

---

# Safe Multi-Agent Navigation guided by Goal-Conditioned Safe Reinforcement Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

Safe navigation is essential for autonomous systems operating in hazardous environments. Traditional planning methods are effective for solving long-horizon tasks but depend on the availability of a graph representation with predefined distance metrics. In contrast, safe Reinforcement Learning (RL) is capable of learning complex behaviors without relying on manual heuristics but fails to solve long-horizon tasks, particularly in goal-conditioned and multi-agent scenarios.

In this paper, we introduce a novel method that integrates the strengths of both planning and safe RL. Our method leverages goal-conditioned RL and safe RL to learn a goal-conditioned policy for navigation while concurrently estimating cumulative distance and safety levels using learned value functions via an automated self-training algorithm. By constructing a graph with states from the replay buffer, our method prunes unsafe edges and generates a waypoint-based plan that the agent then executes by following those waypoints sequentially until their goal locations are reached. This graph pruning and planning approach via the learned value functions allows our approach to flexibly balance the trade-off between faster and safer routes especially over extended horizons.

Utilizing this unified high-level graph and a shared low-level goal-conditioned safe RL policy, we extend this approach to address the multi-agent safe navigation problem. In particular, we leverage Conflict-Based Search (CBS) to create waypoint-based plans for multiple agents allowing for their safe navigation over extended horizons. This integration enhances the scalability of goal-conditioned safe RL in multi-agent scenarios, enabling efficient coordination among agents. Extensive benchmarking against state-of-the-art baselines demonstrates the effectiveness of our method in achieving distance goals safely for multiple agents in complex and hazardous environments.

## 1 Introduction

In many real-world scenarios like search-and-rescue or environmental monitoring, autonomous agents must navigate to distant locations safely as repairs can be costly and time-consuming. This challenge becomes even more critical in multi-agent settings, such as drone swarms or robotic teams, where all agents must coordinate to reach their goals while managing risks throughout the entire mission. Ensuring safety in such complex environments is paramount, as failure to do so could compromise the entire operation. Recent advancements in Multi-Agent Path Finding (MAPF) Stern et al. [2019], Yu and LaValle [2013] demonstrate the ability of planners to efficiently generate path plans for multiple agents by decomposing a multi-agent problem into a series of single-agent problems. However, these approaches assume access to a structured graph representation of the environment, consisting of valid nodes and edges, annotated with utility information like safety levels and estimated travel times,

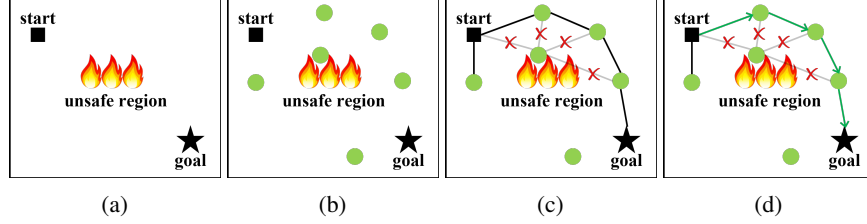


Figure 1: **Safe Navigation via Planning and Safe RL:** (a) Unconstrained goal-conditioned RL does not avoid unsafe regions in the environment. (b) We use observations in the replay buffer (green circles) as potential sub-goals. (c) We construct a graph using states from replay buffer. The distance and safety estimate of each edge is predicted by the agent’s value functions. We prune the edges that are too risky, as illustrated by the crossed out edges (d) We run graph search to find the sequence of sub-goals (green arrows). The goal-conditioned policy then follows the sub-goals to reach the goal. Similar idea is extended to multi-agent setting.

37 limiting their applicability in scenarios where agents only have access to high-level observations,  
38 such as images.

39 In contrast, learning-based approaches such as deep reinforcement learning (RL) excel at learning  
40 state utilities from high-dimensional observations such as images. In particular, goal-conditioned RL  
41 Mirowski et al. [2017], Schaul et al. [2015], Pong et al. [2020] has been effective in enabling agents  
42 to achieve goals of interest in various domains, including navigation and manipulation. However,  
43 learning policies for distant goals remains challenging, as performance often degrades as the task  
44 horizon increases Levy et al. [2019], Nachum et al. [2018] and requires significant reward shaping  
45 Chiang et al. [2019] or expert demonstrations Lynch et al. [2019], Nair et al. [2018], which can  
46 limit the policy’s ability to discover novel solutions. The challenge is compounded when agents  
47 must remain safe while navigating over long horizons. While constrained RL Altman [2021] has  
48 demonstrated the potential to incorporate safety constraints into policy learning, reaching a diverse  
49 set of distant goals under such constraints, particularly in multi-agent contexts, remains as an open  
50 research problem.

51 Prior research has tackled long-horizon problems using goal-conditioned RL combined with graph  
52 search techniques Eysenbach et al. [2019]. By constructing an intermediate graph representation of  
53 the environment from the replay buffer and using search-based methods to generate a waypoint-based  
54 plan, these methods decompose tasks into a series of shorter and easier goal-reaching tasks. However,  
55 existing approaches are limited to single-agent scenarios without considering any safety aspects. This  
56 paper extends these methods to address multi-agent, long-horizon, goal-conditioned tasks ensuring  
57 that the agents reach their goals safely.

58 As shown in Figure 1, our approach leverages goal-conditioned RL and safe RL via a novel self-  
59 training algorithm to learn distance and safety predictors which allows us to dynamically construct an  
60 intermediate graph from the replay buffer. This enables greater planning flexibility that balances the  
61 trade-off between faster and safer routes depending on user-preference. By utilizing this intermediate  
62 graph representation, we extend the range of problems that MAPF approaches can solve, especially in  
63 challenging, image-based environments, effectively bridging the gap between planning and learning  
64 in multi-agent settings. Additionally, we empirically show that our approach can plan a series of  
65 safer sub-goals for multiple agents without requiring new training for each additional agent, thereby  
66 achieving scalability. In summary, we present a scalable and safe approach for navigating multiple  
67 agents to their distance goals in sparse-reward environments, broadening the scope of MAPF through  
68 the integration of goal-conditioned safe RL.

