

Multimodal LLM Guided Exploration and Active Mapping using Fisher Information

Wen Jiang* Boshu Lei* Katrina Ashton Kostas Daniilidis
University of Pennsylvania

Abstract

We present an active mapping system that could plan for long-horizon exploration goals and short-term actions with a 3D Gaussian Splatting (3DGS) representation. Existing methods either did not take advantage of recent developments in multimodal Large Language Models (LLM) or did not consider challenges in localization uncertainty, which is critical in embodied agents. We propose employing multimodal LLMs for long-horizon planning in conjunction with detailed motion planning using our information-based algorithm. By leveraging high-quality view synthesis from our 3DGS representation, our method employs a multimodal LLM as a zero-shot planner for long-horizon exploration goals from the semantic perspective. We also introduce an uncertainty-aware path proposal and selection algorithm that balances the dual objectives of maximizing the information gain for the environment while minimizing the cost of localization errors. Experiments conducted on the Gibson and Habitat-Matterport 3D datasets demonstrate state-of-the-art results of the proposed method.

1. Introduction

Being able to autonomously explore and map an environment while localizing within that map is a core skill for a mobile robot. This ability could empower embodied artificial intelligence systems with effective 3D scene understanding through use in conjunction with open-vocabulary semantic segmentation or vision-language features [55, 68, 87, 91] and can be used as a basis for language-specified robotics tasks [25, 26, 37, 63]. This task is challenging because it requires general knowledge of the typical layout of an environment to identify coarse targets for long-term exploration and the ability to accurately plan each step to ensure localization accuracy while maximizing information gain or ‘surprise’ when traversing the environment.

In this paper, We propose to use 3D Gaussians as our scene representation and tackle this problem in two phases:

*equal contribution

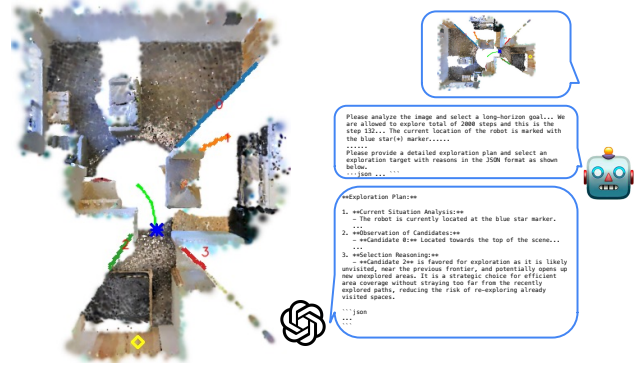


Figure 1. An overview about the interaction with a multimodal LLM. The image on the left is a bird’s-eye view rendering from our 3DGS representation. A detailed example of the dialog can be found in the supplementary.

long-horizon planning with an LLM and detailed motion planning with our information-based approach by quantifying the Fisher Information for 3DGS parameters and current localization states. We prompt a multimodal LLM with our current map state, trajectory, and candidates for long-term exploration frontiers to exploit the LLM’s world knowledge about scene layouts. We then follow the suggestion from the LLM by choosing paths that maximize the Expected Information Gain calculated using the Fisher Information about poses and 3D Gaussians.

Classical active mapping methods often define objectives to reduce the uncertainty of the state estimate based on its covariance [1, 5, 11, 20, 58, 78]. The covariance is readily available for classical systems which use filters to update the state estimate. However, many recent mapping systems use non-linear optimization to update the state estimate, making the covariance difficult to obtain. In particular, systems [27, 32, 44, 83] with 3DGS [33] for the scene representation have been developed for high-fidelity rendering of novel views of the scene. Previous methods attempted to quantify the uncertainty of radiance fields for reconstructing scenes from given data [23, 66, 67, 72] for active view selection [31, 36, 51, 72] and active reconstruction [84] or mapping [85] of scenes with given localization, and for 3D

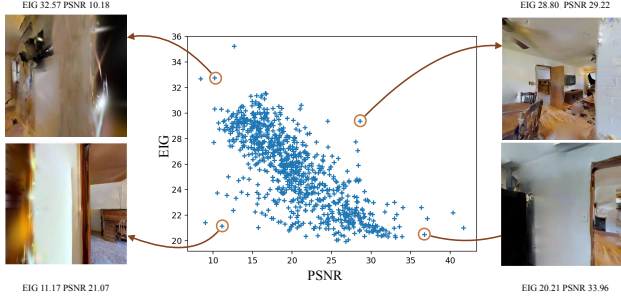


Figure 2. **Scatter Plot of EIG vs. PSNR** We plot the EIG and PSNR at sampled poses in the Cantwell scene of the Gibson dataset. The figure corroborates the intuition that the robot expects to gain little information (low EIG) at well reconstructed region (high PSNR) and gain much information (high EIG) at a poorly reconstructed region (low PSNR).

reconstruction of small scenes with an inward-facing camera [88]. However, all the prior methods only model the uncertainty of the scene representation, ignoring another primary source of error for a practical mapping system: the risk of localization failures. Our proposed system could address not only Information Gain during exploration but also consider localization accuracy when the robot explores unknown and texture-less regions.

In terms of long-horizon planning for robot exploration, classical approaches such as frontier-based exploration and A* algorithms are still used in active mapping systems for their efficiency and simplicity [4, 34, 70, 71, 77]. However, algorithms that use only simple heuristics cannot determine the information gain, nor estimate 3D geometry, limiting their objective to simply improving coverage or having minimal travel distances. To address this issue, learnable approaches have been trained for exploration tasks [22, 38, 46, 59], but these methods are trained for a limited scene distribution and are not natively designed for a 3D Gaussian representation. Recently, LLMs and Vision Language Models (VLM) have demonstrated extraordinary abilities in visual grounding and logical reasoning and have been studied with various robotic tasks [62]. However, these works either focus on other tasks [9, 14, 16], especially tasks with language inputs [26, 74] and navigation [8] or did not take advantage of photo-realistic 3D representation for multimodal LLMs. On the other hand, Neural Radiance Field and 3D Gaussian Splatting have shown the exciting ability to reconstruct accurate 3D scenes just using posed images. Later works [32, 44] enabled online training with localizations of images. However, the 3DGS model reconstructs the scene from a pure optimization perspective without prior knowledge of the scene and only passively works on RGB-D video inputs. In contrast, we can actively select the best paths for exploration to improve the quality of both reconstruction and localization.

We argue that we can leverage the best of both worlds.

We can use high-quality bird’s-eye view rendering from a 3D Gaussian Splatting model as the input of the multimodal LLM to help identify long-horizon exploration goals from frontiers in our scene representation. In this way, we could seamlessly integrate multimodal LLM into a path-planning model to provide knowledge of the typical layout of a scene and potential regions of interest for exploration. Besides, we also alleviate the issue by asking the multimodal LLM to provide precise coordinates, which require precise 3D understanding and are dependent on the scale and relative location of the scene. Given the long-horizon exploration goal from the output of the multimodal LLM, we formulate our path planning problem as an active learning problem and propose multiple paths leading to the exploration goal. We consider the expected information gain (EIG) along the planned trajectory with Fisher Information. As shown in Fig. 2, our EIG metric could reflect the ‘surprise’ and correlate well with the rendering quality on candidate views without actually taking a picture at the candidate location. We also derived localization uncertainty on candidate paths with Cramér–Rao bound to balance the dual problem of exploration and localization. Notably, we compute the full Fisher Information matrix for localization parameters without any approximation with our efficient CUDA implementation. Thus, finding the path that most improves scene reconstruction quality when taking observations during exploration while ensuring the agent still has accurate localization. To validate our approach, we evaluate our method on scenes from the Gibson [81] and Habitat-Matterport 3D [56] datasets quantitatively and qualitatively. We show superior reconstruction quality in various metrics compared to several baselines and recent state-of-the-art methods [7, 22, 53, 82, 85]. Our contributions can be summarized as follows:

- We present an active mapping system that could autonomously explore the environment and extensively compare our system with previous approaches. To the best of our knowledge, we are the first active mapping system with 3D Gaussian representation that is not dependent on ground truth camera pose readings.
- We provide a way to leverage the zero-shot long-horizon planning ability of LLM into our active mapping system seamlessly.
- We introduce localization uncertainty with active mapping systems and effectively balance the information gain for exploration and the cost of possible localization errors.

2. Related work

Active Mapping and Localization Efficiently exploring an environment in order to map it while being able to localize in that map is a fundamental problem in robotics. There are many methods that address aspects of this problem separately – Simultaneous Localization and Mapping (SLAM) methods [3, 17, 18, 30, 47–49, 75] address the mapping and

localization. Many exploration methods [10, 22, 85] address the exploration and mapping aspects while assuming poses are provided. Active SLAM methods [4, 70, 71, 77] consider both of these problems. The exploration in these systems is usually driven by a measure of uncertainty [39, 54]; specific utility functions are often drawn from either Information Theory (IT) [65], or the Theory of Optimal Experimental Design (TOED) [52]. Recently learning-based approaches have been developed for active mapping and localization, Active Neural SLAM [7] learns policies to drive exploration and estimate the agent pose, more similarly to us NARUTO [21] uses a SLAM backbone to estimate the pose and uses uncertainty to drive exploration. However, they only consider the reconstruction uncertainty, whereas we choose long-horizon planning goals and consider the localization uncertainty as well as the reconstruction uncertainty when planning the best path to the long-horizon goal.

Uncertainty quantification for radiance fields The vast majority of previous work on uncertainty quantification for radiance fields has been for post-processing scenes [23, 66, 67, 72], view selection [31, 36, 51, 72] or active view selection [31, 51, 72], all of which assume the input images are posed. Active neural mapping [85] uses neural variability, that is the prediction robustness against random weight perturbations, as an estimate of uncertainty to actively map a scene with ground truth poses provided. Fisher-RF [31] also performs active scene mapping with ground truth poses provided, based on an approximation of the Fisher Information of views along candidate paths. Zhan et al. [88] perform active reconstruction without ground truth camera poses. However, they only evaluate small-scale scenes and limit the camera trajectories to be inwards facing and only model scene uncertainty, not localization uncertainty, which is a key consideration for active mapping and exploration.

Robot Planning with Foundation Models LLMs and VLMs have been used widely in robotics as for language-specified tasks, either as high-level planners in conjunction with other methods for low-level control [2, 15, 16, 26, 28, 40, 57, 69], to create intermediate representations which can be planned over [29] or directly outputting actions [41, 90]. Outside of language-driven tasks, VLMs and LLMs have been used in robotics for providing rewards to drive exploration in Reinforcement Learning [76], and for visual localization [45]. More similar to our task, LLMs and VLMs have also been used for goal-driven navigation, such as searching for specific objects [13, 60, 64, 86].

3. Method

We divide the scene exploration task into two phases: (a) long-horizon planning for the coarse direction of exploration

that leads to better coverage and understanding (Sec. 3.2), and (b) detailed trajectory proposal (Sec. 3.3) and path selection (Sec. 3.4) that aims to improve 3D reconstruction and localization from a geometrical perspective. The former task relies more on semantic information and prior knowledge about the possible layout of the environment, which is suitable for a generalist model with prior distributions on the scene. The latter task, however, involves detailed motion planning and is better addressed by considering information gain on fine-level geometries. Fig. 3 shows an overview of our method.

3.1. Preliminary

In 3D Gaussian Splatting (3DGS) [33], the scene is represented by a set of 3D Gaussians whose color and opacity are learned via a rendering loss. A rendered image can be produced by projecting the Gaussians to 2D and using α -blending for the N ordered points on the 2D splat that overlaps each pixel. The Jacobian of the current camera pose \mathbf{x} with respect to the parameters of 3D Gaussians can be computed by defining the (left) partial derivative on the manifold [44]:

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

where $\tau \in \mathfrak{se}(3)$, \circ is a group composition, and Exp and Log are the exponential and logarithmic mappings between Lie algebra and Lie Group.

