

# Mean Field Control with Envelope $Q$ -learning for Moving Decentralized Agents in Formation

Qiushi Lin and Hang Ma

**Abstract**—We study a decentralized version of Moving Agents in Formation (MAiF), a Multi-Agent Path-Finding variant aiming to plan collision-free paths for multiple agents so that they can reach their goals quickly while staying close to a desired formation. The agents must balance these two objectives under partial observation and limited communication. The maintenance of the formation is determined by the joint state of all agents, whose dimension increases exponentially w.r.t. the number of agents, making the learning process intractable. Furthermore, learning a single policy that can handle different linear preferences for these two objectives adds to the challenge. In this paper, we propose Mean Field Control with Envelope  $Q$ -learning (MFC-EQ), which provides a scalable and adaptable learning framework for this bi-objective multi-agent learning problem. We approximate the dynamic of all agents using the mean field theory while learning a universal preference-agnostic policy with envelope  $Q$ -learning. We empirically evaluate MFC-EQ over numerous instances and demonstrate that it can outperform state-of-the-art centralized MAiF baselines. Furthermore, MFC-EQ tackles more complex scenarios where the desired formation changes dynamically—a challenge that existing MAiF planners cannot handle.

## I. INTRODUCTION

Multi-Agent Path Finding (MAPF) [1] is a widely used technique in various multi-agent systems to find collision-free paths for agents in a shared environment. Applications include warehouse management [2], airport surface operations [3], video games [4], and other multi-agent systems [5]. Additionally, many of these applications require agents to adhere closely to a designated formation to accomplish collaborative tasks or maintain an efficient communication network. For example, in warehouse logistics, multiple robots/vehicles are required to collaborate in transporting large objects. Maintaining a specific formation is critical to optimizing transport efficiency or ensuring reliable communication. Moreover, in video gaming or military strategy simulations, game characters or army personnel must move in formations to safeguard vulnerable members.

To tackle this challenge, [6] has formalized the bi-objective problem of Moving Agents in Formation (MAiF) that combines these two tasks and proposed a centralized MAiF planner based on the leader-follower scheme and a search-based MAPF algorithm. However, existing MAiF planners work only in a centralized setting and do not apply to practical scenarios where agents do not fully observe the environment. Furthermore, centralized MAiF planners suffer from a huge computational burden as the number

of agents increases and are thus not suitable for planning in real time. Additionally, the only scalable MAiF planner, SWARM-MAPF [6], does not have the flexibility to adjust to the particular preferences between two objectives since it balances the two objectives only by setting the makespan allowance between two sets of heuristically determined waypoints, thus not guaranteed to optimize targeted preference. We propose a novel approach to learning a general MAiF solver for decentralized settings that can directly adapt to various preferences of the two objectives.

In the MAPF literature, reinforcement learning and imitation learning [7] have been introduced to solve MAPF in decentralized settings [8], [9], [10]. However, most learning-based MAPF solvers learn one homogeneous policy for any set of agents that treats nearby agents as part of the environment. This learning scheme does not translate seamlessly to decentralized MAiF. Unlike MAPF where the joint action cost can be directly decomposed to action costs of individual agents, the formation in MAiF is determined by the joint state of all agents at any given time. Each agent is thus required to not only avoid colliding with other agents but coordinate with them to maintain proximity to the desired formation. The dimension of the joint state space grows exponentially with the number of agents, incapacitating the scalability. Besides, trading off two objectives merely under partial observation and limited communication make this task even more difficult.

In this paper, we formalize the decentralized MAiF as a bi-objective multi-agent reinforcement learning task. The major contributions of our paper are as follows. We design a practical learning formalization for MAiF, including specifications for observations, actions, rewards, and inter-agent communication. To address the aforementioned challenges of MAiF, we propose a novel approach called **MEAN FIELD CONTROL WITH ENVELOPE  $Q$ -LEARNING** (MFC-EQ), a multi-agent reinforcement learning technique that optimizes towards any linear combination of two objectives for any number of agents, ensuring a stable and efficient learning process. MFC-EQ leverages mean field control to approximate the collective dynamics of the agents, treating the interaction of each agent within the formation as influenced by the collective effect of others. This design choice facilitates seamless scalability to large-scale instances. Furthermore, MFC-EQ extends envelope  $Q$ -learning to a multi-agent setting, enabling the learning of a universal preference-agnostic model adaptable to any linear combinations of the two objectives. To evaluate our method empirically, we extensively test MFC-EQ across various MAiF instances. Our results substantiate

that MFC-EQ consistently produces solutions that dominate those generated by several centralized MAiF planners and scales up well to large numbers of agents without long planning time. Additionally, the learned policy of MFC-EQ can directly adapt to more challenging tasks, including dynamically changing desired formations, which proves to be difficult for centralized MAiF planners.

## II. PROBLEM DEFINITION

In this section, we first describe the standard MAiF formulation in a convenient terminology to better present our learning approach. We then discuss how MAiF can be generalized to a partially observable environment, which is a more practical problem setting. Finally, we define relevant concepts and discuss the bi-objective optimization problem.

### A. Moving Agents in Formation

In the standard formulation, an MAiF instance is defined on an undirected graph  $G = (V, E)$  in a  $d$ -dimensional Cartesian system. Each location in  $V$  can be recognized by its coordinates  $\mathbf{v} = (v_1, \dots, v_d) \in \mathbb{R}^d$ . In this paper, the subscripts represent agents' index numbers and the boldface font denotes multi-dimensional vectors. We also define  $[M] = \{1, \dots, M\}$ . We have a set of  $M$  agents  $I = \{a^i | i \in [M]\}$ . Each agent has a unique start location  $\mathbf{s}^i \in V$  and goal location  $\mathbf{g}^i \in V$ . The time is discretized, and between two consecutive time steps, each agent can choose to wait at the current location or move from  $\mathbf{v}$  to  $\mathbf{v}'$  provided that  $(\mathbf{v}, \mathbf{v}') \in E$ . We also consider two types of collision between agents: a vertex collision  $\langle a^i, a^j, \mathbf{v}, t \rangle$  means agent  $a^i$  and  $a^j$  occupy the same location at the same time step  $t$ , and an edge collision  $\langle a^i, a^j, \mathbf{u}, \mathbf{v}, t \rangle$  happens when  $a^i$  travels from  $\mathbf{u}$  to  $\mathbf{v}$  while  $a^j$  travels backward.

The MAiF problem aims to find a set of  $M$  collision-free paths  $\Pi = \{\Pi^i | i \in [M]\}$  as a solution, where  $\Pi^i = (p_0^i, \dots, p_{T^i}^i)$  represents agent  $i$ 's trajectory. Every solution will be evaluated by two objectives, makespan and formation deviation. The makespan can be defined as  $T = \max_{1 \leq i \leq M} T^i$ , that is, the longest length among all paths. The *formation* at time  $t$  can be represented as an  $M$ -tuple,  $\ell(t) = \langle \mathbf{p}^1(t), \dots, \mathbf{p}^M(t) \rangle$ . The desired formation is the combination of all agents' goal locations,  $\ell_g = \langle \mathbf{g}^1, \dots, \mathbf{g}^M \rangle$ . Following the definition in [6], the *formation distance* between any two formation  $\ell = \langle \mathbf{u}^1, \dots, \mathbf{u}^M \rangle$  and  $\ell' = \langle \mathbf{v}^1, \dots, \mathbf{v}^M \rangle$  indicates the least effort required to transform from  $\ell$  to  $\ell'$ , defined as:

$$\mathcal{F}(\ell, \ell') := \min_{\Delta} \sum_{i=1}^M \|\mathbf{u}^i - (\mathbf{v}^i + \Delta)\|_1. \quad (1)$$

It can be easily proven that this definition is equivalent to:

$$\mathcal{F}(\ell, \ell') = \sum_{i=1}^M \sum_{j=1}^d |(\mathbf{u}_j^i - \mathbf{v}_j^i) - \Delta_j|, \quad (2)$$

where  $j$  indexes the dimension for all the position vectors and  $\Delta_j = \text{median}(\{\mathbf{u}_j^i - \mathbf{v}_j^i\}_{i \in [M]})$  is the median of differences for the  $j$ -th dimension. We can further decompose the

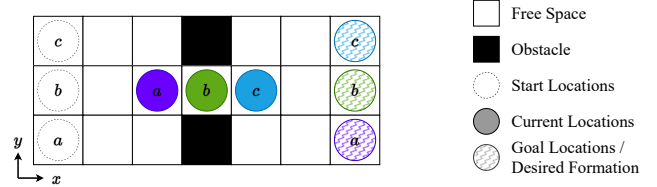


Fig. 1: Example of moving agents in formation.

formation deviation to each agent and denote  $\mathcal{F}^i(\ell, \ell') := \sum_{j=1}^d |(\mathbf{u}_j^i - \mathbf{v}_j^i) - \Delta_j|$  as the subpart only dedicated to the agent  $a^i$ . For the collective objective, unlike the total formation deviation in [6], we consider the average formation deviation per agent across all time steps, and it can be defined as  $\mathcal{F}_{avg} = \frac{1}{M} \sum_{t=0}^T \mathcal{F}(t)$ , where  $\mathcal{F}(t) = \mathcal{F}(\ell(t), \ell_g)$ . We also consider a mix of these two objectives:

$$\text{MIX}(\lambda) = \lambda T + (1 - \lambda) \mathcal{F}_{avg}, \quad (3)$$

which is the linear combination of these two objectives. There could certainly be non-linear combinations of different rewards, but, in this work, we only consider the linear cases, since it has been widely adopted in multi-objective or multi-task reinforcement learning (e.g., [11]).

