# ActiveRMAP: Radiance Field for Active Mapping And Planning

Huangying Zhan[1,2], Jiyang Zheng[3,4], Yi Xu[2], Ian Reid[1], Hamid Rezatofighi[3]

[1]The University of Adelaide [2]InnoPeak Technology, Inc.

[3]Monash University [4]Australian National University

## Abstract

*A high-quality 3D reconstruction of a scene from a collection of 2D images can be achieved through offline/online mapping methods. In this paper, we explore active mapping from the perspective of implicit representations, which have recently produced compelling results in a variety of applications. One of the most popular implicit representations – Neural Radiance Field (NeRF), first demonstrated photorealistic rendering results using multi-layer perceptrons, with promising offline 3D reconstruction as a by-product of the radiance field. More recently, researchers also applied this implicit representation for online reconstruction and localization (i.e. implicit SLAM systems). However, the study on using implicit representation for active vision tasks is still very limited. In this paper, we are particularly interested in applying the neural radiance field for active mapping and planning problems, which are closely coupled tasks in an active system. We, for the first time, present an RGB-only active vision framework using radiance field representation for active 3D reconstruction and planning in an online manner. Specifically, we formulate this joint task as an iterative dual-stage optimization problem, where we alternatively optimize for the radiance field representation and path planning. Experimental results suggest that the proposed method achieves competitive results compared to other offline methods and outperforms active reconstruction methods using NeRFs.*

## 1. Introduction

One of the remarkable successes of computer vision research has been to show that it is possible to generate a high-quality 3D reconstruction of a scene from a collection of 2D images or video of the scene. When done offline in batch, this is usually termed structure from motion. When performed *online*, with the 3D reconstruction built incrementally as more visual data is perceived, it is usually termed *Visual Simultaneous Localisation and Mapping* (vSLAM), in which the 3D reconstruction is the map, and the camera poses are localised with respect to the map. vSLAM
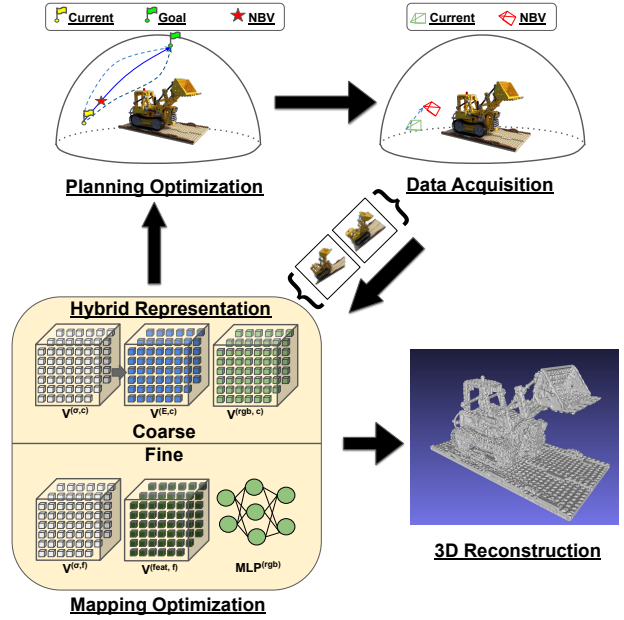


Figure 1. **Overview**. We propose a dual-stage optimization framework for active vision tasks. A hybrid implicit-explicit radiance field is used as the representation of the environment for its speed efficient and remarkable performance on 3D shape representation.

systems have often been researched and developed with the map as the end product. However in many robotic applications the map is simply a means to perform some other tasks such as planning and navigation; the joint tasks of localisation, mapping, planning, and navigation are collectively known as *Active SLAM*. In this paper we explore Active SLAM from the perspective of powerful, learned, implicit representations of geometry (*i.e.* , NeRF), developing methods that can plan and execute camera paths to optimise the fidelity of the NeRF representation.

Implicit neural representations, especially Neural Radiance Fields (NeRF), have been recently applied in a variety of visual geometric applications, including 3D object reconstruction [43], novel view rendering [37, 47, 67, 69], surface reconstruction [3, 28], and generative models [41, 53]. While much of this work assumes well-posed cameras,

1

some work has extended the idea to more general tasks of structure from motion [10, 29, 64] and SLAM [58, 60]. The implicit representation is very attractive as an alternative to point clouds or surface meshes for a variety of reasons, including efficient memory usage, continuous and differentiable representation, and natural ways to leverage prior learned information (see, e.g., PixelNeRF [67]). Nevertheless, there remain questions and challenges associated with using this representation within an active or robotic paradigm. To the best of our knowledge, prior work to exploit NeRF for path planning is limited [1], and active reconstruction using NeRF has been considered by three contemporaneous works [27, 42, 48]. However, [1] assumes a NeRF model of the scene is pre-trained offline using a collection of images from the environment. [27] has to initialize a coarse NeRF representation by pre-training on a small set of collected images. The computation time is very slow, which is not desirable for robotic applications. [42, 48] predict the image reconstruction uncertainty as the criterion for next-best-view selection. However, image reconstruction quality is not always a good proxy for geometric reconstruction fidelity. Moreover, this method requires additional depth supervision to achieve good reconstruction.

To overcome the above-mentioned drawbacks, we propose an iterative dual-stage optimization framework (Fig. 1) for active mapping and planning, which operates in an online fashion, *without* pre-training a NeRF model. Specifically, the mapping stage involves the optimization of radiance field representation given a set of observed images to date; the planning stage involves global and local policies for next-best view selection with a multi-objective function that considers obstacle avoidance, geometric information gain, and path efficiency. We make the following contributions in creating this framework:

- We propose an iterative dual-stage optimization framework for active reconstruction and planning, which (1) operates in an *online* fashion, without the need of pre-training a NeRF model; (2) only requires color images as the visual observation, without the need of a depth sensor, which is usually required in active 3D reconstruction methods.
- We propose a new multi-objective loss function for path planning, which allows the agent to plan a path with the consideration of obstacle avoidance and geometric information gain.
- We propose a differentiable geometric information gain criterion that considers termination uncertainty (entropy) for next-best view selection.
- We showcase the use of the proposed framework in active 3D reconstruction. Extensive experiments on synthetic and real-world datasets have been performed to show the effectiveness and efficiency of the proposed method.

## 2. Related Work

**Active Mapping/SLAM:** An autonomous robotic system requires the robot to form a model of the environment, which would require four essential abilities – including localization, mapping, planning, and motion control [55]. While the former involves estimating the robot's state, creating a representation of the environment, *e.g.* a map, planning a path to safely achieve a goal location, the latter aims to control the movements of the robot according to the planned path. Localization, mapping, and planning are usually investigated solely or in combination, which results in multiple research areas, such as visual odometry [51, 68], stereo matching [21, 61], multi-view stereo [31, 54, 66], structure-from-motion (SfM) [52], simultaneous localization and mapping (SLAM) [8, 14, 17, 62], and active localization/mapping/SLAM. While SfM is an *offline* method that reconstructs the 3D structure from a set of collected images, SLAM is an *online* method that incrementally builds the map of an environment while at the same time localizing the robot within the environment. However, SLAM is still a passive method that requires human operators to scan the environment with a sensor.

It is not concerned with planning, which guides the navigation process. In contrast, *active* methods [2, 4, 12] consider planning problems in the loop. For the cases that aim to improve localization, mapping, and SLAM, the problems are referred to as active localization, active mapping, and active SLAM, respectively. Active mapping, a.k.a. active reconstruction or next best view (NBV), is a longstanding problem [11], which aims to search for the optimal movements to create the best possible representation of an environment. With the assumption of perfect localization, this problem has been primarily addressed to reconstruct scenes and objects of interest from multiple viewpoints [15, 22, 26, 35, 44, 45]. In this paper, we are interested in the problem that the map of the environment is unknown. Active localization aims to improve the robot's pose estimation assuming the map is known. We refer readers to [6, 19, 38, 65] for more relevant works. Active SLAM unifies the three active vision problems (localization, mapping, and planning). It allows robots to autonomously perform localization and mapping to reduce the uncertainty of its localization and the representation of the environment. Before [13] coined the term – Active SLAM –, this problem has been investigated under different names, majorly as known as exploration problems [7, 18, 34, 40, 56, 57, 63]. We refer readers to the survey papers [8, 33, 46] for more relevant works, development of the field, and the comparisons between prior works. In this work, we are more interested in two aspects of the active vision problems – mapping and planning.