## 69 2 Approach

### 70 2.1 Goal-Conditioned Safe Reinforcement Learning

71 Inspired by SoRB Eysenbach et al. [2019] (see Appendix A.1.2), we would like to train an agent that  
72 can predict both distance and cost among state pairs. Meanwhile, the agent should also be capable of  
73 safely navigating amongst them. Since the value functions are conditioned on the associated policy, in

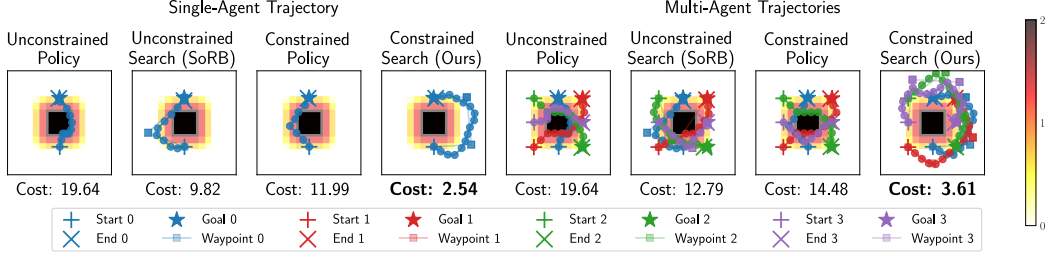


Figure 2: Comparison of different approaches on the 2D navigation problem for a single agent (left) and four agents (right). Baselines tend to traverse high-cost regions, leading to suboptimal paths. In contrast, our approach finds a longer but lower-cost path by avoiding these regions. For both unconstrained and constrained policies, collision avoidance in the multi-agent case is ensured by having agents wait when another agent is nearby.

order to obtain the distance and cost predictions following the shortest path, it is critical to first train an unconstrained agent that includes a goal-conditioned policy  $\pi(s, a, s_g)$ , a  $Q$ -function for cumulative reward  $Q(s, a, s_g)$ , and another  $Q$ -function for cumulative cost  $Q_C(s, a, s_g)$ . For the agent to learn safe navigation, we then fine-tune the unconstrained policy  $\pi$  by applying constrained RL methods which outputs a goal-conditioned safe policy  $\pi_c(s, a, s_g)$  and its associated value functions. As previously noted in Appendix A.1.3, the cost limit from constrained reinforcement learning (CRL) is not directly applicable to the goal-conditioned setting. Therefore, we redefine the cost limit as a hyper-parameter, a soft constraint that penalizes unsafe actions without requiring strict compliance. However, we note that when the choice of cost limit is too low (penalty is too high), the resultant policy would focus too much on reducing costs rather than reaching the goals. In this work, we use the Lagrangian actor critic Ray et al. [2019] thanks to its simplicity.

## 2.2 Self-Sampling and Training

To train a goal-conditioned agent, a set of goal states need to be sampled for the agent to explore and practice reaching them. The design of sampling strategy is crucial for the success of training a goal-conditioned agent. SoRB Eysenbach et al. [2019] accomplishes this by prioritizing nearby goal samples to improve the success rate of short-horizon navigation. However, it relies on an oracle with access to the ground-truth positions rather than the raw observations (e.g., images) to sample goals of desirable distances. This practice is unrealistic as such oracles often do not exist. Instead, we propose Algorithm 1, a self-sampling and training algorithm which, in essence, uses the agent’s  $Q$ -functions to evaluate and select training samples with a diverse set of distances and cost targets. As the  $Q$ -functions become more accurate, the quality of the samples improve, which then generates better training samples to improve the policy.

## 2.3 Categorical Cost Estimate

Training a robust goal-conditioned value function is challenging yet crucial for planning efficient waypoints Eysenbach et al. [2019]. Incorporating distributional RL Bellemare et al. [2017] has shown to be an effective strategy in improving the quality of goal-conditioned value functions thereby enhancing the planning process. Therefore, we follow Bellemare et al. [2017] to model the goal-conditioned value function for costs as a categorical distribution  $Z_\pi \sim \text{Cat}(N, p_\pi(z | s, a, g))$  within a range of  $[V_{\min}^c, V_{\max}^c]$ ,  $V_{\min}^c, V_{\max}^c \in \mathbb{R}^+$ , where  $N \in \mathbb{N}$  is the number of evenly discrete classes,  $p_\pi(z | s, a, g)$  is the probability predicted by the cost value function conditioned on the current policy  $\pi$ . The class supports can be described by a set of atoms  $\{z_i = V_{\min}^c + i\Delta z : 0 \leq i < N\}$ ,  $\Delta z \triangleq (V_{\max}^c - V_{\min}^c)/(N - 1)$ . Consequently, the expected value of the distribution is used estimate the goal-conditioned cost:

$$c_\pi(s_1, g) = \mathbb{E}_{p(z|s,a,g)}[z] \quad (1)$$

Since  $Z_\pi$  is a categorical distribution, standard Bellman update is not applicable. Instead, we follow the Categorical Algorithm Bellemare et al. [2017] to update  $Q_C(s, a, s_g)$ .

## 109 2.4 Graph Construction from Replay Buffer

110 With a trained agent, we use the observations (e.g., images) in the replay buffer  $\mathcal{B}$  to construct a  
 111 weighted, directed graph  $\mathcal{G}$ . Each edge has a predicted distance and cost for the state pair  $(s_i, s_j)$ :

$$d_{sp} \approx d_\pi \leftarrow Q(s, \pi(s_i, a, s_g), s_j) \quad (2)$$

$$c_{sp} \approx c_\pi \leftarrow Q_c(s, \pi(s_i, a, s_g), s_j) \quad (3)$$

112

$$\begin{aligned} \mathcal{G} &\triangleq (\mathcal{V}, \mathcal{E}, \mathcal{W}_d, \mathcal{W}_c) \\ \text{where } \mathcal{V} &= \mathcal{B} \\ \mathcal{E} &= \mathcal{B} \times \mathcal{B} = \{e_{s_i \rightarrow s_j} | s_i, s_j \in \mathcal{B}\} \\ \mathcal{W}_d(e_{s_i \rightarrow s_j}) &= \begin{cases} d_\pi(s_i, s_j) & \text{if } d_\pi(s_i, s_j) < d_{\max} \\ \infty & \text{otherwise} \end{cases} \\ \mathcal{W}_c(e_{s_i \rightarrow s_j}) &= \begin{cases} c_\pi(s_i, s_j) & \text{if } c_\pi(s_i, s_j) < c_{\max} \\ \infty & \text{otherwise} \end{cases} \end{aligned}$$

113 We exclude edges whose distance exceeds a maximum cutoff,  $d_{\max}$ , or whose cost surpasses a  
 114 maximum cutoff,  $c_{\max}$ . Both  $d_{\max}$  and  $c_{\max}$  are treated as hyper-parameters.