Fisher Information is a measurement of the information that a random variable \mathbf{y} carries about an unknown parameter \mathbf{w} of a distribution that models \mathbf{y} . In the problem of novel view synthesis, we are interested in measuring the observed information of a radiance field with parameters \mathbf{w} at a camera pose \mathbf{x} using the negative log-likelihood of the image observation \mathbf{y} taken from that pose:

$$-\log p(\mathbf{y}|\mathbf{x}; \mathbf{w}) = (\mathbf{y} - f(\mathbf{x}, \mathbf{w}))^T (\mathbf{y} - f(\mathbf{x}, \mathbf{w})), \quad (2)$$

where $f(\mathbf{x}, \mathbf{w})$ is the rendering model. Under regularity conditions [61], the Fisher Information of $-\log p(\mathbf{y}|\mathbf{x}; \mathbf{w})$ is the Hessian of Eq. 2 with respect to \mathbf{w} , denoted $\mathbf{H}''[\mathbf{y}|\mathbf{x}, \mathbf{w}]$. In our formulation \mathbf{w} is a tensor for all the 3D Gaussian Parameters $R^{N \times 14}$, where N is the number of Gaussians, and $\mathbf{y} \in R^{h \times w \times c}$ is the RGB-D observation.

FisherRF [31] addressed the active view selection problem that starts with a training set of views D^{train} and aims to select the next best view from a set of candidate $SE(3)$ camera poses $\mathbf{x}_i^{acq} \in D^{pool}$ without obtaining the image \mathbf{y}_i^{acq} at the camera pose \mathbf{x}_i^{acq} . The next best view is chosen by finding:

$$\arg \max_{\mathbf{x}_i^{acq} \in D^{pool}} \text{tr} (\mathbf{H}''[\mathbf{y}_i^{acq}|\mathbf{x}_i^{acq}, \mathbf{w}] \mathbf{H}''[\mathbf{w}|D^{train}]^{-1}), \quad (3)$$

where \mathbf{w} is the initial estimate of model parameters using the current training set D^{train} . $\mathbf{H}''[\mathbf{w}|D^{train}]^{-1}$ can

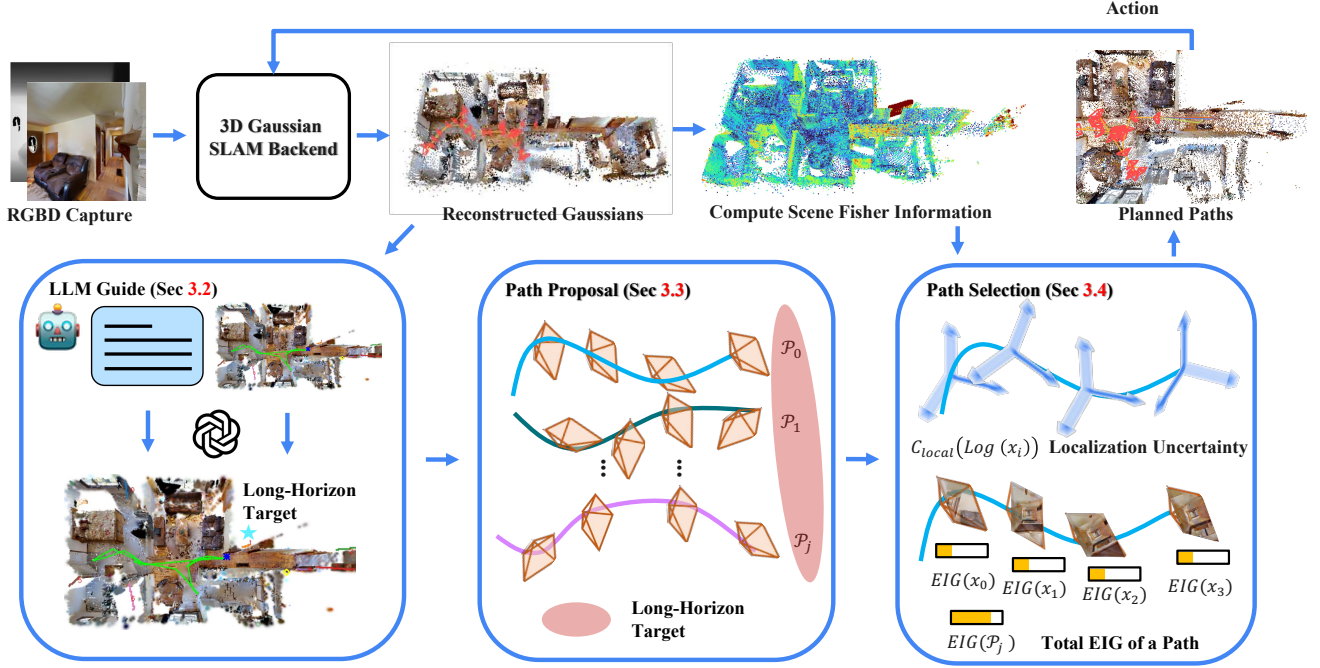


Figure 3. **An Illustration of Our Active Mapping System** Our method first identifies long-horizon targets with a multimodal LLM by using novel-view synthesis from our 3DGS representation. Our information-based path proposal method then proposes and selects short-term action sequences based on the Fisher Information about the 3D Gaussian parameters. The best path and action along the path is selected with respect to both the information gain and localization accuracy.

be computed by summing the Hessians of model parameters across all different views in $\{D_{train}\}$ before inverting. The key of this algorithm is that the Fisher Information $\mathbf{H}''[\mathbf{y}_i^{acq} | \mathbf{x}_i^{acq}, \mathbf{w}]$ does not depend on the label \mathbf{y}_i^{acq} of the acquisition sample \mathbf{x}_i^{acq} . Therefore, it is feasible to compute the Expected Information Gain (EIG) before visiting the potential view candidate \mathbf{x}_i^{acq} . However, the number of optimizable parameters is typically more than 20 million, which means it is impossible to compute without sparsification or approximation. In practice, FisherRF [31] applies a Laplace approximation [12, 42] that approximates the Hessian matrix with its diagonal values plus a log-prior regularizer λI

$$\mathbf{H}''[\mathbf{y}|\mathbf{x}, \mathbf{w}] \simeq \text{diag}(\nabla_{\mathbf{w}} f(\mathbf{x}, \mathbf{w})^T \nabla_{\mathbf{w}} f(\mathbf{x}, \mathbf{w})) + \lambda I. \quad (4)$$

3.2. Long Horizon Exploration with Foundation Model

We propose using a multimodal LLM to perform the task of long-horizon planning as a zero-shot planner for the possible paths and leave the detailed path planning for a closed-form uncertainty-aware motion planning algorithm that computes 3D geometries.

As we consider an agent moving in a 2D action plane (e.g., a ground robot), we leverage the view synthesis ability of 3DGS to create expressive bird’s-eye view renderings that could provide an overview of the environment. An occupancy grid on the motion plane of our robot is created using

our 3D Gaussian representation, which can be used to identify the frontiers of the current environment. The frontiers are defined as points on the boundary between free space and unobserved space (not marked as free or occupied) in our occupancy map. Each frontier can be defined as a set of neighboring points on the action plane $Q = \{\phi(\mathbf{x}_i)\dots\}$, where $\phi(\cdot)$ is the projection from the world frame to the bird’s-eye view. The agent has multiple choices on the frontiers, especially at the early stages of exploration. We apply Chain-of-Thought Prompting [80] to encourage the multimodal LLM to provide analysis on the candidate frontiers first before selection. We provide contextual information such as the description of the task, the total steps allowed for exploration, and our current step in the textual prompt. The multimodal LLM is also allowed to decide whether the robot should not go to a frontier and instead focus on improving existing regions. We embed possible choices of frontiers along with the current location and previously visited trajectory in the bird’s-eye view rendering of the environment along with a text prompt describing the task and contextual information as the input for the multimodal LLM. In this way, we avoid possible infeasible destinations for motion planning. For example, the unexplored regions we have no reconstruction for at this point might be unreachable from the current space. We then parse the long-horizon exploration goal from the reply of LLM. We use GPT-4o [50] for our long-horizon

planning task, but our method is agnostic to the underlying multimodal LLM as the inputs for the multimodal LLM are text prompts and a bird’s-eye view rendering. Fig. 1 provides an overview of our interaction with the multimodal LLM. Detailed examples can be found in the Supplementary.

3.3. Information-driven Path Proposal

After identifying a long horizon with the multimodal LLM, we form an initial set of candidate poses \mathcal{T}_I by sampling points in the coarse region of interest from the LLM. If there are no unvisited boundaries or the multimodal LLM suggests we do not need to explore the frontiers, we sample poses across the free space to form \mathcal{T}_I . We then evaluate the Expected Information Gain (EIG) for each pose $\mathbf{x}_i^{acq} \in \mathcal{T}_I$, given by

$$\text{EIG}(\mathbf{x}_i^{acq}) = \text{tr}(\mathbf{H}''[\mathbf{y}_i^{acq}|\mathbf{x}_i^{acq}, \mathbf{w}] \mathcal{I}(\mathbf{w})^{-1}), \quad (5)$$

as a preliminary selection metric to form our final candidate target poses set \mathcal{T}_F . Thus, we can identify multiple coarse directions for explorations where we can propose multiple paths for detailed path planning and path selection. FisherRF [31] uses $\mathbf{H}''[\mathbf{w}|D^{train}]$ as an approximation for the observed Fisher Information $\mathcal{I}(\mathbf{w})$ by computing the Hessians on the training set. This is also known as empirical Fisher Information, whose limitations have been widely discussed by Kunstner *et al.* [35] and Marten *et al.* [43]. In most scenarios, this is a reluctant design choice because the distribution of $\mathbf{x} \sim p(\mathbf{x})$ is unknown (i.e., the distribution of all possible images). However, $\mathbf{x} \sim p(\mathbf{x})$ in our case is tractable because it represents the possible locations where we can take an observation for the environment, a.k.a. the free space of our map. Therefore, unlike FisherRF [31] as described in Eq. 3, we propose to use Monte-Carlo sampling to compute the Fisher Information of the current model

$$\mathcal{I}(\mathbf{w}) = \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} [\mathbf{H}''[\mathbf{y}|\mathbf{x}, \mathbf{w}]] \simeq \sum_{k=1}^N \mathbf{H}''[\mathbf{y}_k|\mathbf{x}_k, \mathbf{w}], \quad (6)$$

where \mathbf{x}_k is drawn from a uniform distribution of camera poses in the free space of the current map which we use to approximate p . Besides, we also uniformly initialize 3D Gaussians in the space, which will be subsequently updated with rendering losses for visited regions. In Fig. 2, we show the relationship between PSNR and EIG on sampled poses in the Cantwell scene. The result aligns with our expectation that EIG should decrease as PSNR increases. We also show some cases on the lower-left and upper-right of the scatter plot. For the lower left capture, even though the left part is poorly reconstructed, most of the view is occupied by a textureless wall, leading to a low EIG. For the upper right capture, even though the scene has a moderate reconstruction, the content is rich, so our algorithm returns a high EIG.

