

Pathfinding in Stochastic Environments: Learning vs Planning

Alexey Skrynnik^{1, 2, 3}, Anton Andreychuk¹, Konstantin Yakovlev^{2,1}, and
Aleksandr Panov^{2,3}

¹AIRI, Moscow, Russia

²Federal Research Center “Computer Science and Control” of the Russian Academy of
Sciences, Moscow, Russia

³Moscow Institute of Physics and Technology, Moscow, Russia

Corresponding author:

Alexey Skrynnik

Email address: skrynnik@airi.net

ABSTRACT

Among the main challenges associated with navigating a mobile robot in complex environments are partial observability and stochasticity. This work proposes a stochastic formulation of the pathfinding problem, assuming that obstacles of arbitrary shapes may appear and disappear at random moments of time. Moreover, we consider the case when the environment is only partially observable for an agent. We study and evaluate two orthogonal approaches to tackle the problem of reaching the goal under such conditions: planning and learning. Within planning, an agent constantly re-plans and updates the path based on the history of the observations using a search-based planner. Within learning, an agent asynchronously learns to optimize a policy function using recurrent neural networks (we propose an original efficient, scalable approach). We carry on an extensive empirical evaluation of both approaches that show that the learning-based approach scales better to the increasing number of the unpredictably appearing/disappearing obstacles. At the same time, the planning-based one is preferable when the environment is close-to-the-deterministic (i.e., external disturbances are rare).

1 INTRODUCTION

Consider a mobile robot operating in a complex, non-stationary environment, e.g. a service robot that has to transfer documents between the offices in an office building. One of the key challenges that arise when such robot navigates from its current location to the target one is that at no moment of time the robot possesses an accurate model of the environment. Among the main reasons for that the following can be named.

First, the apriori map of the environment (e.g. the floor plan) is either unknown or approximate, i.e. it does not contain the crucial information about the furniture, open/closed doors, etc. This leads to a necessity to invoke the so-called simultaneous localization and mapping (SLAM) (Bresson et al., 2017) pipeline that builds and constantly updates the map, based on the data from the sensors installed on the robot (cameras, lidars, etc.). Due to the measurement errors and to the limited sensors’ range, the robot at each timestep acquires only a local semi-accurate patch of the map. Combining these patches to a global map via SLAM typically results in a more inaccurate map, which is subject to constant changes while the robot is progressing toward its goal.

Indeed, planning a path on the basis of such a map (that constantly changes and does not contain accurate information about a large portion of the environment) is a challenging task. Within a planning framework, it is typically addressed via re-planning, i.e. at each timestep, an agent constructs a new plan taking into account the up-to-date map. This can be done in a straightforward fashion (from scratch) using, for example, a suitable variant of A* algorithm (Hart et al., 1968), or via the more involved techniques, like incremental search, that re-use the search efforts of the previous planning attempts (e.g. D*Lite (Koenig and Likhachev, 2002) algorithm, which is widely used in mobile robotics). Overall, such approaches can

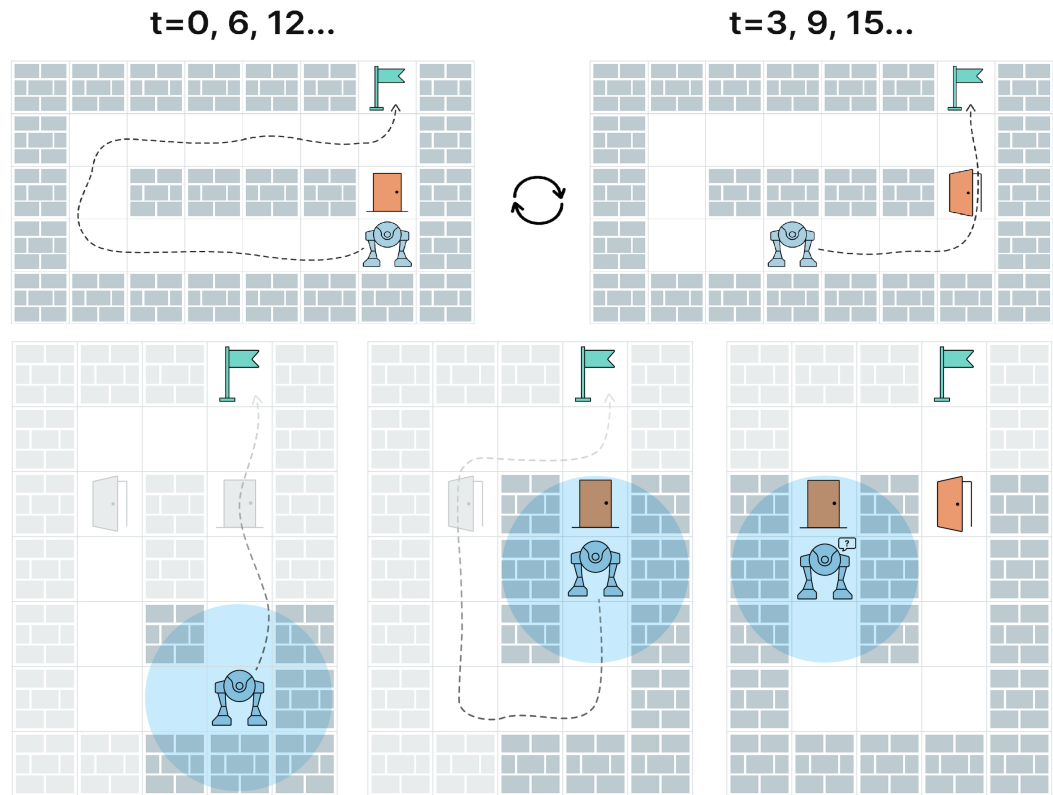


Figure 1. The challenges associated with pathfinding in stochastic environments. Top row: the robot with unlimited field of view is stuck in an oscillating behavior due to blocking/unblocking of a single grid cell (which might correspond to opening/closing the door by humans). Bottom row: due to the partial observability the robot is unaware that the previously blocked cell becomes traversable (the door opens) and, thus, the robot is not able to plan a path to the goal. Both cases can be successfully solved by the learnable policy suggested in the paper, which is empirically shown to learn such behaviors as “wait for an obstacle to disappear,” “keep exploring the environment for additional options of reaching the goal” etc.

46 cope reasonably well with the problems induced by the robot’s internal disturbances like noisy sensor
 47 data and partial observability.

48 There exists, however, a much harder problem, associated with the disturbances that are external to
 49 the robot. The environment (and, thus, its representation as a map) can change due to the actions of
 50 other agents interacting with the environment, e.g. people that close/open doors, move pieces of furniture
 51 from one place to the other etc. These stochastic changes may lead to complete failures of the described
 52 planning/re-planning approach.

53 As an example, consider the environment, depicted at the top of Figure 1. Even if it is fully observable,
 54 the robot can stuck in an oscillating behavior due to the stochastic (from the robot’s perspective) blockage
 55 of one of the grid cells that can correspond to opening/closing a door. When the door is closed the robot
 56 constructs a detouring path and starts moving along it. After three steps, the door opens and a shorter path
 57 is found which is adopted by the robot (Figure 1, top right). However, when the robot comes close to
 58 the door, it might be closed again, and the robot switches back to the detour path. Another challenging
 59 example is depicted on Figure 1 at the bottom row. Due to partial observability and, again, unpredictable
 60 changes in the environment, the robot is stuck being unable to find a path to the goal on a map produced
 61 by accurately combining all observations.

