

RDE: A Hybrid Policy Framework for Multi-Agent Path Finding Problem

Jianqi Gao¹, Yanjie Li¹, *Member, IEEE*, Xiaoqing Yang¹, and Mingshan Tan¹

Abstract—Multi-agent path finding (MAPF) is an abstract model for the navigation of multiple robots in warehouse automation, where multiple robots plan collision-free paths from the start to goal positions. Reinforcement learning (RL) has been employed to develop partially observable distributed MAPF policies that can be scaled to any number of agents. However, RL-based MAPF policies often get agents stuck in deadlock due to warehouse automation’s dense and structured obstacles. This paper proposes a novel hybrid MAPF policy, RDE, based on switching among the RL-based MAPF policy, the Distance heat map (DHM)-based policy and the Escape policy. The RL-based policy is used for coordination among agents. In contrast, when no other agents are in the agent’s field of view, it can get the next action by querying the DHM. The escape policy that randomly selects valid actions can help agents escape the deadlock. We conduct simulations on warehouse-like structured grid maps using state-of-the-art RL-based MAPF policies (DHC and DCC), which show that RDE can significantly improve their performance.

Index Terms—Path Planning for Multiple Mobile Robots or Agents, Planning, Scheduling and Coordination, Collision Avoidance.

I. INTRODUCTION

The widespread adoption of e-commerce and logistics has led to the increased use of warehouse automation. As depicted in Fig. 1, a flexible multi-robot system [1], [2] is employed for this purpose. By providing start and goal positions for a group of warehouse robots, MAPF can plan multiple collision-free paths [3]. Traditional MAPF policies [4]–[9] have been extensively researched and developed, but their real-time performance deteriorates as the number of agents increases, particularly in large-scale warehouse automation.

Researchers have recently used RL to acquire distributed MAPF policies, where each agent’s action is based on its partial observation space [10]–[15]. The RL-based MAPF policy can be applied to any number of agents or map sizes without requiring retraining. When agents conflict, the RL-based MAPF policy will punish them with negative rewards and make them stay at their current positions. The agent retakes new actions based on the observations in the next step. However, if the agent’s observations do not change in the next step, the RL-based policy will make the agent take the same action as before, leading to a path conflict again. When this process continues, the agent will get stuck in a deadlock, which can lead to the failure of MAPF. Moreover, the benefits of cooperative coordination using RL-based policy becomes insignificant if the agent has no other agents within its field

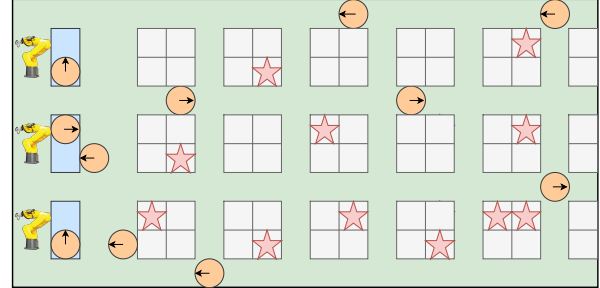


Fig. 1. The two-dimensional schematic of the warehouse automation. The orange circles represent the robots. The grey squares represent the inventory pods (obstacles), and the ones marked with light red stars are to be moved. The inventory stations with manipulators are shown on the left side. The robot enters the bottom of the inventory pod, lifts it, and carries it to its designated location.

of view (FOV). Conversely, RL-based MAPF policy may lead to agent detours.

This paper introduces a hybrid MAPF policy called RDE to address the difficulties outlined earlier. RDE combines the RL-based policy with the DHM-based policy and the escape policy. The RL-based policy mainly manages cooperation and avoids agent conflicts, while the DHM-based policy handles simple scenarios, where no other agents are within an agent’s FOV. We can directly query the DHM¹ to determine the next action where the agent is closest to the goal position and trying to stay straight. Additionally, if an agent is trapped in a deadlock state, the escape policy that relies on randomly selecting actions is activated to help resolve the issue. To verify the effectiveness of RDE, we utilize state-of-the-art RL-based MAPF policies (DHC [14] and DCC [15]) to conduct simulation experiments on warehouse-like structured grid maps, and the results show that RDE can significantly improve their performance.

The remainder of this paper is structured as follows. Section II introduces the research work related to MAPF and lifelong MAPF. Section III presents the problem formulation. Section IV describes the hybrid MAPF framework. Section V presents the simulation experiment results. Section VI provides the conclusions and future work.

