
Learn to Follow: Lifelong Multi-agent Pathfinding with Decentralized Replanning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Multi-agent Pathfinding (MAPF) problem generally asks to find a set of conflict-
2 free paths for a set of agents confined to a graph. In conventional MAPF scenarios,
3 the graph and the agents' start and goal locations are known in advance. Thus,
4 a centralized planning algorithm can be utilized to generate a solution. In this
5 work, we investigate the decentralized MAPF setting, in which the agents can
6 not share the information and must independently navigate toward their goals
7 without knowing the other agents' goals or paths. We focus on the lifelong variant
8 of MAPF, which involves continuously assigning new goals to the agents upon
9 arrival to the previous ones. To address this complex problem, we propose a
10 method that integrates two complementary approaches: planning with heuristic
11 search and reinforcement learning (RL) through policy optimization. Planning
12 is utilized to maintain an individual path, while RL is employed to discover the
13 collision avoidance policies that effectively guide an agent along the path. This
14 decomposition and intrinsic motivation specific for multi-agent scenarios allows
15 leveraging replanning with learnable policies. We evaluate our method on a wide
16 range of setups and compare it to the state-of-the-art competitors (both learnable
17 and search-based). The results show that our method consistently outperforms the
18 competitors in challenging setups when the number of agents is high.

19 1 Introduction

20 Multi-agent pathfinding (MAPF) [1] is a challenging problem that gets increasing attention recently.
21 It is often studied in the AI community with the following assumptions. The agents are confined to a
22 graph, and at each time step an agent can either move to an adjacent vertex or stay at the current one.
23 A central controller possesses information about the graph and the agents' start and goal locations.
24 This unit is in charge of constructing a set of conflict-free plans for all the agents. Thus, a typical
25 setting for MAPF can be attributed as *centralized* and *fully observable*.

26 In many real-world domains, however, the central controller does not exist, or, even if it does, it may
27 not possess full information about the environment. For example, consider a fleet of service robots
28 delivering some items in a human-shared environment, e.g., the robots delivering drugs in the hospital.
29 Each of these robots is likely to have access to the global map of the environment (e.g., the floor
30 plan), possibly refined through the robot's sensors. However, the connection to the central controller
31 may not be consistent. Thus, the latter may not have accurate data on the robots' locations and,
32 consequently, cannot provide valid MAPF solutions. In such scenarios, *decentralized approaches*
33 to the MAPF problems, when the robots themselves have to decide their future paths, are essential.
34 Moreover, decentralized approaches may be preferable due to the poor scalability of the centralized
35 ones. In this work, we aim to develop such an efficient decentralized approach.

36 It is natural to frame the decentralized MAPF as a
 37 sequential decision-making problem where at each
 38 time step, each agent must choose and execute an
 39 action that will advance it to the goal and, at the same
 40 time, will not disallow other agents to reach their
 41 goals as well. The result of solving this problem
 42 is a policy that, at each moment, tells which action
 43 to execute. To form such a policy, learnable meth-
 44 ods are commonly used, for example, reinforcement
 45 learning (RL), which is especially beneficial in tasks
 46 with incomplete information [2, 3, 4]. However, even
 47 state-of-the-art model-free RL methods generally can-
 48 not efficiently solve long-horizon problems with the
 49 involved casual structure [5, 6], and they are often
 50 inferior to the search-based methods when solving
 51 problems with hard combinatorial structure.

52 The additional challenges that make the MAPF prob-
 53 lems challenging for RL are as follows. First, we
 54 want the policy to be highly generalizable to previ-
 55 ously unseen environments, which may differ sig-
 56 nificantly in scale and topology from the ones used
 57 during the learning stage. In MAPF, our primary in-
 58 terest lies not in how well the agents learn to behave
 59 in the environment(s) used for training, but rather how
 60 well they perform in any arbitrary (even out-of-the-
 61 distribution) environment. Second, MAPF problems
 62 are naturally dependent on the goal locations of the
 63 agents, meaning that even in the same environment
 64 (map), the goals may vary significantly. Finally, effectively training in a complex observation and
 65 action spaces poses challenges even for state-of-the-art multi-agent reinforcement learning (MARL)
 66 methods.

67 To this end, in this work we suggest not to solve the MAPF problem directly by RL but rather to
 68 decompose it into a series of sub-tasks utilizing heuristic search algorithms and then solve these
 69 sub-tasks efficiently with a learnable policy, that is obtained through the decentralized training. The
 70 general pipeline of our solution is the following. Each agent plans an individual path to its goal by
 71 the conventional heuristic search algorithm without considering the other agents (we also introduce
 72 an additional technique to penalize paths that are likely to cause deadlocks). Then a waypoint on this
 73 path is chosen in some vicinity of the agent, which becomes its local goal. To reach it, a learnable
 74 policy is utilized, which takes both static obstacles and the locally observable agents into account.
 75 Once a waypoint is reached, or the agent goes too far away from it the cycle repeats.

76 Empirically we compare our method, which we name FOLLOWER, to a range of both learnable and
 77 non-learnable state-of-the-art competitors and show that it *i*) consistently outperforms the competitors
 78 when the number of agents is high; *ii*) better generalizes to unseen environments compared to other
 79 learnable solvers; *iii*) may outperform a centralized search-based solver in certain setups.

80 2 Related Works

81 **Lifelong MAPF** LMAPF is an extension of MAPF when the agents are assigned new goals upon
 82 reaching their current ones. Similarly, in (online) multi-agent pickup and delivery (MAPD), agents
 83 are continuously assigned tasks, comprising two locations that the agent has to visit in a strict order
 84 – pickup location and delivery location. Typically, the assignment problem is not considered in
 85 LMAPF/MAPD. However, there exist works that include the assignment task into the problem,
 86 see [7, 8] for example.