**Example** A simple MAiF example is demonstrated in Fig. 1. The position vectors follow the order of agent  $a$ ,  $b$ , and  $c$ . The start formation is  $\langle (1, 1), (1, 2), (1, 3) \rangle$  and the goal/desired formation is  $\mathbf{u} = \langle (7, 1), (7, 2), (7, 3) \rangle$ . The group of agents cannot go through the fourth column while keeping the formation intact, so they have to change the formation. At the time step  $t$ , the position of agents is  $\mathbf{v} = \langle (3, 2), (4, 2), (5, 2) \rangle$ . We first derive the median of the differences.

$$\begin{cases} \mathbf{u}^a - \mathbf{v}^a = (4, -1) \\ \mathbf{u}^b - \mathbf{v}^b = (3, 0) \\ \mathbf{u}^c - \mathbf{v}^c = (2, 1) \end{cases} \implies \begin{cases} \Delta_x = 3 \\ \Delta_y = 0 \end{cases}$$

Therefore, the formation deviation can be calculated as  $\mathcal{F}(t) = \mathcal{F}^a(t) + \mathcal{F}^b(t) + \mathcal{F}^c(t) = 2 + 0 + 2 = 4$ .

### B. Partially Observable Environments

In this paper, we consider a more practical problem setting where, instead of assuming the full knowledge of the environment, each agent can only have a partial observation of its surroundings. We formulate decentralized MAiF as a decentralized partially observable Markov Decision Process (Dec-POMDP) [12]. A Dec-POMDP can be represented as a 7-tuple  $\langle \mathcal{S}, \mathcal{A}, P_S, \mathcal{O}, P_O, R, \gamma \rangle$ , where  $\mathcal{S}$  is the global state space.  $\mathcal{A} = \prod_{i=1}^M \mathcal{A}^i$  and  $\mathcal{O} = \prod_{i=1}^M \mathcal{O}^i$ , where  $\mathcal{A}^i$  and  $\mathcal{O}^i$  are agent  $i$ 's action and observation space.  $P_S : \mathcal{A} \times \mathcal{S} \rightarrow \mathcal{S}$  describes the state-transition function, and  $P_O : \mathcal{A} \times \mathcal{S} \rightarrow \mathcal{O}$  is the observation-transition function.  $R$  is the reward function with the discount factor  $\gamma$ .

In decentralized MAiF, we assume the observation and the state-transition function are deterministic, in which each agent has full control of its next position and observation by taking a move or the wait action. Following the settings in

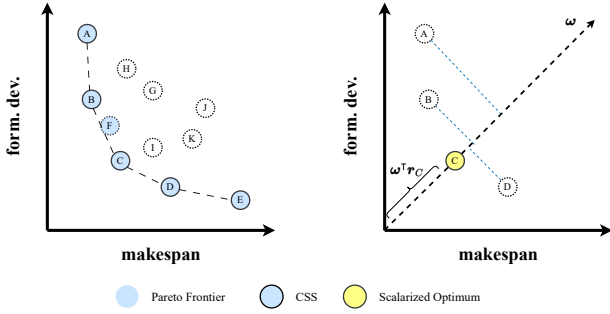


Fig. 2: Illustration of bi-objective optimums.

existing learning-based MAPF methods, we formalize this problem on a 2-dimensional 4-neighbor grid, even though our method can also be easily generalized to other settings. The partial observability limits each agent’s perception to a  $\mathcal{L} \times \mathcal{L}$  square area with it sitting on the center, defined as its FOV. Each agent can take actions merely based on their local observation and limited communication with other agents.

### C. Bi-Objective Optimization

We then formulate the goal of this bi-objective optimization problem. Each MAiF solution is evaluated as  $(v, w)$ , where  $v$  denotes its makespan and  $w$  denotes its average formation deviation per agent. We first define dominance. We say  $\mathbf{r} = (v, w)$  dominates  $\mathbf{r}' = (v', w')$ , denoted as  $\mathbf{r} \preceq \mathbf{r}'$ , iff  $v \leq v'$  and  $w \leq w'$ . A solution is Pareto-optimal if and only if there does not exist any solution that can dominate it. The Pareto-optimal frontier is a set of all Pareto-optimal solutions. In the MAiF setting, we are also interested in evaluating each solution  $\mathbf{r}$  by a scalar function  $f_{\omega}(\mathbf{r}) = \omega^{\top} \mathbf{r}$ , where  $\omega \in \Omega$  is the linear preference and  $\Omega$  is the set of all possible preferences. We let  $\omega = (\lambda, 1 - \lambda)^{\top}$  where  $0 \leq \lambda < 1$ . Our goal is to find the convex convergence set (CSS). The CSS is a subset of the Pareto-optimal frontier, where for each solution in CSS, there exists a preference  $\omega$  such that it minimizes  $f_{\omega}$  among all possible solutions. Intuitively, as shown in Fig. 2, we can regard the scalar function as a projection to the preference  $\omega$ . For example, the solution  $C$  belongs to CSS since it has the smallest projection into  $\omega$  compared to others.

## III. RELATED WORK

We now discuss related work on mean field reinforcement learning and multi-objective reinforcement learning.

### A. Mean Field Reinforcement Learning

Inspired by the mean field theory [13] from the physics world, the mean field reinforcement learning has been proposed in [14] which estimates the dynamic within the entire group of agents as the interaction between each agent and the mean effect of all other agents as a whole. As the dimension of the mean effect is independent of the number of agents, this method does not suffer from the curse of dimensionality, providing a general framework for large-scale multi-agent tasks. This method has been extended to the

partially observable stochastic settings [15], which utilizes certain distributions to sample agents’ actions without the necessity of observing them. The sampling process only serves stochastic games, which does not apply to our task. This mean field framework has also been used to solve multi-type multi-agent tasks [16], where agents are categorized into different types, and a set of mean effects is considered to reflect various types of agents.

### B. Multi-Objective Reinforcement Learning

There exist three major categories of multi-objective reinforcement learning methods. Single-policy methods [17], [18] convert the multi-objective problem into a single-objective optimization by using linear or non-linear functions, but these methods cannot manage unknown preferences. Multi-policy methods [19], [20], [21] update on a set of policies to approximate the real Pareto-optimal frontiers, which requires immense computational resources. These methods are only applicable to problems with limited state and action space. The policy-adaptation methods either train a meta-policy that adapts to different preferences on the fly [22] or learn a policy that conditions on different preference weights [23], [24], [25]. The Envelop  $Q$ -learning [25] has been proposed to increase sample efficiency by introducing a novel envelop operator for updating the multi-objective  $Q$ -function, which has become a standard way to tackle multi-objective problems with linear preferences.

## IV. MFC-EQ

In this section, we show how we design the learning framework for decentralized MAiF. We first design the learning environment with agents’ observation, communication, action, and reward functions. Then, we elaborate on the bi-objective multi-agent learning process based on the mean field theory and the envelope  $Q$ -learning.

### A. Environment and Model Design

1) *Observation*: As most research in the MAPF community [1], we study our problem in the 2-dimensional 4-neighbor grids. To mimic many real-world robotics applications where robots have limited visibility and sense range, each agent, in our grid world, can only observe its field of view (FOV), represented by its surrounding  $\mathcal{L} \times \mathcal{L}$  area. Each agent’s observation is represented by 3-channel feature maps  $\mathcal{F} \in \mathbb{R}^{\mathcal{L} \times \mathcal{L} \times 3}$ . The first two channels indicate obstacles and other neighboring agents’ positions. Inspired by some decentralized MAPF solvers [9], [10], the third channel encompasses the heuristic information where each grid in the FOV is assigned a value proportional to the short-path distance from that to the agent’s goal.

2) *Action*: In 4-neighbor grids, agents can only travel to their cardinally adjacent grids for each step. The action taken by agent  $i$  at time  $t$ , denoted by  $a_t^i \in \mathbb{R}^5$ , is a 5-dimensional one-hot vector with each dimension representing one action from  $\{up, down, left, right, wait\}$ . The first four actions take agents to another location and their observation will shift accordingly. The last action is to have the agent wait at

its current location, and it is especially crucial for formation control as one may have the choice for other agents to catch up for lower formation deviation.

3) *Multi-Agent Communication*: To keep the desired formation, agents not only need to communicate with nearby agents inside FOVs but also have to reach agents outside them. We specifically design the communication message so that it can pass along critical information under low communication bandwidth.

