

# Mean-Field Multi-Agent Reinforcement Learning: A Decentralized Network Approach

Haotian Gu <sup>\*</sup>    Xin Guo <sup>†</sup>    Xiaoli Wei <sup>‡</sup>    Renyuan Xu <sup>§</sup>

February 22, 2022

## Abstract

One of the challenges for multi-agent reinforcement learning (MARL) is designing efficient learning algorithms for a large system in which each agent has only limited or partial information of the entire system. While exciting progress has been made to analyze decentralized MARL with the *network of agents* for social networks and team video games, little is known theoretically for decentralized MARL with the *network of states* for modeling self-driving vehicles, ride-sharing, and data and traffic routing.

This paper proposes a framework of *localized training and decentralized execution* to study MARL with *network of states*. Localized training means that agents only need to collect local information in their neighboring states during the training phase; decentralized execution implies that agents can execute afterwards the learned decentralized policies, which depend only on agents' current states.

The theoretical analysis consists of three key components: the first is the reformulation of the MARL system as a networked Markov decision process with teams of agents, enabling updating the associated team Q-function in a localized fashion; the second is the Bellman equation for the value function and the appropriate Q-function on the probability measure space; and the third is the exponential decay property of the team Q-function, facilitating its approximation with efficient sample efficiency and controllable error.

The theoretical analysis paves the way for a new algorithm LTDE-NEURAL-AC, where the actor-critic approach with over-parameterized neural networks is proposed. The convergence and sample complexity is established and shown to be scalable with respect to the sizes of both agents and states. To the best of our knowledge, this is the first neural network based MARL algorithm with network structure and provably convergence guarantee.

**Keywords:** Multi-Agent Reinforcement Learning, Mean-field Cooperative Games, Neural Network Approximation

## 1 Introduction

Multi-agent reinforcement learning (MARL) has achieved substantial successes in a broad range of cooperative games and their applications, including coordination of robot swarms (Hüttenrauch et al. [28]), self-driving vehicles (Shalev-Shwartz et al. [48], Cabannes et al. [5]), real-time bidding games (Jin et al. [31]), ride-sharing (Li et al. [35]), power management (Zhou et al. [66]) and traffic routing (El-Tantawy et al. [16]). One of the challenges for the development of MARL is designing efficient learning algorithms for a large system, in which each individual agent has only limited or partial information of the entire system. In such a system, it is necessary to design algorithms to learn policies of the decentralized type, i.e., policies that depend only on the *local* information of each agent.

<sup>\*</sup>Department of Mathematics, University of California, Berkeley, USA. **Email:** haotian\_gu@berkeley.edu

<sup>†</sup>Department of Industrial Engineering & Operations Research, University of California, Berkeley, USA.

**Email:** xinguo@berkeley.edu

<sup>‡</sup>Tsinghua-Berkeley Shenzhen Institute, Shenzhen, China. **Email:** xiaoli\_wei@sz.tsinghua.edu

<sup>§</sup>Industrial & Systems Engineering, University of Southern California, Los Angeles, USA. **Email:** renyuanx@usc.edu

In a simulated or laboratory setting, decentralized policies may be learned in a centralized fashion. It is to train a central controller to dictate the actions of all agents. Such paradigm of *centralized training with decentralized execution* has achieved significant empirical successes, especially with the computational power of deep neural networks (Lowe et al. [40], Foerster et al. [17], Chen et al. [14], Rashid et al. [47], Yang et al. [57], Vadori et al. [52]). Such a training approach, however, suffers from the curse of dimensionality as the computational complexity grows exponentially with the number of agents (Zhang et al. [62]); it also requires extensive and costly communications between the central controller and all agents (Rabbat and Nowak [45]). Moreover, policies derived from the centralized training stage may not be robust in the execution phase (Zhang et al. [60]). Most importantly, this approach has not been supported or analyzed theoretically.

An alternative and promising paradigm is to take into consideration the network structure of the system to train decentralized policies. Compared with the centralized training approach, exploiting network structures makes the training procedure more efficient as it allows the algorithm to be updated with parallel computing and reduces communication cost.

There are two distinct types of network structures. The first is the *network of agents*, often found in social networks such as Facebook and Twitter, as well as team video games including StarCraft II. This network describes *interactions and relations among heterogeneous agents*. For MARL systems with such network of agents, Zhang et al. [63] establishes the asymptotic convergence of decentralized-actor-critic algorithms which are scalable in agent actions. Similar ideas are extended to the continuous space where deterministic policy gradient method (DPG) is used (Zhang et al. [61]), with finite-sample analysis for such framework established in the batch setting (Zhang et al. [64]). Qu et al. [44] studies a network of agents where state and action interact in a local manner; by exploiting the network structure and the exponential decay property of the Q-function, it proposes an actor-critic framework scalable in both actions and states. Similar framework is considered for the linear quadratic case with local policy gradients conducted with zero order optimization and parallel updating (Li et al. [36]).

The second type of network, *the network of states*, has been frequently used for modeling self-driving vehicles, ride-sharing, and data and traffic routing. It focuses on the *state* of agents. Compared with the network of agents which is *static* from agent’s perspective (Sunehag et al. [50]), the network of states is *stochastic*: neighboring agents of any given agent may change dynamically. This type of network has been empirically studied in various applications, including packet routing (You et al. [58]), traffic routing (Calderone and Sastry [7], Guérliau and Dusparic [25]), resource allocations (Cao et al. [8]) and social economic systems (Zheng et al. [65]). However, there is no existing theoretical analysis for this type of decentralized MARL. Moreover, the dynamic nature of agents’ relationship makes it difficult to adopt existing methodology from the static network of agents. The goal of this paper is, therefore, to fill this void.

**Our work.** This paper proposes and studies multi-agent systems with network structure of agent states. In this network, homogeneous agents can move from one state to any connecting state, and observe (realistically) only partial information of the entire system in an aggregated fashion. To study this system, we propose a framework of *localized training and decentralized execution* (LTDE). Localized training means that agents only need to collect local information in their neighboring states during the training phase; decentralized execution implies that, agents can execute afterwards the learned decentralized policies which only require knowledge of agents’ current states.

The theoretical analysis consists of three key components. The first is to regroup these homogeneous agents according to their states and reformulate the MARL system as a networked Markov decision process with teams of agents. This reformulation leads to the decomposition of the Q-function and the value function according to the states, enabling the update of the consequent team Q-function in a localized fashion. The second is to establish the Bellman equation for the value function and the appropriate Q-function on the probability measure space, by utilizing the homogeneity of agents. These functions are invariant with respect to the number of agents. The third is to explore the exponential decay property of the team Q-function, enabling its approximation with a truncated version of a much smaller dimension and

yet with a controllable approximation error.

To design an efficient and scalable reinforcement learning algorithm for such framework, the actor-critic approach with over-parameterized neural networks is adopted. The neural networks, representing decentralized policies and localized Q-functions, are much smaller compared with the global one. The convergence and the sample complexity of the proposed algorithm is established and shown to be scalable with respect to the size of both agents and states. To the best of our knowledge, this is the first neural network based MARL algorithm with network structure and provably convergence guarantee.

**Our contribution.** Our work contributes to two lines of research.

The first one is for mean-field control with reinforcement learning, for which existing works require that each agent have the full information of the population distribution (Gu et al. [24], Carmona et al. [10, 11], Motte and Pham [42]) and yet in most applications agents only have access to partial or limited information (Yang et al. [56]). We build a theoretical framework that incorporates network structures in the MARL framework, and provide computationally efficient algorithms where each agent only needs local information of neighborhood states to learn and to execute the policy.

Secondly, our work builds the theoretical foundation for the practically popular scheme of centralized-training-decentralized-execution (CTDE) (Lowe et al. [40], Rashid et al. [47], Vadori et al. [52], Yang et al. [57]). The CTDE framework is first proposed in Lowe et al. [40] to learn optimal policies in cooperative games with two steps: the first step is to train a global policy for the central controller, and the second step is to decompose the central policy (i.e., a large Q-table) into individual policies so that individual agent can apply the decomposed/decentralized policy after training. Despite the popularity of CTDE, however, there has been no theoretical study as to when the Q-table can be decomposed and when the truncation error can be controlled, except for a heuristic argument by Lowe et al. [40] for large  $N$  with local observations. Our paper analyzes for the first time with theoretical guarantee that applying our algorithm to this CTDE paradigm yields a near-optimal sample complexity, when there is a network structure among agent states. Moreover, our algorithm, which is easier to scale-up, improves the centralized training step with a localized training. To differentiate our approach from the CTDE scheme, we call it localized-training-decentralized-execution (LTDE).

**Notation.** For a set  $\mathcal{X}$ , denote  $\mathbb{R}^{\mathcal{X}} = \{f : \mathcal{X} \rightarrow \mathbb{R}\}$  as the set of all real-valued functions on  $\mathcal{X}$ . For each  $f \in \mathbb{R}^{\mathcal{X}}$ , define  $\|f\|_{\infty} = \sup_{x \in \mathcal{X}} |f(x)|$  as the sup norm of  $f$ . In addition, when  $\mathcal{X}$  is finite, denote  $|\mathcal{X}|$  as the size of  $\mathcal{X}$ , and  $\mathcal{P}(\mathcal{X})$  as the set of all probability measures on  $\mathcal{X}$ :  $\mathcal{P}(\mathcal{X}) = \{p : p(x) \geq 0, \sum_{x \in \mathcal{X}} p(x) = 1\}$ , which is equivalent to the probability simplex in  $\mathbb{R}^{|\mathcal{X}|}$ .  $[N] := \{1, 2, \dots, N\}$ . For any  $\mu \in \mathcal{P}(\mathcal{X})$  and a subset  $\mathcal{Y} \subset \mathcal{X}$ , let  $\mu(\mathcal{Y})$  denote the restriction of the vector  $\mu$  on  $\mathcal{Y}$ , and let  $\mathcal{P}(\mathcal{Y})$  denote the set  $\{\mu(\mathcal{Y}) : \mu \in \mathcal{P}(\mathcal{X})\}$ . For  $x \in \mathbb{R}^d$ ,  $d \in \mathbb{N}$ , denote  $\|x\|_2$  as the  $L^2$ -norm of  $x$  and  $\|x\|_{\infty}$  as the  $L^{\infty}$ -norm of  $x$ .

## 2 Mean-field MARL with Local Dependency

The focus of this paper is to study a cooperative multi-agent system with a network of agent states, which consists of nodes representing states of the agents and edges by which states are connected. In this system, every agent is only allowed to move from her present state to its connecting states. Moreover, she is assumed to only observe (realistically) *partial information* of the system on an aggregated level. Mean-field theory provides efficient approximations when agents only observe aggregated information, and has been applied in stochastic systems with large homogeneous agents such as financial markets (Carmona et al. [9], Lacker and Zariphopoulou [34], Hu and Zariphopoulou [27], Casgrain and Jaimungal [12]), energy markets (Germain et al. [21], Aïd et al. [2]), and auction systems (Iyer et al. [29], Guo et al. [26]).

## 2.1 Review of MARL

Let us first recall the cooperative MARL in an infinite time horizon, where there are  $N$  agents whose policies are coordinated by a central controller. We assume that both the state space  $\mathcal{S}$  and the action space  $\mathcal{A}$  are finite.

At each step  $t = 0, 1, \dots$ , the state of agent  $i$  ( $= 1, 2, \dots, N$ ) is  $s_t^i \in \mathcal{S}$  and she takes an action  $a_t^i \in \mathcal{A}$ . Given the current state profile  $\mathbf{s}_t = (s_t^1, \dots, s_t^N) \in \mathcal{S}^N$  and the current action profile  $\mathbf{a}_t = (a_t^1, \dots, a_t^N) \in \mathcal{A}^N$  of  $N$  agents, agent  $i$  will receive a reward  $r^i(\mathbf{s}_t, \mathbf{a}_t)$  and her state will change to  $s_{t+1}^i$  according to a transition probability function  $P^i(\mathbf{s}_t, \mathbf{a}_t)$ . A Markovian game further restricts the admissible policy for agent  $i$  to be of the form  $a_t^i \sim \pi_t^i(\mathbf{s}_t)$ . That is,  $\pi_t^i : \mathcal{S}^N \rightarrow \mathcal{P}(\mathcal{A})$  maps each state profile  $\mathbf{s} \in \mathcal{S}^N$  to a randomized action, with  $\mathcal{P}(\mathcal{A})$  the space of all probability measures on space  $\mathcal{A}$ .

In this cooperative MARL framework, the central controller is to maximize the expected discounted accumulated reward averaged over all agents. That is to find

$$V(\mathbf{s}) = \sup_{\boldsymbol{\pi}} \frac{1}{N} \sum_{i=1}^N v^i(\mathbf{s}, \boldsymbol{\pi}), \quad (2.1)$$

where

$$v^i(\mathbf{s}, \boldsymbol{\pi}) = \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r^i(\mathbf{s}_t, \mathbf{a}_t) \mid \mathbf{s}_0 = \mathbf{s} \right] \quad (2.2)$$

is the accumulated reward for agent  $i$ , given the initial state profile  $\mathbf{s}_0 = \mathbf{s}$  and policy  $\boldsymbol{\pi} = \{\pi_t\}_{t=0}^{\infty}$  with  $\pi_t = (\pi_t^1, \dots, \pi_t^N)$ . Here  $\gamma \in (0, 1)$  is a discount factor,  $a_t^i \sim \pi_t^i(\mathbf{s}_t)$ , and  $s_{t+1}^i \sim P^i(\mathbf{s}_t, \mathbf{a}_t)$ .

The corresponding Bellman equation for the value function (2.1) is

$$V(\mathbf{s}) = \sup_{\mathbf{a} \in \mathcal{A}^N} \left\{ \mathbb{E} \left[ \frac{1}{N} \sum_{i=1}^N r^i(\mathbf{s}, \mathbf{a}) \right] + \gamma \mathbb{E}_{\mathbf{s}' \sim \mathbf{P}(\mathbf{s}, \mathbf{a})} [V(\mathbf{s}')] \right\}, \quad (2.3)$$

with the population transition kernel  $\mathbf{P} = (P^1, \dots, P^N)$ . The value function can be written as

$$V(\mathbf{s}) = \sup_{\mathbf{a} \in \mathcal{A}^N} Q(\mathbf{s}, \mathbf{a}),$$

in which the Q-function is defined as

$$Q(\mathbf{s}, \mathbf{a}) = \mathbb{E} \left[ \frac{1}{N} \sum_{i=1}^N r^i(\mathbf{s}, \mathbf{a}) \right] + \gamma \mathbb{E}_{\mathbf{s}' \sim \mathbf{P}(\mathbf{s}, \mathbf{a})} [V(\mathbf{s}')], \quad (2.4)$$

consisting of the expected reward from taking action  $\mathbf{a}$  at state  $\mathbf{s}$  and then following the optimal policy thereafter. The Bellman equation for the Q-function, defined from  $\mathcal{S}^N \times \mathcal{A}^N$  to  $\mathbb{R}$ , is given by

$$Q(\mathbf{s}, \mathbf{a}) = \mathbb{E} \left[ \frac{1}{N} \sum_{i=1}^N r^i(\mathbf{s}, \mathbf{a}) \right] + \gamma \mathbb{E}_{\mathbf{s}' \sim \mathbf{P}(\mathbf{s}, \mathbf{a})} \left[ \sup_{\mathbf{a}' \in \mathcal{A}^N} Q(\mathbf{s}', \mathbf{a}') \right]. \quad (2.5)$$

One can thus retrieve the optimal (stationary) control  $\pi^*(\mathbf{s}, \mathbf{a})$  (if it exists) from  $Q(\mathbf{s}, \mathbf{a})$ , with  $\pi^*(\mathbf{s}) \in \arg \max_{\mathbf{a} \in \mathcal{A}^N} Q(\mathbf{s}, \mathbf{a})$ .

## 2.2 Mean-field MARL with Local Dependency

In this system, there are  $N$  agents who share a finite state space  $\mathcal{S}$  and take actions from a finite action space  $\mathcal{A}$ . Moreover, there is a network on the state space  $\mathcal{S}$  associated with an underlying undirected graph  $(\mathcal{S}, \mathcal{E})$ , where  $\mathcal{E} \subset \mathcal{S} \times \mathcal{S}$  is the set of edges. The distance between two nodes is defined as the number of edges in a shortest path. For a given  $s \in \mathcal{S}$ ,  $\mathcal{N}_s^1$  denotes the nearest neighbor of  $s$ , which consists of all nodes connected to  $s$  by an edge and includes  $s$

itself; and  $\mathcal{N}_s^k$  denotes the  $k$ -hop neighborhood of  $s$ , which consists of all nodes whose distance to  $s$  is less than or equal to  $k$ , including  $s$  itself. For simplicity, we use  $\mathcal{N}_s := \mathcal{N}_s^1$ . From agent  $i$ 's perspective, agents in her neighborhood  $\mathcal{N}_{s_t^i}$  change stochastically over time.

To facilitate mean-field approximation to this system, assume throughout the paper that the agents are homogeneous and indistinguishable. In particular, at each step  $t = 0, 1, \dots$ , if agent  $i$  at state  $s_t^i \in \mathcal{S}$  takes an action  $a_t^i \in \mathcal{A}$ , then she will receive a stochastic reward which is uniformly upper bounded by  $r_{\max}$  such that

$$r^i(\mathbf{s}_t, \mathbf{a}_t) := r\left(s_t^i, \mu_t(\mathcal{N}_{s_t^i}), a_t^i\right) \leq r_{\max}, \quad i \in [N]; \quad (2.6)$$

and her state will change to a neighboring state  $s_{t+1}^i \in \mathcal{N}_{s_t^i}$  according to a transition probability such that

$$s_{t+1}^i \sim P^i(\mathbf{s}_t, \mathbf{a}_t) := P\left(\cdot \mid s_t^i, \mu_t(\mathcal{N}_{s_t^i}), a_t^i\right), \quad i \in [N], \quad (2.7)$$

where  $\mu_t(\cdot) = \frac{\sum_{i=1}^N \mathbf{1}(s_t^i = \cdot)}{N} \in \mathcal{P}^N(\mathcal{S}) := \{\mu \in \mathcal{P}(\mathcal{S}) : \mu(s) \in \{0, \frac{1}{N}, \frac{2}{N}, \dots, \frac{N-1}{N}, 1\}\}$  for all  $s \in \mathcal{S}$  is the empirical state distribution of  $N$  agents at time  $t$ , with  $N \cdot \mu_t(s)$  the number of agents in state  $s$  at time  $t$ , and  $\mu_t(\mathcal{N}_{s_t^i})$  the restriction of  $\mu_t$  on the 1-hop neighbor of  $s_t^i$ .

(2.6)-(2.7) indicate that the reward and the transition probability of agent  $i$  at time  $t$  depend on both her individual information ( $a_t^i, s_t^i$ ) and the mean-field of her 1-hop neighborhood  $\mu_t(\mathcal{N}_{s_t^i})$ , in an aggregated yet localized format: *aggregated* or *mean-field* meaning that agent  $i$  depends on other agents only through the empirical state distribution; *localized* meaning that agent  $i$  depends on the mean-field information of her 1-hop neighborhood. Intuitive examples of such a setting include traffic-routing, package delivery, data routing, resource allocations, distributed control of autonomous vehicles and social economic systems.

**Policies with partial information.** To incorporate the element of *partial or limited information* into this mean-field MARL system, consider the following *individual-decentralized policies*

$$a_t^i \sim \pi^i(\mathbf{s}_t) := \pi\left(s_t^i, \mu_t(s_t^i)\right) \in \mathcal{P}(\mathcal{A}), \quad i \in [N], \quad (2.8)$$

and denote  $\mathbf{u}$  as the admissible policy set of all such policies.

Note that for a given mean-field information  $\mu_t$ ,  $\pi(\cdot, \mu_t(\cdot)) : \mathcal{S} \rightarrow \mathcal{P}(\mathcal{A})$  maps the agent state to a randomized action. That is, the policy of each agent is executed in a decentralized manner and assumes that each agent only has access to the population information in her own state. This is more realistic than centralized policies which assume full access to the state information of all agents.

**Value function and Q-function.** The goal for this mean-field MARL is to maximize the expected discounted accumulated reward averaged over all agents, i.e.,

$$V(\mu) := \sup_{\pi \in \mathbf{u}} V^\pi(\mu) = \sup_{\pi \in \mathbf{u}} \frac{1}{N} \sum_{i=1}^N \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t^i, \mu_t(\mathcal{N}_{s_t^i}), a_t^i) \mid \mu_0 = \mu \right], \quad (\text{MF-MARL})$$

subject to (2.6)-(2.8) with a discount factor  $\gamma \in (0, 1)$ .