87 In [9], several variants to tackle MAPD were proposed differing in the amount of data the agents
 88 share. Yet, even the decoupled (as attributed by the authors) algorithms based on Token Swapping
 89 rely on global information, i.e., the one provided by the central unit. An enhanced Token Swapping

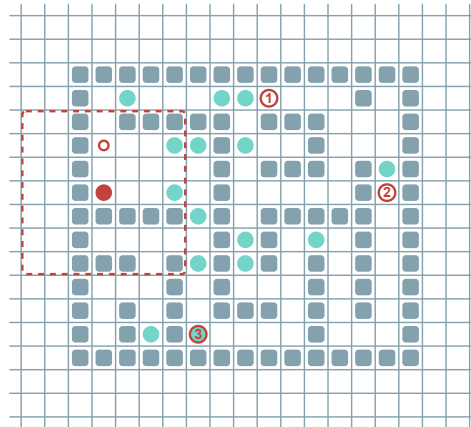


Figure 1: An example of a decentralized LMAPF instance is depicted below. Each agent is represented by a filled circle. The red agent has visibility limited to the positions of other agents within its field of view, indicated by a dotted red line. The red circles with numbers represent the goals that the agent needs to reach. The next goal is revealed only after the previous one is achieved. The small red circle indicates the subgoal the agent needs to accomplish to progress towards its goal.

90 variant that considers kinematic constraints was introduced in [10]. In [11], an efficient rule-based re-
91 planning approach to solve MAPF was introduced that is naturally capable of solving LMAPD/MAPD
92 problems. It did not rely on the several restrictive assumptions of Token Swapping and was empirically
93 shown to outperform the latter.

94 Finally, one of the most recent and effective LMAPF solvers is the RHCR algorithm presented in [12].
95 It relies on the idea of bounded planning, i.e., constructing not a complete plan but rather its initial
96 part. RHCR is a centralized solver that relies on the full knowledge of the agents' locations, their
97 current paths, goals, etc. In this work we empirically compare with RHCR and show that our method
98 is superior when the number of agents is high.

99 **Decentralized MAPF** This setting entails that the paths/actions of the agents are decided not by a
100 central unit but by the agents themselves. Numerous approaches, especially the ones tailored to the
101 robotics applications, boil this problem down to reactive control, see [13, 14, 15] for example. These
102 methods, however, are often prone to deadlocks. Several MAPF algorithms can also be implemented
103 in a decentralized manner. For example, in [16] MAPP algorithm was introduced that relies on the
104 individual pathfinding for each agent and a set of rules to determine priorities and choose actions to
105 avoid conflicts when they happen along the paths. In [11] PIBT algorithm was introduced in which
106 the agents also pick their actions individually (at each time step) based on specific rules. In general,
107 most rule-based MAPF solvers, like [17], can be implemented in such a way that each agent decides
108 its actions. However, in this case, the implicit assumption is that the agents can communicate to share
109 the relevant information (or that they have access to the global MAPF-related data). In contrast, our
110 work assumes that the agents are unable to communicate with one another or a central unit, which
111 significantly increases the complexity of the problem.

112 **Learnable MAPF** This direction has been getting increased attention recently. In [18], a seminal
113 PRIMAL method that utilized reinforcement learning and imitation learning to solve MAPF in a
114 decentralized fashion was introduced. Later in [19], it was also tailored to solve LMAPF. The new
115 version got the name PRIMAL2. Since that, numerous learning-based MAPF solvers emerged, and
116 it became common to compare against PRIMAL/PRIMAL2 (we also compare with it in our work).
117 For example, in [20], another learning-based approach was proposed, tailored explicitly to agents
118 with a non-trivial dynamic model, such as quadrotors. In [21] DHC method that utilized the agents'
119 communications to solve decentralized MAPF efficiently was described. Another communication-
120 based learnable approach, PICO, was presented in [22]. Overall, currently, a wide range of learnable
121 decentralized MAPF solvers exist. However, to the best of our knowledge, they all rely on the
122 communication between the agents or on access to the global MAPF-related data (like in PRIMAL,
123 where each agent knows the goal locations of the others). We lift these assumptions in this work.

124 **MARL** A separate direction in RL can be distinguished that specifically considers the multi-agent
125 setting (MARL) [23]. Mainly these approaches consider game environments (like Starcraft [24])
126 in which pathfinding is not of the primary importance. However, several MARL methods, such as
127 QMIX [3], MAPPO [25], have been adapted specifically for the MAPF task [26]. However they rely
128 on the information sharing between agents.

129 Much attention is paid to multi-agent learnable methods in robotics [27]. Often, the value-based
130 approaches are used to control small groups of agents on simple maps like in [28] where a group
131 of 4 agents is considered. In [29], a combination of Particle Swarm Optimization and Q-Learning
132 controlling up to 100 agents is used. In [30], the model-based DynaQ method is used to learn agents
133 in the knowledge exchange mode. Some works [31, 32] use value-based approaches with prior
134 knowledge of how to interact with other agents. In [33] (MAPPER) an evolutionary reinforcement
135 learning was used for MAPF task. This work also uses a global planner to determine sub-goals in
136 learning one agent. In multi-agent mode, agents using ineffective policies are eliminated and only
137 successful agents continue to be trained.

138 3 Background

139 **Multi-agent Pathfinding** In (Classical) Multi-agent pathfinding [1], the timeline is discretized to
140 the time steps, $T = 0, 1, 2, \dots$ and the workspace, where K agents operate, is discretized to a graph
141 $G = (V, E)$, whose vertices correspond to the locations and the edges to the transitions between

142 these locations. K start and goal vertices are given and each agent i has to reach its goal $g_i \in V$
 143 from the start $s_i \in V$. At each time step, an agent can either stay in its current vertex or move to an
 144 adjacent one. An individual plan for an agent p_i ¹ is a sequence of actions that transfers it between
 145 two designated vertices. The plan’s cost is the time step when the agent reaches the goal.

146 The MAPF problem asks to find a set of K plans s.t. each agent reaches the goal without colliding
 147 with other agents. Formally, two collisions are usually distinguished: vertex collision, when the
 148 agents occupy the same vertex at the same time step, and edge collision, when the agents use the
 149 same edge at the same time step.

