

# POMP++: Pomcp-based Active Visual Search in unknown indoor environments

Francesco Giuliani<sup>1,2</sup>, Alberto Castellini<sup>3</sup>, Riccardo Berra<sup>3</sup>, Alessio Del Bue<sup>2</sup>,  
Alessandro Farinelli<sup>3</sup>, Marco Cristani<sup>3</sup>, Francesco Setti<sup>3</sup> and Yiming Wang<sup>2,4</sup>

**Abstract**—In this paper, we focus on the problem of learning online an optimal policy for Active Visual Search (AVS) of objects in unknown indoor environments. We propose POMP++, a planning strategy that introduces a novel formulation on top of the classic Partially Observable Monte Carlo Planning (POMCP) framework, to allow training-free online policy learning in unknown environments. We present a new belief reinvigoration strategy that enables the use of POMCP with a dynamically growing state space to address the online generation of the floor map. We evaluate our method on two public benchmark datasets, AVD that is acquired by real robotic platforms and Habitat ObjectNav that is rendered from real 3D scene scans, achieving the best success rate with an improvement of  $>10\%$  over the state-of-the-art methods.

## I. INTRODUCTION

Finding a specific object in an indoor environment is a human activity that is trivial to define but complex to encode into an autonomous agent like a robot. This task leverages several basic human skills including self-localisation, visual search, object detection, along all with the necessary motor skills for reaching the object of interest, especially in highly cluttered and dynamic environments.

In such context, this paper focuses on a specific task, Active Visual Search (AVS), where a robot is asked to navigate in indoor environments to search for a given object, with the performance being evaluated mainly on the success rate and the path efficiency (see Fig. 1). To address AVS, a joint solution for both visual perception and motion planning is needed. Recent works mostly tackle this problem by intertwining deep Reinforcement Learning (RL) techniques, *e.g.* deep recurrent Q-network (DRQN), with visual semantics, by either feeding deep visual embeddings to policy learning networks [1], [2] or obtaining 3D scene semantics to guide the planning [3]. These RL-based methods require a large amount of training data, *i.e.* sequences of observations of various lengths, covering successful and unsuccessful episodes from real and simulated scenarios.

<sup>1</sup>Department of Electrical, Electronics and Telecommunication Engineering, University of Genoa, Italy

<sup>2</sup>Visual Geometry and Modelling (VGM) and Pattern Analysis and Computer Vision (PAVIS), Fondazione Istituto Italiano di Tecnologia, Italy.

<sup>3</sup>Department of Computer Science, University of Verona, Italy.

<sup>4</sup>Deep Visual Learning (DVL) Unit, Fondazione Bruno Kessler, Italy.

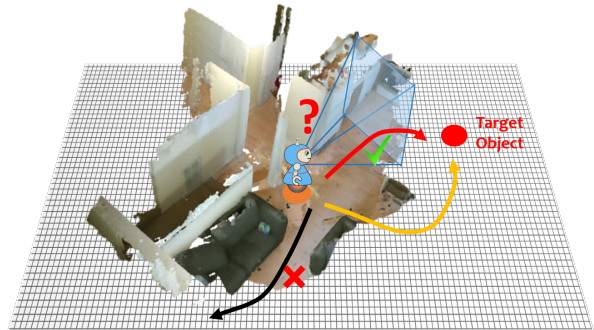


Fig. 1: A robot is navigating in an *unknown* environment with the task of visually searching for a target object (highlighted in red), *i.e.* to have the object captured within the field of view of the camera. The robot is driven by a motion policy based on partially observed scene to visually detect the target object located in the unseen part of scene (highlighted in gray area) with the shortest travelled path (highlighted in red), to avoid longer trajectories (in yellow) or missing entirely the target (in black).

Such pre-collected data may not be representative of the deployed (unknown) environments, thus possibly reducing the performance of a learning-based strategy.

In this paper, we propose a method to specifically address AVS in *unknown* environments without a priori knowledge of the area. While online policy learning has been used in previous approaches [4], our method is the first to address this problem in the *unknown* environment setting. Our approach, named POMP++, is developed within the Partially Observable Monte Carlo Planning (POMCP) framework with two main adaptations introduced. We first present a formal definition of the AVS task in unknown environments following the Partially Observable Markov Decision Process (POMDP) formulation, in particular addressing the problem of a growing search space. Secondly, since the state space is dynamically changing in our problem, the original belief reinvigoration strategy of POMCP does not work as it requires a fixed set of states, we therefore propose a novel belief reinvigoration strategy that allows for an efficient path planning as the robot builds up the 2D map of the environment. Once the POMCP module locates the target, the robot approaches it following the shortest path estimated

by a local planner.

Our main technical contributions can be summarised as:

- 1) We formalise the AVS problem in unknown environments as a POMDP;
- 2) We integrate frontier-based exploration with POMCP planning for the AVS task, allowing the state space of the model to dynamically expand as the robot moves;
- 3) We propose a novel two-step belief reinvigoration strategy to address the dynamic state space within the POMCP framework.

POMP++ is evaluated on public benchmark datasets including Active Vision Dataset (AVD) [5] with sequences acquired by real robotic platforms and the Habitat ObjectNav Challenge [6] with observations rendered from real 3D scans. Our approach outperforms baselines and state-of-the-art (sota) approaches with a significant improvement of >10% in terms of the success rate, without performing any costly offline training.

## II. RELATED WORK

We cover closely related works addressing the AVS task and the development and application of POMCP to exploration or vision domains.