The mean-field assumption leads to the following definition of the corresponding Q-function for (MF-MARL) on the measure space:

$$\begin{aligned} Q(\mu, h) : &= \underbrace{\mathbb{E} \left[ \sum_{i=1}^N \frac{1}{N} r(s_0^i, \mu(\mathcal{N}_{s_0^i}), a_0^i) \mid \mathbf{s}_0, \mathbf{a}_0 \right]}_{\text{Expected reward of taking } \mathbf{a}_0 = (a_0^1, \dots, a_0^N)} \\ &+ \underbrace{\mathbb{E}_{s_1^i \sim P} \left( \cdot \mid s_0^i, \mu(\mathcal{N}_{s_0^i}), a_0^i \right) \left[ \sum_{t=1}^{\infty} \gamma^t \sum_{i=1}^N \frac{1}{N} r(s_t^i, \mu_t(\mathcal{N}_{s_t^i}), a_t^i) \mid a_t^i \sim \pi_t^* \right]}_{\text{Expected reward of playing optimally thereafter } a_t^i \sim \pi_t^*}, \quad (2.9) \end{aligned}$$

where  $\mu(\cdot) = \frac{\sum_{i=1}^N \mathbf{1}(s_0^i = \cdot)}{N}$  is the initial empirical state distribution and  $h(s)(a) = \frac{\sum_{i=1}^N \mathbf{1}(s_0^i = s, a_0^i = a)}{\sum_{i=1}^N \mathbf{1}(s_0^i = s)}$  is a “decentralized” policy representing the proportion of agents in state  $s$  that takes action  $a$ . Specifically, given  $\mu \in \mathcal{P}^N(\mathcal{S})$ ,  $s \in \mathcal{S}$ , and the  $N \cdot \mu(s)$  agents in state  $s$ ,

$$h(s) \in \mathcal{P}^{N \cdot \mu(s)}(\mathcal{A}) := \left\{ \zeta \in \mathcal{P}(\mathcal{A}) : \zeta(a) \in \left\{ 0, \frac{1}{N \cdot \mu(s)}, \dots, \frac{N \cdot \mu(s) - 1}{N \cdot \mu(s)} \right\} \text{ for all } a \in \mathcal{A} \right\} \subset \mathcal{P}(\mathcal{A}),$$

where  $\zeta$  in  $\mathcal{P}^{N \cdot \mu(s)}(\mathcal{A})$  is an empirical action distribution of  $N \cdot \mu(s)$  agents in state  $s$ , and  $\zeta(a)$  is the proportion of agents taking action  $a \in \mathcal{A}$  among all  $N \cdot \mu(s)$  agents in state  $s$ . Furthermore, for a given  $s \in \mathcal{S}$ , denote  $\mathcal{P}^{N \cdot \mu(s)}(\mathcal{A})$  the set of all admissible “decentralized” policies  $h(s)(\cdot)$ ; and for a given  $\mu \in \mathcal{P}^N(\mathcal{S})$ , denote the product of  $\mathcal{P}^{N \cdot \mu(s)}(\mathcal{A})$  over all states by  $\mathcal{H}^N(\mu) := \{h : h(s) \in \mathcal{P}^{N \cdot \mu(s)}(\mathcal{A}) \forall s \in \mathcal{S}\}$ . Here  $\mathcal{H}^N(\mu)$  depends on  $\mu$  and is a subset of  $\mathcal{H} = \{h : \mathcal{S} \rightarrow \mathcal{P}(\mathcal{A})\}$ .

Note that  $Q(\mu, h)$  defined in (2.9) is invariant with respect to the order of the elements in  $\mathbf{s}_0$  and  $\mathbf{a}_0$ . More critically, the input dimension of the Q-function defined in (2.9) is *independent from the number of agents* in the system, hence is easier to scale up in a large population regime. This differs from the the input dimension of the Q-function in (2.4), which grows exponentially with respect to the number of agents, the main culprit of the curse of dimensionality for MARL algorithms.

### 3 Analysis of MF-MARL with Local Dependency

The theoretical study of this mean-field MARL with local dependency (Section 2.2) consists of three key components, which are crucial for subsequent algorithm design and convergence analysis: the first is the reformulation of the MARL system as a networked Markov decision process with teams of agents. This reformulation leads to the decomposition of the Q-function and the value function according to states, facilitating updating the consequent team Q-function in a localized fashion (Section 3.1); the second is the Bellman equation for the value function and the Q-function on the probability measure space (Section 3.2); the third is the exponential decay property of the team Q-function, enabling its approximation with a truncated version of a much smaller dimension and yet with a controllable approximation error (Section 3.3).

#### 3.1 Markov Decision Process (MDP) on Network of States

This section shows that the mean-field MARL (2.6)-(2.8) can be reformulated in an MDP framework by exploiting the network structure of states. This reformulation leads to the decomposition of the Q-function, facilitating more computationally efficient updates.

The key idea is to utilize the homogeneity of the agents in the problem set-up and to regroup these  $N$  agents according to their states. This regrouping translates (MF-MARL) with  $N$  agents into a networked MDP with  $|\mathcal{S}|$  agents teams, indexed by their states.

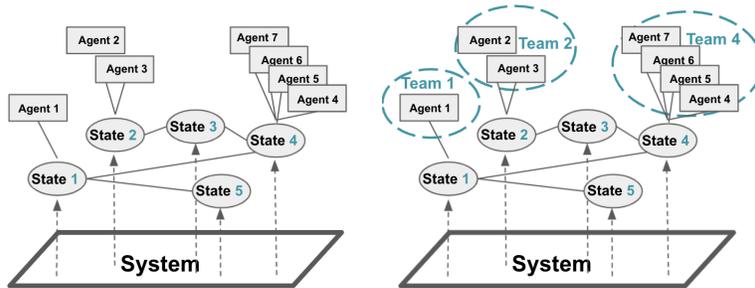


Figure 1: **Left:** MF-MARL problem (2.6)-(2.8). **Right:** Reformulation of team game (3.2)-(3.6).

To see how the policy, the reward function, and the dynamics in this networked Markov decision process are induced by the regrouping approach, recall that there are  $N \cdot \mu(s)$  agents in state  $s$ , each agent  $i$  in state  $s$  will independently choose action  $a_i \sim \pi(s, \mu(s))$  according to the individual-decentralized policy  $\pi(s, \mu(s)) \in \mathcal{P}(\mathcal{A})$  in (2.8). Therefore the empirical action distribution of  $\{a_1, \dots, a_{N \cdot \mu(s)}\}$  is a random variable taking values from  $\mathcal{P}^{N \cdot \mu(s)}(\mathcal{A})$ , the set of empirical action distributions with  $N \cdot \mu(s)$  agents. Moreover, for any  $h(s) \in \mathcal{P}^{N \cdot \mu(s)}(\mathcal{A})$ , we have

$$\begin{aligned} & \mathbb{P} \left( h(s) \text{ is the empirical action distribution of } \{a_1, \dots, a_{N \cdot \mu(s)}\}, a_i \stackrel{i.i.d}{\sim} \pi(s, \mu(s)) \right) \\ &= \mathbb{P} \left( \text{for each } a \in \mathcal{A}, a \text{ appears } N \cdot \mu(s)h(s)(a) \text{ times in } \{a_1, \dots, a_{N \cdot \mu(s)}\}, a_i \stackrel{i.i.d}{\sim} \pi(s, \mu(s)) \right) \\ &= \frac{(N \cdot \mu(s))!}{\prod_{a \in \mathcal{A}} (N \cdot \mu(s)h(s)(a))!} \prod_{a \in \mathcal{A}} \left( \pi(s, \mu(s))(a) \right)^{N \cdot \mu(s)h(s)(a)}. \end{aligned} \quad (3.1)$$

Here  $h(s)(a)$  denotes the proportion of agents taking action  $a$  among all agents in state  $s$ , with last equality derived from the multinomial distribution with parameters  $N \cdot \mu(s)$  and  $\pi(s, \mu(s))$ .

Now, clearly each individual-decentralized policy  $\pi(s, \mu(s)) \in \mathcal{P}(\mathcal{A})$  in (2.8) induces a *team-decentralized policy* of the following form:

$$\Pi_s(h(s) \mid \mu(s)) = \frac{(N \cdot \mu(s))!}{\prod_{a \in \mathcal{A}} (N \cdot \mu(s)h(s)(a))!} \prod_{a \in \mathcal{A}} \left( \pi(s, \mu(s))(a) \right)^{N \cdot \mu(s)h(s)(a)}, \quad (3.2)$$

where  $h(s) \in \mathcal{P}^{N \cdot \mu(s)}(\mathcal{A})$ . Conversely, given a team-decentralized policy  $\Pi_s(\cdot \mid \mu(s))$ , one can recover the individual-decentralized policy  $\pi(s, \mu(s))$  by choosing appropriate  $h(s) \in \mathcal{P}^{N \cdot \mu(s)}(\mathcal{A})$  and querying the value of  $\Pi_s(h(s) \mid \mu(s))$ : let  $h_i(s) = \delta_{a_i}$  be the Dirac measure with  $a_i \in \mathcal{A}$ , which is an action distribution such that all agents in state  $s$  take action  $a_i$ . By (3.2),  $\Pi_s(h_i(s) \mid \mu(s)) = (\pi(s, \mu(s))(a_i))^{N \cdot \mu(s)}$ , implying  $\pi(s, \mu(s))(a_i) = (\Pi(h_i(s) \mid \mu(s)))^{\frac{1}{N \cdot \mu(s)}}$ .

Next, given  $\mu \in \mathcal{P}^N(\mathcal{S})$  and  $h \in \mathcal{H}^N(\mu) = \{h : h(s) \in \mathcal{P}^{N \cdot \mu(s)}(\mathcal{A}), \forall s \in \mathcal{S}\}$ , the set of empirical action distributions on every state, if we define

$$\Pi(h \mid \mu) := \prod_{s \in \mathcal{S}} \Pi_s(h(s) \mid \mu(s)), \quad (3.3)$$

then  $\mathbf{u}$ , the admissible policy set of individual-decentralized policies in the form of (2.8), is now replaced by  $\mathfrak{U}$ , the set of all team-decentralized policies  $\Pi$  induced from  $\pi \in \mathbf{u}$  through (3.2) and (3.3). In addition, denote the set of all state-action distribution pairs as

$$\Xi := \cup_{\mu \in \mathcal{P}^N(\mathcal{S})} \{\zeta = (\mu, h) : h \in \mathcal{H}^N(\mu)\}, \quad (3.4)$$

Moreover, from the team perspective, the transition probability in (2.7) can be viewed as a Markov process of  $\mu_t$  and  $h_t \in \mathcal{H}^N(\mu_t)$  with an induced transition probability  $\mathbf{P}^N$  from (2.7) such that

$$\mu_{t+1} \sim \mathbf{P}^N(\cdot \mid \mu_t, h_t). \quad (3.5)$$

It is easy to verify that for a given state  $s \in \mathcal{S}$ ,  $\mu_{t+1}(s)$  only depends on  $\mu_t(\mathcal{N}_s^2)$ , the empirical distribution in the 2-hop neighborhood of  $s$ , and  $h_t(\mathcal{N}_s)$ .

Finally, given  $\mu(\mathcal{N}_s) \in \mathcal{P}^N(\mathcal{N}_s)$ , an empirical distribution restricted to the 1-hop neighborhood of  $s$ , one can define a *localized team reward function for team  $s$*  from  $\mathcal{P}^{N \cdot \mu(s)}(\mathcal{A})$  to  $\mathbb{R}$  as

$$r_s(\mu(\mathcal{N}_s), h(s)) = \sum_{a \in \mathcal{A}} r(s, \mu(\mathcal{N}_s), a)h(s)(a), \quad (3.6)$$

which depends on the state  $s$  and its 1-hop neighborhood; and define the maximal expected discounted accumulative localized team rewards over all *teams* as

$$\tilde{V}(\mu) := \sup_{\Pi \in \mathfrak{U}} \tilde{V}^\Pi(\mu) = \sup_{\Pi \in \mathfrak{U}} \mathbb{E} \left[ \sum_{t=0}^{\infty} \sum_{s \in \mathcal{S}} \gamma^t r_s(\mu_t(\mathcal{N}_s), h_t(s)) \mid \mu_0 = \mu \right]. \quad (3.7)$$

With all these key elements, one can establish the equivalence between maximizing the reward averaged over all *agents* in (MF-MARL) and maximizing the localized team reward summed over all *teams* in (3.7), and can thus reformulate the (MF-MARL) problem as an equivalent MDP of (3.2)-(3.7) with  $|\mathcal{S}|$  teams, the latter denoted as (MF-DEC-MARL). (The proof is detailed in Appendix A). That is,

**Lemma 3.1** (*Value function and Q-function decomposition*)

$$V(\mu) = \tilde{V}(\mu) = \sup_{\Pi \in \mathfrak{U}} \sum_{s \in \mathcal{S}} \tilde{V}_s^\Pi(\mu), \quad (3.8)$$

where  $h_t \sim \Pi(\cdot | \mu_t)$ ,  $\mu_{t+1} \sim \mathbf{P}^N(\cdot | \mu_t, h_t)$ , and

$$\tilde{V}_s^\Pi(\mu) = \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r_s(\mu_t(\mathcal{N}_s), h_t(s)) \mid \mu_0 = \mu \right] \quad (3.9)$$

is called the value function under policy  $\Pi$  for team  $s$ . Similarly,

$$Q^\Pi(\mu, h) := \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t \sum_{s \in \mathcal{S}} r_s(\mu_t(\mathcal{N}_s), h_t(s)) \mid \mu_0 = \mu, h_0 = h \right] = \sum_{s \in \mathcal{S}} Q_s^\Pi(\mu, h), \quad (3.10)$$

where

$$Q_s^\Pi(\mu, h) = \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r_s(\mu_t(\mathcal{N}_s), h_t(s)) \mid \mu_0 = \mu, h_0 = h \right], \quad (3.11)$$

is the Q-function under policy  $\Pi$  for team  $s$ , called team-decentralized Q-function.

The decomposition for the Q-function in (3.10) is one of the key elements to allow for approximation of  $Q_s^\Pi(\mu, h)$  by a truncated Q-function defined on a smaller space and updated in a localized fashion; it is useful for designing sample-efficient learning algorithms and for parallel computing, as will be clear in the next Section 3.3.

### 3.2 Bellman equation for Q-function.

This section builds the second block for reinforcement learning algorithms, the Bellman equation for Q-function. Indeed, the Bellman equation for  $Q(\mu, h)$  can be derived following a similar argument in Gu et al. [23], after establishing the dynamic programming principle on an appropriate probability measure space.

**Lemma 3.2** (*Bellman Equation for Q-function*) The Q-function defined in (2.9) satisfies:

$$Q(\mu, h) = \mathbb{E} \left[ \sum_{i=1}^N \frac{1}{N} r(s_0^i, \mu(\mathcal{N}_{s_0^i}), a_0^i) \mid \mathbf{s}_0, \mathbf{a}_0 \right] + \gamma \mathbb{E}_{s_1^i \sim P(\cdot | s_0^i, \mu(\mathcal{N}_{s_0^i}), a_0^i)} \left[ \sup_{h' \in \mathcal{H}^N(\mu_1)} Q(\mu_1, h') \right] \quad (3.12)$$

with  $\mu_1(\cdot) = \frac{\sum_{i=1}^N \mathbf{1}(s_1^i = \cdot)}{N}$  the empirical state distribution at time 1.

Note that the Bellman equation (3.12) is for the Q-function defined in (2.9) for general mean-field MARL. In order to enable the *localized-training-decentralized-execution* for computational efficiency, one needs to consider the decomposition of Q-function (3.10) and the updating rule based on the team-decentralized Q-function (3.11). The corresponding Bellman equation for the team-decentralized Q-function (3.11) is:

**Lemma 3.3** Given a policy  $\Pi \in \mathfrak{U}$ ,  $Q_s^\Pi$  defined in (3.11) is the unique solution to the Bellman equation  $Q_s^\Pi = \mathcal{T}_s^\Pi Q_s^\Pi$ , with  $\mathcal{T}_s^\Pi$  the Bellman operator taking the form of

$$\mathcal{T}_s^\Pi Q_s^\Pi(\mu, h) = \mathbb{E}_{\mu' \sim \mathbf{P}^N(\cdot | \mu, h), h' \sim \Pi(\cdot | \mu)} [r_s(\mu, h) + \gamma \cdot Q_s^\Pi(\mu', h')], \forall (\mu, h) \in \Xi. \quad (3.13)$$

These Bellman equations are the basis for general Q-function-based algorithms in mean-field MARL.

### 3.3 Exponential Decay of Q-function

This section will show that the team-decentralized Q-function  $Q_s^\Pi(\mu, h)$  has an *exponential decay* property. This is another key element to enable an approximation to  $Q_s^\Pi$  by a *localized Q-function*  $\widehat{Q}_s^\Pi(\mu(\mathcal{N}_s^k), h(\mathcal{N}_s^k))$ , and to guarantee the scalability and sample efficiency of subsequent algorithm design.

To establish the exponential decay property of the Q-function (3.11), first recall that  $\mathcal{N}_s^k$  is the set of  $k$ -hop neighborhood of state  $s$ , and define  $\mathcal{N}_s^{-k} = \mathcal{S}/\mathcal{N}_s^k$  as the set of states that are outside of  $s$ 'th  $k$ -hop neighborhood. Next, rewrite any given empirical state distribution  $\mu \in \mathcal{P}^N(\mathcal{S})$  as  $(\mu(\mathcal{N}_s^k), \mu(\mathcal{N}_s^{-k}))$ , and similarly,  $h \in \mathcal{H}^N(\mu)$  as  $(h(\mathcal{N}_s^k), h(\mathcal{N}_s^{-k}))$ .

**Definition 3.4** *The  $Q^\Pi$  is said to have  $(c, \rho)$ -exponential decay property, if for any  $s \in \mathcal{S}$  and any  $\Pi \in \mathfrak{U}$ ,  $(\mu, h), (\mu', h') \in \Xi$  with  $\mu(\mathcal{N}_s^k) = \mu'(\mathcal{N}_s^k)$  and  $h(\mathcal{N}_s^k) = h'(\mathcal{N}_s^k)$*

$$\left| Q_s^\Pi(\mu(\mathcal{N}_s^k), \mu(\mathcal{N}_s^{-k}), h(\mathcal{N}_s^k), h(\mathcal{N}_s^{-k})) - Q_s^\Pi(\mu(\mathcal{N}_s^k), \mu'(\mathcal{N}_s^{-k}), h(\mathcal{N}_s^k), h'(\mathcal{N}_s^{-k})) \right| \leq c\rho^{k+1}.$$

Note that the exponential decay property is defined for the team-decentralized Q-function  $Q_s^\Pi$ , instead of the centralized Q-function  $Q^\Pi$ . The following Lemma provides a sufficient condition for the exponential decay property. Its proof is given in Appendix B.

**Lemma 3.5** *When the reward  $r_s$  in (3.6) is uniformly upper bounded by  $r_{max} > 0$ , for any  $s \in \mathcal{S}$ ,  $Q_s^\Pi$  satisfies the  $(\frac{r_{max}}{1-\gamma}, \sqrt{\gamma})$ -exponential decay property.*

The exponential decay property implies that for a given state  $s \in \mathcal{S}$ , the dependence of  $Q_s^\Pi$  on other states decays quickly with respect to its distance from state  $s$ . It motivates and enables the approximation of  $Q_s^\Pi(\mu, h)$  by a truncated function which only depends on  $\mu(\mathcal{N}_s^k)$  and  $h(\mathcal{N}_s^k)$ , especially when  $k$  is large and  $\rho$  is small. Specifically, consider the following class of *localized Q-functions*,

$$\widehat{Q}_s^\Pi(\mu(\mathcal{N}_s^k), h(\mathcal{N}_s^k)) = \sum_{\mu(\mathcal{N}_s^{-k}), h(\mathcal{N}_s^{-k})} \left[ w_s(\mu(\mathcal{N}_s^{-k}), h(\mathcal{N}_s^{-k}); \mu(\mathcal{N}_s^k), h(\mathcal{N}_s^k)) \cdot Q_s^\Pi(\mu(\mathcal{N}_s^k), \mu(\mathcal{N}_s^{-k}), h(\mathcal{N}_s^k), h(\mathcal{N}_s^{-k})) \right],$$

(Local Q-function)

where  $w_s(\mu(\mathcal{N}_s^{-k}), h(\mathcal{N}_s^{-k}); \mu(\mathcal{N}_s^k), h(\mathcal{N}_s^k))$  are any non-negative weights of

$$\sum_{\mu(\mathcal{N}_s^{-k}), h(\mathcal{N}_s^{-k})} w_s(\mu(\mathcal{N}_s^{-k}), h(\mathcal{N}_s^{-k}); \mu(\mathcal{N}_s^k), h(\mathcal{N}_s^k)) = 1$$

for any  $\mu(\mathcal{N}_s^k)$  and  $h(\mathcal{N}_s^k)$ .

Then, direct computation yields the following proposition.

**Proposition 3.6** *Let  $\widehat{Q}_s^\Pi$  be any localized Q-function in the form of (Local Q-function). Assume the  $(c, \rho)$ -exponential decay property in Definition 3.4 holds, then for any  $\mu \in \mathcal{P}^N(\mathcal{S})$  and  $h \in \mathcal{H}^N(\mu)$ ,*

$$\left| \widehat{Q}_s^\Pi(\mu(\mathcal{N}_s^k), h(\mathcal{N}_s^k)) - Q_s^\Pi(\mu, h) \right| \leq c\rho^{k+1}. \quad (3.14)$$

Moreover, (3.14) holds independent of the weights in (Local Q-function).

Note that given a team-decentralized Q-function  $Q_s^\Pi$ , its localized version  $\widehat{Q}_s^\Pi$  only takes  $\mu(\mathcal{N}_s^k), h(\mathcal{N}_s^k)$  as inputs, and  $\widehat{Q}_s^\Pi(\mu(\mathcal{N}_s^k), h(\mathcal{N}_s^k))$  is defined as a weighted average of  $Q_s^\Pi$  over

all  $(\mu, h)$ -pairs which agree with  $(\mu(\mathcal{N}_s^k), h(\mathcal{N}_s^k))$  in the  $k$ -hop neighborhood of  $s$ . Although the localized Q-function  $\widehat{Q}_s^\Pi$  may vary according to different choices of the weights, by the exponential decay property, every  $\widehat{Q}_s^\Pi$  approximates  $Q_s^\Pi$  with uniform error and requires a smaller dimension of input.

**Remark 3.7** (*Exponential Decay Property*) In a discounted reward setting (2.1), the exponential decay property follows directly from the fact that the discount factor  $\gamma \in (0, 1)$  and the local dependency structure in (3.2)-(3.7). For problems of finite-time or infinite horizons with ergodic reward functions, this property can be established by imposing additional Lipschitz condition on the transition kernel. (See Qu et al. [44], Theorem 1 for network of heterogeneous agents and  $\gamma = 1$ ).

It is also worth pointing out that the exponential decay property has been extensively explored in random graphs (e.g., Gamarnik [19], Gamarnik et al. [20]) and for analysis of network of agents in Qu et al. [44] and Lin et al. [37].

## 4 Algorithm Design

The three key analytical components for problem (MF-DEC-MARL) in previous sections pave the way for designing efficient learning algorithms. In this section, we propose and analyze a decentralized neural actor-critic algorithm, called LTDE-NEURAL-AC.