150 *Lifelong MAPF* (LMAPF) is a variant of MAPF where immediately after an agent reaches its goal, it
 151 is assigned to another one (via an external assignment procedure) and has to continue its operation.
 152 Thus, LMAPF generally asks to find not a fixed set of K plans but rather to *i*) find a set of K initial
 153 plans and *ii*) update each agent’s plan when it reaches the current goal and receives a new one. In
 154 extreme cases, when some goal is reached at each step, the plans’ updates are needed constantly (i.e.,
 155 at each time step).

156 **The Considered Decentralized LMAPF Problem** Consider a set of agents operating in the shared
 157 environment, represented as a graph $G = (V, E)$. The timeline is discretized to the time steps
 158 $T = 0, 1, \dots, T_{max}$, where T_{max} is the episode length. Each agent is located initially at the start
 159 vertex and is assigned to the current goal vertex. If it reaches the latter before the episode ends, it is
 160 immediately assigned another goal vertex. We assume that the *goal assignment* unit is external to the
 161 system, and the agents’ behavior does not influence the goal assignments. An agent can reach the
 162 goal by performing the following actions: wait at the current vertex, and move to an adjacent vertex.
 163 The duration of each action is uniform, i.e., 1 time step. We assume that the outcomes of the actions
 164 are deterministic and no inaccuracies occur when executing the actions.

165 Each agent has complete knowledge of the graph G . However, it can observe the other agents only
 166 *locally*. When observing them, *no communication* is happening. Thus an agent does not know the
 167 (current) goals or intended paths of the other agents. It observes only their locations. The observation
 168 function can be defined differently depending on the type of graph. In our experiments, we use
 169 4-connected grids and assume that an agent observes the other agents in the area of the size $m \times m$,
 170 centered at the agent’s current position.

171 Our task is to construct an individual policy π for each agent, i.e., the function that takes as input
 172 a graph (global information) and (a history of) observations (local information) and outputs a
 173 distribution over actions. Equipped with such policy, an agent at each time step samples an action
 174 from the distribution suggested by π and executes it in the environment. This continues until time
 175 step T_{max} is reached when the episode ends. Upon that, we compute the *throughput* as the ratio of
 176 the episode length to the number of goals achieved by all agents. This metric is used to compare
 177 different policies: we say that π_1 outperforms π_2 (in a particular episode) if the throughput of the
 178 former is higher.

179 **Partially Observable Markov Decision Process** We consider a partially observable multi-agent
 180 Markov decision process [34, 35]: $M = \langle S, A, U, P, R, O, \gamma \rangle$. At each timestep, each agent $u \in U$,
 181 where $U = 1, \dots, n$, chooses an action $a^u \in A$, forming a joint action $\mathbf{j} \in \mathbf{J} = J^n$. This joint action
 182 leads to a change in the environment according to the transition function $P(s'|s, \mathbf{j}) : S \times \mathbf{J} \times S \rightarrow [0, 1]$.
 183 After that each agent receives individual observations $o^u \in O$ based on the global observation function
 184 $G(s, a) : S \times A \rightarrow O$. And individual reward $R(s, u, \mathbf{j}) : S \times U \times \mathbf{J} \rightarrow \mathbb{R}$, based on the current
 185 state, agent, and joint action. To make decisions each agent maintains an action-observation history
 186 $\tau^u \in T = (O \times A)^*$, which is used to condition a stochastic policy $\pi^u(a^u|\tau^u) : T \times A \rightarrow [0, 1]$.
 187 The task of the learning process is to optimize the policy π^u for individual each agent in order to
 188 maximize the expected cumulative reward over time.

189 4 Learn to Follow

190 The suggested approach to solve the considered LMAPF problem, which we dub FOLLOWER,
 191 comprises of the two complimentary modules combined into a coherent pipeline shown in Fig. 2.

¹In MAPF literature, a plan is typically denoted with π . However, in RL, this is reserved to denote the policy. As we use both MAPF and RL approaches in this work, we denote a plan as p .

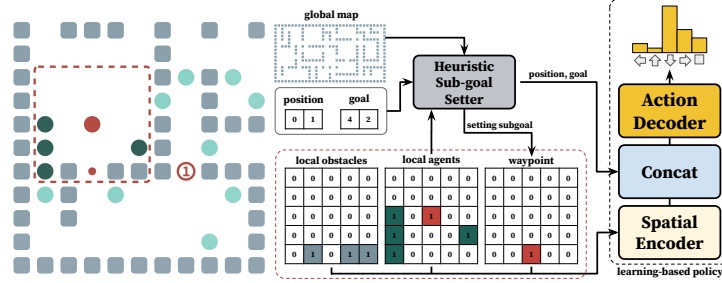


Figure 2: The general pipeline of the FOLLOWER approach. The action selection policy for each agent is decentralized and consists of two modules: Heuristic Sub-goal Decoder, which address long-term path planning problem and Learning-based Policy optimization module, which addresses the short-term conflict resolution task.

192 First, a *Heuristic Sub-goal Decoder* is used to construct an individual path to the goal and choose a
 193 waypoint on this path that becomes the agent’s local goal, which we also call a sub-goal. Second, a
 194 *Learnable Follower* is invoked to reach the sub-goal. This module decides which actions to take at
 195 each time step until the sub-goal is reached or until the agent gets too far away from it. In both cases
 196 the sub-goal decoder is called again and the cycle repeats.

197 4.1 Heuristic Sub-goal Decoder

198 In essence the purpose of this module is to provide a waypoint (sub-goal) in the vicinity of the agent,
 199 pursuing which will allow agent to progress towards its (global) goal. A conventional heuristic
 200 search algorithm, i.e. A*, is used to construct a path to the latter from the current location. Global
 201 information on the locations of the static obstacles, i.e. the map, is used for pathfinding. The other
 202 agents are not taken into account at this stage, thus the constructed path may go through them. Once
 203 the path is built a node located K steps away from the current position is chosen as the current
 204 sub-goal. Here K is the user-specified parameter.