Finally, we compute paths towards each pose in \mathcal{T}_F with the A* algorithm [24] using the occupancy map, selecting which path to follow as described in Section 3.4. The path can be defined as an ordered set of camera poses from the current location x_t at exploration step t to the frontier points x_T^j from the long horizon planning.

$$\mathcal{P}_j = \{x_{t+1}^j, \dots, x_T^j\} \quad (7)$$

3.4. Short-term Path Selection with Localization Uncertainty

The EIG for 3D Gaussian parameters along a path \mathcal{P}_j can be computed as the sum over the path of the following term [31]:

$$\text{EIG}_{\mathcal{P}_j, i}(\mathbf{x}_i) = \text{tr}(\mathbf{H}''[\mathbf{y}_i|\mathbf{x}_i, \mathbf{w}] \mathcal{I}_{\mathcal{P}_j, i}(\mathbf{w})^{-1}) \quad (8)$$

where $\mathcal{I}_{\mathcal{P}_j, i}(\mathbf{w})$ takes the mutual information along the path into account as follows

$$\mathcal{I}_{\mathcal{P}_j, i}(\mathbf{w}) = \mathbf{H}''[\mathbf{w}] + \sum_{\mathbf{x}_t \in \mathcal{P}_j, t < i} \mathbf{H}''[\mathbf{w}|\mathbf{x}_t]. \quad (9)$$

where $\mathbf{H}''[\mathbf{x}]$ is short for $-\nabla^2 \log p(w)$ for clarity. If solely maximizing the EIG, the robot will be more likely to explore unvisited regions. However, exploring regions that have not been well reconstructed also means the agent would have the risk of worse localization accuracy due to noise and ambiguities in the unreconstructed regions during pose optimization. The cost of localization must thus be considered during path planning to balance the importance of exploring new environments with maintaining localization accuracy. We propose to use Fisher Information as a measurement for the localization uncertainty that is also necessary for effective path planning for active Mapping algorithms. During optimization, we essentially optimize on the logarithmic mapping of $\tau_i \triangleq \text{Log}(\mathbf{x}_i)$ of our camera pose. By the Cramér–Rao bound, the covariance of $\tau_i \in \mathfrak{se}(3)$ can be lower-bounded with the inverse of Fisher Information matrix $\mathcal{I}(\tau_i)$:

$$\text{Cov}(T(\hat{\tau}_i)) \geq \mathcal{I}(\tau_i)^{-1} \quad (10)$$

where $T(\tau_i)$ is an unbiased estimator for τ solved by iteratively optimizing photo-metric loss. Hence, we can define the localization cost C_{local} at a pose \mathbf{x}_i in terms of τ_i as:

$$C_{local}(\tau_i) = \log \det(\nabla_{\tau_i} f(\tau_i, \mathbf{w})^T \nabla_{\tau_i} f(\tau_i, \mathbf{w})) \quad (11)$$

Matsuki *et al.* [44] computed the Jacobians of the camera pose with respect to the mean and covariances of each gaussian $\frac{\partial \boldsymbol{\mu}_I}{\partial \mathbf{x}}$ and $\frac{\partial \boldsymbol{\Sigma}_I}{\partial \mathbf{x}}$. However, we need to compute the

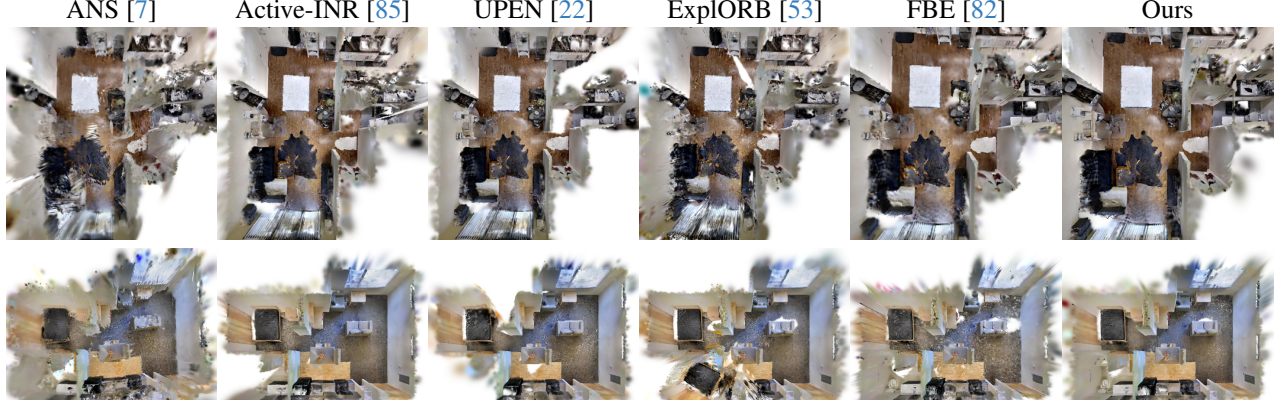


Figure 4. **Qualitative Comparison for Final Scene Reconstruction on Gibson Dataset** Greigsville (top) and Ribera (bottom) scenes. We provide top-down rendering for different methods. Note that UPEN and Active-INR use GT pose in this visualization.

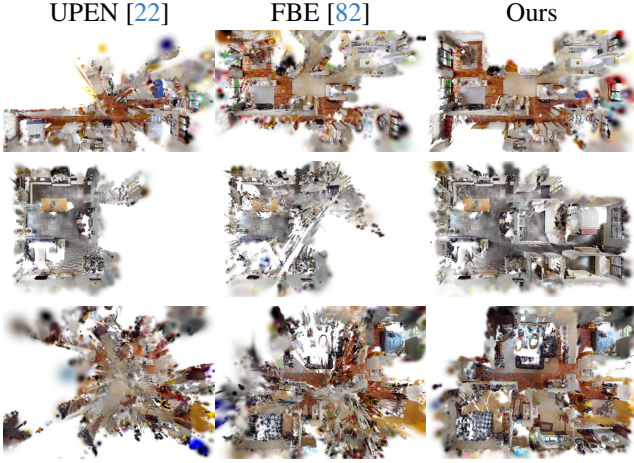


Figure 5. **Qualitative Comparison for Final Scene Reconstruction on Habitat-Matterport 3D Dataset** mscxX4KEBcB (top), oPj9qMxrDEa (middle) and QKGMrrurUVbk (bottom) scenes. We provide top-down rendering for different methods.

Jacobian of τ_i with respect to the rendering output:

$$\nabla_{\tau_i} f(\tau_i, \mathbf{w}) = \frac{\partial f(\tau_i, \mathbf{w})}{\partial \tau_i} = \begin{bmatrix} \frac{\partial f(\tau_i, \mathbf{w})}{\partial \mu_I} & \frac{\partial f(\tau_i, \mathbf{w})}{\partial \Sigma_I} \end{bmatrix} \begin{bmatrix} \frac{\mathcal{D}\mu_C}{\mathcal{D}\tau_i} \\ \frac{\mathcal{D}\mathbf{W}}{\mathcal{D}\tau_i} \end{bmatrix} \quad (12)$$

Unlike 3DGS parameters, our computation of the localization cost can be directly computed on the full Fisher Information matrix without using the Laplace approximation [12, 42]. That is because the full Jacobian $\nabla_{\tau_i} f(\tau_i, \mathbf{w})$ is linear to the number of Gaussians and can be efficiently computed with our CUDA implementation.

Without loss of generality, the path of exploration can be selected by minimizing the total cost for all viewpoints \mathbf{x}_i along a path \mathcal{P}_j :

$$\arg \min_{\mathcal{P}_j} \sum_{\mathbf{x}_i \in \mathcal{P}_j} C_{local}(\text{Log}(\mathbf{x}_i)) - \eta \log(\text{EIG}_{\mathcal{P}_j, i}(\mathbf{x}_i)) \quad (13)$$

where η is a hyper-parameter controlling the importance between EIG and localization accuracy. The agent can then explore the environment with planned path \mathcal{P} . Our active mapping system constantly updates the map, and we replan using our active path planning algorithm if we detect the agent is getting close to a possible obstacle or upon reaching the end of the previously selected path.

4. Experiments

4.1. Experimental Set-up

Dataset Following previous methods [22, 85], our algorithm is evaluated in the Habitat Simulator [73] on the Gibson [81] and Habitat-Matterport 3D (HM3D) [56] datasets comprised of indoor scenes reconstructed from scans of real houses. We use all the scenes in the val split for the Gibson dataset. We also evaluate our method on single-floor scenes whose scene ID can be found in the supplementary. We adopt the default start point in the Habitat Simulator as the starting point for exploration in each scene. The total number of steps for each experiment is 2000. The system takes color and depth images at the resolution of 800x800 and outputs a discrete action at each step. The action space consists of MOVE FORWARD by 5cm, TURN LEFT, and TURN RIGHT by 5°. The field of view (FOV) is set to 90° vertically and horizontally. Please refer to the supplement for more details about the evaluation split and other hyper-parameters.

Metrics We evaluated our method using the Peak-signal-to-noise ratio (PSNR), Structural Similarity Index Measure (SSIM) [79], Learned Perceptual Image Patch Similarity

Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Depth MAE \downarrow	RMSE ATE \downarrow	Coverage (m^2) \uparrow
ANS	16.34	0.6818	0.3923	0.3886	0.1105	10.49
Active-INR (gt)	22.66	0.7652	0.2164	0.1528	-	9.20
UPEN (gt)	21.31	0.7325	0.2714	0.1696	-	8.79
UPEN	16.44	0.6678	0.4134	0.4841	0.5158	8.58
ExplORB	18.99	0.7175	0.3994	0.2664	0.2296	9.00
FBE	21.45	0.7618	0.2126	0.1028	0.1680	10.64
Ours	23.28	0.8067	0.2507	0.0696	0.0226	10.88

Table 1. Quantitative Evaluation on Reconstruction Quality and Tracking Accuracy on Gibson dataset

Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Depth MAE \downarrow	RMSE ATE \downarrow	Coverage (m^2) \uparrow
UPEN (gt)	15.58	0.5175	0.3936	0.4548	-	12.69
UPEN	12.23	0.4795	0.5157	0.7356	0.4393	10.72
FBE	15.80	0.5952	0.4392	0.4085	1.2004	15.69
Ours	19.86	0.7127	0.4122	0.1666	0.0336	18.48

Table 2. Quantitative Evaluation on Reconstruction Quality and Tracking Accuracy on HM3D dataset

(LPIPS) [89] for RGB rendering and mean absolute error (MAE) for depth rendering as metrics for scene reconstruction quality. We calculate these metrics using 2000 points uniformly sampled from the movement plane of the agent in the scene, discarding any points that are not navigable. We argue that the rendering quality reflects both reconstruction quality and pose accuracy because high tracking accuracy would help the training of the 3D Gaussian Splatting model. Meanwhile, misaligned poses will lead to misaligned rendering at test time thus leading to inferior results. Following previous approaches [7, 22], we also use coverage in m^2 as an evaluation metric. To evaluate the pose estimation accuracy, we use the root mean squared average tracking error (RMSE ATE), but as the trajectories for each method are different, the RMSE ATE should only be considered along with other metrics such as coverage.

Baselines We compare to two exploration methods which assume ground truth pose: UPEN [22] and Active Neural Mapping (active-INR) [85]. UPEN chooses trajectories to reduce the epistemic uncertainty of the occupancy prediction produced by an ensemble of models. Active-INR aims to minimize the prediction robustness against random weight perturbation of its signed distance field scene representation. We also compare our method with Active Neural SLAM (ANS) [7], explORB [53] and Frontier Based Exploration (FBE) [82] without ground truth pose provided. ANS [7] learns to predict a map and estimate a global goal based on

it. ExplORB [53] computes the Fisher Information of the Hessian on its pose graph optimization.

To compare the rendering quality fairly, we run all the baselines using the MonoGS [44] backend for reconstruction. We run UPEN and FBE online, but for ANS, active-INR, and ExplORB, we record and playback trajectories obtained using their source code. Because the forward step size for ANS is much larger than for our method, we interpolate the trajectory so that the forward step size matches our method’s to make the steps comparable. For ExplORB, since the official implementation is based on MoveBase, which uses velocity commands, we sample the trajectory at 5 Hz. We also found that ANS, active-INR, and UPEN failed on some scenes due to localization failure. ANS produces a pose estimate (using information from noisy pose sensors not provided to our pipeline), so we set the pose estimate of the MonoGS backend to the one from ANS. As active-INR and UPEN do not produce a pose estimate, we evaluate them using the ground-truth pose.

