 [38] to avoid dynamic distractors during tracking. For dynamic mapping, we represent dynamic Gaussians in a hybrid form with explicit temporal positions. This allows us to leverage prior optical flow for explicit Gaussian propagation and adaptive Gaussian insertion. In parallel, the time-varying opacities and rotations of dynamic Gaussians are modeled with a Gaussian Mixture Model for the reconstruction of more complex dynamic scenes.

3.1. Camera-Induced Motion Decomposition

Category-agnostic Masking. We leverage optical flow and depth to fit a camera ego-motion model. This model predicts the rigid flow induced by camera motion, which allows us to identify and filter out dynamic pixels. The estimated camera motion also provides a coarse initialization of the new camera pose. Specifically, given the new RGB-D frame $\{\mathbf{I}^t, \mathbf{D}^t\}$ at timestamp t , we first feed \mathbf{I}^t and \mathbf{I}^{t-1} into RAFT [35] to obtain optical flow $\mathbf{F}^{t,t-1}(u, v)$ which we denote as $\mathbf{F}(u, v)$ for brevity. Let $Z = \mathbf{D}^t(u, v)$ and $\mathbf{x} = (u, v, Z)^\top$ denote the corresponding 3D point. Under small camera motion

from t to $t-1$, the motion field of a static 3D point [14] is:

$$\mathbf{F}(u, v) = \mathbf{J}(\mathbf{x}) \boldsymbol{\xi}, \quad (1)$$

where the camera twist $\boldsymbol{\xi} = [\boldsymbol{\rho}^\top, \boldsymbol{\theta}^\top]^\top \in \mathbb{R}^6$ contains translation and rotation. The 2×6 image Jacobian \mathbf{J} follows projective geometry [4]:

$$\mathbf{J}(\mathbf{x}) = \begin{bmatrix} -\frac{f_x}{Z} & 0 & \frac{u}{Z} & \frac{uv}{f_y} & -f_x - \frac{u^2}{f_x} & v \\ 0 & -\frac{f_y}{Z} & \frac{v}{Z} & f_y + \frac{v^2}{f_y} & -\frac{uv}{f_x} & -u \end{bmatrix}, \quad (2)$$

where (f_x, f_y) is the camera intrinsics for brevity. The translation-induced flow scales with $1/Z$, while rotation-induced flow is depth-independent. We further obtain a semantic mask \mathcal{M}_s from YOLOv9 [38], where $\mathcal{M}_s(u, v) = 1$ indicates dynamic objects (people). Assuming that most pixels with $\mathcal{M}_s(u, v) = 0$ are static, we stack Eq. (1) over all pixels (u, v) with $Z(u, v) > 0$ and $\mathcal{M}_s(u, v) = 0$, and solve:

$$\hat{\boldsymbol{\xi}} = \arg \min_{\boldsymbol{\xi}} \sum_i w_i \|\mathbf{F}_i - \mathbf{J}_i \boldsymbol{\xi}\|^2, \quad (3)$$

using IRLS with Cauchy weights $\{w_i\}$. The predicted rigid flow is $\hat{\mathbf{F}}(u, v) = \mathbf{J}(u, v, Z) \hat{\boldsymbol{\xi}}$. We compute residuals:

$$r(u, v) = \|\mathbf{F}(u, v) - \hat{\mathbf{F}}(u, v)\|_2. \quad (4)$$

Since dynamic pixels tend to yield large residuals, we obtain the category-agnostic mask using the threshold:

$$\mathcal{M}_{ca}(u, v) = \mathbb{1}(r(u, v) > \text{median}(r) + k \text{MAD}(r)), \quad (5a)$$

$$\text{where } \text{MAD}(r) = \text{median}_i |r_i - \text{median}_j(r_j)|. \quad (5b)$$

The final dynamic mask is $\mathcal{M}_{dy} = \mathcal{M}_s \cup \mathcal{M}_{ca}$.

Camera Initialization. To solve for a more accurate camera motion, we use pixels with $\mathcal{M}_{dy}(u, v) = 0$ and valid depth to perform a second weighted least-squares refinement for a cleaner twist $\hat{\xi}^*$. We map $\hat{\xi}^*$ to SE(3) via the exponential map:

$$\exp_{\text{se}(3)}(\hat{\xi}^*) = \begin{bmatrix} \exp_{\text{so}(3)}(\theta^*) & \mathbf{V}(\theta^*) \rho^* \\ \mathbf{0}^\top & 1 \end{bmatrix}, \quad (6)$$

where $\mathbf{V}(\cdot)$ is the left Jacobian of SO(3). With previous pose $\mathbf{T}_{cw}^{t-1} = [\mathbf{R}^{t-1} \mid \mathbf{t}^{t-1}]$, we initialize the current pose as:

$$\mathbf{T}_{cw}^t = \mathbf{T}_{cw}^{t-1} \exp_{\text{se}(3)}(\hat{\xi}^*). \quad (7)$$

In experiments, we clamp the maximum camera motion with a threshold proportional to the ratio of inlier pixels to improve the robustness to optical flow noise. Our camera-induced motion decomposition leverages prior optical flow and the image Jacobian to efficiently generate a category-agnostic motion mask and provide a better camera pose initialization. As a result, the tracking robustness of dynamic SLAM is significantly improved in complex dynamic scenes.

3.2. Hybrid 4D Gaussian Splatting Representation

Preliminary of 3DGS. Each static 3D Gaussian \mathcal{G}_i^s is defined by a center position $\mathbf{x}_i^s \in \mathbb{R}^3$, a covariance matrix $\Sigma_i^s \in \mathbb{R}^{3 \times 3}$, an opacity $\sigma_i^s \in [0, 1]$, and color parameters \mathbf{c}_i^s . During rendering, each Gaussian is projected onto the image plane and accumulated in a back-to-front order using an α -blending scheme. The rendered color map is:

$$\hat{\mathbf{C}}^s(\mathbf{u}) = \sum_{i=1}^{|\mathcal{G}|} \mathbf{c}_i^s \alpha_i^s(\mathbf{u}) \prod_{j < i} (1 - \alpha_j^s(\mathbf{u})), \quad (8)$$

where $\alpha_i^s(\mathbf{u})$ represents the elliptical footprint of \mathcal{G}_i^s at pixel \mathbf{u} multiplied by σ_i^s . The depth map $\hat{\mathbf{D}}^s(\mathbf{u})$ and opacity map $\hat{\mathbf{O}}^s(\mathbf{u})$ are rendered similarly by replacing \mathbf{c}_i^s in Eq. (8) with depth and 1, respectively.