62 We are unaware of the works that study the pathfinding problem when the environment is both partially
 63 observable and stochastic (when the obstacles may appear/disappear unpredictably). Our work aims to fill
 64 this gap. We develop and study empirically two orthogonal approaches to tackle the problem: the one that
 65 follows the planning/re-planning scheme (Ghallab et al., 2016), and the learnable one, when we utilize a
 66 reinforcement learning approach (Sutton and Barto, 2018; Moerland et al., 2020) to optimize a policy

67 function that maps observations to actions. Within the planning approach, we rely on a heuristic search
68 based method that constantly attempts to find a path on the grid. As expected, this approach often fails due
69 to the reasons explained above. To mitigate this issue to a certain extent we develop an extension of this
70 approach that handles the history of observations in a way to heuristically identify stochastic obstacles and
71 avoid them. Moreover, we apply a learning approach to the problem at hand and optimize a policy (deep
72 neural network) that learns to deal with the stochasticity end-to-end. One can think of this as learning to
73 adaptively change the action selection heuristics online (when solving a pathfinding query). For example,
74 the agent may choose to “wait” as many timesteps as needed when a passage gets temporarily blocked.

75 We carry out an extensive experimental evaluation on a wide range of setups and show that the designed
76 learnable policy outperforms the planning approach both in terms of success rate and the solution cost (i.e.
77 number of actions needed to reach the goal) when the number of stochastically appearing/disappearing
78 obstacles is not zero. Thus, we infer that the learning-based approaches should be the tool of choice when
79 solving the studied kind of problems. Despite the latter may sound obvious, we are unaware of the works
80 that empirically confirm this claim.

81 Overall, our main contributions can be summarized as follows. First, we introduce and study a
82 challenging variant of the single-agent pathfinding problem inspired by the real-world robotic applications.
83 In this setting the obstacles might stochastically appear/disappear in the environment and the latter
84 is only partially-observable to the agent. Second, we propose the planning based and the learning
85 based approaches to solve the problem. The latter utilizes reinforcement learning and is able (as shown
86 empirically) to form an adaptive pathfinding policy that successfully handles a wide range environment’s
87 disturbances. Finally, we empirically show that the designed learnable approach outperforms the planning
88 one in the majority of the setups (on different maps, with different numbers of stochastic obstacles etc.).

89 The rest of the paper is organized as follows. Section 2 provides a brief overview of related works.
90 Section 3 focuses on the formulation of the problem of pathfinding in stochastic environments with partial
91 observability. Section 4 describes the methods of planning and asynchronous reinforcement learning that
92 are being compared. Section 5 is devoted to the experimental study of the limits of applicability of the
93 methods under consideration. In the conclusion, the results obtained are discussed.

94 2 RELATED WORK

95 **Context** The previous work in grid-based pathfinding was mainly focused on the application of the
96 planning-based approaches to solving the problem. It was known that the planners based on the heuristic
97 search are the versatile tools when the environment is static and partially observable. These planners
98 have not been examined (both theoretical and empirical) in the setting with both unpredictably appear-
99 ing/disappearing obstacles and partial observability, thus it was not evident whether planning approaches
100 will succeed in it. On the other hand, pure reinforcement learning (RL) techniques (without hierarchy
101 and model) were known to be very powerful in solving a wide range of problems with simple casual
102 structure (like playing the video game of pong), however, when it comes to the problems that require
103 reasoning about the outcomes of the series of actions (like navigation on a grid) RL shows much worse
104 results. E.g. in one of the most cited papers, that addresses the navigation on a fully-observable grid
105 without stochastic obstacles (Panov et al., 2018), the RL methods demonstrated pure convergence of the
106 learning process in even simplest cases. In some papers considering reinforcement learning in a partially
107 observable stochastic environment, tabular methods are used that do not presuppose scaling to large-sized
108 environments (Pena and Banuti, 2021). In the other recent paper from OpenAI (Cobbe et al., 2020),
109 which describes a state-of-the-art RL benchmark including the grid-based navigation problem with the
110 maximum grid size being 25×25 , the advanced RL methods were not able to generalize easily and fast,
111 which means that navigation queries on large grids remained unsolved. Overall, no clear evidence that RL
112 methods, in general, can handle navigation problems on the large partially-observable grids exist so far.

113 **Relevant papers** There exist a lot of different works relevant to the considered problem. Most of them
114 are related to robotics, as the problem of operating in unknown environments with partial observability,
115 uncertainty, and presence of the dynamic obstacles naturally appears in this field of research. (Fiorini and
116 Shiller, 1998) introduces velocity obstacles – one of the major approaches to avoid dynamic obstacles
117 that is based on the idea of predicting their further movement and choosing such an action that does not
118 lead to a collision with any of the observable dynamic obstacles. This idea was further used in many
119 approaches, including such algorithms as ORCA (Van Den Berg et al., 2011) and ALAN (Godoy et al.,

120 2018), developed for multi-agent pathfinding problems.

121 Neural networks and machine learning are also widely used for planning in dynamic environments. In
122 (Chen et al., 2020), (Zhu et al., 2014), biologically-inspired neural networks were applied for planning in
123 unknown dynamic environments. Recent works, such as (Lei et al., 2018), (Wang et al., 2020), have tried
124 to apply reinforcement learning to solve the problem of planning in dynamic environments. Moreover, it
125 is also worth mentioning such kind of planners as ABT (Kurniawati and Yadav, 2016) or DESPOT (Ye
126 et al., 2017) that solve the POMDP-problem, building a belief-state tree to deal with uncertainties and
127 presence of dynamic obstacles.

128 Another field of research related to this work is heuristic search algorithms. Though having a lot
129 of restrictions and assumptions to be applicable, they provide strong theoretical guarantees such as
130 completeness and even optimality. (Koenig and Likhachev, 2002) intrudes state-of-the-art algorithm
131 called D* Lite for planning in unknown partially-observable environments. In (Van Den Berg et al., 2006),
132 there was presented an anytime version of D* algorithm that works not only in unknown environments,
133 but also can deal with dynamic obstacles. Another well-known approach is SIPP (Phillips and Likhachev,
134 2011) that can be applied for planning in the environments with dynamic obstacles and guarantees to
135 find optimal solutions. However, it assumes that the environment is fully observable and trajectories of
136 dynamic obstacles are known.

137 The stochastic shortest path (SSP) is a generalized version of the classical shortest path problem with
138 a presence of stochasticity. In most cases, stochastic behavior is expressed in a non-deterministic result of
139 the actions' execution. Mainly, it is considered as a Markov Decision Process (MDP) and the solution
140 of such kind of problems is a policy that chooses which action to produce in any state to minimize the
141 solution cost. An overview of different variations of this problem is given in (Randour et al., 2015).
142 Though the SSP problem contains stochasticity, it is assumed that the probability distributions are known,
143 which makes it different from the problem that is considered in this paper.

144 It is also worth noting a direction of research where planning algorithms are combined with reinforce-
145 ment learning. In (Skrynnik et al., 2021; Davydov et al., 2021) authors train RL agents in a centralized
146 (QMIX) and decentralized (PPO) way for solving multi-agent pathfinding tasks. The resulting RL policies
147 are combined with a planning approach (MCTS), which leverages the resulting performance. In (Ferber
148 et al., 2020), RL was used to learn the heuristic function to make it more informative. A similar idea of
149 using reinforcement learning to get a better heuristic was suggested in (Micheli and Valentini, 2021) but
150 for the problem of temporal planning.

151 Overall, there exist a large body of works that study the problems which are similar to ours in some
152 aspects. However they all are different in the set of assumptions. As noted above, some of the works
153 assume that the environment is known and fully observable, some of them – that even trajectories of
154 dynamic obstacles are known. Most assume that the trajectories of the dynamic obstacles can be predicted
155 at least for a short period of time. Approaches that deal with stochastic environments assume that
156 probability distributions are given, so they can use them to build an optimal policy. Contrary to all these
157 assumptions, in the problem that we are considering, the environment is changing in an unpredictable
158 way.

