

NEXTBESTPATH: EFFICIENT 3D MAPPING OF UNSEEN ENVIRONMENTS

Anonymous authors

Paper under double-blind review

ABSTRACT

This work addresses the problem of active 3D mapping, where an agent must find an efficient trajectory to exhaustively reconstruct a new scene. Previous approaches mainly predict the next best view near the agent’s location, which is prone to getting stuck in local areas. Additionally, existing indoor datasets are insufficient due to limited geometric complexity and inaccurate ground truth meshes. To overcome these limitations, we introduce a novel dataset AiMDoom with a map generator for the Doom video game, enabling to better benchmark active 3D mapping in diverse indoor environments. Moreover, we propose a new method we call next-best-path (NBP), which predicts long-term goals rather than focusing solely on short-sighted views. The model jointly predicts accumulated surface coverage gains for long-term goals and obstacle maps, allowing it to efficiently plan optimal paths with a unified model. By leveraging online data collection, data augmentation and curriculum learning, NBP significantly outperforms state-of-the-art methods on both the existing MP3D dataset and our AiMDoom dataset, achieving more efficient mapping in indoor environments of varying complexity.

1 INTRODUCTION

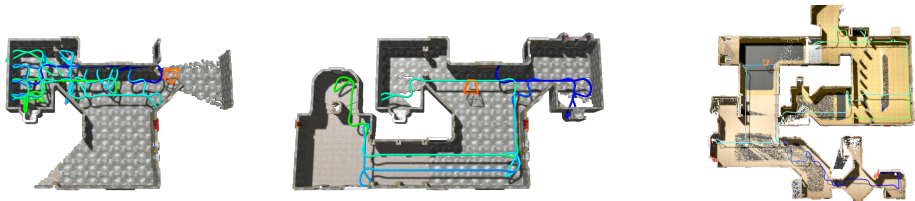
Autonomous 3D mapping of new scenes holds substantial importance for vision, robotics, and graphics communities, with applications including digital twins. In this paper, we focus on the problem of active 3D mapping, where the goal is for an agent to find the shortest possible trajectory to scan the entire surface of a new scene using a depth sensor.

This task is extremely challenging as the agent has to identify an efficient trajectory without knowing the scene in advance. Existing works can be broadly categorized into rule-based and learning-based approaches. Rule-based approaches, such as frontier-based exploration (FBE) (Yamauchi, 1997), utilize heuristic rules to select optimal frontiers at the boundaries of the already-known space for the next movement. Though being simple and generalizable, they fail to leverage data priors to develop more efficient planning strategies. To address this, learning-based methods, often referred to as next-best-view planning (NBV), train parametric policies for action prediction. Although NBV approaches have demonstrated promising results, most of them only are evaluated on single-object datasets or outdoor scenes (Guédon et al., 2022; Chang et al., 2015; Peralta et al., 2020), ignoring a critical but more difficult setting of indoor environments for active 3D mapping applications.

Existing indoor datasets (Xia et al., 2018; Chang et al., 2017), however, offer limited geometry complexity and often include imperfect ground truth meshes, making them inadequate to fully evaluate model performance in complex indoor environments. In this work, we automatically construct a new indoor dataset called AiMDoom for active 3D mapping. AiMDoom is built upon a map generator for the Doom video game, and features a wide range of indoor settings of four difficulty levels: Simple, Normal, Hard and Insane. As illustrated in Figure 1a, even in relatively simple indoor settings of our dataset, the state-of-the-art NBV approach MACARONS (Guédon et al., 2023) is frequently trapped in a limited area and misses substantial portions of the scene. This limitation arises because most NBV methods only look one step ahead to identify the next best view in neighbouring regions, making it difficult to explore under-reconstructed areas at far distances.

Some recent works (Chen et al., 2024; Feng et al., 2024; Zhan et al., 2022; Georgakis et al., 2022) attempt to overcome this limitation by searching for the next optimal view across a broader range. For

054
 055
 056
 057
 058
 059
 060
 061
 062
 063
 064
 065
 066
 067
 068
 069
 070
 071
 072
 073
 074
 075
 076
 077
 078
 079
 080
 081
 082
 083
 084
 085
 086
 087
 088
 089
 090
 091
 092
 093
 094
 095
 096
 097
 098
 099
 100
 101
 102
 103
 104
 105
 106
 107



(a) MACARONS (simple scene). (b) Our NBP (simple scene). (c) Our NBP (hard scene).

Figure 1: Reconstruction results and trajectories of MACARONS (Guédon et al., 2023) and our NBP model. Guédon et al. (2023) fails to fully map the environment in simple scenes (a), while our NBP model manages to capture the full scene (b), even in much more complex geometry (c).

example, Georgakis et al. (2022) utilizes a strategy that relies on averaging predicted uncertainties at each point along every sampled path, and uses a trained point-goal navigation model. However, training separate uncertainty map prediction and navigation models is less efficient, and the scene uncertainty does not directly align with the ultimate objective of 3D mapping.

Therefore, we further propose a novel approach called next-best-path (NBP) planning, which shifts from NBV approaches that predict a single nearby view, to predicting an optimal path in a unified model. Our model is composed of three key components: a mapping progress encoder, a coverage gain decoder and an obstacle map decoder. The mapping progress encoder efficiently encodes the currently reconstructed point cloud along with the agent’s past trajectory. Based on the encoded representation, the coverage gain decoder predicts a value map over a large spatial range centred on the agent’s current location. Each cell in the map represents the surface coverage gain accumulated along the optimal trajectory from the agent’s location to the cell, which corresponds to the final metric for active mapping. The cell with the highest value score is viewed as a long-term goal. The obstacle map decoder predicts obstacles in both seen and unseen regions by leveraging the agent’s current knowledge of the scene. This allows us to compute the shortest path to the long-term goal while avoiding obstacles. To train the model, we collect data online and iteratively improve the model. We also propose a data augmentation method that exploits a property of shortest paths and a combined curriculum and multitask learning strategy to enhance training efficiency.

We evaluate our methods on the existing indoor benchmark MP3D (Chang et al., 2017) and our dataset AiMDoom. The proposed NBP model significantly outperforms state-of-the-art methods on both datasets from simple (Figure 1b) to more complex indoor environments (Figure 1c).

Our key contributions can be summarized as follows:

- We introduce AiMDoom, the first benchmark to systematically evaluate active mapping in indoor scenes of different levels of difficulties.
- We propose a novel next-best-path approach that jointly predicts long-term goals with optimal reconstruction coverage gains, and obstacle maps for trajectory planning.
- Our approach achieved state-of-the-art results on both the AiMDoom and MP3D datasets.

We will release the dataset, codes and trained models publicly upon acceptance.

2 RELATED WORK

Active Mapping. Active mapping aims to exhaustively reconstruct a 3D scene in the shortest possible time with a moving agent. Unlike SLAM (Chaplot et al., 2020; Placed et al., 2023; Matsuki et al., 2024), which addresses both localization and mapping, active mapping focuses on reconstruction, continuously selecting viewpoints to cover the entire scene, assuming the pose is known. Early methods often relied on frontier-based exploration (FBE) approaches (Yamauchi, 1997). The key idea is to move the agent toward a heuristically selected frontier along the boundary between reconstructed and unknown regions of the scene. Among different strategies (Bircher et al., 2016; Cieslewski et al., 2017; Zhou et al., 2021; Tao et al., 2023) for frontier selection, moving to the nearest frontier serves as a strong baseline. Additionally, there are efforts (Cao et al., 2021; Xu et al., 2024) that combine global FBE and local planning strategies within a hierarchical optimiza-

Table 1: **Comparison between AiMDoom and prior indoor 3D datasets.** Navigation complexity is the maximum ratio of geodesic to euclidean distances between any two navigable locations in the scene. Universal accessibility means whether windows and doors are accessible.