Our focus is the localized Q-function  $\widehat{Q}_s^\Pi(\mu(\mathcal{N}_s^k), h(\mathcal{N}_s^k))$ , the approximation to  $Q_s^\Pi$  with a smaller input dimension. First, this localized Q-function  $\widehat{Q}_s^\Pi$  and the team-decentralized policy  $\Pi_s$  will be parameterized by two-layer neural networks with parameters  $\omega_s$  and  $\theta_s$  respectively (Section 4.2). Next, these neural network parameters  $\theta = \{\theta_s\}_{s \in \mathcal{S}}$  and  $\omega = \{\omega_s\}_{s \in \mathcal{S}}$  are updated via an actor-critic algorithm in a *localized fashion* (Section 4.3): the critic aims to find a proper estimate for the localized Q-function under a fixed policy (parameterized by  $\theta$ ), while the actor computes the policy gradient based on the localized Q-function, and updates  $\theta$  by a gradient step.

These networks are updated locally requiring only information of the neighborhood states during the training phase; afterwards agents in the system will execute these learned *decentralized policies* which requires only information of the agent’s current state. This *localized training and decentralized execution* enables efficient parallel computing especially for a large shared state space.

Moreover, over-parameterization of neural networks avoids issues of nonconvexity and divergence associated with the neural network approach, and ensures the global convergence of our proposed LTDE-NEURAL-AC algorithm.

### 4.1 Basic Set-up

**Policy parameterization.** To start, let us assume that at state  $s$  the *team-decentralized policy*  $\Pi_s^{\theta_s}$  is parameterized by  $\theta_s \in \Theta_s$ . Further denote  $\theta := \{\theta_s\}_{s \in \mathcal{S}}$ ,  $\Theta := \prod_{s \in \mathcal{S}} \Theta_s$ ,  $\Pi^\theta := \prod_{s \in \mathcal{S}} \Pi_s^{\theta_s}$ , and  $\mathbf{\Pi} := \{\Pi^\theta : \theta \in \Theta\}$  as the class of admissible policies parameterized by the parameter space  $\{\theta : \theta \in \Theta\}$ .

**Initialization.** Let us also assume that the initial state distribution  $\mu_0$  of  $N$  agents is sampled from a given distribution  $P_0$  over  $\mathcal{P}^N(\mathcal{S})$ , i.e.,  $\mu_0 \sim P_0$ ; and define the expected total reward function  $J(\theta)$  under policy  $\Pi^\theta$  by

$$J(\theta) = \mathbb{E}_{\mu_0 \sim P_0}[\widetilde{V}^{\Pi^\theta}(\mu_0)]. \quad (4.1)$$

**Visitation measure.** Denote  $\nu_\theta$  as the stationary distribution on  $\Xi$  of the Markov process (3.5) induced by  $\Pi^\theta$ .

Similar to the single-agent RL problem (Agarwal et al. [1], Fu et al. [18]), each admissible policy  $\Pi^\theta$  induces a visitation measure  $\sigma_\theta(\mu, h)$  on  $\Xi$  describing the frequency that policy  $\Pi^\theta$  visits  $(\mu, h)$ , with

$$\sigma_\theta(\mu, h) := (1 - \gamma) \cdot \sum_{t=0}^{\infty} \gamma^t \cdot \mathbb{P}(\mu_t = \mu, h_t = h \mid \Pi^\theta), \quad (4.2)$$

where  $\mu_0 \sim P_0$ ,  $h_t \sim \Pi^\theta(\cdot \mid \mu_t)$ , and  $\mu_{t+1} \sim \mathbf{P}^N(\cdot \mid \mu_t, h_t)$ .

**Policy gradient theorem.** In order to find the optimal parameterized policy  $\Pi^\theta$  which maximizes the expected total reward function  $J(\theta)$ , the policy optimization step will search for  $\theta \in \Theta$  along the gradient direction  $\nabla J(\theta)$ . Note that computing the gradient  $\nabla J(\theta)$  depends on both the action selection, which is directly determined by  $\Pi^\theta$ , and the visitation measure  $\sigma_\theta$  in (4.2), which is indirectly determined by  $\Pi^\theta$ .

A simple and elegant result called the policy gradient theorem (Lemma 4.1) proposed in Sutton et al. [51], reformulates the gradient  $\nabla J(\theta)$  in terms of  $Q^{\Pi^\theta}$  in (3.10) and  $\nabla \log \Pi^\theta(h \mid \mu)$  under the visitation measure  $\sigma_\theta$ . This result simplifies the gradient computation significantly, and is fundamental for actor-critic algorithms.

**Lemma 4.1** (Sutton et al. [51])  $\nabla J(\theta) = \frac{1}{1-\gamma} \mathbb{E}_{\sigma_\theta} \left[ Q^{\Pi^\theta}(\mu, h) \nabla \log \Pi^\theta(h \mid \mu) \right]$ .

Now, direct implementation of the actor-critic algorithm with the *centralized* policy gradient theorem in Lemma 4.1 suffers from high sample complexity due to the dimension of the Q-function. Instead, we will show that the exponential decay property of Q-function allows efficient approximation of the policy gradient via *localization* and hence a *scalable* algorithm to solve (MF-MARL).

## 4.2 Neural Policy and Neural Q-function

We now turn to the localized Q-function  $\widehat{Q}_s^\Pi(\mu(\mathcal{N}_s^k), h(\mathcal{N}_s^k))$  (i.e., the approximation of  $Q_s^\Pi$ ) and the team-decentralized policy  $\Pi_s$ , and their parameterization by two-layer neural networks. We emphasize that the parameterization framework in this section can be extended to any neural-based single-agent algorithms with convergence guarantee.

**Two-Layer Neural Network.** For any input space  $\mathcal{X} \subset \mathbb{R}^{d_x}$  with dimension  $d_x \in \mathbb{N}$ , a two-layer neural network  $\tilde{f}(x; W, b)$  with input  $x \in \mathcal{X}$  and width  $M \in \mathbb{N}$  takes the form of

$$\tilde{f}(x; W, b) = \frac{1}{\sqrt{M}} \sum_{m=1}^M b_m \cdot \text{ReLU}(x \cdot [W]_m). \quad (4.3)$$

Here the scaling factor  $\frac{1}{\sqrt{M}}$  called the *Xavier initialization* (Glorot and Bengio [22]) ensures the same input variance and the same gradient variance for all layers; the activation function  $\text{ReLU} : \mathbb{R} \rightarrow \mathbb{R}$ , defined as  $\text{ReLU}(u) = \mathbf{1}\{u > 0\} \cdot u$ ;  $b = \{b_m\}_{m \in [M]}$  and  $W = ([W]_1^\top, \dots, [W]_M^\top)^\top \in \mathbb{R}^{M \times d_x}$  in (4.3) are parameters of the neural network.

Taking advantage of the homogeneity of ReLU (i.e.,  $\text{ReLU}(c \cdot u) = c \cdot \text{ReLU}(u)$  for all  $c > 0$  and  $u \in \mathbb{R}$ ), we adopt the usual trick (Cai et al. [6], Wang et al. [53], Allen-Zhu et al. [3]) to fix  $b$  throughout the training and only to update  $W$  in the sequel. Consequently, denote  $\tilde{f}(x; W, b)$  as  $f(x; W)$  when  $b_m = 1$  is fixed.  $[W]_m$  is initialized according to a multivariate normal distribution  $N(0, I_{d_x}/d_x)$ , where  $I_{d_x}$  is the identity matrix of size  $d_x$ .

**Neural Policy.** For each  $s \in \mathcal{S}$ , denote the tuple  $\zeta_s = (\mu(s), h(s)) \in \mathbb{R}^{d_{\zeta_s}}$  for notational simplicity, where  $d_{\zeta_s} := 1 + |\mathcal{A}|$  is the dimension of  $\zeta_s$ . Given the input  $\zeta_s = (\mu(s), h(s))$  and

parameter  $W = \theta_s$  in the two-layer neural network  $f(\cdot; \theta_s)$  in (4.3), the team-decentralized policy  $\Pi_s^{\theta_s}$ , called the *actor*, is parameterized in the form of an *energy-based policy*,

$$\Pi_s^{\theta_s}(h(s) | \mu(s)) = \frac{\exp[\tau \cdot f((\mu(s), h(s)); \theta_s)]}{\sum_{h'(s) \in \mathcal{P}^{N \cdot \mu(s)}(\mathcal{A})} \exp[\tau \cdot f((\mu(s), h'(s)); \theta_s)]}, \quad (4.4)$$

where  $\tau$  is the temperature parameter and  $f$  is the energy function.

To study the policy gradient for (4.4), let us first define a class of feature mappings that is consistent with the representation of two-layer neural networks. This connection between the gradient of a two-layer ReLU neural network and the feature mapping defined in (4.6) is crucial in the convergence analysis of Theorems 5.4 and 5.10. Specifically, rewrite the two-layer neural network in (4.3) as

$$f(\zeta_s; \theta_s) = \frac{1}{\sqrt{M}} \sum_{m=1}^M \text{ReLU}(\zeta_s^\top [\theta_s]_m) = \frac{1}{\sqrt{M}} \sum_{m=1}^M \mathbb{1}\{\zeta_s^\top [\theta_s]_m > 0\} \cdot \zeta_s^\top [\theta_s]_m := \phi_{\theta_s}(\zeta_s)^\top \theta_s. \quad (4.5)$$

Then the feature mapping  $\phi_{\theta_s} = \left( [\phi_{\theta_s}]_1^\top, \dots, [\phi_{\theta_s}]_M^\top \right)^\top : \mathbb{R}^{d_{\zeta_s}} \rightarrow \mathbb{R}^{M \times d_{\zeta_s}}$  may take the following form:

$$[\phi_{\theta_s}]_m(\zeta_s) = \frac{1}{\sqrt{M}} \cdot \mathbb{1}\{\zeta_s^\top [\theta_s]_m > 0\} \cdot \zeta_s. \quad (4.6)$$

That is, the two-layer neural network  $f(\zeta_s; \theta_s)$  may be viewed as the inner product between the feature  $\phi_{\theta_s}(\zeta_s)$ , and the neural network parameters  $\theta_s$ . Since  $f(\zeta_s; \theta_s)$  is almost everywhere differentiable with respect to  $\theta_s$ , we see  $\nabla_{\theta_s} f(\zeta_s; \theta_s) = \phi_{\theta_s}(\zeta_s)$ .

Furthermore, define a ‘‘centered’’ version of the feature  $\phi_{\theta_s}$  such that

$$\Phi(\theta, s, \mu, h) := \phi_{\theta_s}(\mu(s), h(s)) - \mathbb{E}_{h(s)' \sim \Pi_s^{\theta_s}(\cdot | \mu(s))} [\phi_{\theta_s}(\mu(s), h'(s))]. \quad (4.7)$$

Note that when policy  $\Pi^\theta$  takes the energy-based form (4.4),  $\Phi = \frac{1}{\tau} \nabla_\theta \log \Pi^\theta$ . Therefore,

**Lemma 4.2** *For any  $\theta \in \Theta$ ,  $s \in \mathcal{S}$ ,  $\mu \in \mathcal{P}^N(\mathcal{S})$  and  $h \in \mathcal{H}^N(\mu)$ ,  $\|\Phi(\theta, s, \mu, h)\|_2 \leq 2$ , and*

$$\nabla_{\theta_s} J(\theta) = \frac{\tau}{1-\gamma} \cdot \mathbb{E}_{\sigma_\theta} \left[ Q^{\Pi^\theta}(\mu, h) \cdot \Phi(\theta, s, \mu, h) \right]. \quad (4.8)$$

Moreover, for each  $s \in \mathcal{S}$ , define the following localized policy gradient

$$g_s(\theta) = \frac{\tau}{1-\gamma} \mathbb{E}_{\sigma_\theta} \left[ \left[ \sum_{y \in \mathcal{N}_s^k} \widehat{Q}_y^{\Pi^\theta}(\mu(\mathcal{N}_y^k), h(\mathcal{N}_y^k)) \right] \cdot \Phi(\theta, s, \mu, h) \right], \quad (4.9)$$

with  $\widehat{Q}_s^{\Pi^\theta}$  in (Local Q-function) satisfying the  $(c, \rho)$ -exponential decay property, then there exists a universal constant  $c_0 > 0$  such that

$$\|g_s(\theta) - \nabla_{\theta_s} J(\theta)\| \leq \frac{c_0 \tau |\mathcal{S}|}{1-\gamma} \rho^{k+1}. \quad (4.10)$$

**Neural Q-function.** Note  $\widehat{Q}_s^{\Pi^\theta}$  in (Local Q-function) is unknown *a priori*. To obtain the localized policy gradient (4.9), the neural network (4.3) to parameterize  $\widehat{Q}_s^{\Pi^\theta}$  is taken as:

$$Q_s(\mu(\mathcal{N}_s^k), h(\mathcal{N}_s^k); \omega_s) = f((\mu(\mathcal{N}_s^k), h(\mathcal{N}_s^k)); \omega_s).$$

This  $Q_s$  is called the *critic*. For simplicity, denote  $\zeta_s^k = (\mu(\mathcal{N}_s^k), h(\mathcal{N}_s^k))$ , with  $d_{\zeta_s^k}$  the dimension of  $\zeta_s^k$ .

### 4.3 Actor-Critic

**Critic Update.** For a fixed policy  $\Pi^\theta$ , it is to estimate  $\widehat{Q}_s^{\Pi^\theta}$  of (Local Q-function) by a two-layer neural network  $Q_s(\cdot; \omega_s)$ , where  $\widehat{Q}_s^{\Pi^\theta}$  serves as an approximation to the team-decentralized Q-function  $Q_s^{\Pi^\theta}$ .

To design the update rule for  $\widehat{Q}_s^{\Pi^\theta}$ , note that the Bellman equation (3.13) is for  $Q_s^{\Pi^\theta}$  instead of  $\widehat{Q}_s^{\Pi^\theta}$ . Indeed,  $Q_s^{\Pi^\theta}$  takes  $(\mu, h)$  as the input while  $\widehat{Q}_s^{\Pi^\theta}$  takes the partial information  $(\mu(\mathcal{N}_s^k), h(\mathcal{N}_s^k))$  as the input.

In order to update parameter  $\omega_s$ , we substitute  $(\mu(\mathcal{N}_s^k), h(\mathcal{N}_s^k))$  for the state-action pair in the Bellman equation (3.13). It is therefore necessary to study the error of using  $(\mu(\mathcal{N}_s^k), h(\mathcal{N}_s^k))$  as the input. Specifically, given a tuple  $(\mu_t, h_t, r_s(\mu_t(\mathcal{N}_s), h_t(s)), \mu_{t+1}, h_{t+1})$  sampled from the stationary distribution  $\nu_\theta$  of adopting policy  $\Pi^\theta$ , the parameter  $\omega_s$  will be updated to minimize the error:

$$(\delta_{s,t})^2 = [Q_s(\mu_t(\mathcal{N}_s^k), h_t(\mathcal{N}_s^k); \omega_s) - r_s(\mu_t(\mathcal{N}_s), h_t(s)) - \gamma \cdot Q_s(\mu_{t+1}(\mathcal{N}_s^k), h_{t+1}(\mathcal{N}_s^k); \omega_s)]^2.$$

Estimating  $\delta_{s,t}$  depends only on  $\mu_t(\mathcal{N}_s^k), h_t(\mathcal{N}_s^k)$  and can be *collected locally*. (See Theorem 5.4).

The neural critic update takes the iterative forms of

$$\omega_s(t+1/2) \leftarrow \omega_s(t) - \eta_{\text{critic}} \cdot \delta_{s,t} \cdot \nabla_{\omega_s} Q_s(\mu_t(\mathcal{N}_s^k), h_t(\mathcal{N}_s^k); \omega_s), \quad (4.11)$$

$$\omega_s(t+1) \leftarrow \arg \min_{\omega \in \mathcal{B}_s^{\text{critic}}} \|\omega - \omega_s(t+1/2)\|_2, \quad (4.12)$$

$$\bar{\omega}_s \leftarrow (t+1)/(t+2) \cdot \bar{\omega}_s + 1/(t+2) \cdot \omega_s(t+1), \quad (4.13)$$

in which  $\eta_{\text{critic}}$  is the learning rate. Here (4.11) is the stochastic semigradient step, (4.12) is a projection to the parameter space  $\mathcal{B}_s^{\text{critic}} := \{\omega_s \in \mathbb{R}^{M \times d_{\zeta_s^k}} : \|\omega_s - \omega_s(0)\|_\infty \leq R/\sqrt{M}\}$  for some  $R > 0$ , and (4.13) is the averaging step. This critic update is summarized in Algorithm 1.

---

#### Algorithm 1 Localized-Training-Decentralized-Execution Neural Temporal Difference

---

- 1: **Input:** Width of the neural network  $M$ , radius of the constraint set  $R$ , number of iterations  $T_{\text{critic}}$ , policy  $\Pi^\theta$ , learning rate  $\eta_{\text{critic}}$ , localization parameter  $k$ .
  - 2: **Initialize:** For all  $m \in [M]$  and  $s \in \mathcal{S}$ , sample  $b_m \sim \text{Unif}(\{-1, 1\})$ ,  $[\omega_s(0)]_m \sim N(0, Id_{\zeta_s^k}/d_{\zeta_s^k})$ ,  $\bar{\omega}_s = \omega_s(0)$ .
  - 3: **for**  $t = 0$  to  $T_{\text{critic}} - 2$  **do**
  - 4:   Sample  $(\mu_t, h_t, \{r_s(\mu_t(\mathcal{N}_s), h_t(s))\}_{s \in \mathcal{S}}, \mu_{t+1}, h_{t+1})$  from the stationary distribution  $\nu_\theta$  of  $\Pi^\theta$ .
  - 5:   **for**  $s \in \mathcal{S}$  **do**
  - 6:     Denote  $\zeta_{s,t}^k = (\mu_t(\mathcal{N}_s^k), h_t(\mathcal{N}_s^k))$ ,  $\zeta_{s,t}^{k'} = (\mu_{t+1}(\mathcal{N}_s^k), h_{t+1}(\mathcal{N}_s^k))$ .
  - 7:     Residual calculation:  $\delta_{s,t} \leftarrow Q_s(\zeta_{s,t}^k; \omega_s(t)) - r_s(\mu_t(\mathcal{N}_s), h_t(s)) - \gamma \cdot Q_s(\zeta_{s,t}^{k'}; \omega_s(t))$ .
  - 8:     Temporal difference update:
  - 9:      $\omega_s(t+1/2) \leftarrow \omega_s(t) - \eta_{\text{critic}} \cdot \delta_{s,t} \cdot \nabla_{\omega_s} Q_s(\zeta_{s,t}^k; \omega_s(t))$ .
  - 10:     Projection onto the parameter space:  $\omega_s(t+1) \leftarrow \arg \min_{\omega \in \mathcal{B}_s^{\text{critic}}} \|\omega - \omega_s(t+1/2)\|_2$ .
  - 11:     Averaging the output:  $\bar{\omega}_s \leftarrow \frac{t+1}{t+2} \cdot \bar{\omega}_s + \frac{1}{t+2} \cdot \omega_s(t+1)$ .
  - 12:   **end for**
  - 13: **end for**
  - 14: **Output:**  $Q_s(\cdot; \bar{\omega}_s), \forall s \in \mathcal{S}$ .
- 

**Actor Update.** At the iteration step  $t$ , a neural network estimation  $Q_s(\cdot; \bar{\omega}_s)$  is given for the localized Q-function  $\widehat{Q}_s^{\Pi^{\theta(t)}}$  under the current policy  $\Pi^{\theta(t)}$ . Let  $\{(\mu_l, h_l)\}_{l \in [B]}$  be samples from the state-action visitation measure  $\sigma_{\theta(t)}$  of (4.2), and define an estimator  $\widehat{\Phi}(\theta, s, \mu_l, h_l)$  of  $\Phi(\theta, s, \mu_l, h_l)$  in (4.7):

$$\widehat{\Phi}(\theta, s, \mu_l, h_l) = \phi_{\theta_s}(\mu_l(s), h_l(s)) - \mathbb{E}_{\Pi_s^{\theta_s}}[\phi_{\theta_s}(\mu_l(s), h'(s))].$$

By Lemma 4.2, one can compute the following estimator of  $g_s(\theta(t))$  defined in (4.9),

$$\widehat{g}_s(\theta(t)) = \frac{\tau}{(1-\gamma)B} \sum_{l \in [B]} \left[ \sum_{y \in \mathcal{N}_s^k} Q_y(\mu_l(\mathcal{N}_y^k), h_l(\mathcal{N}_y^k); \bar{\omega}_y) \right] \cdot \widehat{\Phi}(\theta(t), s, \mu_l, h_l). \quad (4.14)$$

This estimator  $\widehat{g}_s$  in (4.14) only depends locally on  $\{(\mu_l, h_l)\}_{l \in [B]}$ . Hence  $\widehat{g}$  and  $\widehat{\Phi}$  can be computed in a *localized fashion* after the samples are collected. Similar to the critic update,  $\theta_s(t)$  is updated by performing a gradient step with  $\widehat{g}_s$ , and then projected onto the parameter space  $\mathcal{B}_s^{\text{actor}} := \{\theta_s \in \mathbb{R}^{M \times d_{\zeta_s}} : \|\theta_s - \theta_s(0)\|_{\infty} \leq R/\sqrt{M}\}$ .

This actor update is summarized in Algorithm 2.

**Sampling from  $\nu_{\theta}$  and the Visitation Measure  $\sigma_{\theta}$ .** In Algorithms 1 and 2, it is assumed that one can sample independently from the stationary distribution  $\nu_{\theta}$  and the visitation measure  $\sigma_{\theta}$ , respectively. Such an assumption of sampling from  $\nu_{\theta}$  can be relaxed by either sampling from a rapidly-mixing Markov chain mixing, with weakly-dependent sequence of samples (Bhandari et al. [4]), or by randomly picking samples from replay buffers consisting of long trajectories, with reduced correlation between samples.