159 3 PROBLEM STATEMENT

160 Consider an agent moving on a 4-connected $M \times N$ grid. At each time step of the discrete timeline
161 $T = 0, 1, \dots, T_{max}$, where T_{max} is the duration (length) of the *episode*, a grid cell can be either occupied or
162 free. The cells that are occupied for all time steps are called static obstacles, the cells that are occupied for
163 some time steps while being free for the others correspond to the stochastic obstacles (e.g. closing doors,
164 chairs that are moved by humans etc.). Indeed, the agent can use only the free cells for movement.

165 The action set for the agent is comprised of five actions: $A = \{up, down, left, right, wait\}$. Being at
166 the grid cell c at timestep t , an agent can opt to either *wait* at the current cell or to move to one of the
167 cardinaly-adjacent cells. Let c' denote the target cell of the action, which is either the same cell or one of
168 the neighboring ones. In case c' is free at $t + 1$, the action is considered valid and the agent is transferred
169 to c' . If, however, the destination cell of the move is blocked in the next time step the agent stays put, i.e.
170 is kept in its current cell. Each action $a \in A$ is associated with a non-negative cost: $cost(a) = w \in \mathbb{R}_{>0}$.
171 We assume this cost to be uniform, i.e. 1, for all actions in the rest of the paper. In case the action chosen
172 by the agent turns to be invalid the agent stays put but it still incurs a +1 cost.

173 The grid topology, its size and status of all the grid cells at a certain/any time step is not known to
 174 the agent. Instead, the agent can observe the grid environment only locally. Different ways to model
 175 this local observability model can be suggested. In this work we adopt the most easy-to-implement one:
 176 when located at the grid cell c at time step t , the agent observes a $(2 \cdot R + 1) \times (2 \cdot R + 1)$ patch of the
 177 grid environment, which is centered at c , where R is the given visibility range. In the example depicted
 178 in Figure 2, $R = 5$, which means that the agent observes a 11×11 patch of the grid at each time step¹.
 179 Within the visibility range the agent is able to observe the blockage status of the cells, however it is not
 180 able to distinguish which cells are blocked only temporarily (due to the appearance of the stochastic
 181 obstacles) and which are blocked constantly (due to the static obstacles).

182 We assume that the agent can not predict the blockage status of the grid cells which are both within or
 183 out of the visibility range. However, it is able to memorize the past observations if necessary.

184 The problem now is formulated as follows. Given the start and the goal location (cell) design a
 185 mapping from the (history of the) observations to actions, i.e. the *policy* π , s.t. the chance of reaching the
 186 goal cell within the T_{max} time steps is maximized. In this work we are not restricting ourselves to design
 187 the policy that minimizes the cost of reaching the goal, i.e. the sum of costs of the actions that led to the
 188 goal cell, however obtaining the lower cost paths is, obviously, preferable.



Figure 2. Examples of the grid environments with different numbers of the stochastic obstacles (shown in orange). Undiscovered static obstacles are transparent. The agent’s field of view is shown by the red square.

189 4 METHODS

190 To solve the considered problem we investigate two approaches. The first one relies on finding a path on
 191 the grid and then applying the first action of this path. The second approach is based on reinforcement
 192 learning. Here the agent optimizes a policy that maps the observations to actions end-to-end and follows
 193 this policy in a reactive fashion. Next, we describe both approaches in more details.

194 4.1 Planning

195 The main idea of planning is to repeatedly *i*) construct a full sequence of actions that reach the goal state
 196 (i.e. the grid cell) from the current one (which is the start cell initially), *ii*) apply the first action of the
 197 plan.

198 When constructing a plan, we rely on the history of the received observations, i.e. we memorize the
 199 observations and construct a map out of the map. At each time step, upon receiving the new observation

¹The introduced observation model allows the agent to “see through the obstacles”. Although this assumption is not realistic in the majority of the real-world cases it, indeed, reflects the property of the local observability and at the same time is very easy to implement and experiment with.

200 we update the map. We use it then to find a path from the current cell to the goal one by applying a
 201 pathfinding algorithm. We then extract the first move action from this path and add it to the resultant
 202 plan. If no path is found, we opt to perform a greedy action that moves the agent closer to the goal (more
 203 details on this will be given below). The high-level pseudocode of the planning algorithm is presented in
 204 Algorithm 1.

Input: start cell s , goal cell g , initial observation obs , maximal number of actions the agent can
 make T_{max}
Output: either *success* or *failure*
 $map \leftarrow obs$
 $c \leftarrow s$
 $t \leftarrow 0$
while $t < T_{max}$ **do**
 | **if** $c = g$ **then**
 | | return *success*;
 | **end**
 | $obs \leftarrow GetCurrentObservation()$
 | $map \leftarrow UpdateMap(map, obs)$
 | $path \leftarrow PathFinding(c, g, map)$
 | **if** $path = \emptyset$ **then**
 | | $a \leftarrow GreedyAction()$
 | **else**
 | | $a \leftarrow FirstActionFromPath(path)$
 | **end**
 | $ApplyAction(a)$
 | $t \leftarrow t + 1$
end
 return *failure*

Algorithm 1: High-level algorithm of reaching the goal in a stochastic environment via planning.

Input: number of training epochs E_{max}
Output: policy θ
 $\mathcal{D} \leftarrow \emptyset; e \leftarrow 0; h_0 \leftarrow \emptyset;$
 $\theta \leftarrow InitializeActor()$
 $\phi \leftarrow InitializeCritic()$
while $e < E_{max}$ **do**
 | $\mathcal{D} \leftarrow GenerateTrajectories(); \theta \leftarrow UpdateActor(\mathcal{D}, \phi)$
 | $\phi \leftarrow UpdateCritic(\mathcal{D})$
end
 return θ

Algorithm 2: A policy optimization algorithm (training phase).

205 The presented algorithm relies on the following intrinsic assumptions. First, it is assumed that the
 206 agent knows the goal's coordinates within a reference frame, initially centered at the start position of
 207 the agent, and is able to localize itself within this frame at each timestep. Second, the agent is able to
 208 combine all the observations into a single representation of the environment, i.e. the map. This map is
 209 built/updated incrementally. The size of the map is unknown.

210 The crucial procedure of the algorithm is *PathFinding*, which finds the path on the map provided
 211 with the start and goal locations, where start constantly changes due to agent moving through the
 212 environment. The most prominent way of solving pathfinding problems in the environments with partial
 213 observability is D*Lite algorithm (Koenig and Likhachev, 2002). Instead of re-planning the path from
 214 scratch after applying each action, it extensively reuses the previously built search tree to speed up
 215 the search. Unexpectedly, our preliminary tests have shown, that the performance of D*Lite is worse
 216 compared to the sequential re-planning with A* from scratch after each move. The main reason for such

Input: initial observation obs , maximal number of actions the agent can make K_{max}
Output: trajectory \mathcal{D}
 $\mathcal{D} \leftarrow \emptyset; k \leftarrow 0; h_0 \leftarrow \emptyset$
while $k < K_{max}$ **do**
 if $EpisodeIsDone()$ **then**
 | return \mathcal{D}
 end
 $action, h_k \leftarrow GetAction(obs, h_{k-1})$
 $ApplyAction(a)$
 $obs \leftarrow GetObservation()$
 $r \leftarrow GetReward()$
 $\mathcal{D}.AddTuple(o, a, r)$
 $k \leftarrow k + 1$
end
return \mathcal{D}

Algorithm 3: An algorithm for generating actions with a learnable policy model (inference phase).

217 phenomenon is that at some (numerous) time steps the feasible path from the agent’s current position to
218 the goal is not existent due to the stochastic obstacles that temporarily block the narrow passages. In case
219 these blockages appear close to the agent, running A* from scratch detects unsolvability notably faster
220 compared to D*Lite which actually plans backwards from the goal state. As such blocking happens often
221 (especially when the number of stochastic obstacles is high) sequential invocation of A* actually proved
222 to be beneficial in our preliminary tests and, thus, we adopted A* to implement *PathFinding* in this work.

