

# Multi-Agent Off-Policy TDC with Near-Optimal Sample and Communication Complexities

Anonymous authors

Paper under double-blind review

## Abstract

The finite-time convergence of off-policy temporal difference (TD) learning has been comprehensively studied recently. However, such a type of convergence has not been established for off-policy TD learning in the multi-agent setting, which covers broader reinforcement learning applications and is fundamentally more challenging. This work develops a decentralized TD with correction (TDC) algorithm for multi-agent off-policy TD learning under Markovian sampling. In particular, our algorithm avoids sharing the actions, policies and rewards of the agents, and adopts mini-batch sampling to reduce the sampling variance and communication frequency. Under Markovian sampling and linear function approximation, we proved that the finite-time sample complexity of our algorithm for achieving an  $\epsilon$ -accurate solution is in the order of  $\mathcal{O}\left(\frac{M \ln \epsilon^{-1}}{\epsilon(1-\sigma_2)^2}\right)$ , where  $M$  denotes the total number of agents and  $\sigma_2$  is a network parameter. This matches the sample complexity of the centralized TDC. Moreover, our algorithm achieves the optimal communication complexity  $\mathcal{O}\left(\frac{\sqrt{M} \ln \epsilon^{-1}}{1-\sigma_2}\right)$  for synchronizing the value function parameters, which is order-wise lower than the communication complexity of the existing decentralized TD(0). Numerical simulations corroborate our theoretical findings.

## 1 Introduction

Multi-agent reinforcement learning (MARL) is an emerging technique that has broad applications in control Yanmaz et al. (2017); Chalaki & Malikopoulos (2020), wireless sensor networks Krishnamurthy et al. (2008); Yuan et al. (2020), robotics Yan et al. (2013), etc. In MARL, agents cooperatively interact with an environment and follow their own policies to collect local rewards. In particular, policy evaluation is a fundamental problem in MARL that aims to learn a multi-agent value function associated with the policies of the agents. This motivates the development of convergent and communication-efficient multi-agent TD learning algorithms.

For single-agent on-policy evaluation (i.e., samples are collected by target policy), the conventional TD(0) algorithm Sutton (1988); Sutton & Barto (2018) and Q-learning algorithm Dayan (1992) have been developed with asymptotic convergence guarantee. Recently, their finite-time (i.e., non-asymptotic) convergence has been established under Markovian sampling and linear approximation Bhandari et al. (2018); Zou et al. (2019). However, these algorithms may diverge in the off-policy setting Baird (1995), where samples are collected by a different behavior policy. To address this important issue, a family of gradient-based TD (GTD) algorithms were developed for off-policy evaluation with asymptotic convergence guarantee Sutton et al. (2008; 2009); Maei (2011). In particular, the TD with gradient correction (TDC) algorithm has been shown to have superior performance and its finite-time convergence has been established recently under Markovian sampling Xu et al. (2019); Gupta et al. (2019); Kaledin et al. (2020).

For multi-agent on-policy evaluation, various decentralized TD learning algorithms have been developed. For example, the finite-time convergence of decentralized TD(0) was established with i.i.d samples Wai et al. (2018); Doan et al. (2019) and Markovian samples Sun et al. (2020), respectively, under linear function approximation, and an improved result is further obtained in Wang et al. (2020) by leveraging gradient tracking. However, these algorithms do not apply to the off-policy setting. In the existing literature, decentralized off-policy TD learning has been studied only in simplified settings, for example, agents obtain

independent MDP trajectories Macua et al. (2014); Stanković & Stanković (2016); Cassano et al. (2020) or share their behavior and target policies with each other Cassano et al. (2020), and the data samples are either i.i.d. or have a finite sample size. These MARL settings either are impractical or reveal the agents' policies that may be sensitive. Therefore, we want to ask the following question:

- *Q1: Can we develop a decentralized off-policy TD algorithm for MARL with interdependent agents and avoids sharing local actions, policies and rewards of the agents?*

In fact, developing such a desired decentralized off-policy TD learning algorithm requires overcoming two major challenges. First, to perform decentralized off-policy TD learning, all the agents need to obtain a global importance sampling ratio (see Section 3.2). In Cassano et al. (2020), the authors obtained this ratio by sharing all local policies among the agents, which may lead to information leakage. Therefore, we aim to develop a safer scheme to synchronize the global importance sampling ratio among the agents without sharing any sensitive local information. Second, the existing decentralized TD-type algorithm achieves a communication complexity (number of communication rounds) of  $\mathcal{O}((\epsilon^{-1} + \frac{\sqrt{M}}{\sqrt{\epsilon(1-\sigma_2)}}) \ln \epsilon^{-1})$  for networks with  $M$  agents and parameter  $\sigma_2$  (See Assumption 5 in Section 4 for the definition of  $\sigma_2$ ) Sun et al. (2020). This induces much communication overhead when the target accuracy  $\epsilon$  is small. Hence, we want to ask the following theoretical question:

- *Q2: Can we develop a decentralized off-policy TD learning algorithm that achieves a near-optimal finite-time sample complexity and a near-optimal communication complexity under Markovian sampling?*

In this work, we provide affirmative answers to these questions by developing a decentralized TDC algorithm that avoids sharing the sensitive information of the agents and achieves the near-optimal sample complexity as well as a significantly reduced communication complexity. We summarize our contributions as follows.

## 1.1 Summary of Contribution

To perform multi-agent off-policy evaluation, we develop a decentralized TDC algorithm with linear function approximation. In every iteration, agents perform two-timescale TDC updates locally and exchange model parameters with their neighbors. In particular, our algorithm adopts the following designs to avoid sharing agents' local sensitive information and reduce communication load.

- We let the agents perform local averaging on their local importance sampling ratios to obtain approximated global importance sampling ratios.
- All the agents use a mini-batch of samples to update their model parameters in each iteration. The mini-batch sampling reduces the sampling variance and the communication frequency, leading to an improved communication complexity over that of the existing decentralized TD(0).
- After the decentralized TDC iterations, our algorithm performs additional local averaging steps to achieve a global consensus on the model parameters. This turns out to be critical for achieving the near-optimal complexity bounds.

Theoretically, we analyze the finite-time convergence of this decentralized TDC algorithm with Markovian samples and show that it attains a fast linear convergence. The overall sample complexity for achieving an  $\epsilon$ -accurate solution is in the order of  $\mathcal{O}(\frac{M \ln \epsilon^{-1}}{\epsilon(1-\sigma_2)^2})$ . When there is a single agent ( $M = 1$ ), this sample complexity result matches that of centralized TDC Xu et al. (2019) and matches the theoretical lower bound  $\mathcal{O}(\epsilon^{-1})$  Kaledin et al. (2020) up to a logarithm factor. In addition, the sample complexity is proportional to  $M$ , which matches the theoretical lower bound of decentralized strongly convex optimization Scaman et al. (2017). Moreover, the communication complexity of our algorithm for synchronizing value function parameters is in the order of  $\mathcal{O}(\frac{\sqrt{M \ln \epsilon^{-1}}}{1-\sigma_2})$ , which is significantly lower than the communication complexity  $\mathcal{O}(\epsilon^{-1} + \frac{\sqrt{M}}{\sqrt{\epsilon(1-\sigma_2)}} \ln \epsilon^{-1})$  of the decentralized TD(0) Sun et al. (2020) and matches the communication complexity lower bound Scaman et al. (2017).

Technically, our analysis is a nontrivial generalization of the analysis of centralized off-policy TDC to the decentralized case. In particular, our analysis establishes tight bounds of the consensus error caused by

synchronizing the global importance sampling ratio, especially under the Markovian sampling where the data samples are correlated. Moreover, we strategically bound the estimation error of the global importance sampling logarithm-ratio. Please refer to the proof sketch at the end of Section 4 for details.

## 1.2 Other Related Work

**Centralized policy evaluation.** TD(0) with linear function approximation Sutton (1988) is popular for on-policy evaluation. The asymptotic and non-asymptotic convergence results of TD(0) have been established in Sutton (1988); Dayan (1992); Jaakkola et al. (1993); Gordon (1995); Baird (1995); Tsitsiklis & Van Roy (1997); Tadić (2001); Hu & Syed (2019) and Korda & La (2015); Liu et al. (2015); Bhandari et al. (2018); Dalal et al. (2018b); Lakshminarayanan & Szepesvari (2018); Wang et al. (2019); Srikant & Ying (2019); Xu et al. (2020b) respectively. Sutton et al. (2009) proposed TDC for off-policy evaluation. The finite-sample convergence of TDC has been established in Dalal et al. (2018a; 2020) with i.i.d. samples and in Xu et al. (2019); Gupta et al. (2019); Kaledin et al. (2020) with Markovian samples.

**Decentralized policy evaluation.** Mathkar & Borkar (2016) proposed the decentralized TD(0) algorithm. The asymptotic and non-asymptotic convergence rate of decentralized TD have been obtained in Borkar (2009) and Sun et al. (2020); Wang et al. (2020) respectively. Existing decentralized off-policy evaluation studies considered simplified settings. Macua et al. (2014); Stanković & Stanković (2016) obtained asymptotic result for decentralized off-policy evaluation where the agents obtained independent MDPs. Cassano et al. (2020) obtained linear convergence rate also with independent MDPs by applying variance reduction and extended to the case where the individual behavior policies and the joint target policy are shared among the agents.

**Decentralized policy control.** Decentralized policy control is also an important MARL problem where the goal is to learn the optimal policy for each agent. Many algorithms have been proposed for decentralized policy control, including policy gradient Chen et al. (2021) and actor-critic Qu et al. (2020); Zhang et al. (2021).

## 2 Policy Evaluation in Multi-Agent RL

In this section, we introduce multi-agent reinforcement learning (MARL) and define the policy evaluation problem. Consider a fully decentralized multi-agent network that consists of  $M$  agents. The network topology is specified by an undirected graph  $\mathcal{G} = (\mathcal{M}, \mathcal{E})$ , where  $\mathcal{M} = \{1, 2, \dots, M\}$  denotes the set of agents and  $\mathcal{E}$  denotes the set of communication links. In MARL, the agents interact with a dynamic environment through a multi-agent Markov decision process (MMDP) specified as  $\{\mathcal{S}, \{\mathcal{A}^{(m)}\}_{m=1}^M, P, \{R^{(m)}\}_{m=1}^M, \gamma\}$ . To elaborate,  $\mathcal{S}$  denotes a global state space that is shared by all the agents,  $\mathcal{A}^{(m)}$  corresponds to the action space of agent  $m$ ,  $P$  is the state transition kernel and  $R^{(m)}$  denotes the reward function of agent  $m$ . All the state and action spaces have finite cardinality.  $\gamma \in (0, 1]$  is a discount factor.

At any time  $t$ , assume that all the agents are in the global state  $s_t \in \mathcal{S}$ . Then, each agent  $m$  takes a certain action  $a_t^{(m)} \in \mathcal{A}^{(m)}$  following its own stationary policy  $\pi^{(m)}$ , i.e.,  $a_t^{(m)} \sim \pi^{(m)}(\cdot | s_t)$ . After all the actions are taken, the global state transfers to a new state  $s_{t+1}$  according to the transition kernel  $P$ , i.e.,  $s_{t+1} \sim P(\cdot | s_t, a_t)$  where  $a_t := \{a_t^{(m)}\}_{m=1}^M$ . At the same time, each agent  $m$  receives a local reward  $R_t^{(m)} := R^{(m)}(s_t, a_t, s_{t+1})$  from the environment for this action-state transition. Throughout the MMDP, each agent  $m$  has access to only the global state  $\{s_t\}_t$  and its own actions  $\{a_t^{(m)}\}_t$  and rewards  $\{R_t^{(m)}\}_t$ . The goal of policy evaluation in MARL is to evaluate the following value function associated with all the local policies  $\pi := \{\pi^{(m)}\}_{m=1}^M$  for any global state  $s$ .

$$V^\pi(s) = \mathbb{E} \left[ \sum_{t=0}^{+\infty} \gamma^t \left( \frac{1}{M} \sum_{m=1}^M R_t^{(m)} \right) \middle| s_0 = s, \pi \right]. \quad (1)$$

A popular algorithm for policy evaluation in MARL is the decentralized TD(0) Sun et al. (2020). Specifically, consider a popular linear function approximation of the value function  $V_\theta(s) := \theta^\top \phi(s)$ , where  $\theta \in \mathbb{R}^d$  contains

the model parameters and  $\phi(s)$  is a feature vector that corresponds to the state  $s$ . In decentralized TD(0), each agent  $m$  collects a single Markovian sample  $\{s_t, a_t^{(m)}, s_{t+1}, R_t^{(m)}\}$  at time  $t$  and updates its own model parameters  $\theta_t^{(m)}$  with learning rate  $\alpha > 0$  as follows.

$$\theta_{t+1}^{(m)} = \sum_{m' \in \mathcal{N}_m} V_{m,m'} \theta_t^{(m')} + \alpha (A_t \theta_t^{(m)} + b_t^{(m)}), \quad (2)$$

where  $\mathcal{N}_m$  denotes the neighborhood of agent  $m$ ,  $V$  corresponds to a doubly stochastic communication matrix and  $A_t = \phi(s_t)(\gamma\phi(s_{t+1}) - \phi(s_t))^\top$ ,  $b_t^{(m)} = R_t^{(m)}\phi(s_t)$ . The above update rule applies the local TD error to update the parameters and synchronize the parameters among neighboring agents through the network.

### 3 Two-Timescale Decentralized TDC for Off-Policy Evaluation

#### 3.1 Centralized TDC

In this subsection, we review the centralized TD with gradient correction (TDC) algorithm Sutton et al. (2009). In RL, the agent may not have enough samples that are collected following the target policy  $\pi$ . Instead, it may have some data samples that are collected under a different behavior policy  $\pi_b$ . Therefore, in this *off-policy* setting, the agent would like to utilize the historical data to help evaluate the value function  $V^\pi$  associated with the target policy  $\pi$ .

In Sutton et al. (2009), a family of gradient-based TD (GTD) learning algorithms have been proposed for off-policy evaluation. In particular, the TDC algorithm has been shown to have superior performance. To explain, consider the linear approximation  $V_\theta(s) = \theta^\top \phi(s)$  and suppose the state space includes states  $s_1, \dots, s_n$ , we can define a total value function  $V_\theta := [V_\theta(s_1), \dots, V_\theta(s_n)]^\top$ . The goal of TDC is to minimize the following mean square projected Bellman error (MSPBE).

$$\text{MSPBE}(\theta) := \mathbb{E}_{\mu_b} \|V_\theta - \Pi T^\pi V_\theta\|^2,$$

where  $\mu_b$  is the stationary distribution induced by  $\pi_b$ ,  $T^\pi$  is the Bellman operator and  $\Pi$  is a projection operator onto the space of linear models. Given the  $i$ -th sample  $(s_i, a_i, s_{i+1}, R_i)$  obtained by the behavior policy, we define the following terms

$$\begin{aligned} \rho_i &:= \frac{\pi(a_i|s_i)}{\pi_b(a_i|s_i)}, \quad b_i := \rho_i R_i \phi(s_i), \quad A_i := \rho_i \phi(s_i)(\gamma\phi(s_{i+1}) - \phi(s_i))^\top, \\ B_i &:= -\gamma \rho_i \phi(s_{i+1})\phi(s_i)^\top, \quad C_i := -\phi(s_i)\phi(s_i)^\top, \end{aligned} \quad (3)$$

where  $\rho_t$  is referred to as the *importance sampling ratio*. Then, with learning rates  $\alpha, \beta > 0$  and initialization parameters  $\theta_0, w_0$ , the two timescale off-policy TDC algorithm takes the following recursive updates for  $t = 0, 1, 2, \dots$

$$\text{(TDC):} \quad \begin{cases} \theta_{t+1} = \theta_t + \alpha(A_t \theta_t + b_t + B_t w_t), \\ w_{t+1} = w_t + \beta(A_t \theta_t + b_t + C_t w_t). \end{cases} \quad (4)$$

#### 3.2 Decentralized Mini-batch TDC

In this subsection, we propose a decentralized TDC algorithm for off-policy evaluation in MARL. In the multi-agent setting, without loss of generality, we assume that each agent  $m$  has a target policy  $\pi^{(m)}$  and its samples are collected by a different behavior policy  $\pi_b^{(m)}$ . In particular, if agent  $m$  is on-policy, then we have  $\pi_b^{(m)} = \pi^{(m)}$ . In this *multi-agent off-policy* setting, the agents aim to utilize the data collected by the behavior policies  $\pi_b = \{\pi_b^{(m)}\}_{m=1}^M$  to help evaluate the value function  $V^\pi$  associated with the target policies  $\pi = \{\pi^{(m)}\}_{m=1}^M$ .

However, directly generalizing the centralized TDC algorithm to the decentralized setting will encounter several challenges. First, the centralized TDC in eq. (4) consumes one sample per-iteration and achieves the sample complexity  $O(\epsilon^{-1} \log \epsilon^{-1})$  Xu et al. (2019). Therefore, the corresponding decentralized TDC would

perform one local communication per-iteration and is expected to have a communication complexity in the order of  $O(\epsilon^{-1} \log \epsilon^{-1})$ , which induces *large communication overhead*. Second, in the multi-agent off-policy setting, every agent  $m$  has a local importance sampling ratio  $\rho_i^{(m)} := \pi^{(m)}(a_i^{(m)}|s_i)/\pi_b^{(m)}(a_i^{(m)}|s_i)$ . However, to correctly perform off-policy updates, every agent needs to know all the other agents' local importance sampling ratios in order to obtain the **global importance sampling ratio**  $\rho_i := \prod_{m=1}^M \rho_i^{(m)}$ . To address these challenges, we next propose a decentralized TDC algorithm that takes mini-batch stochastic updates.