As in many real-world robot applications, each agent can only access the pairwise relative positions between other agents and itself. Assume that the current formation at time step  $t$  is  $\ell_p = \langle p^1, \dots, p^M \rangle$  and the desired formation is  $\ell_g = \langle g^1, \dots, g^M \rangle$ . We define the relative position between agent  $i$  and  $j$  as  $p^{i,j} = p^j - p^i$  (resp.,  $g^{i,j}$ ). Agent  $i$  receives  $\{p^{i,j}\}_{j \in [M]}$  in real time, and it holds the information of the relative positions in the goal formation,  $\{g^{i,j}\}_{j \in [M]}$ , which can be calculated before execution. We show that, only with this information, even without knowing the agent's whereabouts, it can still calculate the formation deviation. As defined in Eq. (2),  $\mathcal{F}(\ell_p, \ell_g) = \min_{\Delta} \sum_{m=1}^M \|p^m - (g^m + \Delta)\|_1 = \sum_{m=1}^M \sum_{n=1}^d |(p_n^m - g_n^m) - \Delta_n|$  where  $\Delta_n$  is the median of  $\{p_n^m - g_n^m\}_{m \in [M]}$ . Recall that  $d$  is the dimension of agents' coordinates. It is easy to verify that  $\mathcal{F}(\ell_p, \ell_g)$  is also equal to  $\sum_{m=1}^M \sum_{n=1}^d |(p_n^m - g_n^m - C_n) - \Delta'_n|$  where  $C$  is any constant  $d$ -dimensional vector and  $\Delta'_n$  is the median of  $\{p_n^m - g_n^m - C_n\}_{m \in [M]}$ , as all the values and the median are shifted by the same margin. Let  $C = p^i - g^i$ . We can rewrite the definition of the formation deviation only using the relative position.  $\sum_{m=1}^M \sum_{n=1}^d |(p_n^m - g_n^m) - \Delta_n| = \sum_{m=1}^M \sum_{n=1}^d |(p_n^m - p_n^i) - (g_n^m - g_n^i) - \Delta_n^*| = \sum_{m=1}^M \sum_{n=1}^d |(p_n^{i,m} - g_n^{i,m} - \Delta_n^*)|$  where  $i$  is the index of the observing agent and  $\Delta_n^*$  is the median of  $\{p_n^{i,m} - g_n^{i,m}\}_{m \in [M]}$ . Therefore, agent  $i$  can calculate the formation deviation merely based on the relative positions, which happens at each decentralized agent during execution with the time complexity of only  $\mathcal{O}(d \cdot M)$ .

We also mentioned that we can infer the mean action based on the relative positions. Given  $p^{i,j}(t)$  and  $p^{i,j}(t+1)$ , we can get that  $p^{i,j}(t+1) - p^{i,j}(t) = (p^j(t+1) - p^i(t)) - (p^i(t+1) - p^i(t)) = a^j(t) - a^i(t)$ . Hence, agent  $i$  can calculate agent  $j$ 's action by  $a^j(t) = p^{i,j}(t+1) - p^{i,j}(t) + a^i(t)$ . Therefore, we pass this to a linear layer to get the relative position encoding. Besides, each agent can infer other agents' actions by simply comparing relative positions in two consecutive time steps, which will later be used to compute the mean action.

4) *Reward*: The reward function for agent  $i$  after taking action  $a$  at time step  $t$ ,  $r_t^i(s_t^i, a_t^i) \in \mathbb{R}^2$ , is represented by a 2-tuple. The first element is designated for the makespan. We modify the individual cost function for makespan from DHC [10] which, instead, intends to minimize the sum of all path lengths (a.k.a., flowtime). The moving cost of agent  $i$  at time step  $t$  with action  $a^i$  is:

$$c_t^i(s^i, a^i) = \begin{cases} -0.075 & \text{collision-free } a^i \\ -0.5 & \text{collision (with obstacles or agents)} \\ 3 & \text{reach goal (first time)} \end{cases}$$

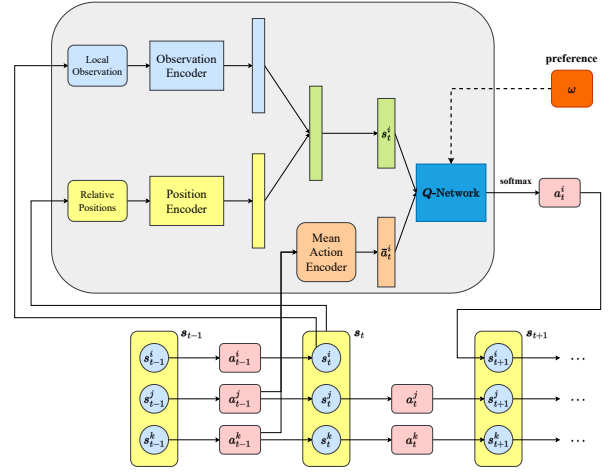


Fig. 3: Illustration of the model architecture of MFC-EQ. The bottom demonstrates the state/observation transition in the partially observable environment. The agent's  $Q$ -network gathers information from the environment through partial observation and limited communication and chooses the next action accordingly.

Each collision-free action, including move (up, down, left, or right) and wait (on goal or away goal), is slightly penalized so that agents are incentivized to approach their goals as quickly as they can. The second element is for the formation deviation. As defined in Eq. (2), we add the individual portion of the collective formation deviation that is dedicated to agent  $j$ , namely  $\mathcal{F}_t^j(\ell_t, \ell^g)$ . We negate the formation deviation so that it will be minimized through maximizing rewards. Hence, the reward function can be represented as:

$$r_t^j(s_t^j, a_t^j) = (c_t^j, -\mathcal{F}_t^j(\ell_t, \ell^g))^\top. \quad (4)$$

5) *Model Architecture*: Given the partially observable multi-agent environment, we further design the  $Q$ -network, whose learning algorithm will be introduced later. As in Fig. 3, we aim to project each agent's observations and communication messages into a corresponding action. Firstly, we feed the local observation and relative positions into two separate encoders. The observation encoder consists of several stacked convolution layers followed by linear layers. The relative position encoder includes two simple linear layers. Then, we concatenate these two encodings and forward them to another linear layer to obtain the final state representation,  $s_t^i$ , for agent  $i$ 's perception at time  $t$ . We then collect other agents' actions from the previous time step to calculate the mean action. Lastly, we use stacked linear layers to project them to the  $Q$ -values which condition on the state, the action, the mean action, and the given preference. The agent will decide its next action that maximizes the  $Q$ -function produced by the  $Q$ -network.

### B. Mean Field and Envelop Optimality

In the rest of this section, we will discuss the details of the learning algorithm. Learning multiple policy networks,  $\pi = [\pi^1, \dots, \pi^M]$ , for this bi-objective multi-agent task can be

extremely challenging. Therefore we simplify it by making some common assumptions.

1) *Mean Field Approximation*: The goal of MAiF is to minimize the makespan and the formation deviation. With the specifically designed reward function, the return, the discounted sum of all rewards from the initial joint state to the goal joint state,  $\sum_t \sum_j \gamma^t r_t^j(s_t, \mathbf{a}_t)$ , can reflect the actual values for the two objectives. Therefore, the goal of the learning is to find a set of policies to maximize the general sum of  $Q$ -values  $\arg \max_{\pi^1, \dots, \pi^M} \sum_{j=1}^M \omega^\top Q^{\pi^j}(s^j, \mathbf{a})$  with the given linear preference  $\omega$ . However, the dimension of  $s$  and  $\mathbf{a}$  grows exponentially w.r.t. the number of agents, rendering it infeasible to learn efficiently.

To tackle this problem, we introduce mean field reinforcement learning. We first lay out two common assumptions of homogeneity and locality that are made in [14] and many other multi-agent reinforcement learning works. The homogeneity assumes each agent shares the same policy, meaning that  $\pi^i = \pi^j$  for all  $i \neq j$ . The locality assumption comes from partial observability, which suggests that agents' actions can only depend on their visible surroundings.

Then, assuming the actions are represented by one-hot vectors, we define the mean action:

$$\bar{a}_t^j = \frac{1}{M} \sum_{j \in [M]} a_t^j, \quad a_t^j \sim \pi^j(\cdot | s^j, \bar{a}_{t-1}^j), \quad (5)$$

where  $\pi^j$  represents agent  $j$ 's policy. As all agents take action based on their observations and evaluate their values based on the joint action, it is infeasible to learn a  $Q$ -network for the joint action. To address this problem, we factorize the global  $Q$ -function using only the pairwise interactions of agents:  $Q^j(s^j, \mathbf{a}) = \frac{1}{M} \sum_{k \neq j} Q^j(s^j, a^j, a^k)$ . With the assumptions of homogeneity and locality, under certain preference  $\omega$ , the local pairwise interactions can be approximated by the interplay of each agent with the mean effect from its neighbors:

$$\frac{1}{|M|} \sum_{k \neq j} \omega^\top Q(s^j, a^j, a^k) = \omega^\top Q(s^j, a^j, \bar{a}^j), \quad (6)$$

where  $a$  is the single-agent action, and  $\bar{a}$  is the mean action. Given this approximated  $Q$ -function, we can derive the agent's policy function with the softmax parameterization:

$$\pi^j(a^j | s^j, \bar{a}^j) = \frac{\exp(\beta \omega^\top Q(s^j, a^j, \bar{a}^j))}{\sum_{a \in A^j} \exp(\beta \omega^\top Q(s^j, a, \bar{a}^j))}, \quad (7)$$

where  $\beta$  is the Boltzmann parameter.