**Neural Radiance Fields:** Implicit neural representations, especially Neural Radiance Field (NeRF) [37], have
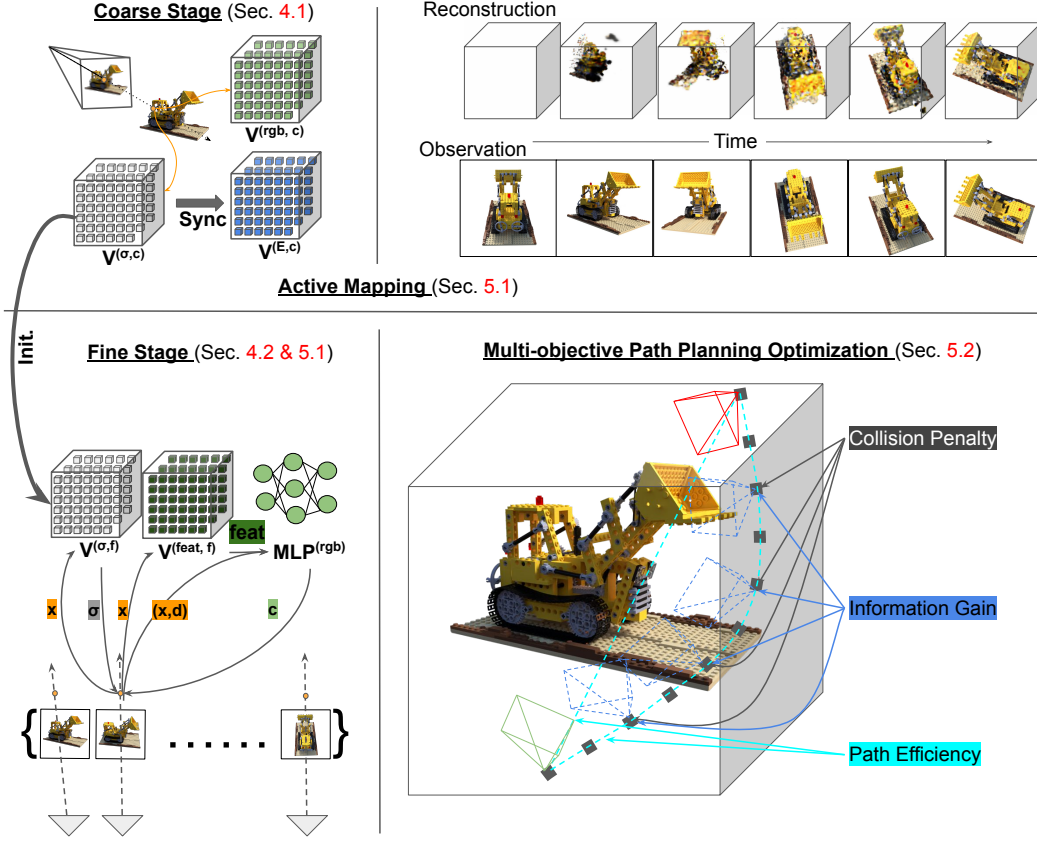
Figure 2. **ActiveRMAP framework**. **Coarse Mapping:** A coarse radiance field representation is optimized with rendering as supervision. **Planning:** Path planning is formulated as an optimization problem that considers collision avoidance and information gain of sampled viewpoints on the predicted trajectory, and path efficiency to reach a goal location. **Iterative framework:** Conditioned on the coarse model, the planning module obtains the next-best view from the optimized trajectory. Thus new observations are acquired incrementally and provided for training the coarse model. **Fine Mapping:** Once the coarse mapping is completed, the optimal coarse model is used for initializing the fine model. Refinement is performed in the Fine Stage.

been proposed as a new representation of the scene. NeRF represents the scene as continuous neural radiance fields using multi-layer perceptrons (MLPs), and it can be simply trained by comparing a set of rendered images with well-posed images. Despite the compelling properties and results demonstrated by NeRFs, training implicit NeRF can be time-consuming and usually takes about a day for a simple scene, due to the nature of volume rendering [23] that requires querying on a significant number of sample points in order to render an image. Though there are some recent works that aim to increase the training and/or inference speed of NeRF [16, 30, 49], these completely implicit representations are still not fast enough for real-time applications. More recently, some works [9, 39, 50, 59] have demonstrated super fast radiance fields with a hybrid representation, *i.e.* implicit for light field and explicit for density field.

NeRF first demonstrated its remarkable ability in novel view rendering [37, 47, 67, 69]. Meanwhile, compelling re-

sults have been demonstrated for 3D object reconstruction [37, 43], surface reconstruction [3, 28], generative models [41, 53], Structure-from-Motion [10, 29, 64], SLAM [58, 60], *etc*. Though implicit representations have been employed in various applications, the study of active vision problems is still very limited.

**NeRF + Active Vision:** To our best knowledge, a prior work [1] shows the use of NeRF for path planning, and three concurrent works [27, 48] show the use of NeRF for active reconstruction. [1] tackles the navigation problem using NeRF representation. It assumes a NeRF model of the scene is pre-trained offline using a collection of images from the environment. Thus, an optimal path is searched in the *NeRF map* by minimizing the collision probability of the path extracted from the NeRF representation. [27, 42, 48] tackle the active mapping problem by finding the next best view (NBV) to optimize the NeRF representation. [27] proposes to use ray-based uncertainty as a criterion for selecting NBV, but a coarse NeRF representation has to be ini-

3

tialized by pre-training on a small set of collected images first. [42,48] predict the image reconstruction uncertainty as the criterion for NBV selection. However, image rendering quality does not directly reflect the geometric reconstruction quality. Moreover, [48] requires additional depth supervision to achieve good reconstruction. More importantly, these methods inherent the speed drawback from implicit NeRF representations. The computation time is very slow, which is not desirable for robotic applications. In this work, we propose an RGB-only active vision framework based on a hybrid implicit-explicit representation, which actively and incrementally builds a representation of the scene without the need for pre-training. It has overcome the drawbacks of the existing approaches.

## 3. Preliminaries

### 3.1. Neural Radiance Field

NeRF [37] was first proposed for solving novel view synthesis. The core of NeRF is the use of multi-layer perceptron networks (MLPs) as an implicit representation of the scene. The MLPs map a 3D position $x$ and a viewing direction $d$ to the corresponding density $\sigma$ and view-dependent color $c$:

$$(\sigma, e) = \text{MLP}_{\Theta_{pos}}^{(pos)}(x); c = \text{MLP}_{\Theta_{rgb}}^{(rgb)}(e, d), \qquad (1)$$

where $\Theta_{pos}$ and $\Theta_{rgb}$ are the learnable MLP parameters, and $e$ is an intermediate embedding to help the shallower $\text{MLP}^{(rgb)}$ to learn the color representation. Given a set of posed images, the training of MLPs relies on the photometric loss between the observed pixel colors $C(r)$ in posed images and the rendered pixel colors $\hat{C}(r)$ in the images viewed from the camera poses , $\mathcal{L}_{photo} = \frac{1}{|R|} \sum_{r \in R} ||C(r) - \hat{C}(r)||$, where $R$ denotes a set of rays in a sampled mini-batch. Volume rendering [23] is employed to render the pixel colors by casting a ray $r$ from the camera center through the pixel. Volume rendering approximates light radiance by integrating the radiance along the ray. $K$ points are sampled along the ray between pre-defined near and far planes. After querying for their densities and colors $(\sigma_i, c_i)_{i=1}^{K}$ via MLPs, the results are accumulated into a single color with the use of quadrature approximation [36] as,

$$\hat{C}(r) = \sum_{i=1}^{K} T_i \alpha_i c_i, \qquad (2)$$

$$\alpha_i = 1 - \exp(-\sigma_i \delta_i)), \qquad (3)$$

$$T_i = \prod_{j=1}^{i-1} (1 - \alpha_j), \qquad (4)$$

where $\alpha_i$ is the light termination probability at the point $i$; $T_i$ is the accumulated transmittance from the near plane to point $i$; $\delta_i$ is the distance between adjacent sampled points.

## 3.2. Direct Voxel Grid Optimization

Training a NeRF with the method presented above can be time-consuming and usually takes hours for a simple scene, due to the querying cost of volume rendering on a significant number of points. More recently, some works [9, 39, 50, 59] have demonstrated super fast radiance fields with a hybrid representation, *i.e.* implicit for color inference and explicit for density inference. With the potential for real-time applications, we choose DVGO [59] among these works as the hybrid representation to use in this paper for the following reasons: (1) a grid-based explicit representation for the density field, which is efficient for computing volumetric information gain in the later stage; (2) the dual-stage design allows a fast convergence in the coarse stage and a high-quality reconstruction in the fine stage. This flexibility in the level of detail suits different levels of active vision task requirements; (3) competitive results and speed when compared to other hybrid representations. However, note that our proposed method can be adapted to the other NeRF methods and is not limited to DVGO only.

We introduce the core of DVGO here and refer readers to [59] for more details. DVGO speeds up NeRF by replacing the density MLP with voxel grids. In the coarse stage, two volumes store densities and colors explicitly, represented by $\mathbf{V}^{(\sigma,c)}$ and $\mathbf{V}^{(rgb,c)}$, respectively. A query of any 3D point $x$ for density $\sigma^{(c)}$ and color $c^{(c)}$ in the coarse stage is efficient with an interpolation:

$$\sigma^{(c)} = \text{softplus}(\ddot{\sigma}^{(c)}) = \log(1 + \exp(\ddot{\sigma}^{(c)} + b) \qquad (5)$$

$$\ddot{\sigma}^{(c)} = \text{interp}(x, \mathbf{V}^{(\sigma,c)}); c^{(c)} = \text{interp}(x, \mathbf{V}^{(rgb,c)}), \quad (6)$$

where shift $b$ is a hyperparameter. The density volume stores a raw density $\ddot{\sigma}^{(c)}$, which is further activated via a softplus operator for the final density. Note that [59] initializes $\mathbf{V}^{(\sigma,c)}$ with zero raw density to prevent suboptimal geometry.

In the fine stage, DVGO refines a higher-resolution density voxel grid $\mathbf{V}^{(\sigma,f)}$ initialized with the optimized coarse geometry $\mathbf{V}^{(\sigma,c)}$. The view-invariant $\mathbf{V}^{(rgb,c)}$ is replaced with a feature voxel grid $\mathbf{V}^{(feat,f)}$ and a shallow MLP for inferring view-dependent color emission. This hybrid implicit-explicit representation provides a good balance between speed and quality. Finally, queries of 3D points $x$ and viewing-direction $d$ at the fine stage are formulated as,

$$\ddot{\sigma}^{(f)} = \text{interp}(x, \mathbf{V}^{(\sigma,f)})) \qquad (7)$$

$$c^{(f)} = \text{MLP}_{\Theta_{rgb}}^{(rgb)} \left( \text{interp}(x, \mathbf{V}^{(feat,f)}, c)), x, d \right). \quad (8)$$

The major supervision for training DVGO models at the coarse and fine stages is the photometric loss. In the following context, we denote the DVGO models at the coarse and fine stage as $\mathbf{M}^{(c)}$ and $\mathbf{M}^{(f)}$ respectively.

# 4. ActiveRMAP

In this section, we present a dual-stage optimization framework, ActiveRMAP (Fig. 2), for active vision tasks, which includes a mapping stage and a planning stage under an assumption of perfect localization and execution. We first present the mapping and planning optimization in Sec. 4.1 and Sec. 4.2. Finally, we introduce the application in Active 3D Reconstruction task in Sec. 4.3.