To elaborate, note that  $\rho_i$  can be rewritten as

$$\rho_i = \exp \left( M \cdot \frac{1}{M} \sum_{m=1}^M \ln \rho_i^{(m)} \right).$$

Therefore, all the agents just need to obtain the average  $\frac{1}{M} \sum_{m=1}^M \ln \rho_i^{(m)}$ , which can be computed via local communication of the logarithm-ratios  $\{\ln \rho_i^{(m)}\}_{m=1}^M$  for  $L$  rounds. Specifically, every agent  $m$  initializes  $\tilde{\rho}_{i,0}^{(m)} = \ln \rho_i^{(m)}$  and for iterations  $\ell = 0, \dots, L-1$  do

$$\tilde{\rho}_{i,\ell+1}^{(m)} = \sum_{m' \in \mathcal{N}_m} V_{m,m'} \tilde{\rho}_{i,\ell}^{(m')}, \quad (5)$$

$$(\text{Output}) : \hat{\rho}_i^{(m)} = \exp(M \cdot \tilde{\rho}_{i,L}^{(m)}). \quad (6)$$

In Corollary 2 (see the appendix), we prove that all of these local estimates  $\{\hat{\rho}_i^{(m)}\}_{m=1}^M$  converge exponentially fast to the desired quantity  $\rho_i$  as  $L$  increases. Then, every agent  $m$  performs the following two-timescale TDC updates

$$\theta_{t+1}^{(m)} = \sum_{m' \in \mathcal{N}_m} V_{m,m'} \theta_t^{(m')} + \frac{\alpha}{N} \sum_{i=tN}^{(t+1)N-1} (A_i^{(m)} \theta_t^{(m)} + \tilde{b}_i^{(m)} + B_i^{(m)} w_t^{(m)}), \quad (7)$$

$$w_{t+1}^{(m)} = \sum_{m' \in \mathcal{N}_m} V_{m,m'} w_t^{(m')} + \frac{\beta}{N} \sum_{i=tN}^{(t+1)N-1} (A_i^{(m)} \theta_t^{(m)} + \tilde{b}_i^{(m)} + C_i w_t^{(m)}), \quad (8)$$

where  $A_i^{(m)}, B_i^{(m)}, \tilde{b}_i^{(m)}$  are defined by replacing the global variables  $\rho_i$  and  $R_i$  involved in  $A_i, B_i, b_i^{(m)}$  (see eq. (3)) with local variables  $\hat{\rho}_i^{(m)}$  and  $R_i^{(m)}$  respectively. To summarize, every TDC iteration of Algorithm 1 consumes  $N$  Markovian samples, and requires two vector communication rounds for synchronizing the parameter vectors  $\theta_t^{(m)}, w_t^{(m)}$ , and  $L$  scalar communication rounds for estimating the global importance sampling ratio ( $\{\rho_{i,\ell}^{(m)} : i = tN, \dots, (t+1)N-1, m \in \mathcal{M}\}$  are shared in the  $\ell$ -th communication round). We summarize these update rules in Algorithm 1. Moreover, after the decentralized TDC updates, the agents perform additional  $T'$  local averaging steps to reach a global consensus on the model parameters.

## 4 Finite-Time Analysis of Decentralized TDC

In this section, we analyze the finite-time convergence of Algorithm 1. Denote  $\mu_{\pi_b}$  as the stationary distribution of the Markov chain  $\{s_t\}_t$  induced by the collection of agents' behavioral policies  $\pi_b$ . Throughout the analysis, we define the following expected quantities.

$$A := \mathbb{E}_{\pi_b}[A_t], \quad B := \mathbb{E}_{\pi_b}[B_t], \quad C := \mathbb{E}_{\pi_b}[C_t], \quad b^{(m)} := \mathbb{E}_{\pi_b}[b_t^{(m)}], \quad \bar{b}_t := \frac{1}{M} \sum_{m=1}^M b_t^{(m)}, \quad \bar{b} := \mathbb{E}_{\pi_b}[\bar{b}_t],$$

where  $\mathbb{E}_{\pi_b}$  denotes the expectation when  $s_t \sim \mu_{\pi_b}$ ,  $a_t^{(m)} \sim \pi_b^{(m)}(s_t)$  and  $s_{t+1} \sim P(\cdot|s_t, a_t)$ . It is well-known that the optimal model parameter is  $\theta^* = -A^{-1}\bar{b}$  Xu et al. (2020b); Xu & Liang (2020). We make the following standard assumptions.

**Assumption 1.** *There exist constants  $\nu > 0$  and  $\delta \in (0, 1)$  such that for all  $t \geq 0$ ,*

$$\sup_{s \in \mathcal{S}} d_{TV}(\mathbb{P}_{\pi_b}(s_t | s_0 = s), \mu_{\pi_b}) \leq \nu \delta^t, \quad (9)$$

where  $d_{TV}$  denotes the total-variation distance.

**Algorithm 1** Decentralized mini-batch TDC.**Input:** Batch size  $N$ , iterations  $T, T'$ , learning rates  $\alpha, \beta$ .**Initialize:**  $\theta_0^{(m)}, w_0^{(m)}$  for all agents  $m \in \mathcal{M}$ .**for** iteration  $t = 0, 1, \dots, T - 1$  **do**    Each agent collects  $N$  Markovian samples and computes their local importance sampling ratios  $\rho_t^{(m)}$ .    **for** agent  $m \in \mathcal{M}$  *in parallel* **do**        Agent  $m$  estimates global importance sampling ratios for the  $N$  samples via eqs. (5) and (6) and performs the updates in eqs. (7) and (8).    **end****end****for** iteration  $t = T, T + 1, \dots, T + T' - 1$  **do**    **for** agent  $m \in \mathcal{M}$  *in parallel* **do**         $\theta_{t+1}^{(m)} = \sum_{m' \in \mathcal{N}_m} V_{m,m'} \theta_t^{(m')}$ .    **end****end****Output:**  $\{\theta_{T+T'}^{(m)}\}_{m=1}^M$ .**Assumption 2.** The matrices  $A$  and  $C$  are invertible.**Assumption 3.** The feature vectors satisfy  $\|\phi(s)\| \leq 1, \forall s$ .**Assumption 4.** There exist  $R_{\max}, \rho_{\max} > 0$  such that for all  $m \in \mathcal{M}$ :  $\max_{s,a,s'} R^{(m)}(s,a,s') < R_{\max}$  and  $\max_{s,a^{(m)}} \rho^{(m)}(s,a^{(m)}) < \rho_{\max}$ .**Assumption 5.** The communication matrix  $V$  is doubly stochastic, and its second largest singular value satisfies  $\sigma_2 \in [0, 1)$ .

Assumption 1 has been widely adopted in the existing literature Bhandari et al. (2018); Xu et al. (2019); Xu & Liang (2020); Shaocong et al. (2020; 2021). It holds for all homogeneous Markov chains with finite state-space and all uniformly ergodic Markov chains. Assumptions 2 – 4 are widely adopted in the analysis of TD learning algorithms Xu et al. (2019); Xu & Liang (2020). In particular, Assumption 2 implies that  $\lambda_1 := -\lambda_{\max}(A^\top C^{-1}A) > 0$ ,  $\lambda_2 := -\lambda_{\max}(C) > 0$ . Assumption 3 can always hold by normalizing the feature vectors. Assumption 4 holds for any uniformly lower bounded behavior policy, i.e.,  $\inf_{(s,a,m)} \pi_b^{(m)}(s,a) > 0$ , which ensures that every state-action pair  $(s,a)$  is visited infinitely often. Assumption 5 is standard in decentralized optimization Singh et al. (2020); Saha et al. (2020) and TD learning Sun et al. (2020); Wang et al. (2020).

We obtain the following finite-time error bound for Algorithm 1 with Markovian samples.

**Theorem 1.** Let Assumptions 1–5 hold. Run Algorithm 1 for  $T$  iterations with learning rates  $\alpha \leq \min\{\mathcal{O}(\beta), \mathcal{O}(\frac{1-\sigma_2}{\sqrt{M}})\}$ ,  $\beta \leq \mathcal{O}(1)$ , batch size  $N \geq \max\{\mathcal{O}(1), \mathcal{O}(\frac{\beta}{\alpha})\}$  and  $L \geq \mathcal{O}(\frac{\ln M + M \ln \rho_{\max}}{\ln \sigma_2^{-1}})$  (see eq. (35)-(38)). Then, we have

$$\mathbb{E}[\|\bar{\theta}_T - \theta^*\|^2] \leq \left(1 - \frac{\alpha \lambda_1}{6}\right)^T (\|\bar{\theta}_0 - \theta^*\|^2 + \|\bar{w}_0 - w_0^*\|^2) + \mathcal{O}\left(\frac{\beta}{N\alpha} + \frac{\beta \sigma_2^{L/4} 2^T}{M}\right). \quad (10)$$

Furthermore, after  $T'$  iterations of local averaging, the local models of all agents satisfy that: for all  $m = 1, \dots, M$ ,

$$\mathbb{E}[\|\theta_{T+T'}^{(m)} - \bar{\theta}_T\|^2] \leq \sigma_2^{2T'} \mathcal{O}\left(1 + \frac{M^4 \beta \alpha}{(1 - \sigma_2)^2} + \frac{M \beta \alpha \sigma_2^{L/4} 2^T}{1 - \sigma_2}\right). \quad (11)$$

Consequently, by choosing  $\alpha = \mathcal{O}(\frac{1-\sigma_2}{\sqrt{M}})$ ,  $\beta = \mathcal{O}(1)$ ,  $T = \mathcal{O}(\frac{\sqrt{M} \ln \epsilon^{-1}}{1-\sigma_2})$ ,  $N = \mathcal{O}(\frac{\sqrt{M}}{\epsilon(1-\sigma_2)})$ ,  $L = \mathcal{O}(\frac{\sqrt{M} \ln \epsilon^{-1}}{(1-\sigma_2)^2} + \frac{M}{1-\sigma_2})$ ,  $T' = \mathcal{O}(\frac{1}{1-\sigma_2} \ln \frac{M}{\epsilon(1-\sigma_2)})$ , we obtain that  $\mathbb{E}(\|\theta_{T+T'}^{(m)} - \theta^*\|^2) \leq \epsilon$  for all  $m$ . The overall communication complexity for synchronizing  $\theta_t^{(m)}$  is  $T + T' = \mathcal{O}(\frac{\sqrt{M} \ln \epsilon^{-1}}{1-\sigma_2})$ , and the total sample complexity is  $NT = \mathcal{O}(\frac{M \ln \epsilon^{-1}}{\epsilon(1-\sigma_2)^2})$ .

The above theorem shows that our decentralized TDC achieves the sample complexity  $\mathcal{O}\left(\frac{M \ln \epsilon^{-1}}{\epsilon(1-\sigma_2)^2}\right)$ , which, in the centralized setting ( $M = 1, \sigma_2 = 0$ ), matches that of centralized TDC for Markovian samples Xu et al. (2019) and matches the theoretical lower bound  $\mathcal{O}(\epsilon^{-1})$  given in Kaledin et al. (2020) up to a logarithm factor. In addition, the sample complexity is proportional to  $M$ , which matches the theoretical lower bound of decentralized strongly convex optimization in Scaman et al. (2017). Importantly, the communication complexity  $\mathcal{O}\left(\frac{\sqrt{M} \ln \epsilon^{-1}}{1-\sigma_2}\right)$  is substantially lower than the communication complexity  $\mathcal{O}\left((\epsilon^{-1} + \frac{\sqrt{M}}{\sqrt{\epsilon(1-\sigma_2)}}) \ln \epsilon^{-1}\right)$  of decentralized TD(0) Sun et al. (2020)<sup>1</sup>. Intuitively, this is because our algorithm adopts mini-batch sampling that significantly reduces the communication frequency, since communication occurs after collecting a mini-batch of samples to compute the mini-batch updates. Moreover, the communication complexity has a logarithm dependence  $\ln \epsilon^{-1}$  on the target accuracy, and this matches the theoretical lower bound of decentralized strongly convex optimization in Scaman et al. (2017).

Taking a deeper look, Theorem 1 shows that the average model  $\bar{\theta}_T$  converges to a small neighborhood of the optimal solution  $\theta^*$  at a fast linear convergence rate (10) that matches the convergence rate of centralized TDC Xu et al. (2019); Xu & Liang (2020). In particular, the convergence error is in the order of  $\mathcal{O}\left(\frac{\beta}{N\alpha} + \frac{\beta\sigma_2^{L/4}2^{2T}}{M}\right)$ , which can be driven arbitrarily close to zero by choosing a sufficiently large mini-batch size  $N$  and communication rounds  $L$  (with  $T$  fixed), and choosing constant-level learning rates  $\alpha, \beta$ . Moreover, the  $T'$  steps of extra local model averaging further help all the agents achieve a small consensus error at a linear convergence rate (11). Eqs. (10) and (11) together ensure the fast convergence of all the local model parameters. We want to point out that the  $T'$  local averaging steps are critical for establishing fast convergence of local model parameters. Specifically, without the  $T'$  local averaging steps, the consensus error  $\mathbb{E}[\|\theta_T^{(m)} - \bar{\theta}_T\|^2]$  would be in the order of at least  $\mathcal{O}(1)$ , which is constant-level and hence cannot guarantee the local model parameters converge arbitrarily close to the true solution.

**Proof sketch of Theorem 1.** The proof of the theorem is a nontrivial generalization of the analysis of centralized off-policy TDC to the decentralized case. Below, we sketch the technical proof and elaborate on the technical novelties.

- **Step 1.** We first consider an ideal case where the agents can access the exact global importance sampling ratio  $\rho_t$  at iteration  $t$ . In this ideal case, every agent  $m$  can replace the estimated global importance sampling ratio  $\hat{\rho}_i^{(m)}$  involved in  $A_i^{(m)}, B_i^{(m)}, \tilde{b}_i^{(m)}$  in the update rules (7) and (8) by the exact value  $\rho_i$ . Then, with the notations defined in Appendix A, the averaged update rules (14) and (15) become

$$\bar{\theta}_{t+1} = \bar{\theta}_t + \alpha(\bar{A}_t \bar{\theta}_t + \bar{\bar{b}}_t + \bar{B}_t \bar{w}_t), \quad (12)$$

$$\bar{w}_{t+1} = \bar{w}_t + \beta(\bar{A}_t \bar{\theta}_t + \bar{\bar{b}}_t + \bar{C}_t \bar{w}_t), \quad (13)$$

which can be seen as one step of centralized TDC. Hence, we can bound its optimization error terms  $\mathbb{E}[\|\bar{w}_{t+1} - w_t^*\|^2]$  and  $\mathbb{E}[\|\bar{\theta}_{t+1} - \theta^*\|^2]$ .

- **Step 2.** We return to Algorithm 1 and bound its optimization error terms  $\mathbb{E}[\|\bar{w}_{t+1} - w_t^*\|^2]$ ,  $\mathbb{E}[\|\bar{w}_{t+1} - w_{t+1}^*\|^2]$  and  $\mathbb{E}[\|\bar{\theta}_{t+1} - \theta^*\|^2]$ . This is done by bounding the gap between the centralized updates (12) and (13) (with exact  $\rho_t$ ) and the decentralized updates (14) and (15) (with inexact  $\rho_t$ ). The key is to establish Corollary 2, which strategically controls the gap between the inexact global importance sampling ratio  $\hat{\rho}_t^{(m)}$  and the exact value  $\rho_t$ .

To elaborate, note that the locally-averaged importance sampling ratios  $\hat{\rho}_{t,\ell}^{(m)}$  in eq. (5) exponentially converges to the value  $\ln \rho_t$ . However, the initial gap  $|\ln \rho_t^{(m)} - \ln \rho_t|$  can be numerically large since  $\rho_t^{(m)}$  may be a numerically small positive number. To avoid such divergence issue, our proof discusses two complementary cases. Case 1: the quantity  $\rho_{\min} := \min_{m \in \mathcal{M}} \rho_t^{(m)} \in [\sigma_2^{L/2}, \rho_{\max}]$ . In this case, the proof is straightforward as the initial gap is bounded. Case 2: the quantity  $\rho_{\min} \in (0, \sigma_2^{L/2}]$ . In this case, we show that the locally-averaged logarithm-ratio  $\hat{\rho}_{t,L}^{(m)}$  is below a large negative number  $\mathcal{O}\left(\frac{\ln \rho_{\min}}{M}\right) \ll 0$  (See eq. (63)). which implies that both the global importance sampling ratio  $\rho_t$  and its estimation  $\hat{\rho}_t^{(m)}$  are close to zero. In both cases,  $\{\hat{\rho}_t^{(m)}\}_{m=1}^M$  converge exponentially fast to  $\rho_t$  as  $L$  increases. This prove eq. (10).

<sup>1</sup>Sun et al. (2020) does not report communication complexity, so we calculated it based on their finite-time error bound.

- **Step 3.** Finally, we prove the consensus error (11). Although the consensus error exponentially decays during the  $T'$  extra local average steps in Algorithm 1, it is non-trivial to bound the initial consensus error  $\|\Delta\Theta_T\|_F$  of the  $T'$  local average iterations (see eq. (34)), which is caused by the  $T$  decentralized TDC steps. To bound this error, note that each decentralized TDC step consists of both local averaging and TDC update, which makes the consensus error  $\|\Delta\Theta_t\|_F$  diminishes geometrically fast with a noise term  $\sum_{m=1}^M \|h_m\|$  (see eq. (30)). Such a noise term is induced by the TDC update and hence its bound depends on both the consensus error and the model estimation error in eq. (10). We need to apply these correlated bounds iteratively for  $T$  iterations to bound the initial consensus error  $\|\Delta\Theta_T\|_F$ .

## 5 Experiments

### 5.1 Simulated Multi-Agent Networks

We simulate a multi-agent MDP with 10 decentralized agents. The shared state space contains 10 states and each agent can take 2 actions. All behavior policies are uniform policies (i.e., each agent takes all actions with equal probability), and the target policies are obtained by first perturbing the corresponding behavior policies with Gaussian noises sampled from  $\mathcal{N}(0, 0.05)$  and then performing a proper normalization. The entries of the transition kernel and the reward functions are independently generated from the uniform distribution on  $[0, 1]$  (with proper normalization for the transition kernel). We generate all state features with dimension 5 independently from the standard Gaussian distribution and normalize them to have unit norm. The discount factor is  $\gamma = 0.95$ .

We consider two types of network topologies: a fully connected network with communication matrix  $V$  having diagonal entries 0.8 and off-diagonal entries  $1/45$ , and a ring network with communication matrix  $V$  having diagonal entries 0.8 and entries 0.1 for adjacent agents. We implement and compare two algorithms in these networks: the decentralized TD(0) with batch size  $N = 1$  (used in Sun et al. (2020)) and our decentralized TDC with batch sizes  $N = 10, 20, 50, 100$ .