II. RELATED WORK

A. Traditional MAPF

Finding an optimal MAPF solution has been proven to be NP-hard [17]. According to whether the optimal solution can

All authors are with the Department of Control Science and Engineering, Harbin Institute of Technology (Shenzhen), Shenzhen 518055, China (e-mail: gaojianqi205a@stu.hit.edu.cn; autolyj@hit.edu.cn; yxqsheep@gmail.com; tanmingshan033@gmail.com).

¹To obtain a DHM with the goal position as the source, we employ the Dijkstra [16] algorithm to compute the minimum distance to all grid positions. Details can be seen in section IV-C.

be obtained, we can divide the traditional MAPF policies into optimal and sub-optimal. The optimal MAPF policies [4], [5], [18], [19] also satisfies completeness. CBS [5] is the most commonly used optimal MAPF policy, and researchers have made many improvements to CBS. However, as the map size or agent's number grows, finding optimal solutions becomes challenging. A number of sub-optimal policies have been proposed, which can be divided into bounded and unbounded. The cost of the bounded sub-optimal policy's solution is no more than $(1 + \epsilon) \times c_{opt}$, where c_{opt} is the optimal cost and $\epsilon > 0$ is the sub-optimality factor [8]. The bounded sub-optimal policies are generally derived from optimal MAPF policies [20]–[22]. Bounded sub-optimal policies can, under certain conditions, improve the speed of solving and provide some guarantee of the quality of the solution. Unbounded sub-optimal MAPF policies [7], [23]–[26] can get solutions faster, but the quality of the solution is not guaranteed.

B. Learning-based MAPF

Learning-based MAPF policies are distributed methods, where every agent takes action based on its partial observation space. PRIMAL [10] combines RL with ODrM* [4]-based imitation learning (IL) to get a distributed MAPF policy. Based on the framework of PRIMAL, many learning-based MAPF policies are proposed. MAPPER [11] and G2RL [12] use A*-based shortest path to guide the policy learning. Agents who deviate from this path will be penalized. However, A*-based shortest paths are not unique and not optimal globally, destabilizing the learning process and the multi-agent implicit coordination. Instead of using a CNN [27] as the encoder of the observed state, [28] proposes a transformer-based [29] policy network for feature extraction. In addition, in the above policies each agent regards other agents as dynamic obstacles, leading to the environment's instability and making the training process challenging to converge. Researchers have recently enabled one agent to communicate with other agents. In [13], a graph neural network is used to help agents communicate with each other. Based on [13] and MAPPER [11], MAGAT [30] and AB-Mapper [31] make use of the attention mechanism to assess the relative importance of agent. DHC [14] uses graph convolution mechanism to communicate between agents. Instead of broadcast communication in FOV, DCC [15] first evaluates the importance of information and then selects some information to aggregate.

C. Hybrid MAPF Solver

One policies cannot effectively address every variant of the classical MAPF problem. Some researchers attempt to combine multiple MAPF policies to solve MAPF problems. In [32], the software is proposed to plan paths for thousands of trains within a few minutes, which incorporates many state-of-the-art MAPF policies. SWARM-MAPF [33] uses swarm-based policy [34] in open areas and CBS [5] in areas with dense obstacles, which improves the efficiency of solving problems. The spatially distributed multi-agent planner [35] identifies high and low-contention areas and

uses different MAPF policies for them. Some studies use hierarchical policies to solve MAPF problems. First, the map is spatially divided into small areas. The high-level performs global planning for each agent, and the low-level solves the MAPF of each small area. HMAPP [36] generates the high-level plan for each agent from a randomly picked shortest path from its start position to its goal position and then uses ECBS [8] to find conflict-free subpaths in every region.

III. PROBLEM FORMULATION

A. MAPF

The input of MAPF includes an undirected connected graph $G(V, E)$ and an agent set $N = \{1, \dots, i, \dots, m\}$. Every agent i has a start vertex, $v_i^s \in V$ and a unique goal vertex, $v_i^g \in V$. At each discrete time step $t = 0, \dots, \infty$, agent i is located in vertex v and takes action a_i^t . Then agent i moves to an adjacent vertex v' that meets $(v, v') \in E$ or stays in its current vertex v . The action can be represented as $a : v \rightarrow v'$ or $a(v) = v'$. The path of agent i can be represented by a sequence of actions $l_i = (a_i^1, a_i^2, \dots, a_i^t)$. A MAPF solution can be represented as $\mathbb{L} = (l_1, l_2, \dots, l_m)$. We assume all agents simultaneously move one grid or stand still at each time step. In this paper, the world is a two-dimensional four-connected unit grid map where every vertex has four adjacent vertices. The path cost of agent i is the discrete-time t_i^g when the agent i reaches the goal vertex v_i^g . Conflict is usually used in MAPF to represent the path plan collision of different agents. The main conflicts in MAPF are edge conflict $\langle i, j, v, v', t \rangle$ and vertex conflict $\langle i, j, v, t \rangle$.