## 115 2.5 Conflict-Based Search (CBS)

116 CBS Sharon et al. [2015] is a widely used MAPF A.1.4 algorithm that addresses the problem’s  
 117 exponential state space by decomposing it into a series of single-agent path planning problems which  
 118 are then combined using efficient conflict-resolution strategies. CBS operates on two levels: the  
 119 high-level search, which manages different constraints on agents’ path plans, and the low-level search,  
 120 where a space-time A\* algorithm is used to find paths for each agent independently while satisfying  
 121 the constraints imposed by the high-level search tree.

122 In CBS, constraints are typically defined in the form  $\langle a_i, v, t \rangle$ , indicating that agent  $a_i$  is prohibited  
 123 from visiting vertex  $v \in \mathcal{V}$  at time step  $t$ . The high level search constructs a *constraint tree* starting  
 124 from a root node with no constraints. The tree expands in a best-first manner where each successor  
 125 node inherits its parent’s constraints and adds a new constraint for a single agent. The goal of CBS is  
 126 to reach a node in the constraint tree where all paths of all agents are conflict-free. This two-level  
 127 structure of CBS effectively manages the complexity of MAPF problems by isolating single-agent  
 128 planning tasks and resolving conflicts incrementally.

## 129 2.6 High-Level Overview

130 In summary, as shown in Algorithm 2 our approach leverages goal-conditioned safe RL augmented  
 131 with self-sampling and training algorithm 1 to build a graph from the replay buffer annotated with  
 132 predicted distances and cost estimates using the  $Q$ -functions  $Q^\pi, Q_c^\pi$  of the *unconstrained* policy  
 133  $\pi$ . CBS then utilizes this intermediate graph to generate effective and safe sequential waypoints for  
 134 the agents, which are followed using the trained low-level goal-conditioned *safe* RL policy  $\pi_c$  for  
 135 enhanced bi-level safety.

## 136 3 Conclusions

137 In conclusion, this work introduces a novel hierarchical framework that combines high-level  
 138 multi-agent planning using MAPF approaches, such as CBS, with low-level control through goal-  
 139 conditioned safe RL policies. The framework leverages learned distance and *cost* critics to automati-  
 140 cally generate and annotate dense graph networks from the replay buffer using self-sampling and  
 141 training. At the high level, our approach utilizes CBS to plan efficient and safe waypoints for multiple  
 142 agents, while at the low level, the RL policy guides the agents to follow these waypoints toward  
 143 distant goal-reaching tasks safely. This bi-level safety control ensures safe multi-agent behavior  
 144 throughout execution, effectively tackling scenarios ranging from simple 2D navigation to more  
 145 complex visual navigation challenges whose results can be found in the Appendix.

## 146 A Appendix

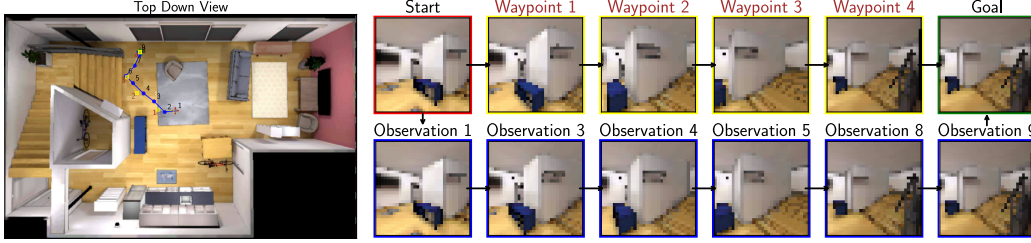


Figure 3: **Safe Visual Navigation:** The top-down map illustrates the agent’s path and waypoints according to our approach. The first row displays the forward-facing start, goal, and suggested waypoint images. Bottom row shows the agent’s forward-facing observations during execution. The waypoints help the agent avoid the couch and blue furniture. Note that the top-down view is not available to the agent.

### 147 A.1 Preliminaries

#### 148 A.1.1 Goal-Conditioned Reinforcement Learning

149 In goal-conditioned RL (GCRL), an agent interacts with an environment modeled as an augmented  
 150 Markov Decision Process (MDP)  $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$ , where  $\mathcal{S}$  is the state space,  $\mathcal{A}$  is the action space,  
 151  $P(s_{t+1}|s_t, a_t, s_g)$  is the transition probability,  $R(s, a, s_g)$  is the reward function, and  $\gamma \in [0, 1]$  is  
 152 the discount factor. Unlike vanilla MDP, GCRL includes a goal  $s_g \in \mathcal{S}$  in the problem formulation.  
 153 The agent’s objective is to learn a policy  $\pi(a|s, s_g)$  that maximizes the expected cumulative reward,  
 154 conditioned on both the current state  $s$  and the goal  $s_g$ . The reward function  $R(s, a, s_g)$  is goal-  
 155 dependent, reflecting how well the actions align with the desired goal. In practice, GCRL excels  
 156 at learning short-horizon skills often leveraging reward shaping Chiang et al. [2019] and human  
 157 demonstrations Lynch et al. [2019], Nair et al. [2018] to enhance performance.

