
Finite-Time Analysis of Fully Decentralized Single-Timescale Actor-Critic

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Decentralized Actor-Critic (AC) algorithms have been widely utilized for multi-
2 agent reinforcement learning (MARL) and have achieved remarkable success.
3 Apart from its empirical success, the theoretical convergence property of decentral-
4 ized AC algorithms is largely unexplored. The existing finite-time convergence
5 results are derived based on either double-loop update or two-timescale step sizes
6 rule, which is not often adopted in real implementation. In this work, we introduce
7 a fully decentralized AC algorithm, where actor, critic, and global reward estimator
8 are updated in an alternating manner with step sizes being of the same order, namely,
9 we adopt the *single-timescale* update. Theoretically, using linear approximation for
10 value and reward estimation, we show that our algorithm has sample complexity of
11 $\tilde{\mathcal{O}}(\epsilon^{-2})$ under Markovian sampling, which matches the optimal complexity with
12 double-loop implementation (here, $\tilde{\mathcal{O}}$ hides a log term). The sample complexity
13 can be improved to $\mathcal{O}(\epsilon^{-2})$ under the i.i.d. sampling scheme. The central to
14 establishing our complexity results is *the hidden smoothness of the optimal critic*
15 *variable* we revealed. We also provide a local action privacy-preserving version
16 of our algorithm and its analysis. Finally, we conduct experiments to show the
17 superiority of our algorithm over the existing decentralized AC algorithms.

18 1 Introduction

19 Multi-agent reinforcement learning (MARL) [16, 30] has been very successful in various models of
20 multi-agent systems, such as robotics [14], autonomous driving [37], Go [25], etc. MARL has been
21 extensively explored in the past decades; see, e.g., [18, 20, 41, 26, 8, 22]. These works either focus
22 on the setting where a central controller is available, or assuming a common reward function for all
23 agents. Among the many cooperative MARL settings, the work [42] proposes the fully decentralized
24 MARL with networked agents. In this setting, each agent maintains a private heterogeneous reward
25 function, and agents can only access local/neighbor information through communicating with its
26 neighboring agents on the network. Then, the objective of all agents is to jointly maximize the average
27 long-term reward through interacting with environment modeled by multi-agent Markov decision
28 process (MDP). They proposed the decentralized Actor-Critic (AC) algorithm to solve this MARL
29 problem, and showed its impressive performance. However, the theoretical convergence properties
30 of such class of decentralized AC algorithms are largely unexplored; see [41] for a comprehensive
31 survey. In this work, our goal is to establish the strong finite-time convergence results under this fully
32 decentralized MARL setting. We first review some recent progresses on this line of research below.

33 **Related works and motivations.** The first fully decentralized AC algorithm with provable con-
34 vergence guarantee was proposed by [42], and they achieved asymptotic convergence results under
35 two-time scale step sizes, which requires actor's step sizes to diminish in a faster scale than the critic's
36 step sizes. The sample complexities of decentralized AC were established recently. In particular, [6]

37 and [11] independently propose two communication efficient decentralized AC algorithms with opti-
 38 mal sample complexity of $\mathcal{O}(\varepsilon^{-2} \log(\varepsilon^{-1}))$ under Markovian sampling scheme. Their analysis are
 39 based on *double-loop* implementation, where each policy optimization step follows a nearly accurate
 40 critic optimization step (a.k.a. policy evaluation), i.e., solving the critic optimization subproblem to
 41 ε -accuracy. Such a double-loop scheme requires careful tuning of two additional hyper-parameters,
 42 which are the batch size and inner loop size. In particular, the batch size and inner loop size need to be
 43 of order $\mathcal{O}(\varepsilon^{-1})$ and $\mathcal{O}(\log(\varepsilon^{-1}))$ in order to achieve their sample complexity results, respectively.
 44 In practice, single-loop algorithmic framework is often utilized, where one updates the actor and
 45 critic in an alternating manner by performing only one algorithmic iteration for both of the two
 46 subproblems; see, e.g., [23, 18, 15, 39]. The work [38] proposes a new decentralized AC algorithm
 47 based on such a single-loop alternative update. Nevertheless, they have to adopt *two-timescale* step
 48 sizes rule to ensure convergence, which requires actor’s step sizes to diminish in a faster scale than
 49 the critic’s step sizes. Due to the separation of the step sizes, the critic optimization sub-problem
 50 is solved exactly when the number of iterations tends to ∞ . Such a restriction on the step size will
 51 slow down the convergence speed of the algorithm. As a consequence, they only obtain sub-optimal
 52 sample complexity of $\mathcal{O}(\varepsilon^{-\frac{5}{2}})$. In practice, most algorithms are implemented with *single-timescale*
 53 step size rule, where the step sizes for actor and critic updates are of the same order. Though there
 54 are some theoretical achievements for single-timescale update in other areas such as TDC [31] and
 55 bi-level optimization [4], similar theoretical understanding under AC setting is largely unexplored.

56 Indeed, even when reducing to single-agent setting, the convergence property of single-timescale
 57 AC algorithm is not well established. The works [9, 10] establish the finite-time convergence result
 58 under a special single-timescale implementation, where they attain the sample complexity of $\mathcal{O}(\varepsilon^{-2})$.
 59 However, their analysis is based on an algorithm where the critic optimization step is formulated as a
 60 least-square temporal difference (LSTD) at each iteration, where they need to sample the transition
 61 tuples for $\tilde{\mathcal{O}}(\varepsilon^{-1})$ times to form the data matrix in the LSTD problem. Then, they solve the LSTD
 62 problem in a closed-form fashion, which requires to invert a matrix of large size. Later, [4] obtains the
 63 same sample complexity using TD(0) update for critic variables under i.i.d. sampling. Nonetheless,
 64 their analysis highly relies on the assumption that the Jacobian of the stationary distribution is
 65 Lipschitz continuous, which is not justified in their work.