205 An crucial design choice for this module is what individual path to build. On the one hand, A* finds
 206 the shortest (individual) path to the goal. On the other, as we noted empirically, when the number
 207 of agents is very high and each agent is following the shortest path, congestion often arise in the
 208 bottleneck parts of the map, such as corridors or doors. This degrades the performance dramatically.
 209 To this end we suggest to search not for the shortest paths but rather the evenly dispersed paths. This
 210 is implemented as follows.

211 At each time step the information on the locations of the locally observed agents is stored, in what we
 212 call a heatmap, and used to compute the additional transition costs for individual pathfinding. The
 213 number of times the other agents were seen in a certain location (grid cell in our experiments) is
 214 multiplied by the user-defined parameter C and added to the transition cost to that location. Intuitively,
 215 if many agents are noticed in particular areas of the map the transition costs of the latter are increased
 216 so A* will avoid them. This balances the distribution of the agents’ paths across the map and
 217 contributes to collision avoidance. Indeed, each agent maintains its own heatmap and performs
 218 pathfinding individually, thus the assumption that the agents do not share any data is not violated.

219 4.2 Learnable Follower

220 This module implements a learnable policy that is tailored to achieve the provided sub-goals while
 221 avoiding collision with the other agents. The policy function is approximated by a (deep) neural
 222 network and, as the agents are assumed to be homogeneous, a single network is utilized during
 223 training (a technique referred to as *policy sharing*). This approach is beneficial for complex tasks
 224 and large maps where it would be infeasible to learn a separate neural network for each agent, as the
 225 number of parameters increases linearly with the number of agents.

226 The input to the neural network represents the local observation of an agent and is comprised of a
 227 $3 \times m \times m$ tensor, where $m \times m$ is the observation range. The channels of the tensor encode the
 228 locations of the static obstacles, other agents and the current sub-goal respectively – see Fig. 2. If the
 229 latter is out of the agent’s field of view, it is projected into the nearest cell (similarly to [19]).

230 The input goes through the *Spatial Encoder* first, then *Concat block* combines both spatial and
 231 non-spatial features (the position of the agent on a map and its global goal). This is followed by an
 232 *Action Decoder* that uses the function f for approximating the state using observation history (the
 233 positions of other agents and the presence of obstacles) to make a decision. The network’s output is a
 234 probability distribution over possible actions.

235 The whole pipeline is trained with a policy optimization algorithm using the reward function separated
 236 into the two components: upon reaching a sub-goal an agent receives a small intrinsic positive reward
 237 of r_s , whose value was determined empirically; upon reaching the global goal a conventional RL
 238 reward $r_g = 1$ is received. If while reaching the current goal the agent goes too far away from it, the
 239 heuristic sub-goal decider is invoked again. This mechanism is helpful in scenarios when to progress
 240 towards the global goal it is actually more beneficial to make a detour to avoid congestion with the
 241 other agents. Practically wise, a goal is recalculated if the agent’s distance from its target exceeds a
 242 certain threshold, which is determined by a hyperparameter H .

243 The task of the learning process is to optimize the shared policy π_θ^u (i.e. the same policy for each
 244 agent) to maximize the expected cumulative reward. During the training process, rollouts (sequences
 245 of observation and action pairs) are gathered asynchronously from multiple environments with
 246 varying numbers of agents. The shared policy π_θ (actor network) is continually updated using the
 247 PPO clipped loss [36]: $\max_\theta \frac{1}{N} \sum_{u=1}^n \sum_j \sum_{\tau^u} \pi_\theta(a^u | \tau^u) \hat{A}_{\text{clip}}(\tau^u, a^u) - \beta H(\pi_\theta(\cdot | \tau^u))$.

248 Here, β is a coefficient that controls the entropy H , and \hat{A} denotes the unclipped advantage function
 249 calculated using returns \hat{R} for each step t with observation history τ^u : $\hat{A}(\tau^u, a^u) = \hat{R}_t^u - V_\phi(\tau^u)$,
 250 with $\hat{R}_t^u = \sum_{k=0}^{T-1} \gamma^k r_{t+k}^u$. Here, we have a shared critic value function V_ϕ , which is optimized using
 251 the following equation: $\min_\phi \frac{1}{N} \sum_{u=1}^n \sum_j \sum_{\tau^u} (V_\phi(\tau^u) - \hat{R}_t^u)^2$.

252 In practice, the observation history τ^u is effectively modeled using a recurrent neural network (RNN)
 253 integrated into the actor and critic networks. The actor network is parameterized by θ , while the critic
 254 network is parameterized by ϕ . In our approach, we specifically utilize the GRU architecture [37].

255 The introduced intrinsic reward function allows the efficient training of an agent using relatively short
 256 rollouts, as evidenced by our experimental results, which demonstrate that a rollout length of 8 is
 257 sufficient for training. This is crucial for ensuring lifelong learning, as episodes may not have a clear
 258 ending point.

259 During the inference phase, each agent uses a copy of the trained weights, and other parameters
 260 remain unchanged. The proposed FOLLOWER scheme, despite its simplicity, allows the agent to
 261 separate the two components of the overall policy transparently and does not require the involvement
 262 of any expert data for training. The learning process is end-to-end and the number of hyperparameters
 263 (such as K and H) that affect the result is relatively small. Finally, the reward function used is simple
 264 and does not require involved manual shaping.

265 5 Experimental Evaluation

266 To evaluate the efficiency of the proposed method², we have conducted a set of experiments, compar-
 267 ing it with the existing learnable and search-based algorithms on different grid maps. The episode
 268 length was set to 512 in all experiments. The agents field-of-view was 11×11 . When training
 269 FOLLOWER we used the following values of the reward components: $r_g = +1$ and $r_s = +0.1$. The
 270 *Spatial Encoder* was realized as ResNet neural model [38], the *Concat block* – as a Multi-Layer
 271 Perceptron (MLP), and the *Action Decoder* – as a recurrent neural network, separated for actor
 272 and critic and based on GRU [37]. In the experiments, values of 2 and 10 were used for K and
 273 H respectively. The weighting coefficient C for sub-goal setter was set to 0.4. More information
 274 about which (hyper) parameters were tuned and how is provided in the Appendix. After fixing all the
 275 parameters, the final policy was trained using a single TITAN RTX GPU in approximately 1 hour.