B. Learning-based MAPF

As shown in Fig. 2, learning-based MAPF can be treated as a partially observable Markov decision process (POMDPs)² [37] solved with RL. At each timestep t , the i th agent has an observation o_i^t , which only provides partial information of $G(V, E)$, and then we independently compute an action a^t by the policy π_θ shared by all agents:

$$a_i^t \sim \pi_\theta(o_i^t), \quad (1)$$

where θ denotes the policy parameters. To find the policy π_θ , we minimize the expectation of the total path cost of all agents, which is defined as:

$$\arg \min_{\pi_\theta} \mathbb{E} \left[\max_{1 \leq i \leq m} |l_i| \mid \pi_\theta \right], \quad (2)$$

where $|l_i|$ is path cost of agent i .

²Formally, a POMDP can be represented as $\langle N, \mathcal{S}, \{\mathcal{A}_i\}_{i=1}^N, \{\mathcal{R}_i\}_{i=1}^N, \{\mathcal{O}_i\}_{i=1}^N, \mathcal{P}, \mathcal{Z}, \gamma \rangle$. N is the number of agents. \mathcal{S} is the state space containing information about agents and environments. \mathcal{A}_i represents the action space of agent i , and $\mathcal{A} = \mathcal{A}_1 \times \dots \times \mathcal{A}_n$ is the space of joint actions. $\mathcal{R}_i : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the reward function for agent i . $\{\mathcal{O}_i\}$ is the observation space of agent i . $\mathcal{P}(s'|s, a) : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ is the state transition function, which is used to describe the probability that agent in state s takes action a and then transits to state s' . $\mathcal{Z} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{O}_i$ is the observation function. $\gamma \in [0, 1]$ is the discount factor.

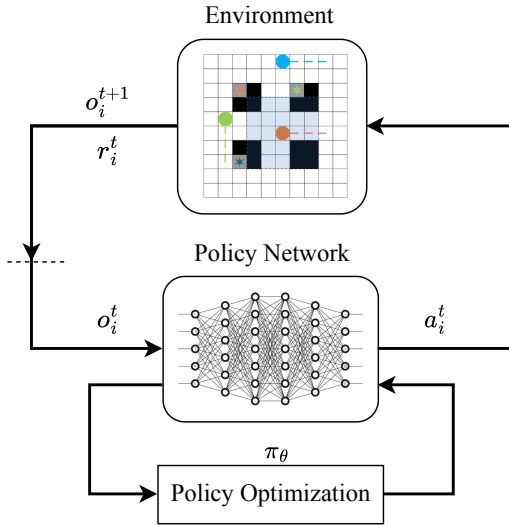


Fig. 2. An overview of RL-based MAPF policy. At each time step t , agent i receives its partial observation space o_i^t from the environment, which is denoted by $o_i^t \in \mathbb{R}^{W_{FOV} \times H_{FOV}}$, where W_{FOV} and H_{FOV} represent the width and height of the FOV. Then an action a_i^t is generated by the shared policy π_θ . When the agent reaches the new position in the next time step $t + 1$, it will receive a reward r_i^t and its new observation space o_i^{t+1} . We repeat the above process until the agent reaches the goal position.

IV. HYBRID MAPF POLICY

This section presents the hybrid MAPF policy, RDE, depicted in Fig. 3. RDE flawlessly combines the RL-based policy π_θ with the DHM-based policy π_H and the escape policy π_s .