To sample from the visitation measure  $\sigma_{\theta}$  and computing the unbiased policy gradient estimator, Konda and Tsitsiklis [33] suggests introducing a new MDP such that the next state is sampled from the transition probability with probability  $\gamma$ , and from the initial distribution with probability  $1 - \gamma$ . Then the stationary distribution of this new MDP is exactly the visitation measure. Alternatively, Liu et al. [39] proposes an importance-sampling-based algorithm which enables off-policy evaluation with low variance.

---

#### Algorithm 2 Localized-Training-Decentralized-Execution Neural Actor-Critic

---

- 1: **Input:** Width of the neural network  $M$ , radius of the constraint set  $R$ , number of iterations  $T_{\text{actor}}$  and  $T_{\text{critic}}$ , learning rate  $\eta_{\text{actor}}$  and  $\eta_{\text{critic}}$ , temperature parameter  $\tau$ , batch size  $B$ , localization parameter  $k$ .
  - 2: **Initialize:** For all  $m \in [M]$  and  $s \in \mathcal{S}$ , sample  $b_m \sim \text{Unif}(\{-1, 1\})$ ,  $[\theta_s(0)]_m \sim N(0, I_{d_{\zeta_s}}/d_{\zeta_s})$ .
  - 3: **for**  $t = 1$  to  $T_{\text{actor}}$  **do**
  - 4:   Define the policy
 
$$\Pi^{\theta}(h | \mu) := \prod_{s \in \mathcal{S}} \Pi_s^{\theta_s}(h(s) | \mu(s)) = \prod_{s \in \mathcal{S}} \frac{\exp[\tau \cdot f((\mu(s), h(s)); \theta_s)]}{\sum_{h'(s) \in \mathcal{H}^N} \exp[\tau \cdot f((\mu(s), h'(s)); \theta_s)]}.$$
  - 5:   Output  $Q_s(\cdot; \bar{\omega}_s)$  using Algorithm 1 with the inputs: policy  $\Pi^{\theta}$ , width of the neural network  $M$ , radius of the constraint set  $R$ , number of iterations  $T_{\text{critic}}$ , learning rate  $\eta_{\text{critic}}$  and localization parameter  $k$ .
  - 6:   Sample  $\{\mu_l, h_l\}_{l \in [B]}$  from the state-action visitation measure  $\sigma_{\theta}$  (4.2) of  $\Pi^{\theta}$ .
  - 7:   **for**  $s \in \mathcal{S}$  **do**
  - 8:     Compute the local gradient estimator  $\widehat{g}_s(\theta(t))$  using (4.14).
  - 9:     Policy update:  $\theta_s(t+1/2) \leftarrow \theta_s(t) + \eta_{\text{actor}} \cdot \widehat{g}_s(\theta(t))$
  - 10:    Projection onto the parameter space:  $\theta_s(t+1) \leftarrow \arg \min_{\theta \in \mathcal{B}_s^{\text{actor}}} \|\theta - \theta_s(t+1/2)\|_2$ .
  - 11:   **end for**
  - 12: **end for**
  - 13: **Output:**  $\{\Pi^{\theta(t)}\}_{t \in [T_{\text{actor}}]}$ .
- 

## 5 Convergence of the Critic and Actor Updates

We now establish the global convergence for LTDE-NEURAL-AC proposed in Section 4.

**Convergence of the Critic Update.** The convergence of the decentralized neural critic update in Algorithm 1 relies on the following assumptions.

**Assumption 5.1** (*Action-Value Function Class*) For each  $s \in \mathcal{S}$ ,  $k \in \mathbb{N}$ , define

$$\mathcal{F}_{R,\infty}^{s,k} = \left\{ f(\zeta_s^k) = Q_s(\zeta_s^k; \omega_s(0)) + \int \mathbf{1}\{v^\top \zeta_s^k > 0\} \cdot (\zeta_s^k)^\top \iota(v) d\mu(v) : \|\iota(v)\|_\infty \leq R \right\},$$

with  $\mu : \mathbb{R}^{d_{\zeta_s^k}} \rightarrow \mathbb{R}$  the density function of Gaussian distribution  $N(0, I_{d_{\zeta_s^k}}/d_{\zeta_s^k})$  and  $Q_s(\zeta_s^k; \omega_s(0))$  the two-layer neural network under the initial parameter  $\omega_s(0)$ . We assume that  $\widehat{Q}_s^{\Pi^\theta} \in \mathcal{F}_{R,\infty}^{s,k}$ .

**Assumption 5.2** (*Regularity of  $\nu_\theta$  and  $\sigma_\theta$* ) There exists a universal constant  $c_0 > 0$  such that for any policy  $\Pi^\theta$ , any  $\alpha \geq 0$ , and any  $v \in \mathbb{R}^{d_\zeta}$  with  $\|v\|_2 = 1$ , the stationary distribution  $\nu_\theta$  and the state visitation measure  $\sigma_\theta$  satisfy

$$\mathbb{P}_{\zeta \sim \nu_\theta} (|v^\top \zeta| \leq \alpha) \leq c_0 \cdot \alpha, \quad \mathbb{P}_{\zeta \sim \sigma_\theta} (|v^\top \zeta| \leq \alpha) \leq c_0 \cdot \alpha.$$

**Remark 5.3** Both Assumption 5.1 and Assumption 5.2 are similar to the standard assumptions in the analysis of single-agent neural actor-critic algorithms (Cai et al. [6], Liu et al. [38], Wang et al. [53], Cayci et al. [13]).

In particular, Assumption 5.1 is a regularity condition for  $\widehat{Q}_s^{\Pi^\theta}$  in (Local Q-function). Here  $\mathcal{F}_{R,\infty}^{s,k}$  is a subset of the reproducing kernel Hilbert space (RKHS) induced by the random feature  $\mathbf{1}\{v^\top \zeta_s^k > 0\} \cdot (\zeta_s^k)$  with  $v \sim N(0, I_{d_{\zeta_s^k}}/d_{\zeta_s^k})$  up to the shift of  $Q_s(\zeta_s^k; \omega_s(0))$  (Rahimi and Recht [46]). This RKHS is dense in the space of continuous functions on any compact set (Micchelli et al. [41], Ji et al. [30]). (See also Section D.1.1 for details of the connection between  $\mathcal{F}_{R,\infty}^{s,k}$  and the linearizations of two-layer neural networks (D.4)).

Assumption 5.2 holds when  $\sigma_\theta$  and  $\nu_\theta$  have uniformly upper bounded probability densities (Cai et al. [6]).

**Theorem 5.4** (*Convergence of Critic Update*) Assume Assumptions 5.1 and 5.2. Set  $T_{\text{critic}} = \Omega(M)$  and  $\eta_{\text{critic}} = \min\{(1-\gamma)/8, (T_{\text{critic}})^{-1/2}\}$  in Algorithm 1. Then  $Q_s(\cdot; \bar{\omega}_s)$  generated by Algorithm 1 satisfies

$$\mathbb{E}_{\text{init}} \left[ \left\| Q_s(\cdot; \bar{\omega}_s) - Q_s^{\Pi^\theta}(\cdot) \right\|_{L^2(\nu_\theta)}^2 \right] \leq \mathcal{O} \left( \frac{R^3 d_{\zeta_s^k}^{3/2}}{M^{1/2}} + \frac{R^{5/2} d_{\zeta_s^k}^{5/4}}{M^{1/4}} + \frac{r_{\max}^2 \gamma^{k+1}}{(1-\gamma)^2} \right), \quad (5.1)$$

where  $\|f\|_{L^2(\nu_\theta)} := (\mathbb{E}_{\zeta \sim \nu_\theta} [f(\zeta)^2])^{1/2}$ , and the expectation (5.1) is taken with respect to the random initialization.

Theorem 5.4 indicates the trade-off between the approximation-optimization error and the localization error. The first two terms in (5.1) correspond to the neural network approximation-optimization error, similar to the single-agent case (Cai et al. [6], Cayci et al. [13]). This approximation-optimization error decreases when the width of the hidden layer  $M$  increases. Meanwhile, the last term in (5.1) represents the additional error from using the localized information in (4.11), unique for the mean-field MARL case. This localization error and  $\gamma^k$  decrease as the number of truncated neighborhood  $k$  increases, with more information from a larger neighborhood used in the update. However, the input dimension  $d_{\zeta_s^k}$  and the approximation-optimization error will increase if the dimension of the problem increases.

In particular, for a relatively sparse network on  $\mathcal{S}$ , one can choose  $k \ll |\mathcal{S}|$  hence  $d_{\zeta_s^k} \ll d_\zeta$ , and Theorem 5.4 indicates the superior performance of the localized training scheme in efficiency over directly approximating the centralized Q-function.

Proof of Theorem 5.4 is presented in Section D.1.

**Convergence of the Actor Update.** This section establishes the global convergence of the actor update. The convergence analysis consists of two steps. The first step proves the convergence to a stationary point  $\tilde{\theta}$ ; the second step controls the gap between the stationary point  $\theta$  and the optimality  $\theta^*$  in the overparametrization regime. The convergence is built under the following assumptions and definition.

**Assumption 5.5** (*Variance Upper Bound*) For every  $t \in [T_{actor}]$  and  $s \in \mathcal{S}$ , denote  $\xi_s(t) = \hat{g}_s(\theta(t)) - \mathbb{E}[\hat{g}_s(\theta(t))]$  with  $\hat{g}_s(\theta(t))$  defined in (4.14). Assume there exists  $\Sigma > 0$  such that  $\mathbb{E}[\|\xi_s(t)\|_2^2] \leq \tau^2 \Sigma^2 / B$ . Here the expectations are taken over  $\sigma_{\theta(t)}$  given  $\{\bar{\omega}_s\}_{s \in \mathcal{S}}$ .

**Assumption 5.6** (*Regularity Condition on  $\sigma_\theta$  and  $\nu_\theta$* ) There exists an absolute constant  $D > 0$  such that for every  $\Pi^\theta$ , the stationary distribution  $\nu_\theta$  and the state-action visitation measure  $\sigma_\theta$  satisfy

$$\left\{ \mathbb{E}_{\nu_\theta} \left[ (d\sigma_\theta / d\nu_\theta(\mu, h))^2 \right] \right\} \leq D^2,$$

where  $d\sigma_\theta / d\nu_\theta$  is the Radon-Nikodym derivative of  $\sigma_\theta$  with respect to  $\nu_\theta$ .

**Assumption 5.7** (*Lipschitz Continuous Policy Gradient*) There exists an absolute constant  $L > 0$ , such that  $\nabla_\theta J(\theta)$  is  $L$ -Lipschitz continuous with respect to  $\theta$ , i.e., for all  $\theta_1, \theta_2$ ,

$$\|\nabla_\theta J(\theta_1) - \nabla_\theta J(\theta_2)\|_2 \leq L \cdot \|\theta_1 - \theta_2\|_2.$$

**Definition 5.8**  $\tilde{\theta} \in \mathcal{B}^{actor}$  is called a stationary point of  $J(\theta)$  if for all  $\tilde{\theta} \in \mathcal{B}^{actor}$ ,

$$\nabla_\theta J(\tilde{\theta})^\top (\theta - \tilde{\theta}) \leq 0. \quad (5.2)$$

**Assumption 5.9** (*Policy Function Class*) Define a function class

$$\mathcal{F}_{R,\infty} = \left\{ f(\zeta) = \sum_{s \in \mathcal{S}} \left[ \phi_{\theta_s(0)}(\zeta_s)^\top \theta_s(0) + \int \mathbb{1}\{v^\top \zeta_s > 0\} \cdot (\zeta_s)^\top \iota(v) d\mu(v) \right] : \|\iota(v)\|_\infty \leq R \right\}$$

where  $\mu : \mathbb{R}^{d_{\zeta_s}} \rightarrow \mathbb{R}$  is the density function of the Gaussian distribution  $N(0, I_{d_{\zeta_s}} / d_{\zeta_s})$  and  $\theta(0)$  is the initial parameter. For any stationary point  $\tilde{\theta}$ , define the function

$$u_{\tilde{\theta}}(\mu, h) := \frac{d\sigma_{\theta^*}}{d\bar{\sigma}_{\tilde{\theta}}}(\zeta) - \frac{d\bar{\sigma}_{\theta^*}}{d\bar{\sigma}_{\tilde{\theta}}}(\mu) + \sum_{s \in \mathcal{S}} \phi_{\tilde{\theta}_s}(\zeta_s)^\top \tilde{\theta}_s,$$

with  $\bar{\sigma}_\theta$  the state visitation measure under policy  $\Pi^\theta$ , and  $\frac{d\sigma_{\theta^*}}{d\bar{\sigma}_{\tilde{\theta}}}, \frac{d\bar{\sigma}_{\theta^*}}{d\bar{\sigma}_{\tilde{\theta}}}$  the Radon-Nikodym derivatives between corresponding measures. We assume that  $u_{\tilde{\theta}} \in \mathcal{F}_{R,\infty}$  for any stationary point  $\tilde{\theta}$ .

A few remarks are in place for these Assumption 5.5 - Assumption 5.9.

**Remark.** All these assumptions are counterparts of standard assumption in the analysis of single-agent policy gradient method (Pirodda et al. [43], Xu et al. [54], Xu et al. [55], Zhang et al. [59], Wang et al. [53]).

In particular, Assumption 5.5 and Assumption 5.6 hold if the Markov chain (3.5) mixes sufficiently fast, and the critic  $Q_s(\cdot; \omega_s)$  has an upper-bounded second moment under  $\sigma_{\theta(t)}$  (Wang et al. [53]). Note that different from Assumption 5.2, where regularity conditions are imposed separately on  $\nu_\theta$  and  $\sigma_\theta$ , Assumption 5.6 imposes the regularity condition directly on the Radon-Nikodym derivative of  $\sigma_\theta$  with respect to  $\nu_\theta$ . This allows the change of measures in the analysis of Theorem 5.10. In general, Assumption 5.2 does not necessarily imply Assumption 5.6.

Assumption 5.7 holds when the transition probability and the reward function are both Lipschitz continuous with respect to their inputs (Pirodda et al. [43]), or when the reward is

uniformly bounded and the score function  $\nabla_{\theta}\Pi^{\theta}$  is uniformly bounded and Lipschitz continuous with respect to  $\theta$  (Zhang et al. [59]).

As for Assumption 5.9, we first emphasize that  $u_{\tilde{\theta}}(\mu, h)$  is a key element in the proof of Theorem 5.10. More specifically, this assumption is motivated by the well-known Performance Difference Lemma (Kakade and Langford [32]) in order to characterize the optimality gap of a stationary point  $\theta$ . In particular, it guarantees that  $u_{\tilde{\theta}}$  can be decomposed into a sum of local functions depending on  $\zeta_s$ , and that each local function lies in a rich RKHS (see the discussion after Assumption 5.1).

With all these assumptions, we now establish the rate of convergence for Algorithm 2.

**Theorem 5.10** *Assume Assumptions 5.1 - 5.9 . Set  $T_{\text{critic}} = \Omega(M)$ ,  $\eta_{\text{critic}} = \min\{(1 - \gamma)/8, (T_{\text{critic}})^{-1/2}\}$ ,  $\eta_{\text{actor}} = (T_{\text{actor}})^{-1/2}$ ,  $R = \tau = 1$ ,  $M = \Omega((f(k)|\mathcal{A}|)^5(T_{\text{actor}})^8)$ ,  $\gamma \leq (T_{\text{actor}})^{-2/k}$ , with  $f(k) := \max_{s \in \mathcal{S}} |\mathcal{N}_s^k|$  the size of the largest  $k$ -neighborhood in the graph  $(\mathcal{S}, \mathcal{E})$ . Then, the output  $\{\theta(t)\}_{t \in [T_{\text{actor}}]}$  of Algorithm 2 satisfies*

$$\min_{t \in [T_{\text{actor}}]} \mathbb{E}[J(\theta^*) - J(\theta(t))] \leq \mathcal{O}\left(|\mathcal{S}|^{1/2}B^{-1/2} + |\mathcal{S}||\mathcal{A}|^{1/4}\left(\gamma^{k/8} + (T_{\text{actor}})^{-1/4}\right)\right). \quad (5.3)$$

The error  $\mathcal{O}(\gamma^{k/8}|\mathcal{S}||\mathcal{A}|^{1/4})$  in Theorem 5.10, coming from the localized training, decays exponentially fast as  $k$  increases and is negligible with a careful choice of  $k$ . According to Theorem 5.10, Algorithm 2 converges at rate  $T_{\text{actor}}^{-1/4}$  with sufficiently large width  $M$  and batch size  $B$ . Technically,  $\{\theta_s(t)\}_{s \in \mathcal{S}}$  in Algorithm 2 are updated in parallel and our analysis extends the single agent actor-critic in Cai et al. [6] to the multi-agent decentralized case.

Detailed proof is provided in Section D.2.

## References

- [1] Agarwal A, Kakade SM, Lee JD, Mahajan G (2021) On the theory of policy gradient methods: Optimality, approximation, and distribution shift. Journal of Machine Learning Research 22(98):1–76.
- [2] Aïd R, Dumitrescu R, Tankov P (2021) The entry and exit game in the electricity markets: a mean-field game approach. Journal of Dynamics & Games 8(4):331.
- [3] Allen-Zhu Z, Li Y, Song Z (2019) A convergence theory for deep learning via over-parameterization. International Conference on Machine Learning, 242–252 (PMLR).
- [4] Bhandari J, Russo D, Singal R (2018) A finite time analysis of temporal difference learning with linear function approximation. Conference on Learning Theory, 1691–1692 (PMLR).
- [5] Cabannes T, Lauriere M, Perolat J, Marinier R, Girgin S, Perrin S, Pietquin O, Bayen AM, Goubault E, Elie R (2021) Solving n-player dynamic routing games with congestion: a mean-field approach. arXiv preprint arXiv:2110.11943 .
- [6] Cai Q, Yang Z, Lee JD, Wang Z (2019) Neural temporal-difference learning converges to global optima. Advances in Neural Information Processing Systems, volume 32, 11315–11326.
- [7] Calderone D, Sastry SS (2017) Markov decision process routing games. International Conference on Cyber-Physical Systems, 273–280 (IEEE).
- [8] Cao Y, Yu W, Ren W, Chen G (2012) An overview of recent progress in the study of distributed multi-agent coordination. IEEE Transactions on Industrial informatics 9(1):427–438.
- [9] Carmona R, Fouque JP, Sun LH (2015) Mean-field games and systemic risk. Communications in Mathematical Sciences 13(4):911–933.
- [10] Carmona R, Laurière M, Tan Z (2019) Linear-quadratic mean-field reinforcement learning: convergence of policy gradient methods. arXiv preprint arXiv:1910.04295 .
- [11] Carmona R, Laurière M, Tan Z (2019) Model-free mean-field reinforcement learning: mean-field MDP and mean-field Q-learning. arXiv preprint arXiv:1910.12802 .
- [12] Casgrain P, Jaimungal S (2020) Mean-field games with differing beliefs for algorithmic trading. Mathematical Finance 30(3):995–1034.

- [13] Cayci S, Satpathi S, He N, Srikant R (2021) Sample complexity and overparameterization bounds for projection-free neural TD learning. [arXiv preprint arXiv:2103.01391](#) .
- [14] Chen T, Zhang K, Giannakis GB, Basar T (2021) Communication-efficient policy gradient methods for distributed reinforcement learning. [IEEE Transactions on Control of Network Systems](#) .
- [15] Dawson D (1993) Measure-valued Markov processes. [École d'été de probabilités de Saint-Flour XXI-1991](#), 1–260 (Springer).
- [16] El-Tantawy S, Abdulhai B, Abdelgawad H (2013) Multi-agent reinforcement learning for integrated network of adaptive traffic signal controllers (MARLIN-ATSC): Methodology and large-scale application on downtown Toronto. [IEEE Transactions on Intelligent Transportation Systems](#) 14(3):1140–1150.
- [17] Foerster J, Farquhar G, Afouras T, Nardelli N, Whiteson S (2018) Counterfactual multi-agent policy gradients. [AAAI Conference on Artificial Intelligence](#), volume 32.
- [18] Fu Z, Yang Z, Wang Z (2020) Single-timescale actor-critic provably finds globally optimal policy. [International Conference on Learning Representations](#).
- [19] Gamarnik D (2013) Correlation decay method for decision, optimization, and inference in large-scale networks. [Theory Driven by Influential Applications](#), 108–121 (INFORMS).
- [20] Gamarnik D, Goldberg DA, Weber T (2014) Correlation decay in random decision networks. [Mathematics of Operations Research](#) 39(2):229–261.
- [21] Germain M, Pham H, Warin X (2021) A level-set approach to the control of state-constrained mckean-vlasov equations: application to renewable energy storage and portfolio selection. [arXiv preprint arXiv:2112.11059](#) .
- [22] Glorot X, Bengio Y (2010) Understanding the difficulty of training deep feedforward neural networks. [International Conference on Artificial Intelligence and Statistics](#), 249–256.
- [23] Gu H, Guo X, Wei X, Xu R (2019) Dynamic programming principles for learning MFCs. [arXiv preprint arXiv:1911.07314](#) .
- [24] Gu H, Guo X, Wei X, Xu R (2021) Mean-field controls with Q-learning for cooperative MARL: convergence and complexity analysis. [SIAM Journal on Mathematics of Data Science](#) 3(4):1168–1196.
- [25] Guériau M, Dusparic I (2018) Samod: Shared autonomous mobility-on-demand using decentralized reinforcement learning. [International Conference on Intelligent Transportation Systems](#), 1558–1563 (IEEE).
- [26] Guo X, Hu A, Xu R, Zhang J (2019) Learning mean-field games. [Advances in Neural Information Processing Systems](#), volume 32, 4966–4976.
- [27] Hu R, Zariphopoulou T (2021) N-player and mean-field games in Itô-diffusion markets with competitive or homophilous interaction. [arXiv preprint arXiv:2106.00581](#) .
- [28] Hüttenrauch M, Šošić A, Neumann G (2017) Guided deep reinforcement learning for swarm systems. [arXiv preprint arXiv:1709.06011](#) .
- [29] Iyer K, Johari R, Sundararajan M (2014) Mean-field equilibria of dynamic auctions with learning. [Management Science](#) 60(12):2949–2970.
- [30] Ji Z, Telgarsky M, Xian R (2020) Neural tangent kernels, transportation mappings, and universal approximation. [International Conference on Learning Representations](#).
- [31] Jin J, Song C, Li H, Gai K, Wang J, Zhang W (2018) Real-time bidding with multi-agent reinforcement learning in display advertising. [ACM International Conference on Information and Knowledge Management](#), 2193–2201.
- [32] Kakade S, Langford J (2002) Approximately optimal approximate reinforcement learning. [International Conference on Machine Learning](#), 267–274 (PMLR).
- [33] Konda VR, Tsitsiklis JN (2000) Actor-critic algorithms. [Advances in Neural Information Processing Systems](#), volume 12, 1008–1014.
- [34] Lacker D, Zariphopoulou T (2019) Mean-field and n-agent games for optimal investment under relative performance criteria. [Mathematical Finance](#) 29(4):1003–1038.
- [35] Li M, Qin Z, Jiao Y, Yang Y, Wang J, Wang C, Wu G, Ye J (2019) Efficient ridesharing order dispatching with mean-field multi-agent reinforcement learning. [The World Wide Web Conference](#), 983–994.