66 The above observations motivate us to ask the following question:

67 *Can we establish finite-time convergence result for decentralized AC algorithm with single-timescale*
 68 *step sizes rule?*¹

69 **Main contributions.** By answering this question positively, we have the following contributions:

- 70 • We design a fully decentralized AC algorithm, which employs a *single-timescale* step sizes
 71 rule and adopts Markovian sampling scheme. The proposed algorithm allows communication
 72 between agents for every K_c iterations with K_c being any integer lies in $[1, \mathcal{O}(\varepsilon^{-\frac{1}{2}})]$, rather
 73 than communicating at each iteration as adopted by previous single-loop decentralized AC
 74 algorithms [38, 42].
- 75 • Using linear approximation for value and reward estimation, we establish the *finite-time*
 76 convergence result for such an algorithm under the standard assumptions. In particular, we
 77 show that the algorithm has the sample complexity of $\tilde{\mathcal{O}}(\varepsilon^{-2})$, which matches the optimal
 78 complexity up to a logarithmic term. In addition, we show that the logarithmic term can be
 79 removed under the i.i.d. sampling scheme. Note that these convergence results are valid for
 80 all the above mentioned choices for K_c .
- 81 • To preserve the privacy of local actions, we propose a variant of our algorithm which utilizes
 82 noisy local rewards for estimating global rewards. We show that such an algorithm will
 83 maintain the optimal sample complexity at the expense of communicating at each iteration.

84 The underlying principle for obtaining the above convergence results is that we reveal *the hidden*
 85 *smoothness of the optimal critic variable*, so that we can derive an approximate descent on the
 86 averaged critic’s optimal gap at each iteration. Consequently, we can resort to the classic convergence
 87 analysis for alternating optimization algorithms to establish the approximate ascent property of the
 88 overall optimization process, which leads to the final sample complexity results.

¹As convention [9], when we use "single-timescale", it means we utilize a single-loop algorithmic framework
 with single-timescale step sizes rule.

89 Another technical highlight is the Lyapunov function we construct for measuring the progress of our
 90 algorithm. Such a construction is motivated by [4], which analyzes bi-level optimization algorithm.
 91 However, our Lyapunov function is different from theirs as it involves the additional optimal gap of
 92 averaged critic and reward estimator, which is necessary for dealing with the decentralized setting.

93 We finish this section by remarking that our convergence results are even new for single agent AC
 94 algorithms under the setting of single-timescale step sizes rule.

95 2 Preliminary

96 In this section, we introduce the problem formulation and the policy gradient theorem, which serves
 97 as the preliminary for the analyzed decentralized AC algorithm.

98 Suppose there are multiple agents aiming to independently optimize a common global objective, and
 99 each agent can communicate with its neighbors through a network. To model the topology, we define
 100 the graph as $\mathcal{G} = (\mathcal{N}, \mathcal{E})$, where \mathcal{N} is the set of nodes with $|\mathcal{N}| = N$ and \mathcal{E} is the set of edges with
 101 $|\mathcal{E}| = E$. In the graph, each node represents an agent, and each edge represents a communication
 102 link. The interaction between agents follows the networked multi-agent MDP.

103 2.1 Markov decision process

104 A networked multi-agent MDP is defined by a tuple $(\mathcal{G}, \mathcal{S}, \{\mathcal{A}^i\}_{i \in \mathcal{N}}, \mathcal{P}, \{r^i\}_{i \in [N]}, \gamma)$. \mathcal{G} denotes the
 105 communication topology (the graph), \mathcal{S} is the finite state space observed by all agents, \mathcal{A}^i represents
 106 the finite action space of agent i . Let $\mathcal{A} := \mathcal{A}^1 \times \dots \times \mathcal{A}^N$ denote the joint action space and
 107 $\mathcal{P}(s'|s, a) : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ denote the transition probability from any state $s \in \mathcal{S}$ to any state
 108 $s' \in \mathcal{S}$ for any joint action $a \in \mathcal{A}$. $r^i : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the local reward function that determines the
 109 reward received by agent i given transition (s, a) ; $\gamma \in [0, 1]$ is the discount factor.

110 For simplicity, we will use $a := [a^1, \dots, a^N]$ to denote the joint action, and $\theta := [\theta^1, \dots, \theta^N] \in$
 111 $\mathbb{R}^{d_\theta \times N}$ to denote joint parameters of all actors, with $\theta^i \in \mathbb{R}^{d_\theta}$. Note that different actors may have
 112 different number of parameters, which is assumed to be the same for our paper without loss of
 113 generality. The MDP goes as follows: For a given state s , each agent make its decision a^i based
 114 on its policy $a^i \sim \pi_{\theta^i}(\cdot|s)$. The state transits to the next state s' based on the joint action of all the
 115 agents: $s' \sim \mathcal{P}(\cdot|s, a)$. Then, each agent will receive its own reward $r^i(s, a)$. For the notation brevity,
 116 we assume that the reward function mapping is deterministic and does not depend on the next state
 117 without loss of generality. The stationary distribution induced by the policy π_θ and the transition
 118 kernel is denoted by $\mu_{\pi_\theta}(s)$.

119 Our objective is to find a set of policies that maximize the accumulated discounted mean reward
 120 received by agents

$$\theta^* = \arg \max_{\theta} J(\theta) := \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k \bar{r}(s_k, a_k) \right]. \quad (1)$$

121 Here, k represents the time step. $\bar{r}(s_k, a_k) := \frac{1}{N} \sum_{i=1}^N r^i(s_k, a_k)$ is the mean reward among agents
 122 at time step k . The randomness of the expectation comes from the initial state distribution $\mu_0(s)$, the
 123 transition kernel \mathcal{P} , and the stochastic policy $\pi_{\theta^i}(\cdot|s)$.