**Effect of Batch size:** We test these algorithms with varying batch size  $N$  and compare their sample and communication complexities. We set learning rate  $\alpha = 0.2$  for the decentralized TD(0) and  $\alpha = 0.2 * N$ ,  $\beta = 0.002 * N$  for our decentralized TDC with varying batch sizes  $N = 10, 20, 50, 100$ . Both algorithms use  $L = 3$  communication rounds for synchronizing  $\hat{\rho}_t^{(m)}$ . All algorithms are repeated 100 times using a fixed set of 100 MDP trajectories, each of which has 20k Markovian samples.

We first implement these algorithms in the fully connected network. Figure 1 plots the relative convergence error  $\|\bar{\theta}_t - \theta^*\|/\|\theta^*\|$  v.s. sample complexity ( $tN$ ) and communication complexity ( $t$ ). For each curve, its upper and lower envelopes denote the 95% and 5% percentiles of the 100 convergence errors, respectively. It can be seen that our decentralized TDC with different batch sizes achieve almost the same sample complexity as that of the decentralized TD(0), demonstrating the sample-efficiency of our algorithms. On the other hand, our decentralized TDC require much less communication complexities than the decentralized TD(0), and the required communication becomes lighter as batch size increases. All these results match our theoretical analysis well.

We further implement these algorithms in the ring network. The comparison results are exactly the same as those in Figure 1, since the update rule of  $\bar{\theta}_t$  does not rely on the network topology under exact global importance sampling.

**Effect of Communication Rounds:** We test our decentralized TDC using varying communication rounds  $L = 1, 3, 5, 7$ . We use a fixed batch size  $N = 100$  and set learning rates  $\alpha = 5$ ,  $\beta = 0.05$ , and repeat each algorithm 100 times using the set of 100 MDP trajectories. We also implement the decentralized TDC with exact global importance sampling ratios as a baseline. Figure 2 plots the relative convergence error v.s. communication complexity in the fully-connected network (Left) and ring network (Right). It can be seen that in both networks, the asymptotic convergence error of the decentralized TDC with inexact  $\rho$  decreases as the number of communication rounds  $L$  for synchronizing the global importance sampling ratio increases. In particular, with  $L = 1$ , decentralized TDC diverges asymptotically due to inaccurate estimation of the



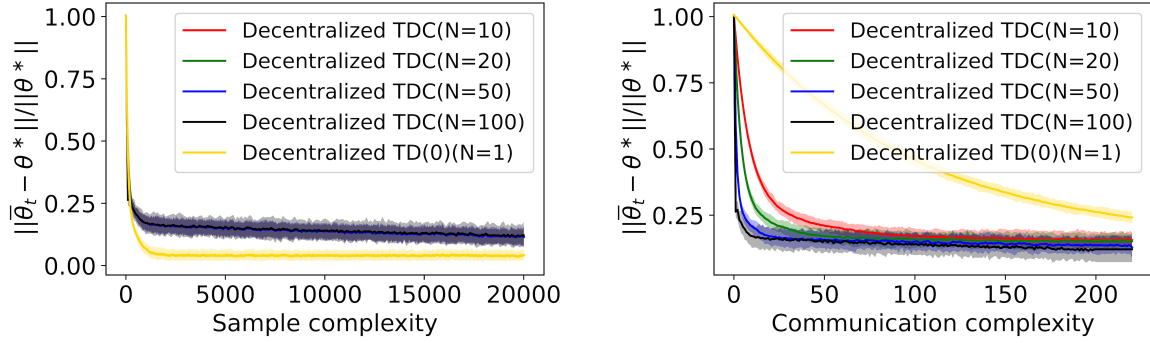


Figure 1: Comparison between decentralized TDC with varying batch sizes and decentralized TD(0).

global importance sampling ratio. As  $L$  increases to more than 5, the convergence error is as small as that under exact global importance sampling.

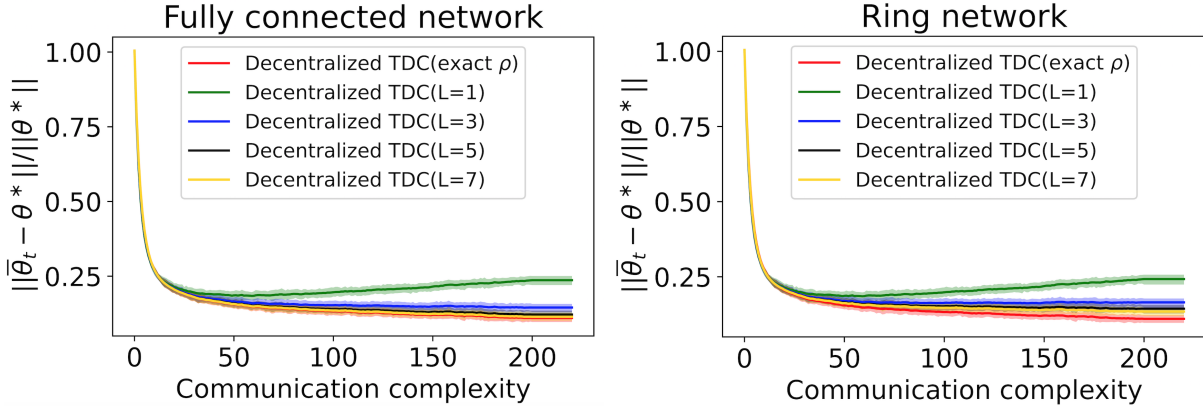


Figure 2: Effect of communication rounds  $L$  on asymptotic convergence error.

We further plot the maximum relative consensus error among all agents  $\max_m \|\theta_t^{(m)} - \bar{\theta}_t\| / \|\bar{\theta}^*\|$  v.s. communication complexity ( $t$ ) in the fully-connected network (Left) and ring network (Right) in Figure 3, where the tails in both figures correspond to the extra  $T' = 20$  local model averaging steps. In both networks, one can see that the consensus error decreases as  $L$  increases, and the extra local model averaging steps are necessary to achieve consensus. Moreover, it can be seen that the consensus errors achieved in the fully connected network are slightly smaller than those achieved in the ring network, as denser connections facilitate achieving the global consensus.

## 5.2 Two-Agent Cliff Navigation Problem

In this subsection, we test our algorithms in solving a two-agent Cliff Navigation problem Qiu et al. (2021) in a grid-world environment. This problem is adapted from its single-agent version (see Example 6.6 of Sutton & Barto (2018)). As illustrated in Figure 4, two agents start from the starting point “S” on a  $3 \times 4$  grid and aim to reach the destination “D”. Here, global state is defined as the joint location of the two agents, and there are in total  $(3 \times 4)^2 = 144$  global states. In most states, an agent can choose to move up, down, left or right by one step and receives  $-1$  reward. However, once an agent falls into the cliff “X”, it will return to the starting point “S” and receive  $-100$  reward. When an agent reaches “D”, it will always stay at “D”, and receives 0 reward if the other agent also reaches/stays at “D”, or receives  $-0.5$  reward otherwise. If an agent

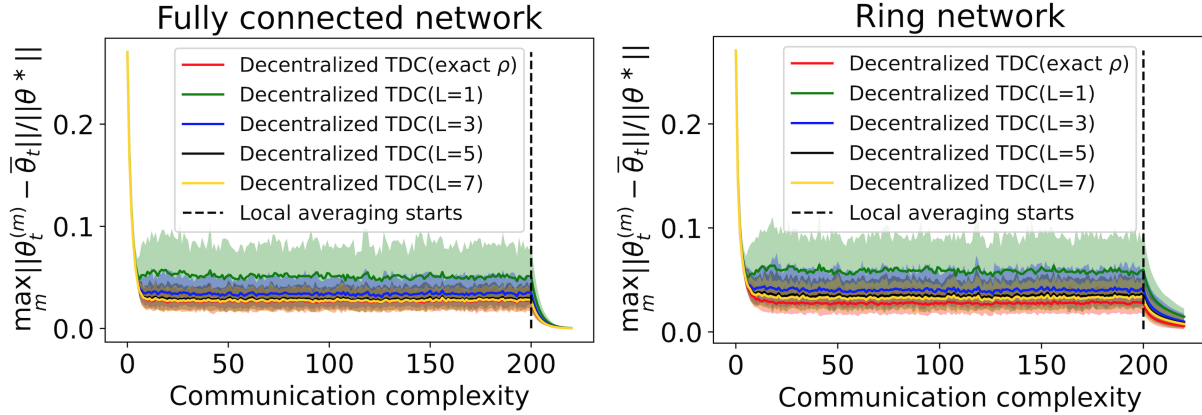
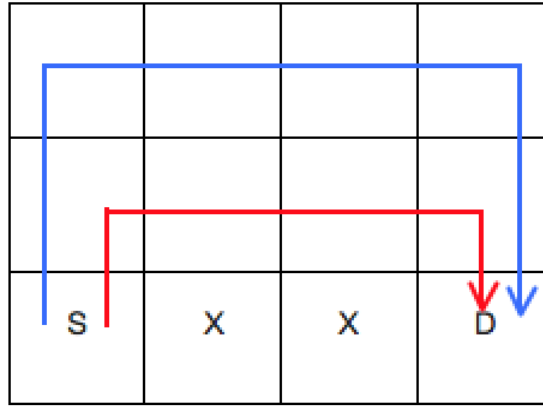
Figure 3: Effect of communication rounds  $L$  on consensus error.

Figure 4: Two-agent cliff navigation. (“S”, “X”, “D” denote starting point, cliff and destination respectively. The optimal path is shown in red.)

is not at “X” or “D” and selects a direction that points outside the grid, then it stays in the previous location and receives  $-1$  reward.

We apply the aforementioned algorithms with different batchsizes  $N$  to solve this problem. The hyperparameters and the ways to generate behavior policy and target policy are the same as the previous simulation experiment, except that  $L = 3$  and the communication matrix  $V$  has diagonal entries 0.7 and off-diagonals 0.3. All algorithms are repeated 100 times using a fixed set of 100 MDP trajectories, each of which has 20k Markovian samples. Figure 5 plots the relative convergence error  $\|\bar{\theta}_t - \theta^*\|/\|\theta^*\|$  v.s. sample complexity ( $tN$ ) and communication complexity ( $t$ ). We can see that compared with the decentralized TD(0), our decentralized TDC with different batch sizes achieve comparable sample complexities and much lower communication complexities. Moreover, the required communication becomes lighter as batch size increases. These properties are similar to those shown in the simulation and thus have generality.

## 6 Conclusion

In this paper, we develop a sample-efficient and communication-efficient decentralized TDC algorithm for multi-agent off-policy evaluation. Our algorithm synchronizes the local importance sampling ratios among the agents and adopts mini-batch stochastic updates to save communication. In particular, it avoids sharing agents’ sensitive local information. We prove that the proposed decentralized TDC algorithms achieve a

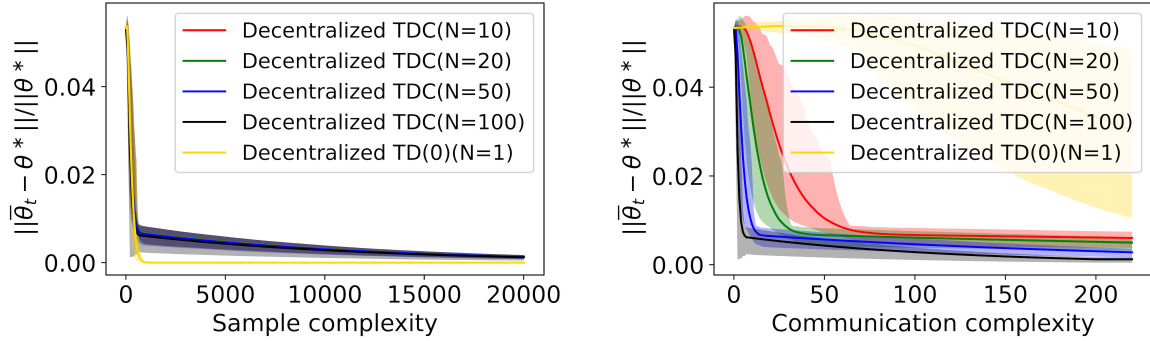


Figure 5: Comparison between decentralized TDC with varying batch sizes and decentralized TD(0).

near-optimal sample complexity as well as an optimal communication complexity that improves over the existing decentralized TD(0). In the future, we expect that our algorithm can serve as a fundamental component in the design of advanced policy optimization algorithms for MARL.

## References

- Leemon Baird. Residual algorithms: Reinforcement learning with function approximation. In *Proc. International Conference on Machine Learning (ICML)*, pp. 30–37, 1995.
- Jalaj Bhandari, Daniel Russo, and Raghav Singal. A finite time analysis of temporal difference learning with linear function approximation. In *Proc. Conference on Learning Theory (COLT)*, volume 75, pp. 1691–1692, 2018.
- Vivek S Borkar. *Stochastic approximation: a dynamical systems viewpoint*, volume 48. 2009.
- Lucas Cassano, Kun Yuan, and Ali H Sayed. Multi-agent fully decentralized value function learning with linear convergence rates. *IEEE Transactions on Automatic Control*, 2020.
- Behdad Chalaki and Andreas A Malikopoulos. A hysteretic q-learning coordination framework for emerging mobility systems in smart cities. *ArXiv:2011.03137*, 2020.
- Tianyi Chen, Kaiqing Zhang, Georgios B Giannakis, and Tamer Basar. Communication-efficient policy gradient methods for distributed reinforcement learning. *IEEE Transactions on Control of Network Systems*, 2021.
- Gal Dalal, Balazs Szorenyi, Gugan Thoppe, and Shie Mannor. Finite sample analysis of two-timescale stochastic approximation with applications to reinforcement learning. In *Proc. Conference on Learning Theory (COLT)*, 2018a.
- Gal Dalal, Balázs Szörényi, Gugan Thoppe, and Shie Mannor. Finite sample analyses for td (0) with function approximation. In *Proc. Association for the Advancement of Artificial Intelligence (AAAI)*, volume 32, 2018b.
- Gal Dalal, Balazs Szorenyi, and Gugan Thoppe. A tale of two-timescale reinforcement learning with the tightest finite-time bound. In *Proc. Association for the Advancement of Artificial Intelligence (AAAI)*, volume 34, pp. 3701–3708, 2020.
- Peter Dayan. The convergence of td ( $\lambda$ ) for general  $\lambda$ . *Machine learning*, 8(3-4):341–362, 1992.
- Thinh Doan, Siva Maguluri, and Justin Romberg. Finite-time analysis of distributed TD(0) with linear function approximation on multi-agent reinforcement learning. In *Proc. International Conference on Machine Learning (ICML)*, volume 97, pp. 1626–1635, 09–15 Jun 2019.

- Geoffrey J Gordon. Stable function approximation in dynamic programming. In *Machine Learning Proceedings 1995*, pp. 261–268, 1995.
- Harsh Gupta, R. Srikant, and Lei Ying. Finite-time performance bounds and adaptive learning rate selection for two time-scale reinforcement learning. In *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, volume 32, pp. 4704–4713, 2019.
- Bin Hu and Usman Ahmed Syed. Characterizing the exact behaviors of temporal difference learning algorithms using markov jump linear system theory. In *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, pp. 8479–8490, 2019.
- Tommi Jaakkola, Michael Jordan, and Satinder Singh. Convergence of stochastic iterative dynamic programming algorithms. In *Proc. Advances in Neural Information Processing Systems (NIPS)*, volume 6, pp. 703–710, 1993.
- Maxim Kaledin, Eric Moulines, Alexey Naumov, Vladislav Tadic, and Hoi-To Wai. Finite time analysis of linear two-timescale stochastic approximation with markovian noise. In *Proc. Conference on Learning Theory (COLT)*, pp. 2144–2203, 2020.
- Nathaniel Korda and Prashanth La. On td (0) with function approximation: Concentration bounds and a centered variant with exponential convergence. In *Proc. International Conference on Machine Learning (ICML)*, pp. 626–634, 2015.
- Vikram Krishnamurthy, Michael Maskery, and George Yin. Decentralized adaptive filtering algorithms for sensor activation in an unattended ground sensor network. *IEEE Transactions on Signal Processing*, 56(12):6086–6101, 2008.
- Chandrashekar Lakshminarayanan and Csaba Szepesvari. Linear stochastic approximation: How far does constant step-size and iterate averaging go? In *Proc. International Conference on Artificial Intelligence and Statistics (AISTATS)*, pp. 1347–1355, 2018.
- Bo Liu, Ji Liu, Mohammad Ghavamzadeh, Sridhar Mahadevan, and Marek Petrik. Finite-sample analysis of proximal gradient td algorithms. In *Proc. Conference on Uncertainty in Artificial Intelligence (UAI)*, pp. 504–513, 2015.
- Sergio Valcarcel Macua, Jianshu Chen, Santiago Zazo, and Ali H Sayed. Distributed policy evaluation under multiple behavior strategies. *IEEE Transactions on Automatic Control*, 60(5):1260–1274, 2014.
- Hamid Reza Maei. *Gradient temporal-difference learning algorithms*. PhD thesis, University of Alberta, 2011.
- Adwaitvedant Mathkar and Vivek S Borkar. Distributed reinforcement learning via gossip. *IEEE Transactions on Automatic Control*, 62(3):1465–1470, 2016.
- Wei Qiu, Xinrun Wang, Runsheng Yu, Rundong Wang, Xu He, Bo An, Svetlana Obraztsova, and Zinovi Rabinovich. Rmix: Learning risk-sensitive policies for cooperative reinforcement learning agents. In *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- Guannan Qu and Na Li. Harnessing smoothness to accelerate distributed optimization. *IEEE Transactions on Control of Network Systems*, 5(3):1245–1260, 2017.
- Guannan Qu, Adam Wierman, and Na Li. Scalable reinforcement learning of localized policies for multi-agent networked systems. In *Learning for Dynamics and Control (L4DC)*, pp. 256–266, 2020.
- Rajarshi Saha, Stefano Rini, Milind Rao, and Andrea Goldsmith. Decentralized optimization over noisy, rate-constrained networks: How to agree by talking about how we disagree. *ArXiv:2010.11292*, 2020.
- Kevin Scaman, Francis Bach, Sébastien Bubeck, Yin Tat Lee, and Laurent Massoulié. Optimal algorithms for smooth and strongly convex distributed optimization in networks. In *Proc. International Conference on Machine Learning (ICML)*, volume 70, pp. 3027–3036, 2017.