- [36] Li Y, Tang Y, Zhang R, Li N (2021) Distributed reinforcement learning for decentralized linear quadratic control: A derivative-free policy optimization approach. IEEE Transactions on Automatic Control .
- [37] Lin Y, Qu G, Huang L, Wierman A (2021) Multi-agent reinforcement learning in stochastic networked systems. Advances in Neural Information Processing Systems, volume 34.
- [38] Liu B, Cai Q, Yang Z, Wang Z (2019) Neural trust region/proximal policy optimization attains globally optimal policy. Advances in Neural Information Processing Systems, volume 32, 10565–10576.
- [39] Liu Y, Swaminathan A, Agarwal A, Brunskill E (2019) Off-policy policy gradient with stationary distribution correction. Conference on Uncertainty in Artificial Intelligence, volume 115, 1180–1190 (PMLR).
- [40] Lowe R, Wu YI, Tamar A, Harb J, Pieter Abbeel O, Mordatch I (2017) Multi-agent actor-critic for mixed cooperative-competitive environments. Advances in Neural Information Processing Systems, volume 30, 6382–6393.
- [41] Micchelli CA, Xu Y, Zhang H (2006) Universal kernels. Journal of Machine Learning Research 7(12):2651–2667.
- [42] Motte M, Pham H (2019) Mean-field markov decision processes with common noise and open-loop controls. arXiv preprint arXiv:1912.07883 .
- [43] Pirodda M, Restelli M, Bascetta L (2015) Policy gradient in lipschitz Markov decision processes. Machine Learning 100(2):255–283.
- [44] Qu G, Wierman A, Li N (2020) Scalable reinforcement learning of localized policies for multi-agent networked systems. Learning for Dynamics and Control, 256–266 (PMLR).
- [45] Rabbat M, Nowak R (2004) Distributed optimization in sensor networks. International Symposium on Information Processing in Sensor Networks, 20–27.
- [46] Rahimi A, Recht B (2008) Uniform approximation of functions with random bases. Annual Allerton Conference on Communication, Control, and Computing, 555–561 (IEEE).
- [47] Rashid T, Samvelyan M, Schroeder C, Farquhar G, Foerster J, Whiteson S (2018) QMIX: Monotonic value function factorisation for deep multi-agent reinforcement learning. International Conference on Machine Learning, 4295–4304 (PMLR).
- [48] Shalev-Shwartz S, Shammah S, Shashua A (2016) Safe, multi-agent, reinforcement learning for autonomous driving. arXiv preprint arXiv:1610.03295 .
- [49] Sra S, Nowozin S, Wright SJ (2012) Optimization for Machine Learning (MIT Press).
- [50] Sunehag P, Lever G, Gruslly A, Czarnecki WM, Zambaldi V, Jaderberg M, Lanctot M, Sonnerat N, Leibo JZ, Tuyls K, Thore G (2018) Value-decomposition networks for cooperative multi-agent learning based on team reward. International Conference on Autonomous Agents and Multi-agent Systems, volume 3, 2085–2087.
- [51] Sutton RS, McAllester DA, Singh SP, Mansour Y (2000) Policy gradient methods for reinforcement learning with function approximation. Advances in Neural Information Processing Systems, volume 99, 1057–1063.
- [52] Vadori N, Ganesh S, Reddy P, Veloso M (2020) Calibration of shared equilibria in general sum partially observable markov games. Advances in Neural Information Processing Systems, volume 33, 14118–14128.
- [53] Wang L, Cai Q, Yang Z, Wang Z (2020) Neural policy gradient methods: Global optimality and rates of convergence. International Conference on Learning Representations.
- [54] Xu P, Gao F, Gu Q (2019) Sample efficient policy gradient methods with recursive variance reduction. International Conference on Learning Representations.
- [55] Xu P, Gao F, Gu Q (2020) An improved convergence analysis of stochastic variance-reduced policy gradient. Uncertainty in Artificial Intelligence, 541–551 (PMLR).
- [56] Yang Y, Hao J, Chen G, Tang H, Chen Y, Hu Y, Fan C, Wei Z (2020) Q-value path decomposition for deep multiagent reinforcement learning. International Conference on Machine Learning, 10706–10715 (PMLR).
- [57] Yang Y, Wen Y, Wang J, Chen L, Shao K, Mguni D, Zhang W (2020) Multi-agent determinantal Q-learning. International Conference on Machine Learning, 10757–10766 (PMLR).

- [58] You X, Li X, Xu Y, Feng H, Zhao J, Yan H (2020) Toward packet routing with fully distributed multiagent deep reinforcement learning. IEEE Transactions on Systems, Man, and Cybernetics: Systems .
- [59] Zhang K, Koppel A, Zhu H, Basar T (2020) Global convergence of policy gradient methods to (almost) locally optimal policies. SIAM Journal on Control and Optimization 58(6):3586–3612.
- [60] Zhang K, Liu Y, Liu J, Liu M, Başar T (2020) Distributed learning of average belief over networks using sequential observations. Automatica 115:108857.
- [61] Zhang K, Yang Z, Basar T (2018) Networked multi-agent reinforcement learning in continuous spaces. Conference on Decision and Control, 2771–2776 (IEEE).
- [62] Zhang K, Yang Z, Başar T (2021) Multi-agent reinforcement learning: A selective overview of theories and algorithms. Handbook of Reinforcement Learning and Control, 321–384 (Springer).
- [63] Zhang K, Yang Z, Liu H, Zhang T, Basar T (2018) Fully decentralized multi-agent reinforcement learning with networked agents. International Conference on Machine Learning, 5872–5881 (PMLR).
- [64] Zhang K, Yang Z, Liu H, Zhang T, Basar T (2021) Finite-sample analysis for decentralized batch multi-agent reinforcement learning with networked agents. IEEE Transactions on Automatic Control .
- [65] Zheng S, Trott A, Srinivasa S, Naik N, Gruesbeck M, Parkes DC, Socher R (2020) The AI economist: Improving equality and productivity with AI-driven tax policies. arXiv preprint arXiv:2004.13332 .
- [66] Zhou Z, Mertikopoulos P, Moustakas AL, Bambos N, Glynn P (2021) Robust power management via learning and game design. Operations Research 69(1):331–345.

# Appendix

## A Proof of Lemma 3.1

The goal is to show that  $V(\mu) = \tilde{V}(\mu)$ , with the former the value function of (MF-MARL) subject to the transition probability  $P$  defined in (2.7) under a given individual policy  $\pi \in \mathfrak{u}$ , and the latter the value function of (3.7) subject to the joint transition probability  $\mathbf{P}^N$  defined in (3.5) under the policy  $\Pi \in \mathfrak{U}$ . The proof consists of two steps. Step 1 shows that  $V(\mu)$  can be reformulated as a *measured-valued* Markov decision problem. Step 2 shows that the measured-valued Markov decision problem from Step 1 is equivalent to  $\tilde{V}(\mu)$  in (3.7).

*Step 1:* Recall that  $\mu_{t+1} := \frac{1}{N} \sum_{i=1}^N \delta_{s_{t+1}^i}$  with  $s_{t+1}^i$  subject to (2.7). First, one can show that  $\mu_t$  is a measure-valued Markov decision process under  $\pi$ . To see this, denote  $\mathcal{F}_t^s = \sigma(s_t^1, \dots, s_t^N)$  as the  $\sigma$ -algebra generated by  $s_t^1, \dots, s_t^N$ . Then it suffices to show

$$\mathbb{P}(\mu_{t+1} \mid \sigma(\mu_t) \vee \mathcal{F}_t^s) = \mathbb{P}(\mu_{t+1} \mid \sigma(\mu_t)), \quad \mathbb{P} - a.s.. \quad (\text{A.1})$$

Following similar arguments for Lemma 2.3.1 and Proposition 2.3.3 in Dawson [15], (A.1) holds due to the exchangeability of the individual transition dynamics (2.7) under  $\pi$ . (A.1) implies that there exists a joint transition probability induced from (2.7) under  $\pi$ , denoted as  $\tilde{\mathbf{P}}^N$  such that

$$\mu_{t+1} \sim \tilde{\mathbf{P}}^N(\cdot \mid \mu_t, \pi). \quad (\text{A.2})$$

Meanwhile, rewrite  $V^\pi(\mu)$  in (MF-MARL) by regrouping the agents according to their states

$$\begin{aligned} V^\pi(\mu) &:= \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t \sum_{i=1}^N \frac{1}{N} r(s_t^i, \mu_t(\mathcal{N}_{s_t^i}), a_t^i) \mid \mu_0 = \mu \right], \\ &= \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t \sum_{s \in \mathcal{S}} \mu_t(s) \sum_{a \in \mathcal{A}} r(s, \mu_t(\mathcal{N}_s), a) \pi(s, \mu_t(s))(a) \mid \mu_0 = \mu \right]. \end{aligned} \quad (\text{A.3})$$

We see (2.7)-(MF-MARL) is reformulated in an equivalent form of (A.2)-(A.3).

*Step 2:* It suffices to show that (A.2) under  $\pi$  is the same as (3.5) under  $\Pi$  and that  $V^\pi$  in (A.3) equals to  $\tilde{V}^\Pi$  in (3.7). To see this, denote  $\langle g, \mu \rangle = \sum_{s \in \mathcal{S}} g(s) \mu(s)$  for any measurable bounded function  $g : \mathcal{S} \rightarrow \mathbb{R}$ , then

$$\begin{aligned} &\mathbb{E}[\langle g, \mu_{t+1} \rangle \mid \sigma(\mu_t)] \\ &= \frac{1}{N} \mathbb{E} \left[ \sum_{i=1}^N \mathbb{E}[g(s_{t+1}^i) \mid \sigma(\mu_t) \vee \mathcal{F}_t^s] \right] \\ &= \frac{1}{N} \sum_{s' \in \mathcal{S}} \sum_{i=1}^N \sum_{a \in \mathcal{A}} g(s') P(s' \mid s_t^i, \mu_t(\mathcal{N}(s_t^i)), a) \pi(s_t^i, \mu_t(s_t^i))(a) \\ &= \frac{1}{N} \sum_{s' \in \mathcal{S}} g(s') \sum_{s \in \mathcal{S}} \sum_{i=1}^N \mathbb{1}(s_t^i = s) \sum_{a \in \mathcal{A}} P(s' \mid s_t^i, \mu_t(\mathcal{N}(s_t^i)), a) \pi(s_t^i, \mu_t(s_t^i))(a) \\ &= \sum_{s' \in \mathcal{S}} g(s') \sum_{s \in \mathcal{S}} \mu_t(s) \sum_{a \in \mathcal{A}} P(s' \mid s, \mu_t(\mathcal{N}(s)), a) \pi(s, \mu_t(s))(a) \\ &= \sum_{s' \in \mathcal{S}} g(s') \sum_{s \in \mathcal{S}} \mu_t(s) \sum_{h \in \mathcal{P}^{N \cdot \mu_t(s)}(\mathcal{A})} \Pi(h \mid \mu_t(s)) \sum_{a \in \mathcal{A}} P(s' \mid s, \mu_t(\mathcal{N}(s)), a) h(s)(a), \end{aligned} \quad (\text{A.4})$$

where in the last step, the expectation of random variable  $h(s)(a)$  with respect to distribution  $\Pi(h \mid \mu)$  is  $\pi(s, \mu_t(s))$ . And from the last equality, clearly  $\mu_{t+1}$  evolves according to transition dynamics  $\mathbf{P}^N(\cdot \mid \mu_t, h_t)$  under  $\Pi(h_t \mid \mu_t)$ . This implies the equivalence of (A.2) and (3.5). As a byproduct, when taking  $g(s') = \mathbb{1}(s' = s^o)$  for any fixed  $s^o \in \mathcal{S}$ , (A.4) becomes

$$\mathbb{E}[\mu_{t+1}(s^o)|\sigma(\mu_t)] = \sum_{s \in \mathcal{N}(s^o)} \mu_t(s) \sum_{h \in \mathcal{P}^{N \cdot \mu_t(s)}(\mathcal{A})} \Pi(h | \mu_t(s)) \sum_{a \in \mathcal{A}} P(s^o | s, \mu_t(\mathcal{N}(s)), a) h(s)(a),$$

where the local structure (2.7) is used. This suggests that  $\mu_{t+1}(s^o)$  only depends on  $\mu_t(\mathcal{N}_{s^o}^2)$  and  $h_t(\mathcal{N}_{s^o})$  since  $\mathcal{N}(s) = \mathcal{N}^2(s^o)$  for  $s \in \mathcal{N}(s^o)$ .

Now we show that  $V^\pi(\mu)$  in (A.3) and  $\tilde{V}^\Pi(\mu)$  in (3.7) are equal. Take  $\tilde{V}^\Pi$  defined in (3.7),

$$\begin{aligned} \tilde{V}^\Pi(\mu) &:= \mathbb{E}_{h_t \sim \Pi(\cdot | \mu_t), \mu_{t+1} \sim \mathbf{P}^N(\cdot | \mu_t, h_t)} \left[ \sum_{t=0}^{\infty} \sum_{s \in \mathcal{S}} \gamma^t r_s(\mu_t(\mathcal{N}_s), h_t) \middle| \mu_0 = \mu \right] \\ &= \mathbb{E}_{\mu_{t+1} \sim \mathbf{P}^N(\cdot | \mu_t, h_t)} \left[ \sum_{t=0}^{\infty} \gamma^t \sum_{s \in \mathcal{S}} \mathbb{E}_{h_t \sim \Pi(\cdot | \mu_t)} [r_s(\mu_t(\mathcal{N}_s), h_t) | \mu_t] \middle| \mu_0 = \mu \right] \\ &= \mathbb{E}_{\mu_{t+1} \sim \mathbf{P}^N(\cdot | \mu_t, h_t)} \left[ \sum_{t=0}^{\infty} \gamma^t \sum_{s \in \mathcal{S}} \sum_{h_t \in \mathcal{P}^{N \cdot \mu_t(s)}(\mathcal{A})} r_s(\mu_t(\mathcal{N}_s), h_t(s)) \Pi(h; \pi) \middle| \mu_0 = \mu \right] \\ &= \mathbb{E}_{\mu_{t+1} \sim \mathbf{P}^N(\cdot | \mu_t, h_t)} \left[ \sum_{t=0}^{\infty} \gamma^t \sum_{s \in \mathcal{S}} \mu_t(s) \sum_{h_t \in \mathcal{P}^{N \cdot \mu_t(s)}(\mathcal{A})} \Pi(h_t | \mu_t) \sum_{a \in \mathcal{A}} r(s, \mu_t(\mathcal{N}_s), a) h(a) \middle| \mu_0 = \mu \right] \\ &= \mathbb{E}_{\mu_{t+1} \sim \tilde{\mathbf{P}}^N(\cdot | \mu_t, \pi)} \left[ \sum_{t=0}^{\infty} \gamma^t \sum_{s \in \mathcal{S}} \mu_t(s) \sum_{a \in \mathcal{A}} r(s, \mu_t(\mathcal{N}_s), a) \pi_t(s, \mu_t(s))(a) \middle| \mu_0 = \mu \right] \\ &= V^\pi(\mu), \end{aligned}$$

where in the last second step,  $\mathbf{P}^N$  under  $\pi$  is equivalent to  $\tilde{\mathbf{P}}^N$  under  $\Pi$ , and the expectation of  $h_t(s)(a)$  with distribution  $\Pi(h_t | \mu_t)$  is  $\pi(s, \mu_t(s))(a)$  such that

$$\begin{aligned} \sum_{h_t \in \mathcal{P}^{N \cdot \mu_t(s)}(\mathcal{A})} \Pi(h_t | \mu_t) \sum_{a \in \mathcal{A}} r(s, \mu_t(\mathcal{N}_s), a) h(a) &= \mathbb{E}_{h \sim \Pi(\cdot | \mu_t)} \left[ \sum_{a \in \mathcal{A}} r(s, \mu_t(\mathcal{N}_s), a) h(a) \right] \\ &= \sum_{a \in \mathcal{A}} r(s, \mu_t(\mathcal{N}_s), a) \pi_t(s, \mu_t(s))(a). \end{aligned}$$

Finally, the decomposition of  $\tilde{V}(\mu)$  and  $Q^{\Pi^\theta}(\mu, h)$  according to the states is straightforward. **Q.E.D.**

## B Proof of Lemma 3.5

Let  $\mathfrak{P}_{t,s}$  and  $\mathfrak{P}'_{t,s}$  be, respectively, distribution of  $(\mu_t(\mathcal{N}_s), h_t(s))$  and  $(\mu'_t(\mathcal{N}_s), h'_t(s))$  under policy  $\Pi^\theta$ . By localized transition kernel (2.7), it is easy to see that for any given  $s \in \mathcal{S}$ ,  $\mu_{t+1}(s)$  only depends on  $\mu_t(\mathcal{N}_s^2)$  and  $h_t(\mathcal{N}_s)$ . Then by the local dependency, (3.5) can be rewritten as

$$\mu_{t+1}(s) \sim \mathbf{P}_s^N(\cdot | \mu_t(\mathcal{N}_s^2), h_t(\mathcal{N}_s)). \quad (\text{B.1})$$

Due to the local structure of dynamics (B.1) and local dependence of  $\Pi^\theta$ , the distribution  $\mathfrak{P}_{t,s}$ ,  $t \leq \lfloor \frac{k}{2} \rfloor$  only depends on the initial value  $(\mu(\mathcal{N}_s^k), h(\mathcal{N}_s^k))$ . Therefore,  $\mathfrak{P}_{t,s} = \mathfrak{P}'_{t,s}$ ,  $t \leq \lfloor \frac{k}{2} \rfloor$ ,

$$\begin{aligned} &\left| Q_s^{\Pi^\theta}(\mu(\mathcal{N}_s^k), \mu(\mathcal{N}_s^{-k}), h(\mathcal{N}_s^k), h(\mathcal{N}_s^{-k})) - Q_s^{\Pi^\theta}(\mu(\mathcal{N}_s^k), \mu'(\mathcal{N}_s^{-k}), h(\mathcal{N}_s^k), h'(\mathcal{N}_s^{-k})) \right| \\ &= \sum_{t=\lfloor \frac{k}{2} \rfloor + 1}^{\infty} \mathbb{E}_{(\mu_t(\mathcal{N}_s), h_t(s)) \sim \mathfrak{P}_{t,s}} [r_s(\mu_t(\mathcal{N}_s), h_t(s))] - \mathbb{E}_{(\mu'_t(\mathcal{N}_s), h'_t(s)) \sim \mathfrak{P}'_{t,s}} [r_s(\mu'_t(\mathcal{N}_s), h'_t(s))] \\ &\leq \sum_{t=\lfloor \frac{k}{2} \rfloor + 1}^{\infty} \gamma^t r_{\max} \text{TV}(\mathfrak{P}_{t,s}, \mathfrak{P}'_{t,s}) \leq \frac{r_{\max}}{1 - \gamma} \gamma^{\lfloor \frac{k}{2} \rfloor + 1}, \end{aligned}$$

where  $\text{TV}(\mathfrak{P}_{t,s}, \mathfrak{P}'_{t,s})$  is total variation between  $\mathfrak{P}_{t,s}$  and  $\mathfrak{P}'_{t,s}$  that is upper bounded by 1. **Q.E.D.**

## C Proof of Lemma 4.2

For any  $\theta \in \Theta$ ,  $s \in \mathcal{S}$ ,  $\mu \in \mathcal{P}^N(\mathcal{S})$  and  $h \in \mathcal{H}^N(\mu)$ , it is easy to verify that  $\|\Phi(\theta, s, \mu, h)\|_2 \leq \|\zeta_s\|_2 \leq 2$ , by the definitions of the feature mapping  $\phi$  in (4.6) and the center feature mapping  $\Phi$  in (4.7).

$$\begin{aligned} \text{To prove (4.8), note that by Lemma 4.1 \& the definition of energy-based policy } \Pi_s^{\theta_s} \text{ (4.4),} \\ \nabla_{\theta_s} \log \Pi_s^{\theta_s}(h(s) \mid \mu(s)) &= \tau \cdot \nabla_{\theta_s} f((\mu(s), h(s)); \theta_s) - \tau \cdot \mathbb{E}_{h(s)' \sim \Pi_s^{\theta_s}(\cdot \mid \mu(s))} [\nabla_{\theta_s} f(\mu(s), h'(s))] \\ &= \tau \cdot \phi_{\theta_s}(\mu(s), h(s)) - \tau \cdot \mathbb{E}_{h(s)' \sim \Pi_s^{\theta_s}(\cdot \mid \mu(s))} [\phi_{\theta_s}(\mu(s), h(s))] \\ &= \tau \cdot \Phi(\theta, s, \mu, h). \end{aligned}$$

The second equality follows from the fact that  $\nabla_{\theta_s} f((\mu(s), h(s)); \theta_s) = \phi_{\theta_s}(\mu(s), h(s))$ . Therefore,

$$\nabla_{\theta_s} J(\theta) = \frac{\tau}{1-\gamma} \mathbb{E}_{\sigma_\theta} \left[ Q^{\Pi^\theta}(\mu, h) \cdot \Phi(\theta, s, \mu, h) \right] = \frac{\tau}{1-\gamma} \mathbb{E}_{\sigma_\theta} \left[ \sum_{y \in \mathcal{S}} Q_y^{\Pi^\theta}(\mu, h) \cdot \Phi(\theta, s, \mu, h) \right],$$

where the second equality is by the decomposition of Q-function in Lemma 3.1.