Hybrid 4DGS. We separate the 3D Gaussians into static and dynamic, which are defined and modeled differently. The static Gaussians are defined similarly as in 3DGS [19]. We define the 4D Gaussian Splatting in a hybrid form, including explicit time-varying trajectories and Gaussian Mixture Model (GMM)-based temporal opacities and rotations. The i -th dynamic Gaussian is defined by a set of static attributes $\{\mathbf{s}_i, \sigma_i, \mathbf{c}_i\}$, where \mathbf{s}_i is its scale, σ_i is its static opacity, and \mathbf{c}_i denotes color features. The dynamic attributes consist of: i) explicit positions $\mathbf{x}_i(t)$ at keyframes, and ii) temporally continuous opacity $\sigma_i(t)$ and rotation $\mathbf{q}_i(t)$.

Explicit Keyframe Positions. To accelerate training, guided by the prior optical flow introduced in Section 3.4, we first discretize the temporal domain into a small number

of keyframes $\{t_k\}$. For the i -th dynamic Gaussian, we learn its 3D center \mathbf{x}_i^k at each keyframe. At arbitrary time t , the position is obtained via linear interpolation, which supports flexible motion and explicit operation for efficient training.

GMM opacity and rotation. To represent the temporal opacity and rotations smoothly without timestamp-specific storage, we model both the opacity $\sigma_i(t)$ and rotation $\mathbf{q}_i(t)$ using a Gaussian mixture model over normalized time $\hat{t} \in [0, 1]$. For each Gaussian we learn K mixture components with learnable weights $w_{i,k}$, temporal means $\mu_{i,k}$ and scales $\tau_{i,k}$. The opacity coefficient is defined as:

$$m_i(t) = 1 - \exp\left(-A_i \sum_{k=1}^K w_{i,k} \mathcal{N}(\hat{t}; \mu_{i,k}, \tau_{i,k}^2)\right), \quad (9)$$

where $A_i > 0$ is a learnable activation amplitude. The final time-varying opacity is $\sigma_i(t) = \sigma_i \cdot m_i(t)$.

Each component k also stores a learnable control quaternion $\mathbf{q}_{i,k}$, and the Gaussian activations serve as blending weights to achieve smooth rotation:

$$\mathbf{q}_i(t) = \frac{\sum_{k=1}^K w_{i,k} \mathcal{N}(\hat{t}; \mu_{i,k}, \tau_{i,k}^2) \mathbf{q}_{i,k}}{\left\| \sum_{k=1}^K w_{i,k} \mathcal{N}(\hat{t}; \mu_{i,k}, \tau_{i,k}^2) \mathbf{q}_{i,k} \right\|}. \quad (10)$$

This formulation produces continuous and smoothly varying rotations. We set $K = 3$ in the experiments.

In our design, explicit keyframe positions are used by the optical-flow-assisted 4D mapping (*cf.* Sec. 3.4) to accelerate training. In addition, GMM-based temporal opacity and rotation ensures smooth temporal behavior with minimal model-size overhead.

3.3. Tracking

Given the previously computed dynamic mask \mathcal{M}_{dy} , we follow 4DGS-SLAM [21] to optimize the camera pose based on L_1 loss between the static-3DGS-rendered color and depth maps and the inputs. Specifically, we first compute an opacity mask and the final valid mask as:

$$\mathcal{M}_o(\mathbf{u}) = \mathbf{1}(\hat{\mathbf{O}}(\mathbf{u}) \geq \alpha), \quad (11a)$$

$$\mathcal{M}_v = (\neg \mathcal{M}_{dy}) \cap \mathcal{M}_o. \quad (11b)$$

The tracking loss is then computed as:

$$\mathcal{L}_{track} = \frac{1}{|\mathcal{V}|} \sum_{\mathbf{u} \in \mathcal{V}} \mathcal{M}_v(\mathbf{u}) \left(\lambda_1 L_1(\hat{\mathbf{C}}(\mathbf{u})) + \lambda_2 L_1(\hat{\mathbf{D}}(\mathbf{u})) \right), \quad (12)$$

where \mathcal{V} is the set of pixels with valid depths. Note that before tracking, the camera pose is initialized with the estimated camera motion from Sec. 3.1 to form a coarse-to-fine pipeline for improved robustness (*cf.* Figure 4).

Table 1. **Trajectory ATE RMSE [cm]**↓ on the TUM RGB-D sequences. Best results are shown in **bold**.

Method	fr3/sit_st	fr3/sit_xyz	fr3/sit_rpy	fr3/walk_st	fr3/walk_xyz	fr3/walk_rpy	Avg.
RoDyn-SLAM [16]	1.5	5.6	5.7	1.7	8.3	8.1	5.1
MonoGS [25]	0.48	1.7	6.1	21.9	30.7	34.2	15.8
SplaTAM [18]	0.52	1.5	11.8	83.2	134.2	142.3	62.2
4DGS-SLAM [21]	0.58	2.9	2.6	0.61	2.7	3.0	2.1
Ours	0.70	1.8	2.4	0.48	2.5	3.6	1.9

Table 2. **Quantitative results on the TUM RGB-D sequences**. Best results are highlighted in **bold**.

