# How NeRFs and 3D Gaussian Splatting are Reshaping SLAM: a Survey

Fabio Tosi[1]    Youmin Zhang[1,2]    Ziren Gong[1]    Erik Sandström[3]
Stefano Mattoccia[1]    Martin R. Oswald[3,4]    Matteo Poggi[1]

[1]University of Bologna, Italy    [2]Rock Universe, China    [3]ETH Zurich, Switzerland    [4]University of Amsterdam, Netherlands

**Abstract**—Over the past two decades, research in the field of Simultaneous Localization and Mapping (SLAM) has undergone a significant evolution, highlighting its critical role in enabling autonomous exploration of unknown environments. This evolution ranges from hand-crafted methods, through the era of deep learning, to more recent developments focused on Neural Radiance Fields (NeRFs) and 3D Gaussian Splatting (3DGS) representations. Recognizing the growing body of research and the absence of a comprehensive survey on the topic, this paper aims to provide the first comprehensive overview of SLAM progress through the lens of the latest advancements in radiance fields. It sheds light on the background, evolutionary path, inherent strengths and limitations, and serves as a fundamental reference to highlight the dynamic progress and specific challenges.

**Index Terms**—Simultaneous Localization and Mapping, SLAM, Deep Learning, Neural Radiance Field, NeRF, 3D Gaussian Splatting
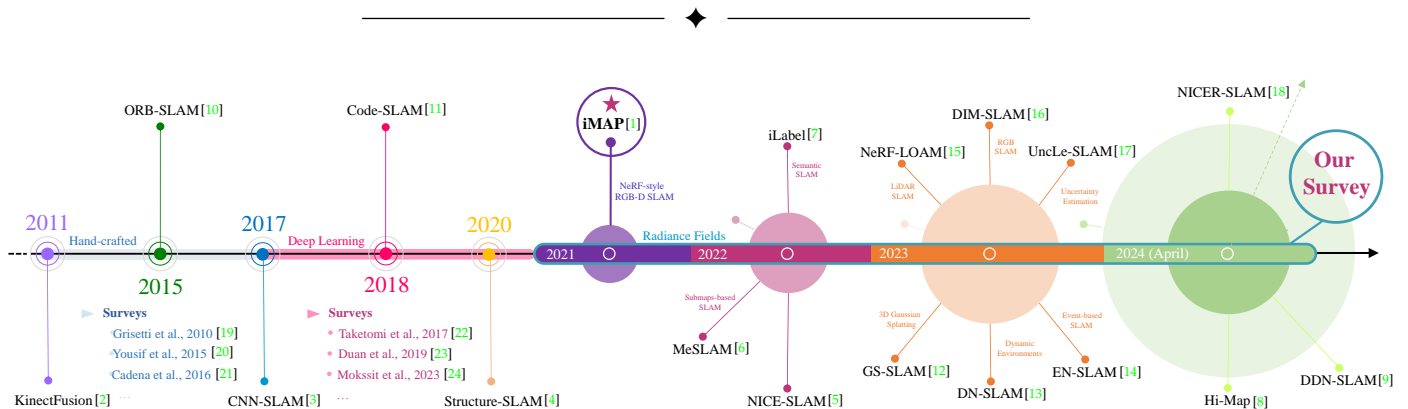
Fig. 1: **Timeline SLAM Evolution.** This timeline begins by illustrating the transition from hand-crafted to deep learning techniques, featuring key surveys from both eras. In 2021, a pivotal shift focuses on radiance-field-based SLAM systems, marked by iMap [1]. The circles on the right side of the figure represent key papers for each year, with size indicating publication volume. The outer circle for 2024 signals a projected surge, highlighting the growing interest in NeRF and 3DGS-inspired SLAM.

## 1 INTRODUCTION

Simultaneous Localization and Mapping (SLAM) is a fundamental concept in the fields of computer vision and robotics. It addresses the challenge of enabling machines to autonomously navigate and incrementally build a map of unknown environments (*mapping*) while simultaneously determining their own position and orientation (*tracking*).

Originally conceived for robotics and automated systems, the demand for SLAM has expanded into a variety of domains, including augmented reality (AR), visual surveillance, medical applications, and beyond. To meet these needs, researchers have focused on developing methods for machines to autonomously construct increasingly highly accurate scene representations, influenced by the convergence of robotics, computer vision, sensor technology, and the recent progress in artificial intelligence (AI).

Typically, SLAM techniques rely on the integration of diverse sensing technologies, including cameras, laser range instruments, inertial devices, and GPS, to effectively accomplish the task at hand. Initially, sonar and LiDAR sensors were prevalent choices due to their high precision, despite being cumbersome and costly. Subsequently, the focus shifted towards visual sensors such as monocular/stereo or RGB-D cameras, which offer advantages in terms of portability, cost-effectiveness, and deployment ease. These visual sensors enable Visual Simultaneous Localization and Mapping (VSLAM) systems to capture more detailed environmental information, improve precise positioning in complex scenarios, and deliver versatile and accessible solutions.

As we outline the ideal SLAM criteria, several key aspects emerge. These include global consistency, robust camera tracking, accurate surface modeling, real-time performance, accurate prediction in unobserved regions, scalability to large scenes, and robustness to noisy data.

Over the years, SLAM methodologies have evolved significantly to meet these specific requirements. At the outset, hand-crafted algorithms [2], [10], [25], [26], [27] demonstrated remarkable real-time performance and scalability. However, they face challenges in strong illumination, radiometric changes, and dynamic/poorly textured environ-

ments, resulting in unsatisfactory performance. The incorporation of advanced techniques, employing deep learning methodologies [3], [4], [11], [28], became crucial in improving the precision and reliability of localization and mapping. This integration takes advantage of the robust feature extraction capabilities of deep neural networks, which are particularly effective in challenging conditions. Nonetheless, their dependence on extensive training data and accurate ground truth annotations limits their ability to generalize to unseen scenarios. Furthermore, both hand-crafted and deep learning-based methods encounter limitations related to using discrete surface representations (point/surfel clouds [29], [30], voxel hashing [31], voxel-grids [2], octrees [32]), which lead to challenges such as sparse 3D modeling, limited spatial resolution and distortion during the reconstruction process. Additionally, accurately estimating geometries in unobserved areas remains an ongoing hurdle.

Driven by the need to overcome existing obstacles and influenced by the success of recent Neural Radiance Fields (NeRF) [33] and 3D Gaussian Splatting (3DGS) [34] representations in high-fidelity view synthesis, along with the introduction of learned representations for modeling geometric fields [35], [36], [37] – extensively discussed in [38] – a revolution is reshaping SLAM systems. Leveraging insights from contemporary research, these approaches offer several advantages over previous methods, including continuous surface modeling, reduced memory requirements, improved noise/outlier handling, and enhanced hole filling and scene inpainting capabilities for occluded or sparse observations. In addition, they have the potential to produce denser and more compact maps that can be reconstructed as 3D meshes at arbitrary resolutions. However, it is important to note that at this early stage, the strengths of each technique coexist with specific challenges and limitations. As such, the field is constantly evolving, and continuous investigation and innovation are required to make further progress.

In response to the lack of SLAM surveys focusing on the latest developments and the growing interest in research exploring this paradigm [39][1], this paper conducts a thorough review of contemporary radiance field-inspired SLAM techniques. Specifically, we undertake an in-depth investigation of 73 SLAM systems that have emerged in the past three years, reflecting the rapid pace of progress in the field. This evolution is illustrated in Figure 1, which provides a visual timeline of the current state of SLAM advancements. Our aim is to fill the existing gap in the survey literature by closely examining and analyzing these cutting-edge techniques, and by highlighting the rapid emergence of innovative solutions aimed at improving their inherent weaknesses. Through a detailed exploration, we intend to categorize these methods, trace their progression, and offer insights that are tailored to the specific requirements of SLAM. By serving as a valuable resource for both the novice and the expert, we believe that this survey represents a significant cornerstone for the future of this paradigm.

The upcoming sections are organized as follows:

- Section 2 provides an overview of existing SLAM surveys (2.1), delves into recent radiance-field volume rendering theory (2.2), introduces prevalent datasets and

---

benchmarks in the field (2.3), and presents the main quality assessment metrics used in this context (2.4).
- Section 3 is the core of our paper, focusing on key NeRF and 3DGS-inspired SLAM techniques and our structured taxonomy for organizing these advancements.
- Section 4 presents quantitative results evaluating SLAM frameworks in tracking, mapping, rendering, and performance analysis across diverse scenarios.
- Sections 5 and 6 focus on limitations, future research directions, and summarize the survey comprehensively.

## 2 BACKGROUND

### 2.1 Existing SLAM Surveys

SLAM has seen significant growth, resulting in a variety range of comprehensive survey papers. In the early stages, Durrant-Whyte and Bailey introduced the probabilistic nature of the SLAM problem and highlighted key methods, alongside implementations [40], [41]. Grisetti et al. [19] further delved into the graph-based SLAM problem, emphasizing its role in navigating in unknown environments. In the field of visual SLAM, Yousif [20] provided an overview of localization and mapping techniques, incorporating basic methods and advances in visual odometry and SLAM. The advent of multiple-robot systems led to Saeedi and Clark [42] reviewing state-of-the-art approaches, with a focus on multiple-robot SLAM challenges and solutions. Cadena et al. [21] presented a comprehensive reflection on the history, robustness, and new frontiers of SLAM, addressing its evolving significance across real-world applications. Taketomi et al. [22] categorized and summarized VSLAM algorithms from 2010 to 2016, classifying them based on feature-based, direct, and RGB-D camera approaches. Saputra et al. [43] addressed the challenge of dynamic environments in VSLAM and Structure from Motion (SfM), presenting a taxonomy of techniques for reconstruction, segmentation and tracking of dynamic objects. The integration of deep learning with SLAM was meticulously examined by Duan et al. [23], highlighting the progression of deep learning methods in visual SLAM. In sensor-specific contexts, Zaffar et al. [44] discussed sensors employed in SLAM, while Yang et al. [45] and Zhao et al. [46] explored the applications of LiDAR and underwater SLAM, respectively. In recent years, deep learning-based VSLAM has gained considerable attention, extensively covered in [47], [48], [49], [50]. Notably, [51] delves into recent advancements in RGB-D scene reconstruction. Ongoing developments in SLAM are explored in surveys like [52], focusing on active SLAM strategies for precise mapping through motion planning.

Despite the extensive body of work describing SLAM systems covering traditional and deep learning-based approaches, there is no comprehensive exploration of the advancing frontiers in SLAM techniques rooted in the latest progress in radiance fields. Nonetheless, within the existing literature of our interest, notably in influential works like [53], two principal SLAM strategies emerge as the *frame-to-frame* and *frame-to-model* tracking approaches, which are influencing the development of new methodologies based on radiance fields. Typically, the former strategy is used in real-time systems, often involving further optimization of the estimated poses through *loop-closure* (LC) or global Bundle
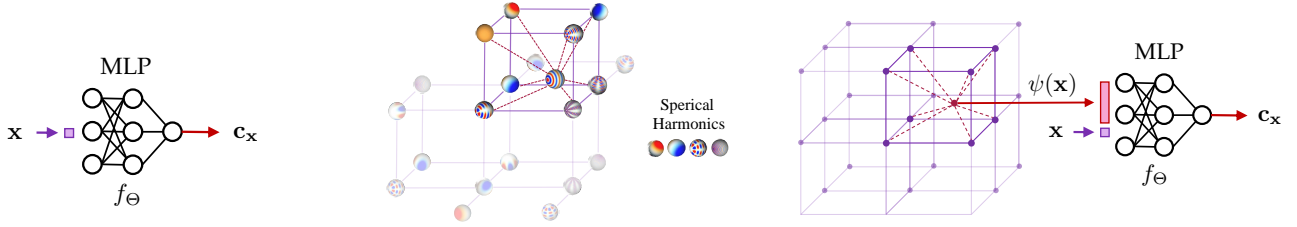
Fig. 2: **Comparison of Scene Representations:** *Implicit*, *Explicit*, **and** *Hybrid*. From left to right: *Implicit* uses a neural network to approximate a radiance field, *explicit* conducts volume rendering directly on learned spatial feature (voxels, hash grids, etc.), excluding neural components, and *hybrid* incorporates learned spatial features $\psi$ with neural networks. Both *hybrid* and *explicit* approaches enable accelerated training and rendering but require additional memory resources.

Adjustment (BA), whereas the latter estimates camera poses from the reconstructed 3D model, often avoiding further optimizations, yet resulting less scalable to large scenes. These strategies, often associated with the concepts of decoupled and coupled methods in recent SLAM research, serve as the foundation for the methodologies we will explore. Decoupled methods employ separate frameworks for tracking and mapping, treating them as independent tasks, while coupled methods utilize a unified representation for both tasks, allowing for a more integrated approach.

## 2.2 Progress in Radiance Field Theory

The term radiance field refers to a representation that describes the behavior and distribution of light within a three-dimensional space. It encapsulates how light interacts with surfaces, materials, and the surrounding environment. It can be represented implicitly, by encoding it entirely within the weights of a neural network or explicitly, by mapping light within a discrete spatial structure such as voxel grids. Explicit representations typically offer faster access but require more memory and have resolution constraints, while implicit representations provide a compact scene encoding with potentially higher rendering computational needs. Hybrid approaches take advantage of both by using a combination of explicitly stored local latent features and shallow neural networks, using various structures such as sparse voxel hashing grids [54], [55], multi-resolution dense voxel grids [56], unordered point sets [57], and more. Figure 2 visually illustrates these representations, which have recently had a significant impact on SLAM methodologies, primarily through the incorporation of models derived from NeRF and more recent explicit methods such as 3DGS. Below, we briefly describe NeRF – for image rendering and surface reconstruction – and 3DGS, essential for understanding the upcoming SLAM approaches.

### 2.2.1 Neural Radiance Field (NeRF)

In 2020, Mildenhall et al. [33] introduced NeRF, an implicit, continuous volumetric representation, setting a new standard for novel view synthesis. In contrast to conventional explicit volumetric models, this method employs a sparse set of input views to optimize a continuous volumetric scene function, representing three-dimensional scenes via a radiance field. To achieve this, the original NeRF implementation requires knowledge of the camera poses and intrinsic parameters corresponding to each input view, which are estimated using the COLMAP structure-from-motion package

[58], [59]. This approach has become the common practice in subsequent research building upon the NeRF framework. Formally expressed as $f(\mathbf{x}, \mathbf{d}) \rightarrow (\mathbf{c}, \sigma)$, the model leverages an MLP (Multi-Layer Perceptron) with weights $\Theta$, denoted as $f_\Theta$, approximating a 5D function of viewing direction $\mathbf{d} = (\theta, \phi)$ and in-scene 3D coordinates $\mathbf{x} = (x, y, z)$. Notably, the representation ensures multi-view consistency by predicting the volume density $\sigma$ independently of viewing direction, while color $\mathbf{c} = (r, g, b)$ depends on both viewing direction and 3D coordinates.

The NeRF workflow for novel view synthesis involves casting camera rays through the scene to generate sampling points per pixel, computing local color and density using the NeRF MLP(s) for each sampling point, and employing volume rendering to synthesize the 2D image. Specifically, the computation of the color $C(\mathbf{r})$ resulting from a camera ray $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$ involves an integral formulation:

$$C(\mathbf{r}) = \int_{t_1}^{t_2} T(t)\sigma(\mathbf{r}(t))\mathbf{c}(\mathbf{r}(t), \mathbf{d})\mathrm{d}t \qquad (1)$$

Here, $\mathrm{d}t$ denotes the differential distance traveled by the ray at each integration step. The terms $\sigma(\mathbf{r}(t))$ and $\mathbf{c}(\mathbf{r}(t), \mathbf{d})$ represent the volume density and color at point $\mathbf{r}(t)$ along the camera ray with viewing direction $\mathbf{d}$, respectively. Additionally, $T(t) = \exp\left(-\int_{t_1}^{t} \sigma(\mathbf{r}(s)), \mathrm{d}s\right)$ is the accumulated transmittance from $t_1$ to $t$.

The integral computation uses quadrature by dividing the ray into $N$ evenly-spaced bins:

$$C(\mathbf{r}) = \sum_{i=1}^{N} \alpha_i T_i \mathbf{c}_i, \quad T_i = \exp\left(-\sum_{j=1}^{i-1} \sigma_j \delta_j\right) \qquad (2)$$

where, $\delta_i$ denotes the interval between consecutive samples $t_i$ and $t_{i+1}$, while $\sigma_i$ and $\mathbf{c}_i$ indicate the density and color evaluated at sample point $i$ along the ray, respectively. Additionally, $\alpha_i = (1 - \exp(-\sigma_i \delta_i))$ characterizes the opacity resulting from alpha compositing at sample point $i$.

The expected depth along a ray, instead, can be calculated using the accumulated transmittance:

$$d(\mathbf{r}) = \int_{t_1}^{t_2} T(t) \cdot \sigma(\mathbf{r}(t)) \cdot t \, dt, \qquad (3)$$

Similarly to Eq. 2, this can be approximated as:

$$\hat{D}(\mathbf{r}) = \sum_{i=1}^{N} \alpha_i t_i T_i. \qquad (4)$$

In this context, some methods propose using expected depth to either impose depth supervision from external priors [60], [61] or apply regularization techniques [62], enhancing scene geometry and enforcing depth smoothness.

For optimization, a square error photometric loss is employed, represented as $L = \sum_{r \in R} \|\hat{C}(\mathbf{r}) - C_{gt}(\mathbf{r})\|_2^2$. Here, $C_{gt}(r)$ denotes the ground truth color for the pixel associated with $\mathbf{r}$ in the training image, and $R$ denotes the batch of rays for synthesizing the target image.

While NeRF achieved success, challenges like slow training/rendering speeds persist. Follow-up methods, comprehensively surveyed in [63], [64], [65], seek to enhance quality or faster training/rendering using techniques such as hashing [54] or sparse 3D grids [66]. However, these methods still struggle to accurately represent empty spaces and face image quality limitations due to structured grids, which significantly impede rendering speeds. Other works [67], [68], [69], instead, aims to reduce the reliance on external tools like COLMAP for camera pose estimation. While these approaches share the goal of joint pose estimation and scene reconstruction with SLAM, they differ in their processing paradigm. SLAM typically processes images sequentially as they are captured, enabling real-time operation. In contrast, these NeRF-based pose estimation methods often require a set of images to be processed simultaneously, limiting their applicability in real-time scenarios. Moreover, they either need a pre-trained neural implicit network or cannot optimize poses and the network concurrently, further constraining their use in SLAM applications.

### 2.2.2  Surface Reconstruction from Neural Fields

Despite the potential of NeRF and its variants to capture the 3D geometry of a scene, these models are implicitly defined in the weights of the neural network. Obtaining an explicit representation of the scene through 3D meshes is desirable for 3D reconstruction applications. Starting with NeRF, a basic approach to achieving coarse scene geometry is to threshold the density predicted by the MLP. More advanced solutions explore three main representations.

**Occupancy.** This representation models free versus occupied space by replacing alpha values $\alpha_i$ along the ray with a learned discrete function $o(x) \in \{0, 1\}$. Specifically, an occupancy probability $\in [0, 1]$ is estimated and surfaces are obtained by running the marching cubes algorithm [70].

**Signed Distance Function (SDF).** An alternative method for scene geometry is the signed distance from any point to the nearest surface, yielding negative values inside objects and positive values outside. NeuS [71] was the first to revisit the NeRF volumetric rendering engine, predicting the SDF with an MLP as $f(\mathbf{r}(t))$ and replacing $\alpha$ with $\rho(t)$, derived from the SDF as follows:

$$\rho(t) = \max\left(\frac{-\frac{d\Phi}{dt}(f(\mathbf{r}(t)))}{\Phi(f(\mathbf{r}(t)))}, 0\right) \qquad (5)$$

with $\Phi$ being the sigmoid function and $\frac{d\Phi}{dt}$ its derivative.

**Truncated Signed Distance Function (TSDF).** Finally, predicting a truncated SDF with the MLP allows for removing the contribution by any SDF value too far from individual surfaces during rendering. In [72], pixel color is obtained as a weighted sum of colors sampled along the ray:
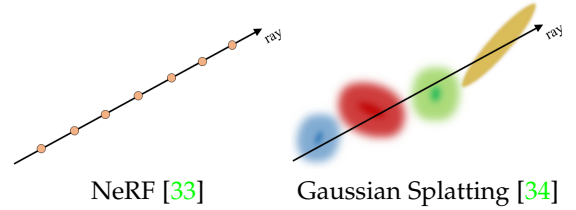


Fig. 3: **NeRF and 3DGS differ conceptually.** (left) NeRF queries an MLP along the ray, while (right) 3DGS blends Gaussians for the given ray.

$$C(\mathbf{r}) = \frac{\sum_{i=1}^{N} w_i \mathbf{c}_i}{\sum_{i=1}^{N} w_i} \qquad (6)$$

with $w_i$ defined, according to truncation distance $t_r$, as

$$w_i = \Phi\left(\frac{f(\mathbf{r}(t))}{t_r}\right) \cdot \Phi\left(-\frac{f(\mathbf{r}(t))}{t_r}\right) \qquad (7)$$

### 2.2.3  3D Gaussian Splatting (3DGS)

Introduced by Kerbl et al. [34] in 2023, 3DGS is an explicit radiance field technique for efficient and high-quality rendering of 3D scenes. Unlike conventional explicit volumetric representations, such as voxel grids, it provides a continuous and flexible representation for modeling 3D scenes in terms of differentiable 3D Gaussian-shaped primitives. These primitives are used to parameterize the radiance field and can be rendered to produce novel views. In addition, in contrast to NeRF, which relies on computationally expensive volumetric ray sampling, 3DGS achieves real-time rendering through a tile-based rasterizer. This conceptual difference is highlighted in Figure 3. This approach offers improved visual quality and faster training without relying on neural components, while also avoiding computation in empty space. More specifically, starting from multi-view images with known camera poses, 3DGS learns a set $\mathcal{G} = \{g_1, g_2, \ldots, g_N\}$ of 3D Gaussians, where $N$ denotes the number of Gaussians in the scene. Each primitive $g_i$, with $1 < i < N$, is parameterized by a full 3D covariance matrix $\mathbf{\Sigma}_i \in \mathbb{R}^{3 \times 3}$, the mean or center position $\boldsymbol{\mu}_i \in \mathbb{R}^3$, the opacity $o_i \in [0, 1]$, and color $\mathbf{c}_i$ represented by spherical harmonics (SH) for view-dependent appearance, where all the properties are learnable and optimized through back-propagation. This allows for the compact expression of the spatial influence of an individual Gaussian primitive as:

$$g_i(\mathbf{x}) = e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_i)^\top \mathbf{\Sigma}_i^{-1}(\mathbf{x}-\boldsymbol{\mu}_i)} \qquad (8)$$

Here, the spatial covariance $\mathbf{\Sigma}$ defines an ellipsoid and it is computed as $\mathbf{\Sigma} = \mathbf{R}\mathbf{S}\mathbf{S}^\top\mathbf{R}^\top$, where $\mathbf{S} \in \mathbb{R}^3$ is the spatial scale and $\mathbf{R} \in \mathbb{R}^{3 \times 3}$ represents the rotation, parameterized by a quaternion. For rendering, 3DGS operates akin to NeRF but diverges significantly in the computation of blending coefficients. Specifically, the process involves first projecting 3D Gaussian points onto a 2D image plane, a process commonly referred to as "splatting". This is done expressing the projected 2D covariance matrix and center as $\mathbf{\Sigma}' = \mathbf{J}\mathbf{W}\mathbf{\Sigma}\mathbf{W}^\mathbf{T}\mathbf{J}^\mathbf{T}$ and $\boldsymbol{\mu}' = \mathbf{J}\mathbf{W}\boldsymbol{\mu}$, where $\mathbf{W}$ represents the viewing transformation, and $\mathbf{J}$ is the Jacobian of the affine approximation of the projective transformation.

Consequently, 3DGS computes the final pixel color $C$ by blending 3D Gaussian splats that overlap at a given pixel, sorted by their depth:

$$C = \sum_{i \in \mathcal{N}} \mathbf{c}_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j) \qquad (9)$$

where the final opacity $\alpha_i$ is the multiplication result of the learned opacity $o_i$ and the Gaussian:

$$\alpha_i = o_i \exp \left( -\frac{1}{2} (\mathbf{x}' - \boldsymbol{\mu}'_i)^\top \boldsymbol{\Sigma}_i'^{-1} (\mathbf{x}' - \boldsymbol{\mu}'_i) \right) \qquad (10)$$

where $\mathbf{x}'$ and $\boldsymbol{\mu}'_i$ are coordinates in the projected space. Similarly, the depth $D$ is rendered as:

$$D = \sum_{i \in \mathcal{N}} d_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j) \qquad (11)$$

Here, $d_i$ refers to the depth of the center of the $i$-th 3D Gaussian, obtained by projecting onto the z-axis in the camera coordinate system. For optimization, instead, the process begins with parameter initialization from SfM point clouds or random values, followed by Stochastic Gradient Descent (SGD) using an L1 and D-SSIM loss function against ground truth and render views. Additionally, periodic adaptive densification handles under- and over-reconstruction by adjusting points with significant gradients and removing low-opacity points, refining scene representation and reducing rendering errors. For more details on 3DGS and related works, refer to [73], [74], [75].

## 2.3 Datasets

This section summarizes datasets commonly used in recent SLAM methodologies, covering various attributes such as sensors, ground truth accuracy, and other key factors, in both indoor and outdoor environments. Figure 4 presents qualitative examples from diverse datasets, which will be introduced in the remainder.

The **TUM RGB-D [76]**[2] dataset comprises RGB-D sequences with annotated camera trajectories, recorded using two platforms: handheld and robot, providing a diverse range of motions. The dataset features 39 sequences, some with loop closures. Core elements include color and depth images from a Microsoft Kinect sensor, captured at 30 Hz and $640 \times 480$ resolution. Ground-truth trajectories are derived from a motion-capture system with eight high-speed cameras operating at 100 Hz. The versatility of the dataset is demonstrated through various trajectories in typical office environments and an industrial hall, encompassing diverse translational and angular velocities.

The **ScanNet [77]**[3] dataset provides a collection of real-world indoor RGB-D acquisitions, featuring 2.5 million images from 1513 scans in 707 unique spaces. In particular, it includes estimated calibration parameters, camera poses, 3D surface reconstructions, textured meshes, detailed semantic segmentations at the object-level, and aligned CAD models.

The development process involved the creation of a user-friendly capture pipeline using a custom RGB-D capture



(a) ETH3D-SLAM [30]   (b) ScanNet [77]
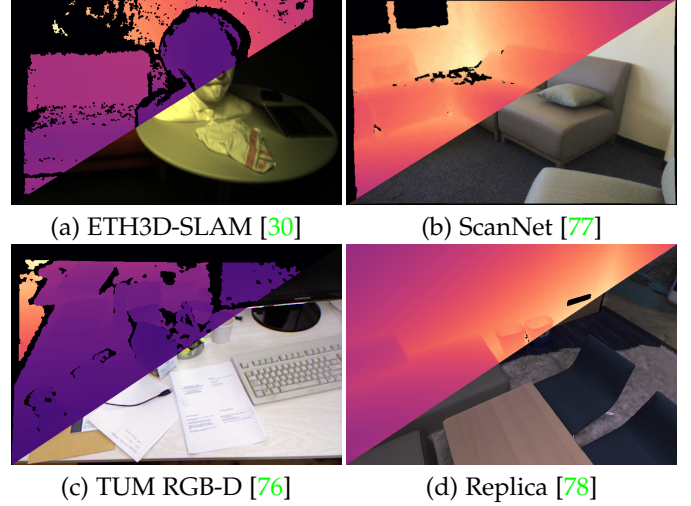
(c) TUM RGB-D [76]   (d) Replica [78]

Fig. 4: **Qualitative Comparison of Key SLAM Datasets.** RGB-D images from: (a) ETH3D-SLAM [30], (b) ScanNet [77], (c) TUM RGB-D [76], and (d) Replica [78].

setup with structure sensors attached to handheld devices such as iPads. The subsequent offline processing phase resulted in comprehensive 3D scene reconstructions, complete with available 6-DoF camera poses and semantic labels. Note that camera poses in ScanNet are derived from the BundleFusion system [53], which may not be as accurate as alternatives such as TUM RGB-D.

