

# Decentralized Multi-AGV Task Allocation based on Multi-Agent Reinforcement Learning with Information Potential Field Rewards

Mengyuan Li, Bin Guo\*, Jiangshan Zhang,  
Jiaqi Liu, Sicong Liu, Zhiwen Yu  
School of Computer Science  
Northwestern Polytechnical University  
Xi'an 710072, China  
guob@nwpu.edu.cn

Zhetao Li  
College of Computer Science John Hopcroft Center for Computer Science  
Xiangtan University  
Xiangtan 411105, China  
liztchina@hotmail.com

Liyao Xiang  
Shanghai Jiao Tong University  
Shanghai 200240, China  
xiangliyao08@sjtu.edu.cn

**Abstract**—Automated Guided Vehicles (AGVs) have been widely used for material handling in flexible shop floors. Each product requires various raw materials to complete the assembly in production process. AGVs are used to realize the automatic handling of raw materials in different locations. Efficient AGVs task allocation strategy can reduce transportation costs and improve distribution efficiency. However, the traditional centralized approaches make high demands on the control center’s computing power and real-time capability. In this paper, we present decentralized solutions to achieve flexible and self-organized AGVs task allocation. In particular, we propose two improved multi-agent reinforcement learning algorithms, MADDPG-IPF (Information Potential Field) and BiCNet-IPF, to realize the coordination among AGVs adapting to different scenarios. To address the reward-sparsity issue, we propose a reward shaping strategy based on information potential field, which provides stepwise rewards and implicitly guides the AGVs to different material targets. We conduct experiments under different settings (3 AGVs and 6 AGVs), and the experiment results indicate that, compared with baseline methods, our work obtains up to 47% task response improvement and 22% training iterations reduction.

**Index Terms**—Multi-agent reinforcement learning, AGVs, decentralized task allocation, information potential field

## I. INTRODUCTION

Driven by the recent advancements in industry 4.0 and industrial artificial intelligence, the use of autonomous systems in manufacturing enterprises has become inevitable [1], [2]. Automated Guided Vehicles (AGVs), as a type of flexible intelligent logistics equipment, have a great degree of freedom and play an essential role in flexibly transporting materials and products. AGVs have been hailed as one of the most promising technologies and have been implemented in a variety of shop floors and warehouse logistics operations for material supply [3], [4].

The multi-variety, small-batch, and customized production mode results in more logistics tasks and higher real-time demands. Using AGVs for cooperative transportation can significantly improve efficiency and cut expenses. How to make

multiple AGVs collaborate to perform material transportation tasks remains a significant topic in intelligent storage systems [5], [6]. The traditional approaches are mostly *centralized* control methods (Fig.1 (a)) and consider task assignment as a path planning problem for single or multiple robots [7], [8]. On one hand, it places extremely high demands on the control center’s computing power and real-time capability. On the other hand, the complexity and dynamic obstacles of the environment can impair the system’s stability and scalability. In comparison to centralized solutions, agent-level *decentralized* task allocation strategies (Fig.1 (b)) evenly distribute computing load and make advantage of agents’ autonomous decision-making ability.

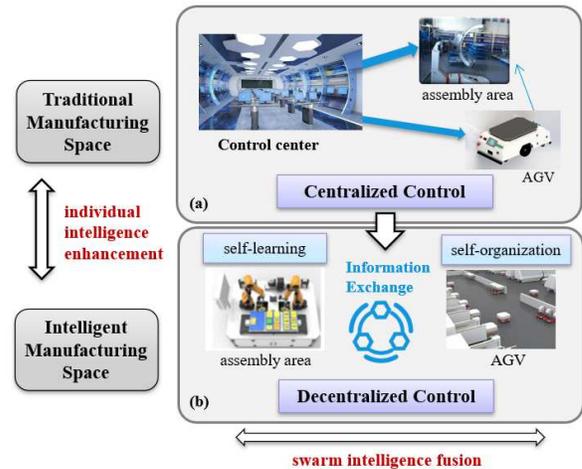


Fig. 1. Centralized control methods and decentralized control methods of AGVs.

With the continuous development of Multi-Agent Reinforcement Learning (MARL) [9], Reinforcement Learning (RL) has developed the capabilities of autonomous learning and distributed computing. Agents generate their own behaviors, modify their own state information, and accomplish the goal efficiently through cooperation with others [10]. For example,

\*Corresponding author.

Lowe et al. [11] propose the Multi-Agent Deep Deterministic Policy Gradient (MADDPG), which extends the DDPG method to MARL by observing the opponent’s behavior. Meanwhile, a global critic function is constructed to evaluate global state action. The Alibaba team proposes the Bidirectionally Coordinated Network (BiCNet) algorithm [12] in the ppsc2 multi-agent scenario [13]. Using Bidirectional Recurrent Neural Networks (BRNN) [14] for implicit communication, BiCNet has demonstrated superior performance in complicated environments.

However, existing MARL approaches have a number of drawbacks that make them unsuitable for decentralized multi-AGV task allocation directly, such as environmental non-stationarity and partial observability. Additionally, the reward mechanism in multi-agent system is more sophisticated than it is in single-agent system, and the reward-sparsity issue frequently makes training progress difficult to converge. A critical question is how to design an effective reward mechanism that will boost performance and expedite convergence. Information Potential Field (IPF) [15] is often utilized to tackle the path planning problem. Using the virtual information gradient diffusion of the target position data, the robot can advance to the target position along a specific gradient direction. By including IPF into reward function, the agents’ status can be assessed more comprehensively, guiding the agents toward the target positions.