2) *Bellman Optimality Operator*: To extend this framework to multi-objective reinforcement learning, we modify the envelop  $Q$ -learning [25] by combining the mean field operator with the envelop optimality operator. We first condition all the  $Q$ -values on the linear preference  $\omega$ , as in  $Q(s, \mathbf{a}, \omega)$ . As the standard  $Q$ -learning [26], we define the bi-objective

---

**Algorithm 1:** Mean Field Control with Envelop  $Q$ -learning

---

```

1 Initialize the  $Q$ -network  $Q_\theta$  and the target  $Q$ -network  $Q_{\bar{\theta}}$ 
2 Initialize the replay buffer  $\mathcal{D}$  and set  $\zeta = 0$ 
3 for episode = 1, ...,  $E$  do
4   Initialize  $\bar{a}_0^j$  for all  $j \in [M]$ 
5   for  $t = 1, \dots, T_{max}$  do
6     Sample  $a_t^j$   $\epsilon$ -greedily from  $Q_\theta$  by Eq. (7) for all  $j \in [M]$ 
7     Compute new mean actions  $\bar{a}_t^j$  by Eq. (5) for all  $j \in [M]$ 
8     Take the joint action  $\mathbf{a}_t = [a_t^1, \dots, a_t^M]$  from the state  $\mathbf{s}$  to
       the next state  $\mathbf{s}_{t+1}$ 
9     Compute the reward  $\mathbf{r}_t = [r_t^1, \dots, r_t^M]$  by Eq. (4)
10    Store the transition,  $\langle \mathbf{s}_t, \mathbf{a}_t, \mathbf{r}_t, \mathbf{s}_{t+1}, \bar{\mathbf{a}} \rangle$ , into  $\mathcal{D}$ , where
        $\bar{\mathbf{a}}_t = [\bar{a}_t^1, \dots, \bar{a}_t^M]$  is the collection of mean actions
11  if update then
12    Sample  $N$  transitions from  $\mathcal{D}$  and  $N_\omega$  preferences from  $\Omega$ 
13    Compute the TD target using the operator in Eq. (8)
14    Update  $Q_\theta$  by minimizing the loss from Eq. (11)
15  Update  $Q_{\bar{\theta}}$  with the learning rate  $\alpha$ :  $\bar{\theta} \leftarrow \alpha \theta + (1 - \alpha) \bar{\theta}$ 
16  Increase  $\zeta$  along the predefined homotopy path

```

---

multi-agent Bellman optimality operator  $\mathcal{T}$  as:

$$\begin{aligned}
(\mathcal{T}Q)(s, \mathbf{a}, \omega) &:= \sum_{j=1}^M r^j(s^j, a^j) \\
&+ \gamma \mathbb{E}_{s'} \sum_{j=1}^M \arg_{Q^j} \left\{ \max_{\omega' \in \Omega} \max_{a^j \in A^j} \omega'^\top Q^j(s'^j, a^j, \bar{a}^j, \omega') \right\}, \quad (8)
\end{aligned}$$

where  $\arg_{Q^j}$  takes out  $Q^j$  that maximizes  $\omega'^\top Q^j$ . This operator resembles the Bellman optimality operator in the standard  $Q$ -learning for single-agent RL and provides the temporal difference (TD) target. The difference is that it also optimizes over the parameter of preference  $\omega$ . By maximizing  $\omega'$  over the next state and its onward trajectory, this approach provides an optimistic perspective for its future rewards. Iteratively applying this operator to the  $Q$ -function, we will be able to reach the convergence of the near-optimal  $Q$ -function, which has been proven in [25].

### C. Double $Q$ -learning

We design our learning algorithm based on the double  $Q$ -learning [27] with two different loss functions and the target network. Algorithm 1 presents the detailed learning framework. During the rollout phase (Line 5-10), we sample the transitions in the multi-agent environment with the homogeneous policy. After we obtain enough transitions in the replay buffer, we enter the learning phase (Line 11-14). Given a mini-batch of  $N$  transitions and  $N_\omega$  preferences, we can estimate the TD target  $\mathbf{y} = (\mathcal{T}Q)(s, \mathbf{a}, \omega)$  via Eq. (8). The first loss function can be computed as the  $L_2$ -norm of the multi-objective TD:

$$L_A(\theta) = \mathbb{E}_{s, \mathbf{a}, \omega} \left[ \left\| \mathbf{y} - \sum_{j=1}^M Q_\theta(s^j, a^j, \bar{a}^j, \omega) \right\|_2^2 \right]. \quad (9)$$

Although this loss function is close to the true expected return, the non-smooth surface makes the learning process



difficult in the early steps. We combine this with another additional loss function with the projected temporal difference:

$$L_B(\theta) = \mathbb{E}_{s, \mathbf{a}, \omega} \left[ \left| \omega^\top \left( \mathbf{y} - \sum_{j=1}^M \mathbf{Q}_\theta(s^j, \mathbf{a}^j, \bar{\mathbf{a}}^j, \omega) \right) \right| \right]. \quad (10)$$

$L_A(\theta)$  provides a closer estimation of the true  $Q$ -function, since it evaluates the  $Q$ -function w.r.t. the optimal frontier which could contain a large number of solutions and result in difficulties for optimization.  $L_B(\theta)$ , on the other hand, makes the landscape of optimization smooth and easy to optimize. We train the  $Q$ -network via the homotopy optimization [28] based on the combination of these two loss functions:

$$L(\theta) = (1 - \zeta)L_A(\theta) + \zeta L_B(\theta), \quad (11)$$

where, in our case, we gradually increases  $\zeta$  from 0 to 1 exponentially as the learning progresses.

## V. EMPIRICAL EVALUATION

In this section, we implement MFC-EQ and experimentally evaluate it with other methods on a server equipped with Intel 2.3GHz 16-Core CPUs and NVIDIA A40 GPUs.

### A. Experimental Setups

We use 4-neighbor grids with two obstacle-free corners in the top-left and the bottom-right. The default obstacle density for grids outside these two corners is set to be 10%. The agents start at the top-left corner and travel toward the bottom-right corner. The formation in the goal position represents the desired formation. We refer to the size of grids as map size and the size of corners as formation size. For each data point, we averaged over 100 samples by crossing 10 random maps and 10 random formations.

### B. Ablatio Study

We first evaluate the ability of our learned policy to adapt to different preferences. We use the environment that has 16 agents with  $48 \times 48$  map size and  $9 \times 9$  formation size. We test different preferences  $\omega = (\lambda, 1 - \lambda)^\top$  by varying  $\lambda$  from  $\{0.1, 0.3, 0.5, 0.7, 0.9\}$ . We generate 5 corresponding solutions and evaluate them under different  $\text{MIX}(\lambda)$  objectives by varying  $\lambda$  from the same set of values. Table I provides 5 different solutions, and every  $\text{MIX}$  column highlights the solution that minimizes the projection onto that particular preference. As we can observe, all  $\text{MIX}(\lambda)$  objectives are minimized by the exact preference  $\omega(\lambda)$  that is given to our policy. This suggests that the learned policy can adapt to different preferences and produce versatile solutions that fit the required objectives.

### C. Centralized Baselines

We compare our method with several methods which have to plan all paths on centralized servers before execution.

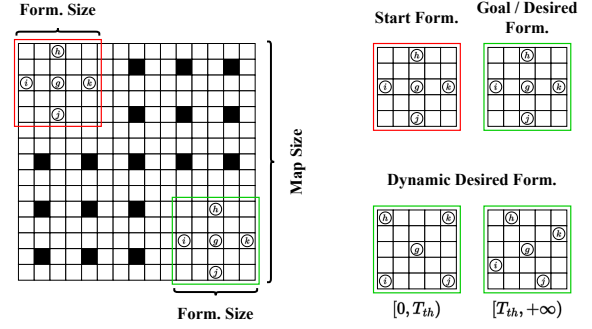


Fig. 4: Demonstration of experiment environments.

$\omega(\lambda)$	Make-span	Form. Dev.	MIX(0.1)	MIX(0.3)	MIX(0.5)	MIX(0.7)	MIX(0.9)
0.1	106.33	14.67	<b>23.84</b>	42.17	60.50	78.83	97.16
0.3	101.14	15.37	23.95	<b>41.10</b>	58.26	75.41	92.56
0.5	98.64	16.84	25.02	41.38	<b>57.74</b>	74.10	90.46
0.7	96.74	19.16	26.92	42.43	57.95	<b>73.47</b>	88.98
0.9	96.42	21.75	29.22	44.15	59.09	74.02	<b>88.95</b>

TABLE I: Results of MFC-EQ with different preferences evaluated by different scalarized objectives.

1) *Scalarized Prioritized Planning (SPP)*: Since it is NP-hard to solve this problem optimally, we come up with an efficient yet suboptimal baseline based on the prioritized planning algorithm [29]. We first give each agent a unique priority, and in that priority order, a low-level A\* search will be invoked to plan the path from the start location to the goal location while respecting the already planned paths of all agents with higher priorities. The low-level A\* search uses a scalarized  $f$ -value which is a mix of the makespan  $f$ -value ( $f_{MS}$ ) and the formation deviation  $f$ -value ( $f_{FD}$ ):

$$f(n) = \underbrace{\lambda [cost(\mathbf{v}^i) + dist(\mathbf{v}^i, \mathbf{g}^i)]}_{f_{MS}(n)} + (1 - \lambda) \underbrace{\sum_{t=1}^{cost(\mathbf{v}^i)} \mathcal{F}_P^i(t_n)}_{f_{FD}(n)}, \quad (12)$$

where  $\mathcal{F}_P^i(t_n)$  is the partial formation deviation among agent  $i$  and all other agents with higher priorities. This baseline is not complete but, in most cases, can find a possible solution much more quickly, albeit the solutions usually have poor quality, especially in congested environments with large numbers of agents. Moreover, this planner, unlike SWARM-MAPF, can target any given linear preference.