4.2. Comparison Against Previous Methods

Table 1 shows the results of our method and the baselines for exploration in scenes from the Gibson dataset [81], and Table 2 shows the results of our method and some baselines on HM3D [6]. Our method outperforms the baselines on all metrics. We further qualitatively compare the reconstruction qualities after active exploration in Fig. 4 and Fig. 5, and the trajectories in Fig. 6. Our method does not have major

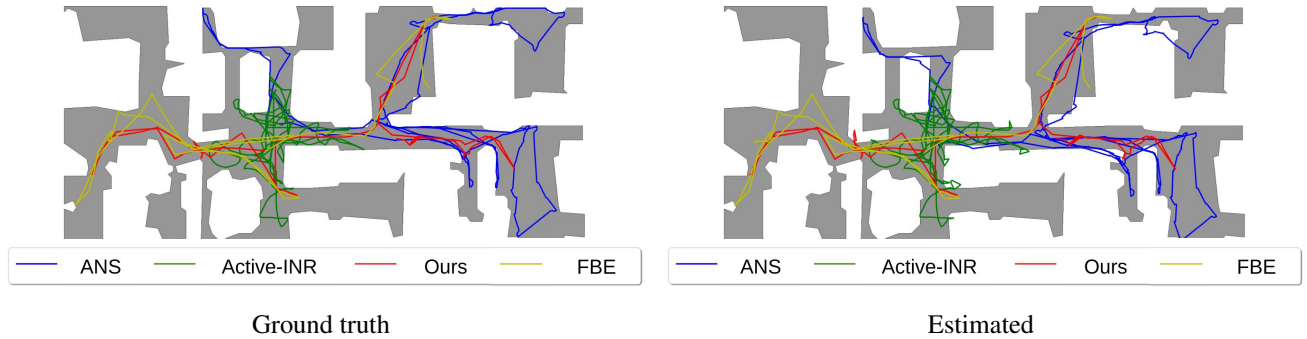


Figure 6. **Qualitative Comparison for Trajectories** on Cantwell from Gibson Dataset. Left is the ground truth trajectory, right is the estimated trajectory.

Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Depth MAE \downarrow	RMSE ATE \downarrow	Coverage (m^2) \uparrow
w.o. LLM & Localization Uncertainty	16.15	0.6550	0.6193	0.3409	0.2478	16.06
w.o. LLM	16.94	0.6799	0.5847	0.2887	0.1694	17.26
Ours	19.86	0.7127	0.4122	0.1666	0.0336	18.48

Table 3. **Ablation Study of Localization Uncertainty Term on Scenes from the HM3D Dataset.** We compare our method with and without the localization uncertainty term to validate that including it provides improvements on both localization and reconstruction

errors from failed localization, and we have fewer gaps in the scenes than other methods. For example, in the Ribera scene, all methods except for us and FBE miss the bathroom at the bottom left, and FBE misses more of the area around the sofa than us. For the trajectories, we show the estimated and ground truth trajectories for the Cantwell scene from the Gibson dataset. Cantwell is a relatively large and challenging scene, so it is suitable for showing the differences between methods. We show only a few baselines to keep the figure legible. We can see that Active-INR stays in a smaller area than the rest of the methods. ANS often goes close to walls, whereas FBE and our method are generally more toward the center of corridors or rooms. ANS also extends further into rooms than these methods. FBE does not go to the region on the bottom right, and ANS does not go to the region on the bottom left, whereas our method efficiently visits most areas of the scene.

4.3. Ablation Study

To validate the effectiveness of localization uncertainty and LLM, we performed an ablative study of each component in our system on the Gibson dataset in Table 3. As can be seen, the average trajectory error is much lower with the localization uncertainty than without. In addition, after adding LLM, the model performs better in the overall rendering. The ATE slightly increases after adding LLM, but this could be attributed to the larger area traversed by the robot.

5. Conclusion

We present an active mapping system that can autonomously explore an environment, the first method to do so using a 3D Gaussian representation without ground truth camera poses. By rendering maps using our scene representation, we can seamlessly query a multimodal Large Language Model for long-horizon planning to leverage its knowledge of scene layouts. This allows the robot to efficiently explore the scene while avoiding setting unreachable navigation goals. We then plan an optimal path to the long-horizon goal using our information-driven path proposal and selection algorithm, which balances the information gain with respect to the map against the localization uncertainty, ensuring that the robot captures the geometric details of the scene for a high-quality reconstruction. We evaluate our method on scenes from the Gibson [81] and Habitat-Matterport 3D [56] datasets in terms of the rendering quality, coverage, and average tracking error. We compare to recent exploration methods with [22, 85] and without [7, 53] ground-truth pose, using the same scene representation for reconstruction as us for fairness, and show that our method has superior performance. To enable our method to support more robotics applications, future work could extend it to take advantage of our 3D scene representation and consider movement with higher degrees of freedom (DOF) than the currently supported 3DOF. Incorporating semantic features [55, 68, 91] to allow for grounding language to the scene would also enable many robotics and computer vision applications.

Acknowledgements Financial support by the grants NSF FRR 2220868, and ONR N00014-22-1-2677 is gratefully acknowledged.

References

- [1] Arash Asgharivaskasi and Nikolay Atanasov. Semantic octree mapping and shannon mutual information computation for robot exploration. *IEEE Transactions on Robotics*, 39(3): 1910–1928, 2023. 1
- [2] Anthony Brohan, Yevgen Chebotar, Chelsea Finn, Karol Hausman, Alexander Herzog, Daniel Ho, Julian Ibarz, Alex Irpan, Eric Jang, Ryan Julian, et al. Do as i can, not as i say: Grounding language in robotic affordances. In *Conference on robot learning*, pages 287–318. PMLR, 2023. 3
- [3] Carlos Campos, Richard Elvira, Juan J Gómez Rodríguez, José MM Montiel, and Juan D Tardós. Orb-slam3: An accurate open-source library for visual, visual-inertial, and multi-map slam. *IEEE Transactions on Robotics*, 37(6):1874–1890, 2021. 2
- [4] Luca Carlone, Jingjing Du, Miguel Kaouk Ng, Basilio Bona, and Marina Indri. Active slam and exploration with particle filters using kullback-leibler divergence. *Journal of Intelligent & Robotic Systems*, 75:291–311, 2014. 2, 3
- [5] Henry Carrillo, Ian Reid, and José A Castellanos. On the comparison of uncertainty criteria for active slam. In *2012 IEEE International Conference on Robotics and Automation*, pages 2080–2087. IEEE, 2012. 1
- [6] Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from rgb-d data in indoor environments. *International Conference on 3D Vision (3DV)*, 2017. 7
- [7] Devendra Singh Chaplot, Dhiraj Gandhi, Saurabh Gupta, Abhinav Gupta, and Ruslan Salakhutdinov. Learning to explore using active neural slam. *arXiv preprint arXiv:2004.05155*, 2020. 2, 3, 6, 7, 8, 20
- [8] Jiaqi Chen, Bingqian Lin, Xinmin Liu, Lin Ma, Xiaodan Liang, and Kwan-Yee K Wong. Affordances-oriented planning using foundation models for continuous vision-language navigation. *arXiv preprint arXiv:2407.05890*, 2024. 2
- [9] Jiaqi Chen, Bingqian Lin, Ran Xu, Zhenhua Chai, Xiaodan Liang, and Kwan-Yee Wong. Mapgpt: Map-guided prompting with adaptive path planning for vision-and-language navigation. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9796–9810, 2024. 2
- [10] Tao Chen, Saurabh Gupta, and Abhinav Gupta. Learning exploration policies for navigation. In *International Conference on Learning Representations*, 2019. 3
- [11] Yongbo Chen, Shoudong Huang, and Robert Fitch. Active slam for mobile robots with area coverage and obstacle avoidance. *IEEE/ASME Transactions on Mechatronics*, 25(3): 1182–1192, 2020. 1
- [12] Erik Daxberger, Agustinus Kristiadi, Alexander Immer, Runa Eschenhagen, Matthias Bauer, and Philipp Hennig. Laplace redux—effortless Bayesian deep learning. In *NeurIPS*, 2021. 4, 6
- [13] Vishnu Sashank Dorbala, James F Mullen Jr, and Dinesh Manocha. Can an embodied agent find your “cat-shaped mug”? IIm-based zero-shot object navigation. *IEEE Robotics and Automation Letters*, 2023. 3
- [14] Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. Palm-e: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*, 2023. 2
- [15] Danny Driess, Fei Xia, Mehdi S. M. Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, Wenlong Huang, Yevgen Chebotar, Pierre Sermanet, Daniel Duckworth, Sergey Levine, Vincent Vanhoucke, Karol Hausman, Marc Toussaint, Klaus Greff, Andy Zeng, Igor Mordatch, and Pete Florence. Palm-e: An embodied multimodal language model. In *arXiv preprint arXiv:2303.03378*, 2023. 3
- [16] Yilun Du, Mengjiao Yang, Pete Florence, Fei Xia, Ayzaan Wahid, Brian Ichter, Pierre Sermanet, Tianhe Yu, Pieter Abbeel, Joshua B Tenenbaum, et al. Video language planning. *arXiv preprint arXiv:2310.10625*, 2023. 2, 3
- [17] Jakob Engel, Thomas Schöps, and Daniel Cremers. Lsd-slam: Large-scale direct monocular slam. In *European conference on computer vision*, pages 834–849. Springer, 2014. 2
- [18] Jakob Engel, Vladlen Koltun, and Daniel Cremers. Direct sparse odometry. *IEEE transactions on pattern analysis and machine intelligence*, 40(3):611–625, 2017. 2
- [19] Martin Ester, Hans-Peter Kriegel, Jörg Sander, Xiaowei Xu, et al. A density-based algorithm for discovering clusters in large spatial databases with noise. In *kdd*, pages 226–231, 1996. 14
- [20] 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. 1
- [21] Ziyue Feng, Huangying Zhan, Zheng Chen, Qingan Yan, Xiangyu Xu, Changjiang Cai, Bing Li, Qilun Zhu, and Yi Xu. Naruto: Neural active reconstruction from uncertain target observations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21572–21583, 2024. 3
- [22] Georgios Georgakis, Bernadette Bucher, Anton Arapin, Karl Schmeckpeper, Nikolai Matni, and Kostas Daniilidis. Uncertainty-driven planner for exploration and navigation. In *ICRA*, 2022. 2, 3, 6, 7, 8, 20, 21
- [23] Lily Goli, Cody Reading, Silvia Sellán, Alec Jacobson, and Andrea Tagliasacchi. Bayes’ Rays: Uncertainty quantification in neural radiance fields. *CVPR*, 2024. 1, 3
- [24] Peter E Hart, Nils J Nilsson, and Bertram Raphael. A formal basis for the heuristic determination of minimum cost paths. *IEEE transactions on Systems Science and Cybernetics*, 4(2): 100–107, 1968. 5
- [25] Jiawei Hou, Tianyu Wang, Tongying Pan, Shouyan Wang, Xiangyang Xue, and Yanwei Fu. TaMMA: Target-driven multi-subscene mobile manipulation. In *8th Annual Conference on Robot Learning*, 2024. 1
- [26] Chenguang Huang, Oier Mees, Andy Zeng, and Wolfram Burgard. Visual language maps for robot navigation. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 10608–10615. IEEE, 2023. 1, 2, 3