To solve the above challenges, this paper proposes a novel multi-agent reinforcement learning algorithm based on information potential field rewards. We model the decentralized multi AGV task allocation as a Partially Observable Markov Decision Process (POMDP). To address reward-sparsity issue, we propose a reward shaping mechanism based on IPF that provides AGV collaboration with stepwise and implicit direction. Additionally, we apply IPF to the state-of-the-art MADDPG and BiCNet algorithms to prove the superiority of this mechanism. Extensive experiments demonstrate that our methodology can result in considerable performance and convergence improvements. The main contributions of this work are summarized as follows.

(1) The traditional centralized task allocation methods place extraordinarily high demands on control center’s computing power and real-time capability. We innovatively formulate the decentralized multi-AGV task allocation problem as a partially observable Markov decision process, and propose two improved multi-agent reinforcement learning algorithms to achieve coordination among AGVs adapting to different scenarios.

(2) We introduce information potential field to address the reward sparsity issue in decentralized multi-AGV task allocation. It can provide implicit direction for autonomous decision-making and improve the AGV system’s cooperation.

(3) We conduct experiments under different settings, and the experiment results show that our strategy obtains up to 47% task response improvement compared with baseline methods. Additionally, we demonstrate the cooperation mechanism of MADDPG-IPF and BiCNet-IPF. The agents establish a dif-

ferential preference for each target in MADDPG, while the agents prefer the closest target in BiCNet.

The paper is organized as follows: In Section II, we discuss related literature on collaborative task allocation and multi-agent reinforcement learning. Section III formulates the task allocation problem. In Section IV, we model the task allocation problem as a partially observable Markov decision process, and the proposed algorithm is demonstrated. The effectiveness of the method is verified by experiments in section V. Section VI gives the conclusions of this study and envisages some future work.

## II. RELATED WORKS

### A. Multi-AGV Task Allocation

Multi-AGV task allocation is a critical part of AGV control, as it seeks to determine the appropriate transit time and equipment for each task. The traditional AGVs task allocation approach is to apply classical optimization algorithms to the production scheduling field, such as genetic algorithm, particle swarm algorithm, ant colony algorithm. Wang et al. [16] optimize the path selection problem using an improved micro-genetic algorithm that takes into account running time, stopping time, and turning time. Zhang et al. [17] employ the makespan of jobs as the goal function and the machine and AGV utilization ratios as the comprehensive evaluation function. An improved particle swarm optimization algorithm is developed to solve a reasonable scheduling scheme. Liu et al. [18] develop a multi-objective mathematical model and integrate with two adaptive genetic algorithms to optimize the task scheduling of AGVs while taking into account the charging task and the AGV’s variable speed. Saidi et al. [19] address the conflict-free AGV path planning problem for job shop scheduling and solve it using a two-stage ant colony algorithm. These algorithms require knowledge of the global environment in order to calculate the optimal policies, and the decision-making capability of a single agent is insufficient in real-world scenarios. The multi-agent system can complete not only a single agent’s goal, but also exceed the efficiency of the single agent, which means that many agents can increase its strength.

### B. Multi-Agent Reinforcement Learning

In multi-agent system, traditional independent Q-learning [20] or DQN based on experience replay [21] cannot be applied to a multi-agent environment directly. Because the experience pool’s samples become old when the environment changes, the method produced from outdated sample training is frequently not ideal. Therefore, Foerster et al. [22] propose two strategies for maintaining the DQN experience replay pool’s stability. The central idea is to augment the experience buffer with additional information and to undertake importance sampling in order to mitigate the influence of unstable surroundings on multi-agent training. Lowe et al. [11] propose MADDPG to train a centralized critic for each agent using all agents’ policies during training in order to reduce variance by eliminating the non-stationarity. The

actor only has local information and the experience buffer records the experiences of all agents. Foerster et al. [23] propose an actor-critic counterfactual multi-agent (COMA) policy gradient method. COMA is intended for use in both the fully centralized and multiagent credit assignment problems. By comparing the current Q value to the counterfactual, an advantage function can be constructed. In contrast to previous approaches, in Bidirectionally Coordinated Network (BiCNet) [12], communication takes place in the latent space, and it also uses parameter sharing. Note that in BiCNet, agents do not explicitly share a message, it might be considered a method for learning cooperation.

Multi-agent reinforcement learning technology provides new ideas for implementing autonomous decision-making of multiple AGVs. Our proposed method utilizes the powerful data representation and decision-making capabilities of deep reinforcement learning to enable self-organizing task assignment of multi-AGV systems.

### C. Information Potential Field

Information Potential Field (IPF) is an effective path planning method. The robot can accomplish the global objective by employing a greedy strategy based on the information gradient. Liu et al. [15] propose two effective algorithms for constructing IPF: the hierarchical skeleton-based construction algorithm and the value estimation replacement algorithm, both of which achieve a trade-off between energy consumption and convergence speed. Wei et al. [24] propose efficient parking navigation via a continuous information ascent method. In the first step, a partial differential equation is used to establish a global potential field. In the second step, a Poisson equation is employed to construct the local potential field in the navigation process. Lin et al. [25] propose an artificial information gradient that is robust and has no local extrema. They use a harmonic function to establish IPF, representing the diffusion of a specific type of event of interest (EoI). Wei et al. [26] offer a novel heat diffusion equation to efficiently and quickly complete the navigation procedure. The strategy assures that a local information field is sufficiently large to encompass many appropriate targets, and that competition conflicts can be addressed concurrently. The majority of current research directly addresses the path planning problem using the information potential field method. In this paper, the information potential field is utilized to design the reward function of multi-agent reinforcement learning. The reward is evaluated in relation to the information potential value of the AGV location to implicitly steer the AGV to the target position.