Method	Metric	fr3/sit_st	fr3/sit_xyz	fr3/sit_rpy	fr3/walk_st	fr3/walk_xyz	fr3/walk_rpy	Avg.
MonoGS [25]	PSNR[dB]↑	19.95	23.92	16.99	16.47	14.02	15.12	17.74
	SSIM↑	0.739	0.803	0.572	0.604	0.436	0.497	0.608
	LPIPS↓	0.213	0.182	0.405	0.355	0.581	0.560	0.382
SplaTAM [18]	PSNR[dB]↑	24.12	22.07	19.97	16.70	17.03	16.54	19.40
	SSIM↑	0.915	0.879	0.799	0.688	0.650	0.635	0.757
	LPIPS↓	0.101	0.163	0.205	0.287	0.339	0.353	0.241
SC-GS [15]	PSNR[dB]↑	27.01	21.45	18.93	20.99	19.89	16.44	20.78
	SSIM↑	0.900	0.686	0.529	0.762	0.590	0.475	0.657
	LPIPS↓	0.182	0.369	0.512	0.291	0.470	0.554	0.396
4DGS-SLAM [21]	PSNR[dB]↑	27.68	24.37	20.71	23.21	19.83	19.47	22.55
	SSIM↑	0.892	0.822	0.746	0.822	0.730	0.713	0.788
	LPIPS↓	0.116	0.179	0.265	0.192	0.281	0.340	0.229
Ours	PSNR[dB]↑	29.70	27.49	26.04	27.64	24.60	23.83	26.55
	SSIM↑	0.905	0.860	0.812	0.852	0.788	0.769	0.831
	LPIPS↓	0.092	0.150	0.223	0.127	0.206	0.263	0.177

3.4. Optical Flow-Guided 4D Mapping

We introduce a scene-flow-based Gaussian propagation module and an adaptive Gaussian insertion module, to accelerate training of dynamic Gaussians for online SLAM and efficiently handle newly appearing dynamic objects. Both modules explicitly operate on Gaussian positions.

Scene Flow Gaussian Propagation. Before mapping keyframe k , we maintain the current set of dynamic Gaussians $\mathcal{G}^{t_{k-1}}$ with 3D centers $\{\mathbf{x}_i^{k-1}\}$. We then initialize the dynamic Gaussians at keyframe k by propagating these centers using the prior scene flow. Let $\mathbf{K} \in \mathbb{R}^{3 \times 3}$ be the camera intrinsic matrix, and $(\mathbf{R}_k, \mathbf{t}_k)$ denote the optimized rotation and translation of keyframe- k camera, we get the projection matrix $\mathbf{P}_k = \mathbf{K}[\mathbf{R}_k \mid \mathbf{t}_k]$.

Each Gaussian center is first projected into the image of keyframe $k-1$: $\bar{\mathbf{u}}_i^{k-1} = \mathbf{P}_{k-1} \begin{bmatrix} \mathbf{x}_i^{k-1} \\ 1 \end{bmatrix}$, $\mathbf{u}_i^{k-1} = \Pi(\bar{\mathbf{u}}_i^{k-1})$, where $\Pi(\cdot)$ performs homogeneous normalization. With the predicted optical flow $\mathbf{F}^{t_{k-1}, t_k}$, we propagate these projections to keyframe k as $\mathbf{u}_i^k = \mathbf{u}_i^{k-1} + \mathbf{F}^{t_{k-1}, t_k}(\mathbf{u}_i^{k-1})$. Subsequently, we obtain a coarse estimate of the 3D deformation from time $k-1$ to k by unprojection:

$$\Delta \mathbf{x}_i^k = \mathbf{R}_k^\top (\mathbf{D}_i^k \mathbf{K}^{-1} \bar{\mathbf{u}}_i^k - \mathbf{t}_k) - \mathbf{x}_i^{k-1}, \quad (13)$$

where $\bar{\mathbf{u}}_i^k = [\mathbf{u}_i^k, 1]^\top$ is the homogeneous pixel location.

We mitigate noise in the flow estimation and enforce local

rigidity by regularizing $\Delta \mathbf{x}_i^k$ with KNN smoothing:

$$\Delta \hat{\mathbf{x}}_i^k = \sum_{j \in \mathcal{N}(i)} w_{ij}^{knn} \Delta \mathbf{x}_j^k, \quad (14)$$

$$w_{ij}^{knn} = \frac{\mathcal{N}(\|\mathbf{x}_j^{k-1} - \mathbf{x}_i^{k-1}\|_2; 0, \tau_{knn}^2)}{\sum_{l \in \mathcal{N}(i)} \mathcal{N}(\|\mathbf{x}_l^{k-1} - \mathbf{x}_i^{k-1}\|_2; 0, \tau_{knn}^2)}, \quad (15)$$

where $\mathcal{N}(i)$ is the nearest neighbors set with a radius of the i -th dynamic Gaussian. Consequently, the corresponding dynamic Gaussian center at keyframe k is initialized as $\mathbf{x}_i^k = \mathbf{x}_i^{k-1} + \Delta \hat{\mathbf{x}}_i^k$. This yields motion-aware and spatially consistent Gaussian initialization, which substantially accelerates convergence in the subsequent 4D optimization.

Adaptive Gaussian Insertion. Although the Gaussian propagation updates existing dynamic Gaussians to a better initialization at the new keyframe, new dynamic objects or previously occluded regions of dynamic objects may appear in keyframe k that was absent in keyframe $k-1$. To ensure complete coverage and fast convergence on such dynamic regions, we initialize new Gaussians based on optical flow back-tracking.

Let \mathcal{M}_{dy}^k and \mathcal{M}_{dy}^{k-1} denote the binary motion masks of keyframes k and $k-1$, respectively. We first warp the current mask into the previous keyframe using the predicted backward flow $\mathbf{F}^{t_k, t_{k-1}}$: $\mathbf{u}_p^{k-1} = \mathbf{u}_p^k + \mathbf{F}^{t_k, t_{k-1}}(\mathbf{u}_p^k)$, for each pixel $\mathbf{u}_p^k \in \mathcal{M}_{dy}^{t_k}$.

Table 3. **Trajectory ATE RMSE [cm]↓ on the BONN sequences.** Best results are highlighted in **bold**.