#### 158 A.1.2 Bridge Planning and RL

159 To address long-horizon problems, Eysenbach et al. [2019] proposed Search on Replay Buffer (SoRB)  
 160 that effectively integrates RL with single-agent planning. Specifically, SoRB defines  $R(s, a, s_g) = -1$   
 161 and  $\gamma = 1$  for the target environment. It then trains a goal-conditioned policy  $\pi(a | s, s_g)$  and the  
 162 corresponding value function  $V(s, s_g)$  using standard off-policy RL algorithms Lillicrap et al. [2019].  
 163 Assuming that  $\pi(a | s, s_g)$  is optimal, the value function  $V(s, s_g)$  can be used to approximate the  
 164 distance between a state  $s$  and a goal state  $s_g$ :

$$V(s, s_g) = -d_{\text{sp}}(s, s_g) \quad (4)$$

165 where  $d_{\text{sp}}$  is the shortest path distance from  $s$  to  $s_g$ , or equivalently, the expected number of steps to  
 166 reach  $s_g$  from  $s$  under the optimal policy  $\pi$ .

167 SoRB then constructs a weighted, directed graph  $\mathcal{G}$  from states that are either randomly sampled or  
 168 drawn from the replay buffer. Each node in the graph corresponds to an observation and edges are  
 169 added based on the predicted distance  $V(s_1, s_2)$  between the node pairs  $(s_1, s_2)$ . Node pairs with  
 170  $V(s_1, s_2) > d_{\text{max}}$  are excluded as they are considered potentially unreachable.

171 Finally, SoRB generates a waypoint-based plan by searching on this graph  $\mathcal{G}$ . The agent sequentially  
 172 navigates to the individual waypoints before heading towards the final goal. This process is executed  
 173 in either an open-loop or closed-loop manner. Notably, SoRB employs a greedy strategy and does not  
 174 take the safety of the agent into account. However, when safety is factored in, longer and safer paths  
 175 may be preferred over shorter yet riskier paths.

#### 176 A.1.3 Constrained Reinforcement Learning

177 To capture the concept of risk and safety, we build our learning algorithm within the framework  
 178 of Constrained RL (CRL). CRL solves a constrained Markov Decision Process (CMDP) Altman

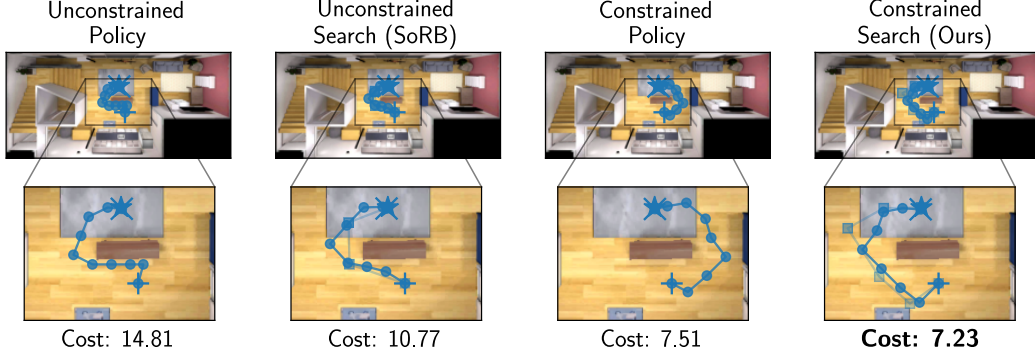


Figure 4: Comparison of different approaches on the SC3 Staging 11 map for single agent shows that unconstrained baselines fail to avoid the furniture, whereas the constrained policy and our approach successfully navigates around it. Our approach goes further by re-routing the agents away from the furniture, ensuring the safest behavior. Note that for both unconstrained and constrained policies, inter-agent collision avoidance is ensured by having agents wait when another agent is nearby.

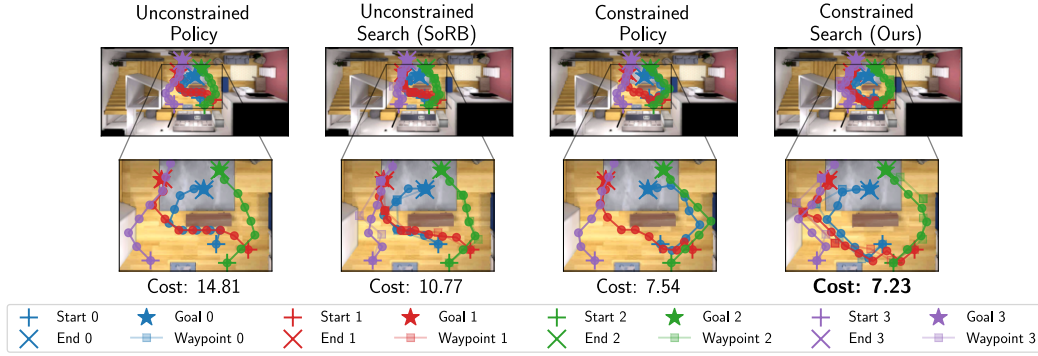


Figure 5: Comparison of different approaches on the SC3 Staging 11 map for four agents shows that unconstrained baselines fail to avoid the furniture, whereas the constrained policy and our approach successfully navigates around it. Our approach goes further by re-routing the agents away from the furniture, ensuring the safest behavior. Note that for both unconstrained and constrained policies, inter-agent collision avoidance is ensured by having agents wait when another agent is nearby.

[2021] problem that augments MDP with extra constraints that restrict the set of allowable policies. CMDP defines an *auxiliary cost* function  $C : S \times A \times S \rightarrow \mathbb{R}^+$  and a cost limit  $\delta \in \mathbb{R}^+$ . For conciseness, we will refer to auxiliary cost as *cost*. The *cost* function serves as a generic safety measure. For example, regions of high environmental uncertainty may be deemed unsafe and induce higher costs. In this work, we do not require the *cost* function to carry any physical interpretation. Let  $J_C(\pi) = \mathbb{E}_{\tau \sim \pi} [\sum_{t=0}^{\infty} \gamma^t C(s_t, a_t, s_{t+1})]$  denote the expected discounted *cost return* following policy  $\pi$ . The set of feasible stationary policies for a CMDP is then

$$\Pi_C \doteq \{\pi \in \Pi : J_C(\pi) \leq d\} \quad (5)$$

and the RL problem in a CMDP is to find a policy such that

$$\pi^* = \operatorname{argmax}_{\pi \in \Pi_C} J(\pi) \quad (6)$$

Standard off-policy methods can be adapted to solve CRL problems by including Lagrange multipliers to penalize unsafe actions that lead to the violation of *cost* constraints Ray et al. [2019]. It is crucial to note that, in CRL, the *cost* function and *cost* limit is task-dependent and is not applicable to the goal-conditioned setting where each goal is akin to a distinct task. For a more detailed reference on CRL, we refer our readers to Schulman et al. [2015, 2017], Sutton and Barto [2018].