## III. PROBLEM FORMULATION AND SYSTEM OVERVIEW

### A. Problem Formulation

In this section, we will formally define the multi-AGV collaborative task allocation problem. In the manufacturing workshop, processing products typically require various raw materials, which are stored in different locations across the warehouse. AGVs must travel to multiple destinations in

order to coordinate transportation tasks. We define the logistic network using  $G = (T, V, L)$ , where  $T$ ,  $V$  and  $L$  denote the set of material targets, vehicles and trajectories, respectively. More specifically,

**Target set  $T$ :** Each cooperative transportation task entails the movement of  $N$  different materials. The material targets  $T_i \in T (1 \leq i \leq N)$  are randomly dispersed in different places, and the position of the target  $T_i$  is represented by  $(x_i^T, y_i^T)$ .

**Vehicle set  $V$ :** we assume that all of  $N$  AGVs are modeled as discs with the same radius  $D$ , i.e., all AGVs are homogeneous. At each timestep  $t$ , utilize the vector  $G_i = \{p_i^t, v_i^t, r_i\}$  to describe the state of the AGV  $i (1 \leq i \leq N)$ , including its position  $p_i^t = (x, y)$ , velocity  $v_i^t = (v_x, v_y)$ , and sensing distance  $r_i$ . The AGV  $i$  obtains an observation  $o_i^t$  within the sensing range  $r_i$ , and then compute the action command  $a_i^t$  according to the policy  $\pi_\theta$ , where  $\theta$  denotes the policy parameters. The calculated action  $a_i^t$  is a velocity  $v_i^t$  that directs the AGV toward the task target while avoiding collisions with other robots.

**Trajectory set  $L$ :** To wrap up the preceding formulation, we define  $L = \{l_i, i = 1, \dots, N\}$  as the set of trajectories of all AGVs, which are subject to the AGV's kinematic constraints, i.e.:

$$\begin{aligned} v_i^t &\sim \pi_\theta(a_i^t | o_i^t) \\ \|v_i^t\| &\leq v_i^{max} \\ p_i^t &= p_i^{t-1} + \Delta t \cdot v_i^t \\ \forall j \in [1, N], j \neq i, &\|p_i^t - p_j^t\| > 2D \end{aligned} \quad (1)$$

To find an optimal policy, we set an objective by minimizing the expectation of the mean arrival time of all AGVs in the same scenario, which is defined as:

$$\operatorname{argmin}_{\pi_\theta} E \left[ \frac{1}{N} \sum_{i=1}^N t_i | \pi_\theta \right] \quad (2)$$

Where  $t_i$  is the travel time of the trajectory  $l_i$  in  $L$  controlled by policy  $\pi_\theta$ .

Decentralized multi-AGV task allocation can be viewed as a special mobile robot moving path planning problem. AGV decides its target and plans a collision-free course based on its surroundings cognition.

### B. System Architecture

We propose improved multi-agent reinforcement learning algorithms to solve this problem, the architecture of which is shown as Fig.2. In real world situations, agents make noisy observations of the true environment state to inform their action selection, typically modeled as a POMDP. Formally, a POMDP can be described as a tuple:  $M = (N, S, A, P, R, O)$ , where  $N$  denotes the number of agents,  $S$  represents the system state space,  $A$  represents the joint action space of all agents,  $P$  is the transition probability function,  $R$  is the reward function, and  $O$  is the observation probability distribution given the system state ( $o \sim O(s)$ ). Specific to the problem scenario of AGV collaborative task allocation, the state space  $S$  and action space  $A$  are specifically designed as follows:

State space  $S$ : For the AGV task assignment problem, the selection of the state space should not only characterize the attributes of the agents and targets, but also not bring too much computational burden. Therefore, we set the state space as  $\{v, p, D_A, D_B\}$ , where  $\{v, p\}$  is the speed and position of the agent itself, and  $\{D_A, D_B\}$  is the relative distance from the targets and other agents.

Action space  $A$ : We set the AGV's action space as a one-dimensional vector  $\{x, y\}$ , the value is  $(-1, 1)$ , representing the acceleration in the left and right directions and the front and back directions. Combined with the weight and damping of the AGV itself, the velocity of the AGV is computed.

Reward  $R$ : Our objective is each AGV avoids collisions and self-organizes to different targets as quickly as possible. A reward function is designed to guide a team of AGVs to achieve this objective. we design a target reward when reaching the target position and a collision penalty when a collision occurs.

When the new tasks arrive, state information is input to the network to determine the action. Following that, the chosen action will be used to route the AGVs to various task targets. The reward function is used to direct model training in this process, allowing the model to learn the ideal strategy.

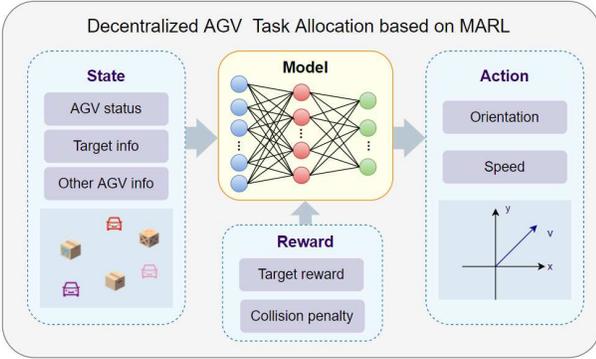


Fig. 2. Architecture of AGVs task allocation approach.

## IV. METHODS

### A. Reward Shaping with IPF

A well-designed reward function can enhance robustness and promote agent collaboration. In the previous section, we discuss a general AGV task allocation framework. In this section, we propose a reward shaping strategy based on information potential field to address the issue of reward sparsity.