Method	ballon	ballon2	ps_track	ps_track2	sync	sync2	p_no_box	p_no_box2	p_no_box3	Avg.
RoDyn-SLAM [16]	7.9	11.5	14.5	13.8	1.3	1.4	4.9	6.2	10.2	7.9
MonoGS [25]	29.6	22.1	54.5	36.9	68.5	0.56	71.5	10.7	3.6	33.1
SplaTAM [18]	32.9	30.4	77.8	116.7	59.5	66.7	91.9	18.5	17.1	56.8
4DGS-SLAM [21]	2.5	4.2	9.8	11.2	0.95	0.64	1.8	1.7	2.6	3.9
Ours	2.6	3.4	7.2	10.8	0.60	0.64	2.2	1.6	2.3	3.5

Table 4. **Quantitative results on the BONN sequences.** Best scores are shown in **bold**. “-” indicates reconstruction failure.

Method	Metric	ballon	ballon2	ps_track	ps_track2	sync	sync2	p_no_box	p_no_box2	p_no_box3	Avg.
MonoGS [25]	PSNR[dB]↑	21.35	20.22	20.53	20.09	22.03	20.55	20.76	19.38	24.81	21.06
	SSIM↑	0.803	0.758	0.779	0.718	0.766	0.841	0.748	0.753	0.857	0.780
	LPIPS↓	0.316	0.354	0.408	0.426	0.328	0.521	0.428	0.372	0.243	0.342
SplaTAM [18]	PSNR[dB]↑	19.65	17.67	18.30	15.57	19.33	19.67	20.81	21.69	21.41	19.34
	SSIM↑	0.781	0.702	0.670	0.606	0.776	0.730	0.824	0.852	0.873	0.757
	LPIPS↓	0.211	0.280	0.283	0.331	0.227	0.258	0.191	0.165	0.152	0.233
SC-GS [15]	PSNR[dB]↑	22.30	21.38	-	-	23.62	22.74	20.60	21.55	19.24	21.63
	SSIM↑	0.737	0.708	-	-	0.788	0.801	0.688	0.722	0.628	0.724
	LPIPS↓	0.448	0.450	-	-	0.427	0.359	0.515	0.491	0.539	0.461
4DGS-SLAM [21]	PSNR[dB]↑	26.19	22.91	21.78	20.65	24.37	25.12	23.14	24.28	25.88	23.81
	SSIM↑	0.875	0.839	0.832	0.820	0.833	0.892	0.845	0.873	0.886	0.855
	LPIPS↓	0.235	0.269	0.289	0.294	0.230	0.173	0.239	0.224	0.207	0.240
Ours	PSNR[dB]↑	30.30	28.36	29.18	28.94	29.85	31.26	29.25	30.03	30.22	29.71
	SSIM↑	0.897	0.877	0.879	0.882	0.874	0.924	0.883	0.896	0.900	0.890
	LPIPS↓	0.193	0.204	0.228	0.211	0.165	0.140	0.204	0.205	0.187	0.193

Pixels of new dynamic regions are identified by:

$$\mathcal{M}_{\text{insert}}^{t_k} = \left\{ \mathbf{u}_p^k \in \mathcal{M}_{dy}^{t_k} \mid \mathbf{u}_p^{k-1} \notin \mathcal{M}_{dy}^{t_{k-1}} \right\}. \quad (16)$$

We randomly sample these newly activated dynamic pixels with a density factor $1/D_{\text{init}}$ and initialize new dynamic Gaussians by unprojecting sampled pixels.

Optimization. We conduct a fast training step after the Gaussian propagation and adaptive gaussian insertion steps. We follow 4DGS-SLAM [21] to train on a sliding window of keyframes and 2 random previous keyframes, but for much fewer iterations (which we set to 50 in the experiments). Finally, the mapping loss is:

$$\mathcal{L}_{\text{map}} = \lambda_1 \mathcal{L}_c + \lambda_2 \mathcal{L}_d + \lambda_f \mathcal{L}_f + \lambda_m \mathcal{L}_m + \lambda_{iso} \mathcal{L}_{iso}, \quad (17)$$

where \mathcal{L}_c , \mathcal{L}_d , \mathcal{L}_f , \mathcal{L}_{iso} are the color, depth, flow, and isotropic loss, respectively. \mathcal{L}_m is a binary mask loss which constrains rendered dynamic Gaussians’ alpha map to be consistent with motion mask. λ_f , λ_m and λ_{iso} are the corresponding weights. The flow loss and isotropic loss are adopted from 4DGS-SLAM [21].

4. Experiments

4.1. Experimental Settings

Datasets. We evaluate our method in two real-world datasets: TUM RGB-D dataset [32] and BONN dataset [29]. We

mainly evaluate on the dynamic sequences following 4DGS-SLAM [21] to form a direct comparison.

Implementation Details. We ran our experiments on single NVIDIA RTX A6000 GPU. For tracking, we follow 4DGS-SLAM [21] to use 100 and 200 maximum iterations on TUM RGB-D [32] and BONN [29] datasets, respectively. For mapping, we train for 50 iterations and use window size of 8. We only apply flow loss in the last 25 iterations to speed up each mapping step. We follow the keyframe selection strategy of 4DGS-SLAM, which selects keyframes based on motion mask differences and at least selects 1 keyframe in every 5 frames. After online training, we also conduct 1500 iterations of color refinement following 4DGS-SLAM [21] and MonoGS [25]. For other detailed hyperparameters, please refer to our supplementary material.

Baseline Methods. We mainly compare with 3DGS-SLAM methods without explicit loop closure. For static SLAM, we compare with MonoGS [25], SplaTAM [18]. For dynamic reconstruction and SLAM, we compare with SC-GS [15] and 4DGS-SLAM [21]. For 4DGS-SLAM, we report our reproduced results using their officially released code.

Metrics. For camera tracking, we follow the standard SLAM methods to evaluate the Root Mean Square Error (RMSE) of the Absolute Trajectory error (ATE) across all frames. For photometric rendering quality, we report the PSNR, SSIM, and LPIPS on all frames.