---

**Algorithm 1** Self-Sampling and Training of Goal-Conditioned Actor Critic

---

**Inputs:** Environment  $E$ , Goal-Conditioned Policy  $\pi$ , Associated  $Q$ -functions for rewards or auxiliary costs  $Q(s, a, s_g)$ , Desirable sample target for rewards or costs  $\tau$ , Population size  $N$ , Number of training problems per batch  $K$

**Outputs:** Updated  $\pi(s, a, s_g)$ ,  $Q(s, a, s_g)$

- 1:  $\mathcal{P} \leftarrow$  Initialize an empty training set
  - 2: **for** each batch **do**
  - 3:    $\{(s_i, s_j)\}_N \leftarrow$  Randomly sample state pairs from  $E$
  - 4:    $\{v_{ij}\}_N \leftarrow Q(s_i, \pi(s_i, a, s_j), s_j)$
  - 5:    $\{l_{ij}\}_N \leftarrow L^2$  distance to  $\tau$
  - 6:   Find  $\{(s_i, s_j)\}$  corresponding to lowest  $K$   $\{l_{ij}\}_N$
  - 7:   Add the selected  $\{(s_i, s_j)\}$  to the training set
  - 8:   Train the agent with  $\mathcal{P}$  until  $\mathcal{P}$  is depleted
- 

#### 192 A.1.4 Multi-Agent Path Finding (MAPF)

193 A Multi-Agent Path Finding (MAPF) problem Stern et al. [2019], Yu and LaValle [2013] involves a  
194 weighted graph  $\mathcal{G}(\mathcal{V}, \mathcal{E})$  and a set of agents  $A = a_1, a_2, \dots, a_k$ , each with a unique start  $s_i \in \mathcal{V}$  and  
195 a goal  $g_i \in \mathcal{V}$ . Agents can either wait or move to adjacent vertices along edges in  $\mathcal{E}$ , with each action  
196 incurring a penalty defined by  $\mathcal{G}$ . Once an agent reaches its goal, no further penalty is incurred.

197 The path length for each agent is the weighted sum of the edges it traverses. A valid MAPF solution  
198 requires that all agents avoid *vertex collision* (two agents occupying the same vertex simultaneously)  
199 and *edge collisions* (two agents traversing the same edge in opposite directions). The typical objective  
200 in MAPF problems is to minimize the Sum of Costs (SoC), which is the total path length for all  
201 agents.

---

**Algorithm 2** Plan and Act

---

**Inputs:** Agents  $N$ , Current States  $s$ , Goal States  $s_g$ , Buffer of observations  $\mathcal{B}$ ,  
Learned constrained policy  $\pi_c$ ,  $Q$ -functions  $Q^\pi$ ,  $Q_c^\pi$  of the unconstrained policy  $\pi$

**Output:** Action  $a$

- 1: **if**  $N > 1$  **then**
  - 2:    $s_{w_1}, \dots \leftarrow \text{CBS}(s, s_g, \mathcal{B}, Q^\pi, Q_c^\pi)$
  - 3: **else**
  - 4:    $s_{w_1}, \dots \leftarrow \text{SHORTEST\_PATH}(s, s_g, \mathcal{B}, Q^\pi, Q_c^\pi)$
  - 5: **if**  $d_\pi(s \rightarrow s_{w_1}) < d_\pi(s \rightarrow s_g)$  or  $d_\pi(s \rightarrow s_g) > d_{\max}$  **then**
  - 6:    $a \leftarrow \pi_c(a \mid s, s_{w_1})$
  - 7: **else**
  - 8:    $a \leftarrow \pi_c(a \mid s, s_g)$
- 

#### 202 A.2 Experiments

203 Our experiments address the following questions to evaluate the effectiveness of our approach:

- 204 **Q1:** Can our approach demonstrate safer agent behavior compared to other methods?
- 205 **Q2:** Does our approach sacrifice its ability to reach distant goals when balancing the accumulated  
206 *cost* compared to other baselines?
- 207 **Q3:** Does the suggested approach effectively demonstrate superior performance when scaled to  
208 larger number of agents compared to other methods?

209 For the single-agent experiments, we compare our approach against a goal-conditioned RL policy,  
210 SoRB and a goal-conditioned safe RL policy. For the multi-agent experiments, we evaluated both  
211 off-policy and on-policy Multi-Agent RL (MARL) approaches, including MADDPG Lowe et al.  
212 [2020] and MAPPO Yu et al. [2022]. However, these methods exhibited the same limitations as  
213 goal-conditioned RL approaches, ultimately failing to enable multiple agents to reach their goals.  
214 Hence, we compare our approach with an extended version of SoRB that utilises CBS to generate



intermediate waypoints for multiple agents. Finally, we compare our approach with these baselines over two different tasks: a 2D navigation problem and a visual navigation problem that involves using images as observations.

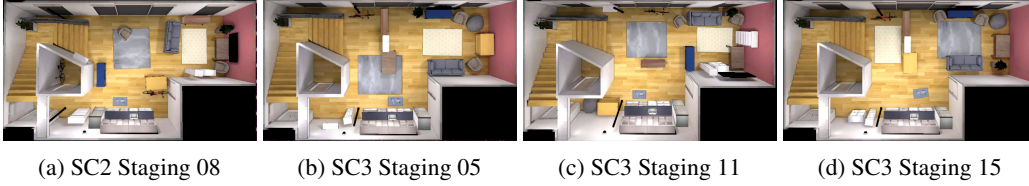


Figure 6: Different interactive digital-twins from the ReplicaCAD dataset we use to generate visual navigation problems.