Information Potential Field (IPF) is introduced to design the reward function  $r_{IPF}$ , as shown in Fig.3. We partition the scenario into a bounded grid map, assign a positive information potential value for the location of the target target, and assign a negative information potential value for the location of other AGVs, which can implicitly guide the AGVs to different targets. The targets are set to a maximum potential value of 5, while the other AGV's positions are set to a minimum potential

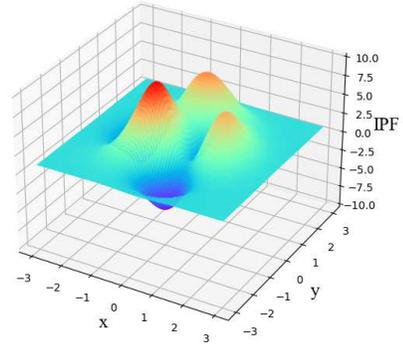


Fig. 3. Information Potential Field.

value of -3. Additionally, we set the information value of some other nodes to 0, often nodes on the network boundary, in order to enforce a gradient throughout the network. The remaining nodes compute the information potential field using Jacobi iterations. Each non-boundary node iterates:

$$\Phi^{k+1}(u) \leftarrow \frac{1}{d(u)} \sum_{v \in N(u)} \Phi^k(u) \quad (3)$$

Where  $\Phi^k(u)$  is the value of node  $u$  in the  $k$ -th iteration.  $N(u)$  signifies the set of  $u$ 's neighbors, while  $d(u)$  denotes the degree of  $u$ . Each position will have a corresponding information potential value after iteration. The AGV obtains the reward value  $r_{IPF}$  according to the information potential value of the position at the time step  $t$ . As illustrated in Fig.4, the IPF value around the target location is high, and the gravitational range grows more vast when several targets are gathered. When another AGV is already in close proximity to the target, the reward is reduced, essentially avoiding multiple AGVs competing for the same target.

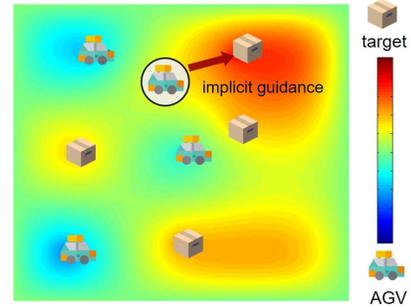


Fig. 4. IPF provides implicit guidance for agent's decision-making.

Along with  $r_{IPF}$  for implicit guidance, we design a target reward  $r_g$  and a collision penalty  $r_c$  for explicit guidance. The target reward  $r_g$  and the collision penalty  $r_c$  are specified as follows:

$$r_g = - \sum_i \min_j(d_{ij}) \quad (4)$$

$$r_c = \begin{cases} -1 & \text{if } \|p_i^t - p_j^t\| \leq 2R \\ 0 & \text{otherwise} \end{cases} \quad (5)$$

Where  $d_{ij}$  is the distance between task target  $j$  and AGV  $i$ . Additionally, when the AGV collides with other AGVs in the environment, it incurs a  $r_c$  penalty.

In general, we hope that when a new handling task arrives, the AGV system can self-organize and complete it in the shortest time possible. Based on the observed information, AGVs must plan a collision-free path to different material targets. We use the sum of  $r_{IPF}$ ,  $r_g$  and  $r_c$  to represent the reward  $r$  acquired by AGV  $i$  at time step  $t$ , as seen in (6), directing the AGV system to achieve self-organizing task assignment.  $r_g$  incentivizes the presence of precisely one agent near each target.  $r_c$  wishes for the fewest potential collisions.  $r_{IPF}$  provides an implicit shove to the AGV, guiding it to the target place in a distributed fashion.

$$r_i^t = (r_{IPF})_i^t + (r_g)_i^t + (r_c)_i^t \quad (6)$$

### B. The Algorithm Design

In multi-agent training, we focus on two algorithms based on the actor-critic framework, MADDPG and BiCNet. These two algorithms offer the following advantages over other MADL algorithms. MADDPG does not require explicit communication rules, is applicable to a wide variety of contexts, including cooperative, competitive, and mixed environments, and is capable of solving the non-stationary problem associated with multi-agent environments. All agents in BiCNet share models and parameters and build communication channels in the hidden layer, enabling any number of agents to cooperate. These two algorithms approach issues differently, and there are clear distinctions in the model structure, loss function, and other factors.

1) *MADDPG-IPF*: MADDPG [11] adopts centralized training with distributed execution method. Each agent trains a critic network that requires global information and an actor network that only requires local knowledge. The actor chooses the best action for a given state by optimizing the neural network parameters  $\theta$ . The critic evaluates the action generated by the actor by computing the temporal difference error. The MADDPG algorithm network structure is shown in Fig.5.

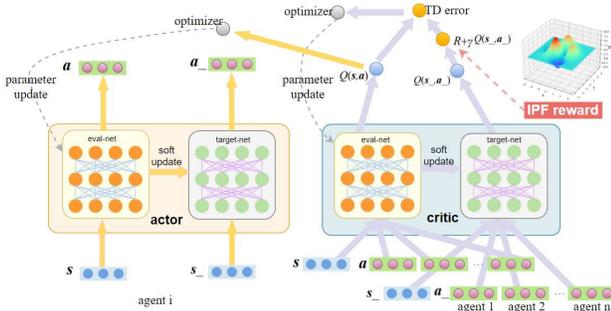


Fig. 5. The structure of MADDPG-IPF.