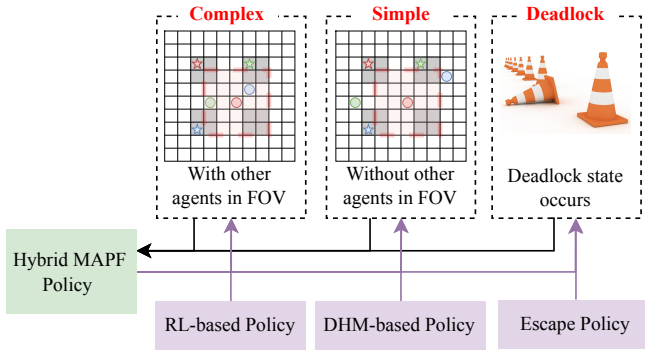


Fig. 3. Overview of RDE. The RL-based policy is utilized in complex scenarios \mathcal{T}_c , while the DHM-based policy is used for simple scenarios \mathcal{T}_s . The escape policy is employed in deadlock state scenarios \mathcal{T}_d .

A. Classification of Scenarios

1) *Complex*: When there are other agents in the field of view of the agent, the agent not only pays attention to the collision with obstacles, but also considers the cooperation with other agents. We call this scenario as *complex* \mathcal{T}_c .

2) *Simple*: When there are no other agents in the FOV of the agent, the agent only pays attention to the collision with obstacles. We call this scenario *simple* \mathcal{T}_s .

3) *Deadlock*: In warehouse automation, encountering deadlock states \mathcal{T}_d like stagnation and oscillation is possible due to the obstacles' high density and structured features [38]. These deadlock states can ultimately fail MAPF.

Stagnation. If the agent stays at a non-goal position for over n steps, we consider the agent in a state of stagnation. When other agents block the warehouse aisle, stagnation happens as the agent cannot find a path. We test different values for n and find that $n = 4$ can improve performance.

Oscillation. When a collision occurs between two agents, they may take the same action to avoid it, causing a collision again in the next step. If this process continues, neither agent can reach the goal position. We call this situation oscillation. As illustrated in Fig. 4, the green agent moves to the right grid first when time $t = 1$, followed by the orange agent. When time $t = 2$ arrives, both agents move to the left grid based on their partial observations. In the next time step $t = 3$, both agents return to the right grid. The single-step oscillation shown in Fig. 4 is the most common in the aisle of the warehouse and has the most significant impact on MAPF. Other larger oscillations occur rarely, so we do not consider them in this paper.

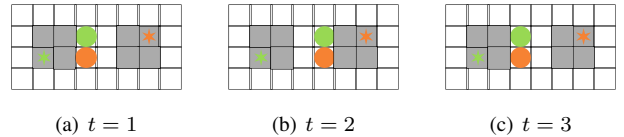


Fig. 4. An example of the agent in a state of oscillation.

B. RL-based Policy

RDE can be adapted to any RL-based MAPF policy. The state-of-the-art RL-based MAPF policies used in this paper are shown in Table I. DHC [14] and DCC [15] perform excellently than other MAPF policies and use DHM-based heuristic observation channels. We can directly leverage the already established DHM to build a DHM-based policy, which makes DHC and DCC ideal for combination with RDE.

TABLE I
COMPARISON OF RDE AND BASELINES.

Solver	Learning	Comm.	Guidance	Encoder	Train Env.
DHC [14]	Dueling DQN [39]	✓	Heuristic	CNN	Random
DCC [15]	Dueling DQN [39]	✓	Heuristic	CNN	Random

C. DHM-Based Policy

If no other agents exist in an agent's FOV, RL-based MAPF policy π_θ loses its inherent collaborative advantage. Additionally, the dense and structured obstacles in warehouse automation often cause agents to get stuck in deadlock using RL-based MAPF policies.

The DHM-based policy π_H is designed in response to the above problems. The policy uses the Dijkstra [16] algorithm to calculate the shortest distance from the goal position to all grid positions, resulting in a DHM with the goal position

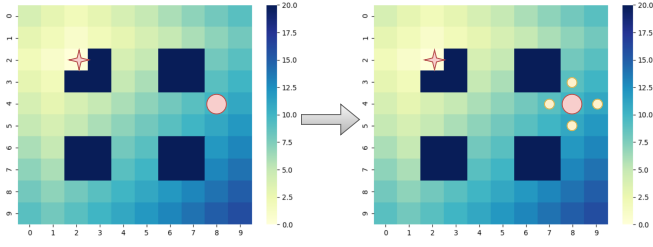


Fig. 5. The 10×10 DHM with a star representing the goal position and the dark blue square indicating the obstacles. On the DHM, the heat value of the obstacle position is infinite. If the goal position is the same, we can call the existing DHM directly.

as the source. As shown in Fig. 5, by querying the DHM, we can determine the shortest distances between the four grids adjacent to the agent’s current and goal positions. The position closest to the goal position is then selected as the agent’s action $a_{i,s}^t$.