223 Additionally, we attempted to adapt the *UpdateMap* procedure to the environments with stochastic
224 obstacles. Recall, that the agent is not able to distinguish which grid cells within its visibility range are
225 blocked temporarily, due to the stochastic obstacles, and which cells are blocked for good by the static
226 obstacles. The basic *UpdateMap* procedure treats all the blocked cells as the static obstacles, adding
227 them to the map. This ignores the fact that some of the temporarily blocked cells might become free
228 in future. To this end, we introduce a modified *UpdateMap* procedure that tries to detect stochastic
229 obstacles leveraging the history of observations. More specifically, in case when the agent observes a cell
230 that is currently blocked but was free according to the preceding observations, it is marked as blocked
231 temporarily. When this cell is within the observation radius its actual blockage status is taken into account
232 while finding a path. When the agent moves away and this cell is no longer within the visibility range it is
233 considered to be free, so the path can go through it. Intuitively, this allows the agent to anticipate that the
234 the previously seen stochastic obstacle might move away. This variant of the planning algorithm with the
235 additional indication of the stochastic obstacles is called Stochastic A* (SA*).

236 Both regular planning algorithm (A*) and the adapted one (SA*) share the following additional
237 techniques: exploration threshold and greedy action. Both these techniques are tailored to handle the case
238 when the path to the goal does not exist (due to the temporal presence of the stochastic obstacles). Recall,
239 that we assume that the size of the map is not known to the agent. This infers that it is not technically
240 possible to detect that the current pathfinding query is unsolvable as the search algorithm will keep on
241 exploring the environment assuming that the portion of the environment, that has not been seen/mapped
242 before, is traversable. To mitigate this issue we impose the threshold on the number of internal iterations
243 of A*/SA*. When this threshold is exceeded the search is aborted with the *no-path-found* result. In
244 such case instead of picking the random action we choose the one that is likely to move the agent closer
245 to the goal – the so-called greedy action. It is chosen as follows. When A*/SA* terminates due to the
246 exploration threshold without finding the path we pick the node in the search tree that corresponds to the
247 cell which is the closest to the goal. We then reconstruct the path to this cell in the tree and pick the first
248 action of the path.

249 4.2 Learning

250 The main idea of the reinforcement learning (RL) approach is to optimize a policy π , which maps the
251 observation to an action. The policy is trained to maximize the cumulative expected reward (cost function)
252 for each interaction episode. We use the partially observable setting since the agent has no access to the
253 global state.

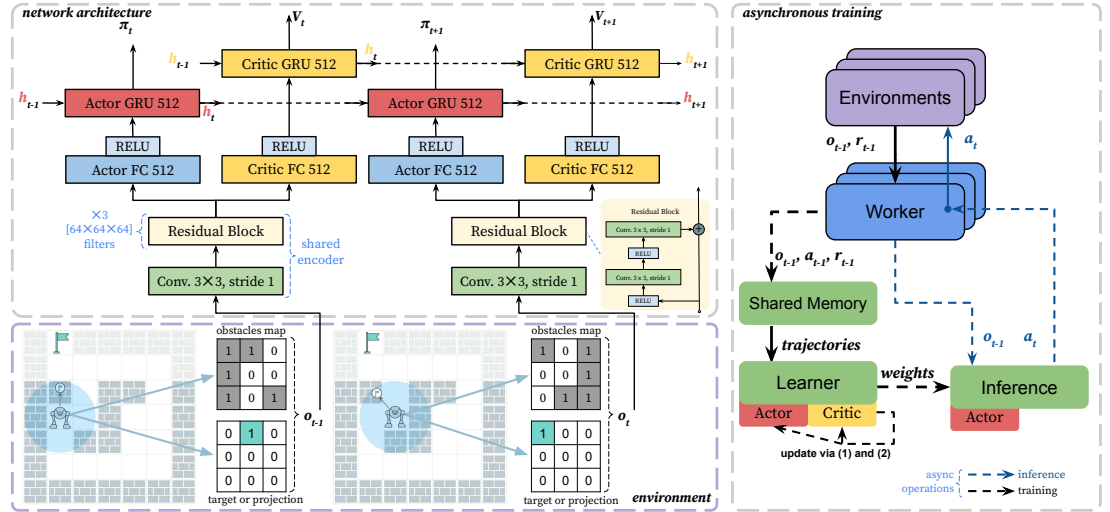


Figure 3. The scheme of the proposed asynchronous learning approach. The bottom-left part shows the environment and observation encoding. The top-left part of the scheme shows the neural network architecture. We use residual layers as a shared encoder and recurrent heads. The recurrent layers (GRU) are responsible for remembering obstacles and environment dynamics. The right part shows the asynchronous learning procedure. The green blocks show GPU computations.

254 The advantage of RL over planning approaches is that the agent can learn the adaptive heuristic of
 255 acting in an environment with stochastic dynamics. At the same time, it should be taken into account
 256 that the work of the learnable approaches is divided into two phases: the actual training on prepared
 257 environment configurations (training phase, see Algorithm 2) and the work of the trained model on any
 258 environments (inference phase, see Algorithm 3).

259 Formally, in RL setting, the interaction of an agent with the environment is described as partially
 260 observable Markov Decision Process (POMDP), which can be described as tuple (S, O, A, P, r, γ) , where
 261 S is the set of environment states, $o \in O$ is a partial observation of the state, $a \in A$ is the set of agent's
 262 actions, $r(s, a) : S \times A \rightarrow \mathbb{R}$ is a reward (cost) function, and γ is the discount factor. The agent has no
 263 access to states S (true coordinates and full obstacle map), and the policy is a mapping from observations
 264 to actions: $\pi(a|o) : A \times O \rightarrow [0, 1]$.

265 We propose and describe an end-to-end architecture to train the agent in grid pathfinding scenarios.
 266 And we believe that our learning approach is applicable for a wide range of pathfinding tasks. The scheme
 267 of the learning approach is presented in Figure 3.

268 As already mentioned in section 3 in the proposed environment the observation space O of the agent is
 269 a multidimensional matrix: $O : 2 \times (2 \times R + 1) \times (2 \times R + 1)$, that represents the part of the environment
 270 around the agent within radius R . It includes the following two matrices.

- 271 • *Obstacle matrix*: 1 encodes an obstacle, and 0 encodes its absence. If any cell of the agent's field of
 272 view, which is outside the grid, is encoded as an obstacle. The agent does not distinguish between
 273 the type of obstacles. Both static and stochastic obstacles are encoded the same.
- 274 • *Target matrix*: if the agent's goal is inside the observation field, then there is 1 in the cell, where it
 275 is located, and 0 in other cells. If the target does not fall into the view, then it is projected onto the
 276 nearest cell of the observation field Skrynnik et al. (2021).

277 At any time step, the agent has five actions available: stay in place, move vertically (up or down), or
 278 move horizontally (right or left). The agent can move to any adjacent free cell.

279 The agent receives a reward of 1.0 when it reaches the goal and 0.0 in all other cases. We have chosen
 280 this function so one can train the agent, which can deal with a wide range of tasks in partially observable
 281 grid environments with stochasticity.

282 To learn a policy in a model-free setting, there are a number of well-known methods in reinforcement
 283 learning, which can be roughly divided into two classes – value-based (Mnih et al., 2015) and policy-

284 based (Schulman et al., 2015; Haarnoja et al., 2018; Lillicrap et al., 2016) methods. The first group of
285 approaches is characterized by the use of replay buffer to store experience and can learn only deterministic
286 policies. The second group of methods is specifically designed for operating stochastic policies. In the
287 problem we are considering, which is characterized by stochastic behavior of the environment itself, it
288 is necessary to provide an opportunity to work with probability distributions on a set of agent actions.
289 On-policy methods (Schulman et al., 2017) that use only current experience are the most promising for
290 the partially observed formulation of the pathfinding problem. This is due to the fact that in some cases
291 of recurrent stochasticity in the environment, it is necessary to use a stochastic policy, which is most
292 effectively learned by on-policy policy gradient methods. Also this makes it possible to improve the
293 quality of state prediction by observation when using recurrent neural network models.