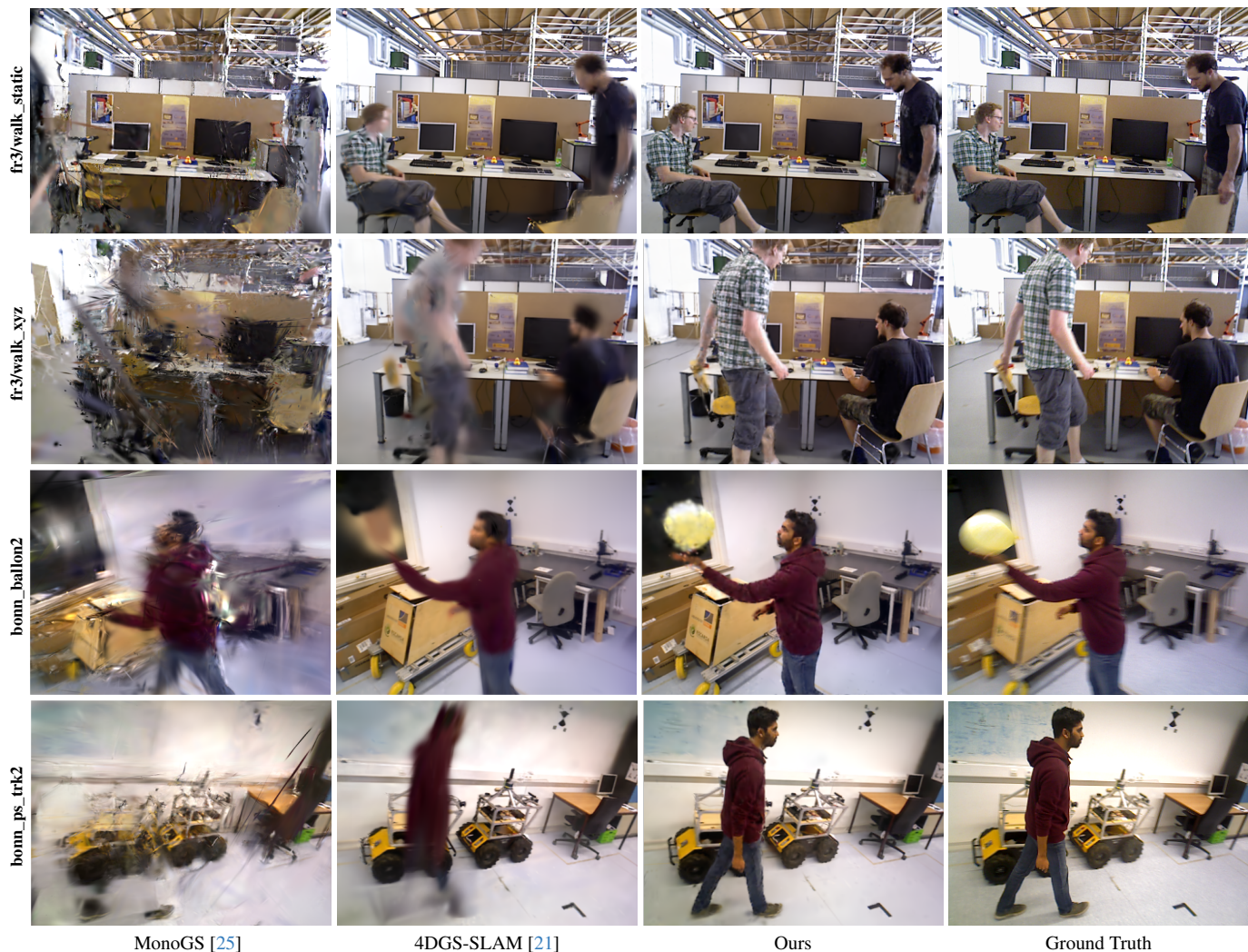


Figure 3. Qualitative Renderings on *non-keyframes* in TUM RGB-D [32] (first two rows) and BONN [29] (last two rows) datasets.

4.2. Main Results

Camera Tracking. The camera tracking results on TUM RGB-D [32] and BONN [29] are reported in Tables 1 and 3. Our method achieves more accurate camera trajectories on both datasets. Notably, 4DGS-SLAM [21] uses prolonged mapping (200 iterations) to reconstruct the scene and refine camera poses within a pose window of length 3, whereas we use far fewer mapping iterations while still obtaining better tracking accuracy. A qualitative comparison is shown in Figure 4. Without the motion decomposition module (w/o Motion Decomp.), our method performs worse than 4DGS-SLAM due to limited camera pose refinement. Adding category-agnostic motion masking (w/o Camera Init) markedly improves the results on *ballon2*, which contains category-agnostic dynamic objects. Furthermore, flow-guided camera initialization improves tracking accuracy on both scenes and mitigates camera trajectory drift.

Photometric Renderings. As shown in Tables 2 and 4, our method achieves substantially higher photometric quality than the baselines. Compared with 4DGS-SLAM [21], the main advantage comes from our hybrid dynamic Gaussian representation, supported by scene-flow based propagation and adaptive insertion for more effective dynamic modeling. Qualitative results are shown in Figure 3. 4DGS-SLAM underperforms in highly dynamic scenes, *e.g.* when people leave and re-enter the view (*fr3/walking_xyz*), because it requires a handcrafted dynamic start time to assign dynamic Gaussians and cannot handle such complex dynamics well. In contrast, our method can adaptively insert new dynamic Gaussians without scene-specific hyperparameters. The BONN [29] results also show that our method preserves more details in fast-moving scenes and reconstructs category-agnostic dynamic objects. We further visualize 2D tracking together with dynamic reconstruction in Figure 1 by sampling 2D points at the start time and incrementally warping