The proof of (4.9) is based on the exponential decay property in Definition 3.4. Notice that

$$\begin{aligned} g_s(\theta) &= \frac{1}{1-\gamma} \mathbb{E}_{\sigma_\theta} \left[ \left[ \sum_{y \in \mathcal{N}_s^k} \widehat{Q}_y^{\Pi^\theta}(\mu(\mathcal{N}_y^k), h(\mathcal{N}_y^k)) \right] \nabla_{\theta_s} \log \Pi_s^{\theta_s}(h(s) \mid \mu(s)) \right] \\ &= \frac{1}{1-\gamma} \mathbb{E}_{\sigma_\theta} \left[ \left[ \sum_{y \in \mathcal{S}} \widehat{Q}_y^{\Pi^\theta}(\mu(\mathcal{N}_y^k), h(\mathcal{N}_y^k)) \right] \nabla_{\theta_s} \log \Pi_s^{\theta_s}(h(s) \mid \mu(s)) \right]. \end{aligned} \quad (\text{C.1})$$

This is because for all  $y \notin \mathcal{N}_s^k$ ,  $\widehat{Q}_y^{\Pi^\theta}(\mu(\mathcal{N}_y^k), h(\mathcal{N}_y^k))$  is independent of  $s$ . Consequently,

$$\mathbb{E}_{\sigma_\theta} \left[ \left[ \sum_{y \notin \mathcal{N}_s^k} \widehat{Q}_y^{\Pi^\theta}(\mu(\mathcal{N}_y^k), h(\mathcal{N}_y^k)) \right] \nabla_{\theta_s} \log \Pi_s^{\theta_s}(h(s) \mid \mu(s)) \right] = 0.$$

Given Lemma 4.1 and (C.1), we have the following bound:

$$\begin{aligned} &\|g_s(\theta) - \nabla_{\theta_s} J(\theta)\|_2 \\ &\leq \frac{1}{1-\gamma} \sum_{y \in \mathcal{S}} \sup_{\substack{\mu \in \mathcal{P}^N(\mathcal{S}), \\ h \in \mathcal{H}^N(\mu)}} \left[ \left| \widehat{Q}_y^{\Pi^\theta}(\mu(\mathcal{N}_y^k), h(\mathcal{N}_y^k)) - Q_y^{\Pi^\theta}(\mu, h) \right| \cdot \|\nabla_{\theta_s} \log \Pi_s^{\theta_s}(h(s) \mid \mu(s))\|_2 \right] \\ &\leq \frac{c_0 \tau |\mathcal{S}|}{1-\gamma} \rho^{k+1}. \end{aligned}$$

The last inequality follows from (3.14) and  $\|\log \Pi_s^{\theta_s}(h(s) \mid \mu(s))\|_2 = \|\Phi(\theta, s, \mu, h)\|_2 \leq 2$  for any  $\mu \in \mathcal{P}^N(\mathcal{S}), h \in \mathcal{H}^N(\mu)$ . **Q.E.D.**

## D Proof of Theorems 5.4 and 5.10

### D.1 Proof of Theorem 5.4: Convergence of Critic Update

This section presents the proof of convergence of the decentralized neural critic update. It consists of several steps. Section D.1.1 introduces necessary notations and definitions. Section D.1.2 proves that the critic update minimizes the projected mean-square Bellman error given a two-layer neural network. Section D.1.3 shows that the global minimizer of the projected mean-square Bellman error converges to the true team-decentralized Q-function as the width of hidden layer  $M \rightarrow \infty$ .

### D.1.1 Notations

Recall that the set of all state-action (distribution) pairs is denoted as  $\Xi := \cup_{\mu \in \mathcal{P}^N(\mathcal{S})} \{\zeta = (\mu, h) : h \in \mathcal{H}^N(\mu)\}$ . For any  $\zeta = (\mu, h) \in \Xi$ , denote the localized state-action (distribution) pair as  $\zeta_s^k = (\mu(\mathcal{N}_s^k), h(\mathcal{N}_s^k))$ . Meanwhile, denote  $\Xi_s^k = \{\zeta_s^k : \zeta \in \Xi\}$  as the set of all possible localized state-action (distribution) pairs. Without loss of generality, assume  $\|\zeta_s^k\|_2 \leq 1$  for any  $\zeta_s^k \in \Xi_s^k$ .

Let  $d_\zeta$  denote the dimension of the space  $\Xi$ . Since  $\mathcal{P}^N(\mathcal{S})$  has dimension  $(|\mathcal{S}| - 1)$  and  $\mathcal{H}^N(\mu)$  has dimension  $|\mathcal{S}|(|\mathcal{A}| - 1)$  for any  $\mu \in \mathcal{P}^N(\mathcal{S})$ , the product space  $\Xi$  has dimension  $d_\zeta = |\mathcal{S}||\mathcal{A}| - 1$ . Similarly, one can see that the dimension of the space  $\Xi_s^k$ , denoted by  $d_{\zeta_s^k}$ , is at most  $f(k)|\mathcal{A}|$ , where  $f(k) := \max_{s \in \mathcal{X}} |\mathcal{N}_s^k|$  is the size of the largest  $k$ -neighborhood in the graph  $(\mathcal{S}, \mathcal{E})$ .

Let  $\mathbb{R}^\Xi$  and  $\mathbb{R}^{\Xi_s^k}$  be the sets of real-valued square-integrable functions (with respect to  $\nu_\theta$ ) on  $\Xi$  and  $\Xi_s^k$ , respectively. Define the norm  $\|\cdot\|_{L^2(\nu_\theta)}$  on  $\mathbb{R}^\Xi$  by

$$\|f\|_{L^2(\nu_\theta)} := (\mathbb{E}_{\zeta \sim \nu_\theta} [f(\zeta)^2])^{1/2}, \quad \forall f \in \mathbb{R}^\Xi. \quad (\text{D.1})$$

Note that for any function  $f \in \mathbb{R}^{\Xi_s^k}$ , a function  $\tilde{f} \in \mathbb{R}^\Xi$  is called a *natural extension* of  $f$  if  $\tilde{f}(\zeta) = f(\zeta_s^k)$  for all  $\zeta \in \Xi$ . Since the natural extension is an injective mapping from  $\mathbb{R}^{\Xi_s^k}$  to  $\mathbb{R}^\Xi$ , one can view  $\mathbb{R}^{\Xi_s^k}$  as a subset of  $\mathbb{R}^\Xi$ . In addition for a function  $f \in \mathbb{R}^{\Xi_s^k}$ , we use the same notation  $f \in \mathbb{R}^\Xi$  to denote the natural extension of  $f$ .

For any closed and convex function class  $\mathcal{F} \subset \mathbb{R}^\Xi$ , define the project operator  $\text{Proj}_{\mathcal{F}}$  from  $\mathbb{R}^\Xi$  onto  $\mathcal{F}$  by

$$\text{Proj}_{\mathcal{F}}(g) := \arg \min_{f \in \mathcal{F}} \|f - g\|_{L^2(\nu_\theta)}. \quad (\text{D.2})$$

This projection operator  $\text{Proj}_{\mathcal{F}}$  is non-expansive in the sense that

$$\|\text{Proj}_{\mathcal{F}}(f) - \text{Proj}_{\mathcal{F}}(g)\|_{L^2(\nu_\theta)} \leq \|f - g\|_{L^2(\nu_\theta)}. \quad (\text{D.3})$$

Recall that for each state  $s \in \mathcal{S}$ , the critic parameter  $\omega_s$  is updated in a localized fashion using information from the  $k$ -hop neighborhood of  $s$ . Without loss of generality, let us omit the subscript  $s$  of  $\omega_s$  in the following presentation, and the result holds for all  $s \in \mathcal{S}$  simultaneously.

Given an initialization  $\omega(0) \in \mathbb{R}^{M \times d_{\zeta_s^k}}$ , define the following function class

$$\mathcal{F}_{R,M} = \left\{ Q_0(\zeta_s^k; \omega) := \frac{1}{\sqrt{M}} \sum_{m=1}^M \mathbb{1} \{ [\omega(0)]_m^\top \zeta_s^k > 0 \} \omega_m^\top \zeta_s^k : \omega \in \mathbb{R}^{M \times d_{\zeta_s^k}}, \|\omega - \omega(0)\|_\infty \leq R/\sqrt{M} \right\}. \quad (\text{D.4})$$

$Q_0(\cdot; \omega)$  locally linearizes the neural network  $Q(\cdot; \omega)$  (with respect to  $\omega$ ) at  $\omega(0)$ . Any function  $Q_0(\cdot; \omega) \in \mathcal{F}_{R,M}$  can be viewed as an inner product between the feature mapping  $\phi_{\omega(0)}(\cdot)$  defined in (4.6) and the parameter  $\omega$ , i.e.  $Q_0(\cdot; \omega) = \phi_{\omega(0)}(\cdot)^\top \omega$ . In addition it holds that  $\nabla_\omega Q_0(\cdot; \omega) = \phi_{\omega(0)}(\cdot)$ . All functions in  $\mathcal{F}_{R,M}$  share the same feature mapping  $\phi_{\omega(0)}(\cdot)$  which only depends on the initialization  $\omega(0)$ .

Recall the Bellman operator  $\mathcal{T}_s^\theta : \mathbb{R}^\Xi \rightarrow \mathbb{R}^\Xi$  defined in (3.13),

$$\mathcal{T}_s^\theta Q_s^{\Pi^\theta}(\mu, h) = \mathbb{E}_{\mu' \sim \mathcal{P}^N(\cdot | \mu, h), h' \sim \Pi^\theta(\cdot | \mu)} \left[ r_s(\mu, h) + \gamma \cdot Q_s^{\Pi^\theta}(\mu', h') \right], \forall (\mu, h) \in \Xi.$$

The team-decentralized Q-function  $Q_s^{\Pi^\theta}$  in (3.10) is the unique fixed point of  $\mathcal{T}_s^\theta$ :  $Q_s^{\Pi^\theta} = \mathcal{T}_s^\theta Q_s^{\Pi^\theta}$ . Now given a general parameterized function class  $\mathcal{F}$ , we aim to learn a  $Q_s(\cdot; \omega) \in \mathcal{F}$  to approximate  $Q_s^{\Pi^\theta}$  by minimizing the following projected mean-squared Bellman error (PMSBE):

$$\min_{\omega} \text{PMSBE}(\omega) = \mathbb{E}_{\zeta \sim \nu_\theta} \left[ \left( Q_s(\zeta_s^k; \omega) - \text{Proj}_{\mathcal{F}} \mathcal{T}_s^\theta Q_s(\zeta_s^k; \omega) \right)^2 \right]. \quad (\text{D.5})$$

In the first step of the convergence analysis, we take  $\mathcal{F} = \mathcal{F}_{R,M}$  (the locally linearized two-layer neural network defined in (D.4)) and consider the following PMSBE:

$$\min_{\omega} \mathbb{E}_{\zeta \sim \nu_{\theta}} \left[ \left( Q_0(\zeta^k; \omega) - \text{Proj}_{\mathcal{F}_{R,M}} \mathcal{T}_s^{\theta} Q_0(\zeta^k; \omega) \right)^2 \right]. \quad (\text{D.6})$$

We will show in Section D.1.2 that the output of Algorithm 1 converges to the global minimizer of (D.6).

### D.1.2 Convergence to the Global Minimizer in $\mathcal{F}_{R,M}$

The following lemma guarantees the existence and the uniqueness of the global minimizer of MSPBE that corresponds to the projection onto  $\mathcal{F}_{R,M}$  in (D.6).

**Lemma D.1** (*Existence and Uniqueness of the Global Minimizer in  $\mathcal{F}_{R,M}$* ) *For any  $b \in \mathbb{R}^M$  and  $\omega(0) \in \mathbb{R}^{M \times d_{\zeta^k}}$ , there exists an  $\omega^*$  such that  $Q_0(\cdot; \omega^*) \in \mathcal{F}_{R,M}$  is unique almost everywhere in  $\mathcal{F}_{R,M}$  and is the global minimizer of MSPBE that corresponds to the projection onto  $\mathcal{F}_{R,M}$  in (D.6).*

*Proof.* Proof of Lemma D.1 We first show that the operator  $\mathcal{T}_s^{\theta} : \mathbb{R}^{\Xi} \rightarrow \mathbb{R}^{\Xi}$  (3.13) is a  $\gamma$ -contraction in the  $L^2(\nu_{\theta})$ -norm.

$$\begin{aligned} & \|\mathcal{T}_s^{\theta} Q_1 - \mathcal{T}_s^{\theta} Q_2\|_{L^2(\nu_{\theta})}^2 = \mathbb{E}_{\zeta \sim \nu_{\theta}} \left[ \left( \mathcal{T}_s^{\theta} Q_1(\zeta) - \mathcal{T}_s^{\theta} Q_2(\zeta) \right)^2 \right] \\ & = \gamma^2 \mathbb{E}_{\zeta \sim \nu_{\theta}} \left[ \left( \mathbb{E} \left[ Q_1(\zeta') - Q_2(\zeta') \mid \zeta' = (\mu', h'), \mu' \sim P^N(\cdot \mid \zeta), h' \sim \Pi^{\theta}(\cdot \mid \mu') \right] \right)^2 \right] \\ & \leq \gamma^2 \mathbb{E}_{\zeta \sim \nu_{\theta}} \left[ \mathbb{E} \left[ (Q_1(\zeta') - Q_2(\zeta'))^2 \mid \zeta' = (\mu', h'), \mu' \sim P^N(\cdot \mid \zeta), h' \sim \Pi^{\theta}(\cdot \mid \mu') \right] \right] \\ & = \gamma^2 \mathbb{E}_{\zeta' \sim \nu_{\theta}} \left[ (Q_1(\zeta') - Q_2(\zeta'))^2 \right] = \gamma^2 \|Q_1 - Q_2\|_{L^2(\nu_{\theta})}^2, \end{aligned}$$

where the first inequality follows from Hölder's inequality for the conditional expectation and the third equality stems from the fact that  $\zeta'$  and  $\zeta$  have the same stationary distribution  $\nu_{\theta}$ .

Meanwhile, the projection operator  $\text{Proj}_{\mathcal{F}_{R,M}} : \mathbb{R}^{\Xi} \rightarrow \mathcal{F}_{R,M}$  is non-expansive. Therefore, the operator  $\text{Proj}_{\mathcal{F}_{R,M}} \mathcal{T}_s^{\theta} : \mathcal{F}_{R,M} \rightarrow \mathcal{F}_{R,M}$  is  $\gamma$ -contraction in the  $L^2(\nu_{\theta})$ -norm. Hence  $\text{Proj}_{\mathcal{F}_{R,M}}$  admits a unique fixed point  $Q_0(\cdot; \omega^*) \in \mathcal{F}_{R,M}$ . By definition,  $Q_0(\cdot; \omega^*)$  is the global minimizer of MSPBE that corresponds to the projection onto  $\mathcal{F}_{R,M}$  in (D.6).  $\square$

### Q.E.D.

We will show that the function class  $\mathcal{F}_{R,M}$  will approximately become  $\mathcal{F}_{R,\infty}^{s,k}$  (defined in Assumption 5.1) as  $M \rightarrow \infty$ , where  $\mathcal{F}_{R,\infty}^{s,k}$  is a rich reproducing kernel Hilbert space (RKHS). Consequently,  $Q_0(\cdot; \omega^*)$  will become the global minimum of the MSPBE (D.6) on  $\mathcal{F}_{R,\infty}^{s,k}$  given Lemma D.1. Moreover, by using similar argument and technique developed in [6, Theorem 4.6], we can establish the convergence of Algorithm 1 to  $Q_0(\cdot; \omega^*)$  as the following.

**Theorem D.2** (*Convergence to  $Q_0(\cdot; \omega^*)$* ) *Set  $\eta_{\text{critic}} = \min\{(1 - \gamma)/8, 1/\sqrt{T_{\text{critic}}}\}$  in Algorithm 1. Then the output  $Q_s(\cdot; \bar{\omega})$  of Algorithm 1 satisfies*

$$\mathbb{E}_{\text{init}} \left[ \|Q_s(\cdot; \bar{\omega}) - Q_0(\cdot; \omega^*)\|_{L^2(\nu_{\theta})}^2 \right] \leq \mathcal{O} \left( \frac{R^3 d_{\zeta^k}^{3/2}}{\sqrt{M}} + \frac{R^{5/2} d_{\zeta^k}^{5/4}}{\sqrt{M}} + \frac{R^2 d_{\zeta^k}}{\sqrt{T_{\text{critic}}}} \right),$$

where the expectation is taken with respect to the random initialization.

The proof of Theorem D.2 is straightforward from [6, Theorem 4.6] and hence omitted.

### D.1.3 Convergence to $Q_s^{\Pi^\theta}$

Next, we analyze the error between the global minimizer of (D.6) and the team-decentralized Q-function  $Q_s^{\Pi^\theta}$  (defined in (3.10)) to complete the convergence analysis. Different from the single-agent case as in Cai et al. [6], we have to bound an additional error from using the localized information in the critic update, in addition to the neural network approximation-optimization error.

*Proof.* Proof of Theorem 5.4 First recall that by Lemma 3.5,  $Q_s^{\Pi^\theta}$  satisfies the  $(c, \rho)$ -exponential decay property in Definition 3.4, with  $c = \frac{r_{\max}}{1-\gamma}$ ,  $\rho = \sqrt{\gamma}$ . Now, let  $\widehat{Q}_s^{\Pi^\theta}$  be any localized Q-function in (Local Q-function), then

$$\left| Q_s^{\Pi^\theta}(\zeta) - \widehat{Q}_s^{\Pi^\theta}(\zeta^k) \right| \leq c\rho^{k+1}, \quad \forall \zeta \in \Xi. \quad (\text{D.7})$$

By the triangle inequality and  $(a+b)^2 \leq 2(a^2 + b^2)$ ,

$$\begin{aligned} \left\| Q_s(\cdot; \bar{\omega}) - Q_s^{\Pi^\theta}(\cdot) \right\|_{L^2(\nu_\theta)}^2 &\leq \left( \left\| Q_s(\cdot; \bar{\omega}) - Q_0(\cdot; \omega^*) \right\|_{L^2(\nu_\theta)} + \left\| Q_s^{\Pi^\theta}(\cdot) - Q_0(\cdot; \omega^*) \right\|_{L^2(\nu_\theta)} \right)^2 \\ &\leq 2 \left( \left\| Q_s(\cdot; \bar{\omega}) - Q_0(\cdot; \omega^*) \right\|_{L^2(\nu_\theta)}^2 + \left\| Q_s^{\Pi^\theta}(\cdot) - Q_0(\cdot; \omega^*) \right\|_{L^2(\nu_\theta)}^2 \right). \end{aligned} \quad (\text{D.8})$$

The first term in (D.8) is studied in Theorem D.2 and it suffices to bound the second term. By interpolating two intermediate terms  $\widehat{Q}_s^{\Pi^\theta}$  and  $\text{Proj}_{\mathcal{F}_{R,M}} \widehat{Q}_s^{\Pi^\theta}$ , we have

$$\begin{aligned} \left\| Q_s^{\Pi^\theta}(\cdot) - Q_0(\cdot; \omega^*) \right\|_{L^2(\nu_\theta)} &\leq \underbrace{\left\| Q_s^{\Pi^\theta}(\cdot) - \widehat{Q}_s^{\Pi^\theta}(\cdot) \right\|_{L^2(\nu_\theta)}}_{\text{(I)}} + \underbrace{\left\| \widehat{Q}_s^{\Pi^\theta}(\cdot) - \text{Proj}_{\mathcal{F}_{R,M}} \widehat{Q}_s^{\Pi^\theta}(\cdot) \right\|_{L^2(\nu_\theta)}}_{\text{(II)}} \\ &\quad + \underbrace{\left\| Q_0(\cdot; \omega^*) - \text{Proj}_{\mathcal{F}_{R,M}} \widehat{Q}_s^{\Pi^\theta}(\cdot) \right\|_{L^2(\nu_\theta)}}_{\text{(III)}}. \end{aligned} \quad (\text{D.9})$$

First, we have (I)  $\leq c\rho^{k+1}$  according to (D.7). To bound (III), we have

$$\begin{aligned} \text{(III)} &= \left\| \text{Proj}_{\mathcal{F}_{R,M}} \mathcal{T}_s^\theta Q_0(\cdot; \omega^*) - \text{Proj}_{\mathcal{F}_{R,M}} \widehat{Q}_s^{\Pi^\theta}(\cdot) \right\|_{L^2(\nu_\theta)} \\ &\leq \left\| \text{Proj}_{\mathcal{F}_{R,M}} \mathcal{T}_s^\theta Q_0(\cdot; \omega^*) - \text{Proj}_{\mathcal{F}_{R,M}} \mathcal{T}_s^\theta Q_s^{\Pi^\theta}(\cdot) \right\|_{L^2(\nu_\theta)} + \left\| \text{Proj}_{\mathcal{F}_{R,M}} \mathcal{T}_s^\theta Q_s^{\Pi^\theta}(\cdot) - \text{Proj}_{\mathcal{F}_{R,M}} \widehat{Q}_s^{\Pi^\theta}(\cdot) \right\|_{L^2(\nu_\theta)} \\ &\leq \gamma \left\| Q_0(\cdot; \omega^*) - Q_s^{\Pi^\theta}(\cdot) \right\|_{L^2(\nu_\theta)} + \left\| \mathcal{T}_s^\theta Q_s^{\Pi^\theta}(\cdot) - \widehat{Q}_s^{\Pi^\theta}(\cdot) \right\|_{L^2(\nu_\theta)} \\ &= \gamma \left\| Q_0(\cdot; \omega^*) - Q_s^{\Pi^\theta}(\cdot) \right\|_{L^2(\nu_\theta)} + \underbrace{\left\| Q_s^{\Pi^\theta}(\cdot) - \widehat{Q}_s^{\Pi^\theta}(\cdot) \right\|_{L^2(\nu_\theta)}}_{\text{(I)}} \\ &\leq \gamma \left\| Q_0(\cdot; \omega^*) - Q_s^{\Pi^\theta}(\cdot) \right\|_{L^2(\nu_\theta)} + c\rho^{k+1}. \end{aligned} \quad (\text{D.10})$$