Dataset	Replica	RoboTHOR	MP3D	Gibson (4+ only)	ScanNet	HM3D	AiMDoom (Ours)			
							Simple	Normal	Hard	Insane
Number of scenes	18	75	90	571 (106)	1613	1000	100	100	100	100
Floor space (m ²)	2.19k	3.17k	101.82k	217.99k (17.74k)	39.98k	365.42k	63.33k	134.84k	321.38k	548.85k
Navigation complexity	5.99	2.06	17.09	14.25 (11.90)	3.78	13.31	11.31	18.38	36.05	45.25
Universal accessibility	✗	✗	✗	✗	✗	✗	✓	✓	✓	✓
Easy expansion	✗	✗	✗	✗	✗	✗	✓	✓	✓	✓

[tion framework to enhance exploration](#). However, these FBE-based approaches are heuristic-based and cannot exploit prior learned from data to explore more efficiently, restricting their performance in complex environments.

To address this limitation, learning-based approaches have been explored to select the next-best views (NBV) for efficient 3D mapping. The NBV-based methods train models to select the optimal pose from nearby camera poses (Guédon et al., 2022; 2023; Lee et al., 2023) or from a limited predefined view space such as a hemisphere (Zhan et al., 2022; Lee et al., 2022; Peralta et al., 2020; Zeng et al., 2020; Mendoza et al., 2020). While these methods show promising results to reconstruct single objects, their performance remains limited in large environments. Due to the narrow search space for the next pose, NBV methods behave like a greedy policy and thus can easily get stuck in local regions. To mitigate this, some works Ramrakhya et al. (2022); Chen et al. (2023) use imitation learning to learn from human demonstrates which prioritize unseen exploration but with the cost of heavy labelling. More recently, efforts have been made to enlarge the search range for the next best view (Chen et al., 2024; Ran et al., 2023; Pan et al., 2022; Georgakis et al., 2022). However, these methods are still primarily evaluated on single-object datasets with small moving steps, and often rely on optimizing indirect metrics like reconstruction uncertainty (Georgakis et al., 2022), which are not directly aligned with the goal of exhaustive 3D reconstruction. In this work, we extend the evaluation to more complex indoor environments and also introduce a new surface coverage gain criterion that optimizes the coverage gain along the best trajectory towards a long-term goal.

3D mapping datasets. Existing datasets for 3D mapping mainly focus on single isolated objects such as those in ShapeNet (Chang et al., 2015) and OmniObject3D (Wu et al., 2023), or outdoor scenes (Lu et al., 2023; Hardouin et al., 2020), where the agent only needs to move around the scene to achieve full reconstruction. These datasets are comparatively less complex than indoor environments where the agent must enter into the scene. The indoor scenes contain unique challenges such as dead ends and tight corners, which often force the agent to backtrack without significantly improving its objective.

While some works (Yan et al., 2023; Georgakis et al., 2022; Ramakrishnan et al., 2020) incorporate indoor scene datasets such as Gibson (Xia et al., 2018) and MP3D (Chang et al., 2017), these often exhibit significant limitations. Existing synthetic datasets (Straub et al., 2019; Deitke et al., 2020) often lack scene complexity, whereas real-world scans (Dai et al., 2017; Ramakrishnan et al., 2021), despite offering greater representational fidelity, are constrained by limited structural and map diversity and often suffer from substantial noise artifacts. This lack of reliable datasets prevents comprehensive evaluation in active 3D mapping tasks. In this work, we propose a new dataset - AiMDoom, designed for benchmarking active mapping in indoor environments of different complexities.

3 THE AiMDOOM DATASET

In this section, we introduce **AiMDoom**, a new dataset for **Active 3D Mapping** in complex indoor environments based on the **Doom** video game ¹. As Doom features a wide variety of indoor settings, we use its map generator to create four sets of maps of increasing geometric complexity: Simple, Normal, Hard, and Insane. In the following, we first detail how we built these maps and then discuss the key challenges presented in our AiMDoom dataset.

¹[https://en.wikipedia.org/wiki/Doom_\(franchise\)](https://en.wikipedia.org/wiki/Doom_(franchise))

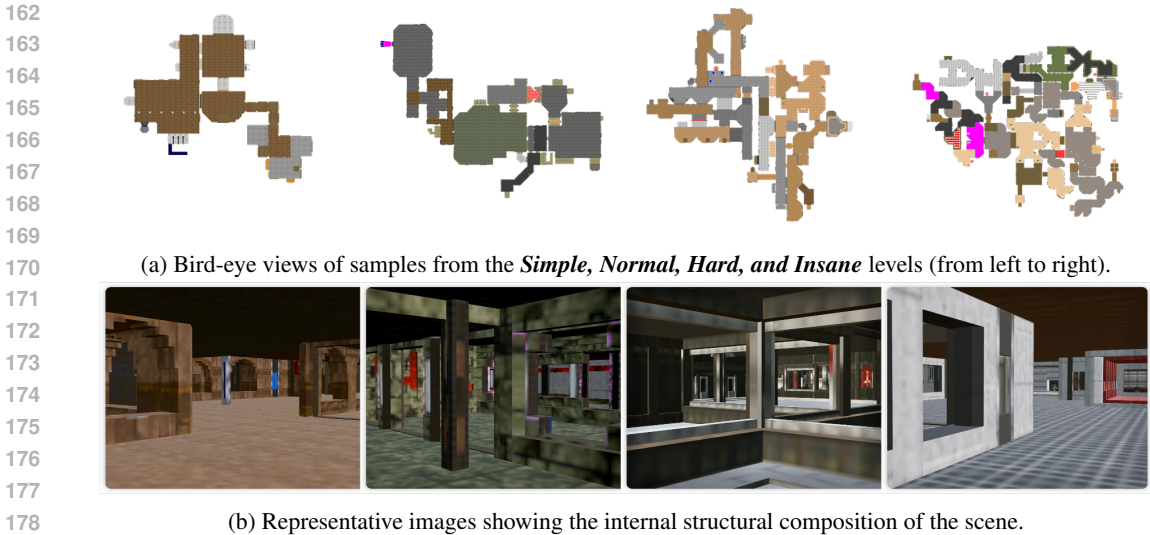


Figure 2: **Maps from our AiMDoom dataset.** The AiMDoom dataset includes four levels of geometric complexity with various textures.

Dataset construction. We used the open-source software Obsidian² to automatically generate Doom maps as our indoor environments. Four sets of hyperparameters are proposed to control architectural complexity and texture styles in Obsidian. By varying these hyperparameters, we produced maps categorized into Simple, Normal, Hard and Insane difficulty levels. Each difficulty level is made of 100 maps with 70 for training and 30 for evaluation.

The maps include doors and windows, all of which are configured to be open. This allows the agent to see and pass through the doors and windows. We converted the maps to the widely used OBJ format, and used Blender (Community, 2018) to consolidate the texture images of each map into a single texture image. This makes the maps compatible with Pytorch3D (Ravi et al., 2020) and Open3D (Zhou et al., 2018). Further details are presented in the supplementary material.