- Ma Shaocong, Zhou Yi, and Zou Shaofeng. Variance-reduced off-policy tdc learning: Non-asymptotic convergence analysis. In *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- Ma Shaocong, Chen Ziyi, Zhou Yi, and Zou Shaofeng. Greedy- $\{gq\}$  with variance reduction: Finite-time analysis and improved complexity. In *Proc. International Conference on Learning Representations (ICLR)*, 2021.
- Navjot Singh, Deepesh Data, Jemin George, and Suhas Diggavi. Squarm-sgd: Communication-efficient momentum sgd for decentralized optimization. *ArXiv:2005.07041*, 2020.
- Rayadurgam Srikant and Lei Ying. Finite-time error bounds for linear stochastic approximation and td learning. In *Proc. Conference on Learning Theory (COLT)*, pp. 2803–2830, 2019.
- Miloš S Stanković and Srdjan S Stanković. Multi-agent temporal-difference learning with linear function approximation: Weak convergence under time-varying network topologies. In *Proc. American Control Conference (ACC)*, pp. 167–172, 2016.
- Jun Sun, Gang Wang, Georgios B Giannakis, Qinmin Yang, and Zaiyue Yang. Finite-sample analysis of decentralized temporal-difference learning with linear function approximation. In *Proc. International Conference on Artificial Intelligence and Statistics (AISTATS)*, pp. 4485–4495, 2020.
- Richard S Sutton. Learning to predict by the methods of temporal differences. *Machine learning*, 3(1):9–44, 1988.
- Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. 2018.
- Richard S Sutton, Csaba Szepesvári, and Hamid Reza Maei. A convergent  $o(n)$  algorithm for off-policy temporal-difference learning with linear function approximation. In *Proc. Advances in Neural Information Processing Systems (NIPS)*, volume 21, pp. 1609–1616, 2008.
- Richard S Sutton, Hamid Reza Maei, Doina Precup, Shalabh Bhatnagar, David Silver, Csaba Szepesvári, and Eric Wiewiora. Fast gradient-descent methods for temporal-difference learning with linear function approximation. In *Proc. International Conference on Machine Learning (ICML)*, pp. 993–1000, 2009.
- Vladislav Tadić. On the convergence of temporal-difference learning with linear function approximation. *Machine learning*, 42(3):241–267, 2001.
- John N Tsitsiklis and Benjamin Van Roy. An analysis of temporal-difference learning with function approximation. *IEEE transactions on automatic control*, 42(5):674–690, 1997.
- Hoi-To Wai, Zhuoran Yang, Zhaoran Wang, and Mingyi Hong. Multi-agent reinforcement learning via double averaging primal-dual optimization. In *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, pp. 9672–9683, 2018.
- Gang Wang, Bingcong Li, and Georgios B Giannakis. A multistep lyapunov approach for finite-time analysis of biased stochastic approximation. *ArXiv:1909.04299*, 2019.
- Gang Wang, Songtao Lu, Georgios Giannakis, Gerald Tesauro, and Jian Sun. Decentralized td tracking with linear function approximation and its finite-time analysis. In *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, volume 33, 2020.
- Tengyu Xu and Yingbin Liang. Sample complexity bounds for two timescale value-based reinforcement learning algorithms. *ArXiv:2011.05053*, 2020.
- Tengyu Xu, Shaofeng Zou, and Yingbin Liang. Two time-scale off-policy td learning: Non-asymptotic analysis over markovian samples. In *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, pp. 10634–10644, 2019.
- Tengyu Xu, Zhe Wang, and Yingbin Liang. Improving sample complexity bounds for (natural) actor-critic algorithms. In *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, volume 33, 2020a.

Tengyu Xu, Zhe Wang, Yi Zhou, and Yingbin Liang. Reanalysis of variance reduced temporal difference learning. In *Proc. International Conference on Learning Representations (ICLR)*, 2020b.

Zhi Yan, Nicolas Jouandeau, and Arab Ali Cherif. A survey and analysis of multi-robot coordination. *International Journal of Advanced Robotic Systems*, 10(12):399, 2013.

Evşen Yanmaz, Markus Quaritsch, Saeed Yahyanejad, Bernhard Rinner, Hermann Hellwagner, and Christian Bettstetter. Communication and coordination for drone networks. In *Proc. International Conference on Ad Hoc Networks*, pp. 79–91, 2017.

Mingqi Yuan, Qi Cao, Man-on Pun, and Yi Chen. Towards user scheduling for 6g: A fairness-oriented scheduler using multi-agent reinforcement learning. *ArXiv:2012.15081*, 2020.

Junyu Zhang, Amrit Singh Bedi, Mengdi Wang, and Alec Koppel. Marl with general utilities via decentralized shadow reward actor-critic. *ArXiv:2106.00543*, 2021.

Shaofeng Zou, Tengyu Xu, and Yingbin Liang. Finite-sample analysis for sarsa with linear function approximation. In *Proc. Advances in Neural Information Processing Systems*, pp. 8665–8675, 2019.

## A Notations and Filtration

### A.1 Notations to rewrite update rules in Algorithm 1

We introduce the following additional notations.

$$\begin{aligned}\bar{\theta}_t &= \frac{1}{M} \sum_{m=1}^M \theta_t^{(m)}, \quad \bar{w}_t = \frac{1}{M} \sum_{m=1}^M w_t^{(m)}, \quad b_t^{(m)} = \rho_t R_t^{(m)} \phi(s_t), \quad \bar{b}_t = \frac{1}{M} \sum_{m=1}^M b_t^{(m)}, \\ \bar{A}_t &= \frac{1}{N} \sum_{i=tN}^{(t+1)N-1} A_i, \quad \bar{B}_t = \frac{1}{N} \sum_{i=tN}^{(t+1)N-1} B_i, \quad \bar{C}_t = \frac{1}{N} \sum_{i=tN}^{(t+1)N-1} C_i, \quad \bar{b}_t^{(m)} = \frac{1}{N} \sum_{i=tN}^{(t+1)N-1} b_i^{(m)}, \\ \bar{\bar{b}}_t &= \frac{1}{N} \sum_{i=tN}^{(t+1)N-1} \bar{b}_i = \frac{1}{M} \sum_{m=1}^M \bar{b}_t^{(m)}, \quad \bar{A}_t^{(m)} = \frac{1}{N} \sum_{i=tN}^{(t+1)N-1} A_i^{(m)}, \quad \bar{B}_t^{(m)} = \frac{1}{N} \sum_{i=tN}^{(t+1)N-1} B_i^{(m)}, \\ \bar{\bar{b}}_t^{(m)} &= \frac{1}{N} \sum_{i=tN}^{(t+1)N-1} \bar{\bar{b}}_i^{(m)}.\end{aligned}$$

Then, by averaging the update rules (7) and (8) over  $m$ , we obtain the following update rules of the model average  $\bar{\theta}_t, \bar{w}_t$ .

$$\bar{\theta}_{t+1} = \bar{\theta}_t + \frac{\alpha}{M} \sum_{m=1}^M (\bar{A}_t^{(m)} \theta_t^{(m)} + \bar{\bar{b}}_t^{(m)} + \bar{B}_t^{(m)} w_t^{(m)}), \quad (14)$$

$$\bar{w}_{t+1} = \bar{w}_t + \frac{\beta}{M} \sum_{m=1}^M (\bar{A}_t^{(m)} \theta_t^{(m)} + \bar{\bar{b}}_t^{(m)} + \bar{C}_t w_t^{(m)}). \quad (15)$$

### A.2 Filtration

Define the filtration  $\mathcal{F}_t = \sigma(\{s_{t'}, a_{t'}\}_{t'=1}^{tN-1} \cup \{s_{tN}\})$ . Then,

$$\bar{A}_t, \bar{B}_t, \bar{C}_t, \bar{b}_t^{(m)}, \bar{\bar{b}}_t, \bar{A}_t^{(m)}, \bar{B}_t^{(m)}, \bar{\bar{b}}_t^{(m)} \in \mathcal{F}_{t+1}/\mathcal{F}_t, \quad \theta_t^{(m)}, \bar{\theta}_t, w_t^{(m)}, \bar{w}_t, w_t^* \in \mathcal{F}_t/\mathcal{F}_{t-1}.$$

## B Proof of Theorem 1

We first bound the following tracking errors of  $\tilde{w}_{t+1}, \tilde{\theta}_{t+1}$  obtained by the centralized update rules (13) and (12) respectively.

$$\mathbb{E}[\|\tilde{w}_{t+1} - w_t^*\|^2 | \mathcal{F}_t]$$

$$= \|\bar{w}_t - w_t^*\|^2 + 2\beta \underbrace{(\bar{w}_t - w_t^*)^\top \mathbb{E}[\bar{A}_t \bar{\theta}_t + \bar{b}_t + \bar{C}_t \bar{w}_t | \mathcal{F}_t]}_{(I)} + \beta^2 \underbrace{\mathbb{E}[\|\bar{A}_t \bar{\theta}_t + \bar{b}_t + \bar{C}_t \bar{w}_t\|^2 | \mathcal{F}_t]}_{(II)}, \quad (16)$$

$$\begin{aligned} & \mathbb{E}[\|\tilde{\theta}_{t+1} - \theta^*\|^2 | \mathcal{F}_t] \\ &= \|\bar{\theta}_t - \theta^*\|^2 + 2\alpha \underbrace{(\bar{\theta}_t - \theta^*)^\top \mathbb{E}[\bar{A}_t \bar{\theta}_t + \bar{b}_t + \bar{B}_t \bar{w}_t | \mathcal{F}_t]}_{(III)} + \alpha^2 \underbrace{\mathbb{E}[\|\bar{A}_t \bar{\theta}_t + \bar{b}_t + \bar{B}_t \bar{w}_t\|^2 | \mathcal{F}_t]}_{(IV)}. \end{aligned} \quad (17)$$

The above four terms (I)-(IV) are respectively bounded below.

$$\begin{aligned} (I) &= (\bar{w}_t - w_t^*)^\top \mathbb{E}[\bar{A}_t \bar{\theta}_t + \bar{b}_t + \bar{C}_t \bar{w}_t | \mathcal{F}_t] \\ &\stackrel{(i)}{=} 2(\bar{w}_t - w_t^*)^\top \mathbb{E}[\bar{C}_t - C | \mathcal{F}_t] (\bar{w}_t - w_t^*) + 2(\bar{w}_t - w_t^*)^\top C (\bar{w}_t - w_t^*) + 2(\bar{w}_t - w_t^*)^\top \\ &\quad \mathbb{E}[(\bar{A}_t - \bar{C}_t C^{-1} A) \bar{\theta}_t + \bar{b}_t - \bar{C}_t C^{-1} b | \mathcal{F}_t] \\ &\stackrel{(ii)}{\leq} 2\|\mathbb{E}[\bar{C}_t - C | \mathcal{F}_t]\|_F \|\bar{w}_t - w_t^*\|^2 - 2\lambda_2 \|\bar{w}_t - w_t^*\|^2 + \lambda_2 \|\bar{w}_t - w_t^*\|^2 + \frac{3}{\lambda_2} \mathbb{E}[\|\bar{b}_t - \bar{C}_t C^{-1} b\|^2 | \mathcal{F}_t] \\ &\quad + \frac{3}{\lambda_2} \mathbb{E}[\|(\bar{A}_t - \bar{C}_t C^{-1} A)\|_F^2 | \mathcal{F}_t] \|\bar{\theta}_t - \theta^*\|^2 + \frac{3}{\lambda_2} \mathbb{E}[\|\bar{A}_t - \bar{C}_t C^{-1} A\|_F^2 | \mathcal{F}_t] \|\theta^*\|^2 \\ &\stackrel{(iii)}{\leq} \left( \frac{4\nu\rho_{\max}}{N(1-\delta)} - \lambda_2 \right) \|\bar{w}_t - w_t^*\|^2 + \frac{96\rho_{\max}^2(\nu+1)}{N\lambda_2(1-\delta)} \left(1 + \frac{1}{\lambda_2}\right)^2 (\|\bar{\theta}_t - \theta^*\|^2 + \|\theta^*\|^2 + R_{\max}^2) \end{aligned} \quad (18)$$

where (i) uses the notation that  $w_t^* = -C^{-1}(A\bar{\theta}_t + b)$ , (ii) uses  $\lambda_2 = -\lambda_{\max}(C)$  and the inequality that  $2a_1^\top a_2 \leq \sigma^{-1}\|a_1\|^2 + \sigma\|a_2\|^2$  for any  $a_1, a_2 \in \mathbb{R}^d$  and  $\sigma > 0$ , and applies Jensen's inequality to convex functions  $\|\cdot\|$  and  $\|\cdot\|^2$  and uses eq. (44), (iv) uses eqs. (47) and (48).

$$\begin{aligned} (II) &= \mathbb{E}[\|\bar{A}_t \bar{\theta}_t + \bar{b}_t + \bar{C}_t \bar{w}_t\|^2 | \mathcal{F}_t] \\ &\stackrel{(i)}{=} \mathbb{E}[\|(\bar{A}_t - \bar{C}_t C^{-1} A)(\bar{\theta}_t - \theta^*) + (\bar{A}_t - \bar{C}_t C^{-1} A)\theta^* + \bar{b}_t + \bar{C}_t(\bar{w}_t - w_t^*) - \bar{C}_t C^{-1} b\|^2 | \mathcal{F}_t] \\ &\stackrel{(ii)}{\leq} 4\mathbb{E}[\|\bar{A}_t - \bar{C}_t C^{-1} A\|_F^2 | \mathcal{F}_t] \|\bar{\theta}_t - \theta^*\|^2 + 4\|\bar{w}_t - w_t^*\|^2 \\ &\quad + 4\mathbb{E}[\|\bar{A}_t - \bar{C}_t C^{-1} A\|_F^2 | \mathcal{F}_t] \|\theta^*\|^2 + 4\mathbb{E}[\|\bar{b}_t - \bar{C}_t C^{-1} b\|^2 | \mathcal{F}_t] \\ &\stackrel{(iii)}{\leq} \frac{128\rho_{\max}^2(\nu+1)}{N(1-\delta)} \left(1 + \frac{1}{\lambda_2}\right)^2 (\|\bar{\theta}_t - \theta^*\|^2 + \|\theta^*\|^2 + R_{\max}^2) + 4\|\bar{w}_t - w_t^*\|^2 \end{aligned} \quad (19)$$

where (i) uses the notation that  $w_t^* = -C^{-1}(A\bar{\theta}_t + b)$ , (ii) uses  $\|\sum_{k=1}^4 a_k\|^2 \leq 4\sum_{k=1}^4 \|a_k\|^2$  for any  $a_1, a_2, a_3, a_4 \in \mathbb{R}^d$  and eq. (41), (iii) uses eqs. (47) and (48).

$$\begin{aligned} (III) &= (\bar{\theta}_t - \theta^*)^\top \mathbb{E}[\bar{A}_t \bar{\theta}_t + \bar{b}_t + \bar{B}_t \bar{w}_t | \mathcal{F}_t] \\ &\stackrel{(i)}{=} 2(\bar{\theta}_t - \theta^*)^\top \mathbb{E}[(\bar{A}_t - \bar{B}_t C^{-1} A - A^\top C^{-1} A)(\bar{\theta}_t - \theta^*) + A^\top C^{-1} A(\bar{\theta}_t - \theta^*) \\ &\quad + (\bar{A}_t - \bar{B}_t C^{-1} A - A^\top C^{-1} A)\theta^* + \bar{B}_t(\bar{w}_t - w_t^*) + \bar{b}_t - b - (\bar{B}_t - B)C^{-1} b | \mathcal{F}_t] \\ &\stackrel{(ii)}{\leq} 2\mathbb{E}[\|\bar{A}_t - \bar{B}_t C^{-1} A - A^\top C^{-1} A\| | \mathcal{F}_t] \|\bar{\theta}_t - \theta^*\|^2 - 2\lambda_1 \|\bar{\theta}_t - \theta^*\|^2 + \lambda_1 \|\bar{\theta}_t - \theta^*\|^2 + \|\bar{B}_t(\bar{w}_t - w_t^*)\|^2 \\ &\quad + \frac{4}{\lambda_1} \mathbb{E}[\|\bar{A}_t - \bar{B}_t C^{-1} A - A^\top C^{-1} A\|^2 | \mathcal{F}_t] \|\theta^*\|^2 + \|\bar{b}_t - b\|^2 + \|(\bar{B}_t - B)C^{-1} b\|^2 | \mathcal{F}_t] \\ &\stackrel{(iii)}{\leq} \left( \frac{64\rho_{\max}^2(\nu+1)}{N(1-\delta)} \left(1 + \frac{\rho_{\max}}{\lambda_2}\right)^2 - \lambda_1 \right) \|\bar{\theta}_t - \theta^*\|^2 + \frac{4\rho_{\max}^2}{\lambda_1} \|\bar{w}_t - w_t^*\|^2 \\ &\quad + \frac{32\rho_{\max}^2(\nu+1)}{N\lambda_1(1-\delta)} \left(4\left(1 + \frac{\rho_{\max}}{\lambda_2}\right)^2 \|\theta^*\|^2 + R_{\max}^2 + \frac{\rho_{\max} R_{\max}}{\lambda_2}\right) \end{aligned} \quad (20)$$

where (i) uses the notations that  $w_t^* = -C^{-1}(A\bar{\theta}_t + b)$  and that  $b = -A\theta^*$ , and the relation that  $C - B = A^\top$ , (ii) uses the notation that  $\lambda_1 = -\lambda_{\max}(A^\top C^{-1} A)$  and the inequality that  $2a_1^\top a_2 \leq \sigma^{-1}\|a_1\|^2 + \sigma\|a_2\|^2$  for

any  $a_1, a_2 \in \mathbb{R}^d$  and  $\sigma > 0$ , and applies Jensen's inequality to the convex functions  $\|\cdot\|$  and  $\|\cdot\|^2$ , (iii) uses eqs. (40), (42), (43), (45), (46) and (49).