## 4.1. Mapping Optimization

Differing from most NeRF methods, which collect a set of images of the scene and thus optimize the radiance field *offline*, we create an *online* system that acquires and processes the observed images incrementally.

**Coarse mapping:** Given an observed image $I_t$ at time-$t$, we first store $I_t$ to a database $\{I_i\}$. Since the database size expands as the agent explores the environment, it is not practical to train the model with all the pixels at each iteration. Therefore, we random sample a set of pixels from $\{I_i\}_{i=0}^{t}$ for training $\mathbf{M}_t^{(c)}$ at times-$t$. Note that rendering a pixel requires sampling points along the ray, which is time-consuming. Instead of re-sampling rays at each iteration, we only sample rays once at the beginning of the optimization until a new observation is acquired.

**Fine mapping:** Once the agent has completed exploring the environment, we can perform fine reconstruction if a high-quality reconstruction is required.

## 4.2. Planning Optimization

One of our major contributions is the proposed planning module, which considers collision avoidance, information gain, and path efficiency via a multi-objective optimization formulation. Given the current agent state (pose) $\mathbf{s}_t$, a goal state $\mathbf{s}_g$, and the DVGO coarse model $\mathbf{M}_t^{(c)}$ at times-$t$, we aim to generate a trajectory that yields various task-dependent objectives. The trajectory is defined as a sequence of agent states $\mathbf{S} = (\mathbf{s}_t, \mathbf{s}_{t+1}..., \mathbf{s}_{t+n}, \mathbf{s}_g)$, where $\mathbf{s}_i$ can be lie-algebra or spherical coordinates representation depending on the desirable representation in the task. This problem of finding the optimal trajectory $S^*$ can therefore be expressed as,

$$\mathbf{S}^* = \underset{\mathbf{S}}{\text{argmin}} \sum_{i=t+1}^{t+n} \mathcal{L}_{plan}(\mathbf{M}_t^{(c)}, \mathbf{s}_i) \qquad (9)$$

$$\mathcal{L}_{plan} = \lambda_{cp}\mathcal{L}_{cp} + \lambda_{oc}\mathcal{L}_{ig} + \lambda_{oc}\mathcal{L}_{pe}, \qquad (10)$$

where $\mathcal{L}_{plan}$ is the overall planning objective, $\mathcal{L}_{cp}$ is the collision penalty objective, $\mathcal{L}_{ig}$ is the information gain objective, $\mathcal{L}_{pe}$ is the path efficiency objective, and $[\lambda_{cp}, \lambda_{ig}, \lambda_{pe}]$ are the weighting terms for each objective. Once the trajectory is optimized, we choose the furthest state $\mathbf{s}_i^*$ that fulfills motion constraints as the next state. We consider a maximum translation (0.5m) and a maximum rotation angle (10) as our motion constraints.

**Bézier Trajectory:** Though we aim to find an optimal sequence of agent states, we do not simply make the agent states the optimizable parameters as this approach can cause infeasible agent motions such as teleporting between distanced viewpoints. To solve this issue, we use parametric curves for trajectory formation. Specifically, we adopt the Quadratic Bézier curve to create a set of discrete waypoints that define a smooth, continuous trajectory. For the Quadratic Bézier curve, the discrete points can be defined as $\mathbf{s}_i = (1-r)[(1-r)\mathbf{s}_t + t\mathbf{s}_c] + r[(1-r)\mathbf{s}_c + r\mathbf{s}_g]$, where $r = (i-t)/(n+1)$, and $\mathbf{s}_c$ is an optimizable Bézier control point, which is initialized as the mid-point of $\mathbf{s}_0$ and $\mathbf{s}_g$.

**Collision Penalty:** The density of the environment is represented as a continuous radiance field, which can be used as a means for computing collision penalty. Eq. (3) computes the light termination probability at a point $x$ based on the density at $x$. [1] assumes that a 3D point $x$ with high density (termination probability for light) is a strong proxy for the probability of terminating a mass particle. We can thus minimize the collision penalty by minimizing the density of the sampled waypoints. Thus, the collision penalty is formulated as

$$\mathcal{L}_{cp} = \sum_{i=t+1}^{t+n} l_{cp}(x_i) = \sum_{i=t+1}^{t+n} \exp(\sigma^{(c)}(x_i)), \qquad (11)$$

where $\sigma^{(c)}(x)$ queries the density at 3D point $x$.

**Uncertainty-based Information Gain:** Active mapping requires the evaluation of the quality of un-visited viewpoints – usually termed information gain. A challenge is to quantify the quality and define a meaningful criterion for evaluation, which has not been well studied with the implicit representation. We address this challenge by using termination uncertainty evaluated from each viewpoint. For each viewpoint $\mathbf{s}_i$, we create a frustum originating from the camera center and sample points within the frustum $\mathbf{F}_i$. Applying the same assumption of termination probability, we further compute the uncertainty of termination probability for the sampled points. Therefore, the entropy of a point $E_x(x)$ and a state $E_s(\mathbf{s}_i)$, and the information gain objective are represented by,

$$E_x(x) = -\alpha(x)\log\alpha(x) - \bar{\alpha}(x)\log\bar{\alpha}(x) \qquad (12)$$

$$\mathcal{L}_{ig} = -\sum_{\mathbf{s}_i \in \mathbf{S}} E_s(\mathbf{s}_i) = -\sum_{\mathbf{s}_i \in \mathbf{S}} \left( \sum_{x_i \in \mathbf{F}_i} E_x(x_i) \right) \qquad (13)$$

where $\bar{\alpha}(x) = 1 - \alpha(x)$. Intuitively, the viewpoint with a higher entropy means higher uncertainty about the termination probability of the associated 3D points. The entropy can be decreased once the model obtains more information about the 3D points. Thus it is more certain about the density (thus termination probability) at the sampled points.

Note that although we can directly compute $\alpha(x_i)$ from the density volume, it comes with an initialization issue. DVGO initializes the density volume with zeros to avoid sub-optimal geometry, which means that the scene has nearly zero entropy at the beginning thus it cannot guide the agent to choose next-best view with the highest information gain. To overcome this issue, we introduce an additional entropy volume $\mathbf{V}^{(E,c)}$ represented by $\alpha^{(E,c)}$ in DVGO coarse model $\mathbf{M}^{(c)}$. We initialize $\alpha^{(E,c)}$ with 0.5, which represents the highest entropy. For the voxels of $\mathbf{V}^{(\sigma,c)}$ that have been updated in the reconstruction stage (Sec. 4.3), we synchronize $\alpha^{(E,c)}$ with $\alpha^{(\sigma,c)}$.

**Path Efficiency:** For some active vision tasks, *e.g.* visual navigation, we want a trajectory not only to avoid potential obstacles but also with the shortest path length. To achieve path efficiency, we regularize the trajectory by minimizing the following objective, $\mathcal{L}_{pe} = \frac{1}{(n)}\sum_{i=t+1}^{t+n}||\mathbf{s}_i - \hat{\mathbf{s}}_i||$, where $\hat{\mathbf{s}}_i$ is a reference state generated by linear interpolating $\mathbf{s}_0$ and $\mathbf{s}_g$.

### 4.3. Active 3D Reconstruction

---

**Algorithm 1** Active 3D Reconstruction

---

**Input:** DVGO coarse model $\mathbf{M}^{(c)}$

 1: **Initialization** agent state $s_t = s_0$; goal state candidates $\mathbf{S}_g$; image database $\{I\}_{i=0}^0$; STOP=False
 2: **Coarse Stage:**
 3:     **while** not STOP **do**
 4:         **Global policy**: obtain a goal state $\mathbf{s}_g$ from $\mathbf{S}_g$
 5:         **Local policy**: $\mathbf{S}^* \leftarrow$ planning optim. (Sec. 4.2)
 6:         **Action**: $\mathbf{s}_{t+1} \leftarrow \mathbf{s}^*$
 7:         **Observation**: obtain a new image observation $I_t$
 8:         **Update database**: $\{I\}_{i=0}^{t+1} \leftarrow \{I\}_{i=0}^t$
 9:         **Mapping Optimization**: $\mathbf{M}_{t+1}^{(c)} \leftarrow \mathbf{M}_t^{(c)}$
10:         **Check STOP**: check STOP condition
11:     **end while**
12: **Fine Stage:**
13:     **Initialization** DVGO coarse model $\mathbf{M}_T^{(c)}$; image database $\{I\}_{i=0}^T$
14:     **for** $k \leftarrow 0$ to N **do**
15:         **if** $k \mod m == 0$ **then**
16:             **Training Sample**: sample training pixels (rays) from $\{I\}_{i=0}^T$
17:         **end if**
18:         **Optimize DVGO**: refine $\mathbf{M}^{(f)}$
19:     **end for**

---

In this section, we present an application of our proposed dual-stage active vision framework for the Active 3D Reconstruction task. For Active 3D Reconstruction, a robot is tasked with autonomously reconstructing a scene or object of interest. We want the agent to create a high-quality and complete 3D reconstruction autonomously without hitting obstacles in the scene while moving around. A summary of the algorithm is presented in Algorithm 1.

**Initialization:** We first initialize a 3D bounding box around the space to be reconstructed. Thus, a set of goal-state candidates are uniformly sampled from the view space. Note that the goal candidates are not necessarily located in empty spaces. We then initialize a DVGO coarse model and an agent state from the current state. The initial image observation is appended to the image database.

**Coarse Stage:** At each time step, path planning involves two steps. First, a global policy step selects the next-best-view $\mathbf{s}_g^*$ from the goal candidates within the whole space, which provides a temporary destination $\mathbf{s}_g$. To evaluate the goal candidates, we adopt the following criterion that considers density and information gain,

$$\mathbf{s}_g^* = \underset{\mathbf{s}_g \in \mathbf{S}_g}{\operatorname{argmin}}(\lambda_{cp}l_{cp}(x_{\mathbf{s}_g}) - \lambda_{ig}E_s(\mathbf{s}_g)), \qquad (14)$$

where $x_{\mathbf{s}_g}$ is the 3D location of the state $\mathbf{s}_g$ Next, we execute the path planning method introduced in Sec. 4.2 to obtain an optimized trajectory and the local next-best view $\mathbf{s}^*$. Once the agent maneuvers to the new location and obtains a new observation, The image database is updated with the new observation. Samples are drawn from the image database for optimizing the DVGO model, including synchronizing the entropy volume $\mathbf{V}^{(E,c)}$. At the end of the iteration, we check the early stop condition and stop the exploration process if the stop criteria are met. Here, we consider two stopping criteria: (1) the change in the entropy volume is below a threshold; and (2) the maximum information gain of the feasible goal candidates is below a threshold.