Key challenges. The AiMDoom dataset presents three key challenges for active 3D mapping. Firstly, the dataset features environments with intricate geometries and layouts as shown in Figure 2, making it challenging to determine the optimal exploration direction for effective mapping. Secondly, the maps have small doors and narrow corridors, requiring careful path planning to navigate. Finally, the map diversity requires the reconstruction system to generalize across different scenes. Table 1 compares AiMDoom with existing indoor 3D datasets (Straub et al., 2019; Deitke et al., 2020; Chang et al., 2017; Dai et al., 2017; Xia et al., 2018; Ramakrishnan et al., 2021), highlighting our dataset’s strengths in scene area and navigation complexity.

We will release the dataset along with a comprehensive toolkit to generate the data, which enables easy expansion of the dataset for future research.

4 LEARNING ACTIVE 3D MAPPING

4.1 OVERVIEW

Problem definition. Active 3D mapping aims to control an agent, such as an unmanned aerial vehicle (UAV) or wheeled robot, to efficiently and exhaustively reconstruct a 3D scene. The agent starts at a random location within the scene, and at each time step t , it receives an RGB-D image I_t and must predict the next one $c_t = (c_t^{\text{pos}}, c_t^{\text{rot}})$ in the immediate surrounding of the agent. Here, c_t^{pos} denotes the position coordinates, and c_t^{rot} represents the orientation angles. The agent continually predicts successive c_t until a predefined time limit T is reached. The final output is the reconstructed 3D point cloud of the explored environment.

²<https://obsidian-level-maker.github.io/>

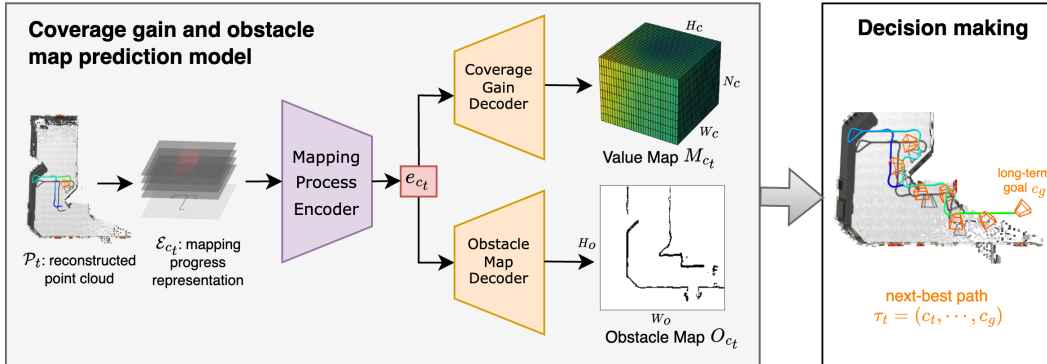


Figure 3: **Overview of the proposed next-best-path (NBP) framework.** The model (left, see Section 4.2) predicts a value map of coverage gain and an obstacle map, which are used for decision making (right, see Section 4.3) to obtain a next-best path.

Overview of our approach. Existing approaches for active mapping (Guédon et al., 2022; 2023) typically predict the next camera pose c_t in a greedy manner, which often suffers from getting stuck in limited areas. To address this limitation, we propose a novel approach that predicts a long-term goal camera pose and uses it to guide the next camera pose selection. Given all past observations and camera poses, our model predicts two key components centred on the agent’s current pose c_t : (1) a value map M_{c_t} , which estimates the surface coverage gain of candidate poses c in the surrounding of c_t , and (2) an obstacle map O_{c_t} , which accounts for both visible and predicted unseen obstacles in the environment. From the value map M_{c_t} , we derive the long-term goal pose c_g and combine it with the obstacle map O_{c_t} to compute an optimal path $\tau_t = (c_t, c_{t+1}, \dots, c_g)$ that navigates the agent from its current pose c_t to the goal pose c_g . This long-term goal-driven strategy helps the model avoid the pitfalls of short-sighted decisions and enhances coverage efficiency.

In the following, we first describe the model for M_{c_t} and O_{c_t} prediction in Section 4.2, then followed by the decision-making process to determine the next best path τ_t in Section 4.3. Finally, in Section 4.4, we introduce the training algorithm for our model.

4.2 COVERAGE GAIN AND OBSTACLE PREDICTION MODEL

Figure 3 depicts the deep model we use to predict the coverage gains and the obstacle map. We detail this model below.

Mapping Progress Encoder. Let’s denote \mathcal{P}_t the reconstructed point cloud at each time step t , obtained by adding the back-projected depth image I_t to the previously accumulated point cloud \mathcal{P}_{t-1} . Directly encoding the point cloud via 3D neural networks can be complex and inefficient. Therefore, we convert the 3D point cloud into multiple 2D images as inputs to a 2D-based encoder.

To be specific, we first centre and crop the point cloud based on the agent’s current position c_t . Centering the input on the agent makes the model invariant to the agent’s position and thus improves generalization. Then, the point cloud is divided into K horizontal layers along the gravity axis. For each layer, we average the occupancy value along the gravity axis to transform each 3D data into a 2D image. In this image, each pixel encodes the density of 3D points within a specific height range. The stack of K point cloud projected images provides a simplified yet informative representation of the 3D structure.

Similarly, we project the 3D trajectory of the agent’s past camera poses onto a 2D plane where each pixel denotes the frequency of visits to that location. This plane serves to mitigate the exploratory value of previously traversed regions. We define \mathcal{E}_{c_t} to include the K point cloud projected images and a single historical trajectory image.

Given the stacked 2D images of \mathcal{E}_{c_t} , we employ an Attention UNet (Oktay et al., 2018) encoder with 4 downsampling convolutional blocks to extract mapping progress features e_{c_t} .

Coverage Gain Decoder. This decoder predicts from e_{c_t} a 3D value map $M_{c_t} \in \mathbb{R}^{H_c \times W_c \times N_c}$ centered on the agent. It is composed of two upsampling convolutional blocks with an attention

mechanism. The first two dimensions of the predicted value map, H_c and W_c , correspond to the camera’s 2D position in the environment, while the third dimension N_c represents different camera orientations. Each value in M_{c_t} quantifies the estimated coverage gain achievable by moving the camera along the shortest trajectory from its current pose to the specific camera pose. The value map M_{c_t} guides the selection of both long-term goal poses c_g and intermediate poses along the trajectory, enabling a two-stage optimization for efficient exploration, which will be discussed in Section 4.3.

Obstacle Map Decoder. This decoder predicts the geometric layout $O_{c_t} \in \mathbb{R}^{H_o \times W_o}$ of the current moving plane, also from the encoder output e_{c_t} . O_{c_t} is a binary map representing potential obstacles around the current agent location, which is used for path planning. To be noted, O_{c_t} includes not only visible obstacles but also anticipated unseen obstacles based on the structure of the partially reconstructed point cloud, providing useful priors for navigation. This decoder is implemented using Attention U-Net with 4 upsampling convolutional blocks, and the output is passed through a sigmoid activation function to generate the binary obstacle map.

4.3 DECISION MAKING FOR NEXT-BEST-PATH PREDICTION

We derive both a long-term goal c_g and next-best-path $\tau_t = (c_t, c_{t+1}, \dots, c_g)$ from the predicted M_{c_t} and O_{c_t} , employing different decision making strategies for training and inference. During training, we balance exploitation and exploration, while we prioritize exploitation during inference.