The first line in (D.10) is due to the fact that  $Q_0(\cdot; \omega^*)$  is the unique fixed point of the operator  $\text{Proj}_{\mathcal{F}_{R,M}} \mathcal{T}_s^\theta$ , (as proved in Lemma D.1); the third line in (D.10) is because the operator  $\text{Proj}_{\mathcal{F}_{R,M}} \mathcal{T}_s^\theta$  is a  $\gamma$ -contraction in the  $L^2(\nu_\theta)$  norm, and  $\text{Proj}_{\mathcal{F}_{R,M}}$  is non-expansive; the fourth line in (D.10) uses the fact that  $Q_s^{\Pi^\theta}$  is the unique fixed point of  $\mathcal{T}_s^\theta$ ; and the last line comes from the fact that (I)  $\leq c\rho^{k+1}$ . Therefore, combining the self-bounding inequality (D.10) with (D.9) and the bound on (I) gives us

$$\left\| Q_s^{\Pi^\theta}(\cdot) - Q_0(\cdot; \omega^*) \right\|_{L^2(\nu_\theta)} \leq \frac{1}{1-\gamma} \left( 2c\rho^{k+1} + \underbrace{\left\| \widehat{Q}_s^{\Pi^\theta}(\cdot) - \text{Proj}_{\mathcal{F}_{R,M}} \widehat{Q}_s^{\Pi^\theta}(\cdot) \right\|_{L^2(\nu_\theta)}}_{\text{(II)}} \right),$$

and consequently,

$$\left\| Q_s^{\Pi^\theta}(\cdot) - Q_0(\cdot; \omega^*) \right\|_{L^2(\nu_\theta)}^2 \leq \frac{1}{(1-\gamma)^2} \left( 8c^2 \rho^{2k+2} + 2 \underbrace{\left\| \widehat{Q}_s^{\Pi^\theta}(\cdot) - \text{Proj}_{\mathcal{F}_{R,M}} \widehat{Q}_s^{\Pi^\theta}(\cdot) \right\|_{L^2(\nu_\theta)}^2}_{\text{(II)}} \right). \quad (\text{D.11})$$

Plugging (D.11) into (D.8) yields

$$\begin{aligned} & \mathbb{E}_{\text{init}} \left[ \left\| Q_s(\cdot; \bar{\omega}) - Q_s^{\Pi^\theta}(\cdot) \right\|_{L^2(\nu_\theta)}^2 \right] \\ & \leq 2 \left( \mathbb{E}_{\text{init}} \left[ \left\| Q_s(\cdot; \bar{\omega}) - Q_0(\cdot; \omega^*) \right\|_{L^2(\nu_\theta)}^2 \right] + \mathbb{E}_{\text{init}} \left[ \left\| Q_s^{\Pi^\theta}(\cdot) - Q_0(\cdot; \omega^*) \right\|_{L^2(\nu_\theta)}^2 \right] \right) \\ & \leq \mathcal{O} \left( \frac{R^3 d_{\zeta_s^k}^{3/2}}{\sqrt{M}} + \frac{R^{5/2} d_{\zeta_s^k}^{5/4}}{\sqrt{M}} + \frac{R^2 d_{\zeta_s^k}}{\sqrt{T}} + c^2 \rho^{2k+2} \right) + \frac{4}{(1-\gamma)^2} \mathbb{E}_{\text{init}} \left[ \underbrace{\left\| \widehat{Q}_s^{\Pi^\theta}(\cdot) - \text{Proj}_{\mathcal{F}_{R,M}} \widehat{Q}_s^{\Pi^\theta}(\cdot) \right\|_{L^2(\nu_\theta)}^2}_{\text{(II)}} \right]. \end{aligned} \quad (\text{D.12})$$

Term (II) measures the distance between  $\widehat{Q}_s^{\Pi^\theta}$  and the class  $\mathcal{F}_{R,M}$ . As discussed in Section D.1.1, the function class  $\mathcal{F}_{R,M}$  converges to  $\mathcal{F}_{R,\infty}^{s,k}$  (defined in Assumption 5.1) as  $M \rightarrow \infty$ . Consequently, term (II) decreases as the neural network gets wider. To quantitatively characterize the approximation error between  $\mathcal{F}_{R,M}$  and  $\mathcal{F}_{R,\infty}^{s,k}$ , one needs the following lemma from Rahimi and Recht [46] and [6, Proposition 4.3]:

**Lemma D.3** *Assume Assumption 5.1, we have*

$$\mathbb{E}_{\text{init}} \left[ \underbrace{\left\| \widehat{Q}_s^{\Pi^\theta}(\cdot) - \text{Proj}_{\mathcal{F}_{R,M}} \widehat{Q}_s^{\Pi^\theta}(\cdot) \right\|_{L^2(\nu_\theta)}^2}_{\text{(II)}} \right] \leq \mathcal{O} \left( \frac{R^2 d_{\zeta_s^k}}{M} \right). \quad (\text{D.13})$$

With this lemma, Theorem 5.4 follows immediately by plugging (D.13) into (D.12), and setting  $c = \frac{r_{\max}}{1-\gamma}$ ,  $\rho = \sqrt{\gamma}$ ,  $T_{\text{critic}} = \Omega(M)$  in (D.12).  $\square$

**Q.E.D.**

## D.2 Proof of Theorem 5.10: Convergence of Actor Update

The proof of Theorem 5.10 consists of two steps: the first step in Section D.2.1 shows that the actor update converges to a stationary point of  $J$  (4.1), and the second step in Section D.2.2 bridges the gap between the stationary point and the optimality.

For the rest of this section, we use  $\eta$  to denote  $\eta_{\text{actor}}$  and  $\mathcal{B}_s$  to denote  $\mathcal{B}_s^{\text{actor}} := \{\theta_s \in \mathbb{R}^{M \times d_{\zeta_s}} : \|\theta_s - \theta_s(0)\|_\infty \leq R/\sqrt{M}\}$  for ease of notation. Meanwhile, define  $\mathcal{B} = \prod_{s \in \mathcal{S}} \mathcal{B}_s$ , the product space of  $\mathcal{B}_s$ 's, which is a convex set in  $\mathbb{R}^{M \times d_\zeta}$ .

### D.2.1 Convergence to Stationary Point

**Definition D.4** *A point  $\tilde{\theta} \in \mathcal{B}$  is called a stationary point of  $J(\cdot)$  if it holds that*

$$\nabla_{\theta} J(\tilde{\theta})^\top (\theta - \tilde{\theta}) \leq 0, \quad \forall \theta \in \mathcal{B}. \quad (\text{D.14})$$

Define the following mapping  $G$  from  $\mathbb{R}^{M \times d_c}$  to itself:

$$G(\theta) := \eta^{-1} \cdot [\text{Proj}_{\mathcal{B}}(\theta + \eta \cdot \nabla_{\theta} J(\theta)) - \theta]. \quad (\text{D.15})$$

It is well-known that (D.14) holds if and only if  $G(\tilde{\theta}) = 0$  (Sra et al. [49]). Now denote  $\rho(t) := G(\theta(t))$ , where  $\theta(t) = \{\theta_s(t)\}_{s \in \mathcal{S}}$  is the actor parameter updated in Algorithm 2 in iteration  $t$ .

To show that Algorithm 2 converges to a stationary point, we focus on analyzing  $\|\rho(t)\|_2$ .

**Theorem D.5** *Assume Assumptions 5.5 - 5.7. Set  $\eta = (T_{\text{actor}})^{-1/2}$  and assume  $1 - L\eta \geq 1/2$ , where  $L$  is the Lipschitz constant in Assumption 5.7. Then the output  $\{\theta(t)\}_{t \in [T_{\text{actor}}]}$  of Algorithm 2 satisfies*

$$\min_{t \in [T_{\text{actor}}]} \mathbb{E} [\|\rho(t)\|_2^2] \leq \frac{8\tau^2 \Sigma^2 |\mathcal{S}|}{B} + \frac{4}{\sqrt{T_{\text{actor}}}} \mathbb{E}[J(\theta(T_{\text{actor}} + 1)) - J(\theta(1))] + \epsilon_Q(T_{\text{actor}}). \quad (\text{D.16})$$

Here  $\epsilon_Q$  measures the error accumulated from the critic steps which is defined as

$$\begin{aligned} \epsilon_Q(T_{\text{actor}}) &= \frac{32\tau D R d_{\xi_s}^{1/2} |\mathcal{S}|}{(1-\gamma)\eta T_{\text{actor}}} \cdot \sum_{t=1}^{T_{\text{actor}}} \sum_{s \in \mathcal{S}} \mathbb{E} \left[ \left\| Q_s(\cdot; \bar{\omega}_s, t) - Q_s^{\Pi^{\theta(t)}}(\cdot) \right\|_{L^2(\nu_{\theta(t)})} \right] \\ &\quad + \frac{16\tau^2 D^2 |\mathcal{S}|^2}{(1-\gamma)^2 T_{\text{actor}}} \cdot \sum_{t=1}^{T_{\text{actor}}} \sum_{s \in \mathcal{S}} \mathbb{E} \left[ \left\| Q_s(\cdot; \bar{\omega}_s, t) - Q_s^{\Pi^{\theta(t)}}(\cdot) \right\|_{L^2(\nu_{\theta(t)})}^2 \right], \end{aligned} \quad (\text{D.17})$$

where  $\{Q_s(\cdot; \bar{\omega}_s, t)\}_{s \in \mathcal{S}}$  is the output of the critic update at step  $t$  in Algorithm 2. All expectations in (D.16) and (D.17) are taken over all randomness in Algorithm 1 and Algorithm 2.

*Proof.* Proof of Theorem D.5

Let  $t \in [T_{\text{actor}}]$ , we first lower bound the difference between the expected total rewards of  $\Pi^{\theta(t+1)}$  and  $\Pi^{\theta(t)}$ . By Assumption 5.7,  $\nabla_{\theta} J(\theta)$  is  $L$ -Lipschitz continuous. Hence by Taylor's expansion,

$$J(\theta(t+1)) - J(\theta(t)) \geq \eta \cdot \nabla_{\theta} J(\theta(t))^{\top} \delta(t) - L/2 \cdot \|\theta(t+1) - \theta(t)\|_2^2, \quad (\text{D.18})$$

where  $\delta(t) = (\theta(t+1) - \theta(t)) / \eta$ . Meanwhile denote  $\xi_s(t) = \hat{g}_s(\theta(t)) - \mathbb{E}[\hat{g}_s(\theta(t))]$ , where  $\hat{g}_s(\theta(t))$  is defined in (4.14) and the expectation is taken over  $\sigma_{\theta(t)}$  given  $\{\bar{\omega}_s\}_{s \in \mathcal{S}}$ . Then

$$\begin{aligned} \nabla_{\theta} J(\theta(t))^{\top} \delta(t) &= \sum_{s \in \mathcal{S}} \nabla_{\theta_s} J(\theta(t))^{\top} \delta_s(t) \\ &= \sum_{s \in \mathcal{S}} \left[ (\nabla_{\theta_s} J(\theta(t)) - \mathbb{E}[\hat{g}_s(\theta(t))])^{\top} \delta_s(t) - \xi_s(t)^{\top} \delta_s(t) + \hat{g}_s(\theta(t))^{\top} \delta_s(t) \right], \end{aligned} \quad (\text{D.19})$$

where  $\delta_s(t) := (\theta_s(t+1) - \theta_s(t)) / \eta$ . The first term in (D.19) represents the error of estimating  $\nabla_{\theta_s} J(\theta(t))$  using

$$\mathbb{E}[\hat{g}_s(\theta(t))] = \frac{1}{1-\gamma} \mathbb{E}_{\sigma_{\theta(t)}} \left[ \left[ \sum_{y \in \mathcal{N}_s^k} Q_y(\mu(\mathcal{N}_y^k), h(\mathcal{N}_y^k); \bar{\omega}_y, t) \right] \nabla_{\theta_s} \log \Pi^{\theta_s}(h(s) | \mu(s)) \right].$$

To bound the first term, first notice that

$$\mathbb{E}[\hat{g}_s(\theta(t))] = \frac{1}{1-\gamma} \mathbb{E}_{\sigma_{\theta(t)}} \left[ \left[ \sum_{y \in \mathcal{S}} Q_y(\mu(\mathcal{N}_y^k), h(\mathcal{N}_y^k); \bar{\omega}_y, t) \right] \nabla_{\theta_s} \log \Pi^{\theta_s}(h(s) | \mu(s)) \right].$$

This is because for all  $y \notin \mathcal{N}_s^k$ ,  $Q_y(\mu(\mathcal{N}_y^k), h(\mathcal{N}_y^k); \bar{\omega}_y)$  is independent of  $s$  and consequently, we can verify that

$$\mathbb{E}_{\sigma_{\theta(t)}} \left[ \left[ \sum_{y \notin \mathcal{N}_s^k} Q_y(\mu(\mathcal{N}^k(y)), h(\mathcal{N}^k(y)); \bar{\omega}_y, t) \right] \nabla_{\theta_s} \log \Pi^{\theta_s}(h(s) | \mu(s)) \right] = 0.$$

Therefore, following the similar computation in Lemma D.2, Cai et al. [6], we have

$$\left| (\nabla_{\theta_s} J(\theta(t)) - \mathbb{E}[\hat{g}_s(\theta(t))])^\top \delta_s(t) \right| \leq \frac{4\tau DRd_{\zeta_s}^{1/2}}{(1-\gamma)\eta} \sum_{s \in \mathcal{S}} \left\| Q_s(\cdot; \bar{\omega}_s, t) - Q_s^{\theta(t)}(\cdot) \right\|_{L^2(\nu_{\theta(t)})}. \quad (\text{D.20})$$

To bound the second term in (D.19), we simply have

$$\xi_s(t)^\top \delta_s(t) \leq \|\xi_s(t)\|_2^2 + \|\delta_s(t)\|_2^2. \quad (\text{D.21})$$

To handle the last term in (D.19), we have

$$\begin{aligned} & \hat{g}_s(\theta(t))^\top \delta_s(t) - \|\delta_s(t)\|_2^2 = \eta^{-1} \cdot (\eta \hat{g}_s(\theta(t)) - (\theta_s(t+1) - \theta_s(t)))^\top \delta_s \\ & = \eta^{-1} \cdot (\theta_s(t+1/2) - \text{Proj}_{\mathcal{B}_s}(\theta_s(t+1/2)))^\top \delta_s(t) \\ & = \eta^{-2} \cdot (\theta_s(t+1/2) - \text{Proj}_{\mathcal{B}_s}(\theta_s(t+1/2)))^\top (\text{Proj}_{\mathcal{B}_s}(\theta_s(t+1/2)) - \theta_s(t)) \geq 0 \end{aligned} \quad (\text{D.22})$$

Here we write  $\theta_s(t) + \eta \hat{g}_s(\theta(t))$  as  $\theta_s(t+1/2)$  to simplify the notation. The last inequality comes from the property of the projection onto a convex set.

Therefore, combining (D.19), (D.20), (D.21) and (D.22) suggests

$$\nabla_{\theta_s} J(\theta(t))^\top \delta_s(t) \geq -\frac{4\tau DRd_{\zeta_s}^{1/2}}{(1-\gamma)\eta} \sum_{s \in \mathcal{S}} \left[ \left\| Q_s(\cdot; \bar{\omega}_s, t) - Q_s^{\theta(t)}(\cdot) \right\|_{L^2(\nu_{\theta(t)})} \right] + \frac{1}{2} (\|\delta_s(t)\|_2^2 - \|\xi_s(t)\|_2^2).$$

Consequently,

$$\nabla_{\theta} J(\theta(t))^\top \delta(t) \geq -\frac{4\tau DRd_{\zeta_s}^{1/2}}{(1-\gamma)\eta} |\mathcal{S}| \sum_{s \in \mathcal{S}} \left[ \left\| Q_s(\cdot; \bar{\omega}_s, t) - Q_s^{\theta(t)}(\cdot) \right\|_{L^2(\nu_{\theta(t)})} \right] + \frac{1}{2} (\|\delta(t)\|_2^2 - \|\xi(t)\|_2^2). \quad (\text{D.23})$$

Thus, by plugging (D.23) into (D.18) and by Assumption 5.5, we have

$$\begin{aligned} \frac{1-L\cdot\eta}{2} \mathbb{E}[\|\delta(t)\|_2^2] & \leq \eta^{-1} \cdot \mathbb{E}[J(\theta(t+1)) - J(\theta(t))] + \frac{\tau^2 \Sigma^2 |\mathcal{S}|}{2B} \\ & \quad + \frac{4\tau DRd_{\zeta_s}^{1/2} |\mathcal{S}|}{(1-\gamma)\eta} \sum_{s \in \mathcal{S}} \left\| Q_s(\cdot; \bar{\omega}_s, t) - Q_s^{\theta(t)}(\cdot) \right\|_{L^2(\nu_{\theta(t)})}. \end{aligned} \quad (\text{D.24})$$

Here the expectation is taken over  $\sigma_{\theta(t)}$  given  $\{\bar{\omega}_s\}_{s \in \mathcal{S}}$ .

Now, in order to bridge the gap between  $\|\delta(t)\|_2$  in (D.24) and  $\|\rho(t)\|_2 = \|G(\theta(t))\|_2$  in (D.15), we next will bound the difference  $\|\delta(t) - \rho(t)\|_2$ . We start with defining a local gradient mapping  $G_s$  from  $\mathbb{R}^{M \times d_\zeta}$  to  $\mathbb{R}^{M \times d_{\zeta_s}}$ :

$$G_s(\theta) := \eta^{-1} \cdot [\text{Proj}_{\mathcal{B}_s}(\theta_s + \eta \cdot \nabla_{\theta_s} J(\theta)) - \theta_s]. \quad (\text{D.25})$$

Since  $\mathcal{B}_s$  is an  $l_\infty$ -ball around the initialization, it is easy to verify that  $G_s(\theta) = (G(\theta))_s$ . Therefore, we can further define  $\rho_s(t) = G_s(\theta(t))$  and the following decomposition holds:

$$\|\delta(t) - \rho(t)\|_2^2 = \sum_{s \in \mathcal{S}} \|\delta_s(t) - \rho_s(t)\|_2^2.$$

From the definitions of  $\bar{\delta}_s(t)$  and  $\rho_s(t)$ ,

$$\begin{aligned}\|\delta_s(t) - \rho_s(t)\|_2 &= \eta^{-1} \cdot \|\text{Proj}_{\mathcal{B}_s}(\theta_s + \eta \cdot \nabla_{\theta_s} J(\theta)) - \theta_s - \text{Proj}_{\mathcal{B}_s}(\theta_s + \eta \cdot \hat{g}_s(\theta)) + \theta_s\|_2 \\ &= \eta^{-1} \cdot \|\text{Proj}_{\mathcal{B}_s}(\theta_s + \eta \cdot \nabla_{\theta_s} J(\theta)) - \text{Proj}_{\mathcal{B}_s}(\theta_s + \eta \cdot \hat{g}_s(\theta))\|_2 \\ &\leq \eta^{-1} \cdot \|\theta_s + \eta \cdot \nabla_{\theta_s} J(\theta) - \theta_s + \eta \cdot \hat{g}_s(\theta)\|_2 = \|\nabla_{\theta_s} J(\theta) - \hat{g}_s(\theta)\|_2\end{aligned}$$

Following similar calculations in [6, Lemma D.3],

$$\begin{aligned}\mathbb{E} [\|\nabla_{\theta_s} J(\theta) - \hat{g}_s(\theta)\|_2^2] &\leq \frac{2\tau^2 \Sigma^2}{B} + \frac{8\tau^2 D^2}{(1-\gamma)^2} \left( \sum_{s \in \mathcal{S}} \left\| Q_s(\cdot; \bar{\omega}_s, t) - Q_s^{\Pi^{\theta(t)}}(\cdot) \right\|_{L^2(\nu_{\theta(t)})} \right)^2 \\ &\leq \frac{2\tau^2 \Sigma^2}{B} + \frac{8\tau^2 D^2 |\mathcal{S}|}{(1-\gamma)^2} \left( \sum_{s \in \mathcal{S}} \left\| Q_s(\cdot; \bar{\omega}_s, t) - Q_s^{\Pi^{\theta(t)}}(\cdot) \right\|_{L^2(\nu_{\theta(t)})}^2 \right).\end{aligned}\tag{D.26}$$

The expectation is taken over  $\sigma_{\theta(t)}$  given  $\{\bar{\omega}_s\}_{s \in \mathcal{S}}$ . Consequently,

$$\mathbb{E} [\|\delta(t) - \rho(t)\|_2^2] \leq \frac{2\tau^2 \Sigma^2 |\mathcal{S}|}{B} + \frac{8\tau^2 D^2 |\mathcal{S}|^2}{(1-\gamma)^2} \left( \sum_{s \in \mathcal{S}} \left\| Q_s(\cdot; \bar{\omega}_s, t) - Q_s^{\Pi^{\theta(t)}}(\cdot) \right\|_{L^2(\nu_{\theta(t)})}^2 \right).\tag{D.27}$$

Set  $\eta = 1/\sqrt{T_{\text{actor}}}$  and take (D.24) and (D.27), we obtain (D.16) from the following estimations:

$$\begin{aligned}\min_{t \in [T_{\text{actor}}]} \mathbb{E} [\|\rho(t)\|_2^2] &\leq \frac{1}{T_{\text{actor}}} \cdot \sum_{t=1}^{T_{\text{actor}}} \|\rho(t)\|_2^2 \leq \frac{2}{T_{\text{actor}}} \cdot \sum_{t=1}^{T_{\text{actor}}} (\mathbb{E} [\|\delta(t) - \rho(t)\|_2^2] + \mathbb{E} [\|\delta(t)\|_2^2]) \\ &\leq \frac{2}{T_{\text{actor}}} \cdot \sum_{t=1}^{T_{\text{actor}}} (\mathbb{E} [\|\delta(t) - \rho(t)\|_2^2] + 2(1-L \cdot \eta) \mathbb{E} [\|\delta(t)\|_2^2]) \\ &\leq \frac{8\tau^2 \Sigma^2 |\mathcal{S}|}{B} + \frac{4}{\sqrt{T_{\text{actor}}}} \mathbb{E} [J(\theta(T_{\text{actor}} + 1)) - J(\theta(1))] + \epsilon_Q(T_{\text{actor}}),\end{aligned}$$

where  $\epsilon_Q$  measures the error accumulated from the critic steps which is defined in (D.17), i.e.,

$$\begin{aligned}\epsilon_Q(T_{\text{actor}}) &= \frac{32\tau D R d_{\zeta_s}^{1/2} |\mathcal{S}|}{(1-\gamma)\eta T_{\text{actor}}} \cdot \sum_{t=1}^{T_{\text{actor}}} \sum_{s \in \mathcal{S}} \mathbb{E} \left[ \left\| Q_s(\cdot; \bar{\omega}_s) - Q_s^{\Pi^{\theta(t)}}(\cdot) \right\|_{L^2(\nu_{\theta(t)})} \right] \\ &\quad + \frac{16\tau^2 D^2 |\mathcal{S}|^2}{(1-\gamma)^2 T_{\text{actor}}} \cdot \sum_{t=1}^{T_{\text{actor}}} \sum_{s \in \mathcal{S}} \mathbb{E} \left[ \left\| Q_s(\cdot; \bar{\omega}_s) - Q_s^{\Pi^{\theta(t)}}(\cdot) \right\|_{L^2(\nu_{\theta(t)})}^2 \right].\end{aligned}$$

Here the expectations in (D.16) and (D.17) are taken over all randomness in Algorithm 1 and Algorithm 2.  $\square$

**Q.E.D.**

## D.2.2 Bridging the gap between Stationarity and Optimality

Recall that  $\sigma_\theta$  in (4.2) denotes the state-action visitation measure under policy  $\Pi^\theta$ . Denote  $\bar{\sigma}_\theta$  as the state visitation measure under policy  $\Pi^\theta$ . Consequently,

$$\bar{\sigma}_\theta(\mu) \Pi^\theta(h | \mu) = \sigma_\theta(\mu, h).$$

Following similar steps in the proof of [6, Theorem 4.8], one can characterize the global optimality of the obtained stationary point  $\bar{\theta} \in \mathcal{B}$  as the following.