### A. Active Visual Search

Earlier works addressing the AVS task often exploit target-specific inferences, such as object co-occurrences [7]–[11], while recent works exploit Deep Reinforcement Learning techniques [1]–[3], where visual neural embeddings are often used for the policy training. EAT [1] performs feature extraction from the current RGB image, and the image crop of the candidate target generated by a Region Proposal Network (RPN). The features are then fed into the action policy network. EAT is trained with twelve scans from the AVD [5] dataset. Similarly, GAPLE [2] uses deep visual features enriched by 3D information (from depth) for training the policy with a large dataset from 100 synthetic scenes. Recent benchmarking efforts offer larger amount of data for training and foster more data-driven AVS solutions. For example, RoboTHOR [12] provides observations rendered from photo-realistic rooms [12], while Habitat ObjectNav Challenge [6] provides renderings from real 3D scans. SemExp [3], the current best-performing method on the ObjectNav Challenge, takes RGB-D frames as input and learns a 3D scene semantics representation, which can then be used for selecting a long-term goal for reaching the target object. SemExp proved that the 3D scene semantics provide more efficient guidance compared to the runner-up solution driven by scene exploration.

In general, learning-based strategies leverage large amount of data to learn a model of the environment together with the motion policy, while online motion planning strategies have the advantage of being general and easy-to-deploy. A recent POMCP-based solution [4] is able to achieve comparable search performance against sota data-driven methods without any offline training. However, this approach is only applicable to known environments, i.e. the 2D floor map of

the scene is required. Instead, our method POMP++ follows the idea of online policy learning, and further proposes novel adaptations in terms of map representation and belief reinvigoration, which address unknown environments.

### B. POMCP

*Partially Observable Markov Decision Process (POMDP)* is an established framework for sequential planning under uncertainty [13]. It allows to model dynamical processes in uncertain environments, where the uncertainty is related to actions and observations. While computing exact solutions for large POMDPs is computationally intractable [14], impressive progress has been made by approximated [15] and online [16] solvers. One of the most used and efficient online solvers for POMDPs is Partially Observable Monte Carlo Planning (POMCP) [17] which combines a Monte Carlo update of the belief state with a *Monte Carlo Tree Search (MCTS)* based policy [18]–[20], thus enabling the scalability to large state spaces. Most recent extensions of POMCP include applications to multi-agent problems [21], reward maximisation with cost constraints [22] and the introduction of prior knowledge about state space to refine the belief space and increase policy performance [23].

In this paper, we advance the state of the art by applying POMCP to AVS in unknown environments by proposing mechanisms to extend the state space, the transition and the observation models step by step. To address unknown environments, we exploit the elements of the frontier-based exploration theory [24]. Frontier-based exploration has been merged with POMDPs [25], but only for the exploration task. We instead focus on the AVS task and propose a novel POMDP formalisation and belief reinvigoration strategy to benefit effective path planning.

## III. PROPOSED METHOD

We consider the scenario where a robot moves in an *unknown* environment, with the task of searching for a specific object. The robot explores the environment to: *i)* observe the target object, *ii)* localise the object in the floor map, and *iii)* finally approach the object, *i.e.* moving close to the object location. Each search episode is terminated once the robot believes the target object is sufficiently near to itself and visually detectable<sup>1</sup>. We formulate the AVS problem as a POMDP and employ POMCP to compute the planning policy online. In the following section, we use the subscript  $t$  to represent the elements at the current time and the subscript  $t + 1$  to represent the updated elements after a new sensor capturing and observation.

Our overall framework is shown in Fig. 2. At each time step  $t$  the robot locates at a pose  $p_t = \{x_t, y_t, \theta_t\}$ , where  $x$  and  $y$  are the coordinates on the 2D floor plane, and  $\theta$  is the orientation. From that pose it receives a new capturing from a RGB-D sensor for updating the environment map  $\mathcal{M}_t$ , which encodes all the information about the environment that the

<sup>1</sup>The exact stop condition depends on the specific application. In the result section we consider standard evaluation protocols used in public benchmarks for AVS evaluation.

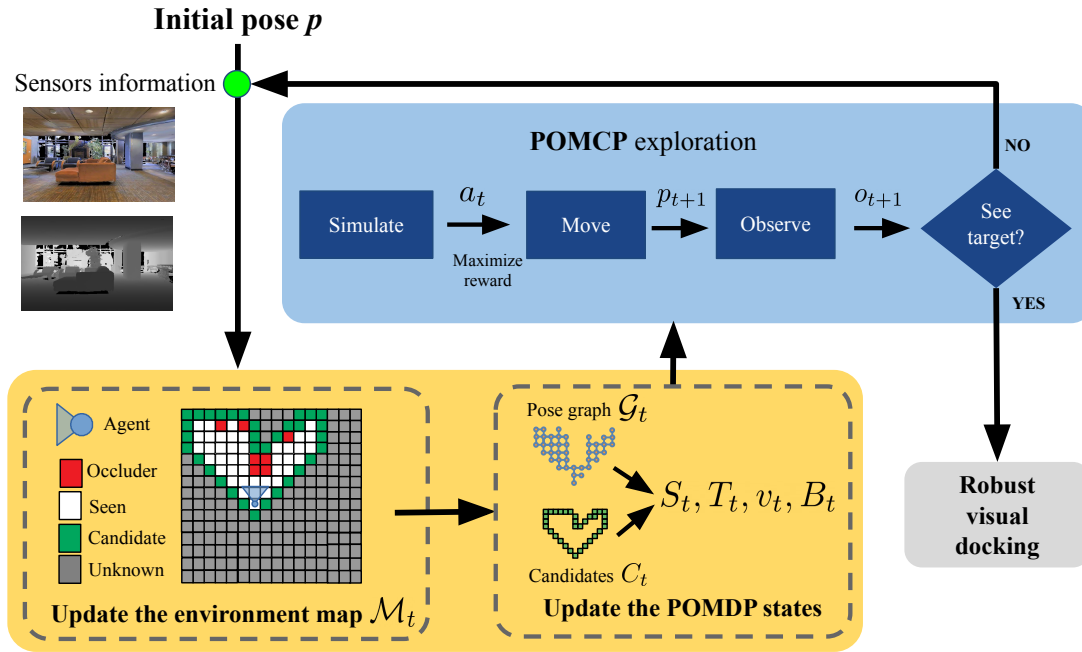


Fig. 2: Overall architecture of POMP++. The yellow block represents the partial knowledge model of the environment that gets updated during the exploration, while the blue block represents the exploration strategy to locate the target. The mathematical notation is defined as: pose  $p_t$ , action  $a_t$ , observation  $o_t$ , POMDP state space  $S_t$ , POMDP transition model  $T_t$ , visibility function  $v_t$ , belief  $B_t$ .