The **Replica [78]**[4] dataset features 18 photorealistic 3D indoor scenes with dense meshes, HDR textures, semantic data, and reflective surfaces. It spans different scene categories, includes 88 semantic classes, and incorporates 6 scans of a single space capturing different furniture arrangements and temporal snapshots. Reconstruction involved a custom-built RGB-D capture rig with synchronized IMU, RGB, IR, and wide-angle grayscale sensors, accurately fusing raw depth data through 6 degrees of freedom (DoF) poses. Although the original data was captured in the real world, the portion of the dataset used for SLAM evaluation is synthetically generated from the accurate meshes produced during reconstruction. Consequently, synthetic sequences lack real-world characteristics like specular highlights, autoexposure, blur, and more.

The **KITTI [79]**[5] dataset serves as a popular benchmark for evaluating stereo, optical flow, visual odometry/SLAM algorithms, among others. Acquired from a vehicle equipped with stereo cameras, Velodyne LiDAR, GPS and inertial sensors, the dataset contains 42,000 stereo pairs and LiDAR pointclouds from 61 scenes representing autonomous driving scenarios. The KITTI odometry dataset, with 22 LiDAR scan sequences, contributes to the evaluation of odometry methods using LiDAR data.

The **Newer College [80]**[6] dataset comprises sensor data captured during a 2.2 km walk around New College, Oxford. It includes information from a stereoscopic-inertial camera, a multi-beam 3D LiDAR with inertial measure-

---

2. https://cvg.cit.tum.de/data/datasets/RGB-D-dataset
3. http://www.scan-net.org/

4. https://github.com/facebookresearch/Replica-Dataset
5. https://www.cvlibs.net/datasets/kitti/
6. https://arxiv.org/pdf/ori.ox.ac.uk/datasets/newer-college-dataset

ments, and a tripod-mounted survey-grade LiDAR scanner, generating a detailed 3D map with around 290 million points. The dataset provides a 6 DoF ground truth pose for each LiDAR scan, accurate to approximately 3 cm. The dataset encompasses diverse environments, including built spaces, open areas, and vegetated zones.

### 2.3.1 Other Datasets

Moreover, we draw attention to less-utilized alternative datasets in recent SLAM research.

The **ETH3D-SLAM** [30][7] dataset includes videos from a custom camera rig, suitable for assessing visual-inertial mono, stereo, and RGB-D SLAM. It features 56 training datasets, 35 test datasets, and 5 independently captured training sequences using SfM techniques for ground truth.

The **EuRoC MAV** [81][8] dataset offers synchronized stereo images, IMU, and accurate ground truth for a micro aerial vehicle. It supports visual-inertial algorithm design and evaluation in diverse conditions, including an industrial setting with millimeter-accurate ground truth and a room for 3D environment reconstruction.

The **7-scenes** [82][9] dataset, created for relocalization performance evaluation, was recorded using a Kinect at $640 \times 480$ resolution. Ground truth poses were obtained through KinectFusion [2]. Sequences from different users were divided into two sets—one for simulating keyframe harvesting and the other for error calculation. The dataset presents challenges such as specularities, motion blur, lighting conditions, flat surfaces, and sensor noise.

The **ScanNet++** [83][10] dataset comprises 460 high-resolution 3D indoor scene reconstructions, dense semantic annotations, DSLR images, and iPhone RGB-D sequences. Captured with a high-end laser scanner at sub-millimeter resolution, each scene includes annotations for over 1,000 semantic classes, addressing label ambiguities and introducing new benchmarks for 3D semantic scene understanding and novel view synthesis.

The **NeuralRGBD** [72][11] dataset is a synthetic dataset and consists of 10 scenes with varying sizes, complexity, and materials. The dataset features trajectories with color and depth images rendered using BlenderProc [84], simulating noise and artifacts characteristic of real-world depth sensors. BundleFusion [53] is used to obtain an initial estimate of the camera trajectory. The camera trajectories were intentionally designed to scan only portions of the scenes, mimicking real-world scenarios.

**Additional Datasets.** For an exhaustive survey of specialized SLAM-related datasets beyond those mentioned, readers can refer to the work by Liu et al. [85]. This paper provides an in-depth exploration of a wide range of datasets designed to facilitate research and benchmarking.

## 2.4 Evaluation Metrics

The evaluation of SLAM systems typically employs several metrics across domains like 3D reconstruction, 2D depth es-

timation, trajectory estimation, and view synthesis to assess the effectiveness of methods against ground truth data.

*A. Mapping.* Metrics assessing the quality of 3D reconstruction and 2D depth estimation include:

- **Accuracy (cm)↓:** Computes the average distance between sampled points from the reconstructed mesh and the nearest ground-truth point.
- **Completion (cm)↓:** Measures the average distance between sampled points from the ground-truth mesh and the nearest reconstructed.
- **Precision (%)↑:** Indicates the proportion of points within the reconstructed mesh with Accuracy under a distance threshold $d$.
- **Recall (%)↑:** Indicates the proportion of points within the reconstructed mesh with Completion under a distance $d$. It is often referred to as *Completion Ratio*.
- **F-Score (%)↑:** An aggregate score defined as the harmonic mean between Precision and Recall.
- **L1-Depth (cm)↓:** Following [5], it computes the absolute difference between depth maps obtained from randomly sampled viewpoints from the reconstructed and the corresponding ground truth meshes respectively.

*B. Tracking.* Metrics for pose estimation, crucial for tracking performance, primarily include:

- **Absolute Trajectory Error (ATE)(cm) ↓:** Evaluates trajectory estimation accuracy by measuring the average Euclidean translation distance between corresponding poses in estimated and ground truth trajectories, often reported in terms of Root Mean Square Error (RMSE). As both trajectories can be specified in arbitrary coordinate frames, alignment is required. Importantly, this metric focuses solely on the translation component.

*C. View Synthesis.* The evaluation of view synthesis relies mainly on three visual quality assessment metrics:

- **Peak Signal to Noise Ratio (PSNR)↑:** Measures image quality by evaluating the ratio between the maximum pixel value and the root mean squared error, usually expressed in terms of the logarithmic decibel scale.
- **Structural Similarity Index Measure (SSIM [86])↑:** Assesses image quality by examining the similarities in luminance, contrast, and structural information among patches of pixels.
- **Learned Perceptual Image Patch Similarity (LPIPS [87])↓ :** Utilizes learned convolutional features to assess image quality based on feature map mean squared error across layers.

*D. Semantic Segmentation.* For SLAM methods that additionally estimate semantic information of the scene, the following metric is included to evaluate the performance of the semantic segmentation:

- **Mean Intersection over Union (mIoU)↑:** mIoU is a widely used metric for evaluating semantic segmentation performance. It is computed by calculating the IoU for each class and then taking the average across all classes. IoU is defined as the ratio of the intersection between the predicted and ground truth segmentation masks to their union. A higher mIoU indicates better semantic segmentation accuracy.

---

7. https://www.eth3d.net/slam_overview
8. https://projects.asl.ethz.ch/datasets/doku.php?id=kmavvisualinertialdatasets
9. http://research.microsoft.com/7-scenes/
10. https://cy94.github.io/scannetpp/
11. https://github.com/dazinovic/neural-rgbd-surface-reconstruction

| | | | (a) | | | | | | (b) | | (c) | | (d) | | | | | (e) | | (f) | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Section | Method | Venue | RGB-D | RGB | D | Stereo | IMU | Event Camera | Scene Encoding | Geometry Representation | Obj/Sem. Segment. | Uncert. | Frame-to Model | Frame-to Frame | External Tracker | Global BA | Loop Closure | Sub-Maps | Dyn. Env. | Extra Priors | Link |
| | | | | | | | | | **RGB-D (Sec. 3.1)** | | | | | | | | | | | | |
| | iMAP [1] | ICCV 2021 | ✓ | | | | | | MLP | Density | | | ✓ | | | | | | | | WebPage |
| | NICE-SLAM [3] | CVPR 2022 | ✓ | | | | | | Hier. Grid + MLP | Occupancy | | | ✓ | | | | | | | | Code |
| | Vox-Fusion [88] | ISMAR 2022 | ✓ | | | | | | Octree Grid + MLP | SDF | | | ✓ | | | | | | | | Code |
| | ESLAM [89] | CVPR 2023 | ✓ | | | | | | Feature Planes + MLP | TSDF | | | ✓ | | | | | | | | Code |
| | Co-SLAM [90] | CVPR 2023 | ✓ | | | | | | Hash Grid + MLP | SDF | | | ✓ | | | | | | | | Code |
| | GO-SLAM [91] | ICCV 2023 | ✓ | ✓ | | ✓ | | | Hash Grid + MLP | SDF | | | | | DROID [92] | ✓ | ✓ | | | | Code |
| | Point-SLAM [93] | ICCV 2023 | ✓ | | | | | | Neural Points + MLP | Occupancy | | | ✓ | | | | | | | | Code |
| Sec. 3.1.1 | ToF-SLAM [94] | ICCV 2023 | ✓ | | | | | | Hier. Grid + MLP | SDF | | | ✓ | | | | | | | | WebPage |
| | ADFP [95] | NeurIPS 2023 | ✓ | | | | | | Hier. Grid + MLP | Occupancy | | | ✓ | | | | | | | | Code |
| | MLM-SLAM [96] | RAL 2023 | ✓ | | | | | | MLP | Occupancy | | | ✓ | | | | | | | | |
| | Plenoxel-SLAM [97] | WACV 2024 | ✓ | | | | | | Plenoxels | Density | | | ✓ | | | | | | | | |
| | Structerf-SLAM [98] | C. & G. 2024 | ✓ | | | | | | Hier. Grid | Occupancy | | | | ✓ | ORB2 [99] | | | | | Super-pixel Segmentation [100] | Code |
| | KN-SLAM [101] | TIM 2024 | ✓ | | | | | | Hier. Grid | Occupancy | | | | ✓ | EPNP [102] | | ✓ | | | HF-Net [103] | |
| | iDF-SLAM [104] | Arx. 09/2022 | ✓ | | | | | | MLP | TSDF | | | | ✓ | URR [105] | | | | | | |
| | NeuV-SLAM [106] | Arx. 02/2024 | ✓ | | | | | | Multi Res. Voxels | SDF | | | ✓ | | | | | | | | Code† |
| | NeSLAM [107] | Arx. 03/2024 | ✓ | | | | | | Hier. Grid + MLP | SDF | | ✓ | ✓ | | | | | | | Depth Completion [61] & SuperPoint [108] | Code† |
| | GSSLAM [109] | CVPR 2024 | ✓ | ✓ | | | | | 3D Gaussians | Density | | | ✓ | | | | | | | | WebPage |
| | Photo-SLAM [110] | CVPR 2024 | ✓ | ✓ | | ✓ | | | 3D Gaussians | Density | | | | | ORB3 [111] | | ✓ | | | | |
| | SplaTAM [112] | CVPR 2024 | ✓ | | | | | | 3D Gaussians | Density | | | ✓ | | | | | | | | Code |
| | GS-SLAM [12] | Arx. 11/2023 | ✓ | | | | | | 3D Gaussians | Density | | | ✓ | | | | | | | | |
| Sec. 3.1.2 | Gaussian-SLAM [113] | Arx. 12/2023 | ✓ | | | | | | 3D Gaussians | Density | | | ✓ | | | | | ✓ | | | WebPage |
| | Compact-GSSLAM [114] | Arx. 03/2024 | ✓ | | | | | | 3D Gaussians | Density | | | ✓ | | | | | | | | |
| | GS-ICP SLAM [115] | Arx. 03/2024 | ✓ | | | | | | 3D Gaussians | Density | | | | | G-ICP [116] | ✓ | | | | | Code |
| | HF-GS SLAM [117] | Arx. 03/2024 | ✓ | | | | | | 3D Gaussians | Density | | | ✓ | | | | | | | | |
| | CG-SLAM [118] | Arx. 03/2024 | ✓ | | | | | | 3D Gaussians | Density | | ✓ | ✓ | | | | | | | NetVLAD [119] | Code |
| | MM3DGS-SLAM [120] | Arx. 04/2024 | ✓ | ✓ | | | ✓ | | 3D Gaussians | Density | | | ✓ | | | | | | | DPT [121] | WebPage |
| | MeSLAM [6] | SMC 2022 | ✓ | | | | | | MLP | Density | | | | | C-ICP [122] | | | ✓ | | | |
| | CP-SLAM [123] | NeurIPS 2023 | ✓ | | | | | | Neural Points + MLP | Occupancy | | | ✓ | | | ✓ | ✓ | ✓ | | | |
| | NISB-Map [124] | RAL 2023 | ✓ | | | | | | MLP | Density | | | | ✓ | Any | | | ✓ | | | |
| | Multiple-SLAM [125] | TIV 2023 | ✓ | | | | | | Octree Grid + MLP | SDF | | | ✓ | | | | | ✓ | | | |
| Sec. 3.1.3 | MIPS-Fusion [126] | TOG 2023 | ✓ | | | | | | MLP | TSDF | ✓ | | | ✓ | | | | ✓ | | | |
| | NGEL-SLAM [127] | ICRA 2024 | ✓ | | | | | | Octree Grid + MLP | Occupancy | ✓ | | | | ORB3 [111] | ✓ | ✓ | ✓ | | | |
| | PLGSLAM [128] | CVPR 2024 | ✓ | | | | | | Feature Planes + MLP | SDF | | | ✓ | | | | | ✓ | | | |
| | Loopy-SLAM [129] | CVPR 2024 | ✓ | | | | | | Neural Points + MLP | Occupancy | | | ✓ | | | | ✓ | ✓ | | | Code |
| | NEWTON [130] | RAL 2024 | ✓ | ✓ | | | | | Hash Grid + MLP | SDF | | | | | ORB2 [99] | | ✓ | ✓ | | | Code† |
| | Vox-Fusion++ [131] | Arx. 03/2024 | ✓ | | | | | | Octree Grid + MLP | SDF | | | ✓ | | | | ✓ | ✓ | | NetVLAD [119] | |
| | MUTE-SLAM [132] | Arx. 03/2024 | ✓ | | | | | | Feature Planes + MLP | TSDF | | | ✓ | | | ✓ | | ✓ | | | |
| | iLabel [7] | RAL 2023 | ✓ | | | | | | MLP | Density | ✓ | ✓ | ✓ | | | | | | | User | WebPage |
| | FR-Fusion [133] | ICRA 2023 | ✓ | | | | | | MLP | Density | ✓ | | ✓ | | | | | | | User & EfficientNet [134]/DINO [135] | WebPage |
| | vMap [136] | CVPR 2023 | ✓ | | | | | | MLP | Occupancy | ✓ | | ✓ | | | | | | | | Code |
| | SNI-SLAM [137] | CVPR 2024 | ✓ | | | | | | Hier. Grid + MLP | TSDF | ✓ | | ✓ | | | | | | | Dinov2 [138] | |
| Sec. 3.1.4 | NIDS-SLAM [139] | Arx. 05/2023 | ✓ | | | | | | Hash Grid + MLP | SDF | ✓ | | | ✓ | ORB3 [111] | | ✓ | | | Mask2Former [140] | |
| | DNS SLAM [141] | Arx. 11/2023 | ✓ | | | | | | Hash Grid + MLP | Occupancy | ✓ | | | ✓ | | ✓ | | | | Dinov2 [138] | |
| | SGS-SLAM [142] | Arx. 02/2024 | ✓ | | | | | | 3D Gaussians | Density | ✓ | | ✓ | | | | ✓ | | | | |
| | SemGauss-SLAM [143] | Arx. 03/2024 | ✓ | | | | | | 3D Gaussians | Density | ✓ | | ✓ | | | | | | | Dinov2 [138] | |
| | NEDS-SLAM [144] | Arx. 03/2024 | ✓ | | | | | | 3D Gaussians | Density | ✓ | | ✓ | | | | | | | Depth Anything [145] | |
| | DN-SLAM [13] | Sensors J. 2023 | ✓ | | | | | | Hash Grid + MLP | Density | | | | | ORB3 [111] | | | | ✓ | SAM [146] | |
| | DynaMoN [147] | Arx. 09/2023 | ✓ | ✓ | | | | | HexPlane + MLP | Density | | | | | DROID [92] | | | | ✓ | DeepLabV3 [148] | Code† |
| Sec. 3.1.5 | DDN-SLAM [9] | Arx. 01/2024 | ✓ | ✓ | | | | | Hash Grid + MLP | SDF | ✓ | | | | ORB3 [111] | ✓ | ✓ | | ✓ | YOLOv5 & ZoeDepth [149] | |
| | NID-SLAM [150] | Arx. 01/2024 | ✓ | | | | | | Hier. Grid + MLP | TSDF | | | ✓ | | | | | | ✓ | | |
| | DVN-SLAM [151] | Arx. 03/2024 | ✓ | | | | | | Feature Planes + MLP | TSDF | | ✓ | ✓ | | | | | | ✓ | | |
| Sec. 3.1.6 | OpenWorld-SLAM [152] | CRV 2023 | ✓ | | | | | ✓ | Hier. Grid + MLP | Occupancy | ✓ | ✓ | ✓ | | | | | | | | |
| | UncLe-SLAM [17] | ICCVW 2023 | ✓ | | ✓ | | | | Hier. Grid + MLP | Occupancy | | ✓ | ✓ | | | | | | | | Code |
| Sec. 3.1.7 | EN-SLAM [14] | CVPR 2024 | ✓ | | | | | ✓ | Hier. Grid + MLP | TSDF | | | ✓ | | | ✓ | | | | | |
| | | | | | | | | | | **RGB (Sec. 3.2)** | | | | | | | | | | | | |
| | DIM-SLAM [16] | ICLR 2023 | | ✓ | | | | | Hier. Grid + MLP | Density | | | ✓ | | | | | | | | Code |
| | Orbeez-SLAM [153] | ICRA 2023 | | ✓ | | | | | Hash Grid + MLP | Density | | | | ✓ | ORB2 [99] | | | | | | Code |
| Sec. 3.2.1 | FMapping [154] | Arx. 06/2023 | | ✓ | | | | | Grid Fact. + MLP | Density | | ✓ | ✓ | | | | | | | | Code† |
| | TT-HO-SLAM [155] | Arx. 12/2023 | | ✓ | | | | | Hier. Grid + MLP | Density | | | ✓ | | | | | | | | |
| | Hi-Map [8] | Arx. 01/2024 | | ✓ | | | | | Grid Fact. + MLP | SDF | | | ✓ | | | | | | | | WebPage |
| | iMode [156] | ICRA 2023 | | ✓ | | | | | MLP | Density | | | | | ORB [10] | | | | | Sparse-to-dense [157] | |
| | Hi-SLAM [158] | RAL 2023 | | ✓ | | | | | Hash Grid + MLP | TSDF | | | | | DROID [92] | ✓ | ✓ | | | Omnidata [159] | |
| | NICER-SLAM [18] | 3DV 2024 | | ✓ | | | | | Hash Grid + MLP | SDF | ✓ | | ✓ | | | | | | | GMFlow [160] & Omnidata [159] | WebPage |
| Sec. 3.2.2 | NeRF-VO [161] | Arx. 12/2023 | | ✓ | | | | | Hash Grid + MLP | SDF | | | | | DPVO [162] | ✓ | | | | Omnidata [159] | |
| | MoD-SLAM [163] | Arx. 02/2024 | ✓ | ✓ | | | | | Hash Grid + MLP | SDF | | | | | DROID [92] | | ✓ | | | DPT [121] & ZoeDepth [149] | |
| | Q-SLAM [164] | Arx. 03/2024 | ✓ | ✓ | | | | | Grid Fact. + MLP | Density | | | | | DROID [92] | | | | | Segmentation Network | |
| | GlORIE-SLAM [165] | Arx. 03/2024 | | ✓ | | | | | Neural Points + MLP | Occupancy | | | | | DROID [92] | ✓ | ✓ | | | Omnidata [159] | |
| Sec. 3.2.3 | RO-MAP [166] | RAL 2023 | | ✓ | | | | | Hash Grid + MLP | Occupancy | ✓ | | | | ORB2 [99] | | | | | YOLOv8 | Code |
| Sec. 3.2.4 | NeRF-SLAM [167] | IROS 2023 | | ✓ | | | | | Hash Grid + MLP | Density | | ✓ | | | DROID [92] | | | | | | Code |
| | | | | | | | | | | **LiDAR (Sec. 3.3)** | | | | | | | | | | | | |
| | NeRF-LOAM [15] | ICCV 2023 | | | ✓ | | | | Octree Grid + MLP | SDF | | | ✓ | | | | | | | | Code |
| Sec. 3.3.1 | LONER [168] | RAL 2023 | | | ✓ | | | | Hier. Grid + MLP | SDF | | | | | P2P-ICP [169] | | | | | | Code |
| | PIN-SLAM [170] | Arx. 01/2024 | | | ✓ | | | | Neural Points + MLP | SDF | ✓ | | ✓ | | | | ✓ | | ✓ | | Code |
| Sec. 3.3.2 | LIV-GaussMap [171] | Arx. 01/2024 | ✓ | | | ✓ | | | 3D Gaussians | Density | | | ✓ | | | | | | | | Code† |
| | MM-Gaussian [172] | Arx. 04/2024 | ✓ | | | | | | 3D Gaussians | Density | | | | | Kiss-ICP [173] | | | | | SuperPoint [108] & LightGlue [174] | |

TABLE 1: **SLAM Systems Overview.** We categorize the different methods into main RGB-D, RGB, and LiDAR-based frameworks. In the leftmost column, we identify sub-categories of methods sharing specific properties, detailed in Sections 3.2.1 to 3.3.2. Then, for each method, we report, from the second leftmost column to the second rightmost, the method name and publication venue, followed by (a) the input modalities they can process: RGB, RGB-D, D (*e.g.* LiDAR, ToF, Kinect, etc.), stereo, IMU, or events; (b) mapping properties: scene encoding and geometry representations learned by the model; (c) additional outputs learned by the method, such as object/semantic segmentation, or uncertainty modeling (Uncert.); (d) tracking properties related to the adoption of a frame-to-frame or frame-to-model approach, the utilization of external trackers, Global Bundle Adjustment (BA), or Loop Closure; (e) advanced design strategies, such as modeling sub-maps or dealing with dynamic environments (Dyn. Env.); (f) the use of additional priors. Finally, we report the link to the project page or source code in the rightmost column. † indicates code not released yet.

# 3 SIMULTANEOUS LOCALIZATION AND MAPPING

This section introduces latest SLAM systems that leverage recent progress in radiance field representations. Organized in a method-based taxonomy, the papers are categorized by their approaches, offering readers a clear and organized presentation. The section begins with a basic classification into RGB-D (3.1), RGB (3.2), and LiDAR (3.3) methodologies, setting the stage for the development of specific subcategories. Each category lists officially published papers in conferences/journals by publication dates, followed by preprints from arXiv arranged by their initial preprint dates.

For a comprehensive understanding, Table 1 offers a detailed overview of the surveyed methods. This table provides an in-depth summary, highlighting key features of each method, and includes references to project pages or source code whenever available. For further details or method specifics, please refer to the original papers.

## 3.1 RGB-D SLAM Approaches

Here we focus on dense SLAM techniques using RGB-D cameras that capture both color images and per-pixel depth information of the environment. These techniques fall into distinct categories: NeRF-style SLAM solutions (3.1.1) and alternatives based on the 3D Gaussian Splatting representation (3.1.2). Specialized solutions derived from both approaches include submap-based SLAM methods for large scenes (3.1.3), frameworks that address semantics (3.1.4), and those tailored for dynamic scenarios (3.1.5). Within this classification, some techniques assess reliability through uncertainty (3.1.6), while others explore the integration of additional sensors like event-based cameras (3.1.7).

### 3.1.1 NeRF-style RGB-D SLAM

Recent advances in implicit neural representations have enabled accurate and dense 3D surface reconstruction. This
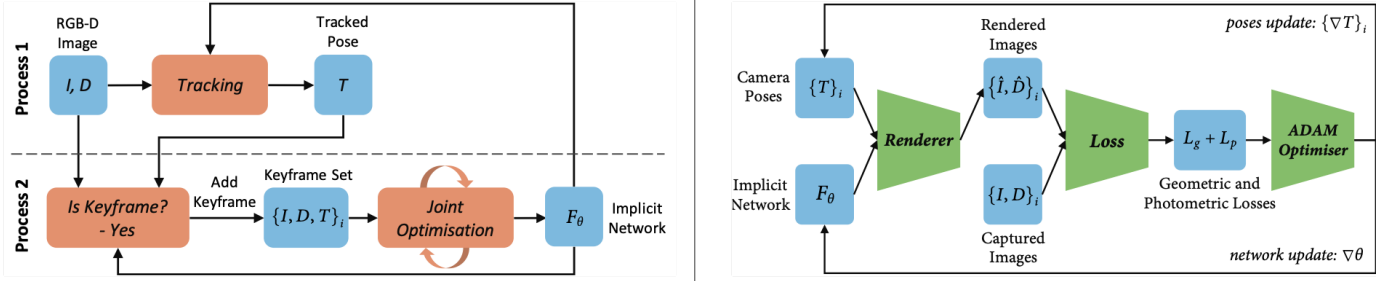
Fig. 5: **Overview of iMap [1], the Pioneering Approach in Neural Implicit-based SLAM.** (Left) The illustration depicts two concurrent processes: *tracking*, optimizing the current frame's pose within the locked network; *mapping*, jointly refining the network and camera poses of selected keyframes. (Right) Jointly optimizing scene network parameters and camera poses for keyframes using differentiable rendering functions. Figure from [1].

has led to novel SLAM systems derived from or inspired by NeRF, initially designed for offline use with known camera poses. In this section, we describe these dense neural VSLAM methods, analyze their main features, and provide a clear overview of their strengths and weaknesses.

**iMAP [1]**. This work marks the first attempt to leverage implicit neural representations for SLAM. This ground-breaking achievement not only pushes the boundaries of SLAM but also establishes a new direction for the field. In particular, iMAP demonstrates the potential of an MLP to dynamically create a scene-specific implicit 3D model. By doing so, this work provides an alternative to conventional techniques, offering efficient geometry representation, automatic detail control, and seamless filling-in of unobserved regions. Specifically, the framework (depicted in Figure 5) uses an MLP to map 3D coordinates to color and volume density, without addressing specularities or considering viewing directions. Differentiable rendering, guided by camera pose and pixel coordinates, generates depth and color images through network queries. Joint optimization of photometric and geometric losses for a fixed set of keyframes refines network parameters and camera poses. A parallel process ensures close-to-frame-rate camera tracking, with dynamic keyframe selection based on information gain. Active sampling optimizes a sparse set of pixels, guided by loss statistics. Integrating keyframe-based structure, multi-processing computation, and dynamic pixel sampling enables real-time tracking and global map updating. However, due to the limited capacity of the model, it leads to less detailed reconstruction and faces challenges with catastrophic forgetting in larger environments.

**NICE-SLAM [5]**. In contrast to iMAP's use of a single MLP as the scene representation, NICE-SLAM adopts a hierarchical strategy that integrates multi-level local data. This approach effectively addresses issues like excessively smoothed reconstructions and scalability limitations in larger scenes. By optimizing a hierarchical representation using pre-trained geometric priors, NICE-SLAM achieves high-quality scene reconstruction. This procedure involves representing geometry by encoding it into three voxel grids of varying resolutions, each associated with its corresponding pre-trained MLP decoder — coarse, mid, and fine levels. Moreover, a dedicated feature grid and decoder are utilized for capturing scene appearance. Notably, compared to iMAP, it updates only the visible grid features at each

step. This strategy significantly enhances the precision and efficiency of optimization processes, overcoming the constraints of iMAP's global updates and addressing its issues with catastrophic forgetting. However, the predictive capabilities are limited to the scale of the coarse representation, and it does not incorporate loop closures.