$$\begin{aligned}
(IV) &= \mathbb{E}[\|\bar{A}_t \bar{\theta}_t + \bar{b}_t + \bar{B}_t \bar{w}_t\|^2 | \mathcal{F}_t] \\
&\stackrel{(i)}{=} \mathbb{E}[\|(\bar{A}_t - \bar{B}_t C^{-1} A)(\bar{\theta}_t - \theta^*) + \bar{b}_t - b + \bar{B}_t(\bar{w}_t - w_t^*) \\
&\quad + (\bar{A}_t - \bar{B}_t C^{-1} A - A^\top C^{-1} A)\theta^* - (\bar{B}_t - B)C^{-1}b\|^2 | \mathcal{F}_t] \\
&\stackrel{(ii)}{\leq} 10\mathbb{E}[\|\bar{A}_t\|_F^2 + \|\bar{B}_t C^{-1} A\|_F^2 | \mathcal{F}_t] \|\bar{\theta}_t - \theta^*\|^2 + 5\mathbb{E}[\|\bar{A}_t - \bar{B}_t C^{-1} A - A^\top C^{-1} A\|_F^2 | \mathcal{F}_t] \|\theta^*\|^2 \\
&\quad + 5\mathbb{E}[\|\bar{b}_t - b\|^2 | \mathcal{F}_t] + 5\mathbb{E}[\|\bar{B}_t\|_F^2 | \mathcal{F}_t] \|\bar{w}_t - w_t^*\|^2 + 5\mathbb{E}[\|\bar{B}_t - B\|_F^2 | \mathcal{F}_t] \|C^{-1}b\|^2 \\
&\stackrel{(iii)}{\leq} 40\rho_{\max}^2 \left(1 + \frac{\rho_{\max}}{\lambda_2}\right)^2 \|\bar{\theta}_t - \theta^*\|^2 + 5\rho_{\max}^2 \|\bar{w}_t - w_t^*\|^2 \\
&\quad + \frac{160\rho_{\max}^2(\nu+1)}{N(1-\delta)} \left(1 + \frac{\rho_{\max}}{\lambda_2}\right)^2 (\|\theta^*\|^2 + R_{\max}^2)
\end{aligned} \tag{21}$$

where (i) uses the notations that  $w_t^* = -C^{-1}(A\bar{\theta}_t + b)$  and that  $b = -A\theta^*$ , and the relation that  $C - B = A^\top$ , (ii) uses the inequality that  $\|\sum_{k=1}^K a_k\|^2 \leq K \sum_{k=1}^K \|a_k\|^2$  for any  $a_k \in \mathbb{R}^d$  and eqs. (40), (45), (46) and (49), (iii) uses eqs. (39), (40), (42) and (43) and  $(1 + \rho_{\max}^2/\lambda_2^2) \leq (1 + \rho_{\max}/\lambda_2)^2$ . Substituting the above terms (18)-(21) into eqs. (16) and (17) gives the following upper bounds of  $\mathbb{E}[\|\tilde{w}_{t+1} - w_t^*\|^2 | \mathcal{F}_t]$  and  $\mathbb{E}[\|\tilde{\theta}_{t+1} - \theta^*\|^2 | \mathcal{F}_t]$ .

$$\begin{aligned}
&\mathbb{E}[\|\tilde{w}_{t+1} - w_t^*\|^2 | \mathcal{F}_t] \\
&\leq \left[1 + 2\beta \left(\frac{4\nu\rho_{\max}}{N(1-\delta)} - \lambda_2\right) + 4\beta^2\right] \|\bar{w}_t - w_t^*\|^2 \\
&\quad + \frac{64\rho_{\max}^2(\nu+1)}{N(1-\delta)} \left(1 + \frac{1}{\lambda_2}\right)^2 \left(\frac{3\beta}{\lambda_2} + 2\beta^2\right) (\|\bar{\theta}_t - \theta^*\|^2 + \|\theta^*\|^2 + R_{\max}^2) \\
&\stackrel{(i)}{\leq} \left(1 - \frac{\beta\lambda_2}{2}\right) \|\bar{w}_t - w_t^*\|^2 + \frac{320\beta\rho_{\max}^2(\nu+1)}{N\lambda_2(1-\delta)} \left(1 + \frac{1}{\lambda_2}\right)^2 (\|\bar{\theta}_t - \theta^*\|^2 + \|\theta^*\|^2 + R_{\max}^2),
\end{aligned} \tag{22}$$

where (i) uses  $N \geq \frac{8\nu\rho_{\max}}{\lambda_2(1-\delta)}$  and  $\beta \leq \min(\frac{\lambda_2}{8}, \frac{1}{\lambda_2})$ .

$$\begin{aligned}
&\mathbb{E}[\|\tilde{\theta}_{t+1} - \theta^*\|^2 | \mathcal{F}_t] \\
&\leq \left(1 + \frac{128\alpha\rho_{\max}^2(\nu+1)}{N(1-\delta)} \left(1 + \frac{\rho_{\max}}{\lambda_2}\right)^2 - 2\alpha\lambda_1 + 40\alpha^2\rho_{\max}^2 \left(1 + \frac{\rho_{\max}}{\lambda_2}\right)^2\right) \|\bar{\theta}_t - \theta^*\|^2 \\
&\quad + \alpha\rho_{\max}^2 \left(\frac{8}{\lambda_1} + 5\alpha\right) \|\bar{w}_t - w_t^*\|^2 + \frac{64\alpha\rho_{\max}^2(\nu+1)}{N\lambda_1(1-\delta)} \left(4\left(1 + \frac{\rho_{\max}}{\lambda_2}\right)^2 \|\theta^*\|^2 + R_{\max}^2 + \frac{\rho_{\max}R_{\max}}{\lambda_2}\right) \\
&\quad + \frac{160\alpha^2\rho_{\max}^2(\nu+1)}{N(1-\delta)} \left(1 + \frac{\rho_{\max}}{\lambda_2}\right)^2 (\|\theta^*\|^2 + R_{\max}^2) \\
&\stackrel{(i)}{\leq} \left(1 - \frac{\alpha\lambda_1}{2}\right) \|\bar{\theta}_t - \theta^*\|^2 + \frac{13\alpha\rho_{\max}^2}{\lambda_1} \|\bar{w}_t - w_t^*\|^2 + \frac{224\alpha\rho_{\max}^2(\nu+1)}{N\lambda_1(1-\delta)} \left(1 + \frac{\rho_{\max}}{\lambda_2}\right)^2 (4\|\theta^*\|^2 + 2R_{\max}^2 + 1),
\end{aligned} \tag{23}$$

where (i) uses  $N \geq \frac{128\rho_{\max}^2(\nu+1)}{\lambda_1(1-\delta)} \left(1 + \frac{\rho_{\max}}{\lambda_2}\right)^2$ ,  $\alpha \leq \min[\frac{\lambda_1}{40\rho_{\max}^2} (1 + \frac{\rho_{\max}}{\lambda_2})^{-2}, \frac{1}{\lambda_1}]$  and  $\frac{\rho_{\max}R_{\max}}{\lambda_2} \leq (1 + \frac{\rho_{\max}}{\lambda_2})^2 + R_{\max}^2$ .

With the above upper bounds (22) and (23), we derive the upper bounds of  $\mathbb{E}[\|\bar{w}_{t+1} - w_t^*\|^2 | \mathcal{F}_t]$ ,  $\mathbb{E}[\|\bar{w}_{t+1} - w_{t+1}^*\|^2 | \mathcal{F}_t]$  and  $\mathbb{E}[\|\bar{\theta}_{t+1} - \theta^*\|^2 | \mathcal{F}_t]$  as follows.

$$\begin{aligned}
&\mathbb{E}[\|\bar{w}_{t+1} - w_t^*\|^2 | \mathcal{F}_t] \\
&\stackrel{(i)}{\leq} \left(1 + \frac{1}{6/(\beta\lambda_2) - 3}\right) \mathbb{E}[\|\tilde{w}_t - w_t^*\|^2 | \mathcal{F}_t] + \left(1 + \frac{6}{\beta\lambda_2} - 3\right) \mathbb{E}[\|\bar{w}_t - \tilde{w}_t\|^2 | \mathcal{F}_t]
\end{aligned}$$



$$\begin{aligned}
&\stackrel{(ii)}{\leq} \frac{6-2\beta\lambda_2}{6-3\beta\lambda_2} \left[ \left(1 - \frac{\beta\lambda_2}{2}\right) \|\bar{w}_t - w_t^*\|^2 + \frac{320\beta\rho_{\max}^2(\nu+1)}{N\lambda_2(1-\delta)} \left(1 + \frac{1}{\lambda_2}\right)^2 (\|\bar{\theta}_t - \theta^*\|^2 + \|\theta^*\|^2 + R_{\max}^2) \right] \\
&\quad + \frac{6}{\beta\lambda_2} \mathbb{E} \left[ \left\| \frac{\beta}{M} \sum_{m=1}^M [(\bar{A}_t^{(m)} - \bar{A}_t) \theta_t^{(m)} + \bar{b}_t^{(m)} - \bar{b}_t^{(m)}] \right\|^2 \middle| \mathcal{F}_t \right] \\
&\stackrel{(iii)}{\leq} \left(1 - \frac{\beta\lambda_2}{3}\right) \|\bar{w}_t - w_t^*\|^2 + \frac{4}{3} \frac{320\beta\rho_{\max}^2(\nu+1)}{N\lambda_2(1-\delta)} \left(1 + \frac{1}{\lambda_2}\right)^2 (\|\bar{\theta}_t - \theta^*\|^2 + \|\theta^*\|^2 + R_{\max}^2) \\
&\quad + \frac{6\beta}{\lambda_2} \mathbb{E} \left[ \frac{1}{M} \sum_{m=1}^M \left( \|(\bar{A}_t^{(m)} - \bar{A}_t) \theta_t^{(m)} + \bar{b}_t^{(m)} - \bar{b}_t^{(m)}\|^2 \right) \middle| \mathcal{F}_t \right] \\
&\stackrel{(iv)}{\leq} \left(1 - \frac{\beta\lambda_2}{3}\right) \|\bar{w}_t - w_t^*\|^2 + \frac{427\beta\rho_{\max}^2(\nu+1)}{N\lambda_2(1-\delta)} \left(1 + \frac{1}{\lambda_2}\right)^2 (\|\bar{\theta}_t - \theta^*\|^2 + \|\theta^*\|^2 + R_{\max}^2) \\
&\quad + \frac{12\beta\sigma_2^{L/4}}{M\lambda_2} \left( 16 \max_{m \in \mathcal{M}} (\|\theta_0^{(m)}\| + \|w_0^{(m)}\| + R_{\max})^2 [1 + \beta(2\rho_{\max} + 3)]^{2t} + R_{\max}^2 \right), \tag{24}
\end{aligned}$$

where (i) uses the inequality that  $\|a_1 + a_2\|^2 \leq (1 + \sigma)\|a_1\|^2 + (1 + \sigma^{-1})\|a_2\|^2$  for any  $a_1, a_2 \in \mathbb{R}^d$  and  $\sigma > 0$ , (ii) uses eqs. (15), (13) and (22), (iii) applies Jensen's inequality to the convex function  $\|\cdot\|^2$  and uses  $\beta \leq \frac{1}{\lambda_2}$ , (iv) uses the condition that  $\beta \leq \frac{1}{\lambda_2}$  which implies that  $1 + \frac{1}{6/(\beta\lambda_2)-3} \leq 2$ , the inequality that  $\|a_1 + a_2\|^2 \leq 2\|a_1\|^2 + 2\|a_2\|^2$  for any  $a_1, a_2 \in \mathbb{R}^d$ , and eqs. (50), (52) and (66).

$$\begin{aligned}
&\mathbb{E} [\|\bar{w}_{t+1} - w_{t+1}^*\|^2 | \mathcal{F}_t] \\
&\stackrel{(i)}{\leq} \left(1 + \frac{1}{2[3/(\beta\lambda_2) - 1]}\right) \mathbb{E} [\|\bar{w}_{t+1} - w_t^*\|^2 | \mathcal{F}_t] + [1 + 2(3/(\beta\lambda_2) - 1)] \mathbb{E} [\|w_{t+1}^* - w_t^*\|^2 | \mathcal{F}_t] \\
&\stackrel{(ii)}{\leq} \frac{6/(\beta\lambda_2) - 1}{2[3/(\beta\lambda_2) - 1]} \left[ \left(1 - \frac{\beta\lambda_2}{3}\right) \|\bar{w}_t - w_t^*\|^2 + \frac{427\beta\rho_{\max}^2(\nu+1)}{N\lambda_2(1-\delta)} \left(1 + \frac{1}{\lambda_2}\right)^2 (\|\bar{\theta}_t - \theta^*\|^2 + \|\theta^*\|^2 + R_{\max}^2) \right. \\
&\quad \left. + \frac{12\beta\sigma_2^{L/4}}{M\lambda_2} \left( 16 \max_{m \in \mathcal{M}} (\|\theta_0^{(m)}\| + \|w_0^{(m)}\| + R_{\max})^2 [1 + \beta(2\rho_{\max} + 3)]^{2t} + R_{\max}^2 \right) \right] \\
&\quad + \frac{6}{\beta\lambda_2} \mathbb{E} [\|C^{-1}A(\bar{\theta}_{t+1} - \bar{\theta}_t)\|^2 | \mathcal{F}_t] \\
&\stackrel{(iii)}{\leq} \left(1 - \frac{\beta\lambda_2}{6}\right) \|\bar{w}_t - w_t^*\|^2 + \frac{534\beta\rho_{\max}^2(\nu+1)}{N\lambda_2(1-\delta)} \left(1 + \frac{1}{\lambda_2}\right)^2 (\|\bar{\theta}_t - \theta^*\|^2 + \|\theta^*\|^2 + R_{\max}^2) \\
&\quad + \frac{15\beta\sigma_2^{L/4}}{M\lambda_2} (17 \max_{m \in \mathcal{M}} (\|\theta_0^{(m)}\| + \|w_0^{(m)}\| + R_{\max})^2 [1 + \beta(2\rho_{\max} + 3)]^{2t}) \\
&\quad + \underbrace{\frac{24\alpha^2\rho_{\max}^2}{\beta\lambda_2^3} \mathbb{E} \left[ \left\| \frac{1}{M} \sum_{m=1}^M (\bar{A}_t^{(m)} \theta_t^{(m)} + \bar{b}_t^{(m)} + \bar{B}_t^{(m)} w_t^{(m)}) \right\|^2 \middle| \mathcal{F}_t \right]}_{(V)}, \tag{25}
\end{aligned}$$

where (i) uses the inequality that  $\|a_1 + a_2\|^2 \leq (1 + \sigma)\|a_1\|^2 + (1 + \sigma^{-1})\|a_2\|^2$  for any  $a_1, a_2 \in \mathbb{R}^d$  and  $\sigma > 0$ , (ii) uses eq. (24) as well as the notation that  $w_t^* = -C^{-1}(A\bar{\theta}_t + b)$ , (iii) uses  $\beta \leq \frac{1}{\lambda_2}$  (this implies that  $\frac{6/(\beta\lambda_2)-1}{2[3/(\beta\lambda_2)-1]} \leq \frac{5}{4}$ ) and eqs. (14), (39) and (43). The above term (V) can be upper bounded as follows.

$$\begin{aligned}
(V) &= \mathbb{E} \left[ \left\| \frac{1}{M} \sum_{m=1}^M (\bar{A}_t^{(m)} \theta_t^{(m)} + \bar{b}_t^{(m)} + \bar{B}_t^{(m)} w_t^{(m)}) \right\|^2 \middle| \mathcal{F}_t \right] \\
&\stackrel{(i)}{\leq} 2\mathbb{E} \left[ \left\| \frac{1}{M} \sum_{m=1}^M ((\bar{A}_t^{(m)} - \bar{A}_t) \theta_t^{(m)} + (\bar{b}_t^{(m)} - \bar{b}_t^{(m)}) + (\bar{B}_t^{(m)} - \bar{B}_t) w_t^{(m)}) \right\|^2 \middle| \mathcal{F}_t \right] \\
&\quad + 2\mathbb{E} \left[ \left\| \frac{1}{M} \sum_{m=1}^M (\bar{A}_t \theta_t^{(m)} + \bar{b}_t^{(m)} + \bar{B}_t w_t^{(m)}) \right\|^2 \middle| \mathcal{F}_t \right]
\end{aligned}$$

$$\begin{aligned}
&\stackrel{(ii)}{\leq} \frac{2}{M} \sum_{m=1}^M \mathbb{E}[\|(\bar{A}_t^{(m)} - \bar{A}_t)\theta_t^{(m)} + (\bar{b}_t^{(m)} - \bar{b}_t) + (\bar{B}_t^{(m)} - \bar{B}_t)w_t^{(m)}\|^2 | \mathcal{F}_t] + 2\mathbb{E}(\|\bar{A}_t\bar{\theta}_t + \bar{b}_t + \bar{B}_t\bar{w}_t\|^2 | \mathcal{F}_t) \\
&\stackrel{(iii)}{\leq} \frac{6}{M} \sum_{m=1}^M \mathbb{E}[\|\bar{A}_t^{(m)} - \bar{A}_t\|^2 \|\theta_t^{(m)}\|^2 + \|\bar{b}_t^{(m)} - \bar{b}_t\|^2 + \|\bar{B}_t^{(m)} - \bar{B}_t\|^2 \|w_t^{(m)}\|^2 | \mathcal{F}_t] \\
&\quad + 2\left(40\rho_{\max}^2\left(1 + \frac{\rho_{\max}}{\lambda_2}\right)^2 \|\bar{\theta}_t - \theta^*\|^2 + 5\rho_{\max}^2 \|\bar{w}_t - w_t^*\|^2 + \frac{160\rho_{\max}^2(\nu+1)}{N(1-\delta)}\left(1 + \frac{\rho_{\max}}{\lambda_2}\right)^2 (\|\theta^*\|^2 + R_{\max}^2)\right) \\
&\stackrel{(iv)}{\leq} \frac{6\sigma_2^{L/4}}{M} \left(20 \max_{m \in \mathcal{M}} (\|\theta_0^{(m)}\| + \|w_0^{(m)}\| + R_{\max})^2 [1 + \beta(2\rho_{\max} + 3)]^{2t} + R_{\max}^2\right) \\
&\quad + 2\left(40\rho_{\max}^2\left(1 + \frac{\rho_{\max}}{\lambda_2}\right)^2 \|\bar{\theta}_t - \theta^*\|^2 + 5\rho_{\max}^2 \|\bar{w}_t - w_t^*\|^2 + \frac{160\rho_{\max}^2(\nu+1)}{N(1-\delta)}\left(1 + \frac{\rho_{\max}}{\lambda_2}\right)^2 (\|\theta^*\|^2 + R_{\max}^2)\right)
\end{aligned}$$