124 2.2 Policy gradient Theorem

125 Under the discounted reward setting, the global state-value function, action-value function, and
 126 advantage function for policy set θ , state s , and action a , are defined as

$$\begin{aligned} V_{\pi_\theta}(s) &:= \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k \bar{r}(s_k, a_k) | s_0 = s \right] \\ Q_{\pi_\theta}(s, a) &:= \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k \bar{r}(s_k, a_k) | s_0 = s, a_0 = a \right] \\ A_{\pi_\theta}(s, a) &:= Q_{\pi_\theta}(s, a) - V_{\pi_\theta}(s). \end{aligned} \quad (2)$$

127 To maximize the objective function defined in (1), the policy gradient [28] can be computed as follow

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{s \sim d_{\pi_{\theta}}, a \sim \pi_{\theta}} \left[\frac{1}{1 - \gamma} A_{\pi_{\theta}}(s, a) \psi_{\pi_{\theta}}(s, a) \right],$$

128 where $d_{\pi_{\theta}}(s) := (1 - \gamma) \sum_{k=0}^{\infty} \gamma^k \mathbb{P}(s_k = s)$ is the discounted state visitation distribution under
129 policy π_{θ} , and $\psi_{\pi_{\theta}}(s, a) := \nabla \log \pi_{\theta}(s, a)$ is the score function.

130 Following the derivation of [42], the policy gradient for each agent under discounted reward setting
131 can be expressed as

$$\nabla_{\theta^i} J(\theta) = \mathbb{E}_{s \sim d_{\pi_{\theta}}, a \sim \pi_{\theta}} \left[\frac{1}{1 - \gamma} A_{\pi_{\theta}}(s, a) \psi_{\pi_{\theta^i}}(s, a^i) \right]. \quad (3)$$

132 3 Decentralized single-timescale actor-critic

Algorithm 1: Decentralized single-timescale AC (reward estimator version)

```

1: Initialize: Actor parameter  $\theta_0$ , critic parameter  $\omega_0$ , reward estimator parameter  $\lambda_0$ , initial state  $s_0$ .
2: for  $k = 0, \dots, K - 1$  do
3:   Option 1: i.i.d. sampling:
4:    $s_k \sim \mu_{\theta_k}(\cdot), a_k \sim \pi_{\theta_k}(\cdot | s_k), s_{k+1} \sim \mathcal{P}(\cdot | s_k, a_k)$ .
5:   Option 2: Markovian sampling:
6:    $a_k \sim \pi_{\theta_k}(\cdot | s_k), s_{k+1} \sim \mathcal{P}(\cdot | s_k, a_k)$ .
7:
8:   Periodical consensus: Compute  $\tilde{\omega}_k^i$  and  $\tilde{\lambda}_k^i$  by (4) and (7).
9:
10:  for  $i = 0, \dots, N$  in parallel do
11:    Reward estimator update: Update  $\lambda_{k+1}^i$  by (8).
12:    Critic update: Update  $\omega_{k+1}^i$  by (5).
13:    Actor update: Update  $\theta_{k+1}^i$  by (6).
14:  end for
15: end for

```

133 We introduce the decentralized single-timescale AC algorithm; see Algorithm 1. In the remaining
134 parts of this section, we will explain the updates in the algorithm in details.

135 In fully-decentralized MARL, each agent can only observe its local reward and action, while trying
136 to maximize the global reward (mean reward) defined in (1). The decentralized AC algorithm solves
137 the problem by performing online updates in an alternative fashion. Specifically, we have N pairs of
138 actor and critic. In order to maximize $J(\theta)$, each critic tries to estimate the *global* state-value function
139 $V_{\pi_{\theta}}(s)$ defined in (2), and each actor then updates its policy parameter based on approximated policy
140 gradient. We now provide more details about the algorithm.

141 **Critics' update.** We will use $\omega^i \in \mathbb{R}^{d_{\omega}}$ to denote the i_{th} critic's parameter and $\bar{\omega} := \frac{1}{N} \sum_{i=1}^N \omega^i$ to
142 represent the averaged parameter of critic. The i_{th} critic approximates the global value function as
143 $V_{\pi_{\theta}}(s) \approx \hat{V}_{\omega^i}(s)$.

144 As we will see, the critic's approximation error can be categorized into two parts, namely, the
145 consensus error $\frac{1}{N} \sum_{i=1}^N \|\omega^i - \bar{\omega}\|$, which measures how close the critics' parameters are; and the
146 approximation error $\|\bar{\omega} - \omega^*(\theta)\|$, which measures the approximation quality of averaged critic.

147 In order for critics to reach consensus, we perform the following update for all critics

$$\tilde{\omega}_k^i = \begin{cases} \sum_{j=1}^N W^{ij} \omega_k^j & \text{if } k \bmod K_c = 0 \\ \omega_k^i & \text{otherwise.} \end{cases} \quad (4)$$

148 where $W \in \mathbb{R}^{n \times n}$ is a weight matrix for communication among agents, whose property will be
149 specified in Assumption 5; K_c denotes the consensus frequency.

150 To reduce the approximation error, we will perform the local TD(0) update [29] as

$$\omega_{k+1}^i = \prod_{R_{\omega}} (\tilde{\omega}_k^i + \beta_k g_c^i(\xi_k, \omega_k^i)), \quad (5)$$

151 where $\xi := (s, a, s')$ represents a transition tuple, $g_c^i(\xi, \omega) := \delta^i(\xi, \omega) \nabla \hat{V}_\omega(s)$ is the update direction,
 152 $\delta^i(\xi, \omega) := r^i(s, a) + \gamma \hat{V}_\omega(s') - \hat{V}_\omega(s)$ is the local temporal difference error (TD-error). β_k is the
 153 step size for critic at iteration k . \prod_{R_ω} projects the parameter into a ball of radius of R_ω containing
 154 the optimal solution, which will be explained when discussing Assumption 1 and 2.

155 **Actors' update.** We will use stochastic gradient ascent to update the policy's parameter, and the
 156 stochastic gradient is calculated based on policy gradient theorem in (3). The advantage function
 157 $A_{\pi_\theta}(s, a)$ can be estimated by

$$\delta(\xi, \theta) := \bar{r}(s, a) + \gamma V(s') - V(s),$$

158 with a sampled from $\pi_\theta(\cdot|s)$. However, to preserve the privacy of each agents, the local reward
 159 cannot be shared to other agents under the fully decentralized setting. Thus, the averaged reward
 160 $\bar{r}(s_k, a_k)$ is not directly attainable. Consequently, we need a strategy to approximate the averaged
 161 reward. In this paper, we will adopt the strategy proposed in [42]. In particular, each agent i will have
 162 a local reward estimator with parameter $\lambda^i \in \mathbb{R}^{d_\lambda}$, which estimates the global averaged reward as
 163 $\bar{r}(s_k, a_k) \approx \hat{r}_{\lambda^i}(s_k, a_k)$.