²We are committed to open-source FOLLOWER.

276 **5.1 Comparison With the Learnable Methods**

277 In the first series of experiments, we have compared FOLLOWER with the two state-of-the-art
 278 learnable MAPF solvers, i.e. PRIMAL2 [19] and PICO [22]. Similarly to FOLLOWER, both are
 279 decentralized and rely on the local observations of the other agents. However, PRIMAL2 assumes
 280 that the local observations contain not only information about the current locations of the agents but
 281 also about their goals on the global map. PICO assumes that the agents can communicate, through
 282 selected central agent. Recall that our solver has access neither to any information about the other
 283 agents except their current locations nor to communication between the agents.

284 As learnable methods assume training on a certain maps topology,
 285 we use the maps suggested by the authors of the respective baselines
 286 for a fair comparison. Specifically, we compare with PRIMAL2 on
 287 the maze-like maps of size 65×65 on which PRIMAL2 was trained,
 288 and we compare with PICO on the maps with random obstacles
 289 described in the PICO paper. The visualizations of the maps are
 290 given in Appendix. We used the readily available weights for PRI-
 291 MAL2 neural network (from the authors’ repository). PICO was
 292 trained by us using the open-source code of its authors. Our method,
 293 FOLLOWER, was trained using the hyperparameters described in
 294 Appendix, and only PRIMAL2 maps were used for training. When
 295 training FOLLOWER, we vary the number of agents in range: 16, 32,
 296 64, 128. After the training phase, we run the solvers on ten different
 297 maze-like/random maps that were not used while training. Each map
 298 was populated with varying numbers of agents: from 2 to 256. The
 299 goals for LMAPF were generated and assigned to agents randomly.

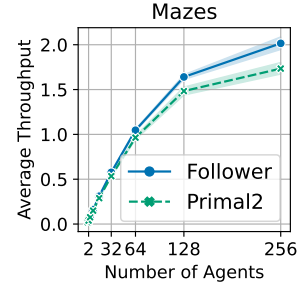


Figure 3: Average throughput on maze-like environments. The shaded area indicates 95% confidence intervals.

300 **FOLLOWER vs. PRIMAL2 results** are depicted on Fig. 3. The OX axis shows the number of
 301 agents and OY axis shows the average throughput. Indeed, when the number of agents is low both
 302 algorithms demonstrate similar results. However, with an increasing number of agents the gap
 303 in performance is getting pronounced. The throughput of the FOLLOWER is 11% better for 128
 304 agents and 17% better for 256 agents. Overall, one can claim that despite having access to less
 305 MAPF-related data FOLLOWER outperforms PRIMAL2 when the number of agents is not low, i.e. in
 306 cases where the potential conflicts between the agents are not rare.

307 **FOLLOWER vs. PICO results** are presented in Table 1. Clearly, FOLLOWER demonstrates a
 308 superior performance across all scenarios. The poor performance of PICO can be attributed to the
 309 inherent difficulties in learning effective communication strategies for prioritizing large number
 310 of agents. Authors of PICO trained their method on 8 agents. We hypothesize that this limited
 311 population size may have impeded the acquisition of knowledge necessary for effective coordination
 312 among a larger number of agents. Both FOLLOWER and PRIMAL2 outperform PICO showing their
 313 ability to generalize (as they were not trained on PICO type of maps), with FOLLOWER being the
 314 ultimate winner.

Table 1: The comparison of FOLLOWER with PICO on random maps with different obstacle densities, taken from PICO evaluation setup.

Algorithm	Agents	Obstacle Density			
		0%	10%	20%	30%
FOLLOWER	8	0.61 (± 0.01)	0.57 (± 0.02)	0.49 (± 0.04)	0.38 (± 0.21)
PRIMAL2	8	0.44 (± 0.03)	0.39 (± 0.04)	0.3 (± 0.05)	0.19 (± 0.11)
PICO	8	0.19 (± 0.01)	0.18 (± 0.03)	0.14 (± 0.04)	0.05 (± 0.05)
FOLLOWER	16	1.1 (± 0.03)	0.96 (± 0.05)	0.85 (± 0.19)	0.56 (± 0.34)
PRIMAL2	16	0.79 (± 0.03)	0.67 (± 0.06)	0.51 (± 0.08)	0.31 (± 0.14)
PICO	16	0.31 (± 0.03)	0.25 (± 0.04)	0.23 (± 0.06)	0.08 (± 0.06)
FOLLOWER	32	1.81 (± 0.05)	1.45 (± 0.15)	1.21 (± 0.27)	0.84 (± 0.39)
PRIMAL2	32	1.25 (± 0.04)	1.02 (± 0.11)	0.71 (± 0.12)	0.4 (± 0.16)
PICO	32	0.46 (± 0.05)	0.35 (± 0.1)	0.28 (± 0.12)	0.12 (± 0.09)
FOLLOWER	64	2.6 (± 0.11)	1.88 (± 0.33)	1.24 (± 0.23)	0.71 (± 0.36)
PRIMAL2	64	1.63 (± 0.05)	1.15 (± 0.15)	0.73 (± 0.13)	0.44 (± 0.18)
PICO	64	0.42 (± 0.07)	0.41 (± 0.13)	0.28 (± 0.12)	0.11 (± 0.1)

315 **Out of the distribution evaluation.**

316 An important attribute of any learnable algorithm is the so-called *generalization*, i.e. the ability to solve
317 problem instances that are not alike
318 to the ones that were used for training.
319 We have already seen that FOLLOWER
320 generalizes better than PRIMAL2 to
321 the PICO type of maps (random ones).
322 Now we run an additional test when
323 we evaluated both algorithms on two
324 (unseen while learning) maps from the
325 well-known in the MAPF community