**Fine Stage:** Once the exploration is completed, we initialize the DVGO fine model with the trained coarse model. Stochastic training is performed by sampling pixels (rays) from the image database.

## 5. Experiments and Results

### 5.1. Datasets and Metrics

We evaluate our framework on synthetic and real-world datasets. **Synthetic-NeRF** [37] includes eight objects with realistic images. Following [37], we set the image resolution to $800 \times 800$ pixels. The dataset has been divided into 100 training views and 200 novel views for testing. **Tanks&Temples** [25] is a real-world dataset captured by an inward-facing camera. We use a set of 5 scenes following [32]. The image resolution is $1920 \times 1080$. 1/8 of the images in each scene are used for testing.

**Evaluation**: There are two categories of evaluation metrics adopted in this work. For evaluating rendering performance, we adopt the metrics used in [37], including *PSNR*, *SSIM*, and *LPIPS* [70]. For reconstruction evaluation, we

| Methods | PSNR ↑ | SSIM ↑ | LPIPS ↓ | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
|---|---|---|---|---|---|---|
| | Discrete (local), 10 views | | | Discrete (local), 20 views | | |
| Random | 16.869 | 0.789 | 0.215 | 20.435 | 0.839 | 0.160 |
| FVS | 18.683 | 0.805 | 0.189 | 23.470 | 0.872 | 0.127 |
| **Entropy** | **19.358** | **0.820** | **0.172** | **23.936** | **0.880** | **0.118** |
| | Discrete (free), 10 views | | | Discrete (free), 20 views | | |
| ActiveNeRF [42] | 20.010 | 0.832 | 0.204 | **26.240** | 0.856 | 0.124 |
| **Entropy** | **21.853** | **0.846** | **0.144** | 24.489 | **0.888** | **0.113** |
| | Continuous | | | | | |
| **Ours (Full)** | **30.181** | **0.942** | **0.063** | - | - | - |

(a) Synthetic-NeRF (Rendering)

| Methods | Accu. ↑ | Comp. ↑ | F1-Score ↑ | Chamfer ↓ |
|---|---|---|---|---|
| | Discrete (local), 10 views | | | |
| Random | 0.315 | 0.348 | 0.323 | 0.070 |
| FVS | 0.377 | 0.376 | 0.372 | 0.048 |
| **Entropy** | **0.407** | **0.424** | **0.412** | **0.026** |
| | Discrete (local), 20 views | | | |
| Random | 0.446 | 0.441 | 0.441 | 0.027 |
| FVS | 0.493 | 0.522 | 0.506 | 0.024 |
| **Entropy** | **0.506** | **0.548** | **0.524** | **0.021** |
| | Continuous | | | |
| **Ours (Full)** | **0.604** | **0.602** | **0.603** | **0.015** |

(b) Synthetic-NeRF (Geometry)

| Methods | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
|---|---|---|---|
| Random | 18.169 | 0.827 | 0.225 |
| FVS | 20.240 | 0.844 | 0.204 |
| **Entropy** | **21.560** | **0.848** | **0.200** |

(c) Tanks & Temple (Rendering)

Table 1. Next best view policy comparison (Rendering & Geometry). "**Discrete (local), N views**" means the camera is limited to local motion and the view space is limited to the discrete viewpoints in the training set. Up to N views can be used for training. "**Discrete (free)**" means teleporting between viewpoints in the training set is allowed. "**Continuous**" means the camera moves within a continuous space.

evaluate the reconstruction quality on a mesh surface where 10,000 points are sampled for evaluation. Four metrics are reported, including (1,2) *Accuracy/Completeness* measures the fraction of predicted/reference points that are closer to the reference/predicted points than a threshold distance, which is set to 1cm; (3) *F1 score* is the harmonic mean of accuracy and completeness, which quantifies the overall reconstruction quality; (4) *Chamfer Distance* measures the distance between nearest neighbor correspondences of two point clouds. Note that we use the DVGO models trained offline with full training set as the reference ground truth for evaluation.

## 5.2. Implementation details

**View space:** NeRFs [37, 59] are trained *offline* with a collection of images. However, as we are interested in active vision tasks, we provide images incrementally while training. We first divide our experiments into two sets, defined by the motion spaces, namely *discrete space* and *con-

tinuous space*. **Discrete space:** viewpoints are limited to the viewpoints in the training set. This is further divided into *discrete (free-view)* where the camera is allowed to teleport between any viewpoints and *discrete (local)* where the camera is only allowed to move to one of the three nearest viewpoints. **Continuous space:** we allow the camera to move freely within a predefined view space. To this end, we need to employ a simulator for rendering novel viewpoints. Since *offline* DVGO [59] has shown compelling rendering and reconstruction results, we employ it for a simulation purpose, which implies that DVGO is our model's performance upper bound.

For **Synthetic-NeRF**, the view space is defined as a hemisphere, where the discrete experiments are on the surface of the hemisphere while the continuous experiments are within the space. For **Tanks&Temples**, the view space is within a 3D bounding box defined by the cameras in the training set.

**Mapping optimization:** In the coarse mapping stage, the model is provided with new observations incrementally in the coarse stage, and all the observations are provided in the fine stage. We sample 8,192 rays from the images to train the coarse models for 5,000 iterations and perform refinement for 20,000 iterations. New observations are inserted for every 50 iterations in the coarse stage. Adam optimizer [24] is used with an initial learning rate $10^{-1}$ for all experiments. Scene parameters, *e.g.* grid size, are dataset-dependent. We adopt the settings from [59].

**Planning optimization:** The path planning optimization is carried out with Adam optimizer [1], with a learning rate $10^{-1}$ and 100 training iterations for each planning. $[\lambda_{cp}, \lambda_{ig}, \lambda_{pe}]$ are set to be [10, 1, 0.1]. For the early stop condition, we set the thresholds to be $5 \times 10^{-5}$. Applying the parametric Bézier curve, we generate 100 samples along the curve for planning optimization.

We ran three times for all sets of the experiment and reported the mean result for each set. For each run, we start with a random initial location, which is shared across all sets of the experiment.

## 5.3. Is termination entropy a good criterion?

In Sec. 4.2, we propose the use of termination uncertainty (entropy) of 3D points observed by a viewpoint to determine the information gain of the viewpoint. To validate the effectiveness of the criterion, we compare the approach with some baselines and a prior work [42]. The comparison is presented in Tab. 1. **Baseline (Random)**: a viewpoint is selected randomly as the next best view (NBV) from the viewpoint candidates. **Baseline (FVS)**: furthest view sampling is applied to select the most distanced viewpoint compared with the current observation set as the NBV. **ActiveNeRF** selects NBV based on image rendering uncertainty. Few-shot active learning experiments are designed to com-

| Methods | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
|---|---|---|---|
| *offline training* | | | |
| NeRF [37] | 31.01 | 0.947 | 0.081 |
| Mip-NeRF [5] | 33.09 | 0.961 | 0.043 |
| FastNeRF [20] | 29.97 | 0.941 | 0.053 |
| DVGO [59] | 31.82 | 0.955 | 0.055 |
| *online training* | | | |
| Ours | 30.18 | 0.942 | 0.063 |

Table 2. Quantitative comparison for novel view synthesis on Synthetic-NeRF. *offline* methods trains the model with all collected images while our method is an *online* method that incrementally adds images into training after the model selected the next best view. Moreover, our application scenario is an even harder setting due to a broader NBV searching space. Note that Ours uses images rendered from DVGO for training, thus DVGO is the performance upper bound of Ours.

pare the performance of the approaches. New observations are added incrementally while training, up to 10/20-views. We evaluate the learned representations based on rendering and 3D reconstruction quality. In the cases the camera is limited to local movements, our proposed entropy-based method consistently outperforms other baselines. Our method shows better performance in most of the metrics when compared to ActiveNeRF [42].

## 5.4. Active 3D Reconstruction

### 5.5. Rendering evaluation

We evaluate our proposed active framework for 3D reconstruction based on rendering and reconstruction results. We compare our method against prior arts on the task of novel view synthesis. To our best knowledge, we are the first RGB-based work that trains radiance field from images incrementally acquired in a continuous space, without the need for depth data [48] or pre-training [27]. Though this *online* setting is more difficult than offline training, we still show competitive results when compared to the state-of-the-art methods.

### 5.6. Reconstruction evaluation

We present a quantitative comparison of our method against baselines in Tab. 1b in Fig. 3. To compare against prior work, we compare against [22] and [27] qualitatively. The results of [22, 27] are obtained from [27]. However, as the F1 threshold is not provided in [27], we cannot compare the quantitative results fairly. From Fig. 3, we can see that our methods provide better details among the active reconstruction methods. More examples are shown in Fig. 4.