164 Thus, the update of the i_{th} actor is given by

$$\theta_{k+1}^i = \theta_k^i + \alpha_k \hat{\delta}(\xi_k, \omega_{k+1}^i, \lambda_{k+1}^i) \psi_{\pi_{\theta_k^i}}(s_k, a_k^i), \quad (6)$$

165 where $\hat{\delta}(\xi, \omega, \lambda) := \hat{r}_\lambda(s, a) + \gamma \hat{V}_\omega(s') - \hat{V}_\omega(s)$ is the approximated advantage function. α_k is the
 166 step size for actor's update at iteration k .

167 **Reward estimators' update.** Similar to critic, each reward estimator's approximation error can be
 168 decomposed into consensus error and the approximation error.

169 For each local reward estimator, we perform the consensus step to minimize the consensus error as

$$\tilde{\lambda}_k^i = \begin{cases} \sum_{j=1}^N W^{ij} \lambda_k^j & \text{if } k \bmod K_c = 0 \\ \lambda_k^i & \text{otherwise.} \end{cases} \quad (7)$$

170 To reduce the approximation error, we perform a local update of stochastic gradient descent.

$$\lambda_{k+1}^i = \prod_{R_\lambda} (\tilde{\lambda}_k^i + \eta_k g_r^i(\xi_k, \lambda_k^i)), \quad (8)$$

171 where $g_r^i(\xi, \lambda) := (r^i(s, a) - \hat{r}_\lambda(s, a)) \nabla \hat{r}_\lambda(s, a)$ is the update direction. η_k is the step size for
 172 reward estimator at iteration k . Note the calculation of $g_r^i(\xi, \lambda)$ does not require the knowledge of s' ;
 173 we use ξ in (8) just for notation brevity. Similar to critic's update, \prod_{R_λ} projects the parameter into a
 174 ball of radius of R_λ containing the optimal solution.

175 In our Algorithm 1, we will use the same order for α_k , β_k , and η_k and hence, our algorithm is in
 176 *single-timescale*.

177 **Linear approximation for analysis.** In our analysis, we will use linear approximation for both critic
 178 and reward estimator variables, i.e. $\hat{V}_\omega(s) := \phi(s)^T \omega$; $\hat{r}_\lambda(s, a) := \varphi(s, a)^T \lambda$, where $\phi(s) : \mathcal{S} \rightarrow$
 179 \mathbb{R}^{d_ω} and $\varphi(s, a) : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^{d_\lambda}$ are two feature mappings, whose property will be specified in the
 180 discussion of Assumption 1.

181 **Algorithm for preserving the local action.** Note that in Algorithm 1, the reward estimators need
 182 the knowledge of joint actions in order to estimate the global rewards. To preserve the privacy of
 183 local actions, we further propose a variant of Algorithm 1, which estimates the global rewards by
 184 communicating noisy local rewards; see [6] for the original idea. However, to maintain the optimal
 185 sample complexity, such an approach requires $\mathcal{O}(\log(\varepsilon^{-1}))$ communication rounds for each iteration.
 186 We postpone the detailed design and analysis of such an algorithm scheme into Appendix B.

187 **Remarks on sampling scheme.** The unbiased update for critic and actor variables requires sampling
 188 from μ_{π_θ} and d_{π_θ} , respectively. However, in practical implementations, states are usually collected
 189 from an online trajectory (Markovian sampling), whose distribution is generally different for μ_{π_θ}
 190 and d_{π_θ} . Such a distribution mismatch will inevitably cause biases during the update of critic and
 191 actor variables. One has to bound the corresponding error terms when analyzing the algorithm. In
 192 this work, we will provide the analysis for both sampling schemes.

193 **4 Main Results**

194 In this section, we first introduce the technical assumptions used for our analysis, which are standard
 195 in the literature. Then, we present the convergence results for both actor and critic variables under
 196 i.i.d. sampling and Markovian sampling.

197 **4.1 Assumptions**

198 **Assumption 1** (bounded rewards and feature vectors). *All the local rewards are uniformly bounded,*
 199 *i.e., there exists a positive constants r_{\max} such that $|r^i(s, a)| \leq r_{\max}$, for all feasible (s, a) and*
 200 *$i \in [N]$. The norm of feature vectors are bounded such that for all $s \in \mathcal{S}$, $a \in \mathcal{A}$, $\|\phi(s)\| \leq$*
 201 *1 , $\|\varphi(s, a)\| \leq 1$.*

202 Assumption 1 is standard and commonly adopted; see, e.g., [3, 35, 38, 24, 21]. This assumption can
 203 be achieved via normalizing the feature vectors.

204 **Assumption 2** (negative definiteness of $A_{\theta, \phi}$ and $A_{\theta, \varphi}$). *There exists two positive constants $\lambda_{\phi}, \lambda_{\varphi}$*
 205 *such that for all policy θ , the following two matrices are negative definite*

$$A_{\theta, \phi} := \mathbb{E}_{s \sim \mu_{\theta}(s)} [\phi(s)(\gamma \phi(s')^T - \phi(s)^T)]$$

$$A_{\theta, \varphi} := \mathbb{E}_{s \sim \mu_{\theta}(s), a \sim \pi_{\theta}(\cdot|s)} [-\varphi(s, a)\varphi(s, a)^T],$$

206 *with $\lambda_{\max}(A_{\theta, \phi}) \leq \lambda_{\phi}$, $\lambda_{\max}(A_{\theta, \varphi}) \leq \lambda_{\varphi}$, where $\lambda_{\max}(\cdot)$ represents the largest eigenvalue.*

207 Assumption 2 can be achieved when the matrices $\Phi_{\phi} := [\phi(s_1), \dots, \phi(s_{|\mathcal{S}|})]$ and $\Phi_{\varphi} :=$
 208 $[\varphi(s_1, a_1), \dots, \varphi(s_{|\mathcal{S}|}, a_{|\mathcal{A}|})]$ have full row rank, which ensures that the optimal critic and reward
 209 estimator are unique; see also [24, 34]. Together with Assumption 1, we can show that the norm of
 210 $\omega^*(\theta)$ and $\lambda^*(\theta)$ are bounded by some positive constant, which justifies the projection steps.