- [27] Huajian Huang, Longwei Li, Cheng Hui, and Sai-Kit Yeung. Photo-slam: Real-time simultaneous localization and photo-realistic mapping for monocular, stereo, and rgb-d cameras. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. 1
- [28] Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. In *International conference on machine learning*, pages 9118–9147. PMLR, 2022. 3
- [29] Wenlong Huang, Chen Wang, Ruohan Zhang, Yunzhu Li, Jiajun Wu, and Li Fei-Fei. Voxposer: Composable 3d value maps for robotic manipulation with language models. *arXiv preprint arXiv:2307.05973*, 2023. 3
- [30] Shahram Izadi, David Kim, Otmar Hilliges, David Molyneaux, Richard Newcombe, Pushmeet Kohli, Jamie Shotton, Steve Hodges, Dustin Freeman, Andrew Davison, et al. Kinectfusion: real-time 3d reconstruction and interaction using a moving depth camera. In *Proceedings of the 24th annual ACM symposium on User interface software and technology*, pages 559–568, 2011. 2
- [31] Wen Jiang, Boshu Lei, and Kostas Daniilidis. Fisherrf: Active view selection and uncertainty quantification for radiance fields using fisher information. In *ECCV*, 2024. 1, 3, 4, 5
- [32] Nikhil Keetha, Jay Karhade, Krishna Murthy Jatavallabhula, Gengshan Yang, Sebastian Scherer, Deva Ramanan, and Jonathon Luiten. Splatmap: Splat, track & map 3d gaussians for dense rgb-d slam. 2024. 1, 2
- [33] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4):1–14, 2023. 1, 3, 14
- [34] Ayoung Kim and Ryan M Eustice. Active visual slam for robotic area coverage: Theory and experiment. *The International Journal of Robotics Research*, 34(4-5):457–475, 2015. 2
- [35] Frederik Kunstner, Philipp Hennig, and Lukas Balles. Limitations of the empirical fisher approximation for natural gradient descent. *Advances in neural information processing systems*, 32, 2019. 5
- [36] 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*, 7(4):12070–12077, 2022. 1, 3
- [37] Peiqi Liu, Yaswanth Orru, Chris Paxton, Nur Muhammad Mahi Shafiullah, and Lerrel Pinto. Ok-robot: What really matters in integrating open-knowledge models for robotics. *arXiv preprint arXiv:2401.12202*, 2024. 1
- [38] Xu Liu, Ankit Prabhu, Fernando Cladera, Ian D Miller, Lifeng Zhou, Camillo J Taylor, and Vijay Kumar. Active metric-semantic mapping by multiple aerial robots. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3282–3288. IEEE, 2023. 2
- [39] Iker Lluvia, Elena Lazkano, and Ander Ansuategi. Active mapping and robot exploration: A survey. *Sensors*, 21(7):2445, 2021. 3
- [40] Yuxing Long, Wenzhe Cai, Hongcheng Wang, Guanqi Zhan, and Hao Dong. Instructnav: Zero-shot system for generic instruction navigation in unexplored environment. *CoRL*, 2024. 3
- [41] Yuxing Long, Xiaoqi Li, Wenzhe Cai, and Hao Dong. Discuss before moving: Visual language navigation via multi-expert discussions. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 17380–17387. IEEE, 2024. 3
- [42] David J. C. MacKay. Bayesian Interpolation. *Neural Computation*, 4(3):415–447, 1992. 4, 6
- [43] James Martens. New insights and perspectives on the natural gradient method. *Journal of Machine Learning Research*, 21(146):1–76, 2020. 5
- [44] Hidenobu Matsuki, Riku Murai, Paul H. J. Kelly, and Andrew J. Davison. Gaussian Splatting SLAM. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. 1, 2, 3, 5, 7, 14
- [45] Reihaneh Mirjalili, Michael Krawez, and Wolfram Burgard. Fm-loc: Using foundation models for improved vision-based localization. In *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1381–1387. IEEE, 2023. 3
- [46] David Morilla-Cabello, Lorenzo Mur-Labadia, Ruben Martinez-Cantin, and Eduardo Montijano. Robust fusion for bayesian semantic mapping. In *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 76–81. IEEE, 2023. 2
- [47] Raul Mur-Artal and Juan D Tardós. Orb-slam2: An open-source slam system for monocular, stereo, and rgb-d cameras. *IEEE transactions on robotics*, 33(5):1255–1262, 2017. 2
- [48] Raul Mur-Artal, Jose Maria Martinez Montiel, and Juan D Tardos. Orb-slam: a versatile and accurate monocular slam system. *IEEE transactions on robotics*, 31(5):1147–1163, 2015.
- [49] Richard A Newcombe, Steven J Lovegrove, and Andrew J Davison. Dtam: Dense tracking and mapping in real-time. In *2011 international conference on computer vision*, pages 2320–2327. IEEE, 2011. 2
- [50] OpenAI. Gpt-4o system card, 2024. <https://openai.com/index/gpt-4o-system-card/>. 4
- [51] Xuran Pan, Zihang Lai, Shiji Song, and Gao Huang. Activenet: Learning where to see with uncertainty estimation. In *European Conference on Computer Vision*, pages 230–246. Springer, 2022. 1, 3
- [52] Andrej Pázmán. Foundations of optimum experimental design. (*No Title*), 1986. 3
- [53] Julio A Placed, Juan J Gómez Rodríguez, Juan D Tardós, and José A Castellanos. Explorb-slam: Active visual slam exploiting the pose-graph topology. In *Iberian Robotics conference*, pages 199–210. Springer, 2022. 2, 6, 7, 8, 20
- [54] 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. *IEEE Transactions on Robotics*, 2023. 3

- [55] Minghan Qin, Wanhua Li, Jiawei Zhou, Haoqian Wang, and Hanspeter Pfister. Langsplat: 3d language gaussian splatting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20051–20060, 2024. 1, 8
- [56] Santhosh Kumar Ramakrishnan, Aaron Gokaslan, Erik Wijmans, Oleksandr Maksymets, Alexander Clegg, John M Turner, Eric Undersander, Wojciech Galuba, Andrew Westbury, Angel X Chang, Manolis Savva, Yili Zhao, and Dhruv Batra. Habitat-matterport 3d dataset (HM3d): 1000 large-scale 3d environments for embodied AI. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2021. 2, 6, 8
- [57] Krishan Rana, Jesse Haviland, Sourav Garg, Jad Abou-Chakra, Ian D Reid, and Niko Sünderhauf. Sayplan: Grounding large language models using 3d scene graphs for scalable task planning. *CoRR*, 2023. 3
- [58] María L Rodríguez-Arévalo, José Neira, and José A Castellanos. On the importance of uncertainty representation in active slam. *IEEE Transactions on Robotics*, 34(3):829–834, 2018. 1
- [59] Julius Rückin, Federico Magistri, Cyrill Stachniss, and Marija Popović. An informative path planning framework for active learning in uav-based semantic mapping. *IEEE Transactions on Robotics*, 39(6):4279–4296, 2023. 2
- [60] Arjun P S, Andrew Melnik, and Gora Chand Nandi. Cognitive planning for object goal navigation using generative AI models. In *NeurIPS 2024 Workshop on Open-World Agents*, 2024. 3
- [61] M.J. Schervish. *Theory of Statistics*. Springer New York, 2012. 3
- [62] Pierre Sermanet, Tianli Ding, Jeffrey Zhao, Fei Xia, Debidatta Dwibedi, Keerthana Gopalakrishnan, Christine Chan, Gabriel Dulac-Arnold, Sharath Maddineni, Nikhil J Joshi, et al. Robovqa: Multimodal long-horizon reasoning for robotics. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 645–652. IEEE, 2024. 2
- [63] Nur Muhammad Mahi Shafiullah, Chris Paxton, Lerrel Pinto, Soumith Chintala, and Arthur Szlam. CLIP-fields: Weakly supervised semantic fields for robotic memory. In *ICRA2023 Workshop on Pretraining for Robotics (PT4R)*, 2023. 1
- [64] Dhruv Shah, Michael Robert Equi, Błażej Osiniński, Fei Xia, Brian Ichter, and Sergey Levine. Navigation with large language models: Semantic guesswork as a heuristic for planning. In *Conference on Robot Learning*, pages 2683–2699. PMLR, 2023. 3
- [65] Claude Elwood Shannon. A mathematical theory of communication. *The Bell system technical journal*, 27(3):379–423, 1948. 3
- [66] Jianxiong Shen, Adria Ruiz, Antonio Agudo, and Francesc Moreno-Noguer. Stochastic neural radiance fields: Quantifying uncertainty in implicit 3d representations. In *2021 International Conference on 3D Vision (3DV)*, pages 972–981. IEEE, 2021. 1, 3
- [67] Jianxiong Shen, Antonio Agudo, Francesc Moreno-Noguer, and Adria Ruiz. Conditional-flow nerf: Accurate 3d modelling with reliable uncertainty quantification. In *European Conference on Computer Vision*, pages 540–557. Springer, 2022. 1, 3
- [68] Jin-Chuan Shi, Miao Wang, Hao-Bin Duan, and Shao-Hua Guan. Language embedded 3d gaussians for open-vocabulary scene understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5333–5343, 2024. 1, 8
- [69] Chan Hee Song, Jiaman Wu, Clayton Washington, Brian M Sadler, Wei-Lun Chao, and Yu Su. Llm-planner: Few-shot grounded planning for embodied agents with large language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2998–3009, 2023. 3
- [70] Cyrill Stachniss, Dirk Hahnel, and Wolfram Burgard. Exploration with active loop-closing for fastslam. In *2004 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)(IEEE Cat. No. 04CH37566)*, pages 1505–1510. IEEE, 2004. 2, 3
- [71] Zezhou Sun, Banghe Wu, Cheng-Zhong Xu, Sanjay E Sarma, Jian Yang, and Hui Kong. Frontier detection and reachability analysis for efficient 2d graph-slam based active exploration. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 2051–2058. IEEE, 2020. 2, 3
- [72] Niko Sünderhauf, Jad Abou-Chakra, and Dmitry Miller. Density-aware nerf ensembles: Quantifying predictive uncertainty in neural radiance fields. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9370–9376. IEEE, 2023. 1, 3
- [73] Andrew Szot, Alexander Clegg, Eric Undersander, Erik Wijmans, Yili Zhao, John Turner, Noah Maestre, Mustafa Mukadam, Devendra Singh Chaplot, Oleksandr Maksymets, et al. Habitat 2.0: Training home assistants to rearrange their habitat. *Advances in neural information processing systems*, 34:251–266, 2021. 6
- [74] Octo Model Team, Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey Hejna, Tobias Kreiman, Charles Xu, et al. Octo: An open-source generalist robot policy. *arXiv preprint arXiv:2405.12213*, 2024. 2
- [75] Zachary Teed and Jia Deng. Droid-slam: Deep visual slam for monocular, stereo, and rgb-d cameras. *Advances in neural information processing systems*, 34:16558–16569, 2021. 2
- [76] Eleftherios Triantafyllidis, Filippos Christianos, and Zhibin Li. Intrinsic language-guided exploration for complex long-horizon robotic manipulation tasks. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 7493–7500. IEEE, 2024. 3
- [77] Darko Trivun, Edin Šalaka, Dinko Osmanković, Jasmin Velagić, and Nedim Osmić. Active slam-based algorithm for autonomous exploration with mobile robot. In *2015 IEEE International Conference on Industrial Technology (ICIT)*, pages 74–79. IEEE, 2015. 2, 3
- [78] Teresa Vidal-Calleja, Andrew J Davison, Juan Andrade-Cetto, and David William Murray. Active control for single camera slam. In *Proceedings 2006 IEEE International Conference on Robotics and Automation, 2006. ICRA 2006.*, pages 1930–1936. IEEE, 2006. 1
- [79] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to

- structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. [6](#)
- [80] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022. [4](#)
- [81] Fei Xia, Amir R. Zamir, Zhi-Yang He, Alexander Sax, Jitendra Malik, and Silvio Savarese. Gibson Env: real-world perception for embodied agents. In *Computer Vision and Pattern Recognition (CVPR), 2018 IEEE Conference on*. IEEE, 2018. [2](#), [6](#), [7](#), [8](#)
- [82] Brian Yamauchi. A frontier-based approach for autonomous exploration. In *Proceedings 1997 IEEE International Symposium on Computational Intelligence in Robotics and Automation CIRA'97.'Towards New Computational Principles for Robotics and Automation'*, pages 146–151. IEEE, 1997. [2](#), [6](#), [7](#), [20](#), [21](#)
- [83] Chi Yan, Delin Qu, Dan Xu, Bin Zhao, Zhigang Wang, Dong Wang, and Xuelong Li. Gs-slam: Dense visual slam with 3d gaussian splatting. In *CVPR*, 2024. [1](#)
- [84] Dongyu Yan, Jianheng Liu, Fengyu Quan, Haoyao Chen, and Mengmeng Fu. Active implicit object reconstruction using uncertainty-guided next-best-view optimization, 2023. [1](#)
- [85] Zike Yan, Haoxiang Yang, and Hongbin Zha. Active neural mapping. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10981–10992, 2023. [1](#), [2](#), [3](#), [6](#), [7](#), [8](#), [14](#), [20](#)
- [86] Naoki Yokoyama, Sehoon Ha, Dhruv Batra, Jiuguang Wang, and Bernadette Bucher. Vlfm: Vision-language frontier maps for zero-shot semantic navigation. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 42–48. IEEE, 2024. [3](#)
- [87] Justin Yu, Kush Hari, Kishore Srinivas, Karim El-Refai, Adam Rashid, Chung Min Kim, Justin Kerr, Richard Cheng, Muhammad Zubair Irshad, Ashwin Balakrishna, et al. Language-embedded gaussian splats (legs): Incrementally building room-scale representations with a mobile robot. [1](#)
- [88] Huangying Zhan, Jiyang Zheng, Yi Xu, Ian Reid, and Hamid Rezatofighi. Activermap: Radiance field for active mapping and planning. *arXiv preprint arXiv:2211.12656*, 2022. [2](#), [3](#)
- [89] 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. [7](#)
- [90] Gengze Zhou, Yicong Hong, and Qi Wu. Navgpt: Explicit reasoning in vision-and-language navigation with large language models. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(7):7641–7649, 2024. [3](#)
- [91] Shijie Zhou, Haoran Chang, Sicheng Jiang, Zhiwen Fan, Zehao Zhu, Dejia Xu, Pradyumna Chari, Suyu You, Zhangyang Wang, and Achuta Kadambi. Feature 3dgs: Supercharging 3d gaussian splatting to enable distilled feature fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21676–21685, 2024. [1](#), [8](#)

A. Appendix

In this supplementary material, we discuss things we left over in our main paper due to page constraints. We provide more background information about 3DGS and the derivative with respect to localization parameters in Sec. A.1. We discussed more implementation details in Sec. A.2 and provided the prompt we use for the Multimodal LLM and an example of the interaction in Sec. A.5 to help reproduce our results. **The source code of this project will be made public soon.** We also include detailed quantitative and qualitative results in Sec. A.3 and Sec. A.4.

A.1. Additional Background Information

For the completeness of our method, we also provide the key definition for the 3D Gaussian Splatting backbone [33] and 3D Gaussian SLAM [44]. In 3D Gaussian Splatting [33], the rendered pixel color is calculated by composing all 3D Gaussians projected in a tile.

$$\hat{C}(\mathbf{r}) = \sum_{i=1}^{N_s} T_i (1 - \exp(-\sigma_i \delta_i)) \mathbf{c}_i \quad (14)$$

$$\alpha_i = \exp \left(- \sum_{j=1}^{i-1} \sigma_j \delta_j \right) (1 - \exp(-\sigma_i \delta_i)) \quad (15)$$

$\delta_i = t_{i+1} - t_i$ represents the distance between adjacent samples, and N_s indicates the number of samples. \mathbf{c}_i is the color of each 3D Gaussian given the current view direction \mathbf{d} and σ_i is given by evaluating a 2D Gaussian with covariance Σ .

The Jacobian of the localization parameters are defined as:

$$\frac{\partial \boldsymbol{\mu}_C}{\partial \mathbf{x}} = [\mathbf{I} \quad -\boldsymbol{\mu}_C^\times] \text{ and } \frac{\partial \mathbf{W}}{\partial \mathbf{x}} = \begin{bmatrix} \mathbf{0} & -\mathbf{W}_{:,1}^\times \\ \mathbf{0} & -\mathbf{W}_{:,2}^\times \\ \mathbf{0} & -\mathbf{W}_{:,3}^\times \end{bmatrix}, \quad (16)$$

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

Unlike PSNR, the EIG can be computed without ground truth images, making it possible to perform view selection during exploration.

A.2. Implementation Details

The 2D occupancy map’s resolution is 5cm. For each single frontier pixel on the 2D map, we add 200 3D Gaussians, which are uniformly distributed in the 3D cube above it. Other parameters like color, opacity, and scale are generated uniformly between 0 and 1. When there are frontiers on the 2D map, we choose the next frontier to be explored by the area of each frontier divided by the distance. When

no frontier exists, we select the top 20% of Gaussians with the highest score. These Gaussians are grouped using DBSCAN [19]. The largest cluster is selected for candidate pose generation. Candidates are uniformly sampled in the range between 0.3m to 1m, facing towards the selected position. Only the poses in free space are kept for path-level selection. The importance factor η in Eq. 13 is set to 5 across all experiments. The source code for this project will be made public no later than the publication of this paper.

We compute the Expected Information Gain (EIG) for each global candidate and use A* to plan a path for each of them. In order to prevent a twisted path, we consider locations 0.15m (3 pixels) away from the current robot position as neighbors and set the robot width to 3 pixels for collision check. However, the path planned by A* might have redundant waypoints, causing unnecessary turns for the robot. Therefore, we smooth the path by finding short-cuts. Specifically, for each waypoint w_i , if the path between waypoint w_{i+2} and w_i is collision-free, then we remove the intermediate waypoint w_{i+1} from the path. Finally, we use a greedy follower for motion planning. If the angle between the heading direction of the robot and the relative next waypoint is larger than 5° , then we turn left or right to decrease the angle. Otherwise, we choose the forward action to approach the next waypoint. In such a way, we get a sequence of actions $\{a_i\}_{i=1}^T$ for each path.

Given a sequence of actions $\{a_i\}_{i=1}^T$ for each path, we use forward dynamics to compute the future camera poses $\{c_i\}_{i=1}^T$. Initially, we use an intermediate variable $\mathbf{H}_{\text{obs}}'' \triangleq \mathbf{H}''[\mathbf{w}^*]$ to help compute expected information gain along the path. For each camera pose x_i , we compute its pose Hessian $\mathbf{H}_{\text{pose}}''$ and the current model Hessian matrix $\mathbf{H}_{\text{cur}}'' \triangleq \mathbf{H}''[\mathbf{y}|x_i, \mathbf{w}^*]$. $\mathbf{H}_{\text{cur}}''$ is then accumulated, and we update $\mathbf{H}_{\text{obs}}''$ to evaluate the remaining poses on the path. We select the path that minimizes the objective given by Eq. 13 for execution.

The compute of Fisher Information in our experiments takes 30 frames per second on our Nvidia L40 GPU. The time on path planning will not significantly impact the robot’s motion if the path is planned ahead of time and runs in parallel as the robot moves. The regularizer λ in Eq. 4 is 10^{-6} across all experiments. The weight parameters η in Eq. 13 is 0.2 across all experiments.

A.3. Results for Each Scene in Gibson and HM3D Dataset

Following previous literature [85], we use the following scenes for Gibson Dataset: Greigsville, Denmark, Cantwell, Eudora, Pablo, Ribera, Swormville, Eastville, Elmira. For HM3D we use the following scenes: DBjEchFg4oq, mscxX4KEBcB, QKGMrurUVbk, oPj9qMxrDEa, CETmJJqkhcK. The detailed results for each scene on

each evaluation metric are presented as bar plots in Fig. 7 for Gibson and Fig. 8 for HM3D. We also present qualitative comparisons on testing views from the Gibson dataset in Fig. 9.

A.4. Detailed version of qualitative results

We provide larger versions of the qualitative rendering comparisons from the main paper, Fig. 11 shows the Gibson scenes and Fig. 12 shows the HM3D scenes.

A.5. Example of Using Multimodal LLM

We provide an example of our interaction with the multimodal LLM in Fig 13.

```

1  <system> You are an AI assistant that can
2    analyze images and plan a long-term goal
3    for the exploration task of a ground
4    robot.
5  You will be given a bird-eye view image of a
6    scene.
7  The goal is to plan a long-term exploration
8    mission for a robot to traverse the area.
9  The robot's task is to explore the terrain
10   efficiently, identifying important
11   areas, potential obstacles, and
12   unvisited areas.
13 Please analyze the image and select a
14   long-term goal from the candidates for
15   the robot to explore the area.
16 Empty space doesn't always mean they are
17   unvisited regions, sometimes it's just
18   outside the floor plan of this scene.
19 We are allowed to explore a total of
20   <TOTAL\_STEPS> steps and this is step
21   <STEP\_ID>.
22 Therefore, it's better to select a space that
23   is close to the visited regions but
24   still unvisited and not behind the walls.
25 The current location of the robot is marked
26   with the blue star(*) marker.
27 The last frontier you selected is marked with
28   a yellow diamond shape.
29 The visited path is painted as green lines in
30   the image.
31 Note that you don't have to select the
32   closest point to the robot, but the
33   point that is most likely to be
34   unvisited and important to explore.
35 As you can see, there are <NUM\_FRAME>
36   candidate points to select from.
37 They are numbered from 0 to <NUM\_FRAME - 1>
38   in red color.
39 If you find all the goals are not necessary
40   to explore and we should instead focus
41   on improving existing reconstruction,
42   please give -1 in the `target` entry of
43   the JSON.
44 Please provide a detailed exploration plan
45   and select an exploration target with
46   reasons in the JSON format as shown
47   below.