## 6. Conclusion

In this work, we tackle the mapping and planning problems in *Active SLAM*, which plays an essential role in au-
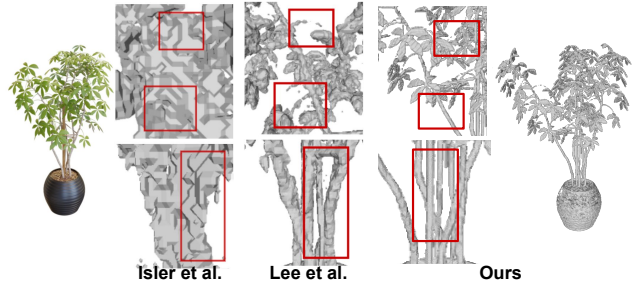


Figure 3. Active reconstruction comparison on Synthetic-NeRF (ficus). Our proposed method is capable of generating high-quality 3D models.
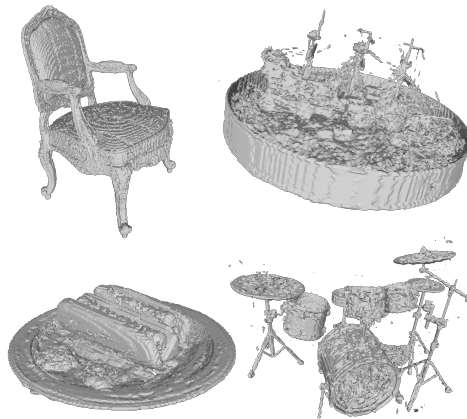


Figure 4. More 3D reconstruction results on Synthetic-NeRF. Our method has shown good reconstruction on *Chair* and *hotdog* but shows minor noises in difficult scenes like *Ship* and *Drums*.

tonomous robotic systems. An iterative dual-stage optimization framework, *ActiveRMAP*, is proposed to perform active mapping and planning with the use of radiance field. Unlike time-consuming NeRFs, we employ an efficient implicit-explicit representation in the mapping stage, which allows our framework to have great potential in robotic applications demanding real-time operations. Photometric difference between observations and volume-rendered images is used as the supervision for mapping optimization. For the planning stage, we propose a new multi-objective loss function that considers collision avoidance, geometric information gain, and path efficiency into consideration. Thanks to the continuity property of the radiance field, we formulate the planning part as a fully differentiable optimization problem. To our best knowledge, though there are concurrent works [27, 48] that address the active 3D reconstruction problem using NeRF, we are the first method without the need for pre-training [27] and depth supervision [48]. Our *online* method has shown a competitive result compared to prior *offline* methods and outperforms other active reconstruction methods using NeRFs.

## 7. Overview

In this supplementary material, we provide more implementation details of ActiveRMAP in Sec. 8. More results and analysis is provided in Sec. 9.

## 8. More ActiveRMAP details

### 8.1. Safe zone modeling

In the main paper Sec. 4.2., we have introduced a collision penalty for calculating the collision cost based on the continuous radiance field. We described the process as minimizing the density of the sampled waypoints. However, in practice, we add a safe zone. The safe zone is defined by a 3D bounding box centered at the waypoints to simulate the agent body and maintain a safe distance away from high-density regions. A set of 3D points $\mathbf{B}_i$ centered at waypoint state $\mathbf{s}_i$ are sampled from the 3D bounding box for computing the penalty. Thus, the full collision penalty is formulated as

$$l_{cp}(x) = \exp(\sigma^{(c)}(x)) \tag{15}$$

$$\mathcal{L}_{cp} = \sum_{i=t+1}^{t+n} \left[ \sum_{x_b \in \mathbf{B}_i} l_{cp}(x_b) \right], \tag{16}$$

where $\sigma^{(c)}(x)$ queries the density at 3D point $x$.

The safe zone setting should be agent-dependent as it should consider the agent's size. As this work performs experiments in simulations, we do not have a specific robotic agent. Thus, we create a 3D bounding box with $1 \times 1 \times 1$ unit as the safe zone. 100 points are sampled uniformly from the bounding box for computing the collision penalty.

### 8.2. View space and Sampling

For **Synthetic-NeRF**, we define the view space as a hemisphere where the maximum radius is 4. We use the spherical coordinate system as the viewpoint representation. However, the spherical coordinate system only represents the 3D coordinates without rotation. Therefore, we constrain the viewpoint to be pointing towards the center of the object of interest. Sampling is performed when generating waypoints from the trajectory and generating 3D points of interest for information gain calculation. We sample 100 waypoints uniformly from the Bézier trajectory. For calculating information gain, we sampled 10 viewpoints from the waypoints to reduce the computation cost. Instead of sampling 3D points from all the rays (pixels) for information gain calculation, we only sample 3D points from $1/100$ of the pixels.

For **Tanks & Temple**, we define the view space as a 3D bounding box whose corners are defined from the raw training set. The viewpoint state is represented by 6DoF $se3$. The sampling process is as same as the one described above.

| Scene | PSNR ↑ | SSIM ↑ | LPIPS ↓ | Accu. ↑ | Comp. ↑ | F1-Score ↑ | Chamfer ↓ |
|---|---|---|---|---|---|---|---|
| chair | 34.100 | 0.976 | 0.027 | - | - | - | - |
| drums | 25.400 | 0.930 | 0.079 | - | - | - | - |
| ficus | 32.500 | 0.977 | 0.025 | - | - | - | - |
| hotdog | 36.700 | 0.980 | 0.035 | - | - | - | - |
| lego | 34.680 | 0.976 | 0.027 | - | - | - | - |
| materials | 29.500 | 0.950 | 0.059 | - | - | - | - |
| mic | 33.100 | 0.982 | 0.018 | - | - | - | - |
| ship | 28.600 | 0.872 | 0.168 | - | - | - | - |
| Avg. | 31.823 | 0.955 | 0.055 | - | - | - | - |

(a) DVGO Baseline (*offline training*)

| Scene | PSNR ↑ | SSIM ↑ | LPIPS ↓ | Accu. ↑ | Comp. ↑ | F1-Score ↑ | Chamfer ↓ |
|---|---|---|---|---|---|---|---|
| chair | 32.770 | 0.968 | 0.037 | 0.639 | 0.637 | 0.638 | 0.011 |
| drums | 24.330 | 0.916 | 0.010 | 0.699 | 0.708 | 0.703 | 0.010 |
| ficus | 31.400 | 0.973 | 0.030 | 0.898 | 0.905 | 0.902 | 0.006 |
| hotdog | 34.200 | 0.970 | 0.054 | 0.424 | 0.411 | 0.418 | 0.016 |
| lego | 32.470 | 0.965 | 0.041 | 0.678 | 0.658 | 0.668 | 0.009 |
| materials | 28.550 | 0.941 | 0.074 | 0.566 | 0.570 | 0.568 | 0.011 |
| mic | 31.830 | 0.976 | 0.026 | 0.803 | 0.808 | 0.806 | 0.007 |
| ship | 25.900 | 0.827 | 0.228 | 0.120 | 0.117 | 0.119 | 0.053 |
| Avg. | 30.181 | 0.942 | 0.063 | 0.604 | 0.602 | 0.603 | 0.015 |

(b) Ours (Full), continuous space

| Scene | PSNR ↑ | SSIM ↑ | LPIPS ↓ | Accu. ↑ | Comp. ↑ | F1-Score ↑ | Chamfer ↓ |
|---|---|---|---|---|---|---|---|
| chair | 20.800 | 0.860 | 0.138 | 0.391 | 0.403 | 0.397 | 0.021 |
| drums | 15.200 | 0.759 | 0.242 | 0.347 | 0.407 | 0.375 | 0.026 |
| ficus | 20.200 | 0.870 | 0.111 | 0.482 | 0.726 | 0.579 | 0.014 |
| hotdog | 24.500 | 0.907 | 0.114 | 0.329 | 0.339 | 0.334 | 0.024 |
| lego | 18.960 | 0.814 | 0.174 | 0.434 | 0.404 | 0.419 | 0.021 |
| materials | 19.600 | 0.833 | 0.148 | 0.522 | 0.482 | 0.501 | 0.013 |
| mic | 22.100 | 0.906 | 0.088 | 0.624 | 0.519 | 0.566 | 0.012 |
| ship | 13.500 | 0.608 | 0.361 | 0.129 | 0.115 | 0.121 | 0.073 |
| Avg. | 19.358 | 0.820 | 0.172 | 0.407 | 0.424 | 0.412 | 0.026 |

(c) Ours (Entropy), Discrete (local), 10 views

| Scene | PSNR ↑ | SSIM ↑ | LPIPS ↓ | Accu. ↑ | Comp. ↑ | F1-Score ↑ | Chamfer ↓ |
|---|---|---|---|---|---|---|---|
| chair | 21.167 | 0.861 | 0.135 | 0.379 | 0.307 | 0.340 | 0.052 |
| drums | 15.030 | 0.748 | 0.254 | 0.379 | 0.308 | 0.339 | 0.052 |
| ficus | 19.667 | 0.863 | 0.119 | 0.489 | 0.711 | 0.579 | 0.013 |
| hotdog | 22.100 | 0.888 | 0.146 | 0.338 | 0.349 | 0.344 | 0.022 |
| lego | 19.630 | 0.816 | 0.171 | 0.430 | 0.390 | 0.409 | 0.023 |
| materials | 18.267 | 0.804 | 0.177 | 0.336 | 0.327 | 0.331 | 0.022 |
| mic | 21.967 | 0.905 | 0.089 | 0.598 | 0.555 | 0.576 | 0.013 |
| ship | 11.640 | 0.56 | 0.423 | 0.064 | 0.059 | 0.061 | 0.184 |
| Avg. | 18.684 | 0.806 | 0.189 | 0.377 | 0.376 | 0.372 | 0.048 |

(d) FVS, Discrete (local), 10 views

| Scene | PSNR ↑ | SSIM ↑ | LPIPS ↓ | Accu. ↑ | Comp. ↑ | F1-Score ↑ | Chamfer ↓ |
|---|---|---|---|---|---|---|---|
| chair | 18.633 | 0.846 | 0.165 | 0.365 | 0.240 | 0.290 | 0.090 |
| drums | 14.700 | 0.756 | 0.261 | 0.356 | 0.386 | 0.371 | 0.029 |
| ficus | 19.133 | 0.854 | 0.133 | 0.413 | 0.681 | 0.514 | 0.017 |
| hotdog | 15.570 | 0.817 | 0.254 | 0.179 | 0.138 | 0.156 | 0.098 |
| lego | 17.960 | 0.801 | 0.189 | 0.351 | 0.331 | 0.341 | 0.031 |
| materials | 17.870 | 0.798 | 0.186 | 0.248 | 0.401 | 0.306 | 0.103 |
| mic | 21.033 | 0.897 | 0.098 | 0.527 | 0.509 | 0.518 | 0.016 |
| ship | 10.050 | 0.545 | 0.431 | 0.079 | 0.099 | 0.088 | 0.176 |
| Avg. | 16.869 | 0.789 | 0.215 | 0.315 | 0.348 | 0.323 | 0.070 |

(e) Random, Discrete (local), 10 views

Table 3. (Part A) Novel view synthesis and 3D reconstruction evaluation result of individual scenes in NeRF-Synthetic dataset.