robot has gathered until  $t$ . The environment map is used to: 1) create a pose graph  $\mathcal{G}_t$ , that includes all the reachable poses of the robot, 2) extract the candidate locations of the target object  $C_t$ , namely the frontier positions bordering the explored and unknown cells, 3) update the internal states of the POMDP, and reinvigorate the belief to make it cover newly discovered candidate locations. Using the planning policy calculated by the POMCP exploration module, the robot reaches a new pose  $p_{t+1}$  by taking an action  $a_t$ . The process repeats until the target object is detected. Once it locates the object, the robot will approach it following the shortest path between the pose of the robot and the estimated position of the object, on the graph  $\mathcal{G}_t$ <sup>2</sup>.

In section III-A we present our new POMDP formulation that addresses the problem of exploring a search space that is initially unknown, and grows as the robot discovers new parts of the environment. In section III-B we detail how we use partial information about the environment to define the search space in POMCP to extract the exploration policy. We also show how the map of the scene can be updated every time when the robot discovers new parts of the environment. Finally, section III-C focuses on a new belief reinvigoration strategy for updating the belief when new states are discovered. This is a key contribution of our approach since the state space in the standard POMCP is fixed and only the probability over states is updated step by step.

<sup>2</sup>We apply the same approaching strategy, i.e. robust visual docking, as in [4] that allows path re-planning with continuous observations to be more robust to poor object detection.

#### A. POMDP formulation for Active Visual Search

The POMDP used in our context is a tuple  $(S, A, O, T, Z, R, \gamma)$  where

- $S$  is a finite set of partially observable *states* representing combinations of robot poses and target object locations, whose cardinality is unknown as the environment is unknown;
- $A$  is the finite set of robot *actions*;
- $O: S \times A \rightarrow \Pi(Z)$  is the *observation model*;
- $T: S \times A \rightarrow \Pi(S)$  is the *state-transition model* (with  $\Pi(S)$  probability distribution over states in  $S$ );
- $Z$  is a finite set of *observations*, namely, 1 if the searched object has been observed, 0 otherwise;
- $R: S \times A \rightarrow R$  is the *reward function*, which is positive when the object is detected, and a negative when the robot moves in position where the object is not detected;
- $\gamma \in [0, 1)$  is a *discount factor*.

The goal is to maximise the expected total discounted reward  $\mathbb{E}[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t)]$ , by choosing the best action  $a_t$  in each state  $s_t$ . The term  $\gamma$  reduces the impact of future rewards and ensures that the (infinite) sum converges to a finite value. The partial observability of the state is dealt with by considering at each time step a probability distribution over all the states, called *Belief*  $B \in \mathcal{B}$ . POMDP solvers compute, in an exact or approximate way, a *policy*, i.e. a function  $\pi: \mathcal{B} \rightarrow A$  that maps beliefs to actions. We use POMCP [17] to compute the policy online. It is a recent algorithm employing Monte-Carlo Tree Search (MCTS) to solve POMDP. Every time an action has to be selected by the robot, POMCP performs controlled simulations using the

Algorithm 1: Map update after each step. We omit the subscript  $t$  or  $t + 1$  for simplicity.

---

```

procedure MAPUPDATE(map  $\mathcal{M}$ , pose  $p$ )
  retrieve image from RGB-D sensor from pose  $p$ 
   $rec \leftarrow$  scene reconstruction from RGBD Image
  for all cells  $c \in \mathcal{M}$  do
    if  $c$  in camera frustum then
      if  $c$  is unoccupied in  $rec$  then
         $c \leftarrow seen$ 
      else
         $c \leftarrow obstacle$ 
    for all cells  $c \in \mathcal{M}$  do
      if  $c$  not camera frustum then
        if  $c$  is unoccupied in  $rec$  then
          if  $c$  is adjacent to  $seen$  cell then
             $c \leftarrow candidate$ 
          else
             $c \leftarrow unknown$ 
  return  $\mathcal{M}$ 

```

---

known transition and observation models, and from those simulations it learns the approximate value of each action. Then it selects the action with the highest value. The algorithm uses a particle filter to represent the (approximated) belief and MCTS to drive the simulation process in an optimal way.

### B. Adaptation to unknown environments

In unknown environments the state space  $S$  can change dynamically because all the possible robot poses and object locations are not known *a priori*, but instead discovered over time as the robot explores the environment. Here we describe how the main elements in the state space  $S$  and, consequently, the transition model  $T$ , the observation model  $O$ , and the belief  $B$  are updated at each time step.

The pose graph  $\mathcal{G}$  is composed of nodes representing possible robot poses and edges connecting poses reachable by the robot with a single action. We use  $V^{\mathcal{G}}$  and  $E^{\mathcal{G}}$  for the set of nodes and edges in  $\mathcal{G}$  respectively. The pose graph  $\mathcal{G}$  initially contains only poses in the field of view of the robot, while new nodes and edges are added as the exploration proceeds. The calculation of the reachable poses from the current camera view is solved by scene reconstruction methods with the depth maps [26].

We represent the *environment map*  $\mathcal{M}$  as a grid that covers the whole environment containing observed and unseen parts of the scene. We assign each cell of the map  $\mathcal{M}$  with one of the four possible values: *i) occluder*, the cell in the real environment contains any occluders, e.g. walls, preventing the robot to see what is behind the cell; *ii) seen*, a location which could contain the target, but has been observed and found empty by the robot; *iii) candidate*, the cell is unknown but adjacent to a *seen* cell, *i.e.* a frontier cell; *iv) unknown*, the cell is unseen and is not a candidate cell. In the beginning, all the cells outside the field of view are initialised as

Algorithm 2: POMCP exploration. We omit subscripts  $t$  and  $t + 1$  for simplicity.