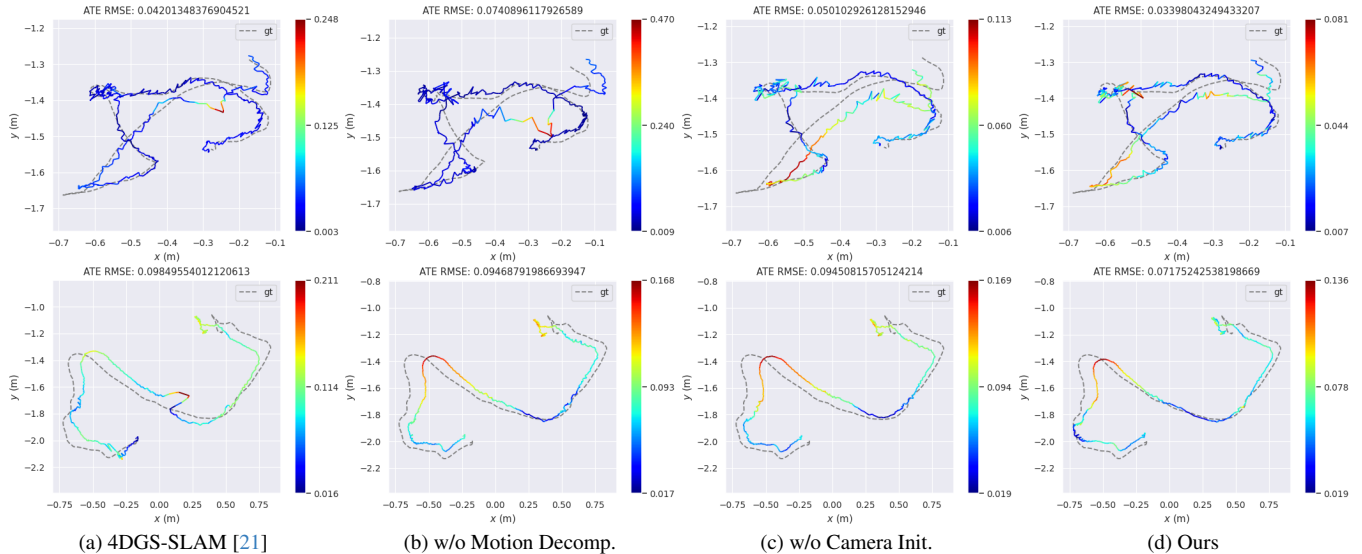


Figure 4. Tracking Results on BONN [29] ballon2 (first row) and person_tracking (second row) scenes.

Table 5. Running time (ms) analysis.

Method	Dynamic Seg.	Tracking	Mapping	FPS
MonoGS [25]	-	476	557	1.93
4DGS-SLAM [21]	16	445	110562	0.04
Ours	68	427	6285	0.50

them with the rendered Gaussian flow maps. This shows that the learned Gaussian flow trajectories are spatially and temporally coherent.

Runtime Analysis. The runtime analysis is reported in Table 5. 4DGS-SLAM [21] has an extremely slow mapping speed due to its computationally expensive design. In each mapping step, it runs 100 iterations to train the deformation MLPs and another 100 iterations to jointly optimize the Gaussian parameters, which substantially reduces FPS. In contrast, our method benefits from explicit dynamic Gaussian centers guided by optical flow, leading to much faster mapping. For dynamic segmentation, our camera-induced motion segmentation takes around 52 ms per frame and plays an indispensable role in our method (*cf.* Table 6, Figure 4).

4.3. Ablation Study

The quantitative ablation results are shown in Table 6. Our camera-induced motion decomposition module improves both camera tracking accuracy and rendering quality, especially on BONN [29], which contains category-agnostic dynamic objects and fast-moving scenes and cameras. The flow propagation and adaptive insertion modules contribute substantially to reconstruction quality in fast-moving scenes (*e.g.* ballon2) and scenes with re-occurring dynamic objects (*e.g.* fr3/walk_xyz). Our GMM-based temporal opacity and rotation further improve representation capabil-

Table 6. Ablation study.

Method	fr3/walk_xyz		ballon2	
	ATE [cm]↓	PSNR↑	ATE [cm]↓	PSNR↑
w/o Motion Decomp.	2.7	24.40	7.4	27.59
w/o Flow Propagate	2.6	23.91	4.2	26.86
w/o Adaptive Insert	3.4	23.53	3.9	27.93
w/o GMM	2.7	24.04	3.7	27.91
w/o KNN smooth	2.5	24.47	3.5	28.14
Ours	2.5	24.60	3.4	28.36

ity in complex dynamic scenes, while the KNN smoothing strategy enhances the local rigidity of propagated Gaussians for fast mapping.

5. Conclusion

In this paper, we proposed Flow4DGS-SLAM, a novel dynamic SLAM system designed to jointly reconstruct static and dynamic regions. We first propose a camera-induced motion decomposition module, which leverages depth and prior optical flow to not only generate category-agnostic motion mask but also provide better camera pose initialization. We accelerate the training of dynamic Gaussians with a hybrid dynamic Gaussian representation combined with optical flow-guided 4D mapping. The empirical results demonstrate the significance of our Flow4DGS-SLAM in dynamic 3DGS SLAM, showing superior performance on camera tracking, photometric rendering quality, and training efficiency.

Acknowledgment. This research / project is supported by the National Research Foundation (NRF) Singapore, under its NRF-Investigatorship Programme (Award ID. NRF-NRFI09-0008), and the Tier 2 grant MOET2EP20124-0015 from the Singapore Ministry of Education.