**Vox-Fusion [88]**. This work combines traditional volumetric fusion methods with neural implicit representations. Specifically, it leverages a voxel-based neural implicit surface representation to encode and optimize the scene within each voxel. While sharing similarities with NICE-SLAM [5], its distinctiveness lies in its adoption of an octree-based structure to enable a dynamic voxel allocation strategy. This dynamic approach empowers real-time tracking and mapping for diverse scenes, eliminating the need for a pre-allocated hierarchical voxel grid, as in NICE-SLAM. Moreover, a fundamental feature of Vox-Fusion involves modeling the local scene geometry within individual voxels using a continuous SDF. This SDF is encoded via a neural implicit decoder along with shared feature embeddings. The use of shared embedding vectors enables a more lightweight decoder, benefiting from its capacity to encapsulate local geometry and appearance knowledge. Furthermore, Vox-Fusion introduces an efficient keyframe selection strategy tailored for sparse voxels, further enhancing its capability for efficient map management.

**ESLAM [89]**. The core of ESLAM is its implementation of multi-scale axis-aligned feature planes, diverging from traditional voxel grids. This approach optimizes memory usage through quadratic scaling, in contrast to the cubic growth exhibited by voxel-based models. These tri-plane architectures store and optimize features on perpendicular axes, enhancing reconstruction quality and addressing the forgetting problem by managing geometry and appearance changes separately. The method employs three coarse and three fine feature planes for both geometry and scene appearance. Furthermore, ESLAM adopts TSDF as the geometric representation. This improves convergence speed and enhances reconstruction quality compared to conventional rendering-based methods like volume density in iMAP and occupancy in NICE-SLAM. The combination of multi-scale feature planes and TSDF representation improves reconstruction and localization while processing frames up to ten times faster than iMAP and NICE-SLAM.

**Co-SLAM [90]**. Systems such as iMAP use coordi-

nate networks for real-time SLAM. However, to ensure interactive operation, they adopt strategies like sparsified ray sampling, potentially leading to oversmoothing issues. In contrast, methods such as NICE-SLAM use parametric embeddings through feature grids, which partially avoid oversmoothing but cannot fully address hole-filling. Co-SLAM fills the gap by combining the smoothness of coordinate encodings (using one-blob encoding [175]) with the fast convergence and local detail advantages of sparse parametric encodings (using a hash grid [54]). Consequently, this results in more robust camera tracking, high-fidelity maps, and improved hole-filling. In addition, unlike previous neural SLAM systems, Co-SLAM performs global bundle adjustment (BA) by sampling few rays from all previous keyframes (around 5% of pixels for each keyframe).

**GO-SLAM** [91]. The absence of global optimization techniques like LC and BA in previous works leads to tracking errors over time. To address these limitations, GO-SLAM is designed for real-time global optimization of camera poses and 3D reconstructions. At the core is a robust pose estimation module, integrating efficient LC and online full BA that utilizes the full history of input frames to ensure accurate trajectory estimation and to maintain a coherent 3D map representation. Specifically, the architecture operates through three parallel threads: *front-end tracking*, responsible for iterative pose and depth updates along with efficient loop closing; *back-end tracking*, focused on generating globally consistent pose and depth predictions via full BA; and *instant mapping*, which updates the 3D reconstruction based on the latest available poses and depths. GO-SLAM's instant mapping draws inspiration from the Instant-NGP framework [54]. It employs a rendering strategy that maps 3D points to multi-resolution hash encodings and predicts both SDF and color using shallow networks. Notably, GO-SLAM supports monocular, stereo, and RGB-D cameras.

**Point-SLAM** [93]. Unlike grid-based or network-based methods, Point-SLAM introduces a dynamic neural point cloud representation, adjusting point density based on input data information, ensuring more points in areas with higher detail and fewer points in less informative regions. The method uses the per pixel input image gradient magnitude to determine the density of the points. Depth and color images are rendered via volume rendering, with each pixel ray extracting geometric and color features from point groups. The features extracted are further processed by specialized decoders, as in [5], to compute occupancy and color values and optimized via gradient descent using an RGB-D re-rendering loss. The mapping process runs alternatively to the tracking process to update the scene and to estimate the location of the camera. A key feature is that its neural point cloud representation expands incrementally during exploration, stabilizing as all relevant regions are incorporated. Unlike voxel-based methods, this strategy optimizes memory usage by only adding points in a region around the surface, removing the need to model free space. Moreover, due to the dynamic resolution of the point cloud, areas with few details are compressed, further saving memory.

**ToF-SLAM** [94]. This work presents the first SLAM system that leverages both a monocular camera and a lightweight ToF sensor, which is limited to providing coarse measurements in the form of low-resolution depth dis-

tributions. To achieve this, a multi-modal feature grid is introduced, offering the ability to perform both zone-level rendering tailored for ToF sensors and pixel-level rendering optimized for other high-resolution signals (*e.g.* RGB). The system optimizes camera poses and scene geometry by comparing these rendered signals to the raw sensor inputs. Additionally, a predicted depth is employed for intermediate supervision, enhancing pose tracking and reconstruction accuracy. The authors also develop a coarse-to-fine optimization strategy to efficiently learn the implicit representation. Furthermore, temporal information is incorporated to handle noisy ToF sensor signals, enhancing system accuracy.

**ADFP** [95]. This work incorporates an attentive depth fusion prior derived from TSDF formed by fusing multiple depth images. This allows neural networks to directly utilize learned geometry and TSDFs during volume rendering, overcoming issues such as incomplete depth at holes and unawareness of occluded structures in the reconstruction process. Through a process involving ray tracing, feature interpolation, occupancy prediction priors, and an attention mechanism to balance the contributions of learned geometry and the depth fusion prior, the methodology significantly enhances accuracy in 3D reconstruction.

**MLM-SLAM** [96]. This work introduces a multi-MLP hierarchical scene representation that utilizes different levels of decoders to extract detailed features, enhancing the reconstruction process without sacrificing scalability. The system employs neural implicit representations, optimizing depth and color estimation through geometric and photometric losses without fixed pre-trained decoders, ensuring better generalization across various scenes. Additionally, it implements a refined tracking strategy and keyframe selection approach, enhancing system reliability, especially in challenging dynamic environments.

**Plenoxel-SLAM** [97]. This work builds upon the Plenoxel radiance field model [66], devoid of neural networks. The paper describes a novel approach: the use of a voxel grid representation and trilinear interpolation within the Plenoxel framework for efficient dense mapping and tracking. The key highlight lies in the analytical derivation of equations essential for both offline RGB-D mapping and online camera pose optimization. Despite the novel approach outlined by Plenoxel-SLAM, it is worth mentioning that no explicit 3D mesh is currently reconstructed from the learned representation.

**Structerf-SLAM** [98]. This methodology uses two-layer feature grids and pre-trained decoders to decode interpolated features into RGB and depth values. During the tracking phase, the use of three-dimensional planar features, based on the Manhattan assumption, improves stability and rapid data association, overcoming the limitations of insufficient texture. Camera pose optimization involves the application of photometric, geometric, and planar feature matching loss terms. In the mapping stage, a planar consistency constraint ensures that the depth predicted by the dual-layer neural radiance field aligns with a plane, resulting in smoother map reconstruction.

**KN-SLAM** [101]. KN-SLAM integrates sparse feature-based localization using HF-Net [103] with neural implicit representations for mapping. It consists of three concurrent threads: 1) tracking, which extracts local and global features

for initial pose estimation and single-frame pose optimization; 2) mapping, which updates the scene representation and jointly optimizes camera poses and implicit mapping; and 3) optimization, which performs loop detection, pose graph optimization, and global refinement of the neural implicit map. KN-SLAM adopts a hierarchical scene representation with feature grids and decoders, similarly to NICE-SLAM [5], and uses a combination of photometric, depth, and reprojection losses in the optimization process to ensure consistency between the implicit scene representation and sparse feature observations.

**iDF-SLAM [104]**. This work integrates a feature-based neural tracker at the front-end for robust camera tracking. The back-end, instead, includes a neural implicit mapper using a single MLP as the map representation and is responsible for estimating TSDF values. Notably, in addition to the 3D position and unlike previous approaches, the MLP takes as input view directions to enhance per-frame rendering outputs. Keyframe updates involve pose optimization and MLP weight tuning, preventing catastrophic forgetting through replay-based keyframe buffering. The selection of keyframes is strategically based on covisibility scores, enhancing map optimization robustness. The front-end features an unsupervised R&R (URR) model [105] for camera tracking, utilizing deep features and point cloud registration. iDF-SLAM employs runtime fine-tuning and point cloud registration for improved tracking accuracy, with the tracker's feature extractor pre-trained on ScanNet.

**NeuV-SLAM [106]**. This framework exploits a voxel-like representation to encode the scene geometry. Specifically, NeuV-SLAM handles multi-resolution voxels using a hash table, allowing for incremental expansion in newly explored areas, named hashMV. This is coupled with a novel implicit scene representation, VDF, that combines the implementation of neural SDF voxels with the SDF activation strategy by directly optimizing color features and SDF values that are anchored within the voxels.

**NeSLAM [107]**. NeSLAM integrates a neural implicit scene representation with hierarchical feature grids, a depth completion/denoising network, and a self-supervised feature tracking method. The depth network provides dense depth images with uncertainty estimates to guide neural point sampling and optimize the implicit representation, while the hierarchical feature grid combined with MLP(s) estimates SDF values for improved geometry. The self-supervised feature tracking network enables online optimization and enhances generalization across different scenes. The system jointly optimizes the implicit scene representation and camera pose using tailored loss functions, including patch-wise depth variance, color, and depth losses, as well as an ICP loss.

### 3.1.2 3DGS-style RGB-D SLAM

Here, we present an overview of pioneering frameworks that use explicit volumetric representations based on 3D Gaussian Splatting for the development of SLAM solutions. These approaches typically exploit the advantages of 3DGS, such as faster and more photorealistic rendering compared to other existing scene representations. They also offer the flexibility to increase map capacity by adding more Gaussian primitives, complete utilization of per-pixel dense photometric losses, and direct parameter gradient flow to facilitate fast optimization. To date, the 3DGS representation has primarily been employed in offline systems dedicated to novel view synthesis from known camera poses. In the following section, we introduce seminal SLAM methodologies that enable the simultaneous optimization of scene geometry along with camera poses.

**GSSLAM [109]**.This work introduces a paradigm shift by leveraging 3D Gaussians as the representation coupled with splatting rendering techniques. Specifically, this system employs 3DGS as its only representation for online 3D reconstruction using a single moving RGB or RGB-D camera. The framework includes several key components, such as tracking and camera pose optimization, Gaussian shape verification and regularization, mapping and keyframing, and resource allocation and pruning. The tracking phase adopts a direct optimization scheme against the 3D Gaussians, providing fast and robust tracking capability with a broad basin of convergence for the camera pose estimation. Meanwhile, geometric verification and regularization techniques are introduced to handle ambiguities in incremental 3D dense reconstruction, with a novel Gaussian shape regularization proposed to ensure geometric consistency. For mapping and keyframing, GSSLAM integrates techniques for efficient online optimization and keyframe management, which involves selecting and maintaining a small window of keyframes based on inter-frame covisibility. Additionally, resource allocation and pruning methods are used to eliminate unstable Gaussians and avoid artifacts in the model.

**Photo-SLAM [110]**. This work integrates explicit geometric features and implicit texture representations within a hyper primitives map. This methodology combines ORB features [176], rotation, scaling, density, and spherical harmonic coefficients to optimize camera poses and mapping accuracy while minimizing a photometric loss. Notably, Photo-SLAM employs a multi-threaded architecture encompassing modules for localization, mapping, photorealistic rendering, and loop closure. This design facilitates efficient factor graph solving, sparse 3D point generation, and progressive optimization of hyper primitives. A key feature lies in the utilization of 3DGS [34] for image rendering from the hyper primitives map. By leveraging advanced techniques, including geometry-based densification and Gaussian-Pyramid-based learning, the framework achieves high-quality rendering, increased mapping accuracy, and real-time operability. The paper evaluates Photo-SLAM with diverse camera setups, including stereo, monocular, and RGB-D.

**SplaTAM [112]**. This method represents the scene as a collection of simplified 3D Gaussians, enabling high-quality color and depth image rendering. The SLAM pipeline encompasses several key steps: *Camera Tracking*: Minimizing re-rendering errors for precise camera pose estimation, focusing on visible silhouette pixels and optimizing within well-structured map regions. *Gaussian Densification*: Adds new Gaussians based on the rendered silhouette and depth information, enhancing the scene representation only where needed for accuracy. *Map Update*: Refines the Gaussian parameters across frames, minimizing RGB and depth errors while optimizing over influential frames to update the geometry of the scene. By adopting this approach, SplaTAM
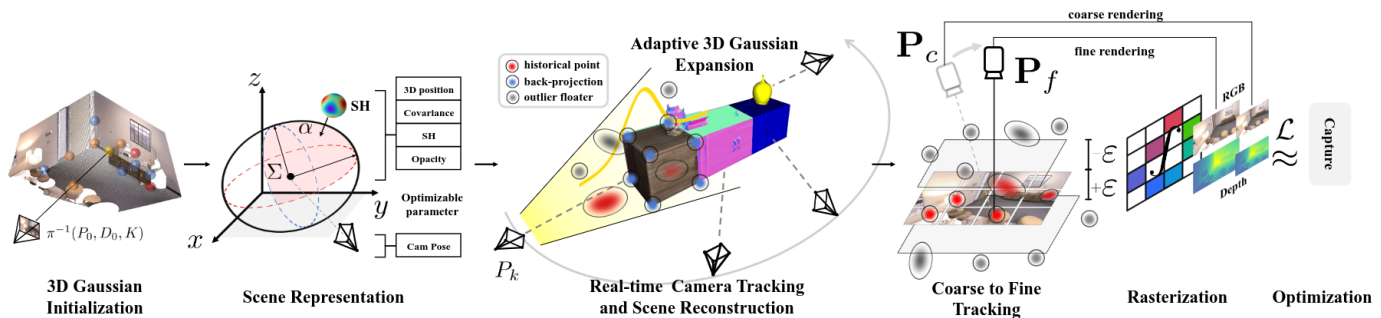
Fig. 6: **Overview of GS-SLAM [12].** This framework leverages the 3D Gaussian scene representation and rendered RGB-D images for inverse camera tracking. Through a novel Gaussian expansion strategy, GS-SLAM achieves real-time tracking, mapping, and rendering on GPUs, enhancing scene reconstruction capabilities. Figure from [12].

fundamentally redefines dense SLAM practices, offering advancements in rendering efficiency, optimization speed, and spatial mapping capabilities. While the study shows remarkable progress, it also acknowledges limitations such as sensitivity to motion blur and depth noise.

**GS-SLAM [12].** In contrast to methods relying on neural implicit representations, GS-SLAM employs a novel approach utilizing 3D Gaussians together with opacity and spherical harmonics to encapsulate both scene geometry and appearance. A key contribution lies in its adaptive expansion strategy, dynamically managing the addition or removal of 3D Gaussians to efficiently reconstruct observed scene geometry while enhancing mapping precision. Additionally, GS-SLAM introduces a robust coarse-to-fine camera tracking technique refining camera pose estimation iteratively through image refinement and reliable 3D Gaussian selection. This is followed by BA, aiming to optimize camera poses and the 3D Gaussian scene representation simultaneously. Despite its strengths, GS-SLAM faces limitations concerning its dependence on high-quality depth information and substantial memory usage in large scenes.

**Gaussian-SLAM [113].** This method efficiently seeds new Gaussians for newly explored areas and optimizes them online by organizing the scene into independently optimized sub-maps, allowing scalability to larger scenes. The sub-maps are initialized based on the camera motion, and new Gaussians are added to the active sub-map by considering the rendered alpha values and color gradients of the keyframes. The Gaussian parameters within each sub-map are jointly optimized using photometric and geometric losses. Frame-to-model camera tracking is performed by minimizing the photometric and geometric losses between the input and rendered frames while using soft alpha and inlier masks to handle occlusions and outliers.

**Compact-GSSLAM [114].** The proposed method introduces a compact 3DGS SLAM system that reduces the number and parameter size of Gaussian ellipsoids, addressing the high memory usage, storage requirements, and slow training speed issues of existing 3D Gaussian-based SLAM approaches. A novel sliding window-based online masking strategy is employed to remove redundant Gaussian ellipsoids during operation, while a geometry codebook compresses the geometric attributes (scale and rotation) of the remaining ellipsoids by exploiting similarities across the scene. This compact representation enables faster rendering

speeds and efficient memory usage. Furthermore, a global bundle adjustment method incorporating reprojection loss is utilized for robust camera pose estimation.

**GS-ICP SLAM [115].** The GS-ICP SLAM approach is a SLAM system that combines two techniques - Generalized Iterative Closest Point (G-ICP) [116] and 3DGS. Unlike previous methods that rely primarily on 2D image-based tracking, this approach actively utilizes 3D information by using G-ICP for the tracking process. This allows the system to directly use the 3D Gaussian map representation for tracking, without the need for additional post-processing. A key novelty is the mutual sharing of covariance information between the tracking and mapping components. The covariances computed during the G-ICP tracking are used to initialize the 3DGS mapping, while the 3D Gaussians in the map are in turn used as 3D points and their covariances for the G-ICP tracking. This bidirectional exchange of information, facilitated by scale alignment techniques, minimizes redundant computations and enables an efficient and high-performance SLAM system.

**HF-GS SLAM [117].** HF-GS SLAM, based on 3DGS, introduces two key novelites. First, it proposes a Gaussian densification strategy guided by the rendering loss to map unobserved areas and refine reobserved regions. Second, it incorporates regularization parameters during mapping to alleviate the forgetting problem. The mapping involves optimizing Gaussian parameters by minimizing the loss between rendered and input images, along with a regularization term to prevent overfitting and preserve details of previously visited areas.

**CG-SLAM [118].** CG-SLAM uses a novel uncertainty-aware 3D Gaussian field for consistent and stable tracking and mapping. The system conducts a mathematical analysis of camera pose derivatives in the EWA (Elliptical Weighted Average) splatting process and develops a CUDA framework to decouple tracking and mapping components. To reduce overfitting and achieve a consistent Gaussian field, CG-SLAM employs techniques such as scale regularization, depth alignment, and a depth uncertainty model to guide the selection of informative Gaussian primitives. The system uses various loss functions to update Gaussian properties and enables efficient and accurate tracking by minimizing a re-rendering loss from low-uncertainty primitives.

**MM3DGS-SLAM [120].** MM3DGS SLAM is the first visual-inertial SLAM framework that utilizes 3D Gaussians
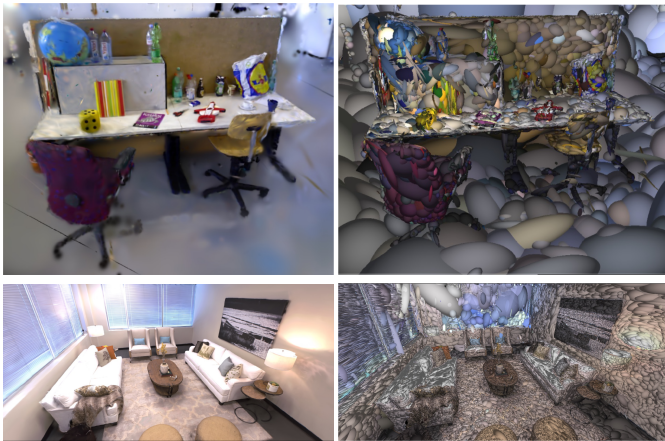
Fig. 7: **3D Gaussian Visualization.** (Left) Rasterized Gaussians, (Right) Gaussians shaded to highlight the underlying geometry. Images adapted from [109].

as the map representation. The framework takes inputs from a monocular camera or RGB-D camera, along with inertial measurements. The camera pose is optimized using a combined tracking loss incorporating depth measurements and IMU pre-integration. Keyframe selection is based on image covisibility and the Naturalness Image Quality Evaluator (NIQE) metric across a sliding window. New 3D Gaussians are initialized for keyframes with low opacity and high depth error, with positions initialized using depth measurements or estimates. The mapping stage optimizes the 3D Gaussian parameters according to a loss that includes photometric rendering quality, structural similarity, and depth correlation. The framework addresses the scale ambiguity of monocular depth estimates by solving for a scaling and shift that fits the depth estimate to the current map. The authors also release a multi-modal dataset, UT-MM, collected from a mobile robot equipped with a camera and an IMU.

### 3.1.3 Submaps-based SLAM

In this category, we focus on methods that address the challenges of catastrophic forgetting and the applicability issues in large environments faced by the previously discussed dense radiance field-inspired SLAM systems.

**MeSLAM [6].** MeSLAM introduces a novel SLAM algorithm for large-scale environment mapping with minimal memory footprint. This is achieved by combining a neural implicit map representation with a novel network distribution strategy. Specifically, by using distributed MLP networks, a global mapping module facilitates the segmentation of the environment into distinct regions and coordinates the stitching of these regions during the reconstruction process. This allows for the creation of a comprehensive global map that captures the entire environment or specific parts thereof. Additionally, a key aspect relies in its integration of external odometry [122] with neural field-based optimization. This combination enhances the system's ability to robustly track poses in regions where maps intersect, leading to improved accuracy and stability. The joint optimization mechanism optimizes both neural network parameters and poses simultaneously, refining the map representation while ensuring precise localization.

**CP-SLAM [123].** This work stands as a collaborative neural implicit SLAM approach, characterized by a unified framework encompassing front-end and back-end modules. At its core, it leverages a neural point-based 3D scene representation associated with keyframes. This allows for seamless adjustments during pose optimization and enhances collaborative mapping capabilities. Leveraging a unique learning strategy, CP-SLAM employs a distributed-to-centralized approach to ensure consistency and cooperation among multiple agents. The method's front-end modules use neural point clouds and differentiable volume rendering to achieve efficient odometry, mapping, and tracking. Additionally, CP-SLAM implements loop detection and submap alignment techniques to mitigate pose drift. The approach concludes with global optimization techniques such as pose graph optimization and map refinement.

**NISB-Map [124].** NISB-Map uses multiple small MLP networks, following the design of iMAP [1], to represent the large-scale environment in compact spatial blocks. Alongside sparse ray sampling with depth priors, this enables scalable indoor mapping with low memory usage. Sparse ray sampling, however, can result in varying density levels among adjacent spatial blocks, leading to inconsistencies in density. To remedy this, a distillation procedure for overlapping Neural Implicit Spatial Block (NISBs) is implemented, effectively minimizing density variations and ensuring geometric consistency. In this process, knowledge from the last trained NISB serves as the teacher and is distilled only within overlapping regions with the current NISB. This ensures continuity while reducing computation and training time compared to training an extra global NISB.

**Multiple-SLAM [125].** This paper introduces a novel collaborative implicit SLAM framework to tackle catastrophic forgetting. By employing multiple SLAM agents to process scenes in blocks, it minimizes trajectory and mapping errors. The system empowers agents in the frontend to operate independently, while also facilitating the sharing and fusion of map information in the backend server. Specifically, the architecture enables complex scene reconstruction through collaborative pose estimation and map fusion processes. The pose estimation process efficiently determines relative poses between agents using a two-stage approach: matching keyframes through a NetVLAD-based [119] global descriptor extraction model and fine-tuning inter-agent poses through an implicit relocalization process. Conversely, the map fusion stage integrates local maps using a floating-point sparse voxel octree for precise alignment. Moreover, for more accurate and efficient map fusion, the method addresses overlapping regions by removing redundant voxels based on observation confidence and a reconstruction loss.

**MIPS-Fusion [126].** This work introduces a divide-and-conquer mapping scheme for online dense RGB-D reconstruction, using a grid-free, purely neural approach with incremental allocation and on-the-fly learning of multiple neural submaps, as depicted in Figure 8. It also incorporates efficient on-the-fly learning through local bundle adjustment, distributed refinement with back-end optimization, and global optimization through loop closure. Moreover, the methodology includes a hybrid tracking scheme, combining gradient-based and randomized optimizations via particle filtering to ensure robust performance, particularly under

fast camera motions. Key features include a depth-to-TSDF loss for efficient fitness evaluation, a lightweight network for classification-based TSDF prediction, and support for parallel submap fine-tuning. Notably, MIPS-Fusion detects loop closures via covisibility thresholds, which does not allow for the correction of large drifts.

**NGEL-SLAM [127].** Utilizing two modules, namely the tracking and mapping modules, this system integrates the robust tracking capabilities of ORB-SLAM3 [111] with the scene representation provided by multiple implicit neural maps. Operating through three concurrent processes—tracking, dynamic local mapping, and loop closing—the system ensures global consistency and low latency. The tracking module, based on ORB-SLAM3, estimates real-time camera poses and identifies keyframes, which are then processed in the dynamic local mapping phase for local BA and efficient scene representation training. Loop closing optimizes poses using global BA, and the system's utilization of multiple local maps minimizes re-training time, ensuring a quick response to significant changes in tracking poses. The system further incorporates uncertainty-based image rendering for optimal sub-map selection, and its scene representation is based on a sparse octree-based grid with implicit neural maps, achieving memory efficiency and accurate representation of the environment.

**PLGSLAM [128].** The progressive scene representation method proposed in this work divides the entire scene into multiple local scene representations, allowing for scalability to larger indoor scenes and improving robustness. As the local scene representation, the system utilizes axis-aligned triplanes for high-frequency features and an MLP for global low-frequency features. This allows for accurate and smooth surface reconstruction. Additionally, it reduces memory growth from cubic to square with respect to the scene size, enhancing scene representation efficiency. Moreover, the system integrates traditional SLAM with an end-to-end pose estimation network, introducing a local-to-global BA algorithm to mitigate cumulative errors in large-scale indoor scenes. Efficiently managing keyframe databases during operation enables seamless BA across all past observations.

**Loopy-SLAM [129].** This system leverages neural point clouds in the form of submaps for local mapping and tracking. The method employs frame-to-model tracking with a data-driven point-based submap generation approach, dynamically growing submaps based on camera motion during scene exploration. Global place recognition triggers loop closures online, enabling robust pose graph optimization for global alignment of submaps and trajectory. The point-based representation facilitates efficient map corrections without storing the entire history of input frames, compared to previous methods such as [91]. The system addresses challenges of error accumulation in camera tracking and avoids visible seams in overlapping regions.

**NEWTON [130].** Most of the neural SLAM systems use a world-centric map representation with a single neural field model. However, this approach faces challenges in capturing dynamic and real-time scenes, as it relies on accurate and fixed prior scene information. This can be particularly problematic in extensive mapping scenarios. In response, NEWTON introduces a view-centric neural field-based mapping method designed to overcome these lim-
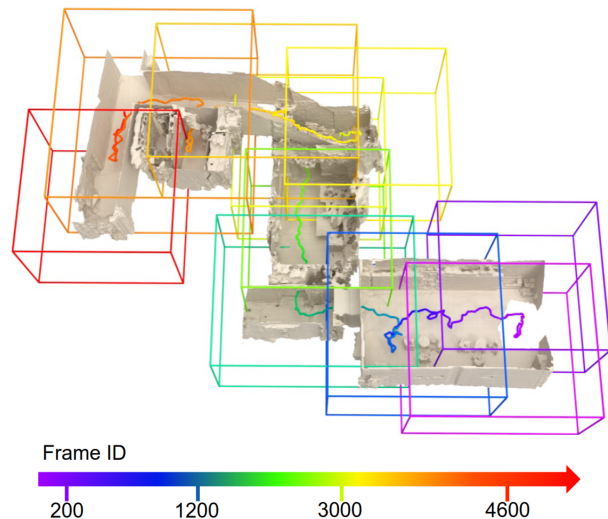


Fig. 8: **Submaps Visualization.** Neural submaps, allocated incrementally along the scanning trajectory, encode precise scene geometry and colors in their dedicated local coordinate frames. Figure from [126].

itations. Unlike existing methods, NEWTON dynamically constructs multiple neural field models, each represented as a multi-resolution feature grid [54] in a spherical coordinate system, based on real-time observations and allowing camera pose updates through loop closures and scene boundary adjustments. This is facilitated by the coordination with the camera tracking component of ORB-SLAM2 [99].

**Vox-Fusion++ [131].** Extending the original Vox-Fusion [88], this system combines sparse voxel embeddings with a neural implicit network. Specifically, it utilizes a dynamic octree structure to manage the voxel embeddings, allowing for efficient on-the-fly expansion of the map. The neural implicit network is trained to regress signed distance values. To handle large-scale scenes, Vox-Fusion++ employs a multi-map approach. Each map is constructed incrementally and optimized independently. Loop detection and hierarchical pose optimization are performed between different maps to reduce long-term drift and eliminate duplicate geometry. This multi-map strategy enables the reconstruction of large indoor scenes and collaborative mapping.