---

```

procedure POMCP-EXP(initial pose  $p$ )
   $\mathcal{M} \leftarrow$  init all cells to unknown
   $V^{\mathcal{G}} \leftarrow \{p\}$ 
   $E^{\mathcal{G}} \leftarrow \emptyset$ 
  repeat
     $\mathcal{M} \leftarrow$  MAPUPDATE( $p$ )
     $V^{\mathcal{G}} \leftarrow V^{\mathcal{G}} \cup \{\text{poses related to } seen \text{ cells in } \mathcal{M}\}$ 
     $E^{\mathcal{G}} \leftarrow$  edges between  $p_i, p_j \in V^{\mathcal{G}}$ 
     $C \leftarrow \{\text{candidate cells in } \mathcal{M}\}$ 
     $S \leftarrow$  states considering updated  $V^{\mathcal{G}}$  and  $C$ 
     $T \leftarrow$  transition model from  $E^{\mathcal{G}}$ 
     $O \leftarrow$  visibility function on  $\mathcal{M}$ 
     $B \leftarrow$  BELIEF REINTEGRATION( $S$ )
     $a \leftarrow$  POMCP( $S, T, O, p, B$ )
     $p \leftarrow$  MOVE( $a$ )
     $o \leftarrow$  DETECTOR  $\triangleright$  true if object is detected
  until  $o = \text{True}$ 

```

---

*unknown*. Cell values are then updated step by step as the robot explores the environment (see Algorithm 1).

With the updated environment map  $\mathcal{M}_{t+1}$ , we can update the pose graph  $\mathcal{G}_{t+1}$ . For each node  $p_i \in \mathcal{G}_t$ , if taking action  $a_t$  leads the robot to a new pose  $p_j$  in the environment which falls into a *seen* cell of the updated map  $\mathcal{M}_{t+1}$ , then a node for pose  $p_j$  is added to  $\mathcal{G}_{t+1}$  with an edge  $(p_i, p_j)$ .

We update the visibility function at each time step  $t$  as  $v_t : V_t^{\mathcal{G}} \times C_t \rightarrow \{0, 1\}$ , where  $V_t^{\mathcal{G}}$  is the set of nodes in  $\mathcal{G}$  at time  $t$  and  $C_t$  is the set of *candidate* cells, *i.e.* candidate object locations at time  $t$ . The function takes a robot pose  $p \in V_t^{\mathcal{G}}$ , and a candidate object location  $c \in C_t$ , and returns 1 if  $c$  is visible from  $p$  (*i.e.* no obstacle lying between  $p$  and  $c$  in the current map  $\mathcal{M}_t$ ), 0 otherwise. Notice that, the aim of function  $v_t$  is to address pose-cell visibility for all robot poses in  $V_t^{\mathcal{G}}$  and all *candidate* cells.

POMP++ uses the elements described above to compute the optimal policy. Poses in  $V_t^{\mathcal{G}}$  along with the *candidate* cells in  $\mathcal{M}_t$  are used to build the state set  $S_t$ ; *edges* in  $E_t^{\mathcal{G}}$  are used in the transition model  $T_t$ , and the visibility function  $v_t$  is used by the observation model  $O_t$  in the simulation phase. Each state in  $S_t$  represents the target object in a specific candidate cell. Algorithm 2 describes how these states are updated during the exploration. At each step  $t$ , POMCP executes  $\beta$  number of simulations to compute the optimal action to perform in the real environment. Each simulation uses a specific state, which assumes the object is in a specific candidate cell. During the simulation the environment is not updated as there is no observation in the real environment. After each step the observation model  $O_t$  predicts if the robot sees the target object from the current pose and a reward is assigned to the action  $a_t$ . The *Reward* for  $a_t$  is a high positive value if the robots moves in a position where it can detect the object, otherwise a small negative value is used to penalise longer paths. Once the optimal action is

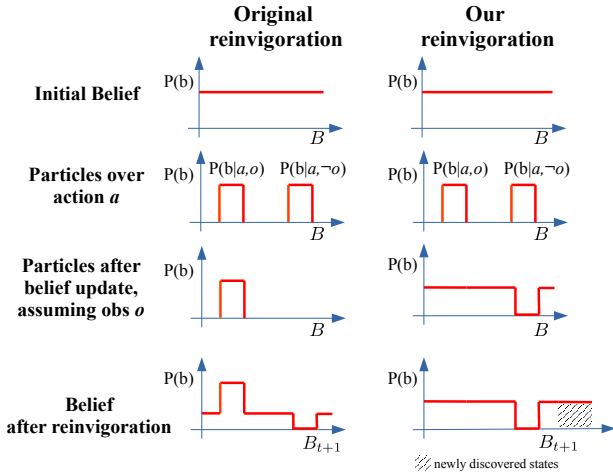


Fig. 3: Difference between the original belief reinvigoriation and our proposed reinvigoriation strategy. The original reinvigoriation penalises particles  $b \in B$  that are not sampled by simulations taking  $a$  as the first action, and receiving the observation  $o$ . Our new strategy, penalise only particles that were sampled by simulations where the simulated observation after the simulated action  $a$  was different from the real world observation  $o$ . Our strategy also add the newly discovered candidate cells to the updated belief  $B_{t+1}$ .

identified by the simulations, the robot moves in the real environment and observes the object using an object detector. If the object is detected, the POMCP exploration terminates and the robot approaches the object with the robust visual docking strategy as in [4]. Otherwise the environment map  $\mathcal{M}_t$  and pose graph  $\mathcal{G}_t$  are updated, and POMCP iterates until the maximum number of allowed steps  $\mu$  is reached.