Training phase. We rely on the Boltzmann exploration strategy (Cesa-Bianchi et al., 2017) to sample a camera pose as the goal c_g based on the value map M_{c_t} . The probability of selecting a camera pose c as the goal is given by:

$$P(c_g = c) = \frac{\exp(M_{c_t}[c]/\beta)}{\sum_{c' \in C} \exp(M_{c_t}[c']/\beta)}, \quad (1)$$

where C represents all possible camera poses within M_{c_t} , β is the temperature parameter that balances exploration and exploitation, and $M_{c_t}[c]$ denotes the value of the cell for candidate c .

Once the long-term goal c_g is sampled, we use the Dijkstra algorithm to find the shortest obstacle-free path from the current position c_t^{pos} to goal position c_g^{pos} with a ground truth obstacle map. To select camera orientation along the path, we also leverage M_{c_t} to sample one orientation from N_c potential orientations at each position. This strategy enhances data diversity and alleviates the risk of converging to local optima.

Inference phase. At inference, we take c_g as the pose with the maximum value in M_{c_t} , and the path planning is based on the predicted obstacle map O_{c_t} instead of ground truth. Each position in the trajectory is assigned the optimal orientation from the heatmap M_{c_t} for its location. In practice, the predicted obstacle map may not be entirely accurate. Encountering an unexpected obstacle requires halting the trajectory and initiating a new decision-making phase.

4.4 MODEL TRAINING

Algorithm 1 outlines the training procedure for our model. We first gather training data from all training scenes using the current model, and then update the model with the new data. This process is repeated iteratively until the model achieves convergence. We detail below the data collection, training objectives to update the model, and the training strategy.

Algorithm 1 Training procedure.

```

 $N$ : number of training iterations
 $N_e$ : number of iterations using easy data
 $S_n$ : the number of trajectories per scene
Initialize memory  $\mathcal{M} \leftarrow \emptyset$  and model parameters  $\theta$ 
for  $n \leftarrow 1$  to  $N$  do
  Initialize training set  $\mathcal{T} \leftarrow \emptyset$ 
  for each scene in training set do
    for  $s \leftarrow 1$  to  $S_n$  do
      Collect training data  $\{d_i\}_{i=1}^L$ 
      if  $n \leq N_e$  then  $\mathcal{T} \leftarrow \mathcal{T} \cup \{d_i : t \geq 10\}_{i=1}^L$ 
      else  $\mathcal{T} \leftarrow \mathcal{T} \cup \{d_i\}_{i=1}^L$  endif
    end for
  end for
   $\mathcal{M} \leftarrow \mathcal{M} \cup \mathcal{T}$ 
   $\mathcal{T} \leftarrow \mathcal{T} \cup \text{RandomSample}(\mathcal{M} \setminus \mathcal{T}, |\mathcal{T}|)$ 
  for  $e \leftarrow 1$  to  $E$  do
    Update  $\theta$  with loss in Eq. (3) over  $\mathcal{T}$ 
  end for
end for
return  $\theta$ 

```

Training data collection. After sampling the goal pose c_g and the trajectory τ_t , we generate ground truth labels to train the value map M_{c_t} and obstacle map O_{c_t} .

For M_{c_t} , we compute the coverage gain for the cell that corresponds to c_g as the ground truth label. Let \mathcal{P}_t and \mathcal{P}_g denote the reconstructed point clouds at pose c_t and c_g respectively, where \mathcal{P}_g is the result of accumulating depth information into \mathcal{P}_t as the agent moves along the trajectory τ_t . By comparing the reconstructed point clouds with the ground truth point cloud \mathcal{P}^{GT} , we can obtain the coverage gain $\Delta\text{Cov}_{c_t \rightarrow c_g}$:

$$\Delta\text{Cov}_{c_t \rightarrow c_g} = \frac{1}{N_{\text{GT}}} \sum_{i=1}^{N_{\text{GT}}} \left[\mathbf{1} \left(\min_{y \in \mathcal{P}_g} \|x_i^{\text{GT}} - y\| < \epsilon \right) - \mathbf{1} \left(\min_{y \in \mathcal{P}_t} \|x_i^{\text{GT}} - y\| < \epsilon \right) \right], \quad (2)$$

where N_{GT} is the number of points in \mathcal{P}^{GT} , $\|\cdot\|$ denotes the Euclidean distance, and ϵ is a predefined distance threshold. Consequently, we set $\Delta\text{Cov}_{c_t \rightarrow c_g}$ as the ground truth value for $M_{c_t}[c_g]$.

For O_{c_t} , we use the 3D mesh of the scene to derive the ground truth obstacle map $O_{c_t}^{\text{GT}}$. Specifically, we intersect the 3D mesh with a plane at the agent’s height, and project this intersection onto a 2D grid. This 2D grid is binarized to distinguish between obstacles and free space. Finally, we centre the 2D grid around the agent’s current position as $O_{c_t}^{\text{GT}}$.

To enhance the efficiency of data generation, we further perform a data augmentation by leveraging the property of Dijkstra’s algorithm, where every sub-path of a shortest path is also a shortest path. From a given path $\tau_t = (c_0 = c_t, \dots, c_m = c_g)$, we compute the coverage gain $\Delta\text{Cov}_{c_i \rightarrow c_j}$ for each segment of the path (c_i, c_j) where $0 \leq i < j \leq m$. More specifically, we update the ground truth values along the Dijkstra path $M_{c_i}^{\text{GT}}[c_j] = \Delta\text{Cov}_{c_i \rightarrow c_j}$. We also collect the input \mathcal{E}_{c_i} and the ground truth of surrounding obstacles $O_{c_i}^{\text{GT}}$ for each $c_i \in \tau_t$. This significantly increases the number of training samples derived from a single trajectory.

We store all augmented pairs $\{d_l\}_{l=1}^L$, $d_l = (\mathcal{E}_{c_i}, M_{c_i}^{\text{GT}}, O_{c_i}^{\text{GT}})$ in memory for training, where L is the length of the trajectory.

Multi-task training. We jointly train the coverage gain and obstacle map prediction using data stored in memory. We use the mean squared error (MSE) loss for training the coverage gain prediction, and the binary cross-entropy (BCE) loss for training the obstacle map prediction. To balance these two tasks effectively, we apply learnable uncertainty weights for each task, following Kendall et al. (2018). Our multi-task loss function for sample d_l is formulated as follows:

$$\mathcal{L}(\theta; d_l) = \frac{1}{2\sigma_1^2} \mathcal{L}_{\text{MSE}}(M_{c_i}^{\text{GT}}, \hat{M}_{c_i}) + \frac{1}{\sigma_2^2} \mathcal{L}_{\text{BCE}}(O_{c_i}^{\text{GT}}, \hat{O}_{c_i}) + \log \sigma_1 + \log \sigma_2, \quad (3)$$

where θ represents the model parameters, σ_1 and σ_2 are learnable uncertainty weights, \hat{M}_{c_i} and \hat{O}_{c_i} are the model’s predictions for the coverage gain and obstacle maps respectively.