211 **Assumption 3** (Lipschitz properties of policy). *There exists constants $C_{\psi}, L_{\psi}, L_{\pi}$ such that for*
 212 *all $\theta, \theta', s \in \mathcal{S}$ and $a \in \mathcal{A}$, we have (1). $|\pi_{\theta}(a|s) - \pi_{\theta'}(a|s)| \leq L_{\pi}\|\theta - \theta'\|$; (2). $\|\psi_{\theta}(s, a) -$*
 213 *$\psi_{\theta'}(s, a)\| \leq L_{\psi}\|\theta - \theta'\|$; (3). $\|\psi_{\theta}(s, a)\| \leq C_{\psi}$.*

214 Assumption 3 is common for analyzing policy-based algorithms; see, e.g., [33, 32, 11]. The assump-
 215 tion ensures the smoothness of objective function $J(\theta)$. It holds for a large range of policy classes
 216 such as tabular softmax policy [1], Gaussian policy [7], and Boltzman policy [13].

217 **Assumption 4** (irreducible and aperiodic Markov chain). *The Markov chain under π_{θ} and transition*
 218 *kernel $\mathcal{P}(\cdot|s, a)$ is irreducible and aperiodic for any θ .*

219 Assumption 4 is a standard assumption, which holds for any uniformly ergodic Markov chains and
 220 any time-homogeneous Markov chains with finite-state space. It ensures that there exists constants
 221 $\kappa > 0$ and $\rho \in (0, 1)$ such that

$$\sup_{s \in \mathcal{S}} d_{TV}(\mathbb{P}(s_k \in \cdot | s_0 = s, \pi_{\theta}), \mu_{\theta}) \leq \kappa \rho^k, \forall k.$$

222 **Assumption 5** (doubly stochastic weight matrix). *The communication matrix W is doubly stochastic,*
 223 *i.e. each column/row sum up to 1. Moreover, the second largest singular value ν is smaller than 1.*

224 Assumption 5 is a common assumption in decentralized optimization and multi-agent reinforcement
 225 learning; see, e.g., [27, 5, 6]. It ensures the convergence of consensus error for critic and reward
 226 estimator variables.

227 **4.2 Sample complexity under i.i.d. sampling**

228 **Theorem 1** (sample complexity under i.i.d. sampling). *Suppose Assumptions 1-5 hold. Consider*
 229 *the update of Algorithm 1 under i.i.d. sampling. Let $\alpha_k = \frac{\bar{\alpha}}{\sqrt{K}}$ for some positive constant $\bar{\alpha}$,*

230 *$\beta_k = \frac{C_9}{2\lambda_{\phi}}\alpha_k$, and $\eta_k = \frac{C_{10}}{2\lambda_{\varphi}}\alpha_k$, $K_c \leq \mathcal{O}(\alpha_k^{-\frac{1}{2}})$, where K denotes the total number of iterations.*
 231 *Then, we have*

$$\frac{1}{K} \sum_{k=1}^K \sum_{i=1}^N \mathbb{E} [\|\omega_k^i - \omega^*(\theta_k)\|^2] \leq \mathcal{O}\left(\frac{1}{\sqrt{K}}\right)$$

$$\frac{1}{K} \sum_{k=1}^K \sum_{i=1}^N \mathbb{E} [\|\nabla_{\theta^i} F(\theta_k)\|^2] \leq \mathcal{O}\left(\frac{1}{\sqrt{K}}\right) + \mathcal{O}(\varepsilon_{app} + \varepsilon_{sp}), \quad (9)$$

232 where C_9, C_{10} are positive constants defined in the proof.

233 The proof of Theorem 1 can be found in Appendix E.1. It establishes the iteration complexity of
 234 $\mathcal{O}(1/\sqrt{K})$, or equivalently, sample complexity of $\mathcal{O}(\varepsilon^{-2})$ for Algorithm 1. Note that actors, critics,
 235 and reward estimators use the step sizes of the same order. The sample complexity matches the
 236 optimal rate of SGD for general non-convex optimization problem. To explain the errors in (9), let us
 237 define the approximation error as the following:

$$\varepsilon_{app} := \max_{\theta, a} \sqrt{\mathbb{E}_{s \sim \mu_\theta} \left[|V_{\pi_\theta}(s) - \hat{V}_{\omega^*(\theta)}(s)|^2 + |\bar{r}(s, a) - \hat{r}_{\lambda^*(\theta)}(s, a)|^2 \right]}.$$

238 The error ε_{app} captures the approximation power of critic and reward estimator. Similar terms
 239 also appear in the literature (see e.g., [35, 1, 21]). Such an approximation error becomes zero in
 240 tabular case. The error ε_{sp} is inevitably caused by the mismatch between discounted state visitation
 241 distribution d_{π_θ} and stationary distribution μ_{π_θ} ; see, e.g., [38, 24]. It is defined as

$$\varepsilon_{sp} := 2C_\theta (\log_\rho \kappa^{-1} + \frac{1}{\rho})(1 - \gamma).$$

242 When γ is close to 1, the error becomes small. This is because d_{π_θ} approaches to μ_{π_θ} when γ goes to
 243 1. In the literature, some works assume that sampling from d_{π_θ} is permitted, thus eliminate this error;
 244 see, e.g., [4].