**MUTE-SLAM [132].** MUTE-SLAM utilizes a multi-map representation with tri-plane hash-encodings for efficient scene reconstruction. The key contributions of the method include the dynamic allocation of sub-maps for newly observed regions and the use of three orthogonal axis-aligned planes for hash-encoding scene properties in each sub-map, which reduces hash collisions and trainable parameters compared to grid-based representations. The TSDF and color features from the sub-maps are decoded using separate multilayer perceptrons and used in a volume rendering process to jointly optimize the sub-maps and camera poses. The optimization strategy concurrently optimizes all sub-maps intersecting with the current camera frustum, ensuring global consistency, while periodic global bundle adjustments are employed to further refine the poses and scene representation.

### 3.1.4 Semantic RGB-D SLAM

Operating as SLAM systems, these methodologies inherently include mapping and tracking processes while also incorporating semantic information to enhance the understanding of the environment. Tailored for tasks such as object recognition or semantic segmentation, these frameworks provide a holistic approach to scene analysis - identifying and classifying objects and/or efficiently categorizing image regions into specific semantic classes (*e.g.* tables, chairs, etc.).

**iLabel** [7]. This framework is a novel system for interactively understanding and segmenting 3D scenes. It uses a neural field representation to map 3D coordinates to color, volumetric density, and semantic values. Specifically, the core of this work is established on the basis of iMAP [1]. User interaction within the framework, instead, involves providing annotations through user-clicks on the scene. The system then employs these annotations to optimize its predictions of semantic labels, which essentially assigns meaningful labels to different parts of the scene. The framework supports two modes of interaction: a *manual mode*, where users provide semantic labels through clicks, and a *hands-free mode*, where the system automatically proposes informative positions for labeling based on semantic uncertainty, thereby reducing user effort. Moreover, iLabel is capable of achieving efficient interactive labeling without relying on pre-existing training data.

**FR-Fusion** [133]. This method seamlessly integrates a neural feature fusion system into the iMAP [1] framework. By incorporating a 2D image feature extractor (either EfficientNet [134] or DINO-based [135]) and augmenting iMAP with a latent volumetric rendering technique, the system efficiently fuses high-dimensional feature maps with low computational and memory requirements. Prioritizing a "feature-realistic" scene representation over "photo-realistic" models, the scene network operates incrementally and enables dynamic open-set semantic segmentation through sparse user interaction. The system's effectiveness is demonstrated on several tasks, including object grouping, object part category specialization, and unreconstructed region exploration, demonstrating its potential for practical applications. This approach is particularly promising in complex and unconventional domains where pre-trained semantic segmentation networks currently have limitations.

**vMap** [136]. This framework introduces a novel approach to object-level dense SLAM. In scenarios where 3D priors are not available, vMAP efficiently constructs a comprehensive scene model by representing each object with a dedicated MLP. This strategy facilitates the creation of watertight and complete object models, even when dealing with partially observed or occluded objects in real-time RGB-D input streams. The methodology revolves around using object masks for segmentation, efficient object association, and tracking via the off-the-shelf ORB-SLAM3 [111] algorithm. Through a seamless integration of these features, the system models the 3D scene in a flexible and efficient manner. This not only ensures precise reconstructions of objects, but also supports the process of re-composing scenes, independently tracking diverse objects, and continually updating objects of interest in real time.

**SNI-SLAM** [137]. SNI-SLAM employs a neural implicit representation and hierarchical semantic encoding
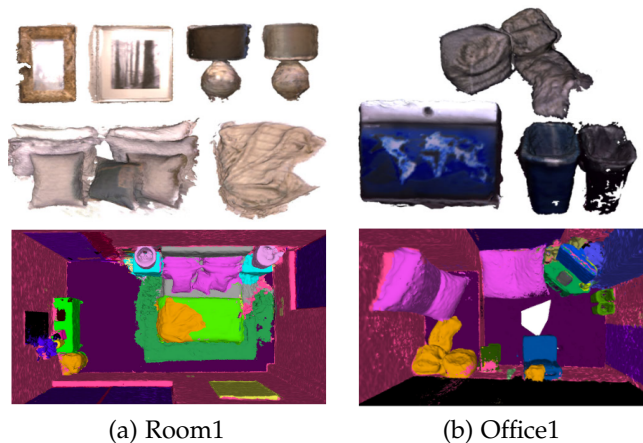


(a) Room1          (b) Office1

Fig. 9: **Semantic Visualization.** 3D semantic mesh (bottom) and its decomposition with RGB colors (top) for two scenes from the Replica [78] dataset. Images from [141].

for multi-level scene comprehension. Key features include cross-attention mechanisms for collaborative integration of appearance, geometry, and semantic features. Hierarchical semantic mapping, an integral component, is adopted as a coarse-to-fine optimization scheme, enabling detailed reconstruction while conserving memory resources. A novel decoder design ensures unidirectional interaction between appearance, geometry, and semantic features, preventing mutual interference and yielding improved results. Additionally, the method incorporates various loss functions (semantic, feature, color, depth) as guidance for the scene representation and network optimization.

**NIDS-SLAM** [139]. This approach enhances scene perception by dynamically learning 3D geometry and semantic segmentation online. The algorithm combines classical tracking and loop closure mechanisms based on ORB-SLAM 3 [111] with neural fields-based mapping, leveraging the strengths of existing methodologies. The core of the method lies in its mapping strategy. A mapping network, built upon the Instant-NGP [54] backbone with adaptations from Neus [71], learns the SDF of the environment. The network optimizes its structure through a combination of photometric, geometric, and semantic losses on selected pixels, with dynamic keyframe selection optimizing computational efficiency. A notable contribution is the novel approach to dense 3D semantic segmentation using 2D semantic color maps from keyframes, as outlined in [177]. The mapping network seamlessly integrates this strategy, enabling precise labeling in scenarios with inconsistent keyframe semantics.

**DNS SLAM** [141]. This work leverages the potential of 2D semantic priors, enabling stable camera tracking while concurrently training class-wise scene representations. The integration of semantic information with multi-view geometry results in a comprehensive and semantically decomposed geometric representation, crucial for improving accuracy and detailing in scene reconstructions. This allows for reconstructing a semantically annotated 3D mesh, as shown in Figure 9. Additionally, by employing image-based feature extraction and imposing multi-view geometry constraints, this framework ensures enhanced color, density, and semantic class information. This approach leads to

superior texture capture and finer geometric details in reconstructions, effectively addressing the over smoothed reconstruction problems seen in other methods. Furthermore, DNS SLAM introduces a novel real-time tracking strategy using a lightweight coarse scene representation, trained through self-supervision. This optimization enhances tracking efficiency while utilizing the multi-class representation as pseudo ground-truth.

**SGS-SLAM [142]**. SGS-SLAM employs the first semantic dense SLAM system using 3D Gaussians. To achieve this, the framework uses multi-channel optimization during the mapping process to integrate visual, geometric, and semantic constraints all at once. Accordingly, semantic information is embedded in the 3D Gaussians in the form of additional color channels that are optimized over 2D semantic segmentation maps corresponding to the color keyframes. During tracking, the camera pose for a new keyframe is estimated from the previous one by assuming a constant velocity camera motion model, and iteratively refined by minimizing the loss of color, depth, and semantics between the rendered view and the ground-truth image.

**SemGauss-SLAM [143]**. SemGauss-SLAM incorporates semantic feature embedding into the 3D Gaussian representation to enable dense semantic mapping. It propagates 2D semantic features directly to the 3D Gaussian initialization for efficient mapping optimization. A feature-level loss is introduced along with the semantic and RGB/depth losses to provide higher-level guidance. The feature-level loss is computed using features extracted from images by a Dinov2 model [138], which are compared to splatted features from the 3D Gaussian representation. To reduce drift, a semantic-informed bundle adjustment is performed for joint optimization of camera poses and the 3D Gaussian representation by exploiting semantic consistency constraints between co-visible frames as extracted by Dinov2.

**NEDS-SLAM [144]**. NEDS-SLAM is a semantic SLAM system based on 3DGS. It employs a Spatially Consistent Feature Fusion model, which combines semantic features from a pre-trained model with appearance features from Depth Anything [145] to reduce the impact of inconsistent semantic estimates. A lightweight encoder-decoder is used to compress high-dimensional semantic features into a compact 3D Gaussian representation, mitigating the burden of excessive memory consumption. NEDS-SLAM leverages the advantages of 3DGS, which enables efficient and differentiable novel view rendering, and proposes a Virtual Camera View Pruning method to eliminate outlier GS points.

### 3.1.5 SLAM in Dynamic Environments

Most of the SLAM methods reviewed so far are based on the fundamental assumption of a static environment characterized by rigid, non-moving objects. While these techniques perform promisingly in static scenes, their performance in dynamic environments faces significant challenges, limiting their applicability in real-world scenarios. Therefore, in this section, we provide an overview of methods that are specifically designed to address the challenges of accurate mapping and localization estimation in dynamic settings.

**DN-SLAM [13]**. This work integrates various components to address challenges in accurate location estimation and map consistency in dynamic environments. Leveraging ORB features for object tracking and employing semantic segmentation, optical flow, and the Segment Anything Model (SAM) [146], DN-SLAM effectively identifies and segregates dynamic objects within the scene while preserving static regions, enhancing SLAM performance. Specifically, the methodology involves utilizing semantic segmentation for object identification, refining dynamic object segmentation via SAM, extracting static features, and employing NeRF for dense map generation.

**DynaMoN [147]** . This framework builds upon DROID-SLAM [28], enhancing it with motion and semantic segmentation. The methodology integrates these elements into a dense BA process, utilizing motion and segmentation masks to weight the optimization process and ignore potentially dynamic pixels. Semantic segmentation, facilitated by a pretrained DeepLabV3 [148] network, aids in refining masks for known object classes and incorporates motion-based filtering for handling unknown dynamic elements. The study also introduces a 4D scene representation using NeRF, employing a combination of implicit and explicit representations for effective 3D reconstruction [178]. Optimization of NeRF involves mean squared error and Total Variation (TV) loss for regularization, enabling the generation of novel views in dynamic scenes.

**DDN-SLAM [9]**. This framework identifies key challenges within dynamic environments, such as dynamic objects, low-texture areas, and significant changes in lighting and viewpoints. To address these problems, the framework comprises semantic perception, sparse flow constraints, a background filling strategy, multi-resolution hash encoding, and tracking. In semantic systems, it detects the static and dynamic feature points using conditional probability fields. Meanwhile, the constraints are created for potential dynamic points and keyframes to improve the tracking performance. To alleviate the problems brought by dynamic objects in the scene, the framework adopts a skip voxel strategy based on selectively updating voxels to perform pixel control on selected keyframes and specific dynamic pixels. Furthermore, it employs semantic mask joint multi-resolution hash encoding to eliminate dynamic interferences and decrease the ghosting artifacts. When achieving the complete bundle adjustment and loop closure detection, DDN-SLAM [150] filters out the feature points validated by semantic and flow threads to gain robust performance.

**NID-SLAM [150]**. This method addresses the dynamic environment by employing a specialized dynamic processing procedure to remove dynamic objects, addressing inaccuracies in depth information. Depth-guided semantic mask enhancement is then introduced to reduce inconsistencies along edge regions and accurately detect dynamic objects, with subsequent background inpainting for occluded areas using static information from prior viewpoints. The keyframe selection strategy prioritizes frames with minimal presence of dynamic objects and limited overlap with previous keyframes for enhanced optimization efficiency. The scene representation involves multi-resolution geometric and color feature grids. The optimization process includes geometric and photometric losses, jointly optimizing scene representation features and camera extrinsic parameters for selected keyframes. In parallel, a tracking process optimizes

camera poses for the current frame. The paper suggests potential improvements to achieve real-time performance by addressing segmentation network speed and exploring the predictive capabilities of neural networks for comprehensive background inpainting.

**DVN-SLAM [151]**. This dynamic SLAM system uses a local-global fusion neural implicit representation, combining the advantages of continuous neural radiance fields for the global representation and discrete feature planes for the local representation. The system employs attention-based feature fusion, result fusion, and an information concentration loss to effectively model both local details and global structure while addressing uncertainties in the rendering process. The local-global fusion representation enables DVN-SLAM to be robust in dynamic scenes by automatically ignoring moving objects and recovering occluded backgrounds while maintaining stable scene structure modeling.

### 3.1.6 Uncertainty Estimation

Analyzing uncertainties in input data, especially depth sensor noise, is critical for robust system processing. This includes tasks such as filtering unreliable sensor measurements or incorporating depth uncertainty into the optimization process. The overall goal is to prevent inaccuracies within the SLAM process that could significantly impact system accuracy. At the same time, acknowledging the intrinsic uncertainty in the neural model reconstruction adds a critical layer for assessing system reliability, especially in challenging scenarios. This section marks the beginning of uncertainty exploration in neural SLAM, emphasizing the integration of both epistemic (knowledge-based) and aleatoric (environmental-noise-based) uncertainty information as essential components to improve overall SLAM system performance.

**OpenWorld-SLAM [152]**. This work improves upon NICE-SLAM [5], addressing its non-real-time execution, limited trajectory estimates, and challenges with adapting to new scenes due to its dependency on a predefined grid. To enhance applicability in open-world scenarios, this work introduces novel improvements, including depth uncertainty integration from RGB-D images for local accuracy refinement, motion information utilization from an inertial measurement unit (IMU), and a division of NeRF a finite foreground grid and a background spherical grid for diverse environment handling. These enhancements result in improved tracking precision and map representation while maintaining NeRF-based SLAM advantages. This work emphasizes the need for specialized datasets supporting NeRF-based SLAM, especially those providing outdoor mesh models, motion data, and well-characterized sensors.

**UncLe-SLAM [17]**. UncLe-SLAM jointly learns the scene geometry and aleatoric depth uncertainties on the fly. This is achieved by employing the Laplacian error distribution associated with the input depth sensor. Unlike existing methods, which lack the integration of depth uncertainty modeling, UncLe-SLAM employs a learning paradigm to adaptively assign weights to different image regions based on their estimated confidence levels, obtained without requiring ground truth depth or 3D. This strategic weighting mechanism allows UncLe-SLAM to prioritize more reliable
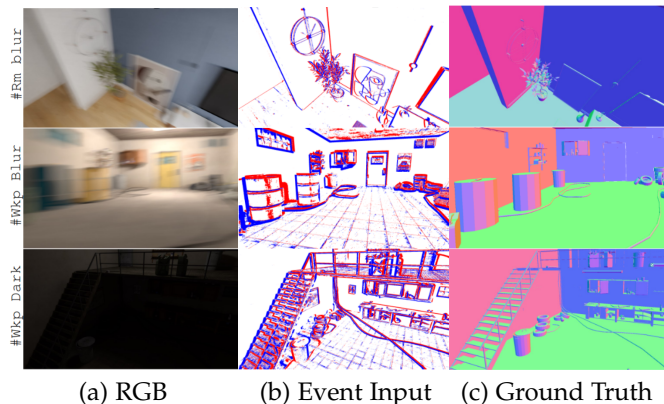


(a) RGB      (b) Event Input     (c) Ground Truth

Fig. 10: **Overview of the DEV-Indoors Dataset [14].** (a) RGB images depicting normal, motion blur, and dark scenes with corresponding (b) event streams and (c) ground truth meshes. Images from [14].

sensor information, leading to refined tracking and mapping results. Furthermore, UncLe-SLAM's adaptability extends to scenarios involving either RGB-D configurations or multiple depth sensors. This enables the system to handle diverse sensor setups, where each sensor might exhibit distinct noise characteristics.

### 3.1.7 Event-based SLAM

While radiance field-inspired VSLAM methods offer advantages in accurate dense reconstruction, practical scenarios involving motion blur and lighting variations pose significant challenges that affect the robustness of the mapping and tracking processes. In this section, we explore a category of systems that make use of data captured by event cameras, to exploit its dynamic range and temporal resolution. The asynchronous event generation mechanism, triggered by a logarithmic change in luminance at a given pixel, shows potential advantages in terms of low latency and high temporal resolution. This has the potential to improve the robustness, efficiency, and accuracy of neural VSLAM in extreme environments. Although event camera-based SLAM systems are still in the early stages of investigation, we believe that ongoing research holds great promise for overcoming the limitations of traditional RGB-based approaches.

**EN-SLAM [14]**. This framework introduces a new paradigm shift by seamlessly integrating event data alongside RGB-D through the implicit neural paradigm. It aims to overcome challenges encountered by existing SLAM methods when operating in non-ideal environments characterized by issues such as motion blur and lighting variation. The methodology centers around a differentiable Camera Response Function (CRF) rendering technique, enabling the mapping of a unified representation from event and RGB cameras. This process entails decoding the scene encoding, establishing a unified representation for geometry and radiance, and decomposing the shared radiance field into color and luminance via differentiable CRF Mappers. Additionally, optimization strategies are implemented for tracking and BA. The paper further proposes the creation of two challenging datasets—DEV-Indoors and DEV-Reals—including scenarios with practical motion blur and lighting changes,

as shown in Fig. 10, to evaluate EN-SLAM's effectiveness and robustness in diverse environments.

## 3.2 RGB-based SLAM Methodologies

This section explores RGB dense SLAM methods, which rely solely on visual cues from color images, eliminating the need for depth sensors, typically light-sensitive, noisy, and, in most cases, only applicable indoors. Hence, RGB-only SLAM using monocular or stereo cameras is gaining attention for scenarios where RGB-D cameras are impractical or costly, making RGB cameras a more viable solution applicable to a broader range of indoor and outdoor environments. However, these methods often face challenges, particularly in monocular setups, as they lack geometric priors, leading to depth ambiguity issues. As a result, they tend to exhibit slower optimization convergence due to less constrained optimization.

Similar to the RGB-D case, we group these methods into categories. We begin by examining the main NeRF-style (3.2.1) SLAM techniques. Subsequently, we cover related RGB SLAM systems that utilize external frameworks for additional supervision signals during optimization (3.2.2), those specifically designed for semantic estimation (3.2.3), and finally, those addressing system uncertainty (3.2.4).

### 3.2.1 NeRF-style RGB SLAM

**DIM-SLAM** [16]. This paper introduces the first RGB SLAM system using a neural implicit map representation. Similar to NICE-SLAM, it combines a learnable multi-resolution volume encoding and an MLP decoder for depth and color prediction. The system dynamically learns scene features and decoders on-the-fly. Moreover, DIM-SLAM optimizes occupancies in a single step by fusing features across scales, improving optimization speed. Notably, it introduces a photometric warping loss inspired by multi-view stereo, enforcing alignment between synthesized and observed images to enhance accuracy by addressing view-dependent intensity changes. Similar to other RGB-D approaches, DIM-SLAM leverages parallel tracking and mapping threads for the concurrent optimization of camera poses and implicit scene representation.

**Orbeez-SLAM** [153]. This approach seamlessly combines dense monocular SLAM strengths with neural radiance field modeling. In accordance with established methodologies, the system operates within a structured process that encompasses both tracking and mapping phases. The tracking process entails the extraction of image features and the estimation of camera poses using visual odometry, specifically derived from ORB-SLAM2 [99]. Conversely, the mapping stage, leveraging the capabilities of Instant-NGP [54], is dedicated to the generation of map points via triangulation and the execution of BA for the optimization of camera poses and 3D points alike. This synergy not only ensures real-time efficiency without the requirement for pre-training but also guarantees precise pose estimations, robust camera tracking, and accurate 3D reconstructions.

**FMapping** [154]. FMapping is a neural field mapping framework for real-time RGB SLAM, aiming to boost efficiency and reduce mapping uncertainty, especially in the absence of depth data. To address this, the authors conduct a thorough theoretical analysis, breaking down the SLAM system into tracking and mapping components, and explicitly defining mapping uncertainty within neural representations. Building on this analysis, the paper presents an innovative factorization scheme for the scene representation. Specifically, this approach employs a factorized neural field to effectively manage uncertainty by decomposing it into a lower-dimensional space, thereby increasing robustness to noise and training efficiency. A complementary sliding window strategy further reduces uncertainty during scene reconstruction by incorporating coherent geometric cues from observed frames. FMapping offers significant advantages, including low memory usage, streamlined computation, and rapid convergence during map initialization.

**TT-HO-SLAM** [155]. Motivated by the observation that existing methods fail to adhere to the binary-type opacity prior for rigid 3D scenes, the paper introduces a novel ternary-type (TT) opacity model for improved optimization during volumetric rendering. The authors empirically observe that frame-wise volumetric rendering and the absence of a binary-type opacity prior contribute to instability in RGB-only NeRF-SLAM. To tackle this, they propose a hybrid odometry scheme that combines volumetric and warping-based image renderings during tracking, leading to a substantial speed-up. Fine adjustments to camera odometry are made during the BA step, which occurs jointly with the mapping process. The methodology involves soft binarization of a decoder network during map initialization and utilizes theoretical insights to optimize opacity, demonstrating superior results in terms of both speed and accuracy.

**Hi-Map** [8] (formerly: FMapping [154]). This work presents a hierarchical factorized representation for monocular mapping. The key idea is to represent the scene as a hierarchical feature grid that encodes and factorizes the radiance into feature planes and vectors. It simplifies the scene's data structure with lower-dimensional elements and allows for fast convergence on changed views where the underlying geometry is unknown. To enhance photometric cues for distant and texture-less regions, Hi-Map employs a dual-path encoding strategy that incorporates absolute coordinates into the appearance encoding. It enables the framework to learn variations in color and lighting caused by changes in viewpoints. As a result, Hi-Map enhances geometric reconstruction and textural details, resulting in higher-quality mapping without external depth priors.

### 3.2.2 Aided Supervision

In this section, we explore RGB-based SLAM methods that use external frameworks to integrate regularization information into the optimization process, referred to as aided supervision. These frameworks include various techniques, such as supervision derived from depth estimates obtained from single or multi-view images, surface normal estimation, optical flow, and more. The incorporation of external signals is crucial for disambiguating the optimization process and helps to significantly improve the performance of SLAM systems using only RGB images as input.

**iMODE** [156]. The system operates through a multi-threaded architecture consisting of three core processes. First, a localisation process utilizes the ORB-SLAM2 [99] sparse SLAM system for real-time camera pose estimation

on a CPU, selecting keyframes for subsequent mapping. Second, inspired by iMAP [1], a semi-dense mapping process enhances reconstruction accuracy by supervising real-time training through depth-rendered geometry. Despite lacking a depth camera, monocular multi-view stereo methods provide depth measurements [179]. Third, the dense reconstruction process, executed on a GPU, optimizes an MLP representing the neural field. This latter differs from iMAP's implementation by incorporating both a view dependency aspect, which addresses photometric consistency by integrating view-dependent effects like specularities, and frequency separation, where a lower frequency embedding is employed for initial input, while a higher frequency embedding is reserved for the color head only. This optimization aligns keyframe images and semi-dense depth maps, minimizing photometric and geometric errors.

**Hi-SLAM [158]**. Three challenges are put forward in this work: low texture and rapid movement, scale ambiguity inherent in depth priors, and lack of global consistency. Hi-SLAM [158] builds on DROID-SLAM [28] to obtain dense pixel correspondences between nearby frames. This allows the methodology to track camera poses with optical flow and create a keyframe graph and a keyframe buffer. Monocular depth priors [159] are incorporated to further improve depth accuracy. Besides, the framework proposes a joint depth and scale adjustment (JDSA) module to achieve depth estimation and avoid scale ambiguity of monocular depth priors. In this module, the scales and offsets of the depth priors are estimated and then incorporated as variables in the BA optimization. As a result, the module maintains the scale consistency of depth priors. To maintain global consistency, Hi-SLAM [158] employs a Sim(3)-based pose graph bundle adjustment (PGBA) approach for online loop closure. When detecting the loop closures, the PGBA process constructs the pose graph using Sim(3) to update the scales. This enables the framework to optimize the pose graph with global consistency despite potential pose and scale drift.

**NICER-SLAM [18]**. This work introduces a unified end-to-end framework that concurrently optimizes tracking and mapping, given an RGB stream of images as input. At the core, hierarchical feature grids provide the structure for accurately representing the scene's geometry and appearance, offering a multi-level framework for modeling SDF and color. Additionally, the system's joint mapping and tracking capabilities are facilitated by a comprehensive suite of loss functions. These include the standard RGB rendering loss, the RGB warping loss for enforcing geometric consistency, and the optical flow loss that aids in addressing ambiguities and imposing smoothness priors. Complementing these are monocular depth and normal losses extracted from an off-the-shelf monocular depth predictor [159], along with the Eikonal loss [180] for regularization.

**NeRF-VO [161]**. This work follows a two-stage approach: firstly, employing a learning-based sparse visual tracking method, DPVO [162], to generate initial poses and sparse depth information. Then, enhancing these sparse cues involves a dense geometry module predicting dense depth maps using the state-of-the-art monocular depth network, DPT [121], [181], and extracting surface normals using Omnidata [159] from monocular RGB input. The alignment of these cues with sparse data is achieved through a scale

alignment procedure. Secondly, the system employs Nerfacto [182] for dense scene representation. It optimizes the representation with camera poses, RGB images, depth maps (with uncertainty-based loss), and surface normals. This joint optimization refines scene geometry and camera poses by minimizing disparities between captured images and rendered views from the neural representation.

**MoD-SLAM [163]**. This monocular framework follows a two-step process: first estimating depth and then refining it for precise scene reconstruction. The depth estimation comprises relative and metric depth modules utilizing DPT [121] and ZoeDepth [149] architectures. To enhance accuracy, depth estimates undergo refinement through a depth distillation module, addressing inaccuracies from the initial depth estimation. MoD-SLAM integrates a multivariate Gaussian encoding and a ray reparameterization technique, facilitating efficient representation of unbounded scenes by capturing detailed 3D space information. The system also employs loop closure to mitigate pose drift, ensuring more precise global pose optimization.

**Q-SLAM [164]**. Q-SLAM integrates quadric representations throughout the pipeline to improve 3D scene geometry modeling. The method uses a tracking module based on DROID-SLAM [92] to predict rough depth maps and camera poses, while a segmentation network estimates masks used in a quadric-based depth correction module to refine noisy depth maps. In the mapping phase, Q-SLAM transforms dense volumetric scenes into a manageable set of quadric surfaces. A quadric-ray transformer employs importance sampling and models interrelationships between points on quadric surfaces and across rays during rendering. Quadric semantics serve as a supervision signal for NeRF network optimization, and an end-to-end joint optimization strategy refines camera poses and scene representation parameters.

**GlORIE-SLAM [165]**. GlORIE-SLAM, instead, utilizes a deformable neural point cloud for mapping and a globally optimized frame-to-frame tracking approach based on optical flow, similarly to GO-SLAM [91]. The tracking module incorporates a monocular depth prior, estimated by an off-the-shelf network [159], and a noisy depth map from the tracker itself [28]. These are fused into a proxy depth map using multi-view consistency checks and depth completion. The Disparity, Scale and Pose Optimization (DSPO) layer is introduced to jointly optimize the pose, depth, and scale of the monocular depth within the bundle adjustment framework. The deformable neural point cloud efficiently adapts to the optimized keyframe poses and depth updates without requiring backpropagation. Loop closure detection and online global bundle adjustment are integrated to maintain global consistency.

### 3.2.3 Semantic RGB SLAM

**RO-MAP [166]**. RO-MAP is a real-time multi-object mapping system that operates without depth priors, utilizing neural radiance fields for object representation. This approach combines a lightweight object-centric SLAM with NeRF models for simultaneous localization and reconstruction of objects from monocular RGB input. The system efficiently trains separate NeRF models for each object, demonstrating real-time performance in semantic object mapping and shape reconstruction. Key contributions include the

development of the first 3D prior-free monocular multi-object mapping pipeline, an efficient loss function tailored for objects, and a high-performance CUDA implementation.