**Lemma D.6** Let  $\tilde{\theta} \in \mathcal{B}$  be a stationary point of  $J(\cdot)$  satisfying condition (D.14) and let  $\theta^* \in \mathcal{B}$  be the global maximum point of  $J(\cdot)$  in  $\mathcal{B}$ . Then the following inequality holds:

$$(1 - \gamma) \left( J(\theta^*) - J(\tilde{\theta}) \right) \leq \frac{2r_{\max}}{1 - \gamma} \inf_{\theta \in \mathcal{B}} \left\| u_{\tilde{\theta}}(\mu, h) - \sum_{s \in \mathcal{S}} \phi_{\tilde{\theta}_s}(\mu(s), h(s))^\top \tilde{\theta}_s \right\|_{L^2(\sigma_{\tilde{\theta}})}, \quad (\text{D.28})$$

where  $u_{\tilde{\theta}}(\mu, h) := \frac{d\sigma_{\theta^*}}{d\sigma_{\tilde{\theta}}}(\mu, h) - \frac{d\bar{\sigma}_{\theta^*}}{d\bar{\sigma}_{\tilde{\theta}}}(\mu) + \sum_{s \in \mathcal{S}} \phi_{\tilde{\theta}_s}(\mu(s), h(s))^\top \tilde{\theta}_s$ , and  $\frac{d\sigma_{\theta^*}}{d\sigma_{\tilde{\theta}}}, \frac{d\bar{\sigma}_{\theta^*}}{d\bar{\sigma}_{\tilde{\theta}}}$  are the Radon-Nikodym derivatives between the corresponding measures.

*Proof.* Proof of Lemma D.6 First recall that by (4.8), for any  $\theta \in \mathcal{B}$ ,

$$\nabla_{\theta} J(\tilde{\theta})^\top (\theta - \tilde{\theta}) = \sum_{s \in \mathcal{S}} \nabla_{\theta_s} J(\tilde{\theta})^\top (\theta_s - \tilde{\theta}_s) = \frac{\tau}{1 - \gamma} \sum_{s \in \mathcal{S}} \mathbb{E}_{\sigma_{\tilde{\theta}}} \left[ Q^{\Pi^{\tilde{\theta}}}(\mu, h) \cdot \Phi(\tilde{\theta}, s, \mu, h)^\top (\theta_s - \tilde{\theta}_s) \right],$$

in which  $\Phi(\theta, s, \mu, h) := \phi_{\theta_s}(\mu(s), h(s)) - \mathbb{E}_{h(s)' \sim \Pi_s^{\theta_s}(\cdot | \mu(s))} [\phi_{\theta_s}(\mu(s), h'(s))]$  is defined in (4.7).

Since  $\tilde{\theta} \in \mathcal{B}$  is a stationary point of  $J(\cdot)$ ,

$$\sum_{s \in \mathcal{S}} \mathbb{E}_{\sigma_{\tilde{\theta}}} \left[ Q^{\Pi^{\tilde{\theta}}}(\mu, h) \cdot \Phi(\tilde{\theta}, s, \mu, h)^\top (\theta_s - \tilde{\theta}_s) \right] \leq 0, \quad \forall \theta \in \mathcal{B}. \quad (\text{D.29})$$

Denote  $A^{\Pi^{\tilde{\theta}}}(\mu, h) := Q^{\Pi^{\tilde{\theta}}}(\mu, h) - V^{\Pi^{\tilde{\theta}}}(\mu)$  as the advantage function under policy  $\Pi^{\tilde{\theta}}$ . It holds from the definition that  $\mathbb{E}_{h \sim \Pi^{\tilde{\theta}}(\cdot | \mu)} [A^{\Pi^{\tilde{\theta}}}(\mu, h)] = V^{\Pi^{\tilde{\theta}}}(\mu) - V^{\Pi^{\tilde{\theta}}}(\mu) = 0$ . Meanwhile,  $\sup_{(\mu, h) \in \Xi} |A^{\Pi^{\tilde{\theta}}}(\mu, h)| \leq 2 \sup_{\mu \in \mathcal{P}^N(\mathcal{S})} |V^{\Pi^{\tilde{\theta}}}(\mu)| \leq \frac{2r_{\max}}{1 - \gamma}$ .

Given that  $\mathbb{E}_{h \sim \Pi^{\tilde{\theta}}(\cdot | \mu)} [A^{\Pi^{\tilde{\theta}}}(\mu, h)] = 0$  and  $\mathbb{E}_{h \sim \Pi^{\tilde{\theta}}(\cdot | \mu)} [\Phi(\tilde{\theta}, s, \mu, h)] = 0$ , we have for any  $s \in \mathcal{S}$ ,

$$\mathbb{E}_{\sigma_{\tilde{\theta}}} \left[ V^{\Pi^{\tilde{\theta}}}(\mu) \cdot \Phi(\tilde{\theta}, s, \mu, h) \right] = 0, \quad \text{and} \quad (\text{D.30})$$

$$\mathbb{E}_{\sigma_{\tilde{\theta}}} \left[ A^{\Pi^{\tilde{\theta}}}(\mu, h) \cdot \mathbb{E}_{h(s)' \sim \Pi_s^{\tilde{\theta}_s}(\cdot | \mu(s))} [\phi_{\tilde{\theta}_s}(\mu(s), h'(s))] \right] = 0. \quad (\text{D.31})$$

Combining (D.29) with (D.30) and (D.31),

$$\sum_{s \in \mathcal{S}} \mathbb{E}_{\sigma_{\tilde{\theta}}} \left[ A^{\Pi^{\tilde{\theta}}}(\mu, h) \cdot \phi_{\tilde{\theta}_s}(\mu(s), h(s))^\top (\theta_s - \tilde{\theta}_s) \right] \leq 0, \quad \forall \theta \in \mathcal{B}. \quad (\text{D.32})$$

Moreover, by the Performance Difference Lemma (Kakade and Langford [32]),

$$(1 - \gamma) \cdot \left( J(\theta^*) - J(\hat{\theta}) \right) = \mathbb{E}_{\sigma_{\theta^*}} \left[ \left\langle A^{\Pi^{\tilde{\theta}}}(\mu, \cdot), \Pi^{\theta^*}(\cdot | \mu) - \Pi^{\tilde{\theta}}(\cdot | \mu) \right\rangle \right]. \quad (\text{D.33})$$

Combining (D.33) with (D.32), it holds that for any  $\theta \in \mathcal{B}$ ,

$$\begin{aligned} & (1 - \gamma) \cdot \left( J(\theta^*) - J(\hat{\theta}) \right) \\ & \leq \mathbb{E}_{\sigma_{\theta^*}} \left[ \left\langle A^{\Pi^{\tilde{\theta}}}(\mu, \cdot), \Pi^{\theta^*}(\cdot | \mu) - \Pi^{\tilde{\theta}}(\cdot | \mu) \right\rangle \right] - \sum_{s \in \mathcal{S}} \mathbb{E}_{\sigma_{\tilde{\theta}}} \left[ A^{\Pi^{\tilde{\theta}}}(\zeta) \cdot \phi_{\tilde{\theta}_s}(\zeta_s)^\top (\theta_s - \tilde{\theta}_s) \right] \\ & = \mathbb{E}_{\sigma_{\tilde{\theta}}} \left[ A^{\Pi^{\tilde{\theta}}}(\mu, h) \cdot \left( \frac{d\sigma_{\theta^*}}{d\sigma_{\tilde{\theta}}}(\mu, h) - \frac{d\bar{\sigma}_{\theta^*}}{d\bar{\sigma}_{\tilde{\theta}}}(\mu) - \sum_{s \in \mathcal{S}} \phi_{\tilde{\theta}_s}(\mu(s), h(s))^\top (\theta_s - \tilde{\theta}_s) \right) \right]. \end{aligned} \quad (\text{D.34})$$

Therefore,

$$\begin{aligned} & (1 - \gamma) \cdot \left( J(\theta^*) - J(\hat{\theta}) \right) \\ & \leq \frac{2r_{\max}}{1 - \gamma} \inf_{\theta \in \mathcal{B}} \left\| \frac{d\sigma_{\theta^*}}{d\sigma_{\tilde{\theta}}}(\mu, h) - \frac{d\bar{\sigma}_{\theta^*}}{d\bar{\sigma}_{\tilde{\theta}}}(\mu) - \sum_{s \in \mathcal{S}} \phi_{\tilde{\theta}_s}(\mu(s), h(s))^\top (\theta_s - \tilde{\theta}_s) \right\|_{L^2(\sigma_{\tilde{\theta}})} \\ & = \frac{2r_{\max}}{1 - \gamma} \inf_{\theta \in \mathcal{B}} \left\| u_{\tilde{\theta}}(\mu, h) - \sum_{s \in \mathcal{S}} \phi_{\tilde{\theta}_s}(\mu(s), h(s))^\top \tilde{\theta}_s \right\|_{L^2(\sigma_{\tilde{\theta}})}, \end{aligned} \quad (\text{D.35})$$

where  $u_{\tilde{\theta}}(\mu, h) := \frac{d\sigma_{\theta^*}}{d\sigma_{\tilde{\theta}}}(\mu, h) - \frac{d\tilde{\sigma}_{\theta^*}}{d\tilde{\sigma}_{\tilde{\theta}}}(\mu) + \sum_{s \in \mathcal{S}} \phi_{\tilde{\theta}_s}(\mu(s), h(s))^\top \tilde{\theta}_s$ , and  $\frac{d\sigma_{\theta^*}}{d\sigma_{\tilde{\theta}}}, \frac{d\tilde{\sigma}_{\theta^*}}{d\tilde{\sigma}_{\tilde{\theta}}}$  are the Radon-Nikodym derivatives between corresponding measures.  $\square$

**Q.E.D.**

To further bound the right-hand-side of (D.28) in Lemma D.6, define the following function class

$$\tilde{\mathcal{F}}_{R,M} = \left\{ f_0(\zeta; \theta) := \underbrace{\sum_{s \in \mathcal{S}} \left[ \frac{1}{\sqrt{M}} \sum_{m=1}^M \mathbb{1} \{ [\theta_s(0)]_m^\top \zeta_s > 0 \} [\theta_s]_m^\top \zeta_s \right]}_{(\star)} : \theta_s \in \mathbb{R}^{M \times d_{\zeta_s}}, \|\theta_s - \theta_s(0)\|_\infty \leq R/\sqrt{M} \right\}, \quad (\text{D.36})$$

given an initialization  $\theta_s(0) \in \mathbb{R}^{M \times d_{\zeta_s}}$ ,  $s \in \mathcal{S}$  and  $b \in \mathbb{R}^M$ .

$\tilde{\mathcal{F}}_{R,M}$  (D.36) is a local linearization of the actor neural network. More specifically, term  $(\star)$  in (D.36) locally linearizes the decentralized actor neural network  $f(\zeta_s; \theta_s)$  (4.4) with respect to  $\theta_s$ . Any  $f_0(\zeta; \theta) \in \tilde{\mathcal{F}}_{R,M}$  is a sum of  $|\mathcal{S}|$  inner products between feature mapping  $\phi_{\theta_s(0)}(\cdot)$  (4.6) and parameter  $\theta_s$ :  $f_0(\zeta; \theta) = \sum_{s \in \mathcal{S}} \phi_{\theta_s(0)}(\zeta_s) \cdot \theta_s$ . As the width of the neural network  $M \rightarrow \infty$ ,  $\tilde{\mathcal{F}}_{R,M}$  converges to  $\mathcal{F}_{R,\infty}$  (defined in Assumption 5.9). The approximation error between  $\tilde{\mathcal{F}}_{R,M}$  and  $\mathcal{F}_{R,\infty}$  is bounded in the following lemma.

**Lemma D.7** *For any function  $f(\zeta) \in \mathcal{F}_{R,\infty}$  defined in Assumption 5.9, we have*

$$\mathbb{E}_{\text{init}} \left[ \left\| f(\cdot) - \text{Proj}_{\tilde{\mathcal{F}}_{R,M}} f(\cdot) \right\|_{L^2(\sigma_{\tilde{\theta}})} \right] \leq \mathcal{O} \left( \frac{|\mathcal{S}| R d_{\zeta_s}^{1/2}}{M^{1/2}} \right). \quad (\text{D.37})$$

Lemma D.7 follows from Rahimi and Recht [46] and [6, Proposition 4.3]. The factor  $|\mathcal{S}|$  stems from the fact that  $\mathcal{F}_{R,\infty}$  can be decomposed into  $|\mathcal{S}|$  independent reproducing kernel Hilbert spaces. With Lemma D.7, we are ready to establish an upper bound for the right-hand-side of (D.28) in the following proposition.

**Proposition D.8** *Under Assumption 5.9, let  $\tilde{\theta} \in \mathcal{B}$  be a stationary point of  $J(\cdot)$  and let  $\theta^* \in \mathcal{B}$  be the global maximum point of  $J(\cdot)$  in  $\mathcal{B}$ . Then the following inequality holds:*

$$(1 - \gamma) \left( J(\theta^*) - J(\tilde{\theta}) \right) \leq \mathcal{O} \left( \frac{|\mathcal{S}| R^{3/2} d_{\zeta_s}^{3/4}}{M^{1/4}} \right). \quad (\text{D.38})$$

*Proof.* Proof of Proposition D.8 First by the triangle inequality,

$$\begin{aligned} \inf_{\theta \in \mathcal{B}} \left\| u_{\tilde{\theta}}(\zeta) - \sum_{s \in \mathcal{S}} \phi_{\tilde{\theta}_s}(\zeta_s)^\top \theta_s \right\|_{L^2(\sigma_{\tilde{\theta}})} &\leq \left\| u_{\tilde{\theta}}(\zeta) - \text{Proj}_{\tilde{\mathcal{F}}_{R,M}} u_{\tilde{\theta}}(\zeta) \right\|_{L^2(\sigma_{\tilde{\theta}})} \\ &+ \inf_{\theta \in \mathcal{B}} \left\| \text{Proj}_{\tilde{\mathcal{F}}_{R,M}} u_{\tilde{\theta}}(\zeta) - \sum_{s \in \mathcal{S}} \phi_{\tilde{\theta}_s}(\zeta_s)^\top \theta_s \right\|_{L^2(\sigma_{\tilde{\theta}})}, \end{aligned} \quad (\text{D.39})$$

where  $\tilde{\mathcal{F}}_{R,M}$  is defined in (D.36). We denote  $\text{Proj}_{\tilde{\mathcal{F}}_{R,M}} u_{\tilde{\theta}}(\zeta) = \sum_{s \in \mathcal{S}} \phi_{\theta_s(0)}(\zeta_s) \cdot \hat{\theta}_s \in \tilde{\mathcal{F}}_{R,M}$  for some  $\hat{\theta} \in \mathcal{B}$ . Therefore, by Lemma D.7, the first term on the right-hand-side of (D.39) is bounded by (D.37):

$$\left\| u_{\tilde{\theta}}(\zeta) - \sum_{s \in \mathcal{S}} \phi_{\theta_s(0)}(\zeta_s) \cdot \hat{\theta}_s \right\|_{L^2(\sigma_{\tilde{\theta}})} \leq \mathcal{O} \left( \frac{|\mathcal{S}| R d_{\zeta_s}^{1/2}}{M^{1/2}} \right).$$

The following Lemma D.9 is a direct application of [53, Lemma E.2], which is used to bound the second term on the right-hand-side of (D.39).

**Lemma D.9** *It holds for any  $\theta_s, \theta'_s \in \mathcal{B}_s = \{\alpha_s \in \mathbb{R}^{M \times d_{\zeta_s}} : \|\alpha_s - \theta_s(0)\|_\infty \leq R/\sqrt{M}\}$  that*

$$\mathbb{E}_{\text{init}} \left[ \|\phi_{\theta_s}(\zeta_s)^\top \theta'_s - \phi_{\theta_s(0)}(\zeta_s)^\top \theta'_s\|_{L^2(\sigma_\theta)} \right] \leq \mathcal{O} \left( \frac{R^{3/2} d_{\zeta_s}^{3/4}}{M^{1/4}} \right), \quad (\text{D.40})$$

where the expectation is taken over random initialization.

Taking  $\theta = \tilde{\theta}$  and  $\theta' = \hat{\theta}$  in Lemma D.9 gives us

$$\sum_{s \in \mathcal{S}} \left\| \phi_{\theta_s(0)}(\zeta_s) \cdot \hat{\theta}_s - \phi_{\tilde{\theta}_s}(\zeta_s)^\top \hat{\theta}_s \right\|_{L^2(\sigma_{\tilde{\theta}})} \leq \mathcal{O} \left( \frac{|\mathcal{S}| R^{3/2} d_{\zeta_s}^{3/4}}{M^{1/4}} \right)$$

Therefore, by Lemma D.1,

$$(1 - \gamma) \left( J(\theta^*) - J(\tilde{\theta}) \right) \leq \inf_{\theta \in \mathcal{B}} \left\| u_{\tilde{\theta}}(\zeta) - \sum_{s \in \mathcal{S}} \phi_{\tilde{\theta}_s}(\zeta_s)^\top \theta_s \right\|_{L^2(\sigma_{\tilde{\theta}})} \leq \mathcal{O} \left( \frac{|\mathcal{S}| R^{3/2} d_{\zeta_s}^{3/4}}{M^{1/4}} \right).$$

□

**Q.E.D.**

Now we are ready to establish Theorem 5.10.

*Proof.* Proof of Theorem 5.10 Following similar calculations as in [53, Section H.3], we obtain that at iteration  $t \in [T_{\text{actor}}]$ ,

$$\nabla_{\theta} J(\theta(t))^\top (\theta - \theta(t)) \leq 2 \left( R + \frac{\eta \cdot r_{\max}}{1 - \gamma} \right) \cdot \|\rho(t)\|_2, \quad \forall \theta \in \mathcal{B}. \quad (\text{D.41})$$

The right-hand-side of (D.41) quantifies the deviation of  $\theta(t)$  from a stationary point  $\tilde{\theta}$ . Having (D.41) and following similar arguments for Lemma D.6 and Proposition D.8, we can show that

$$(1 - \gamma) \min_{t \in [T_{\text{actor}}]} \mathbb{E} [J(\theta^*) - J(\theta(t))] \leq \mathcal{O} \left( \frac{|\mathcal{S}| R^{3/2} d_{\zeta_s}^{3/4}}{M^{1/4}} \right) + 2 \left( R + \frac{\eta \cdot r_{\max}}{1 - \gamma} \right) \cdot \min_{t \in [T_{\text{actor}}]} \mathbb{E} [\|\rho(t)\|_2]. \quad (\text{D.42})$$

Here the last term  $\min_{t \in [T_{\text{actor}}]} \mathbb{E} [\|\rho(t)\|_2]$  is bounded by (D.16) in Theorem D.5, while the term  $\epsilon_Q(T_{\text{actor}})$  in (D.17) can be upper bounded by Theorem 5.4. Finally with the parameters stated in Theorem 5.10, the following statement holds by straightforward calculation:

$$\min_{t \in [T_{\text{actor}}]} \mathbb{E} [J(\theta^*) - J(\theta(t))] \leq \mathcal{O} \left( |\mathcal{S}|^{1/2} B^{-1/2} + |\mathcal{S}| |\mathcal{A}|^{1/4} \left( \gamma^{k/8} + (T_{\text{actor}})^{-1/4} \right) \right).$$

□

**Q.E.D.**