The policy gradient is calculated as:

$$\begin{aligned} \nabla_{\theta_i} J(\mu_i) &= E_{x,a \sim D} [\nabla_{\theta_i} \mu_i(a_i | o_i) \cdot \\ &\quad \nabla_{a_i} Q_i^\mu(x, a_1, \dots, a_n) |_{a_i = \mu_i(o_i)}] \end{aligned} \quad (7)$$

Among them,  $o_i$  represents the observation of the agent  $i$ , and  $x = [o_1, \dots, o_n]$  represents the observation vector.  $Q_i^\mu(x, a_1, \dots, a_n)$  represents the centralized state-action function of the agent  $i$ . The experience replay buffer  $D$  contains  $(x, x', a_1, \dots, a_n, r_1, \dots, r_n)$  these tuples, which acts as the knowledge base of the agent, storing the experience of all agents.

The action-value function  $Q_i^\mu$  is updated based on:

$$\begin{aligned} y &= r_i(s, a) + \lambda Q_i^{\mu'}(x', a'_1, \dots, a'_n) |_{a'_j = \mu'_j(o_j)} \\ L(\theta_i) &= E_{x,a,r,x'} [(Q_i^\mu(x, a_1, \dots, a_n) - y)^2] \end{aligned} \quad (8)$$

Among them,  $Q_i^{\mu'}$  represents the target network, and  $\mu' = [\mu'_1, \mu'_2, \dots, \mu'_n]$  is the parameter  $\theta'_j$  of the target network that has a lagging update.

2) *BiCNet-IPF*: BiCNet [12] is still based on the actor-critic framework, and the network structure as illustrated in Fig.6. The actor and the critic are both constructed using a bidirectional recurrent neural network. Through implicit communication, the actor shares observation and returns action for each agent. Each agent has the ability to retain its own internal state and communicate with other agents.

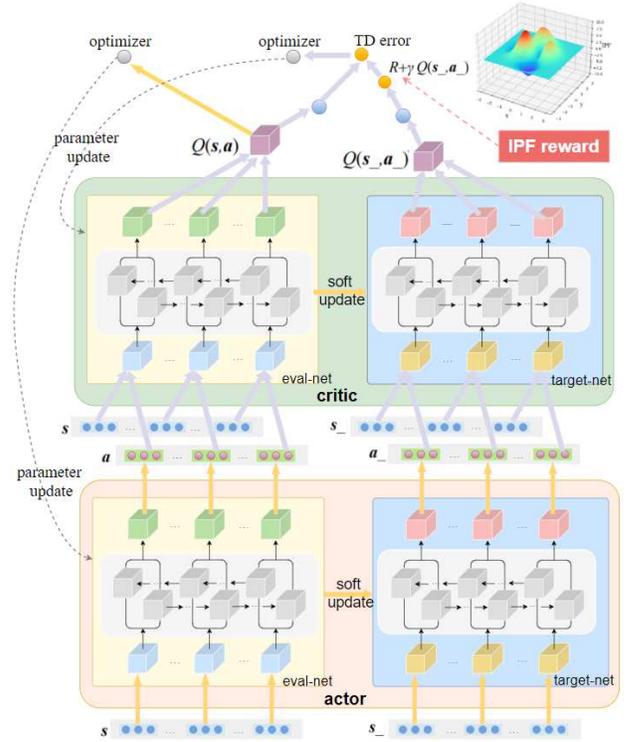


Fig. 6. The structure of BiCNet-IPF.

We denote the objective of a single agent  $i$  by  $J_i(\theta)$ , that is to maximize its expected cumulative individual reward  $r_i$  as  $J_i(\theta) = E_{s \sim \rho_{a_\theta}^\tau} [r_i(s, a_\theta(s))]$ . Therefore, we can get the objective of  $N$  agents denoted by  $J(\theta)$  as follows:

$$J(\theta) = E_{s \sim \rho_{a_\theta}^\tau} \left[ \sum_{i=1}^N r_i(s, a_\theta(s)) \right] \quad (9)$$

Combined with the deterministic policy gradient, we have the policy gradient as follows:

$$\nabla_{\theta} J(\theta) = E_{s \sim \rho_{a_{\theta}}^{\tau}(s)} \left[ \sum_{i=1}^N \sum_{j=1}^N \nabla_{\theta} a_{j,\theta} \cdot \nabla_{a_j} Q_i^{a_{\theta}}(s, a_{\theta}(s)) \right] \quad (10)$$

In training the critic network, using the sum of square loss, the gradient can be written as in (11), where  $\xi$  is the parameter of the Q-network:

$$\begin{aligned} \nabla_{\xi} L(\xi) = E_{s \sim \rho_{a_{\theta}}^{\tau}(s)} & \left[ \sum_{i=1}^N (r_i(s, a_{\theta}(s)) + \lambda Q_i^{\xi}(s', a_{\theta}(s'))) \right. \\ & \left. - Q_i^{\xi}(s, a_{\theta}(s)) \right] \cdot \nabla_{\partial \xi} Q_i^{\xi}(s, a_{\theta}(s)) \end{aligned} \quad (11)$$

In different agents, the parameters are shared, hence the number of parameters is independent of the number of agents. Parameter sharing leads to a compact model that speeds up the learning process.