### 3.2.4 Uncertainty Estimation

**NeRF-SLAM** [167]. By employing real-time implementations of DROID-SLAM [28] as the tracking module and Instant-NGP [54] as the hierarchical volumetric neural radiance field map, this approach successfully achieves real-time operational efficiency given RGB images as input. Moreover, incorporating depth uncertainty estimation addresses inherent noise in depth maps, resulting in improved outcomes through depth loss supervision for neural radiance fields – with weights determined by the depth's marginal covariance. Specifically, the pipeline involves two real-time synchronized threads: tracking and mapping. The tracking thread minimizes BA re-projection errors for a sliding keyframe window. The mapping thread optimizes all keyframes from the tracking thread without a sliding window. Communication only occurs when the tracking thread creates a new keyframe, sharing keyframe data, poses, depth estimates, and covariances.

### 3.3 LiDAR-Based SLAM Strategies

While VSLAM systems discussed so far operate successfully in smaller indoor scenarios where both RGB and dense depth data are available, their limitations become apparent in large outdoor environments where RGB-D cameras are impractical. LiDAR sensors, which provide sparse yet accurate depth information over long distances and in a variety of outdoor conditions, play a critical role in ensuring robust mapping and localization in these settings. However, the sparsity of LiDAR data and the lack of RGB information pose challenges for the application of the previously outlined dense SLAM approaches in outdoor environments. Our focus is now on novel methodologies that exploit the precision of 3D incremental LiDAR data to improve autonomous navigation in outdoor scenarios, while taking advantage of scene representations based on radiance fields, offering the potential to achieve dense, smooth map reconstruction of the environment, even in areas with sparse sensor coverage. Given the limited number of studies addressing this specific setting, we categorize the methodologies into two simple groups: NeRF (3.3.1) and 3DGS-style (3.3.2) LiDAR-based SLAM categories.

### 3.3.1 NeRF-style LiDAR-based SLAM

**NeRF-LOAM** [15]. NeRF-LOAM introduces the first neural implicit approach to jointly determine sensor position and orientation while constructing a comprehensive 3D representation of large-scale environments using LiDAR data. The framework comprises three interconnected modules: neural odometry, neural mapping, and mesh reconstruction. The neural odometry module estimates a 6-DoF pose for each incoming LiDAR scan by minimizing SDF errors through a fixed implicit network. The poses are subsequently optimized via back-projection. In parallel, the neural mapping module employs dynamic voxel embeddings within an octree-based architecture, adeptly capturing local
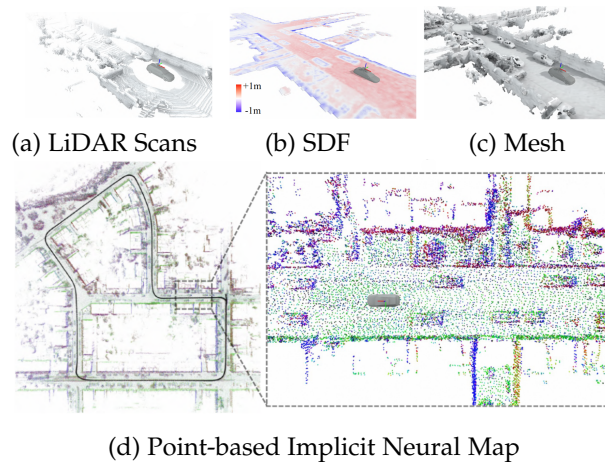


(a) LiDAR Scans     (b) SDF     (c) Mesh

(d) Point-based Implicit Neural Map

Fig. 11: **Overview of PIN-SLAM [170].** Top: (a) LiDAR scans, (b) Implicit SDF , (c) Reconstructed mesh from the SDF. Bottom: (d) Visualization of the Point-based Implicit Neural (PIN) Map. Images from [170].

geometry. This dynamic allocation strategy ensures efficient utilization of computational resources, avoiding the complexities of pre-allocated embeddings or time-intensive hash table searches. The method uses a dynamic voxel embedding look-up table, boosting efficiency and eliminating computational bottlenecks. A key-scans refinement strategy enhances reconstruction quality and addresses catastrophic forgetting during incremental mapping, leading to detailed 3D mesh representations in the final step.

**LONER** [168]. This system employs parallel tracking and mapping threads, with the tracking thread processing incoming scans using ICP for odometry estimation. The mapping thread utilizes selected keyframes to update the neural scene representation. LiDAR scans are transformed using Point-to-Plane ICP [169], and the scene is represented by an MLP with a hierarchical feature grid encoding. The proposed dynamic margin loss function combines Jensen-Shannon Divergence [183], depth loss, and sky loss. The dynamic margin adapts to diverse map regions during online training, enabling the system to learn new areas while preserving acquired geometry. The system incorporates meshing for offline visualization, creating a mesh from the implicit geometry using estimated keyframe poses.

**PIN-SLAM** [170]. PIN-SLAM is a SLAM system designed for LiDAR scans, featuring a point-based implicit neural map representation for building globally consistent maps, depicted in Figure 11. The system leverages sparse, optimizable neural points that exhibit elasticity, allowing for continuous deformation during global pose adjustments. Employing an alternating approach, it performs incremental learning of local implicit signed distance fields and pose estimation using correspondence-free point-to-implicit model registration. The methodology includes efficient odometry estimation, dynamic point filtering, and loop closure detection based on local polar context descriptors. The system corrects drift through optimized neural point maps after loop closure, achieving global consistency. Prominent claims include supporting large-scale mapping, enabling real-time execution thanks to voxel hashing and efficient neural point
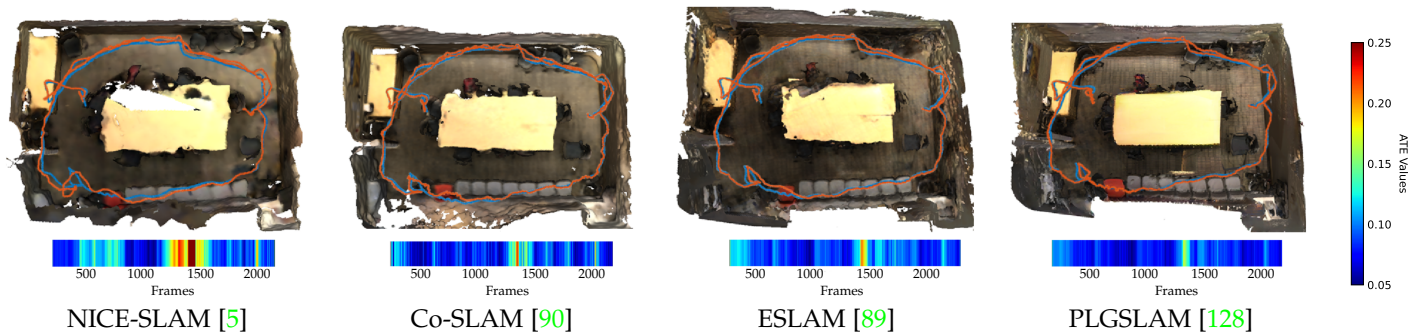
Fig. 12: **SLAM Methods Comparison on the ScanNet [77] Dataset – Surface Reconstruction and Localization Accuracy.** Ground truth trajectory in blue, estimated trajectory in orange. ATE visualized with a color bar.

indexing, and providing a compact map representation suitable for accurate mesh reconstruction.

### 3.3.2 3DGS-style LiDAR-based SLAM

**LIV-GaussMap [171].** The proposed LiDAR-Inertial-Visual (LIV) fused radiance field mapping system integrates LiDAR-inertial sensors with a camera for precise data alignment. The methodology begins with LiDAR-inertial odometry, utilizing size-adaptive voxels to represent planar surfaces. LiDAR point clouds are segmented into voxels, and covariance matrices are computed for initial elliptical splatting estimates. The system is refined by optimizing spherical harmonic coefficients and LiDAR Gaussian structures using visual-derived photometric gradients, enhancing mapping precision and visual realism. The initialization of Gaussians involves size-adaptive voxel partitioning, with further subdivision based on a specified parameter. Adaptive control of the 3D Gaussian map addresses under-reconstruction and over-dense scenarios through structure refinement and photometric gradient optimization.

**MM-Gaussian [172].** MM-Gaussian is a LiDAR-camera fusion system that uses 3DGS for localization and mapping. Specifically, it estimates the LiDAR's pose using point cloud registration and derives the camera's pose, which is further optimized by comparing rendered and captured images. During mapping, LiDAR point clouds are converted into 3D Gaussian points and added to the map, with their attributes updated using a keyframe sequence. A relocalization module detects tracking failures and resets the pose to the correct trajectory, enhancing robustness.

## 4 EXPERIMENTS AND ANALYSIS

In this section, we compare methods across datasets, focusing on tracking (visual in 4.1.1, LiDAR in 4.2.1) and 3D reconstruction (visual in 4.1.2, LiDAR in 4.2.2). Additionally, we explore novel view synthesis (4.1.3), semantic segmentation (4.1.4), and analyze performance in terms of runtime and memory usage (4.3). In each subsequent table, we emphasize the best results withing a subcategory using **bold** and highlight the absolute best in purple. In our analysis, we organized quantitative data from papers with a common evaluation protocol and cross-verified the results. Our priority was to include papers with consistent benchmarks, ensuring a reliable basis for comparison across multiple sources. Although not exhaustive, this approach

| Method | Tracker Based on | Global BA | Loop Closure | fr1/desk | fr2/xyz | fr3/office | Avg (↓) |
|---|---|---|---|---|---|---|---|
| | | | RGB-D | | | | |
| Kintinuous [2] | | | | 3.7 | 2.9 | 3.0 | 3.2 |
| BAD-SLAM [30] | | ✓ | ✓ | 1.7 | 1.1 | 1.7 | 1.5 |
| ORB-SLAM2 [99] | | ✓ | ✓ | 1.6 | 0.4 | 1.0 | 1.0 |
| Vox-Fusion [88] | | | | 3.5 | 1.5 | 26.0 | 10.3 |
| MeSLAM [6] | | | | 6.0 | 6.5 | 7.8 | 6.8 |
| iMAP [1] | | | | 4.9 | 2.0 | 5.8 | 4.2 |
| GS-SLAM [184] | | | | 3.3 | 1.3 | 6.6 | 3.7 |
| SplaTAM [112] | | | | 3.4 | 1.2 | 5.2 | 3.3 |
| HF-GS SLAM [117] | | | | 3.4 | - | 5.1 | - |
| Compact-GSSLAM [114] | | ✓ | | 3.0 | 1.0 | 4.9 | 3.0 |
| MIPS-Fusion [126] | | | ✓ | 3.0 | 1.4 | 4.6 | 3.0 |
| Point-SLAM [93] | | | | 4.3 | 1.3 | 3.5 | 3.0 |
| Gaussian-SLAM [113] | | | | 2.6 | 1.3 | 4.6 | 2.9 |
| Loopy-SLAM [129] | | | ✓ | 3.8 | 1.6 | 3.4 | 2.9 |
| NICE-SLAM [5] | | | | 2.7 | 1.8 | 3.0 | 2.5 |
| GS-ICP SLAM [115] | G-ICP [116] | | | 2.7 | 1.8 | 2.7 | 2.4 |
| vMAP [136] | ORB3 [111] | | | 2.6 | 1.6 | 3.0 | 2.4 |
| Co-SLAM [90] | | ✓ | | 2.4 | 1.7 | 2.4 | 2.2 |
| ESLAM [89] | | | | 2.5 | 1.1 | 2.4 | 2.0 |
| CG-SLAM [118] | | | | 2.4 | 1.2 | 2.5 | 2.0 |
| NeSLAM [107] | | | | 1.8 | 1.0 | 2.1 | 1.6 |
| GSSLAM [109] | | | | 1.5 | 1.6 | 1.7 | 1.6 |
| GO-SLAM [91] | DROID [92] | ✓ | ✓ | 1.5 | 0.6 | 1.3 | 1.1 |
| Q-SLAM [164] | DROID [92] | | | 1.4 | 0.5 | 1.1 | 1.0 |
| NGEL-SLAM [127] | ORB3 [111] | ✓ | ✓ | 1.5 | 0.5 | 1.0 | 1.0 |
| DDN-SLAM [9] | ORB3 [111] | ✓ | ✓ | 1.5 | 0.4 | 0.9 | 0.9 |
| | | | RGB | | | | |
| DROID-SLAM [92] | - | ✓ | | 1.8 | 0.5 | 2.8 | 1.7 |
| ORB-SLAM2 [99] | - | ✓ | ✓ | 1.9 | 0.6 | 2.4 | 1.6 |
| GSSLAM [109] | | | | 4.2 | 4.8 | 4.4 | 4.4 |
| DDN-SLAM [9] | ORB3 [111] | ✓ | ✓ | 1.9 | 2.4 | 2.9 | 2.4 |
| DIM-SLAM [16] | | | | 2.0 | 0.6 | 2.3 | 1.6 |
| GO-SLAM [91] | DROID [92] | ✓ | ✓ | 1.6 | 0.6 | 1.5 | 1.2 |
| Orbeez-SLAM [153] | ORB2 [99] | | | 1.9 | 0.3 | 1.0 | 1.1 |
| MoD-SLAM [163] | DROID [92] | | ✓ | 1.5 | 0.7 | 1.1 | 1.1 |
| GIORIE-SLAM [165] | DROID [92] | ✓ | ✓ | 1.6 | 0.2 | 1.4 | 1.1 |

TABLE 2: **TUM RGB-D [76] Camera Tracking Results**. ATE RMSE [cm] (↓) is used as the evaluation metric.

guarantees the inclusion of methods with verifiable results and a shared evaluation framework in our tables. For performance analysis, we utilized methods with available code to report runtime and memory requirements on a common hardware platform, a single NVIDIA 3090 GPU. For specific implementation details of each method, readers are encouraged to refer to the original papers.

### 4.1 Visual SLAM Evaluation

In line with existing protocols, this section compares SLAM systems using RGB-D or RGB data. We evaluate tracking, 3D reconstruction, rendering, and consider runtime and memory usage. Additionally, for methods that estimate semantic segmentation, we assess the quality of the semantic segmentation using the mIoU metric. Specifically, results are presented on the TUM-RGB-D [76], Replica [78], and Scan-Net [77] datasets. For semantic segmentation evaluation, we focus on the Replica dataset, as it provides ground truth semantic labels, allowing for a comprehensive comparison of the semantic segmentation performance.
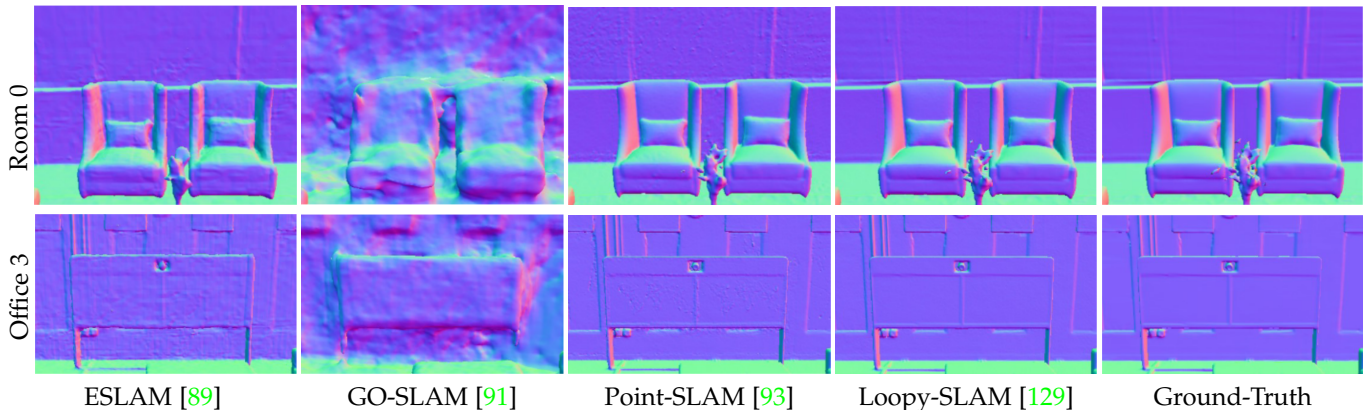
**Fig. 13: SLAM Methods Comparison on the Replica [78] Dataset – Mapping.** Images sourced from [129].

| Method | Tracker Based on | Global BA | Loop Closure | 0000 | 0059 | 0106 | 0169 | 0181 | 0207 | Avg (↓) |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | RGB-D | | | | | |
| DROID-SLAM (VO) [92] | | | | 8.00 | 11.30 | 9.97 | 8.64 | 7.38 | - | - |
| DROID-SLAM [92] | | ✓ | | **5.36** | **7.72** | **7.06** | **8.01** | **6.97** | - | - |
| iMAP [1] | | | | 55.95 | 32.06 | 17.50 | 70.51 | 32.10 | 11.91 | 36.67 |
| iDF-SLAM [104] | | | | 57.7- | 8.2- | 5.8- | 39.9- | 29.1- | 15.4- | 26.01- |
| ADFP [95] | | | | - | 10.50 | 7.48 | 9.31 | - | 5.67 | - |
| Gaussian-SLAM [113] | | | | 24.75 | 8.63 | 11.27 | 14.59 | 18.70 | 14.36 | 15.38 |
| Point-SLAM [93] | | | | 10.24 | 7.81 | 8.65 | 22.16 | 14.77 | 9.54 | 12.19 |
| SplaTAM [112] | | | | 12.83 | 10.10 | 17.72 | 12.08 | 11.10 | 7.46 | 11.88 |
| Compact-GSSLAM [114] | | ✓ | | 11.28 | 9.24 | 16.49 | 11.09 | 10.88 | 6.58 | 10.84 |
| MIPS-Fusion [126] | | | ✓ | 7.9- | 10.7- | 9.7- | 9.7- | 14.2- | 7.8- | 10.0- |
| Vox-Fusion [88] | | | | 8.39 | - | 7.44 | 6.53 | 12.20 | 5.57 | - |
| NEDS-SLAM [144] | | | | 12.34 | - | - | 11.21 | 10.35 | 6.56 | - |
| SemGauss-SLAM [143] | | | | 12.56 | 7.97 | - | 9.05 | 9.78 | 8.97 | - |
| NeuV-SLAM [106] | | | | 12.71 | 9.70 | 8.50 | 8.92 | 12.72 | 5.61 | 9.68 |
| NICE-SLAM [5] | | | | 8.64 | 12.25 | 8.09 | 10.28 | 12.93 | 5.59 | 9.63 |
| Co-SLAM [90] | | ✓ | | 7.18 | 12.29 | 9.57 | 6.62 | 13.43 | 7.13 | 9.37 |
| CG-SLAM [118] | | | | 7.09 | 7.46 | 8.88 | 8.16 | 11.60 | 5.34 | 8.08 |
| MUTE-SLAM [132] | | ✓ | | 7.08 | 9.07 | 8.27 | 6.18 | 10.21 | 7.19 | 8.00 |
| Loopy-SLAM [129] | | | ✓ | **4.2-** | 7.5- | 8.3- | 7.5- | 10.6- | 7.9- | 7.7- |
| ESLAM [89] | | | | 7.3- | 8.5- | 7.5- | 6.5- | 9.0- | 5.7- | 7.4- |
| Structerf-SLAM [98] | ORB2 [99] | | | 7.28 | 6.07 | 8.50 | 7.35 | - | | 7.28 |
| Vox-Fusion++ [131] | | | ✓ | 6.38 | 7.28 | 6.75 | 5.86 | 13.68 | 4.73 | 7.44 |
| NGEL-SLAM [127] | ORB3 [111] | ✓ | ✓ | 7.23 | **6.98** | 7.95 | 6.12 | 10.14 | 6.27 | 7.44 |
| DNS SLAM [141] | | ✓ | | 5.42 | **5.20** | 9.11 | 7.70 | 10.12 | 4.91 | 7.07 |
| NeSLAM [107] | | | | 6.87 | 7.37 | **5.23** | 9.07 | 9.27 | **4.08** | 6.98 |
| SNI-SLAM [137] | | | | 6.90 | 7.38 | 7.19 | **4.70** | - | - | - |
| GO-SLAM [91] | DROID [92] | ✓ | ✓ | 5.35 | 7.52 | 7.03 | 7.74 | 6.84 | 4.78 | 6.54 |
| Q-SLAM [164] | DROID [92] | | | 5.23 | 7.63 | 7.02 | 7.66 | 6.52 | - | - |
| MoD-SLAM [163] | DROID [92] | | ✓ | 5.27 | 7.44 | 6.73 | 6.48 | **6.14** | 5.31 | **6.23** |
| | | | | | RGB | | | | | |
| DROID-SLAM (VO) [92] | - | | | 11.05 | 67.26 | 11.20 | 16.21 | 9.94 | - | - |
| DROID-SLAM [92] | - | ✓ | | **5.48** | **9.00** | **6.76** | **7.86** | **7.41** | - | - |
| Orbeez-SLAM [153] | ORB2 [99] | | | 7.22 | **7.15** | 8.05 | **6.58** | 15.77 | 7.16 | 8.66 |
| GlORIE-SLAM [165] | DROID [92] | ✓ | ✓ | **5.5-** | 9.1- | 7.0- | 8.2- | 8.3- | 7.5- | 7.6- |
| Hi-SLAM [158] | DROID [92] | ✓ | ✓ | 6.40 | 7.20 | **6.50** | 8.50 | 7.60 | 8.40 | 7.40 |
| GO-SLAM [91] | DROID [92] | ✓ | ✓ | 5.94 | 8.27 | 8.07 | 8.42 | 8.29 | **5.31** | 7.38 |
| Q-SLAM [164] | DROID [92] | | | 5.77 | 8.46 | 8.38 | 8.74 | 8.76 | - | - |
| MoD-SLAM [163] | DROID [92] | | ✓ | **5.39** | 7.78 | 7.64 | 6.79 | **6.58** | 5.63 | **6.64** |

TABLE 3: **ScanNet [77] Camera Tracking Results**. ATE RMSE [cm] (↓) is used as the evaluation metric.

### 4.1.1 Tracking

**TUM-RGB-D.** Table 2 provides a thorough analysis of camera tracking results on three scenes of the TUM RGB-D dataset, marked by challenging conditions such as sparse depth sensor information and high motion blur in RGB images. Key benchmarks include established methods like Kintinuous, BAD-SLAM, and ORB-SLAM2, representing traditional hand-crafted baselines.

In the RGB-D setting, it is evident that most methods based on recent radiance field representations generally exhibit lower performance compared to reference methods like BAD-SLAM and ORB-SLAM2. One notable observation, however, is that decoupled methods using external trackers such as ORB3 and DROID, along with advanced strategies such as Global BA and LC, emerge as top performers. Specifically, NGEL-SLAM, DNN-SLAM, Q-SLAM and GO-SLAM demonstrate superior accuracy in RGB-D scenarios, with DDN-SLAM achieving the best average ATE RSME results of 0.9.

When shifting the focus to the RGB scenario, ORB-SLAM2 and DROID-SLAM serve as baselines, with ORB-SLAM2 exhibiting superior tracking accuracy. Among recent SLAM methods, GSSLAM and DDN-SLAM exhibit high ATE RMSE values of 4.4 and 2.4, respectively. This is in contrast to the RGB-D case where DDN-SLAM achieves the best results, indicating the importance of depth information for this methodology to perform well and its greater sensitivity in the RGB-only scenario. Despite this, Orbeez-SLAM, MoD-SLAM and GlORIE-SLAM, jointly with external tracking components, such as ORB-SLAM2 and DROID-SLAM, leads with an ATE RMSE of 1.1.

These results emphasize the varied performance of SLAM frameworks, with approaches based on the latest radiance field representations exhibiting effective results in RGB-D scenarios by separating mapping and tracking processes through external tracking approaches and additional optimization strategies. However, in scenarios where these latter are not applied, most methods still struggle with trajectory drift and sensitivity to noise.

**ScanNet.** Table 3 presents the evaluation of camera tracking methods on six scenes of the ScanNet dataset. In the RGB-D domain, standout performers are the frame-to-frame models MoD-SLAM and GO-SLAM. Both leverage well-crafted visual odometries (such as DROID-SLAM) and LC strategies, with GO-SLAM incorporating also Global BA. Significantly, MoD-SLAM achieves the best average ATE RMSE result of 6.23. A similar trend can be observed in the RGB case, where once again, the best results are achieved by methods employing external trackers. Nevertheless, it is worth noting that these solutions manage to be comparable or even superior to many other SLAM methods that leverage depth information from RGB-D sensors. In Figure 12, we report some qualitative results from selected RGB-D SLAM systems on ScanNet, highlighting recent improvements in trajectory error compared to the seminal systems.

**Replica.** Table 4 evaluates camera tracking across eight scenes from Replica, using higher-quality images compared to challenging counterparts like ScanNet and TUM RGB-D. The evaluation includes the reporting of ATE RMSE results for each individual scene, alongside the averaged outcomes.