245 4.3 Sample complexity under Markovian sampling

246 **Theorem 2** (sample complexity under Markovian sampling). *Suppose Assumptions 1-5 hold. Con-*
 247 *sider the update of Algorithm 1 under Markovian sampling. Let $\alpha_k = \frac{\bar{\alpha}}{\sqrt{K}}$ for some positive constant*

248 *$\bar{\alpha}, \beta_k = \frac{C_9}{2\lambda_\psi} \alpha_k$, and $\eta_k = \frac{C_{10}}{2\lambda_\psi} \alpha_k$, $K_c \leq \mathcal{O}(\alpha_k^{-\frac{1}{2}})$, where K is the total number of iterations. Then,*
 249 *we have*

$$\begin{aligned} \frac{1}{K} \sum_{k=1}^K \sum_{i=1}^N \mathbb{E} [\|\omega_k^i - \omega^*(\theta_k)\|^2] &\leq \mathcal{O}\left(\frac{\log^2 K}{\sqrt{K}}\right) \\ \frac{1}{K} \sum_{k=1}^K \sum_{i=1}^N \mathbb{E} [\|\nabla_{\theta^i} F(\theta_k)\|^2] &\leq \mathcal{O}\left(\frac{\log^2 K}{\sqrt{K}}\right) + \mathcal{O}(\varepsilon_{app} + \varepsilon_{sp}), \end{aligned} \quad (10)$$

250 where C_9, C_{10} are positive constants defined in proof.

251 We put the proof of Theorem 2 in Appendix E.2. In Markovian sampling, the updates are biased for
 252 critics, actors, and reward estimators. The error will decrease as the Markov chain mixes, and the
 253 logarithmic term is due to the cost for mixing.

254 Theorem 2 establishes the iteration complexity of $\mathcal{O}(\log^2 K/\sqrt{K})$, or equivalently, sample complex-
 255 ity of $\tilde{\mathcal{O}}(\varepsilon^{-2})$ for Algorithm 1. It matches the state-of-the-art sample complexity of decentralized AC
 256 algorithms, which are implemented in double-loop fashion [11, 6].

257 4.4 Proof sketch

258 We present the main elements for the proof of Theorem 2, which helps in understanding the difference
 259 between classical two-timescale/double-loop analysis and our single-timescale analysis. The proof of
 260 Theorem 1 follows the same framework with simpler sampling scheme.

261 Under Markovian sampling, it is possible to show the following inequality, which characterizes the
 262 ascent of the objective.

$$\begin{aligned} \mathbb{E}[J(\theta_{k+1})] - J(\theta_k) &\geq \sum_{i=1}^N \left[\frac{\alpha_k}{2} \mathbb{E} \|\nabla_{\theta^i} J(\theta_k)\|^2 + \frac{\alpha_k}{2} \mathbb{E} \|g_a^i(\xi_k, \omega_{k+1}^i, \lambda_{k+1}^i)\|^2 \right. \\ &\quad \left. - 8C_\psi^2 \alpha_k \mathbb{E} \|\omega^*(\theta_k) - \omega_{k+1}^i\|^2 - 4C_\psi^2 \alpha_k \mathbb{E} \|\lambda^*(\theta_k) - \lambda_{k+1}^i\|^2 \right] \\ &\quad - \mathcal{O}(\log^2(K) \alpha_k^2) - \mathcal{O}((\varepsilon_{app} + \varepsilon_{sp}) \alpha_k). \end{aligned} \quad (11)$$

263 To analyze the errors of critic $\|\omega^*(\theta_k) - \omega_{k+1}^i\|^2$ and reward estimator $\|\lambda^*(\theta_k) - \lambda_{k+1}^i\|^2$, the two-
264 timescale analysis requires $\mathcal{O}(\alpha_k) < \min\{\mathcal{O}(\beta_k), \mathcal{O}(\eta_k)\}$ in order for these two errors to converge.
265 The double-loop approach runs lower-level update for $\mathcal{O}(\log(\varepsilon^{-1}))$ times with batch size $\mathcal{O}(\varepsilon^{-1})$
266 to drive these errors below ε and hence, they cannot allow inner loop size and bath size to be $\mathcal{O}(1)$
267 simultaneously. To obtain the convergence result for *single-timescale* update, the idea is to further
268 upper bound these two lower-level errors by the quantity $\mathcal{O}(\alpha_k \mathbb{E}\|g_a^i(\xi_k, \omega_{k+1}^i, \lambda_{k+1}^i)\|^2)$ (through a
269 series of derivations), and then eliminate these errors by the ascent term $\frac{\alpha_k}{2} \mathbb{E}\|g_a^i(\xi_k, \omega_{k+1}^i, \lambda_{k+1}^i)\|^2$.
270 We mainly focus on the analysis of critic's error through the proof sketch. The analysis for reward
271 estimator's error follows similar procedure. We start by decomposing the error of critic as

$$\sum_{i=1}^N \|\omega_{k+1}^i - \omega^*(\theta_k)\|^2 = \sum_{i=1}^N (\|\omega_{k+1}^i - \bar{\omega}_{k+1}\|^2 + \|\bar{\omega}_{k+1} - \omega^*(\theta_k)\|^2). \quad (12)$$

272 The first term represents the consensus error, which can be bounded by the next lemma.

273 **Lemma 1.** *Suppose Assumptions 1 and 5 hold. Consider the sequence $\{\omega_k^i\}$ generated by Algorithm 1,*
274 *then the following holds*

$$\|Q\omega_{k+1}\| \leq \nu^{\frac{k'}{K_c}} \|\omega_0\| + 4 \sum_{t=0}^k \nu^{\lceil \frac{k'-1-t}{K_c} \rceil} \beta_t \sqrt{N} C_\delta,$$

275 where $\omega_0 := [\omega^1, \dots, \omega^N]^T$, $Q := I - \frac{1}{N} \mathbf{1}\mathbf{1}^T$, $k' := \lfloor \frac{k}{K_c} \rfloor * K_c$. The constant $\nu \in (0, 1)$ is the
276 second largest singular value of W .

277 Based on Lemma 1 and follow the step size rule of Theorem 2, it is possible to show $\|Q\omega_{k+1}\|_F^2 =$
278 $\sum_{i=1}^N \|\omega_{k+1}^i - \bar{\omega}_{k+1}\|^2 = \mathcal{O}(K_c^2 \beta_k^2)$. Let $K_c = \mathcal{O}(\beta_k^{-\frac{1}{2}})$, we have $\|Q\omega_{k+1}\|_F^2 = \mathcal{O}(\beta_k)$, which
279 maintains the optimal rate.