### C. Belief reinvigoriation for dynamically growing state space

POMCP assigns the states to simulations by sampling particles  $b$  from Belief  $B$ , which is updated over time when the robot makes an action and perform an observation in the real environment.

In standard POMCP, the belief update is performed considering only particles of the previous belief that are sampled by simulations whose first observation is the one observed in the real environment after performing the optimal action. After the belief update, to avoid particle deprivation, new particles are added to the updated belief  $B_{t+1}$  by sampling from the previous  $B_t$  after executing the chosen action  $a_t$ . In such way, belief  $B_{t+1}$  tends to have a much higher probability in states that were already present before the reinvigoriation, thus discouraging actions towards newly added states, as can be seen in Figure 3.

Moreover, when dealing with an unknown environment where the state space changes over time, the particles at time  $t$  do not account for the part of the environment that will be observed at time  $t + 1$ . Hence the states in the particle filter, after the traditional belief update, cannot represent correctly the state space at time  $t + 1$ . To address the above-mentioned



Fig. 4: Evaluation dataset for AVS methods. The upper row shows sample images in AVD, captured with a real robotic platform in real apartments. The lower row shows sample 3D scans of complex scenes in *Habitat ObjectNav* Challenge.

limitations of the belief update strategy in standard POMCP, we propose a two-step belief reinvigoriation strategy. First, we map all the states in  $B_t$  to  $B_{t+1}$ , and remove the particles in simulations with their simulated observation different from the one made by the robot in the real environment. In this way, states in  $B_t$  are not penalised in  $B_{t+1}$  if they are not used as initial states by the simulations. Second, we add a new set of particles to  $B_{t+1}$  sampled with a uniform distribution over the new candidate cells introduced at time  $t + 1$ . Note that, without adding these states to  $B_{t+1}$ , simulations in next steps will assume that the object could never be in one of these newly discovered parts of the environment. The proposed strategy allows to tune the percentage of new particles introduced in the second step with respect to the previously discovered candidate cells. By doing so, we can tune the attention the robot puts on the newly discovered frontier cells.

## IV. EXPERIMENTS

We evaluate our proposed POMP++ on two public datasets featuring real-world scenarios that have been commonly used for benchmarking methods addressing the AVS task. Active Vision Dataset (AVD) [5] is captured with a RGB-D sensor equipped on a robotic platform moving within real apartments in a grid manner, while Habitat ObjectNav [27] provides a simulator for rendering observations within 3D scans of real large indoor environments provided by Matterport3D dataset [28] (see Fig. 4). We compare POMP++ against baselines and sota methods with both AVD (Section IV-A) and Habitat ObjectNav (Section IV-B) under their corresponding evaluation protocol and performance metrics for fair comparison. Moreover, we conduct an ablation study (Section IV-C) where we show the impact of the proposed belief reinvigoriation strategy and the impact of an imperfect object detector on AVS planning.

### A. Evaluation on AVD

AVD is the largest real-world dataset available for testing AVS, containing 17 scans of 9 real apartments recorded by



TABLE I: Results on the three test scenes from AVD using the GT annotations for object detection. We report the results of EAT [1] and POMP [4] as in their original paper. Results are averaged over multiple search episodes as defined in AVD benchmark, with the standard deviations presented within the parentheses. Results use known floor maps are highlighted in *italic*. Best results of methods without known floor maps are highlighted in **bold**.

Method	AVD											
	Easy			Medium			Hard			Avg.		
	SR $\uparrow$	APL $\downarrow$	ASPPL $\uparrow$	SR $\uparrow$	APL $\downarrow$	ASPPL $\uparrow$	SR $\uparrow$	APL $\downarrow$	ASPPL $\uparrow$	SR $\uparrow$	APL $\downarrow$	ASPPL $\uparrow$
<i>POMP (known env)</i> [4]	<i>0.98</i>	<i>13.6</i>	<i>0.72 (0.26)</i>	<i>0.73</i>	<i>17.1</i>	<i>0.8 (0.23)</i>	<i>0.56</i>	<i>20.5</i>	<i>0.72 (0.39)</i>	<i>0.76</i>	<i>17.1</i>	<i>0.75 (0.29)</i>
Random Walk	0.32	74	0.19 (0.35)	0.11	74.48	0.21 (0.36)	0.10	79.27	0.17 (0.18)	0.18	75.91	0.19 (0.29)
EAT [1]	0.77	<b>12.2</b>	-	0.73	<b>16.2</b>	-	0.58	<b>22.1</b>	-	0.69	<b>16.8</b>	-
<b>POMP++</b>	<b>1.0</b>	14.27	<b>0.70 (0.28)</b>	<b>0.79</b>	19.4	<b>0.76 (0.26)</b>	<b>0.84</b>	29.54	<b>0.61 (0.33)</b>	<b>0.87</b>	21.07	<b>0.68 (0.29)</b>

a robot equipped with a RGB-D camera. For each scene, AVD provides the RGB-D images captured from all possible robot’s poses, as well as the camera intrinsics and extrinsics. The action space  $A$  in AVD is {Forward, Backward, Turn\_Left, Turn\_Right, Stop}, with a translation step of 30 cm and a rotation step of 30°.

**Performance Metrics:** In the AVD Benchmark (AVDB), a run is considered successful when the robot terminates its exploration in one of the ending positions specified in AVDB, which are closest to the object with the object appearing in the camera view. There are three main metrics defined for benchmarking different methods: **Success Rate (SR)**, is the ratio of successful episodes over the total number of episodes; **Average Path Length (APL)** is the average number of poses visited by the robot among the paths that lead to a successful search (a lower value indicates a higher efficiency); and **Average Shortest Predicted Path Length (ASPPL)** is the average ratio between the length of the shortest possible path to reach a valid destination pose provided by AVDB and the length of the path generated by the method (a larger value indicates a higher absolute efficiency). We also compute the standard deviation of ASPPL to investigate the variability of path efficiency.