### A.2.1 Environments

We use a 2D navigation environment with Central Obstacle map, where the agent observations are their locations, represented as  $s = (x, y) \in \mathbb{R}^2$ . The agent’s actions,  $a = (dx, dy) \in [-1, 1]^2$ , are added to the current position at each time step to update their location. For the visual navigation problems, we selected four distinct environments from the rich, interactive digital-twin ReplicaCAD dataset Straub et al. [2019], integrated with Habitat-Sim Szot et al. [2021], Savva et al. [2019]. In these scenarios, the agent’s observations consist of RGB panoramic images created by concatenating  $32 \times 32$  first-person views from each of the four cardinal directions. The action space for the agent in the visual navigation problems is similar to that in the 2D navigation scenario. For all experiments, we choose the *cost* limit  $\delta = 10$ , and we use a *cost* function that linearly decreases from the obstacle boundaries.

$$C(s, a) = \begin{cases} 2 - 2h/r & 0 \leq h \leq r \\ 0 & \text{otherwise} \end{cases} \quad (7)$$

Problem Configurations		Methods (Cumulative Cost (Success Rate))			
Problem Type	Agents	Unconstrained Policy	Unconstrained Search (SoRB)	Constrained Policy	Constrained Search (Ours)
Easy	1	<b>0.46 ± 1.05 (100%)</b>	1.38 ± 2.62 (100%)	0.47 ± 1.07 (100%)	0.49 ± 1.06 (100%)
	5	N/A	6.39 ± 4.66 (100%)	N/A	<b>1.68 ± 1.21 (100%)</b>
	10	N/A	8.69 ± 4.05 (100%)	N/A	<b>2.24 ± 1.03 (100%)</b>
	20	N/A	10.51 ± 3.57 (100%)	N/A	<b>2.72 ± 0.80 (100%)</b>
Medium	1	1.58 ± 2.03 (100%)	2.30 ± 3.02 (100%)	<b>1.58 ± 2.01 (100%)</b>	1.59 ± 1.99 (100%)
	5	N/A	5.12 ± 3.65 (100%)	N/A	<b>3.09 ± 1.50 (100%)</b>
	10	N/A	7.27 ± 3.66 (98%)	N/A	<b>4.11 ± 1.48 (98%)</b>
	20	N/A	8.73 ± 3.73 (98%)	N/A	<b>4.78 ± 1.33 (98%)</b>
Hard	1	3.98 ± 4.40 (100%)	4.05 ± 4.19 (100%)	4.19 ± 4.23 (100%)	<b>3.09 ± 3.87 (100%)</b>
	5	N/A	8.58 ± 3.79 (100%)	N/A	<b>6.01 ± 1.96 (100%)</b>
	10	N/A	10.77 ± 3.51 (100%)	N/A	<b>7.24 ± 1.52 (100%)</b>
	20	N/A	11.96 ± 3.31 (100%)	N/A	<b>8.36 ± 1.58 (100%)</b>

Table 1: **2D Navigation:** Comparison of accumulated cost across different approaches with varying problem difficulty and number of agents. Our approach is consistently able to plan successful and safest paths for multiple agents. Note that we extended SoRB by integrating it with CBS to plan waypoints for multiple agents.

Figure 2 illustrates example trajectories generated by both the baselines and our approach in single-agent and multi-agent settings in the 2D navigation task for a qualitative comparison. Figure 2 shows that, while the baseline methods often fail to generate safer paths when the sampled start and goal locations fall within high-*cost* regions, our approach successfully identifies a path that, while counter-intuitive from the perspective of the RL agent, aligns with human intuition. By navigating around the high-*cost* regions, our approach effectively minimizes the accumulated *cost*.

We now examine the effectiveness of our approach on the more challenging visual navigation problem as shown in Figure 3. Figures 5 present qualitative comparison with example trajectories generated by both the baselines and our approach. The results are consistent with those from the 2D navigation



problem, demonstrating that our approach enables agents to re-route their paths away from obstacles, such as furniture in the room, leading to safer navigation behavior.

Our comprehensive quantitative experiments involved randomly generating problems with varying difficulty levels, placing start and goal positions closer together for easier tasks and further apart for harder ones to increase complexity. Each experiment consisted of 50 different problem instances for each difficulty level. In our multi-agent experiments, we compared our approach with an extended version of SoRB, varying the number of agents from 5 to 20 for the 2D navigation problem, and from 5 to 10 agents for the visual navigation problem.

On the simpler 2D navigation problem, Table 1 shows that our approach consistently generates safer paths with the lowest accumulated *cost* compared to the other baselines, both in single-agent and multi-agent settings. This safer behavior becomes more pronounced as the number of agents increases, with no loss in success rate.