280 To analyze the second term in (12), we first construct the following Lyapunov function

$$\mathbb{V}_k := -J(\theta_k) + \|\bar{\omega}_k - \omega^*(\theta_k)\|^2 + \|\bar{\lambda}_k - \lambda^*(\theta_k)\|^2. \quad (13)$$

281 Then, it remains to derive an approximate descent property of the term $\|\bar{\omega}_k - \omega^*(\theta_k)\|^2$ in (13).
282 Towards that end, our key step lies in establishing the *smoothness of the optimal critic variables*
283 shown in the next lemma.

284 **Lemma 2** (smoothness of optimal critic). *Suppose Assumptions 1-3 hold, under the update of*
285 *Algorithm 1, there exists a positive constant $L_{\mu,1}$ such that for all θ, θ' , it holds that*

$$\|\nabla \omega^*(\theta) - \nabla \omega^*(\theta')\| \leq L_{\mu,1} \|\theta - \theta'\|,$$

286 where $\nabla \omega^*(\theta)$ denotes the Jacobian of $\omega^*(\theta)$ with respect to θ .

287 This smoothness property is essential for achieving our $\tilde{\mathcal{O}}(1/\sqrt{K})$ convergence rate.

288 To the best of our knowledge, the smoothness of $\omega^*(\theta)$ has not been justified in the literature.
289 Equipped with Lemma 2, we are able to establish the following lemma.

290 **Lemma 3** (Error of critic). *Under Assumptions 1-5, consider the update of Algorithm 1. Then, it*
291 *holds that*

$$\begin{aligned} \mathbb{E}[\|\bar{\omega}_{k+1} - \omega^*(\theta_{k+1})\|^2] &\leq (1 + C_9 \alpha_k) \|\bar{\omega}_{k+1} - \omega^*(\theta_k)\|^2 \\ &\quad + \frac{\alpha_k}{4} \sum_{i=1}^N \|\mathbb{E}[g_a^i(\xi_k, \omega_{k+1}^i, \lambda_{k+1}^i)]\|^2 + \mathcal{O}(\alpha_k^2). \end{aligned} \quad (14)$$

$$\begin{aligned} \mathbb{E}[\|\bar{\omega}_{k+1} - \omega^*(\theta_k)\|^2] &\leq (1 - 2\lambda_\phi \beta_k) \|\bar{\omega}_k - \omega^*(\theta_k)\|^2 \\ &\quad + C_{K_1} \beta_k \beta_{k-Z_K} + C_{K_2} \alpha_{k-Z_K} \beta_k. \end{aligned} \quad (15)$$

292 Here, $Z_K := \min\{z \in \mathbb{N}^+ | \kappa \rho^{z-1} \leq \min\{\alpha_k, \beta_k, \eta_k\}\}$, C_9, λ_ϕ are constants specified in appendix,
293 and C_{K_1} and C_{K_2} are of order $\mathcal{O}(\log(K))$ and $\mathcal{O}(\log^2(K))d$ respectively.

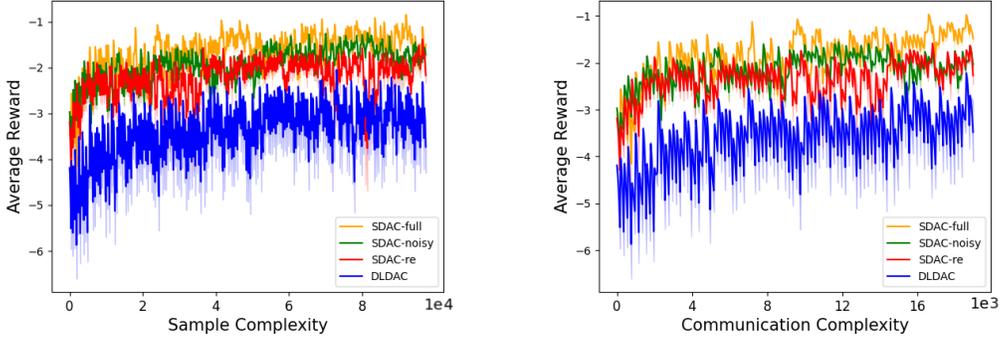


Figure 1: Averaged reward versus sample complexity and communication complexity. The vertical axis is the averaged reward over all the agents.

294 Plug (15) into (14), we can establish the approximate descent property of $\|\bar{\omega}_k - \omega^*(\theta_k)\|^2$ in (13):

$$\begin{aligned}
 \mathbb{E}[\|\bar{\omega}_{k+1} - \omega^*(\theta_{k+1})\|^2] &\leq (1 + C_9\alpha_k)(1 - 2\lambda_\phi\beta_k)\|\bar{\omega}_k - \omega^*(\theta_k)\|^2 \\
 &\quad + \frac{\alpha_k}{4} \sum_{i=1}^N \|\mathbb{E}[g_a^i(\xi_k, \omega_{k+1}^i, \lambda_{k+1}^i)]\|^2 \\
 &\quad + \mathcal{O}(C_{K_1}\beta_k\beta_{k-Z_K} + C_{K_2}\alpha_{k-Z_K}\beta_k). \tag{16}
 \end{aligned}$$

295 Finally, plugging (11), (14), and (16) into (13) gives the ascent of the Lyapunov function, which leads
 296 to our convergence result through steps of standard arguments.

297 5 Numerical results

298 In this section, our objective is to illustrate the empirical sample complexity and communication
 299 complexity of the proposed algorithms. We also implement the algorithm in [6] to serve as a baseline,
 300 which employs double-loop algorithmic framework. Our simulation is based on the grounded
 301 communication environment proposed in [19]; see Appendix A for detailed set up. Through the
 302 discussion, we refer the algorithm in [6] as "DLDAC", the Algorithm 1 as "SDAC-re", the Algorithm 2
 303 as "SDAC-noisy" (see Appendix B). We also provide the result which assumes full reward is available
 304 to serve as baseline, which we refer as "SDAC-full". We set $K_r = 5$ for "SDAC-noisy"; $K_c = 1$
 305 for "SDAC-re", "SDAC-noisy", and "SDAC-full". We choose $T_c = 5$ (loop size), $T'_c = 1$ (critic
 306 consensus number every iteration), $T' = 5$ (reward consensus number every iteration) for "DLDAC".