**Baselines and Comparisons:** We compare POMP++ against:

- **Random Walk:** a standard baseline where the robot move randomly for a certain number of steps.
- **EAT [1]:** a reinforcement learning approach trained on the remaining scenes of AVD. A Region proposal network (RPN) extracting object proposals is combined with a DRQN for the policy learning.
- **POMP [4]:** a POMCP-based solver that only works with known environments. This method requires as input the 2D floor map and the reachable robot poses.

Note that among the RL-based strategies, we consider EAT [1] and GAPLE [2] that both use AVD for evaluation (discussed in Sec. II). However, the evaluation protocol adopted by GAPLE is not documented, while EAT [1] shared their protocol explicitly. For a fair comparison, we follow the evaluation protocol in both EAT [1] and POMP [4], which perform search on a subset of objects under three test scenes with an increasing difficulty level. The test scenes are: *i*) Home\_005\_1, the easy scenario, consists only of a kitchen; *ii*) Home\_001\_2, the medium scenario, the robot has to navigate a living room and an open kitchen; *iii*)

Home\_003\_2 is the hard scenario, where the robot has to explore a living room, a half-open kitchen, a dining room, hallway and bathroom. As in [1] and [4], to only compare the goodness of the motion policy without the nuisances caused by underlying object detectors, we use the ground-truth annotations of the object for the evaluation, and we limit the search to 125 steps. For POMP++ we set the reward value to +1000 for when the object is detected and -1 for when it is not. The impact of a real-world detector on AVD is the subject of another experiment in Section IV-C.

**Results Discussion:** Table I shows the performances achieved by all methods on the three test scenes in AVD. We observe that on average POMP++ outperforms all the other methods in term of Success Rate. Even in complex episodes, where the target object is initialised far away from the robot, POMP++ is able to succeed the search task (see Supplementary Material for more search dynamics). Our path efficiency on the successful searches is lower compared to other methods as POMP++ has to explore more areas to look for the object, while POMP [4] exploits a known 2D floor map and EAT has already encoded the environment knowledge via its offline training. The improvement of POMP++ in terms of SR over POMP is mainly contributed by our new belief reinvigoration strategy. We provide more details in Section IV-C.

### B. Evaluation on Habitat ObjectNav

Habitat [27] is a simulation platform designed for embodied AI, where robots can move and observe within realistic 3D environments, *e.g.* 3D scans of real large scenes from Matterport3D [28] and Gibson [29]. In particular, the Habitat ObjectNav Challenge makes use of 11 scans of complex scenes in Matterport3D for validating methods (see the second row in Fig. 4). While the robot moves in the Habitat environment, the simulator renders the RGB-D images at each step taken from the current pose, as well as the camera intrinsics and extrinsics. The action space  $A$  in Habitat ObjectNav Challenge is the set of four actions {Forward, Turn\_Left, Turn\_Right, Stop}. We set the translation step to 25 cm and the rotation step to 30°.

**Performance Metrics:** In the Habitat ObjectNav Challenge, a success is defined if the distance between the robot and the target object is  $\leq 1$ m, and the object is within the camera’s field of view, without specifying the ending poses as in AVD. There are three main performance metrics defined

for the Habitat ObjectNav Challenge: **Success Rate (SR)**, is the ratio of successful runs over the total number of runs, the same as in AVD metrics; **Success weighted Path Length (SPL)**, measures the path efficiency of successful episodes<sup>3</sup>, defined as  $SPL = \frac{1}{N} \sum_{i=1}^N S_i \frac{l_i}{\max(l_i^*, l_i)}$ , where  $l_i$  is the length of the shortest path between the goal and the target for an episode,  $l_i^*$  is the length of path taken by a robot in an episode, and  $S_i$  is the binary indicator of a success in episode  $i$ . Finally, **Distance To Success (DTS)** is further defined by [3] which measures the distance of the robot from the success threshold boundary when the episode ends. DTS is computed as  $\max(\|x_T - x_G\| - d_s, 0)$ , where  $x_T$  is the robot location at the termination,  $x_G$  is the goal location at the end of the episode, and  $d_s$  is the success threshold.

**Baselines and comparisons:** We compare POMP++ against the winner entry in the Habitat ObjectNav Challenge, **SemExp** [3], as well as baselines reported in the literature:

- **Random walk:** standard baseline where the robot moves randomly for a certain number of steps.
- **RGBD + RL** [27]: a vanilla recurrent RL policy initialised with ResNet18 backbone followed by a GRU.
- **RGBD+Semantics+RL:** an adaptation from [30] which passes semantic segmentation and object detections along with RGBD input to a recurrent RL policy.
- **Classical Mapping + FBE** [31]: classical robotics pipeline for mapping followed by frontier-based exploration (FBE). When the detector finds the object, a local policy reaches the object with an analytical planner.
- **Active Neural SLAM** [32]: an exploration policy trained to maximise the coverage. Whenever the target object is detected, the same local policy as described above is applied to approach the object.
- **SemExp** [3]: a RL policy learner with a semantic mapping of the explored environment, and embedded object-to-object priors.

We follow the evaluation protocol defined in [3], where a subset of object categories that are common between MP3D and MS-COCO is chosen as target:  $\{\textit{chair, couch, potted plant, bed, toilet, tv}\}$ . For object detection, we used MaskRCNN with ResNet50 backbone pretrained on MS-COCO. Experiments are performed with 11 scenes in the validation set of Matterport3D, following the configuration in terms of robot starting pose and target object, provided by the Habitat ObjectNav Challenge. The walls are reconstructed online with the ground-truth semantic segmentation. The reward value is set to +1000 for when the object is detected and -1 for when it is not.