```

```

20   ...
21   {
22     "target": 2, "reason": "The target is located
23       at an unvisited region of the image and
24       seems to be an unvisited bedroom"
25   }
26   ...
27   Do not cut off the JSON and generate the full
28   JSON.
29   </system>
30
31   <user>: I have a bird-eye view image of a
32     scene. The goal is to plan a long-term
33     exploration mission for a robot to
34     traverse the area. Please analyze the
35     attached image and provide the
36     exploration plan first and then an
37     exploration target in the specified JSON
38     format.
39   </user>

```

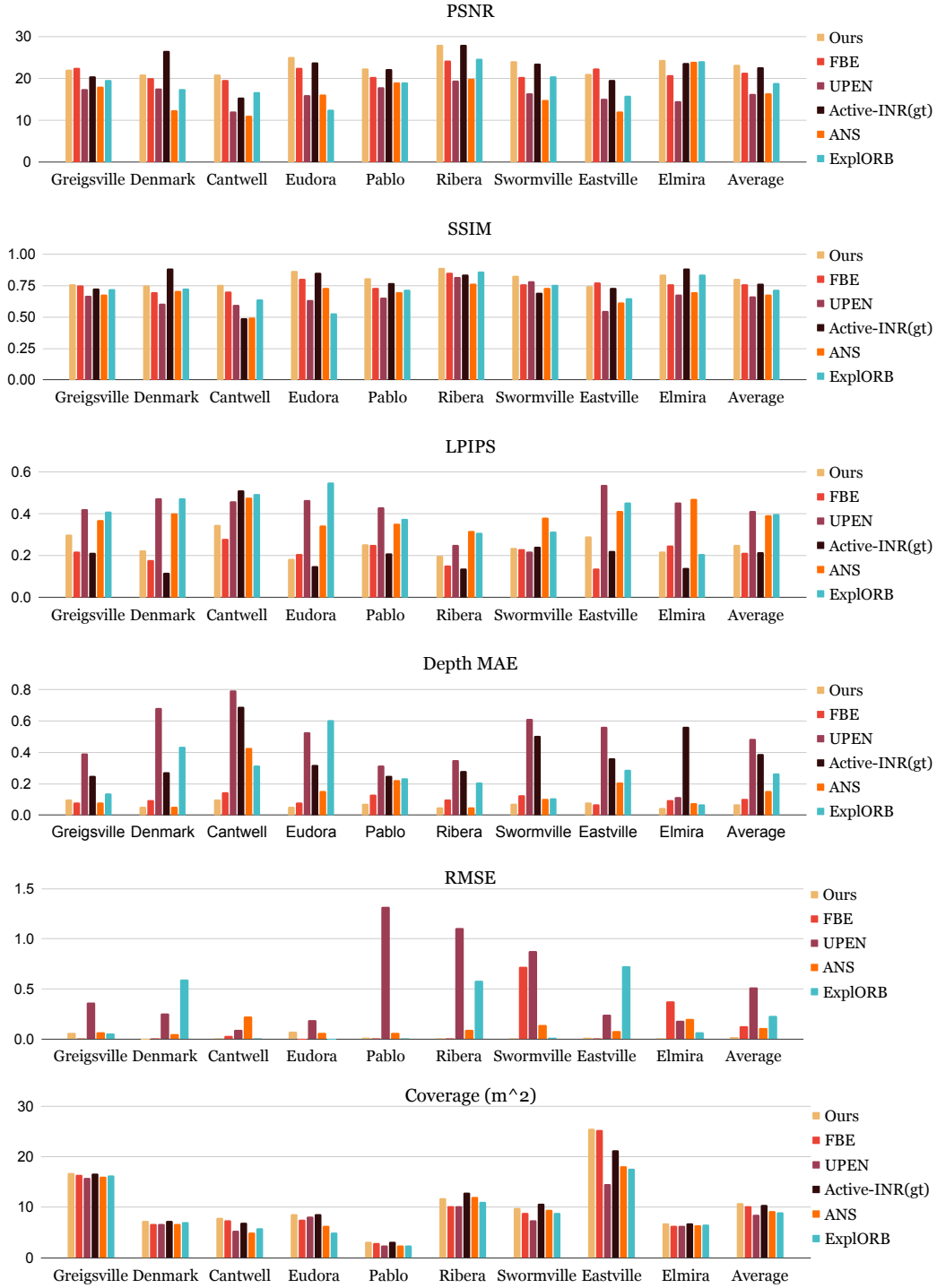
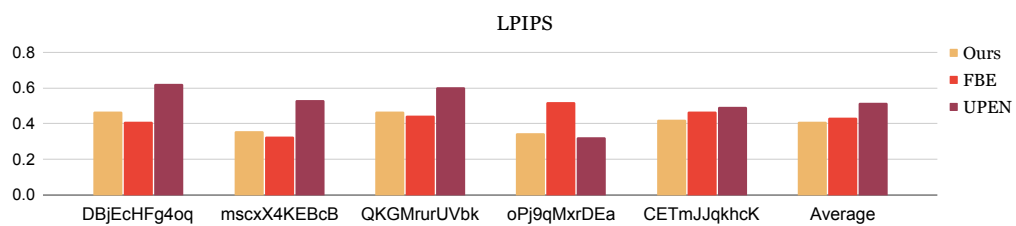
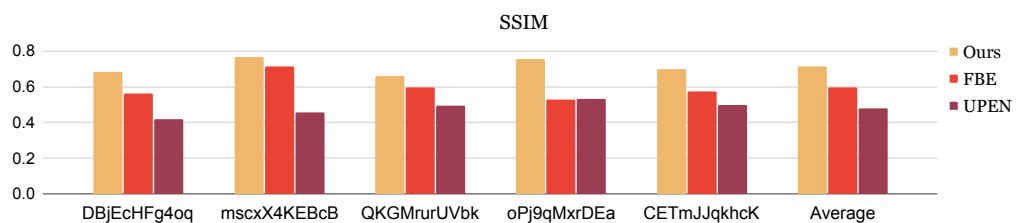
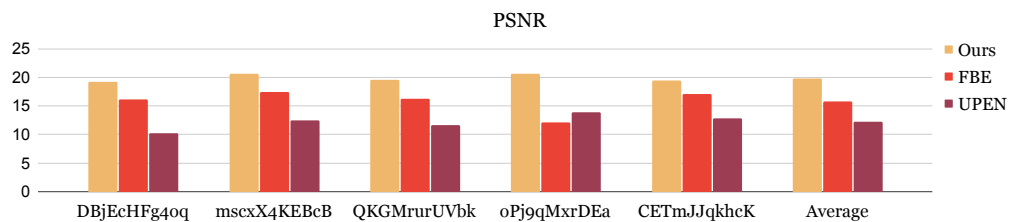
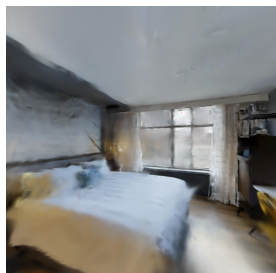


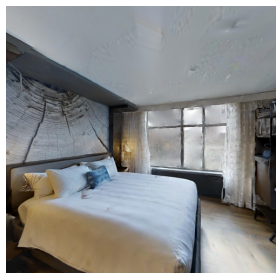
Figure 7. Per-scene results on Gibson Dataset



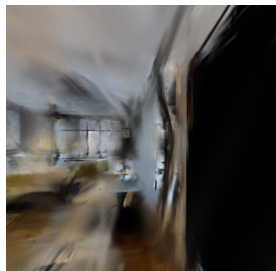
ANS



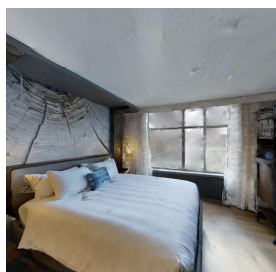
Active-INR (gt)



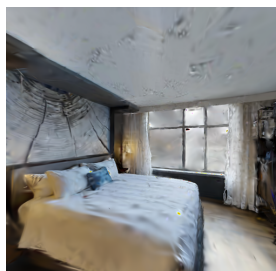
UPEN



ExplORB



FBE



Ours

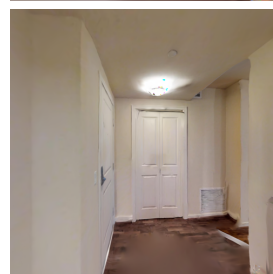
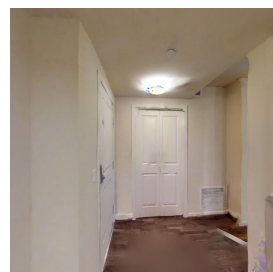
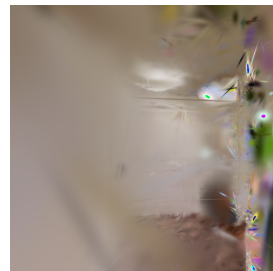
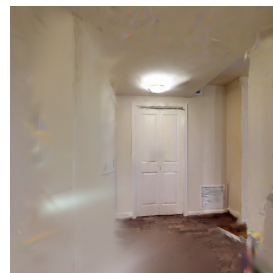
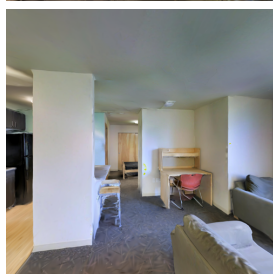
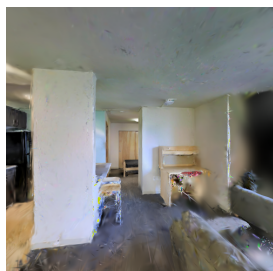
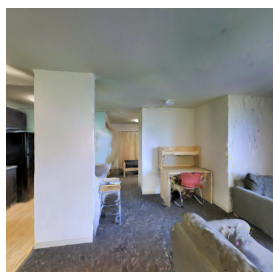
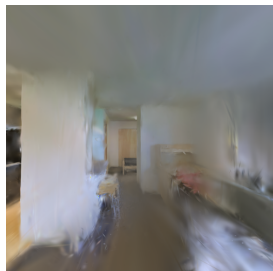
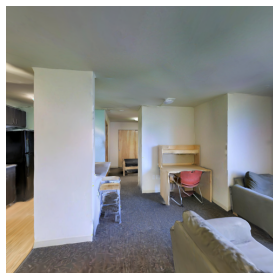
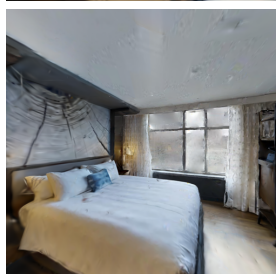


Figure 9. **Test Rendering Qualitative Visualization on Gibson Dataset** All the renderings are from the test view of the Gibson dataset.

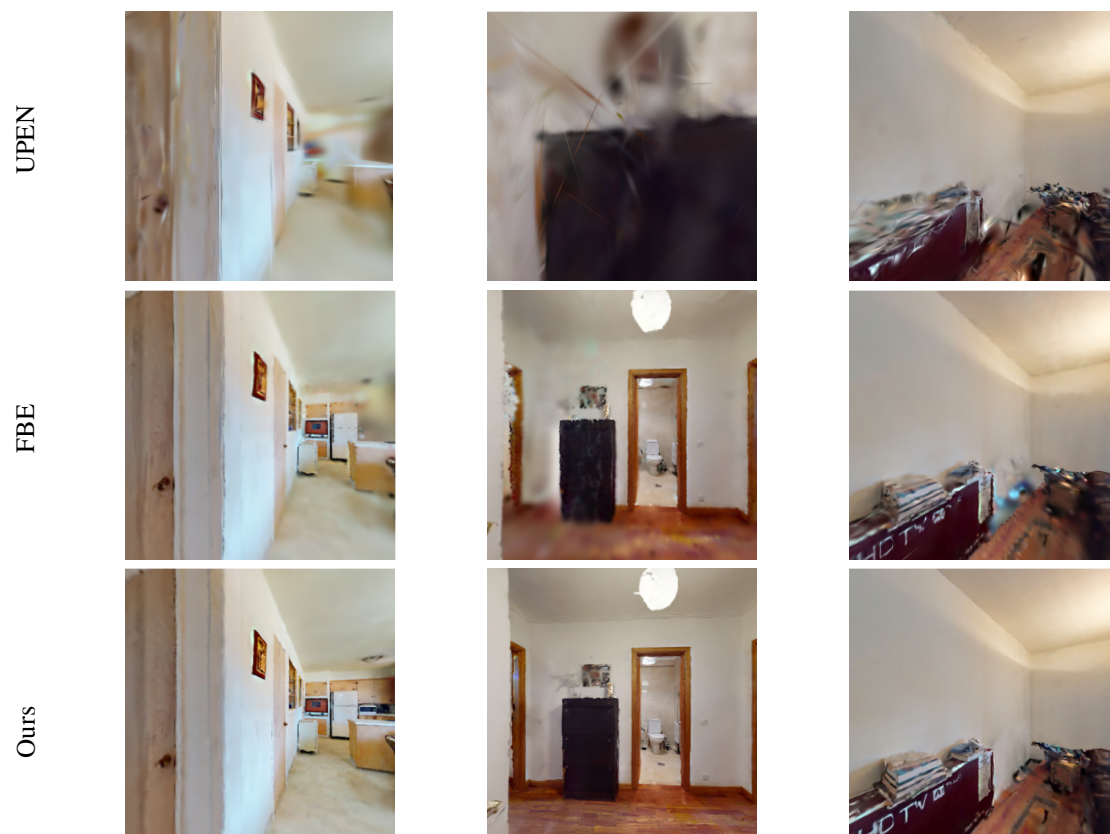


Figure 10. **Test Rendering Qualitative Visualization on Habitat-Matterport 3D Dataset** All the renderings are from the test view of the Habitat-Matterport 3D dataset.