326 MovingAI benchmark [1]: warehouse-10-20-10-2-1 and lak303d. The former map is 63×171
327 in size and represents the warehouse environment. It is similar to a certain extent to the maze-maps
328 on which FOLLOWER and PRIMAL2 were trained. The latter map is a video-game map that was
329 downscaled by us to have size 95×95 . Its topology is quite different from the one of the maps used
330 for training FOLLOWER and PRIMAL2. The results of these experiments are presented in Fig. 4.
331 Note that we did not evaluate PICO on out-of-the-distribution maps due to its poor performance in
332 the previous experiment.
333

334
335 On warehouse-10-20-10-2-1 FOLLOWER and PRIMAL2 demonstrate similar performance, with
336 our method achieving slightly higher throughput. I.e. the FOLLOWER’s throughput is 9.5% higher
337 for 128 agents and 4% higher for the 256 agents. The results on lak303d are quite different. First
338 of all, both algorithms have much lower average throughput compared to maze-like and warehouse
339 environments. This is expected, as the topology of the game map differs a lot from the latter maps.
340 Second, the throughput of FOLLOWER is significantly higher, providing another evidence (in addition
341 to the results on PICO maps) that the generalization ability of our approach is better.

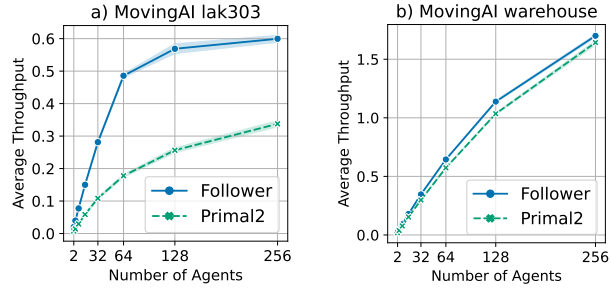


Figure 4: The results on a) lak303d and b) warehouse maps. The shaded area indicates 95% confidence intervals.

342 **5.2 Comparison With the Centralized Search-based Solver**

343 While most of the learnable approaches compare their results only with other learnable methods, we
344 have also compared FOLLOWER with the state-of-the-art search-based algorithm for solving LMAPF
345 – RHCR³ [12]. In contrast to the proposed method, RHCR is a centralized approach that coordinates
346 all the agents and does not restrict the observation and/or communication abilities of the agents. This
347 planner has several parameters that influence its performance. We have varied the planning horizon
348 (2, 5, 10, 20), the re-planning rate (1, 5) and found that the best results are achieved when the first
349 parameter is set to 20 and the second one to 5. The time limit for re-planning was set to either 1 or 10
350 seconds. The MAPF solver used in RHCR was set to PBS [39]. The rest parameters were left default.

351 The comparison was conducted on the same warehouse map as in the original paper [12]. The
352 possible placements of start and goal locations were also restricted in the same way as in the original
353 paper. The number of agents in this experiment reached 192 as no more agents are able to be placed
354 with the given restrictions to start locations. We generated 10 random instances per each number of
355 agents.

356 Besides RHCR we have also evaluated different versions of our solver. First, we want to assess the
357 impact of the learnable component on the FOLLOWER’s performance. To this end, we removed it
358 from FOLLOWER and let each agent simply plan its path with A* and perform the first action. In case
359 the path can not be found the action is selected randomly. We refer to this approach as Randomized
360 A*. We evaluated two versions of Randomized A*: the one that treats the other agents (within
361 field-of-view) as obstacles and the one that does not. Next, we were interested in how weighting the
362 transition costs for A* affects the FOLLOWER’s performance. Thus, we have created a version of
363 FOLLOWER that doesn’t include this weighting (i.e. $C = 0$) and all transitions have uniform cost.

364 The results of this experiment are depicted in Fig. 5. Clearly both version of Randomized A* are
365 outperformed by FOLLOWER. This confirms that the learnable policy is crucial to FOLLOWER.
366 The introduced technique of penalizing the transitions to the areas where the other agents are often
367 observed is also important, as in its absence the results of FOLLOWER are on par with Randomized A*.

³We used an implementation of RHCR from the authors’ repository

368 Only combining the learnable module with the weighting technique we end with the best-performing
 369 method.

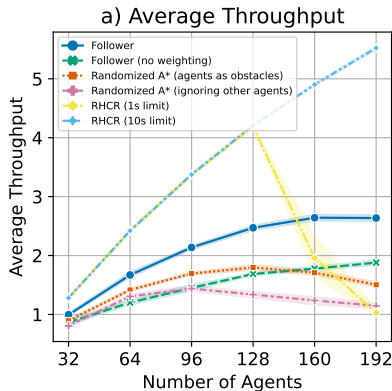


Figure 5: Average throughput on warehouse map. The shaded area indicates 95% confidence intervals.

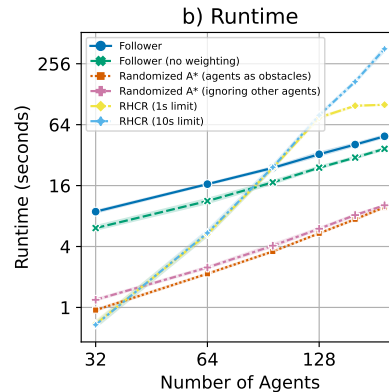


Figure 6: Average runtime on warehouse map.

370 Both versions of RHCR significantly outperform competitors in instances with up to 128 agents.
 371 However, when the number of agents increases to 160 and 192 the performance of RHCR with
 372 1s time cap for re-planning degrades dramatically and it gets outperformed by FOLLOWER. This
 373 pinpoints the principal limitation of the centralized approach – it does not scale well to large number
 374 of agents when the time limit for finding a MAPF solution is imposed. To better understand how the
 375 runtime of the evaluated methods is affected by the increasing number of agents, see Fig. 6. Here
 376 each data point says how much time on average was spent to solve a LMAPF instance (on a single
 377 CPU, 1 thread). Indeed, FOLLOWER scales much better compared to RHCR. Moreover in practice it
 378 can be parallelized, i.e. run on each agent individually, while RHCR – can not.