**Result Discussion:** Table II shows the results evaluated on the 11 test scenes in Habitat ObjectNav dataset. Our approach has a SR 6% higher than the best RL approach SemExp [3] and 10% higher than the Active Neural SLAM [32]. In terms of SPL, we are just slightly better than SemExp, which indicates that our path efficiency could be lower given the higher SR we achieve. Yet, our Distance To Success (DTS) is notably lower than all the other approaches,

<sup>3</sup>SPL is used as the main criterion in the Habitat ObjectNav Challenge

TABLE II: Results of POMP++ and baselines evaluated on 11 test scenes of Matterport3D, with the object detector pretrained on MS-COCO.

Method	Habitat		
	SPL $\uparrow$	SR $\uparrow$	DTS (m) $\downarrow$
Random	0.005	0.005	8.048
RGBD + RL [27]	0.017	0.037	7.654
RGBD + Sem + RL [30]	0.015	0.031	7.612
Classical Mapping + FBE [31]	0.117	0.311	7.102
Active Neural SLAM [32]	0.119	0.321	7.056
SemExp [3]	0.144	0.360	6.733
<b>POMP++</b>	<b>0.148</b>	<b>0.420</b>	<b>3.726</b>

TABLE III: Results of POMP++ and POMP+ on the Hard scenario from AVD, with both GT annotations and an object detector [5].

Method	AVD (Hard)		
	SR $\uparrow$	APL $\downarrow$	ASPPL $\uparrow$
POMP (GT, <i>known env</i> )	0.56	20.05	0.72(0.39)
POMP+ (GT, <i>known env</i> )	0.86	33.71	0.55(0.35)
POMP++ (GT)	0.84	29.54	0.61(0.33)
POMP (Detector, <i>known env</i> )	0.37	22.01	0.63(0.29)
POMP+(Detector, <i>known env</i> )	0.45	23.95	0.63(0.41)
POMP++ (Detector)	0.41	22.41	0.63(0.31)

meaning that our failed episodes tend to terminate nearer to the target location.

It is noticeable that the overall SR achieved by all methods are low due to the unsatisfactory 3D reconstruction quality of scenes (*e.g.* incomplete scans) and sometimes confusing object annotations (*e.g.* "bulletin board" as "tv\_monitor"). The use of a standard object detector without fine-tuning on the rendered images also contributes to the low performance since the detector is prone to errors. The impact of the object detector on AVS task is studied in Section IV-C.

Moreover, we notice that algorithms that use explicit semantic mapping of the environment, such as Classical Mapping+FBE [31], Active Neural SLAM [32], SemExp [3] and POMP++ are able to perform significantly better than RL-based methods [27] that rely on implicit scene representation via visual features extracted from the sequence of observations. However, as the scenes in Habitat ObjectNav varies a lot among each other, in terms of the scene types and dimensions, covering modern/ancient homes, palaces and even a spa. The merits of encoding sophisticated scene semantics can be limited, which allows POMP++ to outperform other methods with only 2D wall maps that are reconstructed online. Please refer to Supplementary Material for the state dynamics of both successful and failed searches as robot moves within the test scenes.

### C. Ablation study

We mainly evaluate the effectiveness of the proposed belief reinvigoration strategy and the impact of imperfect object detector. We compare POMP++ to POMP [4] and POMP+, which is POMP for known environments with our novel belief reinvigoration strategy, with both ground-truth annotation and the object detector provided by AVD [5]. Table III reports the averaged results obtained with the hard

scene Home\_003\_2 in AVD. POMP+ with the novel belief reinvigoration strategy improves the SR by 30% although at the cost of lengthier paths. SR is significantly improved by the new reinvigoration strategy because the original one tends to vanish particles that are not covered by successful simulations, which limits the potential of a successful search, as we explained in Section III-C. POMP++ achieves a slightly worse SR compared to POMP+ due to the lack of a known 2D floor map.

Moreover, the performance of the detector imposes a direct impact on the AVS performance, where we can observe the SR is reduced almost by half compared to the results obtained with GT annotations. This is because during the POMCP exploration, the output of the detector is used to either prune candidate object location in case of no detection; or as a stopping condition for the exploration if the target object is detected.

#### D. Real-time performance analysis

The complexity of POMP++ is  $O(\alpha \times (\beta \times \mu + \delta))$ , where  $\alpha$  is the number of steps performed in the real environment,  $\beta$  is the number of simulations performed for each step,  $\mu$  is the max number of steps per simulation,  $\delta$  is the number of actions to update the state space, the map, the transition and observation models. Note that all parameters can be tuned to satisfy real time performance when scaling the algorithm to large environments. When experimenting on AVD that contains realistic indoor housing environments, e.g. typical homes, we achieve 0.07 seconds per step on average with a standard machine with a 6-cores Intel i7-6800k CPU. Such processing speed is considered sufficient for real-time robotic applications.

### V. CONCLUSION

In this work we proposed a POMCP-based planner, POMP++, that solves the AVS problem in unknown environments, without the need of offline training. We introduced novel modifications to the POMDP's formulation, allowing for the search space to be dynamic, and grow as the robot explores the environment. Following this new formulation, we proposed a novel way to update the belief in the underlying POMDP solver, POMCP, which boosts the correct coverage of the expanding search space. We outperformed the sota methods on two benchmark datasets in terms of the success rate by a significant margin. As future work, we will further integrate scene graphs within the POMCP framework to better inject scene semantic for boosting the search efficiency of our method.