Moreover, when the number of $a_{i,s}^t \in \mathcal{A}_{i,s}^t$ is more than one, we prioritize the action that allows the agent to go straight, which satisfies:

$$\overrightarrow{v_i^{t-1}v_i^t} = \overrightarrow{v_i^t v_i^{t+1}}, \quad (3)$$

where $\overrightarrow{v_i^{t-1}v_i^t}$ represents the current move direction of agent i , and $\overrightarrow{v_i^t v_i^{t+1}}$ represents the future move direction of agent i . If none of the actions $a_{i,s}^t \in \mathcal{A}_{i,s}^t$ can satisfy Eq. (3), the agent randomly selects an action from $\mathcal{A}_{i,s}^t$. As shown in Fig. 6, We try to go straight to reduce robot wear and actual time costs.

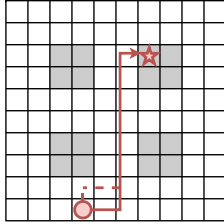


Fig. 6. Two paths with the same cost. Even though the dashed and solid paths have the same cost, robots that follow the dashed paths will take longer and cause more damage to their hardware in real-world situations. Therefore, reducing the number of turns is essential to optimize its performance and reduce wear and tear.

D. Escape Policy

When implementing the RL-based policy in warehouse automation, there is a higher probability for agent to encounter a deadlock state. We introduce an escape policy π_e to enhance the agent’s capability of escaping such a state. As shown in Algorithm 1, if the agent is stuck in a deadlock state while planning, we compel it to randomly choose an action from the available action space, thereby assisting it in breaking free from the deadlock states.

V. SIMULATION EXPERIMENT

In this section, we show the experiment results of RDE. We implemented the experiments with python.

Algorithm 1: The escape policy π_e for MAPF.

Input: Agents set $N = \{1, \dots, i, \dots, m\}$, goals set $V^g = \{v_1, \dots, v_i, \dots, v_m\}$, the position set of each agent for the past five time steps $V_i^p = \{v_i^{t-4}, v_i^{t-3}, v_i^{t-2}, v_i^{t-1}, v_i^t\}, \forall i \in N$, Action set $\mathcal{A}' = \{up, down, left, right\}$

Output: Actions list: $A = \{a_1^t, \dots, a_i^t, \dots, a_m^t\}$

- 1 **for** $i \leftarrow 1$ **to** m **do**
- 2 **if** $v_i^t \neq v_i^g$ **then**
- 3 **if** $v_i^{t-1} = v_i^{t-3}$ & $v_i^{t-2} = v_i^{t-4}$ **then**
- 4 $a_i^t \leftarrow \text{random}(\mathcal{A}')$;
- 5 **end**
- 6 **end**
- 7 **end**

Simulation is carried out on a single desktop computer with 32 GB memory, Intel® Core™ i7-9700 CPU @ 3.00GHz \times 8 processes and GeForce RTX 2060 SUPER/PCIe/SSE2 graphics. The single desktop computer has a Ubuntu 16.04 LTS system.

TABLE II
MAP SETS.

Size	Map Type	ρ_o^\dagger	Agent Density $\rho_a(10^{-2})^*$			
			10	30	50	70
34×34	Sparse	0.419	1.49	4.47	7.44	10.42
	Dense	0.476	1.65	4.95	8.25	11.56

† The density of obstacles ρ_o in the map refers to the percentage of static obstacles.

$*$ The agent density ρ_a is the proportion of all agents to the vacant positions on the map.

A. Setup for Testing

As shown in Fig. 7, we choose two kinds of warehouse-like structured maps for testing: *Sparse* and *Dense*. Table II shows the obstacle density of two kinds of map are both larger than 0.4, much higher than that of other studies [10]–[14] (up to 0.3). There are more path conflicts among agents. In the small-scale scenarios, the map size is 34×34 , and number of agents is 10, 30, 50, and 70, respectively. We randomly generate 1000 test instances with the same map size and the number of agents for the small-scale scenarios. We set the maximum time step size for the small-scale scenario to 150.

As mentioned in Section IV-B, we use state-of-the-art RL-based policies, DHC [14] and DCC [15], to verify the effectiveness of RDE. We directly use the policy network models of DHC and DCC trained on random maps without retraining on warehouse-like structured maps. During the test, we combined DHC, DCC with the DHM-based policy (DHC+DHM, DCC+DHM) and then combined DHC, DCC with the DHM-based and escape policies (DHC+DHM+Escape, DCC+DHM+Escape).