## 9. More results

### 9.1. More results for termination entropy study

In Sec. 5.3, we present a study about using termination entropy as the criterion for calculating information gain. Here we present more details of the study with qualitative comparisons of individual scenes in the dataset, and quantitative comparisons of the scenes. The quantitative comparison is presented in Tab. 3 and Tab. 4 while the qualitative comparison are shown in Fig. 6 and Fig. 7.

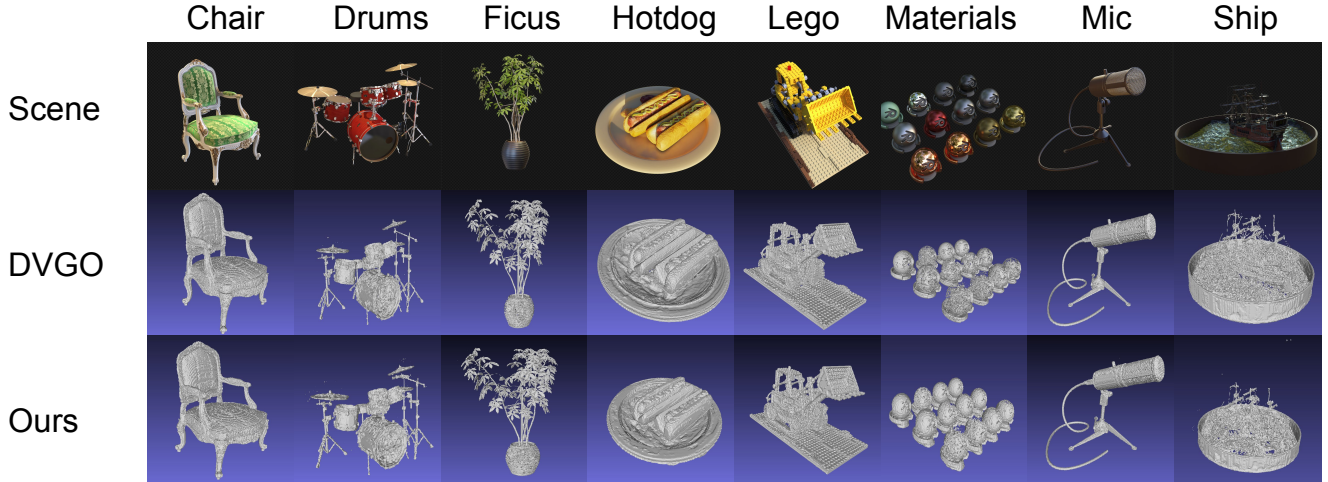For the quantitative result, termination entropy-based

Figure 5. NeRF-Synthetic qualitative comparison between our proposed active 3D reconstruction method and DVGO [59]. Note that, our method uses the images rendered by DVGO as the training images, that implies that DVGO serves as the "ground-truth" and the upper bound of our method. Nevertheless, we show a competitive result when compared to DVGO.
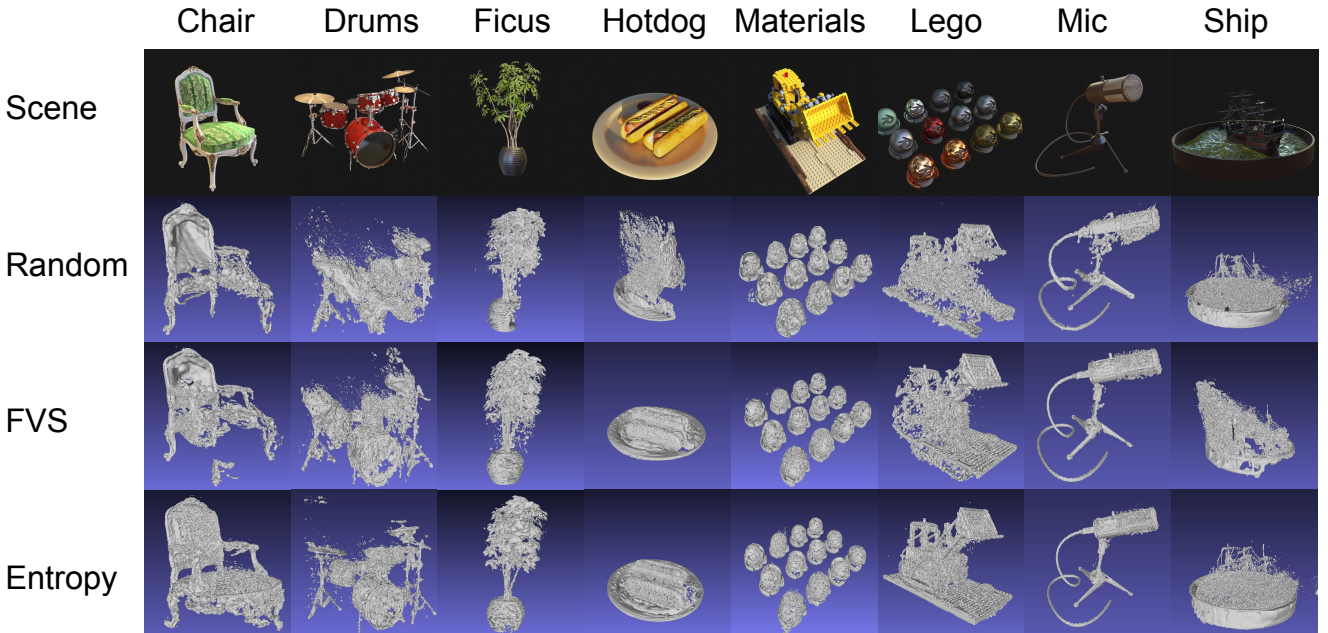


Figure 6. NeRF-Synthetic qualitative comparison for termination entropy study. We compare the 3D reconstruction model of individual scenes. In this set of experiments, the viewpoints are limited to the viewpoints in the raw training set and the movement is contrained locally, *i.e.* the camera can only moved to one of the three nearest viewpoints. The maximum number of viewpoints is 10 in the experiments.

method performs consistently better than the baselines in all the scenes. The qualitative comparison shows that the entropy-based method allows the agent to explore more diverse viewpoints and builds better 3D models with limited number of viewpoints.

## 9.2. More qualitative result

We present a qualitative comparison between our proposed method and DVGO [59] in Fig. 5. As same as most NeRF-based prior work [20,37], DVGO optimizes the scene representation based on a set of collected images in an *offline training* manner. In contrast, we optimize the representation in an *online* fashion by adding new observations
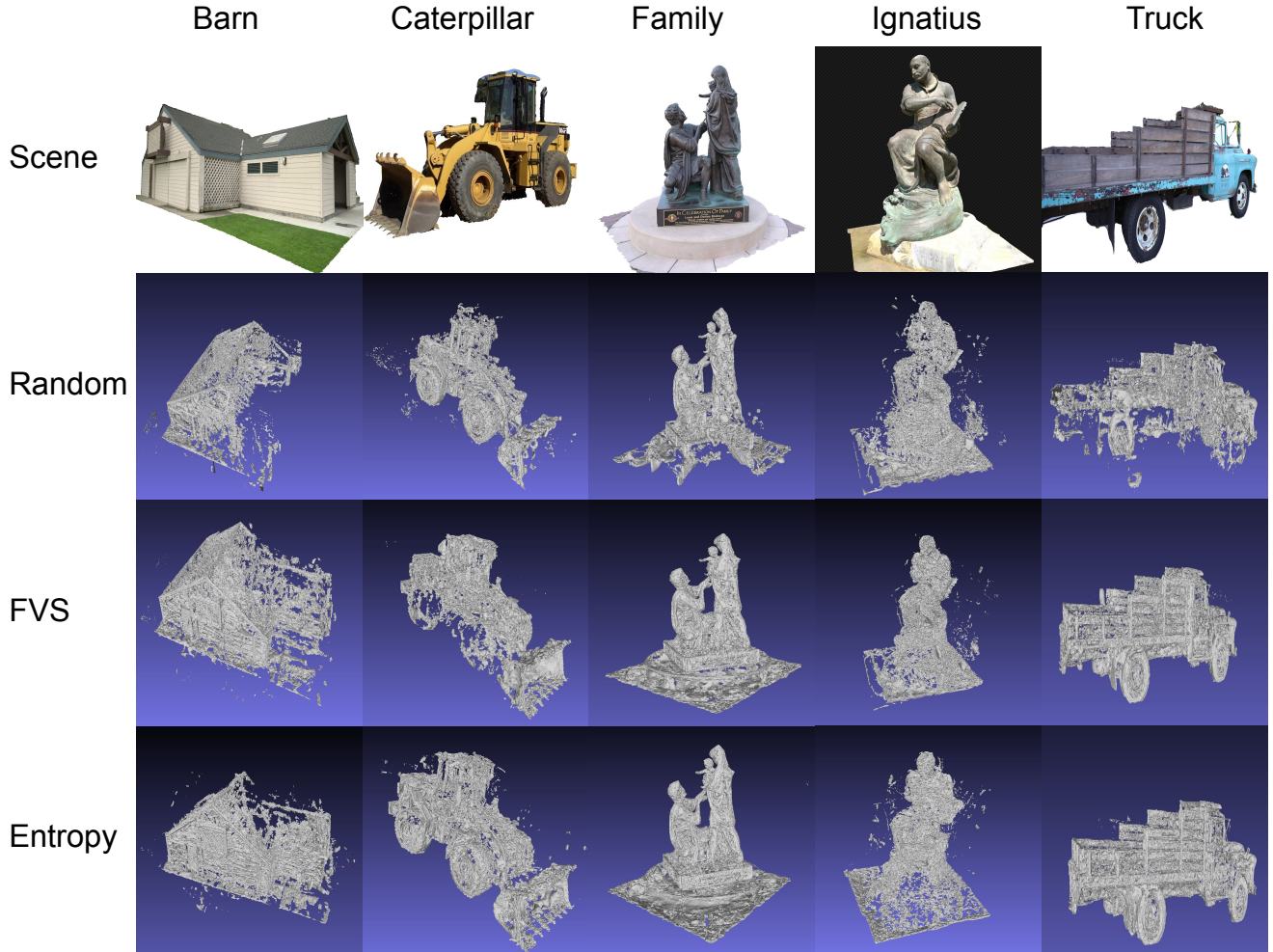
Figure 7. Tanks & Temple qualitative comparison for termination entropy study. We compare the 3D reconstruction model of individual scenes. In this set of experiments, the viewpoints are limited to the viewpoints in the raw training set, and the movement is constrained locally, *i.e.* the camera can only move to one of the three nearest viewpoints. The maximum number of viewpoints is 50 in the experiments.

into the training set incrementally. The agent has to plan for the next best views and perceive new observations based on the current scene representation. Moreover, our method uses the images rendered by DVGO as the training images, which implies that DVGO serves as the "ground-truth" and the upper bound of our method.