Training strategy. We adopt a curriculum training strategy (Liu et al., 2017; Yuan et al., 2022; Yan et al., 2021; De Lange et al., 2021) to train our model, starting with easier-to-predict samples and gradually incorporating the entire dataset. In particular, we consider that the initial steps of a trajectory are more challenging since the agent has limited observations. Therefore, during the first N_e training iterations, we exclude samples from the first 10 steps in a trajectory. After N_e iterations, all samples in a trajectory are used in training.

During each training iteration, we use a balanced combination of previously stored data from the memory and newly collected data generated by the current model (Wulfmeier et al., 2018; Mnih, 2013; Rolnick et al., 2019; Aljundi et al., 2019), which helps prevent catastrophic forgetting. Each training phase is limited to E epochs to balance between enhancing performance and preventing overfitting on sub-optimal data.

5 EXPERIMENTS

5.1 EXPERIMENTAL SETUP

Dataset and simulation setup. We evaluate our model on the Matterport3D (MP3D) dataset (Chang et al., 2017) and our own AiMDoom dataset.

Table 2: **Evaluation results on AiMDoom Dataset.** For each difficulty level, all baseline models, including ours, are trained from scratch on the corresponding training set to ensure a fair comparison.

	Simple		Normal		Hard		Insane	
	Final Cov.	AUCs	Final Cov.	AUCs	Final Cov.	AUCs	Final Cov.	AUCs
Random	0.323 \pm 0.156	0.270 \pm 0.135	0.190 \pm 0.124	0.152 \pm 0.103	0.124 \pm 0.082	0.088 \pm 0.060	0.074 \pm 0.048	0.050 \pm 0.035
FBE	0.760 \pm 0.174	0.605 \pm 0.171	0.565 \pm 0.139	0.415 \pm 0.109	0.425 \pm 0.114	0.311 \pm 0.080	0.330 \pm 0.097	0.239 \pm 0.079
SCONE	0.577 \pm 0.173	0.483 \pm 0.138	0.412 \pm 0.114	0.313 \pm 0.087	0.290 \pm 0.093	0.210 \pm 0.072	0.196 \pm 0.079	0.140 \pm 0.060
MACARONS	0.599 \pm 0.200	0.479 \pm 0.172	0.418 \pm 0.120	0.314 \pm 0.088	0.302 \pm 0.097	0.218 \pm 0.070	0.192 \pm 0.078	0.139 \pm 0.058
NBP (Ours)	0.879\pm0.142	0.692\pm0.156	0.734\pm0.142	0.526\pm0.112	0.618\pm0.153	0.432\pm0.115	0.472\pm0.095	0.312\pm0.073

For MP3D, we use the same setting as prior work (Yan et al., 2023) for fair comparison. The input posed depth images have a resolution of 256×256 with a horizontal field of view (hFOV) of 90° . The mobile agent starts in the traversable space at a height of $1.25m$ and chooses its next camera pose by moving forward by $6.5cm$ or turning left/right by 10° . Depending on the size of each scene, the agent can take a maximum of 1000 or 2000 steps. We focus only on single-floor scenes following Yan et al. (2023) with 10 and 5 scenes in training and evaluation respectively.

For AiMDoom, we utilize a 70/30 train/test split for scenes in each difficulty level. The input RGB-D images are rendered at the resolution of 456×256 with hFOV of 90° . The agent navigates in a traversable space of height $1.65m$. The moving step includes 4 position movements (move forward, backward, left, or right by $1.5m$) and 8 rotation movements (turn left or right by increments of 45° , covering the full 360°). The maximum steps for Simple, Normal, Hard, and Insane levels are set to 100, 200, 400, and 500 respectively, to adapt to their different complexity.

Evaluation metrics. We follow prior works (Chen et al., 2024; Guédon et al., 2023) and adopt two key metrics to evaluate the performance of active 3D reconstruction: (1) *Final Coverage* measures the scene coverage at the end of the trajectory, and (2) *AUCs* evaluates the efficiency of the reconstruction process by calculating the area under the curve of coverage over time. The surface coverage is computed using ground truth meshes, consistent with prior work (Guédon et al., 2023). We evaluate five trajectories per scene using identical random initial camera poses for different methods. We report the mean and standard deviation for each metric across all testing trajectories.

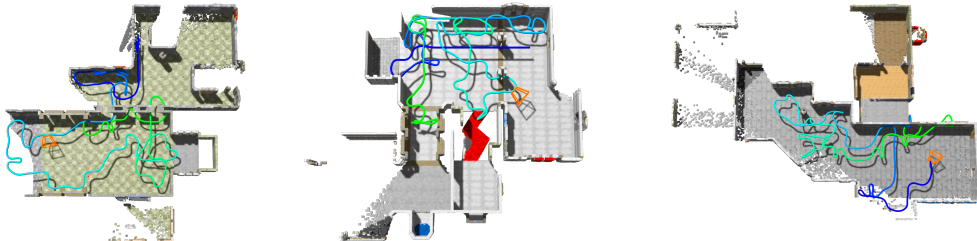
For a fair comparison with prior work in MP3D, we employ another set of metrics to evaluate coverage: (1) *Comp. (%)*, the proportion of ground truth vertices within $5cm$ of any observation, and (2) *Comp. (cm)*, the average minimum distance between ground truth vertices and observations.

Implementation details. Our model takes a stack of $K = 4$ projected 2D images and one previous trajectory projected image as inputs, each with a resolution of 256×256 covering a $40m \times 40m$ exploration area centred on the camera’s current position. The extracted feature e_{c_t} from the encoder is of size $16 \times 16 \times 1024$. The output value map M_{c_t} is of size $64 \times 64 \times 8$ and an obstacle map of $256 \times 256 \times 1$, both representing the same $40m \times 40m$ area. The model is trained for at most $N = 15$ iterations, with the first $N_e = 1$ iterations using easier samples and $S_n = 2$ trajectories per scene. For subsequent iterations, we use all samples and reduce the trajectory count to $S_n = 1$ per scene. Each trajectory has a length of 100 steps and starts at a random location. During the first data collection iteration, we randomly sample 1,000 validation examples from memory and exclude them from training. Gradient accumulation is used in training which results in an effective batch size of 448. The learning rate is set to 0.001 and is decayed by a factor of 0.1 if the validation loss plateaus. We apply early stopping to terminate training when validation loss no longer decreases. The training is performed on a single NVIDIA RTX A6000 GPU, with an average completion time of 25 hours.

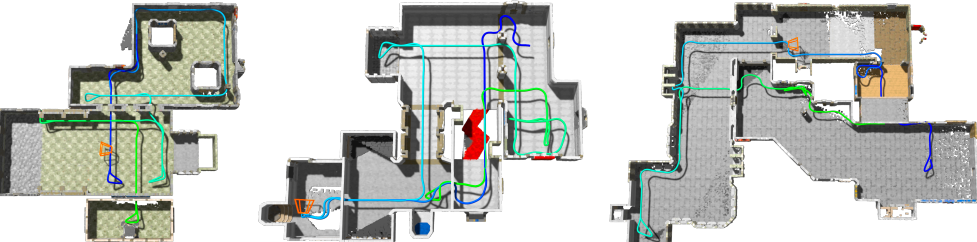
5.2 COMPARISON WITH STATE OF THE ART METHODS

MP3D. We compare our method with five baselines on the MP3D dataset, including: 1) *Random*, which randomly selects a camera pose among all candidates for the next step; 2) *Frontier-based Exploration (FBE)* (Yamauchi, 1997), which heuristically moves the agent to the nearest frontier; 3) *OccAnt* (Ramakrishnan et al., 2020), which predicts the occupancy status of unexplored areas and rewards the agent for accurate predictions; 4) *UPEN* (Georgakis et al., 2022), which utilizes an ensemble of occupancy prediction models to guide the agent towards paths with the highest un-

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485



(a) Results of MACARONS. It generates complicated trajectories and often gets trapped in local areas.