We use the scalarized  $f$ -value to combine these two objectives. The weights for makespan and formation deviation are set to  $\lambda$  and  $1 - \lambda$ , respectively. In the experiments, we vary  $\lambda$  from  $\{0.1, 0.3, 0.5, 0.7, 0.9\}$  to test its performance under different linear preferences.

2) *SWARM-MAPF (SWARM)*: The most effective centralized method, SWARM-MAPF, has been proposed in [6], which combines the swarm-based formation control with the conflict-based MAPF algorithms. The SWARM-MAPF is a two-phase algorithm. In Phase 1, it first calculates the lower bound of the makespan  $B = \max_{1 \leq i \leq M} dist(s_i, g_i)$ .

Map Size	$M$	Success Rate			Makespan			Form. Dev.			MIX(0.5)		
		SPP	SWA-RM	MFC-EQ	SPP	SWA-RM	MFC-EQ	SPP	SWA-RM	MFC-EQ	SPP	SWA-RM	MFC-EQ
32	10	1.00	1.00	1.00	48.30	59.32	60.24	29.79	6.07	4.35	39.05	32.70	<b>32.30</b>
	20	1.00	0.99	0.99	49.03	63.17	60.38	32.42	12.38	10.04	40.73	37.78	<b>35.21</b>
	30	0.79	0.96	0.90	51.54	59.09	54.59	42.25	20.64	20.32	46.90	39.87	<b>37.46</b>
48	10	1.00	0.99	0.99	80.44	98.10	88.07	53.91	8.18	11.05	67.18	53.14	<b>49.56</b>
	20	0.95	0.99	0.96	82.07	108.84	104.28	70.44	23.70	21.49	76.26	66.27	<b>62.89</b>
	30	0.74	0.94	0.88	84.92	101.52	107.42	96.04	36.30	37.18	90.48	68.91	72.30
64	10	1.00	0.99	0.99	113.38	144.54	137.14	97.92	15.00	16.43	105.65	79.77	<b>76.79</b>
	20	1.00	0.97	0.93	114.56	156.03	141.26	113.52	33.24	28.34	114.04	94.64	<b>84.80</b>
	30	0.22	0.98	0.90	115.59	142.65	145.51	107.64	57.31	61.43	111.62	99.98	103.47

TABLE II: Results for MFC-EQ and centralized baselines with different numbers of agents in various sizes of grids.

Map Size	Form Size	Success Rate			Makespan			Form. Dev.			MIX(0.5)		
		SPP	SWA-RM	MFC-EQ	SPP	SWA-RM	MFC-EQ	SPP	SWA-RM	MFC-EQ	SPP	SWA-RM	MFC-EQ
32	7×7	0.94	1.00	0.97	53.81	66.40	68.30	44.66	12.37	9.06	49.24	39.39	<b>38.68</b>
	9×9	1.00	1.00	1.00	48.56	63.20	67.33	29.34	9.10	8.72	38.95	<b>36.15</b>	38.03
	11×11	1.00	1.00	1.00	44.18	57.75	55.12	20.80	7.03	8.32	32.49	32.39	<b>31.72</b>
48	7×7	0.98	1.00	0.87	86.26	110.94	107.37	82.85	19.84	22.40	84.56	65.39	<b>64.89</b>
	9×9	0.93	0.96	0.93	81.49	109.67	98.64	65.74	21.03	16.84	73.62	65.35	<b>57.74</b>
	11×11	1.00	1.00	1.00	77.06	105.06	97.26	60.65	15.12	14.08	68.86	60.09	<b>55.67</b>
64	7×7	0.87	0.99	0.96	118.64	155.36	138.84	115.38	31.07	55.42	117.01	<b>93.22</b>	97.13
	9×9	1.00	1.00	0.96	113.87	153.33	133.92	108.57	25.95	33.20	111.22	89.64	<b>83.56</b>
	11×11	1.00	0.95	1.00	109.43	149.61	131.08	97.68	23.02	27.75	103.56	86.32	<b>79.42</b>

TABLE III: Results for MFC-EQ and centralized baselines under different formation sizes in various sizes of grids.

Given a user-provided parameter  $w \geq 1$ , SWARM-MAPF selects a leader from the group of agents such that its path, whose length is bounded by  $wB$ , can be sufficiently far away from the obstacles and thus others agents can preserve their formation as much as they can. In Phase 2, it will invoke the modified conflict-based search [30] (CBS-M) to minimize the makespan and replanning some critical segments. This planner is complete and suboptimal, but it cannot specifically target any given preference, since we cannot control the trade-off between two objectives based on the parameter  $w$ .

3) *Joint State A\* (JSA\*)*: The joint state A\* directly applies the  $\epsilon$ -constraint search algorithm [31] in the joint state space. The joint state assigns all agents a set of different locations. The operator assigns each agent a set of non-colliding move or wait actions. The OPEN list sorts nodes based on makespan, while the FOCAL list breaks ties based on the formation deviation. Details of this algorithm can be found in [6]. Since the joint state space grows exponentially w.r.t the number of agents, this method can only be applied to instances with relatively small agents (less than 5 agents in our setups). By varying the  $\epsilon$  in the focal search, this method is guaranteed to find the Pareto-optimal frontier.

#### D. Main Results

1) *Number of Agents*: We evaluate our methods under different numbers of agents in different map sizes and compare the results with centralized baselines. The  $\lambda$  is set to 0.5 for SPP and MFC-EQ. The  $w$  is set to 1.0 for SWARM. The runtime limit is only 30 seconds for MFC-EQ and 5 minutes for SWARM and SPP. Due to the partially observable environment, our method naturally does not have perfect success rates, but they are relatively high and acceptable. As we can observe from Table II, SWARM has greater performance in almost all the test cases. We also can

$M$	Success Rate			Makespan			Form. Dev.			MIX(0.5)		
	SPP	SWA-RM	MFC-EQ	SPP	SWA-RM	MFC-EQ	SPP	SWA-RM	MFC-EQ	SPP	SWA-RM	MFC-EQ
10	1.00	0.98	0.96	48.19	59.00	56.33	127.90	172.40	104.61	88.05	115.70	<b>80.47</b>
15	1.00	1.00	1.00	48.29	63.85	57.56	132.71	210.82	114.33	90.50	137.34	<b>85.95</b>
20	0.97	1.00	1.00	49.64	63.60	59.07	141.18	208.28	118.42	95.41	135.94	<b>88.75</b>
25	0.72	1.00	1.00	50.56	62.58	61.29	149.61	204.33	123.20	100.09	133.46	<b>92.25</b>
30	0.90	0.98	0.93	50.00	59.31	62.50	146.82	187.08	129.74	98.41	123.20	<b>96.12</b>
35	0.48	0.94	0.87	51.75	57.87	64.71	163.40	189.64	133.07	107.58	123.76	<b>98.89</b>
40	0.25	0.81	0.74	52.68	54.12	65.29	168.64	165.05	137.33	110.66	109.59	<b>101.31</b>

TABLE IV: Results for MFC-EQ and centralized baselines in scenarios that require dynamic formation.

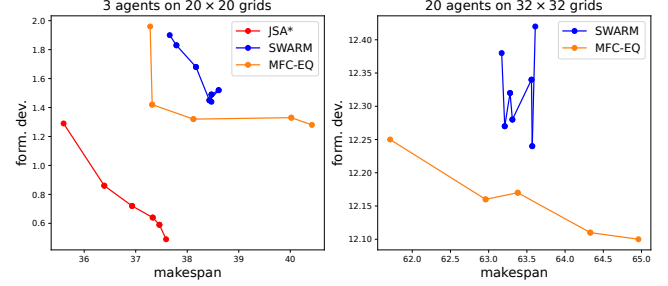


Fig. 5: Trade-off of makespan and formation deviation.

see when projected to the demanded preference, our method can outperform SWARM in most instances. This experiment shows that our method can produce comparable results to the state-of-the-art centralized baseline and scale up well to instances with large numbers of agents in different sizes of maps.

2) *Formation Size*: We repeat the experiment above with various formation sizes. We choose different sizes of obstacle-free corners in which the formation is randomly generated. The larger the corner is, the more spread out the formation will be. The number of agents is fixed at 16. As shown in Table III, we see that smaller formations are usually more difficult to solve, resulting in larger makespan and formation deviation. Compared to SWARM, SPP generally has a better makespan but much worse formation deviation and MIX. MFC-EQ solves most instances with greater solution quality in both objectives when compared to the baselines.

3) *Dynamic Formation*: We also put these methods into more challenging tests where the agents will be asked to adjust to different formations on the fly. We evaluate agents' formations with one desired formation before  $T_{th} = 30$  and another different formation after that. The centralized scheme cannot handle such tasks as agents' paths will have to be planned before execution. In our method, we can simply notify each decentralized agent of the new formation which will result in different ways of calculating relative positions, and therefore the agents can adjust to the new formation seamlessly. The results are shown in Table IV, suggesting that our method has the flexibility to tackle changeable formations, while others result in much larger deviation.