On top, we report the evaluation concerning RGB-D methods. In line with observations from TUM RGB-D and ScanNet datasets, the highest accuracy is achieved by leveraging external tracking and methodologies involving Global BA and/or LC. In particular, GO-SLAM, Loopy-

Fig. 14: **SLAM Methods Comparison on the Replica [78] Dataset– Image Rendering.** Images sourced from [12].

| Method | Tracker Based on | Global BA | Loop Closure | R0 | R1 | R2 | O0 | O1 | O2 | O3 | O4 | Avg (↓) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **RGB-D** | | | | | | | | | | | | |
| iMAP [1] | | | | 3.12 | 2.54 | 2.31 | 1.69 | 1.03 | 3.99 | 4.05 | 1.93 | 2.58 |
| NICE-SLAM [5] | | | | 1.69 | 2.04 | 1.55 | 0.99 | 0.90 | 1.39 | 3.97 | 3.08 | 1.95 |
| ADFP [95] | | | | 1.39 | 1.55 | 2.60 | 1.09 | 1.23 | 1.61 | 3.61 | 1.42 | 1.81 |
| MIPS-Fusion [126] | | | ✓ | 1.10 | 1.20 | 1.10 | 0.70 | 0.80 | 1.30 | 2.20 | 1.10 | 1.19 |
| Co-SLAM [90] | | ✓ | | 0.65 | 1.13 | 1.43 | 0.55 | 0.50 | 0.46 | 1.40 | 0.77 | 0.86 |
| NIDS-SLAM [139] | ORB3 [111] | | ✓ | 0.58 | 0.41 | 0.58 | 0.62 | 0.40 | 1.20 | 0.88 | 1.80 | 0.80 |
| ESLAM [89] | | | | 0.71 | 0.70 | 0.52 | 0.57 | 0.55 | 0.58 | 0.72 | 0.63 | 0.63 |
| GSSLAM [109] | | | | 0.76 | 0.37 | 0.23 | 0.66 | 0.72 | 0.30 | 0.19 | 1.46 | 0.58 |
| Vox-Fusion [88] | | | | 0.40 | 0.54 | 0.54 | 0.50 | 0.46 | 0.75 | 0.50 | 0.60 | 0.54 |
| Point-SLAM [93] | | | | 0.61 | 0.41 | 0.37 | 0.38 | 0.48 | 0.54 | 0.72 | 0.63 | 0.52 |
| GS-SLAM [12] | | | | 0.48 | 0.53 | 0.33 | 0.52 | 0.41 | 0.59 | 0.46 | 0.70 | 0.50 |
| Vox-Fusion++ [131] | | | ✓ | 0.38 | 0.47 | 0.49 | 0.44 | 0.42 | 0.62 | 0.41 | 0.59 | 0.48 |
| SNI-SLAM [137] | | | | 0.50 | 0.55 | 0.45 | 0.35 | 0.41 | 0.33 | 0.62 | 0.50 | 0.46 |
| DNS SLAM [141] | | ✓ | | 0.49 | 0.46 | 0.38 | 0.34 | 0.35 | 0.39 | 0.62 | 0.60 | 0.45 |
| SplaTAM [112] | | | | 0.31 | 0.40 | 0.29 | 0.47 | 0.27 | 0.29 | 0.32 | 0.55 | 0.36 |
| NEDS-SLAM [144] | | | | - | - | - | - | - | - | - | - | 0.35 |
| GO-SLAM [91] | DROID [92] | ✓ | ✓ | 0.32 | 0.30 | 0.39 | 0.39 | 0.46 | 0.34 | 0.29 | 0.29 | 0.34 |
| SemGauss-SLAM [143] | | | | 0.26 | 0.42 | 0.27 | 0.34 | 0.17 | 0.32 | 0.36 | 0.49 | 0.33 |
| Compact-GSSLAM [143] | | ✓ | | 0.27 | 0.38 | 0.25 | 0.44 | 0.27 | 0.27 | 0.29 | 0.50 | 0.33 |
| MoD-SLAM [163] | DROID [92] | ✓ | | - | - | - | - | - | - | - | - | 0.33 |
| Gaussian-SLAM [113] | | | | 0.29 | 0.29 | 0.22 | 0.37 | 0.23 | 0.41 | 0.30 | 0.35 | 0.31 |
| Loopy-SLAM [129] | | | ✓ | 0.24 | 0.24 | 0.28 | 0.26 | 0.40 | 0.29 | 0.22 | 0.35 | 0.29 |
| CG-SLAM [118] | | | | 0.29 | 0.27 | 0.25 | 0.33 | 0.14 | 0.28 | 0.31 | 0.29 | 0.27 |
| HF-GS SLAM [117] | | | | 0.19 | 0.34 | 0.16 | 0.21 | 0.26 | 0.23 | 0.21 | 0.38 | 0.25 |
| GS-ICP SLAM [115] | G-ICP [116] | | | **0.15** | **0.16** | **0.11** | **0.18** | **0.12** | **0.17** | **0.16** | **0.21** | **0.16** |
| **RGB** | | | | | | | | | | | | |
| DROID-SLAM [92] | - | ✓ | | - | - | - | - | - | - | - | - | 0.42 |
| TT-HO-SLAM [155] | | | | 4.51 | 0.91 | 7.49 | 0.59 | 1.74 | 1.70 | 0.81 | 3.47 | 2.65 |
| NICER-SLAM [18] | | | | 1.36 | 1.60 | 1.14 | 2.12 | 3.23 | 2.12 | 1.42 | 2.01 | 1.88 |
| DIM-SLAM [16] | | | | 0.48 | 0.78 | 0.35 | 0.67 | 0.37 | **0.36** | 0.33 | 0.36 | 0.46 |
| GO-SLAM [91] | DROID [92] | ✓ | ✓ | - | - | - | - | - | - | - | - | 0.39 |
| MoD-SLAM [163] | DROID [92] | ✓ | ✓ | **0.28** | **0.29** | 0.30 | 0.40 | 0.45 | 0.50 | **0.31** | **0.27** | **0.35** |
| GlORIE-SLAM [165] | DROID [92] | ✓ | ✓ | 0.31 | 0.37 | **0.28** | **0.29** | **0.28** | 0.45 | 0.45 | 0.44 | **0.35** |

TABLE 4: **Replica [78] Camera Tracking Results**. ATE RMSE [cm] (↓) is used as the evaluation metric.

| Method | L1-Depth ↓ | Acc. [cm] ↓ | Comp. [cm] ↓ | Comp. Ratio [%] ↑ |
|---|---|---|---|---|
| **RGB-D** | | | | |
| COLMAP [59] | - | 8.69 | 12.12 | 67.62 |
| TSDF [185] | 7.57 | **1.60** | **3.49** | **86.08** |
| iMAP [1] | 7.64 | 6.95 | 5.33 | 66.60 |
| NICE-SLAM [5] | 3.53 | 2.85 | 3.00 | 89.33 |
| GO-SLAM* [91] | 3.38 | 2.50 | 3.74 | 88.09 |
| GO-SLAM [91] | 4.68 | 2.50 | 3.74 | 88.09 |
| MoD-SLAM [163] | 3.11 | 2.13 | - | - |
| DNS SLAM [141] | 3.16 | 2.76 | 2.74 | 91.73 |
| NeSLAM [107] | - | 2.57 | 2.46 | 92.66 |
| ADFP [95] | 3.01 | 2.77 | 2.45 | 92.79 |
| NID-SLAM [150] | 2.87 | 2.72 | 2.56 | 91.16 |
| Vox-Fusion [88] | - | 2.37 | 2.28 | **92.86** |
| Vox-Fusion++ [131] | - | 1.44 | **2.43** | 92.37 |
| Q-SLAM [129] | 1.87 | - | - | - |
| Co-SLAM [90] | 1.51 | - | - | - |
| NGEL-SLAM [127] | 1.28 | - | - | - |
| ESLAM [89] | 1.18 | - | - | - |
| SplaTAM [112] | 0.72 | - | - | - |
| Gaussian-SLAM [113] | 0.68 | - | - | - |
| HF-GS SLAM [117] | 0.52 | - | - | - |
| NEDS-SLAM [144] | 0.47 | - | - | - |
| Point-SLAM [93] | 0.44 | 1.41 | 3.10 | 88.89 |
| CG-SLAM [118] | - | **1.01** | 2.84 | 88.51 |
| Loopy-SLAM [129] | **0.35** | - | - | - |
| **RGB** | | | | |
| NeRF-SLAM [167] | 4.49 | - | - | - |
| DIM-SLAM [16] | - | 4.03 | 4.20 | 79.60 |
| GO-SLAM [91] | 4.39 | 3.81 | 4.79 | 78.00 |
| NICER-SLAM [18] | - | 3.65 | 4.16 | 79.37 |
| Hi-SLAM [158] | 3.63 | 3.62 | 4.59 | 80.60 |
| MoD-SLAM [163] | 3.23 | **2.48** | - | - |
| GlORIE-SLAM [165] | 3.24 | 2.96 | **3.95** | **83.72** |
| Q-SLAM [129] | **2.76** | - | - | - |

TABLE 5: **Replica [78] Mapping Results.** L1-Depth (↓), Acc. [cm] (↓), Comp. [cm] (↓) and Comp. Ratio [%] (↑) with 5 cm threshold are used as the evaluation metrics. * evaluates on ground truth poses.

SLAM, and MoD-SLAM (in its RGB-D version) once again stand out on Replica, confirming their effectiveness in optimizing camera tracking accuracy. Additionally, promising results are evident for methods utilizing Gaussian Splatting, with the best results among all achieved with CG-SLAM, HF-GS SLAM and GS-ICP SLAM. This suggests that these approaches struggle with noise and work best in simpler situations, showing less reliability in complex conditions—similar to what was observed in the TUM RGB-D and ScanNet datasets.

At the bottom, we collect results achieved by RGB-only frameworks. Again, we can notice how global BA and LC play a vital role in achieving the highest tracking accuracy; indeed, GO-SLAM and MoD-SLAM yield the best results, outperforming even most RGB-D frameworks not making use of either global BA or LC.

### 4.1.2 Mapping

**Replica.** In Table 5, we provide mapping results according to the evaluation protocol proposed in [5], highlighting the performance in terms of both 3D reconstruction and 2D depth estimation on the Replica dataset. Examining the table, a noticeable progression in both 3D reconstruction and 2D depth estimation metrics is observed, showcasing an improvement from iMap to more recent methods such

as NID-SLAM and ADFP. Notably, Loopy-SLAM leads in the L1-Depth metric, closely followed by Point-SLAM. This suggests that the neural point representation holds significant promise for generating highly accurate scene reconstructions. In terms of 3D error metrics, DNS SLAM, ADFP, NID-SLAM and Vox-Fusion++ outperform other methods, even surpassing hand-crafted approaches like COLMAP and TSDF, with Point-SLAM performing comparably, excelling in the Accuracy metric with a value of 1.41 indicates its superiority in this aspect, with CG-SLAM being the best performer with a score of 1.01. Notably, despite GO-SLAM's notable achievements in tracking, it holds a relatively low position in this ranking, indicating challenges for the mapping process. In Figure 13, qualitatives from a subset of reviewed systems on Replica are presented, emphasizing specific improvements achieved by recent methods in the mapping process.

Shifting focus to RGB methods, NICER-SLAM and Hi-SLAM exhibit a well-balanced performance, showcasing competitive scores in both Accuracy and Completion metrics. However, among these, GlORIE-SLAM stands out as

| Method | Metric | R0 | R1 | R2 | O0 | O1 | O2 | O3 | O4 | Avg |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | RGB-D | | | | | |
| Vox-Fusion [88] | PSNR↑ | 22.39 | 22.36 | 23.92 | 27.79 | 29.83 | 20.33 | 23.47 | 25.21 | 24.41 |
| | SSIM↑ | 0.68 | 0.75 | 0.80 | 0.86 | 0.88 | 0.79 | 0.80 | 0.85 | 0.80 |
| | LPIPS↓ | 0.30 | 0.27 | 0.23 | 0.24 | 0.18 | 0.24 | 0.21 | 0.20 | 0.24 |
| NICE-SLAM [5] | PSNR↑ | 22.12 | 22.47 | 24.52 | 29.07 | 30.34 | 19.66 | 22.23 | 24.94 | 24.42 |
| | SSIM↑ | 0.69 | 0.76 | 0.81 | 0.87 | 0.89 | 0.80 | 0.80 | 0.86 | 0.81 |
| | LPIPS↓ | 0.33 | 0.27 | 0.21 | 0.23 | 0.18 | 0.23 | 0.21 | 0.20 | 0.23 |
| GO-SLAM [91] | PSNR↑ | - | - | - | - | - | - | - | - | 27.38 |
| | SSIM↑ | - | - | - | - | - | - | - | - | 0.851 |
| | LPIPS↓ | - | - | - | - | - | - | - | - | - |
| ESLAM [89] | PSNR↑ | - | - | - | - | - | - | - | - | 27.80 |
| | SSIM↑ | - | - | - | - | - | - | - | - | 0.921 |
| | LPIPS↓ | - | - | - | - | - | - | - | - | 0.25 |
| MoD-SLAM [163] | PSNR↑ | - | - | - | - | - | - | - | - | 29.95 |
| | SSIM↑ | - | - | - | - | - | - | - | - | 0.862 |
| | LPIPS↓ | - | - | - | - | - | - | - | - | - |
| SplaTAM [112] | PSNR↑ | 32.86 | 33.89 | 35.25 | 38.26 | 39.17 | 31.97 | 29.70 | 31.81 | 34.11 |
| | SSIM↑ | 0.98 | 0.97 | 0.98 | 0.98 | 0.98 | 0.97 | 0.95 | 0.95 | 0.97 |
| | LPIPS↓ | 0.07 | 0.10 | 0.08 | 0.09 | 0.09 | 0.10 | 0.12 | 0.15 | 0.10 |
| GS-SLAM [12] | PSNR↑ | 31.56 | 32.86 | 32.59 | 38.70 | 41.17 | 32.36 | 32.03 | 32.92 | 34.27 |
| | SSIM↑ | 0.97 | 0.97 | 0.97 | 0.99 | 0.99 | 0.98 | 0.97 | 0.97 | 0.97 |
| | LPIPS↓ | 0.09 | 0.07 | 0.09 | 0.05 | 0.03 | 0.09 | 0.11 | 0.11 | 0.08 |
| Compact-GSSLAM [114] | PSNR↑ | 32.98 | 34.08 | 35.35 | 38.16 | 39.07 | 32.37 | 31.08 | 32.31 | 34.44 |
| | SSIM↑ | 0.99 | 0.98 | 0.98 | 0.98 | 0.98 | 0.97 | 0.96 | 0.96 | 0.98 |
| | LPIPS↓ | 0.07 | 0.10 | 0.08 | 0.09 | 0.09 | 0.10 | 0.12 | 0.15 | 0.09 |
| NEDS-SLAM [144] | PSNR↑ | 35.23 | 34.86 | 35.16 | 37.53 | 39.71 | 32.68 | 31.07 | 31.82 | 34.76 |
| | SSIM↑ | 0.98 | 0.86 | 0.98 | 0.98 | 0.98 | 0.97 | 0.97 | 0.97 | 0.96 |
| | LPIPS↓ | 0.08 | 0.08 | 0.07 | 0.09 | 0.09 | 0.08 | 0.10 | 0.11 | 0.09 |
| SemGauss-SLAM [143] | PSNR↑ | 32.55 | 33.92 | 35.15 | 39.18 | 39.87 | 32.97 | 31.60 | 35.00 | 35.03 |
| | SSIM↑ | 0.98 | 0.98 | 0.99 | 0.99 | 0.99 | 0.98 | 0.97 | 0.98 | 0.98 |
| | LPIPS↓ | 0.06 | 0.05 | 0.05 | 0.05 | 0.05 | 0.07 | 0.08 | 0.09 | 0.062 |
| Point-SLAM [93] | PSNR↑ | 32.40 | 34.08 | 35.50 | 38.26 | 39.16 | 33.99 | 33.48 | 33.49 | 35.17 |
| | SSIM↑ | 0.97 | 0.98 | 0.98 | 0.98 | 0.99 | 0.96 | 0.96 | 0.98 | 0.97 |
| | LPIPS↓ | 0.11 | 0.12 | 0.11 | 0.10 | 0.12 | 0.16 | 0.13 | 0.14 | 0.12 |
| Q-SLAM [164] | PSNR↑ | 33.24 | 34.81 | 34.16 | 39.32 | 39.51 | 34.08 | 32.65 | 34.93 | 35.34 |
| | SSIM↑ | 0.98 | 0.99 | 0.98 | 0.99 | 0.99 | 0.98 | 0.96 | 0.98 | 0.98 |
| | LPIPS↓ | 0.13 | 0.11 | 0.11 | 0.10 | 0.10 | 0.13 | 0.13 | 0.11 | 0.12 |
| Loopy-SLAM [129] | PSNR↑ | - | - | - | - | - | - | - | - | 35.47 |
| | SSIM↑ | - | - | - | - | - | - | - | - | 0.981 |
| | LPIPS↓ | - | - | - | - | - | - | - | - | 0.109 |
| NIDS-SLAM [139] | PSNR↑ | 33.16 | 35.18 | 36.49 | 40.22 | 38.90 | 34.22 | 34.74 | 33.24 | 35.76 |
| | SSIM↑ | **0.99** | 0.99 | 0.99 | 0.99 | 0.97 | 0.93 | 0.99 | 0.99 | 0.98 |
| | LPIPS↓ | - | - | - | - | - | - | - | - | - |
| HF-GS SLAM [117] | PSNR↑ | 33.06 | 35.74 | 37.21 | 41.12 | 41.11 | 33.56 | 33.21 | 34.48 | 36.19 |
| | SSIM↑ | 0.98 | 0.98 | 0.99 | 0.99 | 0.99 | 0.98 | 0.97 | 0.98 | 0.98 |
| | LPIPS↓ | 0.05 | 0.05 | 0.04 | 0.03 | 0.03 | 0.07 | 0.08 | 0.08 | 0.05 |
| GSSLAM [109] | PSNR↑ | 34.83 | 36.43 | 37.49 | 39.95 | 42.09 | 36.24 | 36.70 | 36.07 | 37.50 |
| | SSIM↑ | 0.95 | 0.96 | 0.96 | 0.97 | 0.98 | 0.96 | 0.96 | 0.96 | 0.96 |
| | LPIPS↓ | 0.07 | 0.08 | 0.07 | 0.07 | 0.06 | 0.08 | 0.07 | 0.10 | 0.07 |
| GS-ICP SLAM [115] | PSNR↑ | 35.37 | 37.80 | 38.50 | 43.13 | 43.26 | 36.93 | 36.90 | 38.75 | 38.83 |
| | SSIM↑ | 0.96 | 0.97 | 0.98 | 0.99 | 0.99 | 0.97 | 0.97 | 0.97 | 0.98 |
| | LPIPS↓ | 0.05 | 0.05 | 0.05 | 0.03 | 0.03 | 0.04 | 0.04 | 0.05 | 0.04 |
| Gaussian-SLAM [113] | PSNR↑ | **38.88** | **41.80** | **42.44** | **46.40** | **45.29** | **40.10** | **39.06** | **42.65** | **42.08** |
| | SSIM↑ | **0.99** | **1.0** | **1.0** | **1.00** | **1.0** | **1.0** | **1.0** | **1.0** | **1.0** |
| | LPIPS↓ | **0.02** | **0.02** | **0.02** | **0.02** | **0.02** | **0.02** | **0.02** | **0.02** | **0.02** |
| | | | | | RGB | | | | | |
| GO-SLAM [91] | PSNR↑ | - | - | - | - | - | - | - | - | 22.13 |
| | SSIM↑ | - | - | - | - | - | - | - | - | 0.73 |
| | LPIPS↓ | - | - | - | - | - | - | - | - | - |
| NICER-SLAM [18] | PSNR↑ | 25.33 | 23.92 | 26.12 | 28.54 | 25.86 | 21.95 | 26.13 | 25.47 | 25.41 |
| | SSIM↑ | 0.75 | 0.77 | 0.83 | 0.87 | 0.85 | 0.82 | 0.86 | 0.87 | 0.83 |
| | LPIPS↓ | 0.25 | 0.22 | 0.18 | 0.17 | 0.18 | 0.20 | 0.16 | 0.18 | 0.19 |
| MoD-SLAM [163] | PSNR↑ | - | - | - | - | - | - | - | - | 27.31 |
| | SSIM↑ | - | - | - | - | - | - | - | - | 0.85 |
| | LPIPS↓ | - | - | - | - | - | - | - | - | - |
| GlORIE-SLAM [165] | PSNR↑ | 28.49 | 30.09 | 29.98 | 35.88 | 37.15 | 28.45 | 28.54 | 29.73 | 31.04 |
| | SSIM↑ | 0.96 | 0.97 | **0.96** | 0.98 | 0.99 | 0.97 | 0.97 | 0.97 | **0.97** |
| | LPIPS↓ | **0.13** | **0.13** | 0.14 | 0.09 | 0.08 | 0.15 | 0.11 | 0.15 | **0.12** |
| Q-SLAM [164] | PSNR↑ | **29.58** | **32.74** | **31.25** | **36.31** | **37.22** | **30.68** | **30.21** | **31.96** | 32.49 |
| | SSIM↑ | 0.83 | 0.91 | 0.87 | 0.94 | 0.94 | 0.90 | 0.88 | 0.89 | 0.89 |
| | LPIPS↓ | 0.18 | 0.16 | 0.15 | 0.13 | 0.15 | 0.20 | 0.19 | 0.18 | 0.17 |
| Photo-SLAM [110] | PSNR↑ | - | - | - | - | - | - | - | - | **33.30** |
| | SSIM↑ | - | - | - | - | - | - | - | - | 0.93 |
| | LPIPS↓ | - | - | - | - | - | - | - | - | - |

TABLE 6: **Replica [78] Train View Rendering Results**. We report PSNR, SSIM, and LPIPS metrics.

| Method | External Priors | R0 | R1 | R2 | O0 | O1 | O2 | O3 | O4 |
|---|---|---|---|---|---|---|---|---|---|
| | | | | | RGB-D | | | | |
| NIDS-SLAM [139] | Mask2Former [140] | 82.45 | 84.08 | 76.99 | 85.94 | – | – | – | – |
| DNS SLAM [141] | | 88.32 | 84.90 | 81.20 | 84.66 | – | – | – | – |
| SNI-SLAM [137] | Dinov2 [138] | 88.42 | 87.43 | 86.16 | 87.63 | 78.63 | 86.49 | 74.01 | 80.22 |
| NEDS-SLAM [144] | | 90.73 | 91.20 | - | 90.42 | - | - | - | - |
| SemGauss-SLAM [143] | Dinov2 [138] | **92.81** | **94.10** | **94.72** | **95.23** | **90.11** | **94.93** | **92.93** | **94.82** |

TABLE 7: **Replica [78] Semantic Results**. Quantitative comparison of input views semantic segmentation performance on the Replica dataset [78] using the mIoU metric.

### 4.1.3 Image Rendering

**Replica.** In Table 6, we show the rendering quality on the training input views of Replica, following the standard evaluation approach of Point-SLAM and NICE-SLAM.

On top, we focus on RGB-D frameworks: recent solutions, particularly those based on Gaussian Splatting or neural points such as Point-SLAM, yield significantly better average metrics in terms of PSNR, SSIM, and LPIPS compared to methodologies proposed in the early stages of the evolution of neural SLAM (showing an improvement of over 10dB in PSNR). These early methods are based on multi-resolution feature grids such as NICE-SLAM or voxel-based neural implicit surface representations such as Vox-Fusion. This suggests that paradigms based on explicit Gaussian primitives or neural points lead to significant improvements in image rendering. Among the latter, Gaussian-SLAM stands out with an impressive average PSNR of 42.08 [113], further underscoring the effectiveness of Gaussian-based approaches in achieving high-quality image rendering.

At the bottom, when considering RGB-only methods, the use of the 3DGS framework allows Photo-SLAM to generate novel view renderings with superior quality compared to other NeRF-style SLAM systems.

In Figure 14, we present qualitative results for image rendering from selected RGB-D SLAM systems on Replica. The latest frameworks demonstrate improved rendering of fine details, with GS-SLAM showing superior rendering quality due to its 3DGS representation.

In our analysis, we agree with the findings presented in the SplaTAM paper. In particular, we share concerns about the relevance of the rendering results on the Replica dataset. Evaluating the same training views used as input raises valid concerns about potential biases introduced by high model capacity and the risk of overfitting to these specific images. Our agreement with the SplaTAM perspective leads us to support the exploration of alternative methods for evaluating novel view rendering in this specific context. It is crucial to emphasize that our agreement with these observations is rooted in a collective understanding of the constraints inherent in the current SLAM benchmarks.

### 4.1.4 Semantic Segmentation Results

**Replica.** Table 7 presents a comparative analysis of state-of-the-art RGB-D semantic SLAM methods on the Replica dataset [78], using the mIoU metric for evaluating the semantic segmentation performance of input views, following the evaluation protocol from SemGauss-SLAM [143]. The methods compared include NIDS-SLAM [139], DNS SLAM [141], SNI-SLAM [137], NEDS-SLAM [144], and SemGauss-SLAM [143]. The table highlights the use of external priors, such as Mask2Former [140] and Dinov2 [138], by some of

the most accurate. Nonetheless, the distinction among different methods in the RGB context is less pronounced compared to RGB-D scenarios. Notably, it becomes evident that, expectedly, methods relying solely on RGB perform less favorably than those leveraging depth sensor information. The exception to this trend is iMAP. This emphasizes the crucial role of depth sensors in SLAM and points towards the potential for advancements in RGB-only methodologies.

| Method | 00 | 01 | 02 | 03 | 04 | 05 | 06 | 07 | 08 | 09 | 10 | Avg. | 11-21 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| LiDAR Odometry Evaluation | | | | | | | | | | | | | |
| MULLS [186] | 0.56 | 0.64 | 0.55 | 0.71 | 0.41 | 0.30 | 0.30 | 0.38 | 0.78 | 0.48 | 0.59 | 0.52 | 0.65 |
| CT-ICP [187] | 0.49 | 0.76 | 0.52 | 0.72 | 0.39 | 0.25 | 0.27 | 0.31 | 0.81 | 0.49 | 0.48 | 0.50 | 0.59 |
| SuMa-LO [189] | 0.72 | 1.71 | 1.06 | 0.66 | 0.38 | 0.50 | 0.41 | 0.55 | 1.02 | 0.48 | 0.71 | 0.75 | 1.39 |
| Litamin-LO [188] | 0.78 | 2.10 | 0.95 | 0.96 | 1.05 | 0.55 | 0.55 | 0.48 | 1.01 | 0.69 | 0.80 | 0.88 | - |
| Nerf-LOAM [15] | 1.34 | 2.07 | - | 2.22 | 1.74 | 1.40 | - | 1.00 | - | 1.63 | 2.08 | 1.69 | - |
| PIN-LO [170] | 0.55 | 0.54 | 0.52 | 0.74 | 0.28 | 0.29 | 0.32 | 0.36 | 0.83 | 0.56 | 0.47 | 0.50 | |
| | 00† | 01 | 02† | 03 | 04 | 05† | 06† | 07† | 08† | 09† | 10 | Avg.† | Avg. |
| LiDAR SLAM Evaluation | | | | | | | | | | | | | |
| MULLS [186] | 1.1 | 1.9 | 5.4 | 0.7 | 0.9 | 1.0 | 0.3 | 0.4 | 2.9 | 2.1 | 1.1 | 1.9 | 1.6 |
| SuMa [189] | 1.0 | 13.8 | 7.1 | 0.9 | 0.4 | 0.6 | 0.6 | 1.0 | 3.4 | 1.1 | 1.3 | 2.1 | 3.2 |
| Litamin2 [188] | 1.3 | 15.9 | 3.2 | 0.8 | 0.7 | 0.6 | 0.8 | 0.5 | 2.1 | 2.1 | 1.0 | 1.5 | 2.4 |
| HLBA [191] | 0.8 | 1.9 | 5.1 | 0.6 | 0.8 | 0.4 | 0.2 | 0.3 | 2.7 | 1.3 | 1.1 | 1.5 | 1.4 |
| PIN-LO [170] | 4.3 | 2.0 | 7.3 | 0.7 | 0.1 | 2.1 | 0.7 | 0.4 | 3.5 | 1.8 | 0.6 | 2.9 | 2.1 |
| PIN-SLAM [170] | 0.8 | 2.0 | 3.3 | 0.7 | 0.1 | 0.2 | 0.4 | 0.3 | 1.7 | 1.0 | 0.6 | 1.1 | 1.0 |

TABLE 8: **KITTI [79] LiDAR Odometry/SLAM Results.** † indicates sequences with loops and Avg.† denotes the average metric for such sequences.

| Method | 01 | 02 | quad | math_e | ug_e | cloister_e | stairs | Avg. |
|---|---|---|---|---|---|---|---|---|
| MULLS [186] | 2.51 | 8.39 | 0.12 | 0.35 | 0.86 | 0.41 | - | 2.11 |
| SuMa [189] | 2.03 | 3.65 | 0.28 | 0.16 | 0.09 | 0.20 | 1.85 | 1.18 |
| PIN-LO [170] | 2.21 | 4.93 | 0.09 | 0.10 | 0.07 | 0.41 | 0.06 | 1.12 |
| PIN-SLAM [170] | 0.43 | 0.30 | 0.09 | 0.09 | 0.07 | 0.18 | 0.06 | 0.19 |

TABLE 9: **Newer College [80] Camera Tracking Results.** ATE RMSE [cm] (↓) is used as the evaluation metric.

| Method | Quad | | Math Institute | |
|---|---|---|---|---|
| | Acc. [cm] ↓ | Comp. [cm] ↓ | Acc. [cm] ↓ | Comp. [cm] ↓ |
| SLAMesh [190] | 19.21 | 48.83 | 12.80 | 23.50 |
| Nerf-LOAM [15] | 12.89 | 22.21 | - | - |
| PIN-SLAM [170] | 11.55 | 15.25 | 13.70 | 21.91 |

TABLE 10: **Newer College [80] Mapping Results.** Acc. [cm] (↓) and Comp. [cm] (↓) are used as the evaluation metrics.

| Method | Scene Encoding | GPU Mem. [G] ↓ | Avg. FPS ↑ |
|---|---|---|---|
| RGB-D | | | |
| iMAP [1] | MLP | 6.44 | 0.13 |
| SplaTAM [112] | 3D Gaussians | 18.54 | 0.14 |
| Point-SLAM [93] | Neural Points + MLP | 7.11 | 0.23 |
| UncLe-SLAM [17] | Hier. Grid + MLP | 8.24 | 0.24 |
| NICE-SLAM [5] | Hier. Grid + MLP | 4.70 | 0.61 |
| ADFP [95] | Hier. Grid + MLP | 3.76 | 0.74 |
| Vox-Fusion [88] | Sparse Voxels + MLP | 21.22 | 0.74 |
| Plenoxel-SLAM [97] | Plenoxels | 13.04 | 1.25 |
| ESLAM [89] | Feature Planes + MLP | 13.04 | 4.62 |
| Co-SLAM [90] | Hash Grid + MLP | 3.56 | 7.97 |
| GO-SLAM [91] | Hash Grid + MLP | 18.50 | 8.36 |
| RGB | | | |
| DIM-SLAM [16] | Hier. Grid + MLP | 4.78 | 3.14 |
| Orbeez-SLAM [153] | Voxels + MLP | 7.55 | 17.70 |
| NeRF-SLAM [167] | Hash Grid + MLP | 9.38 | 20.00 |
| LiDAR | | | |
| Nerf-LOAM [15] | Sparse Voxel + MLP | 11.58 | 0.24 |
| PIN-SLAM [170] | Neural Points + MLP | 6.93 | 6.67 |

TABLE 11: **Performance Evaluation:** GPU memory requirements (GB) and average FPS efficiency on Replica room0 (RGB/RGB-D) and KITTI 00 sequence (LiDAR).

these methods to improve their semantic understanding capabilities. Among the compared methods, SemGauss-SLAM achieves the highest mIoU scores across all eight scenes of the Replica dataset, demonstrating its superior performance in semantic segmentation.