294 We optimize the policy π_θ , which is approximated by a neural network θ , using Proximal Policy
295 Optimization (PPO) method (Schulman et al., 2017). The approach is a variant of the actor-critic algorithm,
296 and proven effective in many challenging domains (Berner et al., 2019; Yu et al., 2021; Cobbe et al.,
297 2020). To adapt PPO for the POMDP setting, we approximate the state s_t using a recurrent neural network
298 (RNN): $s_t \approx f(h_t, o_t)$, where h_t is a hidden state of RNN. PPO uses clipping in objective to improve
299 performance monotonically. The clipped objective penalize the new policy $\pi_{\theta_{k+1}}$ for getting far from the
300 previous one π_{θ_k} .

301 To approximate policy and value we adapt network architecture from IMPALA (Espeholt et al., 2018).
302 As a feature encoder we use residual layers. We have changed the original architecture and removed max
303 pooling, similar to AlphaZero architecture (Silver et al., 2017), which used akin encoding for observations.
304 Removing max pooling is crucial, to prevent losing spatial information of the grid observation. After the
305 shared feature encoder, there are recurrent layers separate for the actor and the critic.

306 We use single GPU asynchronous training setup based on SampleFactory (Petrenko et al., 2020). We
307 called the modified version of the PPO algorithm as asynchronous proximal policy optimization (APPO).
308 The asynchronous training can be divided into two main parts, which run in parallel: accumulating new
309 trajectories and policy updating. The policy used to collect experience may lag behind the current one.
310 That discrepancy is called a policy lag, and it negatively affects the performance. The policy lag especially
311 affects the learning of recurrent architectures. To stabilize training for such case IMPALA (Espeholt et al.,
312 2018) introduced importance sampling for value targets, called V-trace, which we also use in our training.

313 5 EXPERIMENTAL EVALUATION

314 5.1 Environment

315 We designed and implemented the environment simulator that takes all the specifics of the partially
316 observable pathfinding and stochastic obstacles into account. The following parameters are related to
317 the interaction between the agent and the environment: observation radius (agents observe $1 \leq R \leq Size$
318 cells in each direction) and the maximum number of steps in the environment before the episode ends
319 $Horizon \geq 1$.

320 The stochastic part of the environment is described with the following parameters: a number of the
321 stochastic obstacles ≥ 0 ; obstacle size $[x, y]$ (we assume that each stochastic obstacle is a square whose
322 size is in between $[x, y]$); obstacle density $\in (0, 1]$ (the chance of each cell comprising the stochastic
323 obstacle to be blocked, i.e. some of the cells forming the obstacle can be free); obstacle move radius ≥ 0
324 (defines how far from the original placement the obstacle can move); obstacle appear time range $[x, y]$
325 (the range of steps during which a stochastic obstacle can be active); obstacle disappear time range $[x, y]$
326 (the range of steps during which a stochastic obstacle can be inactive). Figure 2 shows examples of the
327 environments with the defined stochastic obstacles.

328 For the experimental evaluation, we have set the following environment parameters of stochastic
329 obstacles: obstacle size = $[5, 10]$, obstacle density = 0.7, obstacle move radius = 5, obstacle appear time
330 range = $[8, 16]$, obstacle disappear time range = $[8, 16]$. All the values for the parameters with ranges are
331 chosen randomly out of the corresponding range. In other words each stochastic obstacle is a square with
332 a side of 5 to 10, and 70% of the cells forming this square are blocked. The structure of each stochastic
333 obstacle does not change with time. However, every tick of time it can move, but not more than 5 cells
334 from the initial position. A stochastic obstacle is present on the map from 8 to 16 time ticks, after which it
335 disappears for a period of 8 to 16 time ticks and then reappears.

336 It is also worth to note that stochastic obstacles can be superimposed on static obstacles, on top of
337 each other, but cannot block the cell in which the agent is located.

338 Moreover, the observation radius was set to 5 and *Horizon* to 512.

339 5.2 Setup

340 During the experiments, we evaluated both planning approaches – the basic one, i.e. A^* , and the improved
341 one, i.e. SA^* . There is only one crucial parameter that is needed to be set for them – a maximum number
342 of iterations, that was set to 10000. Besides the planning approaches, we have also evaluated the learning
343 one, i.e. APPO. The details about its training process and the choice of hyperparameters are described in
344 section 5.3.

345 The experiments were conducted on the maps from three different collections taken from MovingAI
346 (Sturtevant, 2012) – a grid-based pathfinding benchmark: a) *wc3* – 36 maps from Warcraft 3 computer
347 game; b) *sc1* – 75 maps from Starcraft computer game; c) *streets* – 30 maps with real city data taken from
348 OpenStreetMap. The chosen maps represent different landscapes with varying topological structure, i.e.
349 they include maps with small passages, large open areas, prolonged obstacles of non-trivial shapes etc.
350 The original maps can be in size up to 1024×1024 . For our experiments we have scaled them to 64×64
351 or such a size that the lowest side is 64. These collections were divided for training and test subsets in a
352 ratio of 8 : 2, simply using the alphabetic names of the maps. Examples of the maps are shown in Figure 4.
353 The names of all the maps that were used for tests can be found in section 5.4.

354 For each of the evaluated maps, there were generated 200 instances. The instances were generated
355 randomly, but in such a way that the path between start and goal locations is guaranteed to exist for the
356 static obstacles. Moreover, each of these instances was evaluated with a different number of stochastic
357 obstacles: from 0 to 200 with an increment of 25.

358 During the experiments, we evaluated such a parameter as success rate, which shows the ratio between
359 the number of successfully solved instances to their total amount. An instance is considered successfully
360 solved in case the solution was found within less than 512 steps (*Horizon* parameter), i.e. the resulting
361 plan contains less than 512 actions. Besides the success rate, we have also evaluated the average episode
362 length, taking *Horizon* value for the instances that were not solved by an algorithm, and computed how
363 many times each of the algorithms has found a solution with less number of actions.

364 To evaluate the computational efficiency of the approaches, we have measured such a parameter
365 as “steps per second” (SPS), which shows how fast the approach returns an action that the agent needs
366 to perform on the current step. The more actions the approach returns within a second, the faster it
367 works. The evaluation of all algorithms were conducted using AMD Ryzen Threadripper 3970X CPU
368 (single-core). Also, the performance of RL approach, i.e. APPO, were tested using NVIDIA GeForce
369 RTX 3080 Ti.



Figure 4. Examples of the evaluated maps from different collections.

370 5.3 APPO Training

371 First, we make a hyperparameter search to adjust the number of residual blocks and the number of filters
372 in them (see Figure 5). As the environment configuration, we use a map with 30% density of static
373 obstacles and 64 stochastic obstacles in 64×64 grid, obstacle size is $\{5, 10\}$, the density of the stochastic
374 obstacles is 0.7, appear and disappear time ranges are $\{8, 16\}$. The agent was trained for 100 million
375 steps. The best results were shown by a network with 3 residual layers and 60 filters in them, thus we use
376 these settings in all proceedings experiments.

377 Second, we train APPO for one billion steps using training collection. For each episode, the initial
 378 position of the agent and his target, as well as the configuration of stochastic obstacles were sampled
 379 randomly. Solving a task with a large number of stochastic obstacles is difficult in terms of exploration.
 380 Thus, we used curriculum learning to automatically adjust the difficulty of the environment. In each
 381 curriculum phase, the agent is trained until it reaches the average success rate of 0.9 on the 256 consequent
 382 episodes, after that the number of stochastic obstacles is increased by one. The final number of stochastic
 383 obstacles at the end of the training was 46.