379 5.3 Summary

380 The observed results let us infer the following conclusions. First, the suggested approach outperforms
 381 the learnable decentralized competitors, when it comes to a large number of agents and/or maps that
 382 are different from the ones used for training. Second, the learnable component of FOLLOWER is
 383 crucial to its high performance (as well as the introduced weighting technique). Third, when the
 384 number of agents is significantly high, FOLLOWER can outperform the centralized LMAPF solver
 385 when a (reasonable) time cap for the latter is introduced.

386 6 Conclusion

387 This study addresses the challenging problem of decentralized lifelong multi-agent pathfinding. The
 388 proposed FOLLOWER approach utilizes a combination of a planning algorithm for constructing a
 389 long-term plan and reinforcement learning for reaching short-term sub-goals and resolving local
 390 conflicts. The proposed method consistently outperforms decentralized learnable competitors in
 391 challenging scenarios. Moreover, our approach can show better results, compared to state-of-the-art
 392 centralized planner in certain setups. Directions for future research may include: enriching the action
 393 space of the agents, handling uncertain observations and external (stochastic) events.

394 7 Limitations

395 As in many other works on that topic (including the ones we compare with) we rely on the following
 396 assumptions. The map of the environment is accurate and the configuration of the static obstacles
 397 does not change. The agents are assumed to have perfect localization and mapping abilities. The
 398 agents execute actions accurately and their moves are synchronized. All these may be considered as
 399 the limitations as in real world, e.g. in robotic applications, many of the assumptions do not hold.

References

- 400
- 401 [1] Roni Stern, Nathan R Sturtevant, Ariel Felner, Sven Koenig, Hang Ma, Thayne T Walker,
402 Jiaoyang Li, Dor Atzmon, Liron Cohen, TK Satish Kumar, et al. Multi-agent pathfinding:
403 Definitions, variants, and benchmarks. In *Proceedings of the 12th Annual Symposium on*
404 *Combinatorial Search (SoCS 2019)*, pages 151–158, 2019.
- 405 [2] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei a Rusu, Joel Veness, Marc G
406 Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, Stig Pe-
407 tersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan
408 Wierstra, Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement
409 learning. *Nature*, 518(7540):529–533, 2015.
- 410 [3] Tabish Rashid, Mikayel Samvelyan, Christian Schroeder De Witt, Gregory Farquhar, Jakob
411 Foerster, and Shimon Whiteson. QMIX: Monotonic value function factorisation for deep
412 multi-agent reinforcement Learning. In *35th International Conference on Machine Learning,*
413 *ICML 2018*, volume 10, pages 6846–6859, 2018.
- 414 [4] Danijar Hafner, Timothy Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering atari with
415 discrete world models. In *ICLR*, 2021.
- 416 [5] Stephanie Milani, Nicholay Topin, Brandon Houghton, William H. Guss, Sharada P. Mohanty,
417 Oriol Vinyals, and Noboru Sean Kuno. The MineRL Competition on Sample-Efficient Rein-
418 forcement Learning Using Human Priors: A Retrospective. In *NeurIPS Competition track*,
419 2020.
- 420 [6] Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering Diverse Domains
421 through World Models. 2023.
- 422 [7] Minghua Liu, Hang Ma, Jiaoyang Li, and Sven Koenig. Task and path planning for multi-agent
423 pickup and delivery. In *Proceedings of the 18th International Conference on Autonomous*
424 *Agents and Multiagent Systems (AAMAS 2019)*, pages 1152–1160, 2019.
- 425 [8] Zhe Chen, Javier Alonso-Mora, Xiaoshan Bai, Daniel D Harabor, and Peter J Stuckey. Integrated
426 task assignment and path planning for capacitated multi-agent pickup and delivery. *IEEE*
427 *Robotics and Automation Letters*, 6(3):5816–5823, 2021.
- 428 [9] Hang Ma, Jiaoyang Li, TK Kumar, and Sven Koenig. Lifelong multi-agent path finding for
429 online pickup and delivery tasks. In *Proceedings of the 16th Conference on Autonomous Agents*
430 *and MultiAgent Systems (AAMAS 2017)*, pages 837–845, 2017.
- 431 [10] H. Ma, W. Hönl, T. K. S. Kumar, N. Ayanian, and S. Koenig. Lifelong path planning with
432 kinematic constraints for multi-agent pickup and delivery. In *Proceedings of the 33rd AAAI*
433 *Conference on Artificial Intelligence (AAAI 2019)*, pages 7651–7658, 2019.
- 434 [11] Keisuke Okumura, Manao Machida, Xavier Défago, and Yasumasa Tamura. Priority inheritance
435 with backtracking for iterative multi-agent path finding. In *Proceedings of the 28th International*
436 *Joint Conference on Artificial Intelligence (IJCAI 2019)*, pages 535–542, 2019.
- 437 [12] Jiaoyang Li, Andrew Tinka, Scott Kiesel, Joseph W Durham, TK Satish Kumar, and Sven
438 Koenig. Lifelong multi-agent path finding in large-scale warehouses. In *Proceedings of the*
439 *35th AAAI Conference on Artificial Intelligence (AAAI 2021)*, pages 11272–11281, 2021.
- 440 [13] Vladimir J. Lumelsky and KR Harinarayan. Decentralized motion planning for multiple mobile
441 robots: The cocktail party model. *Autonomous Robots*, 4(1):121–135, 1997.
- 442 [14] Jur Van den Berg, Ming Lin, and Dinesh Manocha. Reciprocal velocity obstacles for real-time
443 multi-agent navigation. In *Proceedings of The 2008 IEEE International Conference on Robotics*
444 *and Automation (ICRA 2008)*, pages 1928–1935. IEEE, 2008.
- 445 [15] Hai Zhu, Bruno Brito, and Javier Alonso-Mora. Decentralized probabilistic multi-robot collision
446 avoidance using buffered uncertainty-aware voronoi cells. *Autonomous Robots*, 46(2):401–420,
447 2022.