We also include a video showing the active 3D reconstruction process.

## References

[1] Michal Adamkiewicz, Timothy Chen, Adam Caccavale, Rachel Gardner, Preston Culbertson, Jeannette Bohg, and Mac Schwager. Vision-only robot navigation in a neural radiance world. *IEEE Robotics and Automation Letters*, 7(2):4606–4613, 2022. 2, 3, 5, 7

[2] John Aloimonos, Isaac Weiss, and Amit Bandyopadhyay. Active vision. *International journal of computer vision*, 1(4):333–356, 1988. 2

[3] Dejan Azinović, Ricardo Martin-Brualla, Dan B Goldman, Matthias Nießner, and Justus Thies. Neural rgb-d surface reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6290–6301, 2022. 1, 3

| Scene | PSNR ↑ | SSIM ↑ | LPIPS ↓ | Accu. ↑ | Comp. ↑ | F1-Score ↑ | Chamfer ↓ |
|---|---|---|---|---|---|---|---|
| chair | 23.470 | 0.875 | 0.117 | - | - | - | - |
| drums | 17.750 | 0.789 | 0.198 | - | - | - | - |
| ficus | 20.070 | 0.87 | 0.105 | - | - | - | - |
| hotdog | 25.000 | 0.902 | 0.122 | - | - | - | - |
| lego | 22.700 | 0.838 | 0.136 | - | - | - | - |
| materials | 20.330 | 0.836 | 0.148 | - | - | - | - |
| mic | 25.500 | 0.941 | 0.055 | - | - | - | - |
| ship | 20.000 | 0.72 | 0.273 | - | - | - | - |
| Avg. | 21.853 | 0.846 | 0.144 | - | - | - | - |

(a) Ours (Entropy), Discrete (free), 10 views

| Scene | PSNR ↑ | SSIM ↑ | LPIPS ↓ | Accu. ↑ | Comp. ↑ | F1-Score ↑ | Chamfer ↓ |
|---|---|---|---|---|---|---|---|
| chair | 26.690 | 0.919 | 0.081 | - | - | - | - |
| drums | 20.200 | 0.847 | 0.154 | - | - | - | - |
| ficus | 21.520 | 0.892 | 0.095 | - | - | - | - |
| hotdog | 28.840 | 0.944 | 0.079 | - | - | - | - |
| lego | 25.220 | 0.893 | 0.099 | - | - | - | - |
| materials | 22.360 | 0.866 | 0.124 | - | - | - | - |
| mic | 27.480 | 0.954 | 0.046 | - | - | - | - |
| ship | 23.600 | 0.786 | 0.222 | - | - | - | - |
| Avg. | 24.489 | 0.888 | 0.113 | - | - | - | - |

(b) Ours (Entropy), Discrete (free), 20 views

| Scene | PSNR ↑ | SSIM ↑ | LPIPS ↓ | Accu. ↑ | Comp. ↑ | F1-Score ↑ | Chamfer ↓ |
|---|---|---|---|---|---|---|---|
| chair | 27.700 | 0.932 | 0.067 | 0.489 | 0.557 | 0.521 | 0.016 |
| drums | 18.730 | 0.822 | 0.179 | 0.512 | 0.561 | 0.536 | 0.014 |
| ficus | 22.800 | 0.905 | 0.083 | 0.744 | 0.867 | 0.801 | 0.008 |
| hotdog | 28.430 | 0.938 | 0.085 | 0.319 | 0.461 | 0.377 | 0.045 |
| lego | 24.730 | 0.885 | 0.105 | 0.553 | 0.534 | 0.543 | 0.012 |
| materials | 22.300 | 0.869 | 0.117 | 0.572 | 0.534 | 0.552 | 0.012 |
| mic | 25.800 | 0.942 | 0.055 | 0.670 | 0.707 | 0.688 | 0.009 |
| ship | 21.000 | 0.752 | 0.255 | 0.186 | 0.166 | 0.175 | 0.053 |
| Avg. | 23.936 | 0.881 | 0.118 | 0.506 | 0.548 | 0.524 | 0.021 |

(c) Ours (Entropy), Discrete (local), 20 views

| Scene | PSNR ↑ | SSIM ↑ | LPIPS ↓ | Accu. ↑ | Comp. ↑ | F1-Score ↑ | Chamfer ↓ |
|---|---|---|---|---|---|---|---|
| chair | 28.100 | 0.934 | 0.065 | 0.449 | 0.497 | 0.472 | 0.019 |
| drums | 18.967 | 0.825 | 0.174 | 0.484 | 0.528 | 0.505 | 0.015 |
| ficus | 22.770 | 0.904 | 0.086 | 0.723 | 0.835 | 0.775 | 0.009 |
| hotdog | 29.830 | 0.947 | 0.075 | 0.408 | 0.425 | 0.416 | 0.017 |
| lego | 25.070 | 0.89 | 0.102 | 0.551 | 0.529 | 0.540 | 0.012 |
| materials | 21.870 | 0.858 | 0.13 | 0.477 | 0.487 | 0.482 | 0.013 |
| mic | 25.330 | 0.939 | 0.059 | 0.716 | 0.752 | 0.733 | 0.008 |
| ship | 15.830 | 0.679 | 0.327 | 0.132 | 0.122 | 0.127 | 0.094 |
| Avg. | 23.471 | 0.872 | 0.127 | 0.493 | 0.522 | 0.506 | 0.024 |

(d) FVS, Discrete (local), 20 views

| Scene | PSNR ↑ | SSIM ↑ | LPIPS ↓ | Accu. ↑ | Comp. ↑ | F1-Score ↑ | Chamfer ↓ |
|---|---|---|---|---|---|---|---|
| chair | 21.667 | 0.871 | 0.127 | 0.387 | 0.331 | 0.357 | 0.045 |
| drums | 17.530 | 0.805 | 0.199 | 0.417 | 0.473 | 0.444 | 0.020 |
| ficus | 21.370 | 0.887 | 0.105 | 0.631 | 0.796 | 0.704 | 0.010 |
| hotdog | 21.270 | 0.884 | 0.151 | 0.358 | 0.337 | 0.347 | 0.030 |
| lego | 22.767 | 0.859 | 0.13 | 0.431 | 0.405 | 0.417 | 0.023 |
| materials | 20.267 | 0.838 | 0.145 | 0.486 | 0.435 | 0.459 | 0.015 |
| mic | 24.630 | 0.932 | 0.064 | 0.696 | 0.612 | 0.651 | 0.011 |
| ship | 13.980 | 0.634 | 0.357 | 0.163 | 0.137 | 0.149 | 0.065 |
| Avg. | 20.435 | 0.839 | 0.160 | 0.446 | 0.441 | 0.441 | 0.027 |

(e) Random, Discrete (local), 20 views

Table 4. (Part B) Novel view synthesis and 3D reconstruction evaluation result of individual scenes in NeRF-Synthetic dataset.

[4] Ruzena Bajcsy. Active perception. *Proceedings of the IEEE*, 76(8):966–1005, 1988. 2

[5] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5855–5864, 2021. 8

[6] Giuseppe Borghi and Vincenzo Caglioti. Minimum uncertainty explorations in the self-localization of mobile robots. *IEEE Transactions on Robotics and Automation*, 14(6):902–911, 1998. 2

[7] Frederic Bourgault, Alexei A Makarenko, Stefan B Williams, Ben Grocholsky, and Hugh F Durrant-Whyte. Information based adaptive robotic exploration. In *IEEE/RSJ international conference on intelligent robots and systems*, volume 1, pages 540–545. IEEE, 2002. 2

[8] Cesar Cadena, Luca Carlone, Henry Carrillo, Yasir Latif, Davide Scaramuzza, José Neira, Ian Reid, and John J Leonard. Past, present, and future of simultaneous localization and mapping: Toward the robust-perception age. *IEEE Transactions on robotics*, 32(6):1309–1332, 2016. 2

[9] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. In *European Conference on Computer Vision (ECCV)*, 2022. 3, 4

[10] Shin-Fang Chng, Sameera Ramasinghe, Jamie Sherrah, and Simon Lucey. Gaussian activated neural radiance fields for high fidelity reconstruction and pose estimation. In *European Conference on Computer Vision*, pages 264–280. Springer, 2022. 2, 3

[11] Cl Connolly. The determination of next best views. In *Proceedings. 1985 IEEE international conference on robotics and automation*, volume 2, pages 432–435. IEEE, 1985. 2

[12] Cregg K. Cowan and Peter D Kovesi. Automatic sensor placement from vision task requirements. *IEEE Transactions on Pattern Analysis and machine intelligence*, 10(3):407–416, 1988. 2