B. Metrics

Success rate (*SSR*). One MAPF instance is deemed unsuccessful when the time step surpasses the maximum value and not all agents reach their goals. We usually use

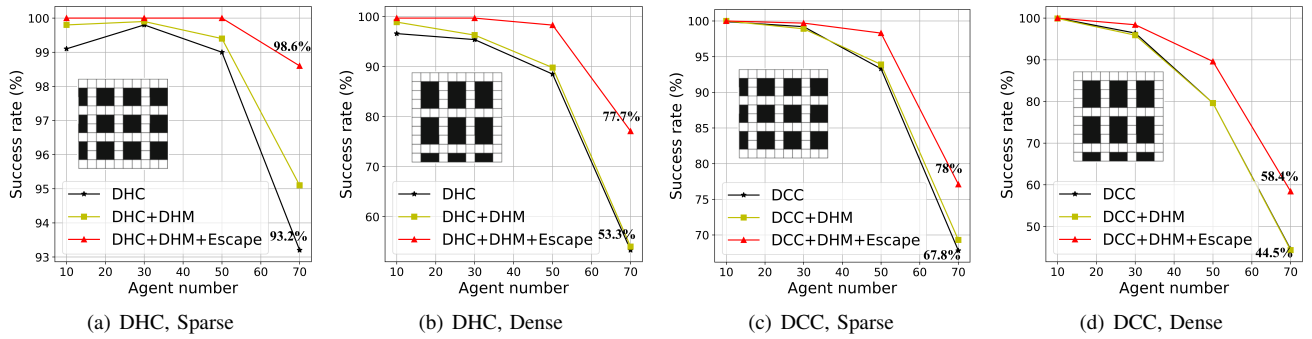


Fig. 7. The comparison of SSR for all policies in the small-scale scenarios.

the SSR of solving multiple MAPF instances to evaluate the performance of the MAPF policy. SSR can be defined as:

$$SSR = \frac{n_s}{n_t}, \quad (4)$$

where n_s is the number of successfully solved MAPF instances \mathcal{I}_s and n_t is the total number of instances \mathcal{I}_t .

C. Results

As shown in Fig. 7, RDE can effectively improve the success rate of DHC and DCC on two different warehouse-like structured grid maps. As the number of agents increases, the success rate of all policies decreases. In particular, when the number of agents exceeds 50, the success rates of DHC and DCC drop sharply. By combining DHC, DCC with DHM and escape policies, the sharp decline in success rate has been effectively alleviated. As shown in Fig. 7(a) and 7(b), when the number of agents is less than 30, the DHM-based policy can further improve the success rate of DHC. Similarly, by comparing Fig. 7(a) with 7(b), and Fig. 7(c) with 7(d), we can also find that the DHM-based policy improves the success rate of the RL-based policy more significantly on sparse maps than on dense maps. It is because when the agent density or obstacle density is low, the possibility of there being no agent in the agent’s FOV becomes higher. The DHM-based policy can help the agent take optimal actions to approach the target position. From Fig. 7, we can see that the addition of the escape policy has significantly improved the success rates of DHC and DCC on two different maps. Especially on dense maps, when the number of agents reaches 70, the success rates of both DHC and DCC increase by 15%. Deadlock is the key reason for the failure of MAPF solution, and escape policies based on random action selection can effectively help agents escape from deadlock.

VI. CONCLUSIONS AND FUTURE WORK

This paper introduces a hybrid MAPF policy framework called RDE, which combines RL-based policy with DHM-based and escape policy for warehouse automation. The RL-based policy can handle cooperative coordination between agents, while the DHM-based policy is suitable for situations where no other agents are within the FOV. In a deadlock, agents can use escape policy to escape from the deadlock. RDE can be adapted to any RL-based MAPF policy. This paper combines RDE with state-of-the-art RL-based MAPF

policies, DHC and DCC. Simulation results show that RDE can further improve the success rate of DHC and DCC. In RDE, we use an event-based heuristic method for policy switching. However, it may be challenging to cope with more complex scenarios. In the future, we will seek more intelligent ways to switch policies.