where (i) uses the inequality that  $\|a_1 + a_2\|^2 \leq 2\|a_1\|^2 + 2\|a_2\|^2$  for any  $a_1, a_2 \in \mathbb{R}^d$ , (ii) applies Jensen's inequality to the convex function  $\|\cdot\|^2$ , (iii) uses eq. (21) and the inequality that  $\|a_1 + a_2 + a_3\|^2 \leq 3\sum_{k=1}^3 \|a_k\|^2$  for any  $a_1, a_2, a_3 \in \mathbb{R}^d$ , (iv) uses eqs. (50), (51), (52) and (66). Substituting the above inequality into eq. (25) yields that

$$\begin{aligned}
&\mathbb{E}[\|\bar{w}_{t+1} - w_{t+1}^*\|^2 | \mathcal{F}_t] \\
&\leq \left(1 - \frac{\beta\lambda_2}{6} + \frac{120\alpha^2\rho_{\max}^4}{\beta\lambda_2^3}\right) \|\bar{w}_t - w_t^*\|^2 + \left[\frac{534\beta\rho_{\max}^2(\nu+1)}{N\lambda_2(1-\delta)}\left(1 + \frac{1}{\lambda_2}\right)^2 + \frac{1920\alpha^2\rho_{\max}^4}{\beta\lambda_2^3}\left(1 + \frac{\rho_{\max}}{\lambda_2}\right)^2\right] \|\bar{\theta}_t - \theta^*\|^2 \\
&\quad + \frac{3840\rho_{\max}^2(\nu+1)}{N(1-\delta)} \left[\frac{\beta}{\lambda_2}\left(1 + \frac{1}{\lambda_2}\right)^2 + \frac{\alpha^2\rho_{\max}^2}{\beta\lambda_2^3}\left(1 + \frac{\rho_{\max}}{\lambda_2}\right)^2\right] (\|\theta^*\|^2 + R_{\max}^2) \\
&\quad + \frac{255\sigma_2^{L/4}}{M\lambda_2} \left(\beta + \frac{12\alpha^2\rho_{\max}^2}{\beta\lambda_2^2}\right) \max_{m \in \mathcal{M}} (\|\theta_0^{(m)}\| + \|w_0^{(m)}\| + R_{\max})^2 [1 + \beta(2\rho_{\max} + 3)]^{2t}. \tag{26}
\end{aligned}$$

$$\begin{aligned}
&\mathbb{E}[\|\bar{\theta}_{t+1} - \theta^*\|^2 | \mathcal{F}_t] \\
&\stackrel{(i)}{\leq} \left(1 + \frac{1}{6/(\alpha\lambda_1) - 3}\right) \mathbb{E}[\tilde{\theta}_{t+1} - \theta^* | \mathcal{F}_t] + \left(1 + \frac{6}{\alpha\lambda_1} - 3\right) \mathbb{E}[\bar{\theta}_{t+1} - \tilde{\theta}_{t+1} | \mathcal{F}_t] \\
&\stackrel{(ii)}{\leq} \frac{6-2\alpha\lambda_1}{6-3\alpha\lambda_1} \left[\left(1 - \frac{\alpha\lambda_1}{2}\right) \|\bar{\theta}_t - \theta^*\|^2 + \frac{13\alpha\rho_{\max}^2}{\lambda_1} \|\bar{w}_t - w_t^*\|^2\right. \\
&\quad \left.+ \frac{224\alpha\rho_{\max}^2(\nu+1)}{N\lambda_1(1-\delta)} (4\|\theta^*\|^2 + 2R_{\max}^2 + 1) \left(1 + \frac{\rho_{\max}}{\lambda_2}\right)^2\right] \\
&\quad + \frac{6}{\alpha\lambda_1} \mathbb{E}\left[\left\|\frac{\alpha}{M} \sum_{m=1}^M [(\bar{A}_t^{(m)} - \bar{A}_t)\theta_t^{(m)} + \bar{b}_t^{(m)} - \bar{b}_t + (\bar{B}_t^{(m)} - \bar{B}_t)w_t^{(m)}]\right\|^2 | \mathcal{F}_t\right] \\
&\stackrel{(iii)}{\leq} \left(1 - \frac{\alpha\lambda_1}{3}\right) \|\bar{\theta}_t - \theta^*\|^2 + \frac{18\alpha\rho_{\max}^2}{\lambda_1} \|\bar{w}_t - w_t^*\|^2 + \frac{300\alpha\rho_{\max}^2(\nu+1)}{N\lambda_1(1-\delta)} (4\|\theta^*\|^2 + 2R_{\max}^2 + 1) \left(1 + \frac{\rho_{\max}}{\lambda_2}\right)^2 \\
&\quad + \frac{6\alpha}{M\lambda_1} \sum_{m=1}^M \mathbb{E}[\|(\bar{A}_t^{(m)} - \bar{A}_t)\theta_t^{(m)} + \bar{b}_t^{(m)} - \bar{b}_t + (\bar{B}_t^{(m)} - \bar{B}_t)w_t^{(m)}\|^2 | \mathcal{F}_t] \\
&\stackrel{(iv)}{\leq} \left(1 - \frac{\alpha\lambda_1}{3}\right) \|\bar{\theta}_t - \theta^*\|^2 + \frac{18\alpha\rho_{\max}^2}{\lambda_1} \|\bar{w}_t - w_t^*\|^2 + \frac{300\alpha\rho_{\max}^2(\nu+1)}{N\lambda_1(1-\delta)} (4\|\theta^*\|^2 + 2R_{\max}^2 + 1) \left(1 + \frac{\rho_{\max}}{\lambda_2}\right)^2 \\
&\quad + \frac{18\alpha\sigma_2^{L/4}}{M\lambda_1} \left[20 \max_{m \in \mathcal{M}} (\|\theta_0^{(m)}\| + \|w_0^{(m)}\| + R_{\max})^2 [1 + \beta(2\rho_{\max} + 3)]^{2t} + R_{\max}^2\right], \tag{27}
\end{aligned}$$

where (i) uses the inequality that  $\|a_1 + a_2\|^2 \leq (1 + \sigma)\|a_1\|^2 + (1 + \sigma^{-1})\|a_2\|^2$  for any  $a_1, a_2 \in \mathbb{R}^d$  and  $\sigma > 0$ , (ii) uses eqs. (14), (12) and (23), (iii) uses  $\alpha \leq 1/\lambda_1$  and applies Jensen's inequality to the convex function  $\|\cdot\|^2$ , and (iv) uses the inequality that  $\|a_1 + a_2 + a_3\|^2 \leq 3\sum_{k=1}^3 \|a_k\|^2$  for any  $a_1, a_2, a_3 \in \mathbb{R}^d$  and then uses eqs. (50)-(52).

Taking expectation on both sides of eqs. (26) and (27) and summing up the two inequalities yields that

$$\begin{aligned}
& \mathbb{E}(\|\bar{\theta}_{t+1} - \theta^*\|^2) + \mathbb{E}(\|\bar{w}_{t+1} - w_t^*\|^2) \\
& \stackrel{(i)}{\leq} \left(1 - \frac{\beta\lambda_2}{6} + \frac{120\alpha^2\rho_{\max}^4}{\beta\lambda_2^3} + \frac{18\alpha\rho_{\max}^2}{\lambda_1}\right) \mathbb{E}(\|\bar{w}_t - w_t^*\|^2) \\
& \quad + \left[1 - \frac{\alpha\lambda_1}{3} + \frac{534\beta\rho_{\max}^2(\nu+1)}{N\lambda_2(1-\delta)}\left(1 + \frac{1}{\lambda_2}\right)^2 + \frac{1920\alpha^2\rho_{\max}^4}{\beta\lambda_2^3}\left(1 + \frac{\rho_{\max}}{\lambda_2}\right)^2\right] \mathbb{E}(\|\bar{\theta}_t - \theta^*\|^2) \\
& \quad + \frac{3840\rho_{\max}^2(\nu+1)}{N(1-\delta)} \left[\frac{\beta}{\lambda_2}\left(1 + \frac{1}{\lambda_2}\right)^2 + \frac{\alpha^2\rho_{\max}^2}{\beta\lambda_2^3}\left(1 + \frac{\rho_{\max}}{\lambda_2}\right)^2 + \frac{\alpha}{\lambda_1}\right] (\|\theta^*\|^2 + R_{\max}^2 + 1) \\
& \quad + \frac{255\sigma_2^{L/4}}{M\lambda_2} \left(\beta + \frac{12\alpha^2\rho_{\max}^2}{\beta\lambda_2^2} + \frac{\alpha\lambda_2}{\lambda_1}\right) \max_{m \in \mathcal{M}} (\|\theta_0^{(m)}\| + \|w_0^{(m)}\| + R_{\max})^2 [1 + \beta(2\rho_{\max} + 3)]^{2t} \\
& \stackrel{(i)}{\leq} \left(1 - \frac{\alpha\lambda_1}{6}\right) [\mathbb{E}(\|\bar{\theta}_t - \theta^*\|^2) + \mathbb{E}(\|\bar{w}_t - w_t^*\|^2)] + \frac{6000\beta\rho_{\max}^2(\nu+1)}{N\lambda_2(1-\delta)} \left(1 + \frac{\rho_{\max}}{\lambda_2}\right)^2 (\|\theta^*\|^2 + R_{\max}^2 + 1) \\
& \quad + \frac{574\beta\sigma_2^{L/4}}{M\lambda_2} \max_{m \in \mathcal{M}} (\|\theta_0^{(m)}\| + \|w_0^{(m)}\| + R_{\max})^2 [1 + \beta(2\rho_{\max} + 3)]^{2t},
\end{aligned}$$

where (i) uses the conditions that  $N \geq \frac{6408\beta\rho_{\max}^2(\nu+1)}{\alpha\lambda_1\lambda_2(1-\delta)} \left(1 + \frac{1}{\lambda_2}\right)^2$ ,

$\alpha \leq \min\left(\frac{\beta\lambda_1\lambda_2^3}{23040\rho_{\max}^4}, \frac{\beta\lambda_2^2}{53\rho_{\max}^2}, \frac{\beta\lambda_1\lambda_2}{432\rho_{\max}^2}, \frac{\beta\lambda_2}{2\lambda_1}, \frac{\beta\lambda_1}{2\lambda_2}, \frac{\beta\lambda_2}{4\rho_{\max}}\right)$ . Iterating the inequality above yields that

$$\begin{aligned}
& \mathbb{E}(\|\bar{\theta}_T - \theta^*\|^2 + \|\bar{w}_T - w^*\|^2) \\
& \leq \left(1 - \frac{\alpha\lambda_1}{6}\right)^T (\|\bar{\theta}_0 - \theta^*\|^2 + \|\bar{w}_0 - w_0^*\|^2) \\
& \quad + \sum_{t=0}^{T-1} \left(1 - \frac{\alpha\lambda_1}{6}\right)^{T-1-t} \left[\frac{6000\beta\rho_{\max}^2(\nu+1)}{N\lambda_2(1-\delta)} \left(1 + \frac{\rho_{\max}}{\lambda_2}\right)^2 (\|\theta^*\|^2 + R_{\max}^2 + 1) \right. \\
& \quad \left. + \frac{574\beta\sigma_2^{L/4}}{M\lambda_2} \max_{m \in \mathcal{M}} (\|\theta_0^{(m)}\| + \|w_0^{(m)}\| + R_{\max})^2 [1 + \beta(2\rho_{\max} + 3)]^{2t}\right] \\
& \stackrel{(i)}{\leq} \left(1 - \frac{\alpha\lambda_1}{6}\right)^T (\|\bar{\theta}_0 - \theta^*\|^2 + \|\bar{w}_0 - w_0^*\|^2) + \frac{36000\beta\rho_{\max}^2(\nu+1)}{\alpha N\lambda_1\lambda_2(1-\delta)} \left(1 + \frac{\rho_{\max}}{\lambda_2}\right)^2 (\|\theta^*\|^2 + R_{\max}^2 + 1) \\
& \quad + \frac{574\beta\sigma_2^{L/4}2^T}{M\lambda_2} \max_{m \in \mathcal{M}} (\|\theta_0^{(m)}\| + \|w_0^{(m)}\| + R_{\max})^2 \\
& = \left(1 - \frac{\alpha\lambda_1}{6}\right)^T (\|\bar{\theta}_0 - \theta^*\|^2 + \|\bar{w}_0 - w_0^*\|^2) + \mathcal{O}\left(\frac{\beta}{N\alpha} + \frac{\beta\sigma_2^{L/4}2^T}{M}\right), \tag{28}
\end{aligned}$$

where (i) uses the definitions of  $c_{10}$ ,  $c_{11}$  in Appendix A and the conditions that  $\alpha \leq \frac{1}{\lambda_1}$  and  $\beta \leq \frac{2}{5(2\rho_{\max}+3)}$  which respectively imply that  $1 - \frac{\alpha\lambda_1}{6} \geq \frac{5}{6}$  and that  $1 + \beta(2\rho_{\max} + 3) \leq \sqrt{2}$ . This proves eq. (10).

Next, we prove eq. (11). Note that the local model averaging iterations can be rewritten into the matrix-vector form as  $\Theta_{t+1} = V\Theta_t$  where  $T \leq t \leq T + T'$  and  $\Theta_t \triangleq [\theta_t^{(1)}; \theta_t^{(2)}; \dots; \theta_t^{(M)}]^\top$ . Hence, it can be derived from Lemma C.3 that

$$\|\Delta\Theta_{T+T'}\|_F = \|\Delta V^{T'}\Theta_T\|_F = \|V^{T'}\Delta\Theta_T\|_F \leq \sigma_2^{T'} \|\Delta\Theta_T\|_F. \tag{29}$$

Hence, we only need to obtain an upper bound of  $\mathbb{E}\|\Delta\Theta_T\|^2$ . Subtracting eq. (14) from the local update rule yields that for any  $0 \leq t \leq T - 1$ ,

$$\begin{aligned}
\theta_{t+1}^{(m)} - \bar{\theta}_{t+1} &= \sum_{m' \in \mathcal{N}_m} V_{m,m'} (\theta_t^{(m')} - \bar{\theta}_t) + \frac{M-1}{M} \alpha (\bar{A}_t^{(m)} \theta_t^{(m)} + \bar{\tilde{b}}_t^{(m)} + \bar{B}_t^{(m)} w_t^{(m)}) \\
&\quad - \frac{\alpha}{M} \sum_{m'=1, m' \neq m}^M (\bar{A}_t^{(m')} \theta_t^{(m')} + \bar{\tilde{b}}_t^{(m')} + \bar{B}_t^{(m')} w_t^{(m')}).
\end{aligned}$$

This can be rewritten into the following matrix-vector form,

$$\Delta\Theta_{t+1} = V\Delta\Theta_t + [h_1; h_2; \dots; h_M]^\top,$$

where  $h_m \triangleq \frac{M-1}{M}\alpha(\bar{A}_t^{(m)}\theta_t^{(m)} + \bar{b}_t^{(m)} + \bar{B}_t^{(m)}w_t^{(m)}) - \frac{\alpha}{M}\sum_{m'=1, m' \neq m}^M (\bar{A}_t^{(m')}\theta_t^{(m')} + \bar{b}_t^{(m')} + \bar{B}_t^{(m')}w_t^{(m')})$ .

The item 2 of Lemma C.3 implies that for any  $0 \leq t \leq T-1$ ,

$$\|\Delta\Theta_{t+1}\|_F \leq \sigma_2\|\Delta\Theta_t\|_F + \sqrt{\sum_{m=1}^M \|h_m\|^2} \leq \sigma_2\|\Delta\Theta_t\|_F + \sum_{m=1}^M \|h_m\|. \quad (30)$$

Then, using eqs. (53)-(55) yields that

$$\begin{aligned} \sum_{m=1}^M \|h_m\| &\leq \frac{M-1}{M}(2\alpha)(2\rho_{\max}+2) \sum_{m=1}^M (\|\theta_t^{(m)}\| + R_{\max} + \|w_t^{(m)}\|) \\ &\stackrel{(i)}{\leq} 4\alpha(\rho_{\max}+1) \sum_{m=1}^M (\|\theta_t^{(m)} - \bar{\theta}_t\| + \|w_t^{(m)} - \bar{w}_t\| + \|\bar{\theta}_t - \theta^*\| + \|\bar{w}_t - w_t^*\| + \|\theta^*\| + \|C^{-1}(A\bar{\theta}_t + b)\|) \\ &\stackrel{(ii)}{\leq} 4\alpha(\rho_{\max}+1) \left( \sum_{m=1}^M (\|\theta_t^{(m)} - \bar{\theta}_t\| + \|w_t^{(m)} - \bar{w}_t\|) \right. \\ &\quad \left. + M\left(1 + \frac{2\rho_{\max}}{\lambda_2}\right) \|\bar{\theta}_t - \theta^*\| + M\|\bar{w}_t - w_t^*\| + \frac{M\rho_{\max}}{\lambda_2}(2\|\theta^*\| + R_{\max}) \right), \end{aligned}$$

where (i) uses the notations that  $w_t^* = -C^{-1}(A\bar{\theta}_t + b)$ , (ii) uses eqs. (39), (42) and (43). Hence, we obtain that

$$\begin{aligned} \mathbb{E}(\|\Delta\Theta_{t+1}\|_F^2) &\stackrel{(i)}{\leq} \left(1 + \frac{\sigma_2^{-2}-1}{2}\right) \sigma_2^2 \mathbb{E}(\|\Delta\Theta_t\|_F^2) + \left(1 + \frac{2}{\sigma_2^{-2}-1}\right) \mathbb{E}\left[\left(\sum_{m=1}^M \|h_m\|\right)^2\right], \\ &\stackrel{(ii)}{\leq} \frac{1+\sigma_2^2}{2} \mathbb{E}(\|\Delta\Theta_t\|_F^2) + \frac{48\alpha^2(1+\sigma_2^2)}{1-\sigma_2^2} (\rho_{\max}+1)^2 \mathbb{E}\left[2M \sum_{m=1}^M (\|\theta_t^{(m)} - \bar{\theta}_t\|^2 + \|w_t^{(m)} - \bar{w}_t\|^2) \right. \\ &\quad \left. + M^2\left(1 + \frac{2\rho_{\max}}{\lambda_2}\right)^2 (\|\bar{\theta}_t - \theta^*\|^2 + \|\bar{w}_t - w_t^*\|^2) + \frac{4M^2\rho_{\max}^2}{\lambda_2^2} (\|\theta^*\| + R_{\max})^2\right], \end{aligned} \quad (31)$$

where (i) uses eq. (30) and the fact that  $(u+v)^2 \leq (1+\sigma)u^2 + (1+\sigma^{-1})v^2$  for any  $u, v, \sigma \geq 0$ , (ii) uses  $(\sum_{i=1}^n q_i)^2 \leq n \sum_{i=1}^n q_i^2$  for any  $q_i \in \mathbb{R}$  and  $n \in \mathbb{N}^+$ . Similarly, we obtain from the update rule of  $W_t = [w_t^{(1)}; w_t^{(2)}; \dots; w_t^{(M)}]^\top \in \mathbb{R}^{M \times d}$  and eq. (15) that

$$\begin{aligned} \mathbb{E}(\|\Delta W_{t+1}\|_F^2) &\leq \frac{1+\sigma_2^2}{2} \mathbb{E}(\|\Delta W_t\|_F^2) + \frac{48\alpha^2(1+\sigma_2^2)}{1-\sigma_2^2} (\rho_{\max}+1)^2 \mathbb{E}\left[2M \sum_{m=1}^M (\|\theta_t^{(m)} - \bar{\theta}_t\|^2 + \|w_t^{(m)} - \bar{w}_t\|^2) \right. \\ &\quad \left. + M^2\left(1 + \frac{2\rho_{\max}}{\lambda_2}\right)^2 (\|\bar{\theta}_t - \theta^*\|^2 + \|\bar{w}_t - w_t^*\|^2) + \frac{4M^2\rho_{\max}^2}{\lambda_2^2} (\|\theta^*\| + R_{\max})^2\right], \end{aligned} \quad (32)$$