4) *Makespan and Formation Trade-off*: We further compare our method with others under different preferences. Due to the limited scalability of JSA\*, we first use the environment with  $20 \times 20$  map size,  $3 \times 3$  formation size, 15%

obstacle density, and 3 agents. We vary  $\epsilon$  of JSA\* from 1.0 to 1.8 and  $w$  of SWARM from 1.0 to 1.6. The value of  $\lambda$  for  $\omega$  in MFC-EQ is varied from 0.1 to 0.9. JSA\* can provide the Pareto-optimal frontier only for small-scale instances. We then repeat this experiment in larger instances with  $32 \times 32$  map size,  $9 \times 9$  formation size, and 20 agents. Fig. 5 shows the results. Although SPP is also tested, it only gives solutions that have near-optimal makespan but significantly larger formation deviation than the shown results. In large-scale cases, SWARM tends to fluctuate, meaning that, even given more makespan allowance, it may result in solutions with worse formation deviation. The envelope generated by our method is near-convex and can cover all solutions from SWARM, albeit still suboptimal. It also has a wider range of makespan with greater solution variety.

## VI. CONCLUSIONS

We proposed MFC-EQ, a general  $Q$ -learning framework for solving decentralized MAiF with partial observation and limited communication. MFC-EQ utilizes the mean field approximation to simplify the complex multi-agent interaction and employs the envelope  $Q$ -learning to enable the adaptability to various preferences for this bi-objective task. Our theoretical proofs further show that the combination of these two operators can still converge to a fixed optimum. Empirical results demonstrate that MFC-EQ outperforms existing centralized baselines in most cases and is more versatile in handling dynamically changing desired formations. Moreover, MFC-EQ is not limited to solving MAiF and has great potential to be generalized to other multi-objective tasks in multi-agent systems.

## REFERENCES

- [1] R. Stern, N. Sturtevant, A. Felner, S. Koenig, H. Ma, T. Walker, J. Li, D. Atzmon, L. Cohen, T. Kumar, *et al.*, "Multi-agent pathfinding: Definitions, variants, and benchmarks," in *Proceedings of the International Symposium on Combinatorial Search*, vol. 10, no. 1, 2019, pp. 151–158.
- [2] P. R. Wurman, R. D'Andrea, and M. Mountz, "Coordinating hundreds of cooperative, autonomous vehicles in warehouses," *AI magazine*, vol. 29, no. 1, pp. 9–9, 2008.
- [3] R. Morris, C. S. Pasareanu, K. S. Luckow, W. Malik, H. Ma, T. S. Kumar, and S. Koenig, "Planning, scheduling and monitoring for airport surface operations," in *AAAI Workshop: Planning for Hybrid Systems*, 2016, pp. 608–614.
- [4] H. Ma, J. Yang, L. Cohen, T. Kumar, and S. Koenig, "Feasibility study: Moving non-homogeneous teams in congested video game environments," in *Proceedings of the AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment*, vol. 13, no. 1, 2017, pp. 270–272.
- [5] A. Gautam and S. Mohan, "A review of research in multi-robot systems," in *ICIIS*. IEEE, 2012, pp. 1–5.
- [6] J. Li, K. Sun, H. Ma, A. Felner, T. Kumar, and S. Koenig, "Moving agents in formation in congested environments," in *Proceedings of the International Symposium on Combinatorial Search*, vol. 11, no. 1, 2020, pp. 131–132.
- [7] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.
- [8] 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 Robotics and Automation Letters*, vol. 4, no. 3, pp. 2378–2385, 2019.
- [9] 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 *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2020, pp. 11 748–11 754.
- [10] Z. Ma, Y. Luo, and H. Ma, "Distributed heuristic multi-agent path finding with communication," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2021, pp. 8699–8705.
- [11] A. Barreto, W. Dabney, R. Munos, J. J. Hunt, T. Schaul, H. P. van Hasselt, and D. Silver, "Successor features for transfer in reinforcement learning," *Advances in neural information processing systems*, vol. 30, 2017.
- [12] M. L. Littman, "Markov games as a framework for multi-agent reinforcement learning," in *Machine learning proceedings 1994*. Elsevier, 1994, pp. 157–163.
- [13] H. E. Stanley, *Phase transitions and critical phenomena*. Clarendon Press, Oxford, 1971, vol. 7.
- [14] Y. Yang, R. Luo, M. Li, M. Zhou, W. Zhang, and J. Wang, "Mean field multi-agent reinforcement learning," in *International conference on machine learning*. PMLR, 2018, pp. 5571–5580.
- [15] S. G. Subramanian, M. E. Taylor, M. Crowley, and P. Poupart, "Partially observable mean field reinforcement learning," *arXiv preprint arXiv:2012.15791*, 2020.
- [16] S. G. Subramanian, P. Poupart, M. E. Taylor, and N. Hegde, "Multi type mean field reinforcement learning," *arXiv preprint arXiv:2002.02513*, 2020.
- [17] Z. Gábor, Z. Kalmár, and C. Szepesvári, "Multi-criteria reinforcement learning," in *ICML*, vol. 98, 1998, pp. 197–205.
- [18] S. Mannor and N. Shimkin, "The steering approach for multi-criteria reinforcement learning," *Advances in Neural Information Processing Systems*, vol. 14, 2001.
- [19] S. Natarajan and P. Tadepalli, "Dynamic preferences in multi-criteria reinforcement learning," in *Proceedings of the 22nd international conference on Machine learning*, 2005, pp. 601–608.
- [20] S. Parisi, M. Pirodda, N. Smacchia, L. Bascetta, and M. Restelli, "Policy gradient approaches for multi-objective sequential decision making," in *2014 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2014, pp. 2323–2330.
- [21] K. Van Moffaert and A. Nowé, "Multi-objective reinforcement learning using sets of pareto dominating policies," *The Journal of Machine Learning Research*, vol. 15, no. 1, pp. 3483–3512, 2014.
- [22] X. Chen, A. Ghadirzadeh, M. Björkman, and P. Jensfelt, "Meta-learning for multi-objective reinforcement learning," in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2019, pp. 977–983.
- [23] A. Castelletti, F. Pianosi, and M. Restelli, "Multi-objective fitted q-iteration: Pareto frontier approximation in one single run," in *2011 International Conference on Networking, Sensing and Control*. IEEE, 2011, pp. 260–265.
- [24] A. Abels, D. Roijers, T. Lenaerts, A. Nowé, and D. Steckelmacher, "Dynamic weights in multi-objective deep reinforcement learning," in *International conference on machine learning*. PMLR, 2019, pp. 11–20.
- [25] R. Yang, X. Sun, and K. Narasimhan, "A generalized algorithm for multi-objective reinforcement learning and policy adaptation," *Advances in neural information processing systems*, vol. 32, 2019.
- [26] C. J. Watkins and P. Dayan, "Q-learning," *Machine learning*, vol. 8, pp. 279–292, 1992.
- [27] H. Hasselt, "Double q-learning," *Advances in neural information processing systems*, vol. 23, 2010.
- [28] L. T. Watson and R. T. Haftka, "Modern homotopy methods in optimization," *Computer Methods in Applied Mechanics and Engineering*, vol. 74, no. 3, pp. 289–305, 1989.
- [29] D. Silver, "Cooperative pathfinding," in *Proceedings of the aaai conference on artificial intelligence and interactive digital entertainment*, vol. 1, no. 1, 2005, pp. 117–122.
- [30] G. Sharon, R. Stern, A. Felner, and N. R. Sturtevant, "Conflict-based search for optimal multi-agent pathfinding," *Artificial Intelligence*, vol. 219, pp. 40–66, 2015.
- [31] Y. Haimes, "On a bicriterion formulation of the problems of integrated system identification and system optimization," *IEEE transactions on systems, man, and cybernetics*, no. 3, pp. 296–297, 1971.
- [32] A. W. Naylor and G. R. Sell, *Linear operator theory in engineering and science*. Springer Science & Business Media, 1982.



In this appendix, we provide some theoretical analysis for the proposed method.

#### A. Mean Field Approximation in Multi-Objective Reinforcement Learning

In Section IV-B, we have Eq. 6 that gives an approximation for all pairwise interactions among agents. Here, we provide a theorem to analyze the legitimacy of such approximation in the context of multiple objectives. We adopt the proof from [14] and extend it to the multi-objective settings.

*Theorem 1 (Mean Field Approximation):* With the assumptions of homogeneity and locality, under certain preference  $\omega$ , the local pairwise interactions can be approximated by the interplay of each agent with the mean effect from its neighbors:

$$\frac{1}{M} \sum_{k \neq j} \omega^\top Q(s_t^j, a_t^j, a_t^k) \approx \omega^\top Q(s_t^j, a_t^j, \bar{a}_t^j).$$

*Proof:* First, we denote the difference between each neighboring agent's action and the mean action as  $\delta_t^{j,k} = a_t^k - \bar{a}_t^j$  for agent  $j$  where  $\bar{a}_t^j$  is the mean action. According to the definition of the mean action,  $\sum_{k \neq j} \delta_t^{j,k} = 0$ . Assuming the given  $Q$ -function is twice-differentiable, we then approximate each pairwise  $Q$ -function at  $\bar{a}_t^j$  w.r.t  $a_t^k$  using Taylor's theorem:

$$\begin{aligned} & \frac{1}{M} \sum_{k \neq j} \omega^\top Q^j(s_t^j, a_t^j, a_t^k) \\ &= \frac{1}{M} \sum_{k \neq j} \left[ \omega^\top Q^j(s_t^j, a_t^j, \bar{a}_t^j) + \omega^\top \nabla_{\bar{a}_t^j} Q^j(s_t^j, a_t^j, \bar{a}_t^j) \cdot \delta_t^{j,k} + \frac{1}{2} \omega^\top (\delta_t^{j,k})^\top \nabla_{\bar{a}_t^j}^2 Q^j(s_t^j, a_t^j, \bar{a}_t^j) \cdot \delta_t^{j,k} \right] \\ &= \omega^\top Q^j(s_t^j, a_t^j, \bar{a}_t^j) + \frac{1}{M} \omega^\top \nabla Q^j(s_t^j, a_t^j, \bar{a}_t^j) \sum_{k \neq j} \delta_t^{j,k} + \frac{1}{2M} \sum_{k \neq j} \omega^\top (\delta_t^{j,k})^\top \nabla^2 Q^j(s_t^j, a_t^j, \bar{a}_t^j) \cdot \delta_t^{j,k} \end{aligned}$$