### REFERENCES

- [1] J. F. Schmid, M. Lauri, and S. Frintrop, "Explore, approach, and terminate: Evaluating subtasks in active visual object search based on deep reinforcement learning," in *IROS*, 2019.
- [2] X. Ye, Z. Lin, J. Lee, J. Zhang, S. Zheng, and Y. Yang, "GAPLE: Generalizable Approaching Policy LEarning for Robotic Object Searching in Indoor Environment," *IEEE Robotics and Automation Letters*, vol. 4, no. 4, pp. 4003–4010, 2019.
- [3] D. S. Chaplot, D. Gandhi, A. Gupta, and R. Salakhutdinov, "Object goal navigation using goal-oriented semantic exploration," in *NeurIPS*, 2020.
- [4] Y. Wang, F. Giuliani, R. Berra, A. Castellini, A. Del Bue, A. Farinelli, M. Cristani, and F. Setti, "POMP: Pomcp-based online motion planning for active visual search in indoor environments," in *BMVC*, 2020.
- [5] P. Ammirato, P. Poirson, E. Park, J. Košecká, and A. C. Berg, "A dataset for developing and benchmarking active vision," in *ICRA*, 2017.
- [6] D. Batra, A. Gokaslan, A. Kembhavi, O. Maksymets, R. Mottaghi, M. Savva, A. Toshev, and E. Wijnmans, "ObjectNav Revisited: On Evaluation of Embodied Agents Navigating to Objects," in *arXiv:2006.13171*, 2020.
- [7] T. D. Garvey, "Perceptual strategies for purposive vision." Ph.D. dissertation, Stanford University, Stanford, CA, USA, 1976, aAI7613006.
- [8] L. E. Wixson and D. H. Ballard, "Using intermediate objects to improve the efficiency of visual search," *International Journal of Computer Vision*, vol. 12, no. 2, pp. 209–230, Apr 1994.
- [9] T. Kollar and N. Roy, "Utilizing object-object and object-scene context when planning to find things," in *ICRA*, 2009.
- [10] K. Sjöö, A. Aydemir, and P. Jensfelt, "Topological spatial relations for active visual search," *Robotics and Autonomous Systems*, vol. 60, no. 9, pp. 1093–1107, 2012.
- [11] A. Aydemir, A. Pronobis, M. Göbelbecker, and P. Jensfelt, "Active visual object search in unknown environments using uncertain semantics," *IEEE Trans. on Robotics*, vol. 29, no. 4, pp. 986–1002, 2013.
- [12] M. Deitke, W. Han, A. Herrasti, A. Kembhavi, E. Kolve, R. Mottaghi, J. Salvador, D. Schwenk, E. VanderBilt, M. Wallingford, L. Weihs, M. Yatskar, and A. Farhadi, "RoboTHOR: An Open Simulation-to-Real Embodied AI Platform," in *CVPR*, 2020.
- [13] L. Kaelbling, M. Littman, and A. Cassandra, "Planning and Acting in Partially Observable Stochastic Domains," *Artificial Intelligence*, vol. 101, no. 1-2, pp. 99–134, 1998.
- [14] C. Papadimitriou and J. Tsitsiklis, "The Complexity of Markov Decision Processes," *Math. Oper. Res.*, vol. 12, no. 3, pp. 441–450, 1987.
- [15] M. Hauskrecht, "Value-Function Approximations for Partially Observable Markov Decision Processes," *JAIR*, vol. 13, pp. 33–94, 2000.
- [16] S. Ross, J. Pineau, S. Paquet, and B. Chaib-draa, "Online Planning Algorithms for POMDPs," *JAIR*, vol. 32, pp. 663–704, 2008.
- [17] D. Silver and J. Veness, "Monte-Carlo Planning in Large POMDPs," in *NeurIPS*, 2010.
- [18] R. Coulom, "Efficient Selectivity and Backup Operators in Monte-Carlo Tree Search," in *CG*, 2006.
- [19] L. Kocsis and C. Szepesvári, "Bandit Based Monte-Carlo Planning," in *ECML*, 2006.
- [20] C. Browne, E. Powley, D. Whitehouse, S. Lucas, P. Cowling, P. Rohlfshagen, S. Tavener, D. Perez, S. Samothrakis, and S. Colton, "A Survey of Monte Carlo Tree Search Methods," *IEEE Trans. Comp. Intell. AI Games*, vol. 4, no. 1, pp. 1–43, 2012.
- [21] C. Amato and F. A. Oliehoek, "Scalable Planning and Learning for Multiagent POMDPs," in *AAAI*, 2015.
- [22] J. Lee, G.-H. Kim, P. Poupart, and K.-E. Kim, "Monte-Carlo Tree Search for Constrained POMDPs," in *NeurIPS*, 2018.
- [23] A. Castellini, G. Chalkiadakis, and A. Farinelli, "Influence of State-Variable Constraints on Partially Observable Monte Carlo Planning," in *IJCAI*, 2019.
- [24] B. Yamauchi, "A frontier-based approach for autonomous exploration," in *CIRA*, 1997.
- [25] M. Lauri and R. Ritala, "Planning for robotic exploration based on forward simulation," *Robotics and Autonomous Systems*, vol. 83, pp. 15–31, 2016.
- [26] S. Choi, Q. Zhou, and V. Koltun, "Robust reconstruction of indoor scenes," in *CVPR*, 2015.
- [27] M. Savva, A. Kadian, O. Maksymets, Y. Zhao, E. Wijnmans, B. Jain, J. Straub, J. Liu, V. Koltun, J. Malik, D. Parikh, and D. Batra, "Habitat: A platform for embodied AI research," in *ICCV*, 2019.
- [28] A. Chang, A. Dai, T. Funkhouser, M. Halber, M. Niessner, M. Savva, S. Song, A. Zeng, and Y. Zhang, "Matterport3d: Learning from rgb-d data in indoor environments," *3DV*, 2017.
- [29] F. Xia, A. R. Zamir, Z.-Y. He, A. Sax, J. Malik, and S. Savarese, "Gibson env: real-world perception for embodied agents," in *CVPR*, 2018.
- [30] A. Mousavian, A. Toshev, M. Fišer, J. Košecká, A. Wahid, and J. Davidson, "Visual representations for semantic target driven navigation," in *ICRA*, 2019.



- [31] B. Yamauchi, "A frontier-based approach for autonomous exploration," in *CIRA*, 1997.
- [32] D. S. Chaplot, D. Gandhi, S. Gupta, A. Gupta, and R. Salakhutdinov, "Learning to explore using active neural slam," in *ICLR*, 2020.