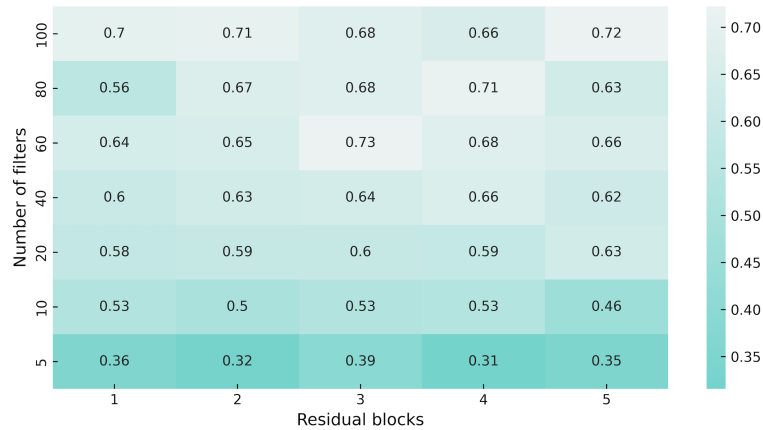


Figure 5. Heatmap of the network architecture parameters. We use version with 60 filters and 3 residual blocks as encoder in APPO model for all experiments, since it showed better performance in our hyperparameter sweep.

384 5.4 Results

385 The aggregated results of the experiments are shown in Table 1. As one can see in most of the cases,
 386 APPO solves more instances than planning approaches. The only case where APPO shows a worse
 387 success rate is the one that doesn't contain any stochastic obstacles, where planning approaches have
 388 solved 100% of instances. The success rate of A* algorithm, which doesn't detect stochastic obstacles and
 389 remembers all the obstacles it has seen, is very poor. While APPO and SA* successfully solved 45% of
 390 the instances with 200 stochastic obstacles, A* has shown even worse success rate with only 25 stochastic
 391 obstacles. Such behavior explains by the fact that A* in most cases can't pass stochastic obstacles and
 392 fails to find a path to the goal. This behavior also explains its high SPS – A* makes a very few expansions
 393 before it makes a conclusion that the path cannot be found.

394 A more detailed view of the success rates of the approaches is presented in Figure 6. It shows the
 395 success rates of the approaches on each of the maps separately. As one can see, there is not a single map
 396 where SA* significantly outperforms APPO. Despite the points with 0 stochastic obstacles, in all the
 397 cases they either show very close results, or APPO outperforms SA*. There are some maps, for example,
 398 *timbermawhold* or *swampofsorrows*, that show significantly higher success rates compared to
 399 other maps. Such behavior explains by the fact that these maps contain several relatively small disjoint
 400 areas. As a result, the instances on these maps are much simpler than on the other maps, as the distance
 401 between start and goal location is much lower. On the other hand, there are some maps, where success
 402 rates of all approaches drop to less than 20% on the instances with the highest amount of stochastic
 403 obstacles. There are actually two reasons for such behavior. First, there are maps such as *ThinIce* (see
 404 Figure 4(a)) or *Typhoon*, that contains a difficult structure for partially observable environments, i.e.
 405 they contain “trap” areas, that are on the way to the goal but do not actually lead to it. Second, some maps
 406 of *wc3* collection in the original MovingAI benchmark contain huge borders of obstacles, that affected
 407 the scaled maps and reduced their actual size with traversable areas to 48x48 (see Figure 4(c)). Thus, the
 408 density of stochastic obstacles on such maps is actually much higher than on other maps.

409 To get some insight about the quality of the found solutions, we have computed how many times
 410 APPO or SA* has found a solution with less number of actions. We have excluded the results of A* in
 411 this comparison, as its success rate is very low. However, there actually were some instances where A*

412 has found a solution with the least number of actions – 220 out of 56,000. The results were aggregated
413 among the collections and are presented in Figure 7. As one can see, there actually presents two more
414 lines called “Equal” and “Failed”. “Equal” line indicates the portion of instances that were successfully
415 solved by both approaches with an equal number of actions, while “Failed” one indicates the portion of
416 instances that were solved neither by APPO nor by SA*. The results on *scl* and *street* collections are
417 very similar. The portion of instances that were solved with equal number of steps on *scl* and *street*
418 collections is relatively small, while on *wc3* it’s much higher. This behavior on *wc3* is explained by the
419 presence of such maps as *timbermawhold* and *swampofsorrows* with isolated parts, where the
420 instances are easy and thus solved by both approaches with equal cost. About 70% of all the instances on
421 these maps were successfully solved by both approaches with an equal number of actions. The behavior
422 of SA* and APPO methods shows that the planning approach outperforms the others when the number
423 of stochastic obstacles is very small. The learning approach outperforms the others when the number of
424 stochastic obstacles rises to about 100. When the number of stochastic obstacles rises to the maximum,
425 both approaches show close results. Generally, these plots indicate the same trends as the ones that show
426 the success rates of the algorithms.

Algorithm	Obstacles	SPS (CPU)	SPS (GPU)	Episode length	Success rate
A*		1127.87	-	57	1
APPO	0	105.58	812.3	95.64	0.93
SA*		1052.73	-	57	1
A*		2226.9	-	331.14	0.38
APPO	25	100.71	754.58	120.11	0.93
SA*		847.81	-	135.39	0.88
A*		2589.9	-	410.22	0.21
APPO	50	100.39	731.58	157.29	0.9
SA*		866.21	-	189.96	0.8
A*		2438.91	-	438.91	0.15
APPO	75	98.23	712.26	200.41	0.82
SA*		909.81	-	232.68	0.72
A*		2226.41	-	453.42	0.12
APPO	100	100.04	702.45	242.72	0.74
SA*		955.12	-	271.56	0.65
A*		2061.56	-	462.88	0.1
APPO	125	98.49	693.69	279.28	0.65
SA*		976.84	-	297.88	0.59
A*		1911.56	-	468.22	0.09
APPO	150	98.23	686.08	309.25	0.57
SA*		988.06	-	320.56	0.54
A*		1793.01	-	474.08	0.08
APPO	175	97.75	667.67	334.54	0.51
SA*		1021.84	-	338.94	0.49
A*		1631.31	-	475.79	0.07
APPO	200	98.39	659.66	353.37	0.45
SA*		1014.4	-	354.45	0.45

Table 1. Averaged results of A*, SA*, and APPO aggregated over all the evaluated maps. Bold values highlight better performance (episode length – lower better, success rate – higher better) for each number of stochastic obstacles. SPS denotes steps per second.