Summing up eqs. (31) and (32) yields that

$$\begin{aligned} &\mathbb{E}(\|\Delta\Theta_{t+1}\|_F^2 + \|\Delta W_{t+1}\|_F^2) \\ &\stackrel{(i)}{\leq} \frac{1+\sigma_2^2}{2} \mathbb{E}(\|\Delta\Theta_t\|_F^2 + \|\Delta W_t\|_F^2) + \frac{96\alpha^2(1+\sigma_2^2)}{1-\sigma_2^2} (\rho_{\max}+1)^2 \mathbb{E}\left[2M(\|\Delta\Theta_t\|_F^2 + \|\Delta W_t\|_F^2) \right. \\ &\quad \left. + M^2\left(1 + \frac{2\rho_{\max}}{\lambda_2}\right)^2 \left(\left(1 - \frac{\alpha\lambda_1}{6}\right)^t (\|\bar{\theta}_0 - \theta^*\|^2 + \|\bar{w}_0 - w_0^*\|^2) \right. \right. \\ &\quad \left. \left. + \frac{36000\beta\rho_{\max}^2(\nu+1)}{\alpha N\lambda_1\lambda_2(1-\delta)} \left(1 + \frac{\rho_{\max}}{\lambda_2}\right)^2 (\|\theta^*\|^2 + R_{\max}^2 + 1) + \frac{574\beta\sigma_2^{L/4}2^t}{M\lambda_2} \max_{m \in \mathcal{M}} (\|\theta_0^{(m)}\| + \|w_0^{(m)}\| + R_{\max})^2 \right) \right] \end{aligned}$$

$$\begin{aligned}
& + \frac{4M^2\rho_{\max}^2}{\lambda_2^2}(\|\theta^*\| + R_{\max})^2 \Big] \\
& \stackrel{(ii)}{\leq} \frac{2 + \sigma_2^2}{3} \mathbb{E}(\|\Delta\Theta_t\|_F^2 + \|\Delta W_t\|_F^2) + \frac{192M^2\alpha^2}{1 - \sigma_2^2}(\rho_{\max} + 1)^2 \left(1 + \frac{2\rho_{\max}}{\lambda_2}\right)^2 \left[ (\|\bar{\theta}_0 - \theta^*\|^2 + \|\bar{w}_0 - w_0^*\|^2) \right. \\
& \quad \left. + \left( \frac{282\beta}{\alpha\lambda_2} + \frac{8M^2\rho_{\max}^2}{\lambda_2^2} \right) (\|\theta^*\|^2 + R_{\max}^2 + 1) + \frac{574\beta\sigma_2^{L/4}2^t}{M\lambda_2} \max_{m \in \mathcal{M}} (\|\theta_0^{(m)}\| + \|w_0^{(m)}\| + R_{\max})^2 \right], \quad (33)
\end{aligned}$$

where (i) uses eq. (28), and (ii) uses  $(\|\theta^*\| + R_{\max})^2 \leq 2\|\theta^*\|^2 + 2R_{\max}^2$ ,  $\alpha \leq \frac{1 - \sigma_2}{50\sqrt{M}(\rho_{\max} + 1)}$  and  $N \geq \frac{128\rho_{\max}^2(\nu+1)}{\lambda_1(1-\delta)}(1 + \frac{\rho_{\max}}{\lambda_2})^2$ .

Iterating eq. (33) yields that

$$\begin{aligned}
\mathbb{E}(\|\Delta\Theta_T\|_F^2) & \leq \mathbb{E}(\|\Delta\Theta_T\|_F^2 + \|\Delta W_T\|_F^2) \\
& \leq \left(\frac{2 + \sigma_2^2}{3}\right)^T (\|\Delta\Theta_0\|_F^2 + \|\Delta W_0\|_F^2) + \frac{192M^2\alpha^2}{1 - \sigma_2^2}(\rho_{\max} + 1)^2 \left(1 + \frac{2\rho_{\max}}{\lambda_2}\right)^2 \\
& \quad + \left[ \frac{3}{1 - \sigma_2} (\|\bar{\theta}_0 - \theta^*\|^2 \|\bar{w}_0 - w_0^*\|^2 + \left( \frac{282\beta}{\alpha\lambda_2} + \frac{8M^2\rho_{\max}^2}{\lambda_2^2} \right) (\|\theta^*\|^2 + R_{\max}^2 + 1)) \right. \\
& \quad \left. + \frac{574\beta\sigma_2^{L/4}2^T}{M\lambda_2} \max_{m \in \mathcal{M}} (\|\theta_0^{(m)}\| + \|w_0^{(m)}\| + R_{\max})^2 \right] \\
& \leq \|\Delta\Theta_0\|_F^2 + \|\Delta W_0\|_F^2 + \mathcal{O}\left[ \frac{M^2\alpha^2}{(1 - \sigma_2)^2} \left( \frac{\beta}{\alpha} + M^2 \right) + \frac{M\beta\alpha^2\sigma_2^{L/4}2^T}{1 - \sigma_2} \right] \\
& \stackrel{(i)}{\leq} \mathcal{O}\left( 1 + \frac{M^4\beta\alpha}{(1 - \sigma_2)^2} + \frac{M\beta\alpha\sigma_2^{L/4}2^T}{1 - \sigma_2} \right), \quad (34)
\end{aligned}$$

where (i) uses  $\alpha \leq \frac{\beta}{2\rho_{\max} + 2} = \mathcal{O}(\beta)$ . Substituting the above inequality into eq. (29) proves eq. (11).

To summarize, the following conditions of the hyperparameters are used in the proof of Theorem 1, including those required by Corollary 2 and lemma C.4.

$$\begin{aligned}
\alpha & \leq \min \left( \frac{\beta}{2\rho_{\max} + 2}, \frac{\lambda_1}{40\rho_{\max}^2} \left(1 + \frac{\rho_{\max}}{\lambda_2}\right)^{-2}, \frac{1}{\lambda_1}, \frac{\beta\lambda_2}{4\rho_{\max}}, \frac{\beta\lambda_1\lambda_2^3}{23040\rho_{\max}^4}, \frac{\beta\lambda_2^2}{53\rho_{\max}^2}, \frac{\beta\lambda_1\lambda_2}{432\rho_{\max}^2}, \right. \\
& \quad \left. \frac{1 - \sigma_2}{50\sqrt{M}(\rho_{\max} + 1)}, \frac{\beta\lambda_2}{2\lambda_1}, \frac{\beta\lambda_1}{2\lambda_2} \right) = \min\{\mathcal{O}(\beta), \mathcal{O}(M^{-1/2}(1 - \sigma_2))\} \quad (35)
\end{aligned}$$

$$\beta \leq \min \left( \frac{1}{\lambda_2}, \frac{2}{5(2\rho_{\max} + 3)} \right) = \mathcal{O}(1) \quad (36)$$

$$N \geq \max \left( \frac{8c_{\text{sd}}}{\lambda_2}, \frac{8c_{\text{var},3}\Omega_A^2}{\lambda_1^2}, \frac{12\beta c_8}{\alpha\lambda_1} \right) = \max\{\mathcal{O}(1), \mathcal{O}(\beta/\alpha)\}. \quad (37)$$

$$L \geq \frac{12 \ln M + (8M + 10) \ln \rho_{\max}}{\ln \sigma_2^{-1}} = \mathcal{O}\left( \frac{M}{1 - \sigma_2} \right) \quad (38)$$

Under the above conditions, we choose the following hyperparameter values.

$$\begin{aligned}
\alpha & = \mathcal{O}(M^{-1/2}(1 - \sigma_2)), \beta = \mathcal{O}(1) \\
T & = \frac{6}{\alpha\lambda_1} \ln \epsilon^{-1} = \mathcal{O}\left( \frac{\sqrt{M} \ln \epsilon^{-1}}{1 - \sigma_2} \right) \\
N & = \frac{\beta}{\alpha\epsilon} = \mathcal{O}\left( \frac{\sqrt{M}}{\epsilon(1 - \sigma_2)} \right) \\
L & = \frac{4}{\ln \sigma_2^{-1}} \left( \ln \left( \frac{\beta}{M\epsilon} \right) + T \ln 2 \right) + \frac{12 \ln M + (8M + 10) \ln \rho_{\max}}{\ln \sigma_2^{-1}} = \mathcal{O}\left( \frac{\sqrt{M} \ln \epsilon^{-1}}{(1 - \sigma_2)^2} + \frac{M}{1 - \sigma_2} \right) \leq \mathcal{O}\left( \frac{M \ln \epsilon^{-1}}{(1 - \sigma_2)^2} \right) \\
T' & = \frac{1}{\ln \sigma_2^{-1}} \ln \left( \epsilon^{-1} \mathcal{O}\left( 1 + \frac{M^4\beta\alpha}{(1 - \sigma_2)^2} + \frac{M\beta\alpha\sigma_2^{L/4}2^T}{1 - \sigma_2} \right) \right)
\end{aligned}$$

$$= \frac{1}{\ln \sigma_2^{-1}} \ln \left( \epsilon^{-1} \mathcal{O} \left( 1 + \frac{M^{3.5}}{1 - \sigma_2} + \frac{\alpha M^2 \epsilon}{1 - \sigma_2} \right) \right) = \mathcal{O} \left( \frac{1}{1 - \sigma_2} \ln \left( \frac{M}{\epsilon(1 - \sigma_2)} \right) \right).$$

Substituting these hyperparameters into eqs. (10) and (11) implies  $\mathbb{E}(\|\bar{\theta}_T - \theta^*\|^2), \mathbb{E}(\|\theta_{T+T'}^{(m)} - \bar{\theta}_T\|^2) \leq \mathcal{O}(\epsilon)$ , so  $\mathbb{E}(\|\theta_{T+T'}^{(m)} - \theta^*\|^2) \leq 2\mathbb{E}(\|\theta_{T+T'}^{(m)} - \bar{\theta}_T\|^2) + 2\mathbb{E}(\|\bar{\theta}_T - \theta^*\|^2) \leq \mathcal{O}(\epsilon)$ . Therefore, the overall communication complexity for synchronizing  $\theta_t^{(m)}$  is  $T + T' = \mathcal{O}(\frac{\sqrt{M} \ln \epsilon^{-1}}{1 - \sigma_2})$ , and the total sample complexity is  $NT = \mathcal{O}(\frac{M \ln \epsilon^{-1}}{\epsilon(1 - \sigma_2)^2})$ .

## C Supporting Lemmas

In this section, we prove some supporting lemmas that are used throughout the analysis of Algorithm 1.

**Lemma C.1.** *Regarding the terms defined in Appendix A, their norms have the following upper bounds.*

$$\|A_t\|_F, \|\bar{A}_t\|_F, \|A\|_F \leq 2\rho_{\max}, \quad (39)$$

$$\|B_t\|_F, \|\bar{B}_t\|_F, \|B\|_F \leq \rho_{\max}, \quad (40)$$

$$\|C_t\|_F, \|\bar{C}_t\|_F, \|C\|_F \leq 1, \quad (41)$$

$$\|b_t^{(m)}\|, \|\bar{b}_t^{(m)}\|, \|\bar{b}_t\|, \|\bar{\bar{b}}_t\|, \|b\| \leq \rho_{\max} R_{\max}, \quad (42)$$

$$\|C^{-1}\| = \lambda_2^{-1}. \quad (43)$$

*Proof.* Consider any two vectors  $u, v \in \mathbb{R}^d$ , we have that  $\|uv^\top\|_F = \sqrt{\text{tr}(vu^\top uv^\top)} = \|u\| \|v\|$ . Therefore, by Assumption 3, we obtain that

$$\|A_t\|_F \leq \rho_t \|\phi(s_t)\| \|\gamma \phi(s_{t+1}) - \phi(s_t)\| \leq \rho_{\max} [\gamma \|\phi(s_{t+1})\| + \|\phi(s_t)\|] \leq 2\rho_{\max},$$

$$\|B_t\|_F \leq \gamma \rho_t \|\phi(s_{t+1})\| \|\phi(s_t)\| \leq \rho_{\max},$$

$$\|C_t\|_F \leq \|\phi(s_t)\|^2 \leq 1,$$

$$\|b_t^{(m)}\| \leq \rho_t R_t^{(m)} \|\phi(s_t)\| \leq \rho_{\max} R_{\max}.$$

The proof for  $\|A_t^{(m)}\|_F, \|\tilde{b}_t^{(m)}\|$ , etc. is similar.

On the other hand, by Jensen's inequality, we obtain that

$$\|A\|_F = \|\mathbb{E}_{\pi_b}[A_t]\|_F \leq \mathbb{E}_{\pi_b}\|A_t\|_F \leq 2\rho_{\max}, \quad \|\bar{A}_t\|_F \leq \frac{1}{N} \sum_{i=tN}^{(t+1)N-1} \|A_i\|_F \leq 2\rho_{\max}.$$

The proof for the other remaining matrices in eqs. (39)-(42) is similar by using the Jensen's inequality. Finally, we prove eq. (43). Note that  $-C = \mathbb{E}_{\pi_b}(\phi(s_t)\phi(s_t)) \succ 0$  with  $\lambda_{\min}(-C) = -\lambda_{\max}(C) = \lambda_2$ . Hence,  $\|C^{-1}\| = \lambda_{\max}(-C^{-1}) = \lambda_{\min}^{-1}(-C) = \lambda_2^{-1}$ .  $\square$

**Lemma C.2.** *Suppose the MDP trajectory  $\{s_t, a_t\}_{t \geq 0}$  is generated following a behavioral policy  $\pi_b$  where  $a_t \triangleq \{a_t^{(m)}\}_m$ . For any deterministic mappings  $Y : \mathcal{S} \times \mathcal{A}_1 \times \dots \times \mathcal{A}_M \times \mathcal{S} \rightarrow \mathbb{R}^{p \times q}$  such that  $\|Y(s, a, s')\|_F \leq C_y, \forall s, s' \in \mathcal{S}, a^{(m)} \in \mathcal{A}_m$  where  $a = \{a^{(m)}\}_m$ , we have*

$$\begin{aligned} \left\| \mathbb{E} \left[ \frac{1}{N} \sum_{i=tN}^{(t+1)N-1} Y(s_i, a_i, s_{i+1}) \middle| \mathcal{F}_t \right] - \bar{Y} \right\| &\leq \frac{2\nu C_y}{N(1-\delta)}, \\ \mathbb{E} \left[ \left\| \frac{1}{N} \sum_{i=tN}^{(t+1)N-1} Y(s_i, a_i, s_{i+1}) - \bar{Y} \right\|_F^2 \middle| \mathcal{F}_t \right] &\leq \frac{8C_y^2(\nu+1)}{N(1-\delta)}, \end{aligned}$$

where  $\bar{Y} = \mathbb{E}Y(s_i, a_i, s_{i+1})$ .

**Note:** A simplified version of the above lemma has been proposed and proved in Xu et al. (2020a), where  $a_i$  and  $s_{i+1}$  are omitted in the above inequality. We add  $a_i$  and  $s_{i+1}$  so that this lemma can be better applied to the quantities  $A_i$ ,  $B_i$ ,  $C_i$  and  $b_i^{(m)}$  which rely on  $s_i$  as well as  $a_i$  and  $s_{i+1}$ . The proof logic is very similar to that of Xu et al. (2020a) and thus omitted here.

**Corollary 1.** *Regarding the terms defined in Appendix A, they have the following upper bounds.*

$$\mathbb{E}[\|\bar{C}_t - C\|_F | \mathcal{F}_t] \leq \frac{2\nu\rho_{\max}}{N(1-\delta)} \quad (44)$$

$$\mathbb{E}[\|\bar{B}_t - B\|_F^2 | \mathcal{F}_t] \leq \frac{8\rho_{\max}^2(\nu+1)}{N(1-\delta)} \quad (45)$$

$$\mathbb{E}[\|\bar{b}_t - b\|^2 | \mathcal{F}_t] \leq \frac{8\rho_{\max}^2 R_{\max}^2(\nu+1)}{N(1-\delta)} \quad (46)$$

$$\mathbb{E}[\|\bar{A}_t - \bar{C}_t C^{-1} A\|_F^2 | \mathcal{F}_t] \leq \frac{32\rho_{\max}^2(\nu+1)}{N(1-\delta)} \left(1 + \frac{1}{\lambda_2}\right)^2 \quad (47)$$

$$\mathbb{E}[\|\bar{b}_t - \bar{C}_t C^{-1} b\|^2 | \mathcal{F}_t] \leq \frac{8\rho_{\max}^2 R_{\max}^2(\nu+1)}{N(1-\delta)} \left(1 + \frac{1}{\lambda_2}\right)^2 \quad (48)$$

$$\mathbb{E}[\|\bar{A}_t - \bar{B}_t C^{-1} A - A^\top C^{-1} A\|_F^2 | \mathcal{F}_t] \leq \frac{32\rho_{\max}^2(\nu+1)}{N(1-\delta)} \left(1 + \frac{\rho_{\max}}{\lambda_2}\right)^2, \quad (49)$$

*Proof.* Let  $Y(s, a, s') = -\gamma\rho(s, a)\phi(s')\phi(s)^\top$  in Lemma C.2. Then it can be checked that  $Y(s_t, a_t, s_{t+1}) = B_t$ ,  $C_y = \rho_{\max}$ ,  $\frac{1}{N} \sum_{i=tN}^{(t+1)N-1} Y(s_i, a_i, s_{i+1}) = \bar{B}_t$ , and  $\bar{Y} = \mathbb{E}_{\pi_b} Y(s_i, a_i, s_{i+1}) = B$ .