We first drop the second term since  $\sum_{k \neq j} \delta_t^{j,k} = 0$ . Given that the  $Q$ -functions are  $L$ -smooth, we know that for each dimension (objective)  $i$  of  $Q^j$ , we have that  $\nabla^2 Q_i^j \preceq LI_{|A|}$ , meaning that  $\sigma_{\max}(\nabla^2 Q_i^j) \leq L$  and  $\sigma_{\min}(\nabla^2 Q_i^j) \geq -L$ , where  $\sigma_{\max}$  is the largest eigenvalue and  $\sigma_{\min}$  is the smallest eigenvalue. Without the loss of generality, we then limit the preference  $\omega$  such that  $|\omega| = 1$ . As mentioned above, in our task, we define  $\omega = (\lambda, 1 - \lambda)^\top$ , which satisfies this condition. Then the largest eigenvalue of the Hessian matrix of  $\omega^\top Q^j$  is:

$$\begin{aligned} \sigma_{\max}(\omega^\top \nabla^2 Q^j(s_t^j, a_t^j, \bar{a}_t^j)) &\leq \sigma_{\max}(\nabla^2 \left[ |\omega| \cdot |Q^j(s_t^j, a_t^j, \bar{a}_t^j)| \right]) && \text{(Cauchy-Schwarz inequality)} \\ &= |\omega| \cdot \sigma_{\max}(\nabla^2 |Q^j(s_t^j, a_t^j, \bar{a}_t^j)|) \leq L. && (|\omega| = 1 \text{ and } \sigma_{\max} \leq L) \end{aligned}$$

Similarly, we can also have  $\sigma_{\min}(\omega^\top \nabla^2 Q^j(s_t^j, a_t^j, \bar{a}_t^j)) \geq -L$ . Since the Hessian matrix is symmetric and diagonalizable, we can apply the orthogonal decomposition:  $\omega^\top \nabla^2 Q^j(s_t^j, a_t^j, \bar{a}_t^j) = U^\top \Sigma U$ , where  $U$  is an orthogonal matrix and  $\Sigma = \text{diag}[\sigma_1, \dots, \sigma_{|A|}]$  is the eigenmatrix of  $\nabla^2 \omega^\top Q^j(s_t^j, a_t^j, \bar{a}_t^j)$ . It can be shown that

$$\begin{aligned} (\delta_t^{j,k})^\top \nabla^2 \omega^\top Q^j(s_t^j, a_t^j, \bar{a}_t^j) (\delta_t^{j,k}) &= (\delta_t^{j,k})^\top U^\top \Sigma U (\delta_t^{j,k}) \\ &\leq \sigma_{\max}(U^\top \delta_t^{j,k})^2 \\ &= L \cdot (\delta_t^{j,k})^\top (\delta_t^{j,k}) && (U^\top U = 1 \text{ and } \sigma_{\max} \leq L) \\ &= L \cdot (a_t^k - \bar{a}_t^j)^\top (a_t^k - \bar{a}_t^j) && \text{(definition of } \delta_t^{j,k}) \\ &\leq L \cdot [(a_t^k)^\top (a_t^k) + (\bar{a}_t^j)^\top (\bar{a}_t^j)] \\ &\leq 2L. && \text{(one-hot encoding)} \end{aligned}$$

Similarly, we can show that  $(\delta_t^{j,k})^\top \nabla^2 \omega^\top Q^j(s_t^j, a_t^j, \bar{a}_t^j) (\delta_t^{j,k}) \geq -2L$ . Therefore, we prove that this term is bounded by  $[-2L, 2L]$ . In this symmetrical range, with  $L$  being relatively small and the assumptions of homogeneity and locality, these terms across the neighborhood tend to cancel each other, resulting in the third term close to 0. Finally, we prove that

$$\frac{1}{M} \sum_{k \neq j} \omega^\top Q^j(s_t^j, a_t^j, a_t^k) \approx \omega^\top Q^j(s_t^j, a_t^j, \bar{a}_t^j)$$

■

### B. Convergence of MFC-EQ

Here, we provide some theoretical insights for the convergence of MFC-EQ. The multi-object multi-agent Bellman operator can be regarded as a combination of the mean field operator and the envelop optimality operator. However, it is not obvious that the combined operator can still guarantee the convergence of  $Q$ -function to any fixed optimum. We adopt the proof from [25] and extend it to homogeneous multi-agent systems.

To analyze the convergence, we first define the metric of distance between any two considered  $Q$ -function.

*Definition 1:* The distance between any two  $Q$ -functions  $Q$  and  $Q'$  in  $\mathcal{Q} \subseteq S^j \times A^j \times \bar{A}^j \times \Omega \rightarrow \mathbb{R}^2$  is defined as:

$$d(Q, Q') := \max_{\substack{s \in S, a \in A \\ \omega \in \Omega}} \left| \sum_{j=1}^M \omega^\top (Q(s^j, a^j, \bar{a}^j, \omega) - Q'(s^j, a^j, \bar{a}^j, \omega)) \right|,$$

where  $\bar{a}$  is the mean action defined by Eq. (5).

This metric forms a complete pseudo-metric space [32]. We then generalize the theorems from [25] to multi-agent environments over the defined metric space  $\langle \mathcal{Q}, d \rangle$ . The proofs follow the framework of the well-known Banach's Fixed-Point Theorem. We first define the optimal  $Q$ -function and prove that it is a fixed point of the operator  $\mathcal{T}$  in Eq. (8).

To prove that MFC-EQ is still guaranteed to converge in multi-agent systems, we follow the proof framework of the well-known Banach's Fixed-Point Theorem. We first define the multi-objective multi-agent Bellman optimality operator  $\mathcal{T}$ . Then we prove the following theorems. Theorem 2 states that there exists an optimal  $Q$ -function that is a fixed point under the multi-objective multi-agent operator  $\mathcal{T}$ . Theorem 3 states that  $\mathcal{T}$  is a contraction map, meaning that iterative applying it to the  $Q$ -function will lead to a closer distance to the fixed point. Lastly, Theorem 4 concludes that in the pseudo-metric space  $\langle \mathcal{Q}, d \rangle$ ,  $\mathcal{T}$  results in convergence to the optimal point.

*Theorem 2 (Fixed Point):* Let  $Q^*$  be the optimal  $Q$ -function under the preference  $\omega$ , defined as:

$$Q^* = \arg_Q \max_{\pi} \mathbb{E}_{\tau \sim (s, a, \pi, P_S)} \left[ \sum_{t=0}^{\infty} \gamma^t \omega^\top \sum_{j=1}^M \mathbf{r}^j(s_t^j, a_t^j) \right],$$

where  $\tau$  denotes the trajectory under  $\pi$ . Then,  $Q^*$  is a fixed point under  $\mathcal{T}$ , that is,  $Q^* = \mathcal{T}Q^*$ .

*Proof:* We first expand  $\omega^\top \mathcal{T}Q^*$ :

$$\begin{aligned} \omega^\top \mathcal{T}Q^* &= \omega^\top \mathcal{T} \arg_Q \max_{\pi} \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t \omega^\top \sum_{j=1}^M \mathbf{r}^j(s_t^j, a_t^j) \right] && \text{(definition of } Q^*) \\ &= \omega^\top \sum_{j=1}^M \mathbf{r}^j(s_t^j, a_t^j) + \gamma \mathbb{E}_{s_{t+1}} \sum_{j=1}^M \max_{a^j \in A^j} \max_{\omega' \in \Omega} \omega^\top \arg_{Q^j} \max_{\pi} \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t \omega'^\top \sum_{j=1}^M \mathbf{r}^j(s_t^j, a_t^j) \right] && \text{(definition of } \mathcal{T}) \end{aligned}$$

Assume that  $\omega'_*$  is the maximum and the corresponding policy is  $\pi_{\omega'_*}$ . It can be shown that  $\omega^\top \arg_{Q^j} \max_{\pi} Q(\omega) \leq \max_{\omega'} \omega^\top \arg_Q \max_{\pi} Q(\omega')$  due to the definition of  $\max_{\omega'}$ . We also know that  $\max_{\omega'} \omega^\top \arg_Q \max_{\pi} Q(\omega') = \omega^\top Q^{\pi_{\omega'_*}} \leq \omega^\top \arg_Q \max_{\pi} Q(\omega)$  due to the definition of  $\arg_Q$  and  $\max_{\pi}$ . Hence, the two terms have to be equal to avoid contradictory, as in  $\omega^\top \arg_Q \max_{\pi} Q(\omega) = \max_{\omega'} \omega^\top \arg_Q \max_{\pi} Q(\omega')$ . We plug that into the equation above.