## V. EVALUATION

### A. Experimental Settings

In order to conduct experiments, we build an AGV task allocation simulator based on a multi-agent environment [11], which comprises of  $N$  AGVs and  $N$  tasks inhabiting a two-dimensional world with continuous space and discrete time (see Fig.7). For MARL algorithms, as the number of agents increases, the joint state-action space increases exponentially, which makes the task intractable. Therefore we verify the robustness of the proposed methods under two scenarios: a 3 AGVs and 3 tasks simple scenario and a 6 AGVs and 6 tasks complex scenario, referred to as 3V3 scenario and 6V6 scenario. In each scenario, the position of the AGV and the position of the task are randomly generated. Taking into account the actual scenario, we define boundaries around the simulator, within which the agent can only move. We hope that the AGV can learn to disperse to different task targets in the shortest time and avoid collisions as much as feasible. Performance is measured by average task response rate, average reward, and average time:

*Average task response rate:* the number of tasks completed by  $N$  AGVs in the entire test epochs divided by the total number of tasks generated.

*Average reward:* the rewards obtained by  $N$  AGVs at each time step, calculated using the formula  $R = -\sum_i \min_j (d_{ij}) - C$ .

*Average time:* the total time required for  $N$  AGVs to execute all tasks (for example, in the 3V3 scenario, the three AGVs have reached the three task targets correctly).

### B. Performance Comparison

In this subsection, the performance of following methods is extensively evaluated by the simulation.

**MADDPG-MiniDist:** The MiniDist is a global reward that sums the distance between each task target and its nearest

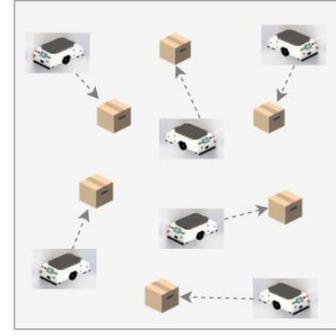


Fig. 7. Cooperative task allocation.

TABLE I  
MODEL PERFORMANCE IN 3V3 SCENARIO

	Average task response rate	Average reward	Average time
MADDPG-MiniDist	88.64%	-101.1	14.3
MADDPG-Greedy	88.67%	-116.61	11.8
MADDPG-IPF	95.00%	-85.8	11.1
BiCNet-MiniDist	93.03%	-71.5	10.4
BiCNet-Greedy	73.56%	-143.8	10.2
BiCNet-IPF	97.58%	-65.8	9.7

agent. The shorter the distance between two targets, the larger the reward.

**MADDPG-Greedy:** The Greedy is an individual reward. When an agent approaches the task target, it receives a positive reward, which rises as the distance between the agent and the task target decreases.

**MADDPG-IPF:** The IPF as we discussed in Section 4.

Additionally, **BiCNet-MiniDist**, **BiCNet-Greedy** and **BiCNet-IPF** are similar to the above. The Q-network and policy network in MADDPG are parameterized by three fully connected layers. The Q-network and policy network in BiCNet are based on the bi-directional RNN structure. Both the input and output modules are made up of four fully connected layers.

Each model is trained for 30k epochs in the 3V3 scenario. For the 6V6 scenario, the action space and state space dimensions are greatly increased, necessitating the use of additional rounds. As a result, each model based on MADDPG is trained for 50k epochs, and each model based on BiCNet, a more complicated network structure, is trained for 90k epochs. Finally, we execute 300 epochs for testing on each model in the two scenarios, and the results are presented in Table I and Table II.

MADDPG-IPF achieves a task response rate of 95% in the 3V3 scenario, an increase of approximately 6% over the other MADDPG models. Comparing the results of MADDPG-IPF and BiCNet-IPF, the BiCNet-IPF consistently outperforms MADDPG-IPF, possibly because of implicit communication, which enables better decision-making with more information. In the more complex 6V6 scenario, BiCNet-IPF achieves a task response rate of 91.61%, a significant advantage over all

TABLE II  
MODEL PERFORMANCE IN 6V6 SCENARIO

	Average task response rate	Average reward	Average time
MADDPG-MiniDist	69.22%	-438.5	17.7
MADDPG-Greedy	46.06%	-675.0	17.1
MADDPG-IPF	80.22%	-371.5	16.2
BiCNet-MiniDist	80.44%	-249.8	16.0
BiCNet-Greedy	44.56%	-664.5	17.5
BiCNet-IPF	91.61%	-241.1	15.6

other models. Although MADDPG-IPF is not as good as the best approach, it still achieves an 80.22% task response rate. In general, the global reward (MiniDist) assigns the same reward to all agents without regard of their contributions, which may encourage slothful agents. In comparison, the local reward (Greedy) only provides different local rewards to each agent based on individual behavior, leading to selfish agents. IPF reward incorporates global and local information and gives ongoing rewards at each step, allowing the agent to improve its performance on various task targets.

### C. The Effectiveness of IPF

Along with the performance comparisons mentioned above, we examine the task completion of each round of three AGVs under different reward designs in 3V3 scenario. As shown in the Fig.8, after applying the IPF reward mechanism, the agents can complete all tasks mostly in a distributed manner. IPF can significantly reduce the likelihood of multiple AGVs competing for the same target by offering implicit guidance. Global rewards may lead to laziness, so the agents inspired by MiniDist sometimes reach the target nearby but stagnate, resulting in worse task response than IPF. The Greedy reward frequently motivates agents to fight for a single task target, resulting in suboptimal performance.

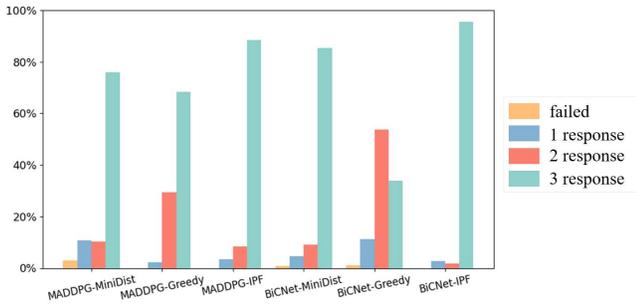


Fig. 8. Task response rate per round. For 3 response, all three tasks are completed.