Applying Lemma C.2 to these equations proves eq. (45). The eqs. (44) and (46) can be proved in a similar way.

Let  $Y(s, a, s') = \rho(s, a)\phi(s)[\gamma\phi(s') - \phi(s)]^\top + \gamma\rho(s, a)\phi(s')\phi(s)^\top C^{-1}A$ . Then, it can be checked that  $Y(s_t, a_t, s_{t+1}) = A_t - B_t C^{-1}A$ ,  $\frac{1}{N} \sum_{i=tN}^{(t+1)N-1} Y(s_i, a_i, s_{i+1}) = \bar{A}_t - \bar{B}_t C^{-1}A$ . Moreover,

$$\begin{aligned} \|Y(s, a, s')\|_F &\leq \rho_{\max}(\gamma+1) + \gamma\rho_{\max}\|C^{-1}\| \|A\| \leq 2\rho_{\max} + \rho_{\max}(\lambda_2^{-1})(2\rho_{\max}) = 2\rho_{\max}(1 + \rho_{\max}/\lambda_2) := C_y, \\ \bar{Y} &= \mathbb{E}_{\pi_b} Y(s_i, a_i, s_{i+1}) = A - BC^{-1}A = A^\top C^{-1}A. \end{aligned}$$

Applying Lemma C.2 to these equations proves eq. (49). The equations (47) and (48) can be proved in a similar way.  $\square$

**Lemma C.3.** *The doubly stochastic matrix  $V$  and the difference matrix  $\Delta = I - \frac{1}{M}\mathbf{1}\mathbf{1}^\top$  have the following properties:*

1.  $\Delta V = V\Delta = V - \frac{1}{M}\mathbf{1}\mathbf{1}^\top$
2. For any  $x \in \mathbb{R}^M$  and  $n \in \mathbb{N}^+$ ,  $\|V^n \Delta x\| \leq \sigma_2^n \|\Delta x\|$  ( $\sigma_2$  is the second largest singular value of  $V$ ). Hence, for any  $H \in \mathbb{R}^{M \times M}$ ,  $\|V^n \Delta H\|_F \leq \sigma_2^n \|\Delta H\|_F$

*Proof.* The first item can be proved by the following two equalities.

$$\begin{aligned} \Delta V &= \left(I - \frac{1}{M}\mathbf{1}\mathbf{1}^\top\right)V = V - \frac{1}{M}\mathbf{1}\mathbf{1}^\top V = V - \frac{1}{M}\mathbf{1}\mathbf{1}^\top \\ V\Delta &= V\left(I - \frac{1}{M}\mathbf{1}\mathbf{1}^\top\right) = V - \frac{1}{M}V\mathbf{1}\mathbf{1}^\top = V - \frac{1}{M}\mathbf{1}\mathbf{1}^\top \end{aligned}$$

The proof of the item 2 can be seen in page 3 of Qu & Li (2017).  $\square$

Based on Lemma C.3, we obtain the following inexactness of importance sampling ratio estimation  $\hat{\rho}_t^{(m)} \approx \rho_t$ .

**Corollary 2.** *Under Assumption 4 and choosing  $L \geq \mathcal{O}(\frac{\ln M + M \ln \rho_{\max}}{\ln \sigma_2^{-1}})$ , the estimation error of the inexact global importance sampling ratio  $\hat{\rho}_t^{(m)}$  satisfies  $\sum_{m=1}^M (\hat{\rho}_t^{(m)} - \rho_t)^2 \leq \sigma_2^{L/4}$ . Therefore, the following inequalities hold.*

$$\sum_{m=1}^M \|A_t^{(m)} - A_t\|_F^2, \sum_{m=1}^M \|\bar{A}_t^{(m)} - \bar{A}_t\|_F^2 \leq 4\sigma_2^{L/4} \quad (50)$$

$$\sum_{m=1}^M \|B_t^{(m)} - B_t\|_F^2, \sum_{m=1}^M \|\bar{B}_t^{(m)} - \bar{B}_t\|_F^2 \leq \sigma_2^{L/4}, \quad (51)$$

$$\sum_{m=1}^M \|\tilde{b}_t^{(m)} - b_t^{(m)}\|^2, \sum_{m=1}^M \|\bar{\tilde{b}}_t^{(m)} - \bar{b}_t^{(m)}\|^2 \leq \sigma_2^{L/4} R_{\max}^2. \quad (52)$$

As a result, the following upper bounds hold.

$$\|A_t^{(m)}\|_F, \|\bar{A}_t^{(m)}\|_F \leq 2\rho_{\max} + 2 \quad (53)$$

$$\|B_t^{(m)}\|_F, \|\bar{B}_t^{(m)}\|_F \leq \rho_{\max} + 1 \quad (54)$$

$$\|\tilde{b}_t^{(m)}\|, \|\bar{\tilde{b}}_t^{(m)}\| \leq R_{\max}(\rho_{\max} + 1) \quad (55)$$

*Proof.* Eq. (5) can be rewritten into the following matrix form.

$$[\tilde{\rho}_{t,L}^{(1)}; \dots; \tilde{\rho}_{t,L}^{(M)}] = V^L [\tilde{\rho}_{t,0}^{(1)}; \dots; \tilde{\rho}_{t,0}^{(M)}].$$

Hence, the item 1 of Lemma C.3 yields that

$$\Delta[\tilde{\rho}_{t,L}^{(1)}; \dots; \tilde{\rho}_{t,L}^{(M)}] = V^L \Delta[\tilde{\rho}_{t,0}^{(1)}; \dots; \tilde{\rho}_{t,0}^{(M)}].$$

Then the item 2 of Lemma C.3 yields that

$$\|\Delta[\tilde{\rho}_{t,L}^{(1)}; \dots; \tilde{\rho}_{t,L}^{(M)}]\|^2 \leq \sigma_2^{2L} \|\Delta[\tilde{\rho}_{t,0}^{(1)}; \dots; \tilde{\rho}_{t,0}^{(M)}]\|^2. \quad (56)$$

Denote  $\rho_{\min} := \min_{m \in \mathcal{M}} \rho_t^{(m)}$ . Then Assumption 4 implies that  $\tilde{\rho}_{t,0}^{(m)} = \ln \rho_t^{(m)} \in [\ln \rho_{\min}, \ln \rho_{\max}]$ . Then it can be proved by iterating eq. (5) that  $\tilde{\rho}_{t,L}^{(m)} \in [\ln \rho_{\min}, \ln \rho_{\max}]$ . Hence,

$$\frac{1}{M} \ln \rho_t = \frac{1}{M} \sum_{m=1}^M \ln \rho_t^{(m)} \in [\ln \rho_{\min}, \ln \rho_{\max}] \quad (57)$$

Then eqs. (56) and (57) imply that

$$\sum_{m=1}^M \left( \tilde{\rho}_{t,L}^{(m)} - \frac{1}{M} \ln \rho_t \right)^2 \leq \sigma_2^{2L} \sum_{m=1}^M \left( \tilde{\rho}_{t,0}^{(m)} - \frac{1}{M} \ln \rho_t \right)^2 \leq M \sigma_2^{2L} \ln^2(\rho_{\max}/\rho_{\min}). \quad (58)$$

Hence,

$$\left| \tilde{\rho}_{t,L}^{(m)} - \frac{1}{M} \ln \rho_t \right| \leq \sqrt{M} \sigma_2^L \ln \left( \frac{\rho_{\max}}{\rho_{\min}} \right) \stackrel{(i)}{\leq} \frac{1}{2M} \ln \left( \frac{\rho_{\max}}{\rho_{\min}} \right), \quad (59)$$

where (i) uses the conditions that  $L \geq \frac{12 \ln M + (8M+10) \ln \rho_{\max}}{\ln(\sigma_2^{-1})}$  and  $\sigma_2 \in [0, 1)$ .

Hence, eqs. (57) and (58) imply that

$$\tilde{\rho}_{t,L}^{(m)} \leq \ln \rho_{\max} + \frac{1}{2M} \ln(\rho_{\max}/\rho_{\min}). \quad (60)$$



Therefore, we obtain that

$$\begin{aligned}
\sum_{m=1}^M (\hat{\rho}_t^{(m)} - \rho_t)^2 &\stackrel{(i)}{=} \sum_{m=1}^M (e^{M\tilde{\rho}_{t,L}^{(m)}} - e^{\ln \rho_t})^2 \\
&\stackrel{(ii)}{\leq} \sum_{m=1}^M [\max(e^{M\tilde{\rho}_{t,L}^{(m)}}, e^{\ln \rho_t})]^2 (M\tilde{\rho}_{t,L}^{(m)} - \ln \rho_t)^2 \\
&\stackrel{(iii)}{\leq} M^3 \sigma_2^{2L} \rho_{\max}^M \sqrt{\rho_{\max}/\rho_{\min}} \ln^2(\rho_{\max}/\rho_{\min}) \\
&\stackrel{(iv)}{\leq} M^3 \sigma_2^{2L} (\rho_{\max}^{M+2.5}/\rho_{\min}^{2.5}), \tag{61}
\end{aligned}$$

where (i) uses eq. (6), (ii) uses the Lagrange's Mean Value Theorem, (iii) uses eqs. (57), (58) and (60), (iv) uses the inequality that  $\ln x < x$  for  $x = \rho_{\max}/\rho_{\min} \geq 1$ .

Since at least one of  $\{\tilde{\rho}_{t,0}^{(m)}\}_{m \in \mathcal{M}}$  equals  $\ln \rho_{\min}$ , we have

$$\ln \rho_t = \sum_{m=1}^M \tilde{\rho}_{t,0}^{(m)} \leq \ln \rho_{\min} + (M-1) \ln \rho_{\max}. \tag{62}$$

Then, eqs. (59) and (62) imply that

$$\tilde{\rho}_{t,L}^{(m)} \leq \frac{1}{2M} \ln \rho_{\min} + \left(1 - \frac{1}{2M}\right) \ln \rho_{\max} \tag{63}$$

Hence, we conclude that

$$\sum_{m=1}^M (\hat{\rho}_t^{(m)} - \rho_t)^2 \stackrel{(i)}{=} \sum_{m=1}^M (e^{M\tilde{\rho}_{t,L}^{(m)}} - e^{\ln \rho_t})^2 \leq \sum_{m=1}^M \max(e^{2M\tilde{\rho}_{t,L}^{(m)}}, e^{2\ln \rho_t}) \stackrel{(ii)}{\leq} M \rho_{\min} \rho_{\max}^{2M-1} \tag{64}$$

where (i) uses eq. (6), (ii) uses eqs. (62) and (63).

When  $\rho_{\min} \geq \sigma_2^{L/2}$ , eq. (61) implies that  $\sum_{m=1}^M (\hat{\rho}_t^{(m)} - \rho_t)^2 \leq M^3 \rho_{\max}^{M+2.5} \sigma_2^{0.75L}$ . When  $\rho_{\min} < \sigma_2^{L/2} < 1$ , eq. (64) implies that  $\sum_{m=1}^M (\hat{\rho}_t^{(m)} - \rho_t)^2 \leq M \rho_{\max}^{2M-1} \sigma_2^{L/2}$ . Both imply  $\sum_{m=1}^M (\hat{\rho}_t^{(m)} - \rho_t)^2 \leq \sigma_2^{L/4}$  since  $L \geq \frac{12 \ln M + (8M+10) \ln \rho_{\max}}{\ln(\sigma_2^{-1})}$ .

Then, eq. (50) can be proved as follows.

$$\sum_{m=1}^M \|A_t^{(m)} - A_t\|_F^2 \leq \|\phi(s_t)[\gamma\phi(s_{t+1}) - \phi(s_t)]^\top\|_F^2 \sum_{m=1}^M (\hat{\rho}_t^{(m)} - \rho_t)^2 \leq (1+\gamma)^2 \sigma_2^{L/4} \leq 4\sigma_2^{L/4} \tag{65}$$

The above inequality implies that  $\sum_{m=1}^M \|\bar{A}_t^{(m)} - \bar{A}_t\|_F^2 = \sum_{m=1}^M \|\frac{1}{N} \sum_{i=tN}^{(t+1)N-1} (A_i^{(m)} - A_i)\|_F^2 \leq 4\sigma_2^{L/4}$ , where  $\leq$  applies Jensen's inequality to the convex function  $\|\cdot\|_F^2$ . Eqs. (51) and (52) can be proved similarly.

Eq. (50) implies that  $\|A_t^{(m)} - A_t\|_F, \|\bar{A}_t^{(m)} - \bar{A}_t\|_F \leq 2$ . Hence, eq. (53) can be proved using triangle inequality and eq. (39). Eqs. (54) and (55) can be proved similarly.  $\square$

The proof of Corollary 2 introduces a new technique, which includes discussion of two cases:  $\rho_{\min} := \min_{m \in \mathcal{M}} \rho_t^{(m)}$  lies in  $[\sigma_2^{L/2}, \rho_{\max}]$  and  $(0, \sigma_2^{L/2}]$ . This is necessary as the local average is applied to  $\ln \hat{\rho}_t^{(m)}$ , which may be a large negative number that cannot ensure a small consensus error for a fixed number of local average steps  $L$ .

**Lemma C.4.** *Under the update rules of Algorithm 1 and choosing  $L \geq \frac{12 \ln M + (8M+10) \ln \rho_{\max}}{\ln(\sigma_2^{-1})}$ ,  $\alpha \leq \frac{\beta}{2\rho_{\max}+2}$ , the parameters have the following upper bound.*

$$\max_{m \in \mathcal{M}} \|\theta_T^{(m)}\| + \max_{m \in \mathcal{M}} \|w_T^{(m)}\| \leq 2 \max_{m \in \mathcal{M}} (\|\theta_0^{(m)}\| + \|w_0^{(m)}\| + R_{\max}) [1 + \beta(2\rho_{\max} + 3)]^T. \tag{66}$$

*Proof.* Since  $L \geq \frac{12 \ln M + (8M+10) \ln \rho_{\max}}{\ln(\sigma_2^{-1})}$ , eqs. (53)-(55) hold. Hence, these equations and the update rule imply that

$\|\theta_{t+1}^{(m)}\| \leq \sum_{m' \in \mathcal{N}_m} V_{m,m'} \|\theta_t^{(m')}\| + \alpha(\rho_{\max} + 1)(2\|\theta_t^{(m)}\| + R_{\max} + \|w_t^{(m)}\|)$ . Taking maximum with respect to  $m$  yields that

$$\begin{aligned} \max_{m \in \mathcal{M}} \|\theta_{t+1}^{(m)}\| &\leq \max_{m \in \mathcal{M}} \sum_{m' \in \mathcal{N}_m} V_{m,m'} \max_{m'' \in \mathcal{M}} \|\theta_t^{(m'')}\| + \alpha(\rho_{\max} + 1)(2 \max_{m \in \mathcal{M}} \|\theta_t^{(m)}\| + R_{\max} + \max_{m \in \mathcal{M}} \|w_t^{(m)}\|). \\ &\leq \alpha(\rho_{\max} + 1)(2 \max_{m \in \mathcal{M}} \|\theta_t^{(m)}\| + R_{\max} + \max_{m \in \mathcal{M}} \|w_t^{(m)}\|) + \max_{m \in \mathcal{M}} \|w_t^{(m)}\|. \end{aligned} \quad (67)$$

Similarly, it can be obtained that

$$\max_{m \in \mathcal{M}} \|w_{t+1}^{(m)}\| \leq 2\beta(\rho_{\max} + 1) \max_{m \in \mathcal{M}} \|\theta_t^{(m)}\| + (1 + \beta) \max_{m \in \mathcal{M}} \|w_t^{(m)}\| + \beta R_{\max}(\rho_{\max} + 1). \quad (68)$$

Adding up eqs. (67) and (68) yields that

$$\begin{aligned} &\max_{m \in \mathcal{M}} \|\theta_{t+1}^{(m)}\| + \max_{m \in \mathcal{M}} \|w_{t+1}^{(m)}\| \\ &\leq 2(\alpha + \beta)(\rho_{\max} + 1) \max_{m \in \mathcal{M}} \|\theta_t^{(m)}\| + (\alpha\rho_{\max} + \alpha + \beta + 1) \max_{m \in \mathcal{M}} \|w_t^{(m)}\| + R_{\max}(\alpha + \beta)(\rho_{\max} + 1) \\ &\stackrel{(i)}{\leq} \beta(2\rho_{\max} + 3) \max_{m \in \mathcal{M}} \|\theta_t^{(m)}\| + (1.5\beta + 1) \max_{m \in \mathcal{M}} \|w_t^{(m)}\| + 0.5\beta R_{\max}(2\rho_{\max} + 3) \\ &\leq [1 + \beta(2\rho_{\max} + 3)] \left( \max_{m \in \mathcal{M}} \|\theta_t^{(m)}\| + \max_{m \in \mathcal{M}} \|w_t^{(m)}\| \right) + 0.5\beta R_{\max}(2\rho_{\max} + 3), \end{aligned}$$

where (i) uses the condition that  $\alpha \leq \frac{\beta}{2\rho_{\max} + 2}$  and (ii) uses  $\rho_{\max} \geq 1$ . By iterating the inequality above and using  $\max_{m \in \mathcal{M}} \|\theta_0^{(m)}\| + \max_{m \in \mathcal{M}} \|w_0^{(m)}\| \leq 2 \max_{m \in \mathcal{M}} (\|\theta_0^{(m)}\| + \|w_0^{(m)}\|)$ , we prove eq. (66).  $\square$