$$\begin{aligned} &= \omega^\top \sum_{j=1}^M \mathbf{r}^j(s_t^j, a_t^j) + \gamma \mathbb{E}_{s_{t+1}} \sum_{j=1}^M \max_{a^j \in A^j} \omega^\top \arg_{Q^j} \max_{\pi^j} \mathbb{E}_{\tau \sim \pi^j} \left[ \sum_{t=0}^{\infty} \gamma^t \omega^\top \sum_{j=1}^M \mathbf{r}^j(s_t^j, a_t^j) \right] \\ &= \omega^\top \sum_{j=1}^M \mathbf{r}^j(s_t^j, a_t^j) + \gamma \omega^\top \sum_{j=1}^M \arg_{Q^j} \max_{\pi^j} \mathbb{E}_{\tau \sim \pi^j} \left[ \sum_{t=0}^{\infty} \gamma^t \sum_{j=1}^M \mathbf{r}^j(s_t^j, a_t^j) \right] && \text{(rearrange)} \\ &= \omega^\top \sum_{j=1}^M \mathbf{r}^j(s_t^j, a_t^j) + \gamma \sum_{j=1}^M \arg_Q \max_{\pi} \omega^\top \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t \mathbf{r}^j(s_t^j, a_t^j) \right] && \text{(homogeneity)} \\ &= \omega^\top \sum_{j=1}^M \arg_Q \max_{\pi} \omega^\top \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t \sum_{j=1}^M \mathbf{r}^j(s_t^j, a_t^j) \right] && \text{(rearrange)} \\ &= \omega^\top Q^*. && \text{(definition of } Q^*) \end{aligned}$$

For any arbitrary  $\omega$ , we have  $\omega^\top Q^* = \omega^\top \mathcal{T}Q^* \Rightarrow \omega^\top (Q^* - \mathcal{T}Q^*) = 0$ . Since this holds for any  $\omega$ , we can finally conclude that  $Q^* = \mathcal{T}Q^*$ .  $\blacksquare$

We then state that the operator  $\mathcal{T}$  forms a contraction mapping under the metric  $d$  in  $Q$ -function space  $\mathcal{Q}$ .

*Theorem 3 (Contraction):*  $\mathcal{T}$  forms a contraction mapping on the pseudo-metric space  $\langle \mathcal{Q}, d \rangle$ , that is, for any  $Q, Q' \in \mathcal{Q}$ ,  $d(\mathcal{T}Q, \mathcal{T}Q') \leq \gamma d(Q, Q')$ , where  $\gamma$  ( $0 \leq \gamma < 1$ ) is the discount factor.

*Proof:* We first expand the expression of  $d(\mathcal{T}Q, \mathcal{T}Q')$ :

$$\begin{aligned} d(\mathcal{T}Q, \mathcal{T}Q') &= \max_{\substack{s \in \mathcal{S}, a \in \mathcal{A} \\ \omega \in \Omega}} \left| \sum_{j=1}^M \omega^\top ([\mathcal{T}Q(s^j, a^j, \bar{a}^j, \omega)]_j - [\mathcal{T}Q'(s^j, a^j, \bar{a}^j, \omega)]_j) \right| \\ &= \max_{\substack{s \in \mathcal{S}, a \in \mathcal{A} \\ \omega \in \Omega}} \left| \sum_{j=1}^M \gamma \cdot \omega^\top \left[ \mathbb{E}_{s'} \arg_{Q^j} \left\{ \max_{\omega'_{Q^j} \in \Omega, a^j_{Q^j} \in A^j} \omega^\top Q((s')^j, a^j_{Q^j}, (\bar{a}')^j, \omega'_{Q^j}) \right\} \right. \right. \\ &\quad \left. \left. - \mathbb{E}_{s'} \arg_{Q'^j} \left\{ \max_{\omega'_{Q'^j} \in \Omega, a^j_{Q'^j} \in A^j} \omega^\top Q'((s')^j, a^j_{Q'^j}, (\bar{a}')^j, \omega'_{Q'^j}) \right\} \right] \right| \end{aligned}$$

We first apply the Cauchy–Schwarz inequality to pull the sum over all the agents out. We then loosen the expectation to the maximum over the next state. Next, we drop the  $\omega^\top$  inside  $\arg_Q$  as we did in Theorem 2.

$$\leq \gamma \cdot \max_{s' \in \mathcal{S}, \omega \in \Omega} \sum_{j=1}^M \left| \max_{\omega'_{Q^j} \in \Omega, a^j_{Q^j} \in A^j} \omega^\top Q((s')^j, a^j_{Q^j}, (\bar{a}')^j, \omega'_{Q^j}) - \max_{\omega'_{Q'^j} \in \Omega, a^j_{Q'^j} \in A^j} \omega^\top Q'((s')^j, a^j_{Q'^j}, (\bar{a}')^j, \omega'_{Q'^j}) \right|$$

Let  $a_Q^*$  and  $\omega_Q^*$  be the optimal vectors that maximizes  $Q$ . Without the loss of generality, we assume that  $\omega^\top Q(s^j, a^j, \bar{a}^j, \omega') \geq \omega^\top Q'(s^j, a^j, \bar{a}^j, \omega')$ . Then we continue the expansion and rearrangement.

$$\begin{aligned} &\leq \gamma \cdot \max_{s' \in \mathcal{S}, \omega \in \Omega} \sum_{j=1}^M \left| \omega^\top Q((s')^j, (a^*)^j_{Q^j}, (\bar{a}')^j, \omega_Q^*) - \max_{\omega'_{Q'^j} \in \Omega, a^j_{Q'^j} \in A^j} \omega^\top Q'((s')^j, a^j_{Q'^j}, (\bar{a}')^j, \omega'_{Q'^j}) \right| \\ &\leq \gamma \cdot \max_{s' \in \mathcal{S}, \omega \in \Omega} \sum_{j=1}^M \left| \omega^\top Q((s')^j, (a^*)^j_{Q^j}, (\bar{a}')^j, \omega_Q^*) - \omega^\top Q'((s')^j, (a^*)^j_{Q^j}, (\bar{a}')^j, \omega_Q^*) \right| \\ &\leq \gamma \cdot \max_{\substack{s \in \mathcal{S}, a \in \mathcal{A} \\ \omega \in \Omega}} \sum_{j=1}^M \left| \omega^\top Q(s^j, a^j, \bar{a}^j, \omega_Q) - \omega^\top Q'(s^j, a^j, \bar{a}^j, \omega_Q) \right| \\ &= \gamma \cdot d(Q, Q') \end{aligned}$$

Therefore, we prove that  $d(\mathcal{T}Q, \mathcal{T}Q') \leq \gamma \cdot d(Q, Q')$ , and thus  $\mathcal{T}$  is a contraction mapping over  $\mathcal{Q}$ .  $\blacksquare$

Finally, we show that by iterating  $\mathcal{T}$  to any  $Q$ -function, it will asymptotically converge to the fixed optimum under  $d$ .

**Theorem 4 (Convergence):** Given the defined  $Q^*$  in Theorem 2 and the contraction mapping  $\mathcal{T}$  on the pseudo-metric space  $\langle \mathcal{Q}, d \rangle$  in Theorem 3, iterating  $\mathcal{T}$  leads to convergence to  $Q^*$ , for any  $Q \in \mathcal{Q}$ , that is,  $\lim_{n \rightarrow \infty} d(\mathcal{T}^n Q, Q^*) = 0$ .

*Proof:* We prove that  $\{\mathcal{T}^n Q\}$  is a Cauchy sequence under the metric  $d$ . We first look at

$$d(\mathcal{T}^{m+1} Q, \mathcal{T}^m Q) \leq \gamma d(\mathcal{T}^m Q, \mathcal{T}^{m-1} Q^*) \leq \dots \leq \gamma^m d(\mathcal{T} Q, Q).$$

For any  $m$  and  $l$  (w.l.o.g.,  $m > l$ ), by the triangular inequality we have

$$\begin{aligned} d(\mathcal{T}^m Q, \mathcal{T}^l Q) &\leq d(\mathcal{T}^m Q, \mathcal{T}^{m-1} Q) + \dots + d(\mathcal{T}^{l+1} Q, \mathcal{T}^l Q) \\ &\leq (\gamma^{m-1} + \dots + \gamma^l) \cdot d(\mathcal{T} Q, Q) \\ &= \frac{\gamma^m - \gamma^l}{1 - \gamma} \cdot d(\mathcal{T} Q, Q) \leq \frac{\gamma^m}{1 - \gamma} \cdot d(\mathcal{T} Q, Q) \end{aligned}$$

Since  $0 \leq \gamma < 1$ , we know that  $d(\mathcal{T}^m Q, \mathcal{T}^l Q)$  will converge to 0 as  $m, l \rightarrow \infty$ , and hence  $\{\mathcal{T}^n Q\}$  is a Cauchy sequence. We can further conclude that  $\lim_{n \rightarrow \infty} d(\mathcal{T}^n Q, Q_\infty) = 0$  such that  $d(Q_\infty, Q^*) = 0$  which yields that  $\lim_{n \rightarrow \infty} d(\mathcal{T}^n Q, Q^*) = 0$ .  $\blacksquare$

These theorems guarantee that through applying our proposed operator repeatedly, any initialization of the  $Q$ -function can converge to a point that is inherently equivalent to the fixed optimum. Even with one homogeneous  $Q$ -function distributed to multiple agents, the  $Q$ -learning can be tractable and sample-efficient. This generalization is beneficial not only to our problem but to other homogeneous multi-agent systems.