References

- [1] Berta Bescos, José M Fácil, Javier Civera, and José Neira. Dynaslam: Tracking, mapping, and inpainting in dynamic scenes. 2018. [2](#)
- [2] Carlos Campos, Richard Elvira, Juan J Gómez Rodríguez, José MM Montiel, and Juan D Tardós. Orb-slam3: An accurate open-source library for visual, visual-inertial, and multi-view slam. 2021. [1](#), [2](#)
- [3] David Charatan, Sizhe Lester Li, Andrea Tagliasacchi, and Vincent Sitzmann. pixelsplat: 3d gaussian splats from image pairs for scalable generalizable 3d reconstruction. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 19457–19467, 2024. [2](#)
- [4] François Chaumette, Seth Hutchinson, and Peter Corke. Visual servoing. In *Springer handbook of robotics*, pages 841–866. Springer, 2016. [3](#)
- [5] Hanlin Chen, Chen Li, and Gim Hee Lee. Neusg: Neural implicit surface reconstruction with 3d gaussian splatting guidance. *arXiv preprint arXiv:2312.00846*, 2023. [2](#)
- [6] Hanlin Chen, Fangyin Wei, Chen Li, Tianxin Huang, Yunsong Wang, and Gim Hee Lee. Vcr-gaus: View consistent depth-normal regularizer for gaussian surface reconstruction. *Advances in Neural Information Processing Systems*, 37:139725–139750, 2024. [2](#)
- [7] Yuedong Chen, Haofei Xu, Chuanxia Zheng, Bohan Zhuang, Marc Pollefeys, Andreas Geiger, Tat-Jen Cham, and Jianfei Cai. Mvsplat: Efficient 3d gaussian splatting from sparse multi-view images. In *European Conference on Computer Vision*, pages 370–386. Springer, 2024. [2](#)
- [8] Brian Curless and Marc Levoy. A volumetric method for building complex models from range images. In *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, pages 303–312, 1996. [2](#)
- [9] Jan Czarnowski, Tristan Laidlow, Ronald Clark, and Andrew J Davison. Deepfactors: Real-time probabilistic dense monocular slam. *IEEE Robotics and Automation Letters*, 5(2):721–728, 2020. [2](#)
- [10] 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](#)
- [11] Yuanxing Duan, Fangyin Wei, Qiyu Dai, Yuhang He, Wenzheng Chen, and Baoquan Chen. 4d-rotor gaussian splatting: towards efficient novel view synthesis for dynamic scenes. In *ACM SIGGRAPH 2024 Conference Papers*, pages 1–11, 2024. [1](#)
- [12] Antoine Guédon and Vincent Lepetit. Sugar: Surface-aligned gaussian splatting for efficient 3d mesh reconstruction and high-quality mesh rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5354–5363, 2024. [2](#)
- [13] Zhiyang Guo, Wengang Zhou, Li Li, Min Wang, and Houqiang Li. Motion-aware 3d gaussian splatting for efficient dynamic scene reconstruction. *IEEE Transactions on Circuits and Systems for Video Technology*, 2024. [1](#), [2](#)
- [14] Berthold K.P. Horn and Brian G. Schunck. Determining optical flow. *Artificial Intelligence*, 17(1-3):185–203, 1981. [3](#)
- [15] Yi-Hua Huang, Yang-Tian Sun, Ziyi Yang, Xiaoyang Lyu, Yan-Pei Cao, and Xiaojuan Qi. Sc-gs: Sparse-controlled gaussian splatting for editable dynamic scenes. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4220–4230, 2024. [1](#), [2](#), [5](#), [6](#)
- [16] Haochen Jiang, Yueming Xu, Kejie Li, Jianfeng Feng, and Li Zhang. Rodyn-slam: Robust dynamic dense rgb-d slam with neural radiance fields. *IEEE Robotics and Automation Letters*, 2024. [5](#), [6](#)
- [17] Mohammad Mahdi Johari, Camilla Carta, and François Fleuret. Eslam: Efficient dense slam system based on hybrid representation of signed distance fields. In *CVPR*, 2023. [1](#)
- [18] Nikhil Keetha, Jay Karhade, Krishna Murthy Jatavallabhula, Gengshan Yang, Sebastian Scherer, Deva Ramanan, and Jonathon Luiten. Splatam: Splat, track & map 3d gaussians for dense rgb-d slam. In *CVPR*, 2024. [1](#), [2](#), [5](#), [6](#)
- [19] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Trans. Graph.*, 42(4):139–1, 2023. [1](#), [2](#), [4](#)
- [20] Mangyu Kong, Jaewon Lee, Seongwon Lee, and Euntai Kim. Dgs-slam: Gaussian splatting slam in dynamic environment. *arXiv preprint arXiv:2411.10722*, 2024. [1](#)
- [21] Yanyan Li, Youxu Fang, Zunjie Zhu, Kunyi Li, Yong Ding, and Federico Tombari. 4d gaussian splatting slam, 2025. [1](#), [2](#), [4](#), [5](#), [6](#), [7](#), [8](#)
- [22] Zhan Li, Zhang Chen, Zhong Li, and Yi Xu. Spacetime gaussian feature splatting for real-time dynamic view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8508–8520, 2024. [1](#), [2](#)
- [23] Yiqing Liang, Numair Khan, Zhengqin Li, Thu Nguyen-Phuoc, Douglas Lanman, James Tompkin, and Lei Xiao. Gafre: Gaussian deformation fields for real-time dynamic novel view synthesis. In *2025 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 2642–2652. IEEE, 2025. [1](#), [2](#)
- [24] Jonathon Luiten, Georgios Kopanas, Bastian Leibe, and Deva Ramanan. Dynamic 3d gaussians: Tracking by persistent dynamic view synthesis. In *2024 International Conference on 3D Vision (3DV)*, pages 800–809. IEEE, 2024. [2](#)
- [25] Hidenobu Matsuki, Riku Murai, Paul HJ Kelly, and Andrew J Davison. Gaussian splatting slam. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18039–18048, 2024. [1](#), [2](#), [5](#), [6](#), [7](#), [8](#)
- [26] Raul Mur-Artal and Juan D Tardós. Orb-slam2: An open-source slam system for monocular, stereo, and rgb-d cameras. 2017. [1](#), [2](#)
- [27] Richard A Newcombe, Shahram Izadi, Otmar Hilliges, David Molyneaux, David Kim, Andrew J Davison, Pushmeet Kohi, Jamie Shotton, Steve Hodges, and Andrew Fitzgibbon. Kinectfusion: Real-time dense surface mapping and tracking. In *2011 10th IEEE international symposium on mixed and augmented reality*, pages 127–136. Ieee, 2011. [2](#)
- [28] Richard A Newcombe, Steven J Lovegrove, and Andrew J Davison. Dtam: Dense tracking and mapping in real-time.