- 448 [16] Ko-Hsin Cindy Wang and Adi Botea. Mapp: a scalable multi-agent path planning algorithm
449 with tractability and completeness guarantees. *Journal of Artificial Intelligence Research*,
450 42:55–90, 2011.
- 451 [17] Boris de Wilde, Adriaan W ter Mors, and Cees Witteveen. Push and rotate: cooperative multi-
452 agent path planning. In *Proceedings of the 12th International Conference on Autonomous*
453 *Agents and Multiagent Systems (AAMAS 2013)*, pages 87–94, 2013.
- 454 [18] Guillaume Sartoretti, Justin Kerr, Yunfei Shi, Glenn Wagner, TK Satish Kumar, Sven Koenig,
455 and Howie Choset. Primal: Pathfinding via reinforcement and imitation multi-agent learning.
456 *IEEE Robotics and Automation Letters*, 4(3):2378–2385, 2019.
- 457 [19] Mehul Damani, Zhiyao Luo, Emerson Wenzel, and Guillaume Sartoretti. Primal _2: Pathfinding
458 via reinforcement and imitation multi-agent learning-lifelong. *IEEE Robotics and Automation*
459 *Letters*, 6(2):2666–2673, 2021.
- 460 [20] Benjamin Riviere, Wolfgang Hönig, Yisong Yue, and Soon-Jo Chung. Glas: Global-to-local
461 safe autonomy synthesis for multi-robot motion planning with end-to-end learning. *IEEE*
462 *Robotics and Automation Letters*, 5(3):4249–4256, 2020.
- 463 [21] Ziyuan Ma, Yudong Luo, and Hang Ma. Distributed heuristic multi-agent path finding with
464 communication. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*,
465 pages 8699–8705. IEEE, 2021.
- 466 [22] Wenhao Li, Hongjun Chen, Bo Jin, Wenzhe Tan, Hong Zha, and Xiangfeng Wang. Multi-
467 agent path finding with prioritized communication learning. *2022 International Conference on*
468 *Robotics and Automation (ICRA)*, pages 10695–10701, 2022.
- 469 [23] Annie Wong, Thomas Bäck, Anna V. Kononova, and Aske Plaat. Deep multiagent reinforcement
470 learning: challenges and directions. *Artificial Intelligence Review*, (0123456789), oct 2022.
- 471 [24] Mikayel Samvelyan, Tabish Rashid, Christian Schroeder De Witt, Gregory Farquhar, Nantas
472 Nardelli, Tim G.J. Rudner, Chia Man Hung, Philip H.S. Torr, Jakob Foerster, and Shimon
473 Whiteson. The StarCraft multi-agent challenge. In *Proceedings of the International Joint*
474 *Conference on Autonomous Agents and Multiagent Systems, AAMAS*, volume 4, pages 2186–
475 2188, 2019.
- 476 [25] Chao Yu, Akash Velu, Eugene Vinitsky, Yu Wang, Alexandre Bayen, and Yi Wu. The Surprising
477 Effectiveness of PPO in Cooperative Multi-Agent Games. 2021.
- 478 [26] Alexey Skrynnik, Alexandra Yakovleva, Vasilii Davydov, Konstantin Yakovlev, and Aleksandr I.
479 Panov. Hybrid Policy Learning for Multi-Agent Pathfinding. *IEEE Access*, 9:126034–126047,
480 2021.
- 481 [27] Mike Wesselhöft, Johannes Hinckeldeyn, and Jochen Kreutzfeldt. Controlling fleets of au-
482 tonomous mobile robots with reinforcement learning: A brief survey. *Robotics*, 11(5), 2022.
- 483 [28] Jixuan Zhi and Jyh-Ming Lien. Learning to herd agents amongst obstacles: Training robust
484 shepherding behaviors using deep reinforcement learning. *IEEE Robotics and Automation*
485 *Letters*, 6(2):4163–4168, 2021.
- 486 [29] Syed Irfan Ali Meerza, Moinul Islam, and Md. Mohiuddin Uzzal. Q-learning based particle
487 swarm optimization algorithm for optimal path planning of swarm of mobile robots. In *2019*
488 *1st International Conference on Advances in Science, Engineering and Robotics Technology*
489 *(ICASERT)*, pages 1–5, 2019.
- 490 [30] Emanuele Vitolo, Alberto San Miguel, Javier Civera, and Cristian Mahulea. Performance
491 evaluation of the dyna-q algorithm for robot navigation. In *2018 IEEE 14th International*
492 *Conference on Automation Science and Engineering (CASE)*, pages 322–327, 2018.
- 493 [31] Bo Li and Hongbin Liang. Multi-robot path planning method based on prior knowledge and
494 q-learning algorithms. *Journal of Physics: Conference Series*, 1624(4):042008, oct 2020.

- 495 [32] Hyansu Bae, Gidong Kim, Jonguk Kim, Dianwei Qian, and Sukgyu Lee. Multi-robot path
496 planning method using reinforcement learning. *Applied Sciences*, 9(15), 2019.
- 497 [33] Zuxin Liu, Baiming Chen, Hongyi Zhou, Guru Koushik, Martial Hebert, and Ding Zhao.
498 Mapper: Multi-agent path planning with evolutionary reinforcement learning in mixed dynamic
499 environments. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems*
500 (*IROS*), pages 11748–11754, 2020.
- 501 [34] Daniel S Bernstein, Robert Givan, Neil Immerman, and Shlomo Zilberstein. The complexity
502 of decentralized control of markov decision processes. *Mathematics of operations research*,
503 27(4):819–840, 2002.
- 504 [35] Leslie Pack Kaelbling, Michael L Littman, and Anthony R Cassandra. Planning and acting in
505 partially observable stochastic domains. *Artificial Intelligence*, 101(1-2):99–134, may 1998.
- 506 [36] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal
507 Policy Optimization Algorithms. pages 1–12, 2017.
- 508 [37] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation
509 of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*,
510 2014.
- 511 [38] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
512 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
513 pages 770–778, 2016.
- 514 [39] Hang Ma, Daniel Harabor, Peter J Stuckey, Jiaoyang Li, and Sven Koenig. Searching with
515 consistent prioritization for multi-agent path finding. In *Proceedings of the AAAI Conference*
516 *on Artificial Intelligence*, volume 33, pages 7643–7650, 2019.