[13] Andrew J Davison and David W. Murray. Simultaneous localization and map-building using active vision. *IEEE transactions on pattern analysis and machine intelligence*, 24(7):865–880, 2002. 2

[14] Andrew J Davison, Ian D Reid, Nicholas D Molton, and Olivier Stasse. Monoslam: Real-time single camera slam. *IEEE transactions on pattern analysis and machine intelligence*, 29(6):1052–1067, 2007. 2

[15] Jeffrey Delmerico, Stefan Isler, Reza Sabzevari, and Davide Scaramuzza. A comparison of volumetric information gain metrics for active 3d object reconstruction. *Autonomous Robots*, 42(2):197–208, 2018. 2

[16] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised nerf: Fewer views and faster training for free. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12882–12891, 2022. 3

[17] Hugh Durrant-Whyte and Tim Bailey. Simultaneous localization and mapping: part i. *IEEE robotics & automation magazine*, 13(2):99–110, 2006. 2

[18] Hans Jacob S Feder, John J Leonard, and Christopher M Smith. Adaptive mobile robot navigation and mapping. *The International Journal of Robotics Research*, 18(7):650–668, 1999. 2

[19] Dieter Fox, Wolfram Burgard, and Sebastian Thrun. Active markov localization for mobile robots. *Robotics and Autonomous Systems*, 25(3-4):195–207, 1998. 2

[20] Stephan J Garbin, Marek Kowalski, Matthew Johnson, Jamie Shotton, and Julien Valentin. Fastnerf: High-fidelity neural rendering at 200fps. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14346–14355, 2021. 8, 10

[21] Heiko Hirschmuller. Accurate and efficient stereo processing by semi-global matching and mutual information. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, volume 2, pages 807–814. IEEE, 2005. 2

[22] Stefan Isler, Reza Sabzevari, Jeffrey Delmerico, and Davide Scaramuzza. An information gain formulation for active volumetric 3d reconstruction. In *2016 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3477–3484. IEEE, 2016. 2, 8

[23] James T Kajiya and Brian P Von Herzen. Ray tracing volume densities. *ACM SIGGRAPH computer graphics*, 18(3):165–174, 1984. 3, 4

[24] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 7

[25] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Transactions on Graphics*, 36(4), 2017. 6

[26] Simon Kriegel, Christian Rink, Tim Bodenmüller, and Michael Suppa. Efficient next-best-scan planning for autonomous 3d surface reconstruction of unknown objects. *Journal of Real-Time Image Processing*, 10(4):611–631, 2015. 2

[27] Soomin Lee, Le Chen, Jiahao Wang, Alexander Liniger, Suryansh Kumar, and Fisher Yu. Uncertainty guided policy for active robotic 3d reconstruction using neural radiance fields. *IEEE Robotics and Automation Letters*, 2022. 2, 3, 8

[28] Kejie Li, Yansong Tang, Victor Adrian Prisacariu, and Philip HS Torr. Bnv-fusion: Dense 3d reconstruction using bi-level neural volume fusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 1, 3

[29] Chen-Hsuan Lin, Wei-Chiu Ma, Antonio Torralba, and Simon Lucey. Barf: Bundle-adjusting neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5741–5751, 2021. 2, 3

[30] Celong Liu, Zhong Li, Junsong Yuan, and Yi Xu. Neulf: Efficient novel view synthesis with neural 4d light field. *arXiv preprint arXiv:2105.07112*, 2021. 3

[31] Jiachen Liu, Pan Ji, Nitin Bansal, Changjiang Cai, Qingan Yan, Xiaolei Huang, and Yi Xu. Planemvs: 3d plane reconstruction from multi-view stereo. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8665–8675, 2022. 2

[32] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. *Advances in Neural Information Processing Systems*, 33:15651–15663, 2020. 6

[33] Iker Lluvia, Elena Lazkano, and Ander Ansuategi. Active mapping and robot exploration: A survey. *Sensors*, 21(7):2445, 2021. 2

[34] Alexei A Makarenko, Stefan B Williams, Frederic Bourgault, and Hugh F Durrant-Whyte. An experiment in integrated exploration. In *IEEE/RSJ international conference on intelligent robots and systems*, volume 1, pages 534–539. IEEE, 2002. 2

[35] Jasna Maver and Ruzena Bajcsy. Occlusions as a guide for planning the next view. *IEEE transactions on pattern analysis and machine intelligence*, 15(5):417–433, 1993. 2

[36] Nelson Max. Optical models for direct volume rendering. *IEEE Transactions on Visualization and Computer Graphics*, 1(2):99–108, 1995. 4

[37] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021. 1, 2, 3, 4, 6, 7, 8, 10

[38] Christian Mostegel, Andreas Wendel, and Horst Bischof. Active monocular localization: Towards autonomous monocular exploration for multirotor mavs. In *2014 IEEE international conference on robotics and automation (ICRA)*, pages 3848–3855. IEEE, 2014. 2

[39] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Trans. Graph.*, 41(4):102:1–102:15, July 2022. 3, 4

[40] Paul Newman, Michael Bosse, and John Leonard. Autonomous feature-based exploration. In *2003 IEEE International Conference on Robotics and Automation (Cat. No. 03CH37422)*, volume 1, pages 1234–1240. IEEE, 2003. 2

[41] Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11453–11464, 2021. 1, 3

[42] Xuran Pan, Zihang Lai, Shiji Song, and Gao Huang. Activenerf: Learning where to see with uncertainty estimation. In *European Conference on Computer Vision*, pages 230–246. Springer, 2022. 2, 3, 4, 7, 8

[43] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 165–174, 2019. 1, 3

[44] Daryl Peralta, Joel Casimiro, Aldrin Michael Nilles, Justine Aletta Aguilar, Rowel Atienza, and Rhandley Cajote. Next-best view policy for 3d reconstruction. In *European Conference on Computer Vision*, pages 558–573. Springer, 2020. 2

[45] Richard Pito. A solution to the next best view problem for automated surface acquisition. *IEEE Transactions on pattern analysis and machine intelligence*, 21(10):1016–1030, 1999. 2

[46] Julio A Placed, Jared Strader, Henry Carrillo, Nikolay Atanasov, Vadim Indelman, Luca Carlone, and José A Castellanos. A survey on active simultaneous localization and mapping: State of the art and new frontiers. *arXiv preprint arXiv:2207.00254*, 2022. 2

[47] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10318–10327, 2021. 1, 3

[48] Yunlong Ran, Jing Zeng, Shibo He, Lincheng Li, Yingfeng Chen, Gimhee Lee, Jiming Chen, and Qi Ye. Neurar: Neural uncertainty for autonomous 3d reconstruction. *arXiv preprint arXiv:2207.10985*, 2022. 2, 3, 4, 8

[49] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. In *International Conference on Computer Vision (ICCV)*, 2021. 3

[50] Sara Fridovich-Keil and Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *CVPR*, 2022. 3, 4

[51] Davide Scaramuzza and Friedrich Fraundorfer. Visual odometry [tutorial]. *IEEE robotics & automation magazine*, 18(4):80–92, 2011. 2

[52] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4104–4113, 2016. 2

[53] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. Graf: Generative radiance fields for 3d-aware image synthesis. *Advances in Neural Information Processing Systems*, 33:20154–20166, 2020. 1, 3

[54] Steven M Seitz, Brian Curless, James Diebel, Daniel Scharstein, and Richard Szeliski. A comparison and evaluation of multi-view stereo reconstruction algorithms. In *2006 IEEE computer society conference on computer vision and pattern recognition (CVPR'06)*, volume 1, pages 519–528. IEEE, 2006. 2

[55] Roland Siegwart, Illah Reza Nourbakhsh, and Davide Scaramuzza. *Introduction to autonomous mobile robots*. MIT press, 2011. 2

[56] Cyrill Stachniss. *Robotic mapping and exploration*, volume 55. Springer, 2009. 2

[57] Cyrill Stachniss, Dirk Hahnel, and Wolfram Burgard. Exploration with active loop-closing for fast-slam. In *2004 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)(IEEE Cat. No. 04CH37566)*, volume 2, pages 1505–1510. IEEE, 2004. 2

[58] Edgar Sucar, Shikun Liu, Joseph Ortiz, and Andrew J Davison. imap: Implicit mapping and positioning in real-time. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6229–6238, 2021. 2, 3

[59] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *CVPR*, 2022. 3, 4, 7, 8, 10

[60] Jiaming Sun, Yiming Xie, Linghao Chen, Xiaowei Zhou, and Hujun Bao. Neuralrecon: Real-time coherent 3d reconstruction from monocular video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15598–15607, 2021. 2, 3

[61] Jian Sun, Nan-Ning Zheng, and Heung-Yeung Shum. Stereo matching using belief propagation. *IEEE Transactions on pattern analysis and machine intelligence*, 25(7):787–800, 2003. 2

[62] Sebastian Thrun. Probabilistic robotics. *Communications of the ACM*, 45(3):52–57, 2002. 2

[63] Sebastian B Thrun and Knut Möller. Active exploration in dynamic environments. *Advances in neural information processing systems*, 4, 1991. 2

[64] Zirui Wang, Shangzhe Wu, Weidi Xie, Min Chen, and Victor Adrian Prisacariu. Nerf–: Neural radiance fields without known camera parameters. *arXiv preprint arXiv:2102.07064*, 2021. 2, 3

[65] Qingyong Xie and Yongcai Wang. A survey of filtering based active localization methods. In *Proceedings of the 2020 the 4th International Conference on Big Data and Internet of Things*, pages 69–73, 2020. 2

[66] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. In *Proceedings of the European conference on computer vision (ECCV)*, pages 767–783, 2018. 2

[67] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4578–4587, 2021. 1, 2, 3

[68] Huangying Zhan, Chamara Saroj Weerasekera, Jia-Wang Bian, and Ian Reid. Visual odometry revisited: What should be learnt? In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 4203–4210. IEEE, 2020. 2

[69] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*, 2020. 1, 3

[70] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018. 6