## 4.2 LiDAR SLAM/Odometry Evaluation

### 4.2.1 Tracking

**KITTI**. Table 8 presents the evaluation of LiDAR SLAM strategies on the KITTI dataset, detailing odometry accuracy at the top and SLAM performance metrics at the bottom. The odometry section reports the average relative translational drift error (%) and highlights the performance of PIN-LO, a variant of PIN-SLAM that disables the loop closure detection correction and pose graph optimization modules. PIN-LO outperforms several LiDAR odometry systems using different map representations (feature points [186], denser voxel downsampling points [187], normal distribution transformation [188], surfels [189] and triangle meshes [190]), achieving an impressive translation error of 0.5%, competing with KISS-ICP and CT-ICP, and outperforming the neural implicit approach Nerf-LOAM due to improved SDF training and robust point-to-SDF registration.

In the LiDAR SLAM evaluation at the bottom of the table 8, the ATE RMSE [m] is used as the evaluation metric. As a representative of implicit LiDAR-based SLAM strategies, PIN-SLAM consistently outperforms state-of-the-art LiDAR SLAM systems. Specifically, PIN-SLAM achieves an average RMSE of 1.1 m on sequences with loop closure and 1.0 m over all eleven sequences. The results of PIN-LO underscore the significant improvement of PIN-SLAM in ensuring global trajectory consistency.

**Newer College**. Table 9 reports the tracking accuracy on the Newer College dataset, measured in terms of ATE RMSE [cm]. Again, we can observe how PIN-SLAM consistently outperforms PIN-LO, with an average RMSE of 0.19 cm over the whole set of sequences, which is 5× lower compared to

PIN-LO. This further confirms the superiority of PIN-SLAM at global trajectory tracking.

### 4.2.2 Mapping

**Newer College**. Table 10 collects the results concerning the quality of 3D reconstruction on the New College dataset – specifically, on *Quad* and *Math Institute* sequences. Accuracy and Completeness scores are used to assess the effectiveness of Nerf-LOAM and PIN-SLAM, with the latter confirming again as the best LiDAR-based SLAM system among those evaluating on this dataset. In particular, on *Quad* we can appreciate a large margin in terms of completeness between PIN-SLAM and Nerf-LOAM – *i.e.*, about 7 cm.

## 4.3 Performance Analysis

We conclude the experimental studies by considering the efficiency of the SLAM systems reviewed so far. For this purpose, we run methods with source code publicly available and measure 1) the GPU memory requirements (as the peak memory use in GB) and 2) the average FPS (computed as the total time required to process a single sequence, divided by the total amount of frames in it) achieved on a single NVIDIA RTX 3090 board. Table 11 collects the outcome of our benchmark for RGB-D and RGB systems running on Replica, sorted in increasing order of average FPS. On top, we consider RGB-D frameworks: we can notice how SplaTAM, despite its high efficiency at rendering images, is however much slower at processing both tracking and mapping simultaneously. This is also the case for hybrid methods using hierarchical feature grids, on the other hand require much less GPU memory – 4 to 5× lower compared to SplaTAM. Finally, the use of more advanced representations such as hash grids or point features allows for much

faster processing. This is confirmed also by the studies on the RGB-only methods, in the middle, with NeRF-SLAM resulting 6× faster than DIM-SLAM. Finally, concerning LiDAR SLAM systems, we can observe how PIN-SLAM is much more efficient than Nerf-LOAM, requiring as few as 7 GB of GPU memory while running at nearly 7 FPS, compared to the nearly 12 GB and 4 seconds per frame required by Nerf-LOAM.

This analysis highlights how, despite the great promise brought by this new generation of SLAM systems, most of them are still unsatisfactory in terms of hardware and runtime requirements, making them not yet ready for real-time applications.

## 5 DISCUSSION

In this section we focus on highlighting the key findings of the survey. We will outline the main advances achieved through the most recent methodologies examined, while identifying ongoing challenges and potential avenues for future research in this area.

**Scene Representation.** The choice of scene representation is critical in current SLAM solutions, significantly affecting mapping/tracking accuracy, rendering quality, and computation. Early approaches, such as iMAP [1], used network-based methods, implicitly modeling scenes with coordinate-based MLP(s). While these provide compact, continuous modeling of the scene, they struggle with real-time reconstruction due to challenges in updating local regions and scaling for large scenes. In addition, they tend to produce over-smoothed scene reconstructions. Subsequent research has explored grid-based representations, such as multi-resolution hierarchical [5], [90] and sparse octree grids [88], [125], which have gained popularity. Grids allow for fast neighbor lookups, but require a pre-specified grid resolution, resulting in inefficient memory use in empty space and a limited ability to capture fine details constrained by the resolution. Recent advances, such as Point-SLAM [93] and Loopy-SLAM [129], favor hybrid neural point-based representations. Unlike grids, point densities vary naturally and need not be pre-specified. Points concentrate efficiently around surfaces while assigning higher density to details, facilitating scalability and local updates compared to network-based methods. At present, point-based methods have demonstrated superior performance in 3D reconstruction, yielding highly accurate 3D surfaces, as evidenced by experiments conducted on the Replica dataset. However, similar to other NeRF-style approaches, volumetric ray sampling significantly restricts its efficiency.

Promising techniques include explicit representations based on the 3D Gaussian Splatting paradigm. Explicit representations based on 3DGS have been shown to achieve state-of-the-art rendering accuracy compared to other representations while also exhibiting faster rendering. However, these methods have several limitations, including a heavy reliance on initialization and a lack of control over primitive growth in unobserved regions. Furthermore, the original 3DGS-based scene representation requires a large number of 3D Gaussian primitives to achieve high-fidelity reconstruction, resulting in substantial memory consumption.

Despite significant progress over the past three years, ongoing research is still actively engaged in overcoming existing scene representation limitations and finding ever more effective alternatives to improve accuracy and real-time performance in SLAM.

**Catastrophic Forgetting.** Existing methods often exhibit a tendency to forget previously learned information, particularly in large scenarios or extended video sequences. In the case of network-based methods, this is attributed to their reliance on single neural networks or global models with fixed capacity, which are affected by global changes during optimization. One common approach to alleviate this problem is to train the network using sparse ray sampling with current observations while replaying keyframes from historical data. However, in large-scale incremental mapping, such a strategy results in a cumulative increase in data, requiring complex resampling procedures for memory efficiency. The forgetting problem extends to grid-based approaches. Despite efforts to address this issue, obstacles arise due to quadratic or cubic spatial complexity, which poses scalability challenges. Similarly, while explicit representations, such as 3DGS-style solutions, offer a practical workaround for catastrophic forgetting, they face challenges due to increased memory requirements and slow processing, especially in large scenes. Some methods attempt to mitigate these limitations by employing sparse frame sampling, but this leads to inefficient information sampling across 3D space, resulting in slower and less uniform model updates compared to approaches that integrate sparse ray sampling.

Eventually, some strategies recommend dividing the environment into submaps and assigning local SLAM tasks to different agents. However, this introduces additional challenges in handling multiple distributed models and devising efficient strategies to manage overlapping regions while preventing the occurrence of map fusion artifacts.

**Real-Time Constraints**. Many of the techniques reviewed face challenges in achieving real-time processing, often failing to match the sensor frame rate. This limitation is mainly due to the chosen map data structure or the computationally intensive ray-wise rendering-based optimization, which is especially noticeable in NeRF-style SLAM methods. In particular, hybrid approaches using hierarchical grids require less GPU memory but exhibit slower runtime performance. On the other hand, advanced representations such as hash grids or sparse voxels allow for faster computation, but with higher memory requirements. Finally, despite their advantages in fast image rendering, current 3DGS-style methods still struggle to efficiently handle simultaneous tracking and mapping processing, preventing their effective use in real-time applications.

**Global Optimization.** Implementing LC and global BA requires significant computational resources, risking performance bottlenecks, especially in real-time applications. Many reviewed frame-to-model methods (*e.g.*, iMap [1], NICE-SLAM [5], etc.) face challenges with loop closure and global bundle adjustment due to the prohibitive computational complexity of updating the entire 3D model. In contrast, frame-to-frame techniques (*e.g.*, GO-SLAM [91], etc.) facilitate global correction by executing the global BA in a background thread, which significantly improves tracking accuracy, as demonstrated in the reported experiments,

although at a slower computational speed compared to real-time rates. For both approaches, the computational cost is largely due to the lack of flexibility of latent feature grids to accommodate pose corrections from loop closures. Indeed, this requires re-allocating feature grids and retraining the entire map once a loop is corrected and poses are updated. However, this challenge becomes more pronounced as the number of frames processed increases, leading to the accumulation of camera drift errors and eventually either an inconsistent 3D reconstruction or a rapid collapse of the reconstruction process.

Overall, decoupled methods, which separate the mapping and tracking processes, tend to achieve better tracking performance compared to coupled approaches. By allowing the tracking module to focus solely on camera pose estimation without the added complexity of simultaneously updating the map representation, decoupled methods can achieve more accurate and robust tracking. However, this improved accuracy and robustness come at the cost of increased computational overhead, as the independent mapping and tracking stages require separate processing pipelines and memory allocation, which may impact the overall efficiency of the SLAM system.

**NeRF vs. 3DGS in SLAM.** NeRF-style SLAM, which relies mostly on MLP(s), is well suited for novel view synthesis, mapping and tracking but faces challenges such as oversmoothing, susceptibility to catastrophic forgetting, and computational inefficiency due to its reliance on per-pixel ray marching. 3DGS bypasses per-pixel ray marching and exploits sparsity through differentiable rasterization over primitives. This benefits SLAM with an explicit volumetric representation, fast rendering, rich optimization, direct gradient flow, increased map capacity, and explicit spatial extent control. Thus, while NeRF shows a remarkable ability to synthesize novel views, its slow training speed and difficulty in adapting to SLAM are significant drawbacks. 3DGS, with its efficient rendering, explicit representation, and rich optimization capabilities, emerges as a powerful alternative. Despite its advantages, current 3DGS-style SLAM approaches have limitations. These include scalability issues for large scenes, the lack of a direct mesh extraction algorithm (although recent methods such as [192] have been proposed), the inability to accurately encode precise geometry and, among others, the potential for uncontrollable Gaussian growth into unobserved areas, causing artifacts in rendered views and the underlying 3D structure. Moreover, the computational complexity of 3DGS-based SLAM systems is significantly higher than NeRF-based methods, which can hinder real-time performance and practical deployment, especially on resource-constrained devices. In order to mitigate these issues, recent research efforts, such as Compact-GSSLAM [114], have focused on developing compact 3D Gaussian scene representations that optimize storage efficiency while maintaining high-quality reconstruction, rapid training convergence, and real-time rendering capabilities.

**Evaluation Inconsistencies.** The lack of standardized benchmarks or online servers with well-defined evaluation protocols results in inconsistent evaluation methods, making it difficult to conduct fair comparisons between approaches and introducing inconsistencies within the methodologies presented in different research papers.

This is exemplified by challenges in datasets such as ScanNet [77], where ground truth poses are derived from BundleFusion [53], raising concerns about the reliability and generalizability of evaluation results. Xu et al. [193] and Hua et al. [194] both acknowledge these inconsistencies and propose solutions to address them. Xu et al. [193] introduce a comprehensive taxonomy of perturbations for SLAM in dynamic and unstructured environments, along with the Robust-SLAM dataset[12], created using 3D scene models sourced from Replica, which includes diverse perturbation types and offers a consistent evaluation protocol. Similarly, Hua et al. [194] establish an open-source benchmark framework[13] to evaluate the performance of a wide spectrum of commonly used implicit neural representations and geometric rendering functions for examining their effectiveness in mapping and localization. They propose a novel RGB-D SLAM benchmark framework, featuring a unified evaluation protocol to assess different NeRF components effectively. These works highlight the importance of standardized benchmarks and evaluation protocols in mitigating inconsistencies and enabling more reliable and generalizable research outcomes in SLAM. However, to further address these issues, we believe it is crucial to establish online evaluation platforms with well-defined protocols, error metrics, and leaderboards for tracking, mapping, similar in spirit to the ETH3D benchmark [30]. These online benchmarks should provide high-quality ground truth data for both mapping and tracking, ensuring that the proposed methods are evaluated against reliable and accurate reference data. Moreover, we consider that they should include a dedicated evaluation protocol for novel view rendering to address overfitting risks and promote more generalizable rendering methods. In summary, by adopting standardized benchmarks, well-defined protocols, and high-quality ground truth data, we believe that the research community can make more informed and fair comparisons between different approaches.

**Additional Challenges.** SLAM approaches, whether traditional, deep learning based, or influenced by radiance field representations, face common challenges. One notable obstacle is the handling of *dynamic scenes*, which proves difficult due to the underlying assumption of a static environment, leading to artifacts in the reconstructed scene and errors in the tracking process. While some approaches attempt to address this issue, there is still significant room for improvement, especially in highly dynamic environments.

Another challenge is sensitivity to *sensor noise*, which includes motion blur, depth noise, and aggressive rotation, all of which affect tracking and mapping accuracy. This is further compounded by the presence of *non-Lambertian objects* in the scene, such as glass or metal surfaces, which introduce additional complexity due to their varying reflective properties. In the context of these challenges, it is noteworthy that many approaches often overlook explicit *uncertainty estimation* across input modalities, hindering a comprehensive understanding of system reliability.

Additionally, the absence of external sensors, especially depth information, poses a fundamental problem to *RGB-*

---

12. https://github.com/Xiaohao-Xu/SLAM-under-Perturbation/
13. https://vlis2022.github.io/nerf-slam-benchmark/

*only* SLAM, leading to depth ambiguity and 3D reconstruction optimization convergence issues.

A less critical but specific issue is the quality of rendered images of the scene. Reviewed techniques often struggle with *view-dependent appearance* elements, such as specular reflections, due to the lack of modeling of view directions in the model, which affects rendering quality.

# 6 CONCLUSION

In summary, this overview pioneers the exploration of SLAM methods influenced by recent advances in radiance field representations. Ranging from seminal works such as iMap [1] to the latest advances, the review reveals a substantial body of literature that has emerged in just three years. Through structured classification and analysis, it highlights key limitations and innovations, providing valuable insights with comparative results across tracking, mapping, and rendering. It also identifies current open challenges, providing interesting avenues for future exploration.

As a result, this survey is intended to serve as an essential guide for both novices and seasoned experts, establishing itself as a comprehensive reference in this rapidly evolving field.

## REFERENCES

[1] E. Sucar, S. Liu, J. Ortiz, and A. J. Davison, "imap: Implicit mapping and positioning in real-time," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 6229–6238.

[2] R. A. Newcombe, S. Izadi, O. Hilliges, D. Molyneaux, D. Kim, A. J. Davison, P. Kohi, J. Shotton, S. Hodges, and A. Fitzgibbon, "Kinectfusion: Real-time dense surface mapping and tracking," in *IEEE ISMAR*. IEEE, 2011, pp. 127–136.

[3] K. Tateno, F. Tombari, I. Laina, and N. Navab, "Cnn-slam: Real-time dense monocular slam with learned depth prediction," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 6243–6252.

[4] Y. Li, N. Brasch, Y. Wang, N. Navab, and F. Tombari, "Structure-slam: Low-drift monocular slam in indoor environments," *IEEE Robotics and Automation Letters*, vol. 5, no. 4, pp. 6583–6590, 2020.

[5] Z. Zhu, S. Peng, V. Larsson, W. Xu, H. Bao, Z. Cui, M. R. Oswald, and M. Pollefeys, "Nice-slam: Neural implicit scalable encoding for slam," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 12786–12796.

[6] E. Kruzhkov, A. Savinykh, P. Karpyshev, M. Kurenkov, E. Yudin, A. Potapov, and D. Tsetserukou, "Meslam: Memory efficient slam based on neural fields," in *2022 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*. IEEE, 2022, pp. 430–435.

[7] S. Zhi, E. Sucar, A. Mouton, I. Haughton, T. Laidlow, and A. J. Davison, "ilabel: Revealing objects in neural fields," *IEEE Robotics and Automation Letters*, vol. 8, no. 2, pp. 832–839, 2022.

[8] T. Hua, H. Bai, Z. Cao, M. Liu, D. Tao, and L. Wang, "Hi-map: Hierarchical factorized radiance field for high-fidelity monocular dense mapping," *arXiv preprint arXiv:2401.03203*, 2024.

[9] M. Li, J. He, G. Jiang, and H. Wang, "Ddn-slam: Real-time dense dynamic neural implicit slam with joint semantic encoding," *arXiv preprint arXiv:2401.01545*, 2024.

[10] R. Mur-Artal, J. M. M. Montiel, and J. D. Tardos, "Orb-slam: a versatile and accurate monocular slam system," *IEEE transactions on robotics*, vol. 31, no. 5, pp. 1147–1163, 2015.

[11] M. Bloesch, J. Czarnowski, R. Clark, S. Leutenegger, and A. J. Davison, "Codeslam—learning a compact, optimisable representation for dense visual slam," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 2560–2568.

[12] C. Yan, D. Qu, D. Wang, D. Xu, Z. Wang, B. Zhao, and X. Li, "Gs-slam: Dense visual slam with 3d gaussian splatting," *arXiv preprint arXiv:2311.11700*, 2023.

[13] C. Ruan, Q. Zang, K. Zhang, and K. Huang, "Dn-slam: A visual slam with orb features and nerf mapping in dynamic environments," *IEEE Sensors Journal*, 2023.

[14] D. Qu, C. Yan, D. Wang, J. Yin, D. Xu, B. Zhao, and X. Li, "Implicit event-rgbd neural slam," *arXiv preprint arXiv:2311.11013*, 2023.

[15] J. Deng, Q. Wu, X. Chen, S. Xia, Z. Sun, G. Liu, W. Yu, and L. Pei, "Nerf-loam: Neural implicit representation for large-scale incremental lidar odometry and mapping," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 8218–8227.

[16] H. Li, X. Gu, W. Yuan, L. Yang, Z. Dong, and P. Tan, "Dense rgb slam with neural implicit maps," in *Proceedings of the International Conference on Learning Representations*, 2023. [Online]. Available: https://openreview.net/forum?id=QUK1ExlbbA

[17] E. Sandström, K. Ta, L. V. Gool, and M. R. Oswald, "Uncle-SLAM: Uncertainty learning for dense neural SLAM," in *International Conference on Computer Vision Workshops (ICCVW)*, 2023.

[18] Z. Zhu, S. Peng, V. Larsson, Z. Cui, M. R. Oswald, A. Geiger, and M. Pollefeys, "Nicer-slam: Neural implicit scene encoding for rgb slam," in *International Conference on 3D Vision (3DV)*, March 2024.

[19] G. Grisetti, R. Kümmerle, C. Stachniss, and W. Burgard, "A tutorial on graph-based slam," *IEEE Intelligent Transportation Systems Magazine*, vol. 2, no. 4, pp. 31–43, 2010.

[20] K. Yousif, A. Bab-Hadiashar, and R. Hoseinnezhad, "An overview to visual odometry and visual slam: Applications to mobile robotics," *Intelligent Industrial Systems*, vol. 1, no. 4, pp. 289–311, 2015.

[21] C. Cadena, L. Carlone, H. Carrillo, Y. Latif, D. Scaramuzza, J. Neira, I. Reid, and J. J. Leonard, "Past, present, and future of simultaneous localization and mapping: Toward the robust-perception age," *IEEE Transactions on robotics*, vol. 32, no. 6, pp. 1309–1332, 2016.

[22] T. Taketomi, H. Uchiyama, and S. Ikeda, "Visual slam algorithms: A survey from 2010 to 2016," *IPSJ Transactions on Computer Vision and Applications*, vol. 9, no. 1, pp. 1–11, 2017.

[23] C. Duan, S. Junginger, J. Huang, K. Jin, and K. Thurow, "Deep learning for visual slam in transportation robotics: A review," *Transportation Safety and Environment*, vol. 1, no. 3, pp. 177–184, 2019.

[24] S. Mokssit, D. B. Licea, B. Guermah, and M. Ghogho, "Deep learning techniques for visual slam: A survey," *IEEE Access*, vol. 11, pp. 20026–20050, 2023.

[25] R. A. Newcombe, S. J. Lovegrove, and A. J. Davison, "Dtam: Dense tracking and mapping in real-time," in *2011 international conference on computer vision*. IEEE, 2011, pp. 2320–2327.

[26] R. F. Salas-Moreno, R. A. Newcombe, H. Strasdat, P. H. Kelly, and A. J. Davison, "Slam++: Simultaneous localisation and mapping at the level of objects," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2013, pp. 1352–1359.

[27] T. Whelan, S. Leutenegger, R. Salas-Moreno, B. Glocker, and A. Davison, "Elasticfusion: Dense slam without a pose graph." Robotics: Science and Systems, 2015.

[28] Z. Teed and J. Deng, "Droid-slam: Deep visual slam for monocular, stereo, and rgb-d cameras," *Advances in neural information processing systems*, vol. 34, pp. 16558–16569, 2021.

[29] Q.-Y. Zhou, S. Miller, and V. Koltun, "Elastic fragments for dense scene reconstruction," in *Proceedings of the IEEE International Conference on Computer Vision*, 2013, pp. 473–480.

[30] T. Schops, T. Sattler, and M. Pollefeys, "Bad slam: Bundle adjusted direct rgb-d slam," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 134–144.

[31] M. Nießner, M. Zollhöfer, S. Izadi, and M. Stamminger, "Real-time 3d reconstruction at scale using voxel hashing," *ACM Transactions on Graphics (ToG)*, vol. 32, no. 6, pp. 1–11, 2013.

[32] F. Steinbrucker, C. Kerl, and D. Cremers, "Large-scale multi-resolution surface reconstruction from rgb-d sequences," in *Proceedings of the IEEE International Conference on Computer Vision*, 2013, pp. 3264–3271.

[33] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "Nerf: Representing scenes as neural radiance fields for view synthesis," *Communications of the ACM*, vol. 65, no. 1, pp. 99–106, 2021.

[34] B. Kerbl, G. Kopanas, T. Leimkühler, and G. Drettakis, "3d gaussian splatting for real-time radiance field rendering," *ACM Transactions on Graphics*, vol. 42, no. 4, 2023.

[35] L. Mescheder, M. Oechsle, M. Niemeyer, S. Nowozin, and A. Geiger, "Occupancy networks: Learning 3d reconstruction in function space," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 4460–4470.

[36] Z. Chen and H. Zhang, "Learning implicit fields for generative shape modeling," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.

[37] J. J. Park, P. Florence, J. Straub, R. Newcombe, and S. Lovegrove, "Deepsdf: Learning continuous signed distance functions for shape representation," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 165–174.

[38] Y. Xie, T. Takikawa, S. Saito, O. Litany, S. Yan, N. Khan, F. Tombari, J. Tompkin, V. Sitzmann, and S. Sridhar, "Neural fields in visual computing and beyond," in *Computer Graphics Forum*, vol. 41, no. 2. Wiley Online Library, 2022, pp. 641–676.

[39] D. Li, "awesome-implicit-nerf-slam," 2022. [Online]. Available: https://github.com/DoongLi/awesome-Implicit-NeRF-SLAM

[40] H. Durrant-Whyte and T. Bailey, "Simultaneous localization and mapping: part i," *IEEE robotics & automation magazine*, vol. 13, no. 2, pp. 99–110, 2006.

[41] T. Bailey and H. Durrant-Whyte, "Simultaneous localization and mapping (slam): Part ii," *IEEE robotics & automation magazine*, vol. 13, no. 3, pp. 108–117, 2006.

[42] S. Saeedi, M. Trentini, M. Seto, and H. Li, "Multiple-robot simultaneous localization and mapping: A review," *Journal of Field Robotics*, vol. 33, no. 1, pp. 3–46, 2016.

[43] M. R. U. Saputra, A. Markham, and N. Trigoni, "Visual slam and structure from motion in dynamic environments: A survey," *ACM Computing Surveys (CSUR)*, vol. 51, no. 2, pp. 1–36, 2018.

[44] M. Zaffar, S. Ehsan, R. Stolkin, and K. M. Maier, "Sensors, slam and long-term autonomy: A review," in *2018 NASA/ESA Conference on Adaptive Hardware and Systems (AHS)*. IEEE, 2018, pp. 285–290.

[45] J. Yang, Y. Li, L. Cao, Y. Jiang, L. Sun, and Q. Xie, "A survey of slam research based on lidar sensors," *International Journal of Sensors*, vol. 1, no. 1, p. 1003, 2019.

[46] W. Zhao, T. He, A. Y. M. Sani, and T. Yao, "Review of slam techniques for autonomous underwater vehicles," in *Proceedings of the 2019 International Conference on Robotics, Intelligent Control and Artificial Intelligence*, 2019, pp. 384–389.

[47] C. Chen, B. Wang, C. X. Lu, N. Trigoni, and A. Markham, "A survey on deep learning for localization and mapping: Towards the age of spatial machine intelligence," *arXiv preprint arXiv:2006.12567*, 2020.

[48] W. Chen, G. Shang, A. Ji, C. Zhou, X. Wang, C. Xu, Z. Li, and K. Hu, "An overview on visual slam: From tradition to semantic," *Remote Sensing*, vol. 14, no. 13, p. 3010, 2022.

[49] I. A. Kazerouni, L. Fitzgerald, G. Dooly, and D. Toal, "A survey of state-of-the-art on visual slam," *Expert Systems with Applications*, vol. 205, p. 117734, 2022.

[50] Y. Tang, C. Zhao, J. Wang, C. Zhang, Q. Sun, W. X. Zheng, W. Du, F. Qian, and J. Kurths, "Perception and navigation in autonomous systems in the era of learning: A survey," *IEEE Transactions on Neural Networks and Learning Systems*, 2022.

[51] M. Zollhöfer, P. Stotko, A. Görlitz, C. Theobalt, M. Nießner, R. Klein, and A. Kolb, "State of the art on 3d reconstruction with rgb-d cameras," in *Computer graphics forum*, vol. 37, no. 2. Wiley Online Library, 2018, pp. 625–652.

[52] J. A. Placed, J. Strader, H. Carrillo, N. Atanasov, V. Indelman, L. Carlone, and J. A. Castellanos, "A survey on active simultaneous localization and mapping: State of the art and new frontiers," *IEEE Transactions on Robotics*, 2023.

[53] A. Dai, M. Nießner, M. Zollhöfer, S. Izadi, and C. Theobalt, "Bundlefusion: Real-time globally consistent 3d reconstruction using on-the-fly surface reintegration," vol. 36, no. 4, p. 1, 2017.

[54] T. Müller, A. Evans, C. Schied, and A. Keller, "Instant neural graphics primitives with a multiresolution hash encoding," *ACM Trans. Graph.*, vol. 41, no. 4, pp. 102:1–102:15, Jul. 2022. [Online]. Available: https://doi.org/10.1145/3528223.3530127

[55] Z. Li, T. Müller, A. Evans, R. H. Taylor, M. Unberath, M.-Y. Liu, and C.-H. Lin, "Neuralangelo: High-fidelity neural surface reconstruction," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 8456–8465.

[56] S. Peng, M. Niemeyer, L. Mescheder, M. Pollefeys, and A. Geiger, "Convolutional occupancy networks," in *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part III 16*. Springer, 2020, pp. 523–540.

[57] Q. Xu, Z. Xu, J. Philip, S. Bi, Z. Shu, K. Sunkavalli, and U. Neumann, "Point-nerf: Point-based neural radiance fields," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 5438–5448.

[58] J. L. Schönberger and J.-M. Frahm, "Structure-from-motion revisited," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

[59] J. L. Schönberger, E. Zheng, M. Pollefeys, and J.-M. Frahm, "Pixelwise view selection for unstructured multi-view stereo," in *European Conference on Computer Vision (ECCV)*, 2016.