Problem Configurations			Methods (Cumulative Cost (Success Rate))			
Map	Problem Type	Agents	Unconstrained Policy	Unconstrained Search (SoRB)	Constrained Policy	Constrained Search (Ours)
SC2 Staging 08	Easy	1	3.16 $\pm$ 9.40 (100%)	2.14 $\pm$ 2.22 (98%)	1.85 $\pm$ 1.80 (100%)	<b>1.28 <math>\pm</math> 1.69 (100%)</b>
		5	N/A	4.32 $\pm$ 2.29 (98%)	N/A	<b>3.35 <math>\pm</math> 1.48 (100%)</b>
		10	N/A	5.36 $\pm$ 2.60 (96%)	N/A	<b>4.02 <math>\pm</math> 1.15 (100%)</b>
	Medium	1	4.98 $\pm$ 3.46 (98%)	5.15 $\pm$ 3.32 (98%)	4.77 $\pm$ 3.25 (100%)	<b>3.35 <math>\pm</math> 3.23 (100%)</b>
		5	N/A	9.74 $\pm$ 11.13 (90%)	N/A	<b>6.18 <math>\pm</math> 1.64 (98%)</b>
		10	N/A	14.23 $\pm$ 23.53 (84%)	N/A	<b>7.30 <math>\pm</math> 1.69 (98%)</b>
	Hard	1	13.41 $\pm$ 4.00 (96%)	12.35 $\pm$ 5.10 (90%)	10.43 $\pm$ 3.07 (100%)	<b>9.62 <math>\pm</math> 3.66 (100%)</b>
		5	N/A	21.23 $\pm$ 14.64 (84%)	N/A	<b>13.02 <math>\pm</math> 2.38 (98%)</b>
		10	N/A	27.37 $\pm$ 19.05 (62%)	N/A	<b>13.87 <math>\pm</math> 2.17 (92%)</b>
SC3 Staging 05	Easy	1	1.87 $\pm$ 1.65 (98%)	2.23 $\pm$ 1.85 (98%)	1.80 $\pm$ 1.58 (100%)	<b>1.25 <math>\pm</math> 1.60 (100%)</b>
		5	N/A	3.71 $\pm$ 1.96 (98%)	N/A	<b>2.76 <math>\pm</math> 1.36 (100%)</b>
		10	N/A	4.20 $\pm$ 1.73 (98%)	N/A	<b>3.32 <math>\pm</math> 1.11 (100%)</b>
	Medium	1	5.06 $\pm$ 3.27 (100%)	4.87 $\pm$ 3.15 (100%)	4.38 $\pm$ 2.91 (100%)	<b>2.66 <math>\pm</math> 2.32 (100%)</b>
		5	N/A	7.58 $\pm$ 2.16 (98%)	N/A	<b>5.88 <math>\pm</math> 1.38 (100%)</b>
		10	N/A	8.83 $\pm$ 1.75 (98%)	N/A	<b>7.28 <math>\pm</math> 1.53 (96%)</b>
	Hard	1	14.86 $\pm$ 4.9 (98%)	13.16 $\pm$ 4.24 (100%)	8.59 $\pm$ 2.21 (98%)	<b>6.47 <math>\pm</math> 3.30 (98%)</b>
		5	N/A	17.79 $\pm$ 2.68 (96%)	N/A	<b>15.10 <math>\pm</math> 2.19 (96%)</b>
		10	N/A	19.56 $\pm$ 2.66 (94%)	N/A	<b>16.71 <math>\pm</math> 2.28 (94%)</b>
SC3 Staging 11	Easy	1	1.20 $\pm$ 1.56 (100%)	1.27 $\pm$ 1.35 (100%)	1.20 $\pm$ 1.55 (100%)	<b>0.78 <math>\pm</math> 1.12 (100%)</b>
		5	N/A	2.32 $\pm$ 1.11 (100%)	N/A	<b>1.88 <math>\pm</math> 0.90 (100%)</b>
		10	N/A	3.00 $\pm$ 1.07 (100%)	N/A	<b>2.54 <math>\pm</math> 1.15 (100%)</b>
	Medium	1	3.66 $\pm$ 2.95 (98%)	4.19 $\pm$ 3.23 (98%)	3.40 $\pm$ 2.83 (98%)	<b>1.95 <math>\pm</math> 2.24 (100%)</b>
		5	N/A	8.33 $\pm$ 2.37 (98%)	N/A	<b>5.93 <math>\pm</math> 2.40 (96%)</b>
		10	N/A	10.23 $\pm$ 2.27 (96%)	N/A	<b>8.13 <math>\pm</math> 2.18 (94%)</b>
	Hard	1	16.41 $\pm$ 6.50 (100%)	14.75 $\pm$ 5.43 (100%)	11.59 $\pm$ 3.60 (100%)	<b>10.93 <math>\pm</math> 3.94 (100%)</b>
		5	N/A	18.67 $\pm$ 2.51 (96%)	N/A	<b>15.78 <math>\pm</math> 2.16 (98%)</b>
		10	N/A	20.37 $\pm$ 2.41 (92%)	N/A	<b>16.94 <math>\pm</math> 1.66 (92%)</b>
SC3 Staging 15	Easy	1	1.90 $\pm$ 1.76 (100%)	2.08 $\pm$ 1.73 (100%)	1.87 $\pm$ 1.77 (100%)	<b>1.47 <math>\pm</math> 1.39 (100%)</b>
		5	N/A (100%)	3.82 $\pm$ 3.53 (100%)	N/A	<b>2.77 <math>\pm</math> 1.32 (100%)</b>
		10	N/A	4.83 $\pm$ 3.68 (100%)	N/A	<b>3.63 <math>\pm</math> 1.39 (100%)</b>
	Medium	1	3.76 $\pm$ 3.42 (100%)	4.16 $\pm$ 2.79 (100%)	3.23 $\pm$ 2.51 (100%)	<b>2.01 <math>\pm</math> 1.75 (100%)</b>
		5	N/A	6.70 $\pm$ 1.81 (100%)	N/A	<b>5.18 <math>\pm</math> 2.32 (100%)</b>
		10	N/A	10.21 $\pm$ 14.89 (96%)	N/A	<b>6.66 <math>\pm</math> 1.59 (100%)</b>
	Hard	1	13.83 $\pm$ 4.35 (100%)	14.37 $\pm$ 11.06 (98%)	9.75 $\pm$ 2.74 (100%)	<b>7.61 <math>\pm</math> 3.10 (100%)</b>
		5	N/A	19.23 $\pm$ 11.03 (88%)	N/A	<b>11.63 <math>\pm</math> 2.70 (100%)</b>
		10	N/A	19.70 $\pm$ 3.73 (84%)	N/A	<b>12.85 <math>\pm</math> 2.12 (98%)</b>

Table 2: **Safe Visual Navigation:** Comparison of accumulated cost across different approaches with varying problem difficulty and number of agents. Our approach is consistently able to plan successful and safest paths for multiple agents. Note that we extended SoRB by integrating it with CBS to plan waypoints for multiple agents.

For the challenging visual navigation problems, Table 2 demonstrates that our approach again outperforms the baselines by generating safer paths for the agents in both single and multi-agent settings. Notably, we observed a significant drop in success rate for SoRB with an increased number of agents, particularly in the SC2 Staging 08 map. This is likely due to longer re-routing times as agents avoid each other, leading to violation of the time limits. However, as shown in Table 2, our approach consistently maintains safer paths for the agents, with the safety margin improving as the number of agents increases.