(b) Results of our NBP method. It efficiently travels in the scene and reconstructs the scene well.

Figure 4: **Comparison of our NBP method with the state-of-the-art MACARONS method.** Both methods start from the same initial pose, marked in deep blue. [We also include a demonstration video of active mapping using our method in the supplementary materials.](#)

certainty; 5) ANM (Yan et al., 2023), which guides exploration through a continually-learned neural scene representation. Table 3 presents the results. Our NBP method achieves the best performance, with a 6.23 absolute gain for the completion ratio compared to the previous best ANM (Yan et al., 2023) model.

Table 3: Comparison on the MP3D dataset.

Method	Comp. (%) \uparrow	Comp. (cm) \downarrow
Random	45.67	26.53
FBE	71.18	9.78
UPEN	69.06	10.60
OccAnt	71.72	9.40
ANM	73.15	9.11
NBP (ours)	79.38	6.78

AiMDoom. The proposed AiMDoom dataset is more challenging than MP3D dataset for active 3D mapping. We mainly benchmark our approach against state-of-the-art Next-Best-View (NBV) approaches, including: *SCONE* (Guédon et al., 2022) which employs volumetric integration to sum the potential visibility points for each candidate camera pose in the subsequent step and is trained using supervised learning; and *MACARONS* (Guédon et al., 2023) which quantifies the coverage gains of potential next camera poses to select

the best one and utilizes a self-supervised online learning paradigm. Both approaches select the next camera pose in a greedy manner. Unfortunately, we were unable to include UPEN (Georgakis et al., 2022) and ANM (Yan et al., 2023) in our comparison. These methods rely on the navigation policy DD-PPO (Wijmans et al., 2019) trained on their environments (Savva et al., 2019), which requires extensive GPU hours and thus is infeasible to retrain it on our dataset. However, we implemented FBE (Yamauchi, 1997) on our dataset, a recognized strong baseline in reconstruction and exploration tasks.

As shown in Table 2, our method significantly outperforms the baselines across all metrics on four levels of AiMDoom. While NBV approaches such as SCONE and MACARONS excel in outdoor or single-object scenarios, their performance deteriorates in complex indoor environments. As illustrated in Figure 4, MACARONS struggles to escape local areas due to its short-term focus. It only selects the next best pose in nearby regions, and once these areas - such as the interior of a single room - are fully reconstructed, it has difficulty moving out of the room to explore under-explored, distant regions. In contrast, our approach overcomes this limitation by incorporating long-term goal guidance to determine the next-best path. In addition, our method surpasses the strong baseline FBE. Although FBE enables better exploration compared to state-of-the-art NBV methods on our dataset, its simple heuristic of moving to the nearest frontier leads to sub-optimal scene reconstruction as it lacks strategic planning for efficient coverage.

486 Despite the superior performance of our model, the results in hard and insane environments are still
 487 unsatisfactory, highlighting the significant challenges posed in our dataset.
 488

489 5.3 ABLATION STUDY
 490

491 In this section, we perform ablation experi-
 492 ments to demonstrate the effectiveness of dif-
 493 ferent components in our model. All the ex-
 494 periments below are conducted on the Normal
 495 level of AiMDoom.

496 **Spatial range of long-term goal.** We compare
 497 the impact of different spatial ranges for the
 498 prediction of the value map M_{c_t} and obsta-
 499 cle map O_{c_t} , which in turn determines the
 500 maximum distance of the long-term goals c_g .
 501 Specifically, we experiment with map sizes of
 502 $20m \times 20m$ to $50m \times 50m$. The results are pre-
 503 sented in Figure 5. When the value map covers
 504 a smaller area, the goal c_g is close to the agent’s
 505 current position, leading to behaviour similar to
 506 existing NBV methods that struggle with explo-
 507 ration. On the other hand, if the map size is too large, predicting M_{c_t} and O_{c_t} becomes much more
 508 challenging. Our findings demonstrate that selecting an appropriate spatial range for the value map
 509 is crucial for balancing exploration efficiency and prediction accuracy.

510 **Oracle obstacle map.** In Table 4, we replace
 511 the predicted obstacle map with the ground truth
 512 obstacle map for path planning during inference,
 513 while maintaining to use the predicted value map
 514 for long-term goals. Using the oracle obstacle
 515 map improves the performance by 0.074 on final
 516 coverage and 0.054 on AUCs, but is far from per-
 517 fect. This suggests that the major bottleneck is
 518 the value map prediction.

519 Table 5: Comparison of single-task and
 520 multi-task learning for the value map
 521 and obstacle map prediction.

Strategy	Final Cov.	AUCs
Single-task	0.712 ± 0.136	0.501 ± 0.101
Multi-task	0.734 ± 0.142	0.526 ± 0.112

526 complement each other to enhance learning.
 527

528
 529 6 CONCLUSION
 530

531 In this paper, we tackle the challenging problem of active 3D mapping of unknown environments.
 532 We introduce a new dataset, AiMDoom, designed to benchmark active mapping in indoor scenes
 533 with four difficulty levels. Our evaluations of existing methods on the AiMDoom dataset reveal
 534 shortcomings of short-sighted next-best-view prediction in complex large indoor environments.
 535 Hence, we propose the next-best-path (NBP) method, which integrates a mapping progress encoder,
 536 a coverage gain decoder and an obstacle map decoder. The NBP model can efficiently reconstruct
 537 unseen environments guided by predicted long-term goals, achieving state-of-the-art performance
 538 on both the MP3D and AiMDoom datasets. However, we observe considerable room for improve-
 539 ment in more difficult levels of our dataset, and the major limitation lies in long-term goal prediction
 rather than obstacle map prediction.

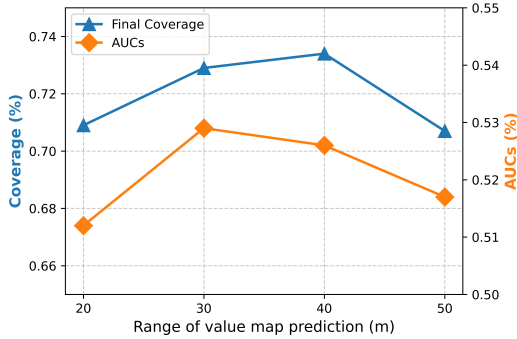


Figure 5: Comparisons of different spatial ranges for value map prediction.

Table 4: Ablation study on using the oracle map for obstacle avoidance at inference.