REFERENCES

- [1] P. R. Wurman, R. D’Andrea, and M. Mountz, “Coordinating hundreds of cooperative, autonomous vehicles in warehouses,” *AI Mag.*, vol. 29, no. 1, pp. 9–9, 2008.
- [2] J. Berger and N. Lo, “An innovative multi-agent search-and-rescue path planning approach,” *Computers & Operations Research*, vol. 53, pp. 24–31, 2015.
- [3] H. Ma, “Graph-based multi-robot path finding and planning,” *Current Robot. Reports*, vol. 3, pp. 77–84, 2022.
- [4] C. Ferner, G. Wagner, and H. Choset, “ODrM* optimal multirobot path planning in low dimensional search spaces,” in *Proc. IEEE Int. Conf. Robot. Autom.*, 2013, pp. 3854–3859.
- [5] G. Sharon, R. Stern, A. Felner, and N. R. Sturtevant, “Conflict-based search for optimal multi-agent pathfinding,” *Artif. Intell.*, vol. 219, pp. 40–66, 2015.
- [6] L. Cohen, M. Greco, H. Ma, C. Hernandez, A. Felner, T. S. Kumar, and S. Koenig, “Anytime focal search with applications,” in *Proc. Int. Joint Conf. Artif. Intell.*, 2018, pp. 1434–1441.
- [7] D. Silver, “Cooperative pathfinding,” in *Proc. AAAI Conf. Artif. Intell. Interact. Digit. Entertain.*, 2005, pp. 117–122.
- [8] M. Barer, G. Sharon, R. Stern, and A. Felner, “Suboptimal variants of the conflict-based search algorithm for the multi-agent pathfinding problem,” in *Proc. Int. Symp. Comb. Search*, 2014, pp. 19–27.
- [9] P. Surynek, “Towards optimal cooperative path planning in hard setups through satisfiability solving,” in *Pacific Rim Int. Conf. Artif. Intell.*, 2012, pp. 564–576.
- [10] G. Sartoretti, J. Kerr, Y. Shi, G. Wagner, T. S. Kumar, S. Koenig, and H. Choset, “PRIMAL: Pathfinding via reinforcement and imitation multi-agent learning,” *IEEE Robot. Autom. Lett.*, vol. 4, no. 3, pp. 2378–2385, 2019.
- [11] Z. Liu, B. Chen, H. Zhou, G. Koushik, M. Hebert, and D. Zhao, “MAPPER: Multi-agent path planning with evolutionary reinforcement learning in mixed dynamic environments,” in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.*, 2020, pp. 11 748–11 754.
- [12] B. Wang, Z. Liu, Q. Li, and A. Prorok, “Mobile robot path planning in dynamic environments through globally guided reinforcement learning,” *IEEE Robot. Autom. Lett.*, vol. 5, no. 4, pp. 6932–6939, 2020.
- [13] Q. Li, F. Gama, A. Ribeiro, and A. Prorok, “Graph neural networks for decentralized multi-robot path planning,” in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.*, 2020, pp. 11 785–11 792.
- [14] Z. Ma, Y. Luo, and H. Ma, “Distributed heuristic multi-agent path finding with communication,” in *Proc. IEEE Int. Conf. Robot. Autom.*, 2021, pp. 8699–8705.
- [15] Z. Ma, Y. Luo, and J. Pan, “Learning selective communication for multi-agent path finding,” *IEEE Robot. Autom. Lett.*, vol. 7, no. 2, pp. 1455–1462, 2021.
- [16] H. Wang, Y. Yu, and Q. Yuan, “Application of dijkstra algorithm in robot path-planning,” in *Proc. IEEE Int. Conf. Mech. Automat. Control Eng.*, 2011, pp. 1067–1069.