- In *2011 international conference on computer vision*, pages 2320–2327. IEEE, 2011. 2
- [29] Emanuele Palazzolo, Jens Behley, Philipp Lottes, Philippe Giguere, and Cyrill Stachniss. Refusion: 3d reconstruction in dynamic environments for rgb-d cameras exploiting residuals. 2019. 6, 7, 8
- [30] Thomas Schops, Torsten Sattler, and Marc Pollefeys. Bad slam: Bundle adjusted direct rgb-d slam. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 134–144, 2019. 2
- [31] Colton Stearns, Adam Harley, Mikaela Uy, Florian Dubost, Federico Tombari, Gordon Wetzstein, and Leonidas Guibas. Dynamic gaussian marbles for novel view synthesis of casual monocular videos. In *SIGGRAPH Asia 2024 Conference Papers*, pages 1–11, 2024. 1
- [32] Jürgen Sturm, Nikolas Engelhard, Felix Endres, Wolfram Burgard, and Daniel Cremers. A benchmark for the evaluation of rgb-d slam systems. 2012. 6, 7
- [33] Edgar Sucar, Shikun Liu, Joseph Ortiz, and Andrew Davison. iMAP: Implicit mapping and positioning in real-time. In *ICCV*, 2021. 1, 2
- [34] Zachary Teed and Jia Deng. Deepv2d: Video to depth with differentiable structure from motion. *arXiv preprint arXiv:1812.04605*, 2018. 2
- [35] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *ECCV*, 2020. 3
- [36] Zachary Teed and Jia Deng. DROID-SLAM: Deep Visual SLAM for Monocular, Stereo, and RGB-D Cameras. 2021. 2
- [37] Benjamin Ummenhofer, Huizhong Zhou, Jonas Uhrig, Nikolaus Mayer, Eddy Ilg, Alexey Dosovitskiy, and Thomas Brox. Demon: Depth and motion network for learning monocular stereo. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5038–5047, 2017. 2
- [38] Chien-Yao Wang and Hong-Yuan Mark Liao. YOLOv9: Learning what you want to learn using programmable gradient information. 2024. 3
- [39] Hengyi Wang, Jingwen Wang, and Lourdes Agapito. Co-slam: Joint coordinate and sparse parametric encodings for neural real-time slam. In *CVPR*, 2023. 1, 2
- [40] Qianqian Wang, Vickie Ye, Hang Gao, Weijia Zeng, Jake Austin, Zhengqi Li, and Angjoo Kanazawa. Shape of motion: 4d reconstruction from a single video. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9660–9672, 2025. 1, 2
- [41] Shizun Wang, Xingyi Yang, QiuHong Shen, Zhenxiang Jiang, and Xinchao Wang. Gflow: Recovering 4d world from monocular video. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 7862–7870, 2025. 2
- [42] Yunsong Wang, Tianxin Huang, Hanlin Chen, and Gim Hee Lee. Freesplat: Generalizable 3d gaussian splatting towards free view synthesis of indoor scenes. *Advances in Neural Information Processing Systems*, 37:107326–107349, 2024. 2
- [43] Yunsong Wang, Tianxin Huang, Hanlin Chen, and Gim Hee Lee. Freesplat++: Generalizable 3d gaussian splatting for efficient indoor scene reconstruction. *arXiv preprint arXiv:2503.22986*, 2025. 2
- [44] Thomas Whelan, Stefan Leutenegger, Renato F Salas-Moreno, Ben Glocker, and Andrew J Davison. Elasticfusion: Dense slam without a pose graph. In *Robotics: science and systems*, page 3. Rome, Italy, 2015. 2
- [45] Wenhua Wu, Chenpeng Su, Siting Zhu, Tianchen Deng, Zhe Liu, and Hesheng Wang. Add-slam: Adaptive dynamic dense slam with gaussian splatting. *arXiv preprint arXiv:2505.19420*, 2025. 2
- [46] Yueming Xu, Haochen Jiang, Zhongyang Xiao, Jianfeng Feng, and Li Zhang. Dg-slam: Robust dynamic gaussian splatting slam with hybrid pose optimization. *Advances in Neural Information Processing Systems*, 37:51577–51596, 2025. 1, 2
- [47] 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
- [48] 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 *ISMAR*, 2022. 1
- [49] Zeyu Yang, Hongye Yang, Zijie Pan, Xiatian Zhu, and Li Zhang. Real-time photorealistic dynamic scene representation and rendering with 4d gaussian splatting. *arXiv preprint arXiv:2310.10642*, 2023. 1, 2
- [50] Ziyi Yang, Xinyu Gao, Wen Zhou, Shaohui Jiao, Yuqing Zhang, and Xiaogang Jin. Deformable 3d gaussians for high-fidelity monocular dynamic scene reconstruction. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 20331–20341, 2024. 1, 2
- [51] Chao Yu, Zuxin Liu, Xin-Jun Liu, Fugui Xie, Yi Yang, Qi Wei, and Qiao Fei. Ds-slam: A semantic visual slam towards dynamic environments. 2018. 2
- [52] Yuanwen Yue, Anurag Das, Francis Engelmann, Siyu Tang, and Jan Eric Lenssen. Improving 2d feature representations by 3d-aware fine-tuning. In *European Conference on Computer Vision*, pages 57–74. Springer, 2024. 2
- [53] Vladimir Yugay, Yue Li, Theo Gevers, and Martin R. Oswald. Gaussian-slam: Photo-realistic dense slam with gaussian splatting, 2023. 1
- [54] Jun Zhang, Mina Henein, Robert Mahony, and Viorela Ila. Vdo-slam: a visual dynamic object-aware slam system. *arXiv preprint*, 2020. 2
- [55] Jianhao Zheng, Zihan Zhu, Valentin Bieri, Marc Pollefeys, Songyou Peng, and Iro Armeni. Wildgs-slam: Monocular gaussian splatting slam in dynamic environments. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 11461–11471, 2025. 1, 2
- [56] Huizhong Zhou, Benjamin Ummenhofer, and Thomas Brox. Deeptam: Deep tracking and mapping. In *Proceedings of the European conference on computer vision (ECCV)*, pages 822–838, 2018. 2
- [57] Zihan Zhu, Songyou Peng, Viktor Larsson, Weiwei Xu, Hujun Bao, Zhaopeng Cui, Martin R. Oswald, and Marc Pollefeys. Nice-slam: Neural implicit scalable encoding for slam. In *CVPR*, 2022. 1, 2