[60] K. Deng, A. Liu, J.-Y. Zhu, and D. Ramanan, "Depth-supervised nerf: Fewer views and faster training for free," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 12 882–12 891.

[61] B. Roessle, J. T. Barron, B. Mildenhall, P. P. Srinivasan, and M. Nießner, "Dense depth priors for neural radiance fields from sparse input views," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 12 892–12 901.

[62] M. Niemeyer, J. T. Barron, B. Mildenhall, M. S. Sajjadi, A. Geiger, and N. Radwan, "Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 5480–5490.

[63] K. Gao, Y. Gao, H. He, D. Lu, L. Xu, and J. Li, "Nerf: Neural radiance field in 3d vision, a comprehensive review," *arXiv preprint arXiv:2210.00379*, 2022.

[64] A. S. A. Rabby and C. Zhang, "Beyondpixels: A comprehensive review of the evolution of neural radiance fields," *arXiv e-prints*, pp. arXiv–2306, 2023. [Online]. Available: https://arxiv.org/abs/2306.03000

[65] A. Tewari, J. Thies, B. Mildenhall, P. Srinivasan, E. Tretschk, W. Yifan, C. Lassner, V. Sitzmann, R. Martin-Brualla, S. Lombardi *et al.*, "Advances in neural rendering," in *Computer Graphics Forum*, vol. 41, no. 2. Wiley Online Library, 2022, pp. 703–735.

[66] S. Fridovich-Keil, A. Yu, M. Tancik, Q. Chen, B. Recht, and A. Kanazawa, "Plenoxels: Radiance fields without neural networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 5501–5510.

[67] L. Yen-Chen, P. Florence, J. T. Barron, A. Rodriguez, P. Isola, and T.-Y. Lin, "inerf: Inverting neural radiance fields for pose estimation. in 2021 ieee," in *RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 1323–1330.

[68] Z. Wang, S. Wu, W. Xie, M. Chen, and V. A. Prisacariu, "Nerf–: Neural radiance fields without known camera parameters," *arXiv preprint arXiv:2102.07064*, 2021.

[69] C.-H. Lin, W.-C. Ma, A. Torralba, and S. Lucey, "Barf: Bundle-adjusting neural radiance fields," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 5741–5751.

[70] W. E. Lorensen and H. E. Cline, "Marching cubes: A high resolution 3d surface construction algorithm," in *Seminal graphics: pioneering efforts that shaped the field*, 1998, pp. 347–353.

[71] P. Wang, L. Liu, Y. Liu, C. Theobalt, T. Komura, and W. Wang, "Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction," *arXiv preprint arXiv:2106.10689*, 2021.

[72] D. Azinović, R. Martin-Brualla, D. B. Goldman, M. Nießner, and J. Thies, "Neural rgb-d surface reconstruction," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 6290–6301.

[73] G. Chen and W. Wang, "A survey on 3d gaussian splatting," *arXiv preprint arXiv:2401.03890*, 2024.

[74] T. Wu, Y.-J. Yuan, L.-X. Zhang, J. Yang, Y.-P. Cao, L.-Q. Yan, and L. Gao, "Recent advances in 3d gaussian splatting," *arXiv preprint arXiv:2403.11134*, 2024.

[75] B. Fei, J. Xu, R. Zhang, Q. Zhou, W. Yang, and Y. He, "3d gaussian as a new vision era: A survey," *arXiv preprint arXiv:2402.07181*, 2024.

[76] J. Sturm, N. Engelhard, F. Endres, W. Burgard, and D. Cremers, "A benchmark for the evaluation of rgb-d slam systems," in *2012 IEEE/RSJ international conference on intelligent robots and systems*. IEEE, 2012, pp. 573–580.

[77] A. Dai, A. X. Chang, M. Savva, M. Halber, T. Funkhouser, and M. Nießner, "Scannet: Richly-annotated 3d reconstructions of indoor scenes," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 5828–5839.

[78] J. Straub, T. Whelan, L. Ma, Y. Chen, E. Wijmans, S. Green, J. J. Engel, R. Mur-Artal, C. Ren, S. Verma, A. Clarkson, M. Yan, B. Budge, Y. Yan, X. Pan, J. Yon, Y. Zou, K. Leon, N. Carter, J. Briales, T. Gillingham, E. Mueggler, L. Pesqueira, M. Savva, D. Batra, H. M. Strasdat, R. D. Nardi, M. Goesele, S. Lovegrove,

and R. Newcombe, "The Replica dataset: A digital replica of indoor spaces," *arXiv preprint arXiv:1906.05797*, 2019.

[79] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012.

[80] M. Ramezani, Y. Wang, M. Camurri, D. Wisth, M. Mattamala, and M. Fallon, "The newer college dataset: Handheld lidar, inertial and vision with ground truth," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2020, pp. 4353–4360.

[81] M. Burri, J. Nikolic, P. Gohl, T. Schneider, J. Rehder, S. Omari, M. W. Achtelik, and R. Siegwart, "The euroc micro aerial vehicle datasets," *The International Journal of Robotics Research*, vol. 35, no. 10, pp. 1157–1163, 2016.

[82] B. Glocker, S. Izadi, J. Shotton, and A. Criminisi, "Real-time rgb-d camera relocalization," in *2013 IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*. IEEE, 2013, pp. 173–179.

[83] C. Yeshwanth, Y.-C. Liu, M. Nießner, and A. Dai, "Scannet++: A high-fidelity dataset of 3d indoor scenes," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 12–22.

[84] M. Denninger, M. Sundermeyer, D. Winkelbauer, Y. Zidan, D. Olefir, M. Elbadrawy, A. Lodhi, and H. Katam, "Blenderproc," *arXiv preprint arXiv:1911.01911*, 2019.

[85] Y. Liu, Y. Fu, F. Chen, B. Goossens, W. Tao, and H. Zhao, "Simultaneous localization and mapping related datasets: A comprehensive survey," *arXiv preprint arXiv:2102.04036*, 2021.

[86] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE transactions on image processing*, vol. 13, no. 4, pp. 600–612, 2004.

[87] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 586–595.

[88] X. Yang, H. Li, H. Zhai, Y. Ming, Y. Liu, and G. Zhang, "Voxfusion: Dense tracking and mapping with voxel-based neural implicit representation," in *2022 IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*. IEEE, 2022, pp. 499–507.

[89] M. M. Johari, C. Carta, and F. Fleuret, "Eslam: Efficient dense slam system based on hybrid representation of signed distance fields," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2023, pp. 17408–17419.

[90] H. Wang, J. Wang, and L. 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)*, June 2023, pp. 13293–13302.

[91] Y. Zhang, F. Tosi, S. Mattoccia, and M. Poggi, "Go-slam: Global optimization for consistent 3d instant reconstruction," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 3727–3737.

[92] Z. Teed and J. Deng, "Droid-slam: Deep visual slam for monocular, stereo, and rgb-d cameras," *NeurIPS*, vol. 34, pp. 16558–16569, 2021.

[93] E. Sandström, Y. Li, L. Van Gool, and M. R. Oswald, "Point-slam: Dense neural point cloud-based slam," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023.

[94] L. Xinyang, L. Yijin, T. Yanbin, B. Hujun, Z. Guofeng, Z. Yinda, and C. Zhaopeng, "Multi-modal neural radiance field for monocular dense slam with a light-weight tof sensor," in *International Conference on Computer Vision (ICCV)*, 2023.

[95] P. Hu and Z. Han, "Learning neural implicit through volume rendering with attentive depth fusion priors," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.

[96] M. Li, J. He, Y. Wang, and H. Wang, "End-to-end rgb-d slam with multi-mlps dense neural implicit representations," *IEEE Robotics and Automation Letters*, 2023.

[97] A. L. Teigen, Y. Park, A. Stahl, and R. Mester, "Rgb-d mapping and tracking in a plenoxel radiance field," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2024, pp. 3342–3351.

[98] H. Wang, Y. Cao, X. Wei, Y. Shou, L. Shen, Z. Xu, and K. Ren, "Structerf-slam: Neural implicit representation slam for structural environments," *Computers & Graphics*, p. 103893, 2024.

[99] R. Mur-Artal and J. D. Tardós, "Orb-slam2: An open-source slam system for monocular, stereo, and rgb-d cameras," *IEEE transactions on robotics*, vol. 33, no. 5, pp. 1255–1262, 2017.

[100] P. F. Felzenszwalb and D. P. Huttenlocher, "Efficient graph-based image segmentation," *International journal of computer vision*, vol. 59, pp. 167–181, 2004.

[101] X. Wu, Z. Liu, Y. Tian, Z. Liu, and W. Chen, "Kn-slam: Keypoints and neural implicit encoding slam," *IEEE Transactions on Instrumentation and Measurement*, vol. 73, pp. 1–12, 2024.

[102] V. Lepetit, F. Moreno-Noguer, and P. Fua, "Ep n p: An accurate o (n) solution to the p n p problem," *International journal of computer vision*, vol. 81, pp. 155–166, 2009.

[103] P.-E. Sarlin, C. Cadena, R. Siegwart, and M. Dymczyk, "From coarse to fine: Robust hierarchical localization at large scale," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 12716–12725.

[104] Y. Ming, W. Ye, and A. Calway, "idf-slam: End-to-end rgb-d slam with neural implicit mapping and deep feature tracking," *arXiv preprint arXiv:2209.07919*, 2022.

[105] M. El Banani, L. Gao, and J. Johnson, "Unsupervisedr&r: Unsupervised point cloud registration via differentiable rendering," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 7129–7139.

[106] W. Guo, B. Wang, and L. Chen, "Neuv-slam: Fast neural multiresolution voxel optimization for rgbd dense slam," 2024.

[107] T. Deng, Y. Wang, H. Xie, H. Wang, J. Wang, D. Wang, and W. Chen, "Neslam: Neural implicit mapping and self-supervised feature tracking with depth completion and denoising," *arXiv preprint arXiv:2403.20034*, 2024.

[108] D. DeTone, T. Malisiewicz, and A. Rabinovich, "Superpoint: Self-supervised interest point detection and description," in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2018, pp. 224–236.

[109] H. Matsuki, R. Murai, P. H. J. Kelly, and A. J. Davison, "Gaussian Splatting SLAM," 2024.

[110] H. Huang, L. Li, H. Cheng, and S.-K. Yeung, "Photo-slam: Real-time simultaneous localization and photorealistic mapping for monocular, stereo, and rgb-d cameras," *arXiv preprint arXiv:2311.16728*, 2023.

[111] C. Campos, R. Elvira, J. J. G. Rodríguez, J. M. Montiel, and J. D. Tardós, "Orb-slam3: An accurate open-source library for visual, visual–inertial, and multimap slam," *IEEE Transactions on Robotics*, vol. 37, no. 6, pp. 1874–1890, 2021.

[112] N. Keetha, J. Karhade, K. M. Jatavallabhula, G. Yang, S. Scherer, D. Ramanan, and J. Luiten, "Splatam: Splat, track & map 3d gaussians for dense rgb-d slam," *arXiv preprint arXiv:2312.02126*, 2023.

[113] V. Yugay, Y. Li, T. Gevers, and M. R. Oswald, "Gaussian-slam: Photo-realistic dense slam with gaussian splatting," *arXiv preprint arXiv:2312.10070*, 2023.

[114] T. Deng, Y. Chen, L. Zhang, J. Yang, S. Yuan, D. Wang, and W. Chen, "Compact 3d gaussian splatting for dense visual slam," *arXiv preprint arXiv:2403.11247*, 2024.

[115] S. Ha, J. Yeon, and H. Yu, "Rgbd gs-icp slam," *arXiv preprint arXiv:2403.12550*, 2024.

[116] A. Segal, D. Haehnel, and S. Thrun, "Generalized-icp." in *Robotics: science and systems*, vol. 2, no. 4. Seattle, WA, 2009, p. 435.

[117] S. Sun, M. Mielle, A. J. Lilienthal, and M. Magnusson, "High-fidelity slam using gaussian splatting with rendering-guided densification and regularized optimization," *arXiv preprint arXiv:2403.12535*, 2024.

[118] J. Hu, X. Chen, B. Feng, G. Li, L. Yang, H. Bao, G. Zhang, and Z. Cui, "Cg-slam: Efficient dense rgb-d slam in a consistent uncertainty-aware 3d gaussian field," *arXiv preprint arXiv:2403.16095*, 2024.

[119] R. Arandjelovic, P. Gronat, A. Torii, T. Pajdla, and J. Sivic, "Netvlad: Cnn architecture for weakly supervised place recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 5297–5307.

[120] Lisong C. Sun, Neel Bhatt, Jonathan C. Liu, Zhiwen Fan, Zhangyang Wang, Todd E. Humphreys, and Ufuk Topcu, "Mm3dgs slam: Multi-modal 3d gaussian splatting for slam using vision, depth, and inertial measurements," 2024.

[121] R. Ranftl, A. Bochkovskiy, and V. Koltun, "Vision transformers for dense prediction," *ICCV*, 2021.

[122] J. Park, Q.-Y. Zhou, and V. Koltun, "Colored point cloud registration revisited," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 143–152.

[123] J. Hu, M. Mao, H. Bao, G. Zhang, and Z. Cui, "CP-SLAM: Collaborative neural point-based SLAM system," in *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. [Online]. Available: https://openreview.net/forum?id=dFSeZm6dTC

[124] B. Xiang, Y. Sun, Z. Xie, X. Yang, and Y. Wang, "Nisb-map: Scalable mapping with neural implicit spatial block," *IEEE Robotics and Automation Letters*, 2023.

[125] S. Liu and J. Zhu, "Efficient map fusion for multiple implicit slam agents," *IEEE Transactions on Intelligent Vehicles*, 2023.

[126] Y. Tang, J. Zhang, Z. Yu, H. Wang, and K. Xu, "Mips-fusion: Multi-implicit-submaps for scalable and robust online neural rgb-d reconstruction," *ACM Transactions on Graphics (TOG)*, vol. 42, no. 6, pp. 1–16, 2023.

[127] Y. Mao, X. Yu, K. Wang, Y. Wang, R. Xiong, and Y. Liao, "Ngel-slam: Neural implicit representation-based global consistent low-latency slam system," *arXiv preprint arXiv:2311.09525*, 2023.

[128] T. Deng, G. Shen, T. Qin, J. Wang, W. Zhao, J. Wang, D. Wang, and W. Chen, "Plgslam: Progressive neural scene represenation with local to global bundle adjustment," *arXiv preprint arXiv:2312.09866*, 2023.

[129] L. Liso, E. Sandström, V. Yugay, L. V. Gool, and M. R. Oswald, "Loopy-slam: Dense neural slam with loop closures," 2024.

[130] H. Matsuki, K. Tateno, M. Niemeyer, and F. Tombari, "Newton: Neural view-centric mapping for on-the-fly large-scale slam," *IEEE Robotics and Automation Letters*, 2024.

[131] H. Zhai, H. Li, X. Yang, G. Huang, Y. Ming, H. Bao, and G. Zhang, "Vox-fusion++: Voxel-based neural implicit dense tracking and mapping with multi-maps," *arXiv preprint arXiv:2403.12536*, 2024.

[132] Y. Yan, R. He, and Z. Liu, "Mute-slam: Real-time neural slam with multiple tri-plane hash representations," *arXiv preprint arXiv:2403.17765*, 2024.

[133] K. Mazur, E. Sucar, and A. J. Davison, "Feature-realistic neural fusion for real-time, open set scene understanding," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023, pp. 8201–8207.

[134] M. Tan and Q. Le, "Efficientnet: Rethinking model scaling for convolutional neural networks," in *International conference on machine learning*. PMLR, 2019, pp. 6105–6114.

[135] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin, "Emerging properties in self-supervised vision transformers," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 9650–9660.

[136] X. Kong, S. Liu, M. Taher, and A. J. Davison, "vmap: Vectorised object mapping for neural field slam," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2023, pp. 952–961.

[137] S. Zhu, G. Wang, H. Blum, J. Liu, L. Song, M. Pollefeys, and H. Wang, "Sni-slam: Semantic neural implicit slam," *arXiv preprint arXiv:2311.11016*, 2023.

[138] M. Oquab, T. Darcet, T. Moutakanni, H. Vo, M. Szafraniec, V. Khalidov, P. Fernandez, D. Haziza, F. Massa, A. El-Nouby *et al.*, "Dinov2: Learning robust visual features without supervision," *arXiv preprint arXiv:2304.07193*, 2023.

[139] Y. Haghighi, S. Kumar, J. P. Thiran, and L. Van Gool, "Neural implicit dense semantic slam," *arXiv preprint arXiv:2304.14560*, 2023.

[140] B. Cheng, A. Schwing, and A. Kirillov, "Per-pixel classification is not all you need for semantic segmentation," *Advances in neural information processing systems*, vol. 34, pp. 17 864–17 875, 2021.

[141] K. Li, M. Niemeyer, N. Navab, and F. Tombari, "Dns slam: Dense neural semantic-informed slam," *arXiv preprint arXiv:2312.00204*, 2023.

[142] M. Li, S. Liu, and H. Zhou, "Sgs-slam: Semantic gaussian splatting for neural dense slam," 2024.

[143] S. Zhu, R. Qin, G. Wang, J. Liu, and H. Wang, "Semgauss-slam: Dense semantic gaussian splatting slam," *arXiv preprint arXiv:2403.07494*, 2024.

[144] Y. Ji, Y. Liu, G. Xie, B. Ma, and Z. Xie, "Neds-slam: A novel neural explicit dense semantic slam framework using 3d gaussian splatting," *arXiv preprint arXiv:2403.11679*, 2024.

[145] L. Yang, B. Kang, Z. Huang, X. Xu, J. Feng, and H. Zhao, "Depth anything: Unleashing the power of large-scale unlabeled data," *arXiv preprint arXiv:2401.10891*, 2024.

[146] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo *et al.*, "Segment anything," *arXiv preprint arXiv:2304.02643*, 2023.

[147] M. A. Karaoglu, H. Schieber, N. Schischka, M. Görgülü, F. Grötzner, A. Ladikos, D. Roth, N. Navab, and B. Busam, "Dynamon: Motion-aware fast and robust camera localization for dynamic nerf," *arXiv preprint arXiv:2309.08927*, 2023.

[148] L.-C. Chen, G. Papandreou, F. Schroff, and H. Adam, "Rethinking atrous convolution for semantic image segmentation," *arXiv preprint arXiv:1706.05587*, 2017.

[149] S. F. Bhat, R. Birkl, D. Wofk, P. Wonka, and M. Müller, "Zoedepth: Zero-shot transfer by combining relative and metric depth," *arXiv preprint arXiv:2302.12288*, 2023.

[150] Z. Xu, J. Niu, Q. Li, T. Ren, and C. Chen, "Nid-slam: Neural implicit representation-based rgb-d slam in dynamic environments," *arXiv preprint arXiv:2401.01189*, 2024.

[151] W. Wu, G. Wang, T. Deng, S. Aegidius, S. Shanks, V. Modugno, D. Kanoulas, and H. Wang, "Dvn-slam: Dynamic visual neural slam based on local-global encoding," *arXiv preprint arXiv:2403.11776*, 2024.

[152] D. Lisus, C. Holmes, and S. Waslander, "Towards open world nerf-based slam," in *2023 20th Conference on Robots and Vision (CRV)*, 2023, pp. 37–44.

[153] C.-M. Chung, Y.-C. Tseng, Y.-C. Hsu, X.-Q. Shi, Y.-H. Hua, J.-F. Yeh, W.-C. Chen, Y.-T. Chen, and W. H. Hsu, "Orbeez-slam: A real-time monocular visual slam with orb features and nerf-realized mapping," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023, pp. 9400–9406.

[154] T. Hua, H. Bai, Z. Cao, and L. Wang, "Fmapping: Factorized efficient neural field mapping for real-time dense rgb slam," *arXiv preprint arXiv:2306.00579*, 2023.

[155] J. Lin, A. Nachkov, S. Peng, L. Van Gool, and D. P. Paudel, "Ternary-type opacity and hybrid odometry for rgb-only nerf-slam," *arXiv preprint arXiv:2312.13332*, 2023.

[156] H. Matsuki, E. Sucar, T. Laidow, K. Wada, R. Scona, and A. J. Davison, "imode: Real-time incremental monocular dense mapping using neural field," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023, pp. 4171–4177.

[157] F. Ma and S. Karaman, "Sparse-to-dense: Depth prediction from sparse depth samples and a single image," 2018.

[158] W. Zhang, T. Sun, S. Wang, Q. Cheng, and N. Haala, "Hi-slam: Monocular real-time dense mapping with hybrid implicit fields," *IEEE Robotics and Automation Letters*, 2023.

[159] A. Eftekhar, A. Sax, J. Malik, and A. 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*, 2021, pp. 10 786–10 796.

[160] H. Xu, J. Zhang, J. Cai, H. Rezatofighi, and D. Tao, "Gmflow: Learning optical flow via global matching," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 8121–8130.

[161] J. Naumann, B. Xu, S. Leutenegger, and X. Zuo, "Nerf-vo: Real-time sparse visual odometry with neural radiance fields," *arXiv preprint arXiv:2312.13471*, 2023.

[162] Z. Teed, L. Lipson, and J. Deng, "Deep patch visual odometry," *arXiv preprint arXiv:2208.04726*, 2022.

[163] H. Zhou, Z. Guo, S. Liu, L. Zhang, Q. Wang, Y. Ren, and M. Li, "Mod-slam: Monocular dense mapping for unbounded 3d scene reconstruction," 2024.

[164] C. Peng, C. Xu, Y. Wang, M. Ding, H. Yang, M. Tomizuka, K. Keutzer, M. Pavone, and W. Zhan, "Q-slam: Quadric representations for monocular slam," *arXiv preprint arXiv:2403.08125*, 2024.

[165] G. Zhang, E. Sandström, Y. Zhang, M. Patel, L. Van Gool, and M. R. Oswald, "Glorie-slam: Globally optimized rgb-only implicit encoding point cloud slam," *arXiv preprint arXiv:2403.19549*, 2024.

[166] X. Han, H. Liu, Y. Ding, and L. Yang, "Ro-map: Real-time multi-object mapping with neural radiance fields," *IEEE Robotics and Automation Letters*, vol. 8, no. 9, pp. 5950–5957, 2023.

[167] A. Rosinol, J. J. Leonard, and L. Carlone, "Nerf-slam: Real-time dense monocular slam with neural radiance fields," in *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2023, pp. 3437–3444.

[168] S. Isaacson, P.-C. Kung, M. Ramanagopal, R. Vasudevan, and K. A. Skinner, "Loner: Lidar only neural representations for real-time slam," *IEEE Robotics and Automation Letters*, 2023.

[169] S. Rusinkiewicz and M. Levoy, "Efficient variants of the icp algorithm," in *Proceedings third international conference on 3-D digital imaging and modeling*. IEEE, 2001, pp. 145–152.

[170] Y. Pan, X. Zhong, L. Wiesmann, T. Posewsky, J. Behley, and C. Stachniss, "Pin-slam: Lidar slam using a point-based implicit neural representation for achieving global map consistency," *arXiv preprint arXiv:2401.09101*, 2024.

[171] S. Hong, J. He, X. Zheng, H. Wang, H. Fang, K. Liu, C. Zheng, and S. Shen, "Liv-gaussmap: Lidar-inertial-visual fusion for real-time 3d radiance field map rendering," *arXiv preprint arXiv:2401.14857*, 2024.

[172] C. Wu, Y. Duan, X. Zhang, Y. Sheng, J. Ji, and Y. Zhang, "Mm-gaussian: 3d gaussian-based multi-modal fusion for localization and reconstruction in unbounded scenes," *arXiv preprint arXiv:2404.04026*, 2024.

[173] I. Vizzo, T. Guadagnino, B. Mersch, L. Wiesmann, J. Behley, and C. Stachniss, "Kiss-icp: In defense of point-to-point icp–simple, accurate, and robust registration if done the right way," *IEEE Robotics and Automation Letters*, vol. 8, no. 2, pp. 1029–1036, 2023.

[174] P. Lindenberger, P. Sarlin, and M. Pollefeys, "Lightglue: Local feature matching at light speed. arxiv 2023," *arXiv preprint arXiv:2306.13643*.

[175] T. Müller, B. McWilliams, F. Rousselle, M. Gross, and J. Novák, "Neural importance sampling," *ACM Transactions on Graphics (ToG)*, vol. 38, no. 5, pp. 1–19, 2019.

[176] E. Rublee, V. Rabaud, K. Konolige, and G. Bradski, "Orb: An efficient alternative to sift or surf," in *2011 International conference on computer vision*. Ieee, 2011, pp. 2564–2571.

[177] B. Cheng, I. Misra, A. G. Schwing, A. Kirillov, and R. Girdhar, "Masked-attention mask transformer for universal image segmentation," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 1290–1299.

[178] A. Cao and J. Johnson, "Hexplane: A fast representation for dynamic scenes," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 130–141.

[179] R. Mur-Artal and J. D. Tardós, "Probabilistic semi-dense mapping from highly accurate feature-based monocular slam." in *Robotics: Science and Systems*, vol. 2015. Rome, 2015.

[180] A. Gropp, L. Yariv, N. Haim, M. Atzmon, and Y. Lipman, "Implicit geometric regularization for learning shapes," *arXiv preprint arXiv:2002.10099*, 2020.

[181] R. Ranftl, K. Lasinger, D. Hafner, K. Schindler, and V. Koltun, "Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 3, 2022.

[182] M. Tancik, E. Weber, E. Ng, R. Li, B. Yi, T. Wang, A. Kristoffersen, J. Austin, K. Salahi, A. Ahuja *et al.*, "Nerfstudio: A modular framework for neural radiance field development," in *ACM SIGGRAPH 2023 Conference Proceedings*, 2023, pp. 1–12.

[183] J. Lin, "Divergence measures based on the shannon entropy," *IEEE Transactions on Information theory*, vol. 37, no. 1, pp. 145–151, 1991.

[184] C. Yan, D. Qu, D. Wang, D. Xu, Z. Wang, B. Zhao, and X. Li, "Gs-slam: Dense visual slam with 3d gaussian splatting," *arXiv preprint arXiv:2311.11700*, 2023.

[185] A. Zeng, S. Song, M. Nießner, M. Fisher, J. Xiao, and T. Funkhouser, "3dmatch: Learning local geometric descriptors from rgb-d reconstructions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1802–1811.

[186] Y. Pan, P. Xiao, Y. He, Z. Shao, and Z. Li, "Mulls: Versatile lidar slam via multi-metric linear least square," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2021, pp. 11 633–11 640.

[187] P. Dellenbach, J.-E. Deschaud, B. Jacquet, and F. Goulette, "Ct-icp: Real-time elastic lidar odometry with loop closure," in *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022, pp. 5580–5586.

[188] M. Yokozuka, K. Koide, S. Oishi, and A. Banno, "Litamin2: Ultra light lidar-based slam using geometric approximation applied with kl-divergence," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2021, pp. 11 619–11 625.

[189] J. Behley and C. Stachniss, "Efficient surfel-based slam using 3d laser range data in urban environments." in *Robotics: Science and Systems*, vol. 2018, 2018, p. 59.

[190] J. Ruan, B. Li, Y. Wang, and Y. Sun, "Slamesh: Real-time lidar simultaneous localization and meshing," *arXiv preprint arXiv:2303.05252*, 2023.

[191] X. Liu, Z. Liu, F. Kong, and F. Zhang, "Large-scale lidar consistent mapping using hierarchical lidar bundle adjustment," *IEEE Robotics and Automation Letters*, vol. 8, no. 3, pp. 1523–1530, 2023.

[192] A. Guédon and V. Lepetit, "Sugar: Surface-aligned gaussian splatting for efficient 3d mesh reconstruction and high-quality mesh rendering," *arXiv preprint arXiv:2311.12775*, 2023.

[193] X. Xu, T. Zhang, S. Wang, X. Li, Y. Chen, Y. Li, B. Raj, M. Johnson-Roberson, and X. Huang, "Customizable perturbation synthesis for robust slam benchmarking," *arXiv preprint arXiv:2402.08125*, 2024.

[194] T. Hua and L. Wang, "Benchmarking implicit neural representation and geometric rendering in real-time rgb-d slam," *arXiv preprint arXiv:2403.19473*, 2024.