Convergence is assessed by examining the average task completion rate of BiCNet during the training phase under the challenging 6V6 scenario. As illustrated in Fig.9, the approach using IPF can achieve a 40% task response rate after 40k epochs and 60% task response rate after 60k epochs. Due to the fact that BiCNet-Greedy is an individual reward network, its convergence rate is slower. The agents inspired by MiniDist

are unable to acquire vital knowledge in the first 60k epochs, but there performance improves significantly after 70k epochs.

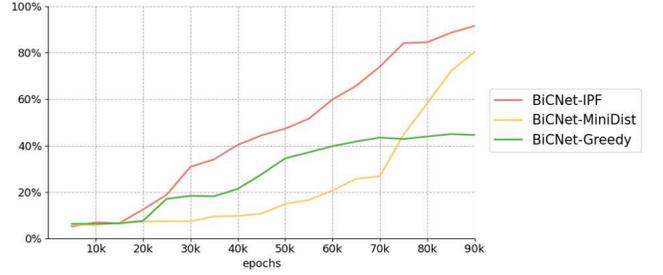


Fig. 9. Convergence comparison. Average task response rate under different reward mechanisms during the training phase.

### D. Implicit Cooperation Mechanism Analysis

In the 6V6 scenario, by numbering each AGV and each task, we observe an interesting phenomenon: the 2-th AGV and 6-th AGV directed by MADDPG always arrive at the identical 3-th task, resulting in no AGV reaching the 1-th task. However, this will not occur in BiCNet. Thus, we count the task targets achieved by each agent of MADDPG-IPF and BiCNet-IPF in the 3V3 scenario and 6V6 scenario, and investigate the cooperation mechanism of the two methods MADDPG and BiCNet, as shown in Fig. 10.

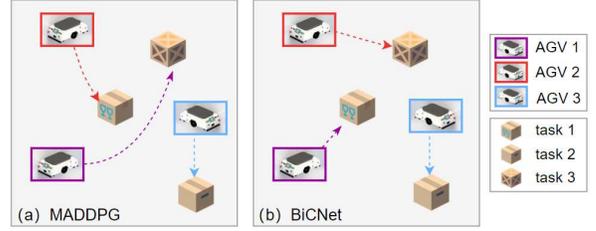


Fig. 10. The cooperation mechanism of MADDPG and BiCNet. In (a), the agents develop a differential preference for each task, e.g. the 1-th agent prefers task 1. In (b), the agents tend to complete the nearest task.

We discover that what MADDPG learned is each agent's preference for a certain fixed task. As illustrated in Fig.11, while training 30k epochs in the 3V3 scenario, 97% of the epochs of 1-th AGV chooses the 1-th task. What BiCNet learned is the choice of each agent for the closest task targets. As shown in Fig.12, the reach rate of 1-th AGV for the three tasks in 30k epochs is approximately 30%, and it does not show exceptional performance for a particular task.

In terms of this phenomenon, we argue that under MADDPG framework, each agent has an independent network structure and takes decisions based on local observations. Therefore, by continuously strengthening the rewards obtained at a particular task target during initial training, the agent will prefer it. While all agents in BiCNet share parameters and communicate implicitly via the bi-directional RNN, each agent coordinates with others and moves toward the nearest task target.

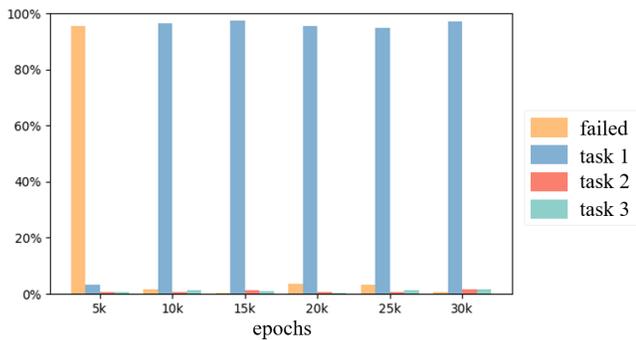


Fig. 11. The 1-th agent's preference in MADDPG. After fully training the model, the 1-th agent tends to complete task 1, but rarely chooses task 2 and task 3.

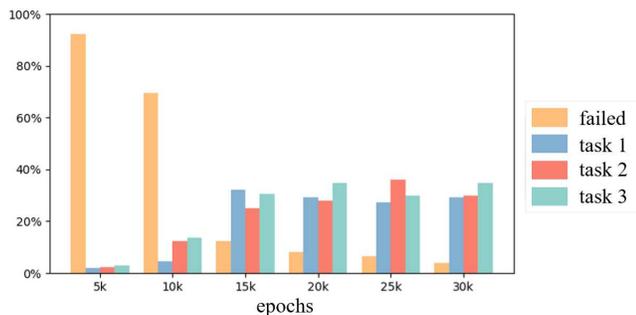


Fig. 12. The 1-th agent's preference in BiCNet. After fully training the model, the probabilities of the 1-th agent choosing three tasks are similar.

## VI. CONCLUSION

In this paper, we first formulated the AGVs task allocation problem in logistics networks as a partially observable Markov decision process. Given this setting, we introduced the information potential field optimization reward mechanism and proposed two cooperative multi-agent reinforcement learning algorithms to solve the problem. Extensive experiments demonstrate that our new approach can stimulate cooperation among agents and give rise to a significant improvement in both performance and convergence. For future work, we will create more multi-agent coordination and communications scenarios considering complex operation situations and uncertainties. Another interesting and practical direction to develop is to use a heterogeneous agent setting with individual specific feature to improve collaboration.