Obstacle Map	Final Cov.	AUCs
Predicted	0.734 ± 0.142	0.526 ± 0.112
Oracle	0.808 ± 0.115	0.580 ± 0.105

Multi-task training. We also explore the influence of multi-task learning in predicting the value map M_{c_t} and the obstacle map O_{c_t} . For comparison, we train two separate models that use the same input to predict M_{c_t} and O_{c_t} respectively. The results show that multi-task learning improved the precision of obstacle prediction to 0.805, exceeding the 0.754 achieved by single-task learning. Table 5 further demonstrates that multi-task learning achieves better performance, indicating that the two tasks

REFERENCES

- 540
541
542 Rahaf Aljundi, Eugene Belilovsky, Tinne Tuytelaars, Laurent Charlin, Massimo Caccia, Min Lin,
543 and Lucas Page-Caccia. Online continual learning with maximal interfered retrieval. *Advances*
544 *in neural information processing systems*, 32, 2019.
- 545
546 Andreas Bircher, Mina Kamel, Kostas Alexis, Helen Oleynikova, and Roland Siegwart. Receding
547 horizon” next-best-view” planner for 3d exploration. In *2016 IEEE international conference on*
548 *robotics and automation (ICRA)*, pp. 1462–1468. IEEE, 2016.
- 549
550 Chao Cao, Hongbiao Zhu, Howie Choset, and Ji Zhang. Tare: A hierarchical framework for effi-
551 ciently exploring complex 3d environments. In *Robotics: Science and Systems*, volume 5, pp. 2,
2021.
- 552
553 Nicolò Cesa-Bianchi, Claudio Gentile, Gábor Lugosi, and Gergely Neu. Boltzmann Exploration
554 Done Right. In *Advances in Neural Information Processing Systems*, 2017.
- 555
556 Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva,
557 Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3D: Learning from RGB-D Data in Indoor
Environments. In *International Conference on 3D Vision*, 2017.
- 558
559 Angel X. Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li,
560 Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, and Others. Shapenet: An Information-
561 Rich 3D Model Repository. In *arXiv Preprint*, 2015.
- 562
563 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.
- 564
565 Shizhe Chen, Thomas Chabal, Ivan Laptev, and Cordelia Schmid. Object Goal Navigation with
566 Recursive Implicit Maps. In *International Conference on Intelligent Robots and Systems*, 2023.
- 567
568 Xiao Chen, Quanyi Li, Tai Wang, Tianfan Xue, and Jiangmiao Pang. GenNBV: Generalizable Next-
569 Best-View Policy for Active 3D Reconstruction. In *Conference on Computer Vision and Pattern*
Recognition, 2024.
- 570
571 Titus Cieslewski, Elia Kaufmann, and Davide Scaramuzza. Rapid exploration with multi-rotors: A
572 frontier selection method for high speed flight. In *2017 IEEE/RSJ International Conference on*
Intelligent Robots and Systems (IROS), pp. 2135–2142. IEEE, 2017.
- 573
574 Blender Online Community. *Blender - A 3D Modelling and Rendering Package*, 2018. URL <http://www.blender.org>.
575
- 576
577 Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias
578 Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proceedings of the*
IEEE conference on computer vision and pattern recognition, pp. 5828–5839, 2017.
- 579
580 Matthias De Lange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Aleš Leonardis, Gregory
581 Slabaugh, and Tinne Tuytelaars. A continual learning survey: Defying forgetting in classification
582 tasks. *IEEE transactions on pattern analysis and machine intelligence*, 44(7):3366–3385, 2021.
- 583
584 Matt Deitke, Winson Han, Alvaro Herrasti, Aniruddha Kembhavi, Eric Kolve, Roozbeh Mot-
585 taghi, Jordi Salvador, Dustin Schwenk, Eli VanderBilt, Matthew Wallingford, Luca Weihs, Mark
586 Yatskar, and Ali Farhadi. Robothor: An open simulation-to-real embodied ai platform. In
IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), June 2020.
- 587
588 Ziyue Feng, Huangying Zhan, Zheng Chen, Qingan Yan, Xiangyu Xu, Changjiang Cai, Bing Li,
589 Qilun Zhu, and Yi Xu. Naruto: Neural active reconstruction from uncertain target observations.
590 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.
591 21572–21583, 2024.
- 592
593 Georgios Georgakis, Bernadette Bucher, Anton Arapin, Karl Schmeckpeper, Nikolai Matni, and
Kostas Daniilidis. Uncertainty-Driven Planner for Exploration and Navigation. In *International*
Conference on Robotics and Automation, 2022.

- 594 Antoine Guédon, Pascal Monasse, and Vincent Lepetit. Scone: Surface Coverage Optimization in
595 Unknown Environments by Volumetric Integration. In *Advances in Neural Information Process-*
596 *ing Systems*, 2022.
- 597
- 598 Antoine Guédon, Tom Monnier, Pascal Monasse, and Vincent Lepetit. Macarons: Mapping and
599 Coverage Anticipation with RGB Online Self-Supervision. In *Conference on Computer Vision*
600 *and Pattern Recognition*, 2023.
- 601
- 602 Guillaume Hardouin, Julien Moras, Fabio Morbidi, Julien Marzat, and El Mustapha Mouaddib.
603 Next-best-view planning for surface reconstruction of large-scale 3d environments with multiple
604 uavs. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp.
605 1567–1574. IEEE, 2020.
- 606 Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-Task Learning Using Uncertainty to Weigh
607 Losses for Scene Geometry and Semantics. In *Conference on Computer Vision and Pattern Recog-*
608 *nition*, 2018.
- 609
- 610 Keifer Lee, Shubham Gupta, Sungyoung Kim, Bhargav Makwana, Chao Chen, and Chen Feng. So-
611 nerf: Active view planning for nerf using surrogate objectives. *arXiv preprint arXiv:2312.03266*,
612 2023.
- 613
- 614 Soomin Lee, Le Chen, Jiahao Wang, Alexander Liniger, Suryansh Kumar, and Fisher Yu. Uncer-
615 tainty Guided Policy for Active Robotic 3D Reconstruction Using Neural Radiance Fields. *IEEE*
616 *Robotics and Automation Letters*, 7(4), 2022.
- 617
- 618 Weiwei Liu, Ivor W. Tsang, and Klaus-Robert Müller. An Easy-To-Hard Learning Paradigm for
619 Multiple Classes and Multiple Labels. *Journal of Machine Learning Research*, 18(94), 2017.
- 620
- 621 Chongshan Lu, Fukun Yin, Xin Chen, Wen Liu, Tao Chen, Gang Yu, and Jiayuan Fan. A large-
622 scale outdoor multi-modal dataset and benchmark for novel view synthesis and implicit scene
623 reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,
624 pp. 7557–7567, 2023.
- 625
- 626 Hidenobu Matsuki, Riku Murai, Paul HJ Kelly, and Andrew J Davison. Gaussian splatting slam.
627 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.
628 18039–18048, 2024.
- 629
- 630 Miguel Mendoza, J Irving Vasquez-Gomez, Hind Taud, L Enrique Sucar, and Carolina Reta. Su-
631 pervised learning of the next-best-view for 3d object reconstruction. *Pattern Recognition Letters*,
632 133:224–231, 2020.
- 633
- 634 Volodymyr Mnih. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*,
635 2013.
- 636
- 637 Ozan Oktay, Jo Schlemper, Loic Le Folgoc, Matthew Lee, Mattias Heinrich, Kazunari Misawa,
638 Kensaku Mori, Steven Mcdonagh, Nils Y. Hammerla, and Bernhard Kainz. Attention U-Net:
639 Learning Where to Look for the Pancreas. In *Medical Imaging with Deep Learning*, 2018.
- 640
- 641 Xuran Pan, Zihang Lai, Shiji Song, and Gao Huang. ActiveNeRF: Learning Where to See with
642 Uncertainty Estimation. In *European Conference on Computer Vision*, 2022.
- 643
- 644 Daryl Peralta, Joel Casimiro, Aldrin Michael Nilles, Justine Aletta Aguilar, Rowel Atienza, and
645 Rhandley Cajote. Next-Best View Policy for 3D Reconstruction. In *Workshop at European*
646 *Conference on Computer Vision*, 2020.
- 647
- 648 Julio A Placed, Jared Strader, Henry Carrillo, Nikolay Atanasov, Vadim Indelman, Luca Carlone,
649 and José A Castellanos. A survey on active simultaneous localization and mapping: State of the
650 art and new frontiers. *IEEE Transactions on Robotics*, 39(3):1686–1705, 2023.
- 651
- 652 Santhosh K. Ramakrishnan, Ziad Al-Halah, and Kristen Grauman. Occupancy Anticipation for
653 Efficient Exploration and Navigation. In *European Conference on Computer Vision*, 2020.