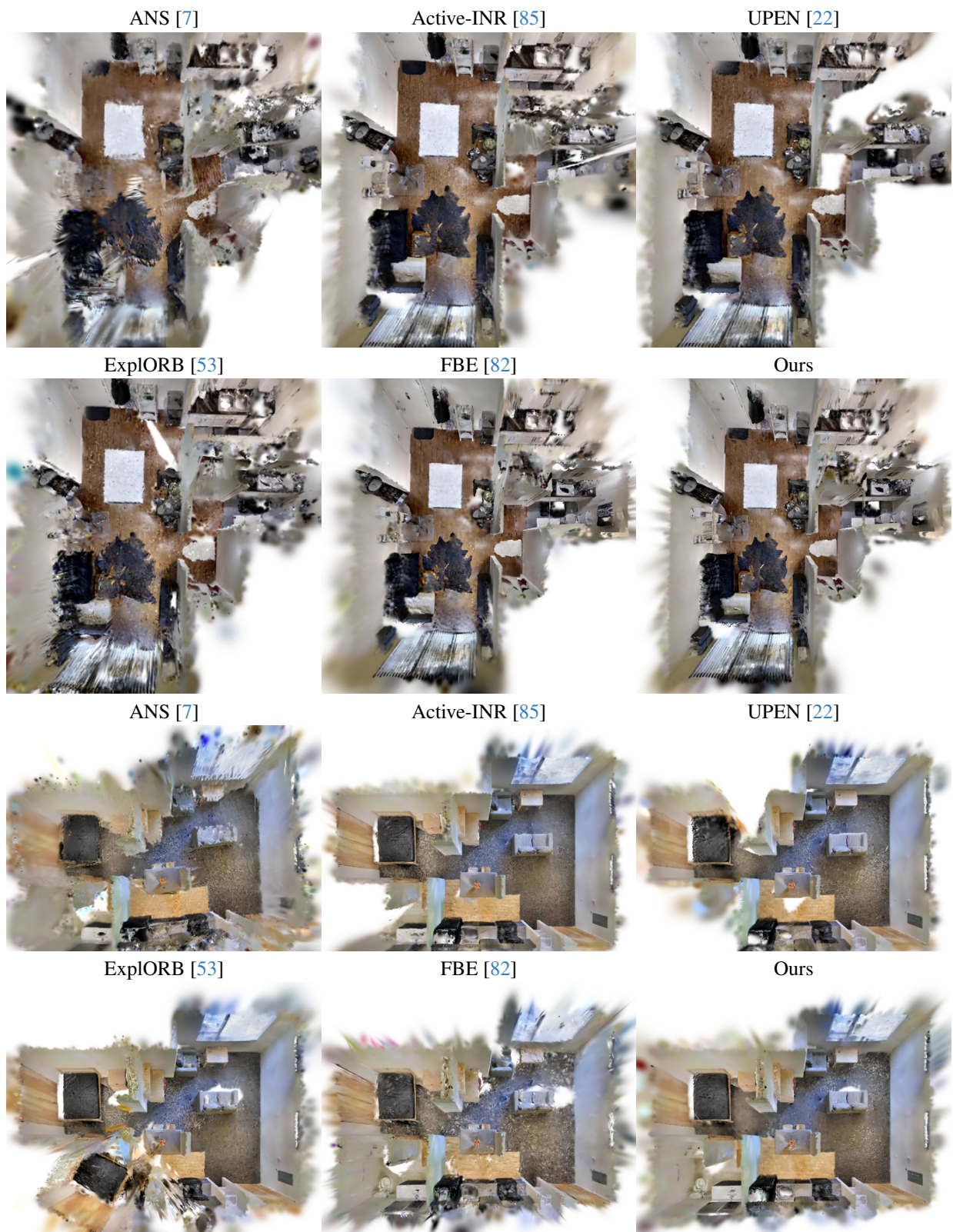


Figure 11. **Qualitative Comparison for Final Scene Reconstruction on Gibson Dataset** Greigsville (top) and Ribera (bottom) scenes. We provide top-down rendering for different methods. Note that UPEN and Active-INR use GT pose in this visualization.

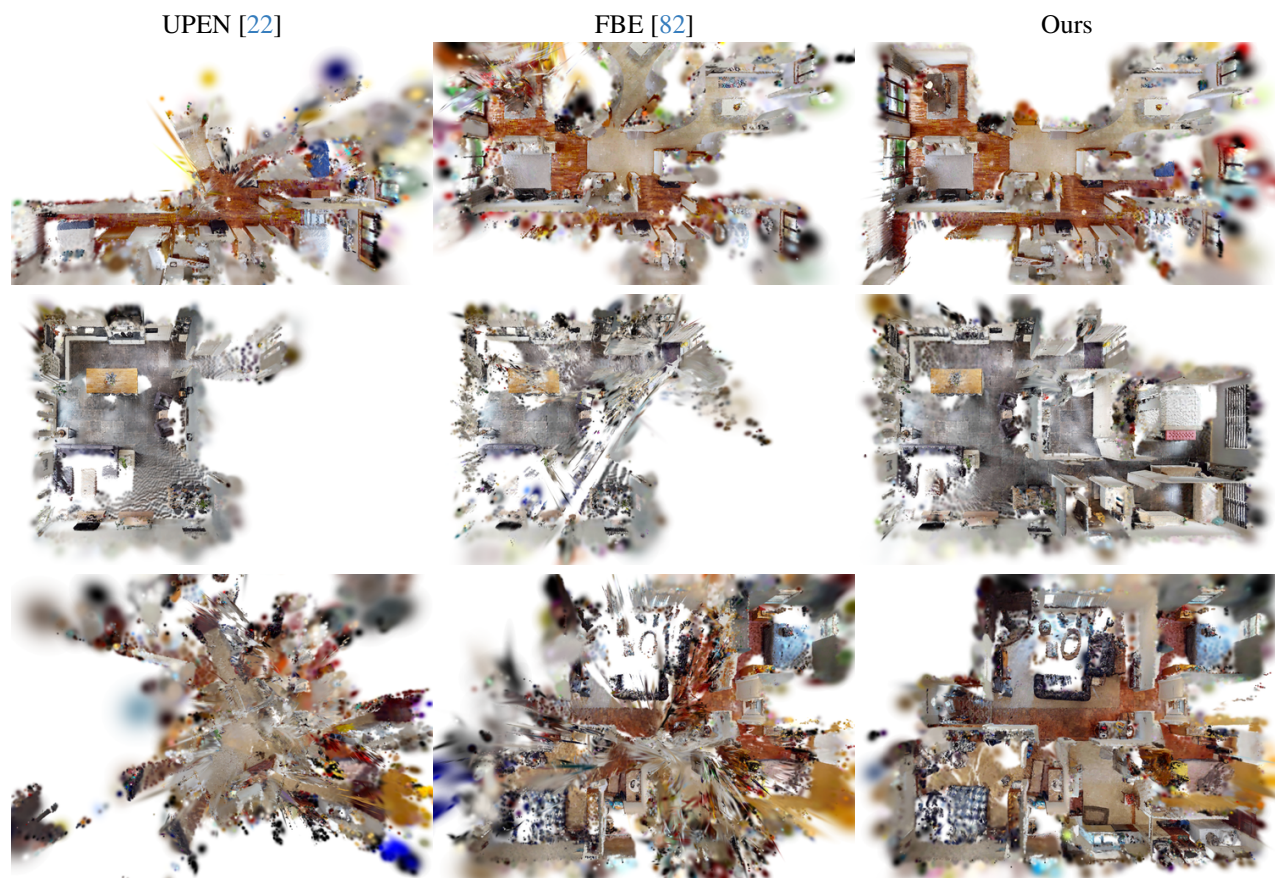



Figure 12. **Qualitative Comparison for Final Scene Reconstruction on Habitat-Matterport 3D Dataset** mscxX4KEBcB (top), oPj9qMxrDEa (middle) and QKGMrurUVbk (bottom) scenes. We provide top-down rendering for different methods.



System: You are an AI assistant that can analyze images and plan a long-term goal for the exploration task of a ground robot.

You will be given a bird-eye view image of a scene.

The goal is to plan a long-term exploration mission for a robot to traverse the area.

The robot's task is to explore the terrain efficiently, identifying important areas, potential obstacles, and unvisited areas.

Please analyze the image and select a long-term goal from the candidates for the robot to explore the area.

Empty space doesn't always mean they are unvisited regions, sometimes it's just outside the floor plan of this scene.

We are allowed to explore a total of 2000 steps and this is step 2.

Therefore, it's better to select a space that is close to the visited regions but still unvisited and not behind the walls.

The current location of the robot is marked with the blue star(*) marker.

The last frontier you selected is marked with a yellow diamond(◆) shape.

The visited path is painted as green lines in the image.

Note that you don't have to select the closest point to the robot, but the point that is most likely to be unvisited and important to explore.

As you can see, there are 4 candidate points to select from.

They are numbered from 0 to 3 in red color.

If you find all the goals are not necessary to explore and we should instead focus on improving existing reconstruction, please give -1 in the 'target' entry of the JSON.

Please provide a detailed exploration plan and select an exploration target with reasons in the JSON format as shown below.


```

'''
{
  "target": 2, "reason": "The target is located at an unvisited region of the
image and seems to be an unvisited bedroom"
}
'''

```

Do not cut off the JSON and generate the full JSON.

User: I have a bird-eye view image of a scene. The goal is to plan a long-term exploration mission for a robot to traverse the area. Please analyze the attached image and provide the exploration plan first and then an exploration target in the specified JSON format.




****Exploration Plan:****

- **Current Situation Analysis:****
 - The robot is currently located at the blue star marker.
 - The previously selected frontier, marked by a yellow diamond, might still have unexplored neighboring areas.
 - So far, the robot has traversed 132 steps and can take a total of 2000 steps.
- **Observation of Candidates:****
 - ****Candidate 0:**** Located towards the top of the scene. It appears to be in a separated area from the main path by walls and could be challenging to access directly.
 - ****Candidate 1:**** Positioned near the robot's previous path, to the left. The space seems to be partly blocked but may lead to a new room or area.
 - ****Candidate 2:**** Situated to the bottom right, close to the previously marked yellow diamond. This area seems unvisited with potential space for exploration.
 - ****Candidate 3:**** On the right side, aligned with the robot's current heading, showing open space but appearing more distant from the recent paths.
- **Selection Reasoning:****
 - ****Candidate 2**** is favored for exploration as it is likely unvisited, near the previous frontier, and potentially opens up new unexplored areas. It is a strategic choice for efficient area coverage without straying too far from the recently explored paths, reducing the risk of re-exploring already visited spaces.

```

'''json
{
  "target": 2,
  "reason": "Candidate 2 is favored as it is near the previously marked frontier, indicating potential unexplored area relevant to the mission objectives. Continuing exploration here aligns with efficient coverage of the space."
}
'''

```

Figure 13. An Example about a Dialog with the Multi-modal LLM. The chat box on the top is the question from our method, and the bottom is the response from GPT-4o, which we parsed to get a long-term goal for our path proposal. We provided the contextual information in both the textual prompt and the image.