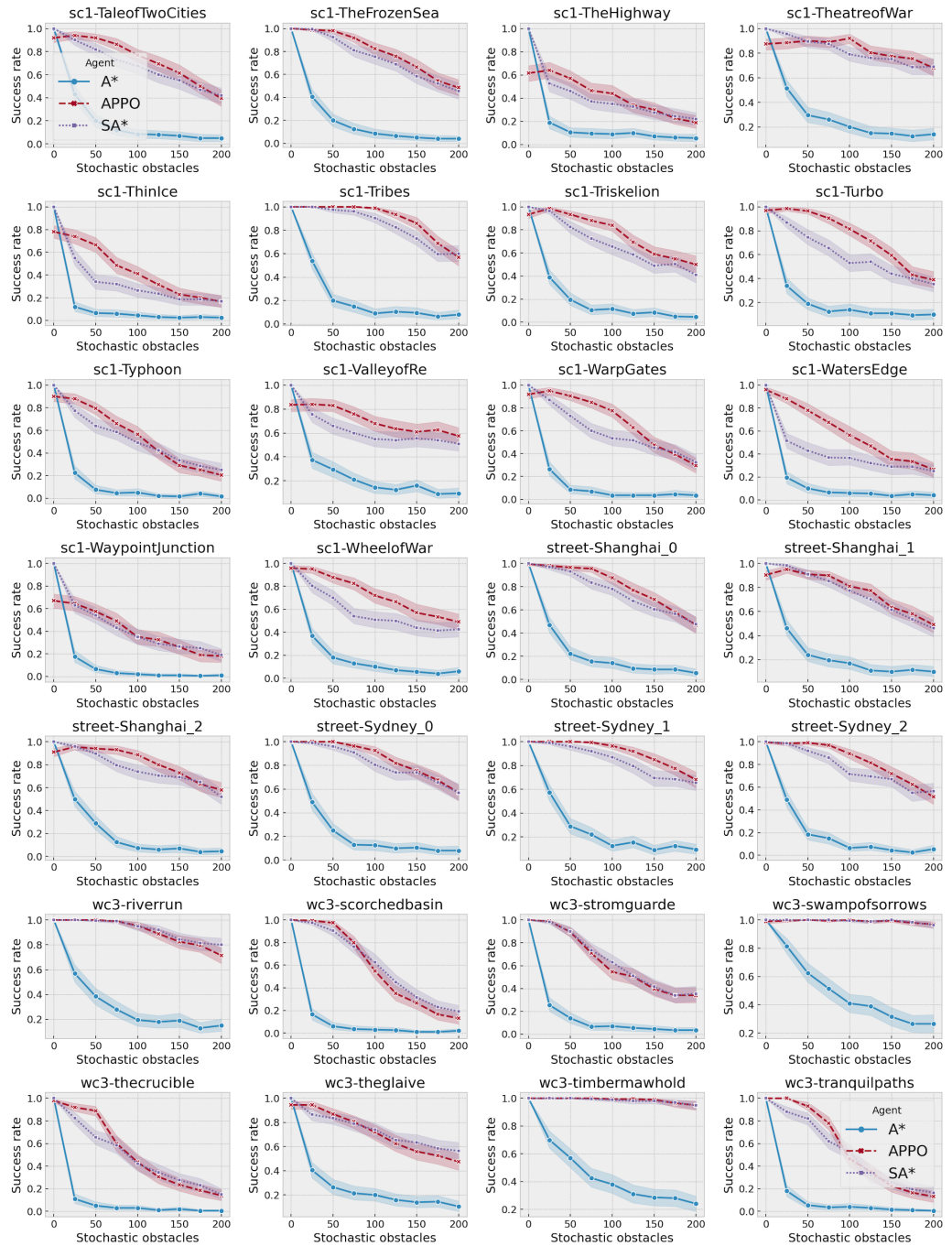


Figure 6. Success rates of A*, SA* and APPO on each of the evaluated maps depending on the number of stochastic obstacles. The shaded area shows 95% confidence interval.

427 6 CONCLUSION AND DISCUSSION

428 In this paper, we have introduced and studied a challenging variant of the single-agent pathfinding
 429 problem inspired by real-world robotic applications. In this setting, some of the obstacles unpredictably
 430 appear/disappear in the environment, and the latter is only partially observable to the agent. We designed
 431 two orthogonal approaches to solve this problem: planning-based and learning-based. For the former,
 432 we utilized the well-known A* algorithm and suggested its modification, called Stochastic A* (SA*),
 433 that differs in the way how the incoming observations are processed; for the latter, we have proposed an

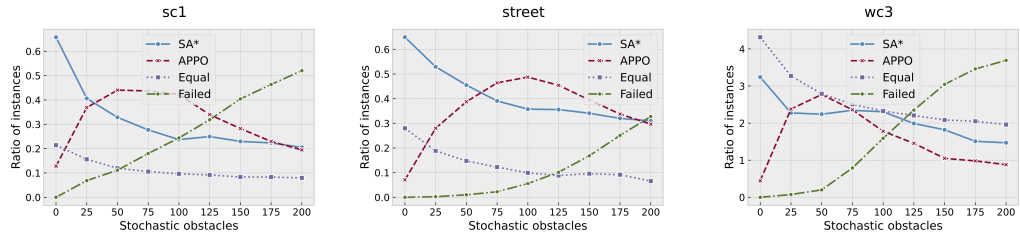


Figure 7. Comparison of the solutions found by SA* and APPO depending on the collection of the maps and number of stochastic obstacles. Points on “SA*” and “APPO” lines indicates the portions of instances for which the corresponding approach has found a solution of a better quality than other one. “Equal” correspond to the portion of instances for which both approaches have found solution with equal quality. “Failed” correspond to the portion of instances that were solved neither by SA* nor by APPO.

434 original asynchronous policy optimization method (APPO) based on the established actor-critic neural
 435 network architecture. Both approaches were experimentally evaluated on a range of setups involving
 436 different maps and degrees of stochasticity (i.e. numbers of the appearing/disappearing obstacles). The
 437 results indicate that both of the suggested approaches has their own pros and contras. SA* is evidently
 438 faster than APPO but its success rate is generally lower compared to APPO. The only case when SA*
 439 performs better/on par with APPO is either when no stochastic obstacles are present at all or when
 440 this number is very high. Both cases can be seen as the outliers. For all other configuration APPO,
 441 indeed, is able to successfully solve more instances than SA*. We believe that this happens due to
 442 the ability of APPO to adaptively adjust the heuristic of choosing actions, which is learned rather than
 443 hard-coded. We also would like to note that low computational efficiency of our implementation of APPO
 444 is not a fundamental problem, as there is a room for a significant speed-up via using specialized code
 445 implementations (e.g. TensorRT).

446 One of the perspective avenues for future research is investigating the analogous problem statements,
 447 but when certain predictions about the dynamics of the environment can be made. One of such settings
 448 that is of particular interest is the decentralized multi-agent pathfinding setting, when each agent can
 449 distinguish between the static obstacles and the moving agents and is able to predict, to a certain extent,
 450 the future moves of the latter. We assume that in such settings, the planning-based approaches may
 451 exhibit a better performance due to the additional knowledge that they can take into account, i.e. the
 452 locations that will be blocked at the next time step due to the moves of the other agents. In such case,
 453 pathfinding algorithms can be straightforwardly extended to reason about the temporal dimension and to
 454 build plans that will avoid future collisions with the other agents. As for the learning-based approaches,
 455 modifying (and learning) them for such settings might be more problematic. Indeed, such approaches for
 456 decentralized multi-agent pathfinding do exist currently, see (Sartoretti et al., 2019; Riviere et al., 2020) and
 457 others, but mainly they rely on the accurate long-horizon predictions of how the other agents will behave,
 458 i.e. they rely on the ability to acquire/accurately reconstruct the full paths of the other agents to their goals.
 459 In case this ability is limited, e.g. only the next action can be inaccurately predicted, their performance
 460 might get worse. So the question whether the learning-based approaches will beat the planning-based
 461 ones in such settings is still to be answered.

462 REFERENCES

- 463 Berner, C., Brockman, G., Chan, B., Cheung, V., Debiak, P., Dennison, C., Farhi, D., Fischer, Q., Hashme,
 464 S., Hesse, C., et al. (2019). Dota 2 with large scale deep reinforcement learning. *arXiv preprint*
 465 *arXiv:1912.06680*.
- 466 Bresson, G., Alsayed, Z., Yu, L., and Glaser, S. (2017). Simultaneous localization and mapping: A survey
 467 of current trends in autonomous driving. *IEEE Transactions on Intelligent Vehicles*, 2(3):194–220.
- 468 Chen, Y., Liang, J., Wang, Y., Pan, Q., Tan, J., and Mao, J. (2020). Autonomous mobile robot path planning
 469 in unknown dynamic environments using neural dynamics. *Soft Computing*, 24(18):13979–13995.