- 648 Santhosh Kumar Ramakrishnan, Aaron Gokaslan, Erik Wijmans, Oleksandr Maksymets, Alexander
649 Clegg, John M. Turner, Eric Undersander, Wojciech Galuba, Andrew Westbury, Angel X.
650 Chang, Manolis Savva, Yili Zhao, and Dhruv Batra. Habitat-Matterport 3D Dataset (HM3D):
651 1000 Large-Scale 3D Environments for Embodied AI. In *Thirty-fifth Conference on Neural In-*
652 *formation Processing Systems Datasets and Benchmarks Track*, 2021.
- 653 Ram Ramrakhya, Eric Undersander, Dhruv Batra, and Abhishek Das. Habitat-web: Learning em-
654 bodied object-search strategies from human demonstrations at scale. In *International Conference*
655 *on Computer Vision*, pp. 5173–5183, 2022.
- 657 Yunlong Ran, Jing Zeng, Shibo He, Jiming Chen, Lincheng Li, Yingfeng Chen, Gimhee Lee, and
658 Qi Ye. Neurar: Neural Uncertainty for Autonomous 3D Reconstruction with Implicit Neural
659 Representations. *IEEE Robotics and Automation Letters*, 8(2), 2023.
- 660 Nikhila Ravi, Jeremy Reizenstein, David Novotny, Taylor Gordon, Wan-Yen Lo, Justin Johnson, and
661 Georgia Gkioxari. Accelerating 3D Deep Learning with PyTorch3D. In *arXiv Preprint*, 2020.
- 663 David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy Lillicrap, and Gregory Wayne. Experience
664 replay for continual learning. *Advances in neural information processing systems*, 32, 2019.
- 665 Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain,
666 Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, Devi Parikh, and Dhruv Batra. Habitat: A
667 Platform for Embodied AI Research. In *International Conference on Computer Vision*, 2019.
- 669 Julian Straub, Thomas Whelan, Lingni Ma, Yufan Chen, Erik Wijmans, Simon Green, Jakob J.
670 Engel, Raul Mur-Artal, Carl Ren, Shobhit Verma, Anton Clarkson, Mingfei Yan, Brian Budge,
671 Yajie Yan, Xiaqing Pan, June Yon, Yuyang Zou, Kimberly Leon, Nigel Carter, Jesus Briales, Tyler
672 Gillingham, Elias Mueggler, Luis Pesqueira, Manolis Savva, Dhruv Batra, Hauke M. Strasdat,
673 Renzo De Nardi, Michael Goesele, Steven Lovegrove, and Richard Newcombe. The Replica
674 Dataset: A Digital Replica of Indoor Spaces. In *arXiv Preprint*, 2019.
- 675 Yuezhan Tao, Yuwei Wu, Beiming Li, Fernando Cladera, Alex Zhou, Dinesh Thakur, and Vijay
676 Kumar. Seer: Safe efficient exploration for aerial robots using learning to predict information
677 gain. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 1235–
678 1241. IEEE, 2023.
- 679 Erik Wijmans, Abhishek Kadian, Ari Morcos, Stefan Lee, Irfan Essa, Devi Parikh, Manolis Savva,
680 and Dhruv Batra. Dd-ppo: Learning near-perfect pointgoal navigators from 2.5 billion frames.
681 *arXiv preprint arXiv:1911.00357*, 2019.
- 682 Tong Wu, Jiarui Zhang, Xiao Fu, Yuxin Wang, Liang Pan Jiawei Ren, Wayne Wu, Lei Yang, Jiaqi
683 Wang, Chen Qian, Dahua Lin, and Ziwei Liu. Omniobject3d: Large-vocabulary 3d object dataset
684 for realistic perception, reconstruction and generation. In *IEEE/CVF Conference on Computer*
685 *Vision and Pattern Recognition (CVPR)*, 2023.
- 687 Markus Wulfmeier, Alex Bewley, and Ingmar Posner. Incremental adversarial domain adaptation
688 for continually changing environments. In *2018 IEEE International conference on robotics and*
689 *automation (ICRA)*, pp. 4489–4495. IEEE, 2018.
- 690 Fei Xia, R. Zamir, Amir, Zhi-Yang He, Alexander Sax, Jitendra Malik, and Silvio Savarese. Gib-
691 son Env: Real-World Perception for Embodied Agents. In *Conference on Computer Vision and*
692 *Pattern Recognition*, 2018.
- 694 Zhefan Xu, Christopher Suzuki, Xiaoyang Zhan, and Kenji Shimada. Heuristic-based incremen-
695 tal probabilistic roadmap for efficient uav exploration in dynamic environments. In *2024 IEEE*
696 *International Conference on Robotics and Automation (ICRA)*, pp. 11832–11838. IEEE, 2024.
- 697 Brian Yamauchi. A Frontier-Based Approach for Autonomous Exploration. In *IEEE International*
698 *Symposium on Computational Intelligence in Robotics and Automation*, 1997.
- 700 Zike Yan, Yuxin Tian, Xuesong Shi, Ping Guo, Peng Wang, and Hongbin Zha. Continual Neural
701 Mapping: Learning an Implicit Scene Representation from Sequential Observations. In *Internat-*
ional Conference on Computer Vision, 2021.

702 Zike Yan, Haoxiang Yang, and Hongbin Zha. Active Neural Mapping. In *International Conference*
703 *on Computer Vision*, 2023.
704

705 Zhenghang Yuan, Lichao Mou, Qi Wang, and Xiao Xiang Zhu. From easy to hard: Learning
706 language-guided curriculum for visual question answering on remote sensing data. *IEEE trans-*
707 *actions on geoscience and remote sensing*, 60:1–11, 2022.

708 Rui Zeng, Wang Zhao, and Yong-Jin Liu. Pc-Nbv: A Point Cloud Based Deep Network for Efficient
709 Next Best View Planning. In *International Conference on Intelligent Robots and Systems*, 2020.
710

711 Huangying Zhan, Jiyang Zheng, Yi Xu, Ian Reid, and Hamid Rezatofighi. Activermap: Radiance
712 Field for Active Mapping and Planning. In *arXiv Preprint*, 2022.

713 Boyu Zhou, Yichen Zhang, Xinyi Chen, and Shaojie Shen. Fuel: Fast uav exploration using incre-
714 mental frontier structure and hierarchical planning. *IEEE Robotics and Automation Letters*, 6(2):
715 779–786, 2021.

716 Qian-Yi Zhou, Jaesik Park, and Vladlen Koltun. Open3D: A Modern Library for 3D Data Process-
717 ing. In *arXiv Preprint*, 2018.
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755