- [17] J. Yu and S. LaValle, "Structure and intractability of optimal multi-robot path planning on graphs," in *Proc. AAAI Conf. Artif. Intell.*, 2013, pp. 1443–1449.
- [18] G. Sharon, R. Stern, M. Goldenberg, and A. Felner, "The increasing cost tree search for optimal multi-agent pathfinding," *Artif. Intell.*, vol. 195, pp. 470–495, 2013.
- [19] J. Yu and S. M. LaValle, "Optimal multirobot path planning on graphs: Complete algorithms and effective heuristics," *IEEE Trans. Robot.*, vol. 32, no. 5, pp. 1163–1177, 2016.
- [20] G. Wagner and H. Choset, "Subdimensional expansion for multirobot path planning," *Artif. Intell.*, vol. 219, pp. 1–24, 2015.
- [21] F. Aljaloud and N. Sturtevant, "Finding bounded suboptimal multi-agent path planning solutions using increasing cost tree search," in *Int. Symp. Comb. Search, SoCS*, 2013, pp. 203–204.
- [22] M. Rahman, M. A. Alam, M. M. Islam, I. Rahman, M. M. Khan, and T. Iqbal, "An adaptive agent-specific sub-optimal bounding approach for multi-agent path finding," *IEEE Access*, vol. 10, pp. 22 226–22 237, 2022.
- [23] H. Ma, C. Tovey, G. Sharon, T. S. Kumar, and S. Koenig, "Multi-agent path finding with payload transfers and the package-exchange robot-routing problem," in *Proc. AAAI Conf. Artif. Intell.*, 2016, pp. 3166–3173.
- [24] K. Okumura, M. Machida, X. Défago, and Y. Tamura, "Priority inheritance with backtracking for iterative multi-agent path finding," in *Proc. Int. Joint Conf. Artif. Intell.*, 2019, pp. 535–542.
- [25] J. Li, Z. Chen, D. Harabor, P. J. Stuckey, and S. Koenig, "Anytime multi-agent path finding via large neighborhood search," in *Proc. Int. Joint Conf. Auton. Agents Multiagent Syst.*, 2021, pp. 1581–1583.
- [26] —, "Mapf-Ins2: Fast repairing for multi-agent path finding via large neighborhood search," in *Proc. AAAI Conf. Artif. Intell.*, 2022, pp. 10 256–10 265.
- [27] Z. Li, F. Liu, W. Yang, S. Peng, and J. Zhou, "A survey of convolutional neural networks: Analysis, applications, and prospects," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 33, no. 12, pp. 6999–7019, 2022.
- [28] L. Chen, Y. Wang, Z. Miao, Y. Mo, M. Feng, Z. Zhou, and H. Wang, "Transformer-based imitative reinforcement learning for multirobot path planning," *IEEE Trans. Industr. Inform.*, vol. 19, no. 10, pp. 10 233–10 243, 2023.
- [29] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," 2017.
- [30] Q. Li, W. Lin, Z. Liu, and A. Prorok, "Message-aware graph attention networks for large-scale multi-robot path planning," *IEEE Robot. Autom. Lett.*, vol. 6, no. 3, pp. 5533–5540, 2021.
- [31] H. Guan, Y. Gao, M. Zhao, Y. Yang, F. Deng, and T. L. Lam, "Ab-mapper: Attention and bicnet based multi-agent path planning for dynamic environment," in *Proc. IEEE/RSSJ Int. Conf. Intell. Robots Syst.*, 2022, pp. 13 799–13 806.
- [32] J. Li, Z. Chen, Y. Zheng, S.-H. Chan, D. Harabor, P. J. Stuckey, H. Ma, and S. Koenig, "Scalable rail planning and replanning: Winning the 2020 flatland challenge," in *Proc. Int. Conf. Automated Plan. Sched.*, 2021, pp. 477–485.
- [33] J. Li, K. Sun, H. Ma, A. Felner, T. S. Kumar, and S. Koenig, "Moving agents in formation in congested environments," in *Proc. Int. Joint Conf. Auton. Agents Multiagent Syst.*, 2020, pp. 726–734.
- [34] A. Jain, D. Ghose, and P. P. Menon, "Achieving a desired collective centroid by a formation of agents moving in a controllable force field," in *Indian Control Conference*, 2016, pp. 182–187.
- [35] C. Wilt and A. Botea, "Spatially distributed multiagent path planning," in *Proc. Int. Conf. Automated Plan. Sched., ICAPS*, 2014, pp. 332–340.
- [36] H. Zhang, M. Yao, Z. Liu, J. Li, L. Terr, S.-H. Chan, T. S. Kumar, and S. Koenig, "A hierarchical approach to multi-agent path finding," in *Int. Symp. Comb. Search, SoCS*, 2021, pp. 209–211.
- [37] M. L. Littman, "Markov games as a framework for multi-agent reinforcement learning," in *Machine learning proceedings*, 1994, pp. 157–163.
- [38] Y. Xu, Y. Li, Q. Liu, J. Gao, Y. Liu, and M. Chen, "Multi-agent pathfinding with local and global guidance," in *2021 IEEE International Conference on Networking, Sensing and Control (ICNSC)*, vol. 1, 2021, pp. 1–7.
- [39] Z. Wang, T. Schaul, M. Hessel, H. Hasselt, M. Lanctot, and N. Freitas, "Duelling network architectures for deep reinforcement learning," in *International conference on machine learning*, 2016, pp. 1995–2003.