- 470 Cobbe, K., Hesse, C., Hilton, J., and Schulman, J. (2020). Leveraging procedural generation to benchmark
471 reinforcement learning. In *International conference on machine learning*, pages 2048–2056. PMLR.
- 472 Davydov, V., Skrynnik, A., Yakovlev, K., and Panov, A. (2021). Q-Mixing Network for Multi-agent
473 Pathfinding in Partially Observable Grid Environments. In Kovalev, S. M., Kuznetsov, S. O., and Panov,
474 A. I., editors, *Artificial Intelligence. RCAI 2021. Lecture Notes in Computer Science*, volume 12948,
475 pages 169–179. Springer.
- 476 Espeholt, L., Soyer, H., Munos, R., Simonyan, K., Mnih, V., Ward, T., Doron, Y., Firoiu, V., Harley, T.,
477 Dunning, I., et al. (2018). Impala: Scalable distributed deep-rl with importance weighted actor-learner
478 architectures. In *International Conference on Machine Learning*, pages 1407–1416. PMLR.
- 479 Ferber, P., Helmert, M., and Hoffmann, J. (2020). Reinforcement learning for planning heuristics.
480 In *Proceedings of the 1st Workshop on Bridging the Gap Between AI Planning and Reinforcement*
481 *Learning (PRL)*, pages 119–126. University_of_Basel.
- 482 Fiorini, P. and Shiller, Z. (1998). Motion planning in dynamic environments using velocity obstacles. *The*
483 *International Journal of Robotics Research*, 17(7):760–772.
- 484 Ghallab, M., Nau, D., and Traverso, P. (2016). *Automated planning and acting*. Cambridge University
485 Press.
- 486 Godoy, J., Chen, T., Guy, S. J., Karamouzas, I., and Gini, M. (2018). Alan: adaptive learning for
487 multi-agent navigation. *Autonomous Robots*, 42(8):1543–1562.
- 488 Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. (2018). Soft actor-critic: Off-policy maximum entropy
489 deep reinforcement learning with a stochastic actor. In *International conference on machine learning*,
490 pages 1861–1870. PMLR.
- 491 Hart, P. E., Nilsson, N. J., and Raphael, B. (1968). A formal basis for the heuristic determination of
492 minimum cost paths. *IEEE transactions on Systems Science and Cybernetics*, 4(2):100–107.
- 493 Koenig, S. and Likhachev, M. (2002). D* lite. In *Proceedings of the 18th AAI Conference on Artificial*
494 *Intelligence (AAAI 2002)*, pages 476–483.
- 495 Kurniawati, H. and Yadav, V. (2016). An online pomdp solver for uncertainty planning in dynamic
496 environment. In *Robotics Research*, pages 611–629. Springer.
- 497 Lei, X., Zhang, Z., and Dong, P. (2018). Dynamic path planning of unknown environment based on deep
498 reinforcement learning. *Journal of Robotics*, 2018.
- 499 Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., and Wierstra, D. (2016).
500 Continuous control with deep reinforcement learning. In Bengio, Y. and LeCun, Y., editors, *4th*
501 *International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4,*
502 *2016, Conference Track Proceedings*.
- 503 Micheli, A. and Valentini, A. (2021). Synthesis of search heuristics for temporal planning via reinforce-
504 ment learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages
505 11895–11902.
- 506 Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. a., Veness, J., Bellemare, M. G., Graves, A., Riedmiller,
507 M., Fidjeland, A. K., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H.,
508 Kumaran, D., Wierstra, D., Legg, S., and Hassabis, D. (2015). Human-level control through deep
509 reinforcement learning. *Nature*, 518(7540):529–533.
- 510 Moerland, T. M., Broekens, J., and Jonker, C. M. (2020). A framework for reinforcement learning and
511 planning. *ArXiv*.
- 512 Panov, A. I., Yakovlev, K. S., and Suvorov, R. (2018). Grid Path Planning with Deep Reinforcement
513 Learning: Preliminary Results. In Klimov, V. and Samsonovich, A., editors, *Procedia Computer*
514 *Science*, volume 123, pages 347–353. Elsevier.

- 515 Pena, B. D. and Banuti, D. T. (2021). Reinforcement learning for pathfinding with restricted observation
516 space in variable complexity environments. In *AIAA Scitech 2021 Forum*, page 1755.
- 517 Petrenko, A., Huang, Z., Kumar, T., Sukhatme, G., and Koltun, V. (2020). Sample factory: Egocentric
518 3d control from pixels at 100000 fps with asynchronous reinforcement learning. In *International
519 Conference on Machine Learning*, pages 7652–7662. PMLR.
- 520 Phillips, M. and Likhachev, M. (2011). Sipp: Safe interval path planning for dynamic environments. In
521 *2011 IEEE International Conference on Robotics and Automation*, pages 5628–5635. IEEE.
- 522 Randour, M., Raskin, J.-F., and Sankur, O. (2015). Variations on the stochastic shortest path problem.
523 In *International Workshop on Verification, Model Checking, and Abstract Interpretation*, pages 1–18.
524 Springer.
- 525 Riviere, B., Hönig, W., Yue, Y., and Chung, S.-J. (2020). Glas: Global-to-local safe autonomy synthesis
526 for multi-robot motion planning with end-to-end learning. *IEEE Robotics and Automation Letters*,
527 5(3):4249–4256.
- 528 Sartoretti, G., Kerr, J., Shi, Y., Wagner, G., Kumar, T. S., Koenig, S., and Choset, H. (2019). Primal:
529 Pathfinding via reinforcement and imitation multi-agent learning. *IEEE Robotics and Automation
530 Letters*, 4(3):2378–2385.
- 531 Schulman, J., Levine, S., Abbeel, P., Jordan, M., and Moritz, P. (2015). Trust region policy optimization.
532 In *International conference on machine learning*, pages 1889–1897. PMLR.
- 533 Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. (2017). Proximal policy optimization
534 algorithms. *arXiv preprint arXiv:1707.06347*.
- 535 Silver, D., Hubert, T., Schrittwieser, J., Antonoglou, I., Lai, M., Guez, A., Lanctot, M., Sifre, L., Kumaran,
536 D., Graepel, T., et al. (2017). Mastering chess and shogi by self-play with a general reinforcement
537 learning algorithm. *arXiv preprint arXiv:1712.01815*.
- 538 Skrynnik, A., Yakovleva, A., Davydov, V., Yakovlev, K., and Panov, A. I. (2021). Hybrid Policy Learning
539 for Multi-Agent Pathfinding. *IEEE Access*, 9.
- 540 Sturtevant, N. R. (2012). Benchmarks for grid-based pathfinding. *IEEE Transactions on Computational
541 Intelligence and AI in Games*, 4(2):144–148.
- 542 Sutton, R. S. and Barto, A. G. (2018). *Reinforcement learning: An introduction*. Bradford Books, 2nd
543 edition.
- 544 Van Den Berg, J., Ferguson, D., and Kuffner, J. (2006). Anytime path planning and replanning in dynamic
545 environments. In *Proceedings 2006 IEEE International Conference on Robotics and Automation, 2006.
546 ICRA 2006.*, pages 2366–2371. IEEE.
- 547 Van Den Berg, J., Guy, S. J., Lin, M., and Manocha, D. (2011). Reciprocal n-body collision avoidance. In
548 *Robotics research*, pages 3–19. Springer.
- 549 Wang, B., Liu, Z., Li, Q., and Prorok, A. (2020). Mobile robot path planning in dynamic environments
550 through globally guided reinforcement learning. *IEEE Robotics and Automation Letters*, 5(4):6932–
551 6939.
- 552 Ye, N., Somani, A., Hsu, D., and Lee, W. S. (2017). Despot: Online pomdp planning with regularization.
553 *Journal of Artificial Intelligence Research*, 58:231–266.
- 554 Yu, C., Velu, A., Vinitzky, E., Wang, Y., Bayen, A., and Wu, Y. (2021). The surprising effectiveness of
555 ppo in cooperative multi-agent games.
- 556 Zhu, D., Li, W., Yan, M., and Yang, S. X. (2014). The path planning of auv based on ds information
557 fusion map building and bio-inspired neural network in unknown dynamic environment. *International
558 Journal of Advanced Robotic Systems*, 11(3):34.