Regarding success rate, our approach manages to maintain near-perfect scores, though there is a slight drop in performance when scaling to larger numbers of agents. This behavior is expected, as with more agents, the traversable regions become narrower, making it more challenging for CBS to find collision-free paths. Consequently, agents may need to wait longer at some assigned waypoints, resulting in increased steps to reach their goals and a higher likelihood of exceeding time limits.

Based on these results, we conclude that our approach not only plans safer paths for agents but does so with minimal loss in their ability to reach their goals. It also demonstrates consistent trends when scaling to larger numbers of agents, effectively addressing **Q1**, **Q2**, and **Q3**.

## References

- Roni Stern, Nathan Sturtevant, Ariel Felner, Sven Koenig, Hang Ma, Thayne Walker, Jiaoyang Li, Dor Atzmon, Liron Cohen, TK Kumar, et al. Multi-agent pathfinding: Definitions, variants, and benchmarks. In *Proceedings of the International Symposium on Combinatorial Search*, volume 10, pages 151–158, 2019.
- Jingjin Yu and Steven LaValle. Structure and intractability of optimal multi-robot path planning on graphs. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 27, pages 1443–1449, 2013.
- Piotr Mirowski, Razvan Pascanu, Fabio Viola, Hubert Soyer, Andrew J. Ballard, Andrea Banino, Misha Denil, Ross Goroshin, Laurent Sifre, Koray Kavukcuoglu, Dharshan Kumaran, and Raia Hadsell. Learning to navigate in complex environments, 2017.
- Tom Schaul, Daniel Horgan, Karol Gregor, and David Silver. Universal value function approximators. In Francis Bach and David Blei, editors, *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 1312–1320, Lille, France, 07–09 Jul 2015. PMLR. URL <https://proceedings.mlr.press/v37/schaul15.html>.
- Vitchyr Pong, Shixiang Gu, Murtaza Dalal, and Sergey Levine. Temporal difference models: Model-free deep rl for model-based control, 2020. URL <https://arxiv.org/abs/1802.09081>.
- Andrew Levy, Robert Platt, and Kate Saenko. Hierarchical reinforcement learning with hindsight, 2019.
- Ofir Nachum, Shixiang Gu, Honglak Lee, and Sergey Levine. Data-efficient hierarchical reinforcement learning, 2018.
- Hao-Tien Lewis Chiang, Aleksandra Faust, Marek Fiser, and Anthony Francis. Learning navigation behaviors end-to-end with autorl, 2019.
- Corey Lynch, Mohi Khansari, Ted Xiao, Vikash Kumar, Jonathan Tompson, Sergey Levine, and Pierre Sermanet. Learning latent plans from play, 2019.
- Ashvin Nair, Bob McGrew, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel. Overcoming exploration in reinforcement learning with demonstrations, 2018.
- Eitan Altman. *Constrained Markov decision processes*. Routledge, 2021.
- Ben Eysenbach, Russ R Salakhutdinov, and Sergey Levine. Search on the replay buffer: Bridging planning and reinforcement learning. *Advances in Neural Information Processing Systems*, 32, 2019.
- Alex Ray, Joshua Achiam, and Dario Amodei. Benchmarking Safe Exploration in Deep Reinforcement Learning. 2019.
- Marc G Bellemare, Will Dabney, and Rémi Munos. A distributional perspective on reinforcement learning. In *International conference on machine learning*, pages 449–458. PMLR, 2017.
- Guni Sharon, Roni Stern, Ariel Felner, and Nathan R Sturtevant. Conflict-based search for optimal multi-agent pathfinding. *Artificial Intelligence*, 219:40–66, 2015.

303 Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa,  
304 David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning, 2019.  
305 URL <https://arxiv.org/abs/1509.02971>.

306 John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region  
307 policy optimization. In *International conference on machine learning*, pages 1889–1897. PMLR,  
308 2015.

309 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy  
310 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

311 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.

312 Ryan Lowe, Yi Wu, Aviv Tamar, Jean Harb, Pieter Abbeel, and Igor Mordatch. Multi-agent actor-  
313 critic for mixed cooperative-competitive environments, 2020. URL <https://arxiv.org/abs/1706.02275>.

315 Chao Yu, Akash Velu, Eugene Vinitzky, Jiaxuan Gao, Yu Wang, Alexandre Bayen, and Yi Wu.  
316 The surprising effectiveness of ppo in cooperative, multi-agent games, 2022. URL <https://arxiv.org/abs/2103.01955>.

318 Julian Straub, Thomas Whelan, Lingni Ma, Yufan Chen, Erik Wijnmans, Simon Green, Jakob J. Engel,  
319 Raul Mur-Artal, Carl Ren, Shobhit Verma, Anton Clarkson, Mingfei Yan, Brian Budge, Yajie Yan,  
320 Xiaqing Pan, June Yon, Yuyang Zou, Kimberly Leon, Nigel Carter, Jesus Briales, Tyler Gillingham,  
321 Elias Mueggler, Luis Pesqueira, Manolis Savva, Dhruv Batra, Hauke M. Strasdat, Renzo De Nardi,  
322 Michael Goesele, Steven Lovegrove, and Richard Newcombe. The Replica dataset: A digital  
323 replica of indoor spaces. *arXiv preprint arXiv:1906.05797*, 2019.

324 Andrew Szot, Alex Clegg, Eric Undersander, Erik Wijnmans, Yili Zhao, John Turner, Noah Maestre,  
325 Mustafa Mukadam, Devendra Chaplot, Oleksandr Maksymets, Aaron Gokaslan, Vladimir Vondrus,  
326 Sameer Dharur, Franziska Meier, Wojciech Galuba, Angel Chang, Zsolt Kira, Vladlen Koltun,  
327 Jitendra Malik, Manolis Savva, and Dhruv Batra. Habitat 2.0: Training home assistants to rearrange  
328 their habitat. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.

329 Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijnmans, Bhavana Jain,  
330 Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, Devi Parikh, and Dhruv Batra. Habitat: A  
331 Platform for Embodied AI Research. In *Proceedings of the IEEE/CVF International Conference  
332 on Computer Vision (ICCV)*, 2019.