## ACKNOWLEDGMENT

This work was partially supported by the National Science Fund for Distinguished Young Scholars(62025205), National Key R&D Program of China(2019YFB1703901), and the National Natural Science Foundation of China (No. 62032020, 61960206008, 61725205).

## REFERENCES

[1] Zhang Y, Ma S, Yang H, et al. A big data driven analytical framework for energy-intensive manufacturing industries[J]. *Journal of Cleaner Production*, 2018, 197: 57-72.

[2] Wang J, Zhang Y, Liu Y, et al. Multiagent and bargaining-game-based real-time scheduling for internet of things-enabled flexible job shop[J]. *IEEE Internet of Things Journal*, 2018, 6(2): 2518-2531.

[3] Demesure G, Defoort M, Bekrar A, et al. Decentralized motion planning and scheduling of AGVs in an FMS[J]. *IEEE Transactions on Industrial Informatics*, 2017, 14(4): 1744-1752.

[4] Wang W, Zhang Y, Zhong R Y. A proactive material handling method for CPS enabled shop-floor[J]. *Robotics and Computer-Integrated Manufacturing*, 2020, 61: 101849.

[5] Khamis A, Hussein A, Elmogy A. Multi-robot task allocation: A review of the state-of-the-art[J]. *Cooperative Robots and Sensor Networks 2015*, 2015: 31-51.

[6] Wu G, Sun X. AGV Task Distribution Study[C]//*Journal of Physics: Conference Series*. IOP Publishing, 2020, 1486(7): 072016

[7] Zhu Z, Tang B, Yuan J. Multirobot task allocation based on an improved particle swarm optimization approach[J]. *International Journal of Advanced robotic systems*, 2017, 14(3): 1729881417710312.

[8] Mousavi M, Yap H J, Musa S N, et al. Multi-objective AGV scheduling in an FMS using a hybrid of genetic algorithm and particle swarm optimization[J]. *PloS one*, 2017, 12(3): e0169817.

[9] Oroojlooyjadid A, Hajinezhad D. A review of cooperative multi-agent deep reinforcement learning[J]. *arXiv preprint arXiv:1908.03963*, 2019.

[10] Nguyen T T, Nguyen N D, Nahavandi S. Deep reinforcement learning for multiagent systems: A review of challenges, solutions, and applications[J]. *IEEE transactions on cybernetics*, 2020, 50(9): 3826-3839.

[11] Lowe R, Wu Y, Tamar A, et al. Multi-agent actor-critic for mixed cooperative-competitive environments[J]. *arXiv preprint arXiv:1706.02275*, 2017.

[12] Peng P, Wen Y, Yang Y, et al. Multiagent bidirectionally-coordinated nets: Emergence of human-level coordination in learning to play starcraft combat games[J]. *arXiv preprint arXiv:1703.10069*, 2017.

[13] Vinyals O, Ewalds T, Bartunov S, et al. Starcraft ii: A new challenge for reinforcement learning[J]. *arXiv preprint arXiv:1708.04782*, 2017.

[14] Schuster M, Paliwal K K. Bidirectional recurrent neural networks[J]. *IEEE transactions on Signal Processing*, 1997, 45(11): 2673-2681.

[15] Liu S, Du J, Liu H, et al. Energy-efficient algorithm to construct the information potential field in WSNs[J]. *IEEE Sensors Journal*, 2017, 17(12): 3822-3831.

[16] Guo Z, Sun F. Research on integrated navigation method for AUV[J]. *Journal of Marine Science and Application*, 2005, 4(2): 34-38.

[17] Zhang F, Li J. An improved particle swarm optimization algorithm for integrated scheduling model in AGV-served manufacturing systems[J]. *Journal of Advanced Manufacturing Systems*, 2018, 17(03): 375-390.

[18] Liu S, Ji S, Su Z, et al. Multi-objective AGV scheduling in an automatic sorting system of an unmanned (intelligent) warehouse by using two adaptive genetic algorithms and a multi-adaptive genetic algorithm[J]. *PloS one*, 2019, 14(12): e0226161.

[19] Saidi-Mehrabad M, Dehnavi-Arani S, Evazabadian F, et al. An Ant Colony Algorithm (ACA) for solving the new integrated model of job shop scheduling and conflict-free routing of AGVs[J]. *Computers & Industrial Engineering*, 2015, 86: 2-13.

[20] Tan M. Multi-agent reinforcement learning: Independent vs. cooperative agents[C]//*Proceedings of the tenth international conference on machine learning*. 1993: 330-337.

[21] Mnih V, Kavukcuoglu K, Silver D, et al. Human-level control through deep reinforcement learning[J]. *nature*, 2015, 518(7540): 529-533.

[22] Foerster J, Nardelli N, Farquhar G, et al. Stabilising experience replay for deep multi-agent reinforcement learning[C]//*Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR. org, 2017: 1146-1155.

[23] Foerster J, Farquhar G, Afouras T, et al. Counterfactual multi-agent policy gradients[C]//*Proceedings of the AAAI Conference on Artificial Intelligence*. 2018, 32(1).

[24] Wei W, Song H, Li W, et al. Gradient-driven parking navigation using a continuous information potential field based on wireless sensor network[J]. *Information Sciences*, 2017, 408: 100-114.

[25] Lin H, Lu M, Milosavljevic N, et al. Composable information gradients in wireless sensor networks[C]//*2008 International Conference on Information Processing in Sensor Networks (ipsn 2008)*. IEEE, 2008: 121-132.

[26] Wei W, Qi Y. Information potential fields navigation in wireless Ad-Hoc sensor networks[J]. *Sensors*, 2011, 11(5): 4794-4807.