307 The sample complexity and communication complexity are shown in Figure 1. The results are
 308 averaged over 10 Monte Carlo runs. As we can see, the proposed two algorithms achieve significantly
 309 higher reward than "DLDAC" in terms of both sample complexity and communication complexity.
 310 Moreover, their performances approach the baseline "SDAC-full", where the global reward is assumed
 311 to be available, indicating that the reward approximation is nearly accurate. Due to space limit, we
 312 will put additional experiments on the comparison with existing decentralized AC algorithms and the
 313 ablation study of hyper-parameters to Appendix A.

314 6 Conclusion and future direction

315 In this paper, we studied the convergence of fully decentralized AC algorithm under practical single-
 316 timescale update for the first time. We designed such an algorithm which maintains the optimal
 317 sample complexity of $\tilde{\mathcal{O}}(\varepsilon^{-2})$ under less communications. We also proposed a variant to preserve the
 318 privacy of local actions by communicating noisy rewards. Extensive simulation results demonstrate
 319 the superiority of our algorithms' empirical performance over existing decentralized AC algorithms.
 320 One limitation of our work is that we only study the convergence to stationary point. Thus, we leave
 321 the research on the avoidance of saddle points and convergence to global optimum as promising
 322 future directions.

323 **References**

- 324 [1] A. Agarwal, S. M. Kakade, J. D. Lee, and G. Mahajan. Optimality and approximation with
325 policy gradient methods in markov decision processes. In *Conference on Learning Theory*
326 (*COLT*), pages 64–66, 2020.
- 327 [2] J. Baxter and P. L. Bartlett. Infinite-horizon policy-gradient estimation. *Journal of Artificial*
328 *Intelligence Research*, 15:319–350, 2001.
- 329 [3] J. Bhandari, D. Russo, and R. Singal. A finite time analysis of temporal difference learning with
330 linear function approximation. In *Conference on Learning Theory (COLT)*, pages 1691–1692,
331 2018.
- 332 [4] T. Chen, Y. Sun, and W. Yin. Closing the gap: Tighter analysis of alternating stochastic gradient
333 methods for bilevel problems. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan,
334 editors, *Advances in Neural Information Processing Systems*, 2021.
- 335 [5] Z. Chen, Y. Zhou, and R. Chen. Multi-agent off-policy td learning: Finite-time analysis with near-
336 optimal sample complexity and communication complexity. *arXiv preprint arXiv:2103.13147*,
337 2021.
- 338 [6] Z. Chen, Y. Zhou, R.-R. Chen, and S. Zou. Sample and communication-efficient decentralized
339 actor-critic algorithms with finite-time analysis. In *International Conference on Machine*
340 *Learning*, pages 3794–3834. PMLR, 2022.
- 341 [7] K. Doya. Reinforcement learning in continuous time and space. *Neural Computation*, 12(1):219–
342 245, 2000.
- 343 [8] L. Espeholt, H. Soyer, R. Munos, K. Simonyan, V. Mnih, T. Ward, Y. Doron, V. Firoiu, T. Harley,
344 I. Dunning, et al. Impala: Scalable distributed deep-rl with importance weighted actor-learner
345 architectures. In *International Conference on Machine Learning*, pages 1407–1416. PMLR,
346 2018.
- 347 [9] Z. Fu, Z. Yang, and Z. Wang. Single-timescale actor-critic provably finds globally optimal
348 policy. In *International Conference on Learning Representations*, 2021.
- 349 [10] H. Guo, Z. Fu, Z. Yang, and Z. Wang. Decentralized single-timescale actor-critic on zero-
350 sum two-player stochastic games. In *International Conference on Machine Learning*, pages
351 3899–3909. PMLR, 2021.
- 352 [11] F. Hairi, J. Liu, and S. Lu. Finite-time convergence and sample complexity of multi-agent actor-
353 critic reinforcement learning with average reward. In *International Conference on Learning*
354 *Representations*, 2022.
- 355 [12] S. M. Kakade. A natural policy gradient. In *Proc. Advances in Neural Information Processing*
356 *Systems (NIPS)*, pages 1531–1538, 2002.
- 357 [13] V. R. Konda and V. S. Borkar. Actor-critic-type learning algorithms for Markov decision
358 processes. *SIAM Journal on Control and Optimization*, 38(1):94–123, 1999.
- 359 [14] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra.
360 Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- 361 [15] Y. Lin, K. Zhang, Z. Yang, Z. Wang, T. Başar, R. Sandhu, and J. Liu. A communication-efficient
362 multi-agent actor-critic algorithm for distributed reinforcement learning. In *2019 IEEE 58th*
363 *Conference on Decision and Control (CDC)*, pages 5562–5567, 2019.
- 364 [16] M. L. Littman. Markov games as a framework for multi-agent reinforcement learning. In
365 *Machine learning proceedings 1994*, pages 157–163. Elsevier, 1994.
- 366 [17] Y. Liu, K. Zhang, T. Basar, and W. Yin. An improved analysis of (variance-reduced) policy
367 gradient and natural policy gradient methods. *Advances in Neural Information Processing*
368 *Systems*, 33:7624–7636, 2020.

- 369 [18] R. Lowe, Y. I. Wu, A. Tamar, J. Harb, O. Pieter Abbeel, and I. Mordatch. Multi-agent actor-critic
370 for mixed cooperative-competitive environments. *Advances in neural information processing*
371 *systems*, 30, 2017.
- 372 [19] I. Mordatch and P. Abbeel. Emergence of grounded compositional language in multi-agent
373 populations. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- 374 [20] S. Omidshafiei, J. Papis, C. Amato, J. P. How, and J. Vian. Deep decentralized multi-task
375 multi-agent reinforcement learning under partial observability. In *International Conference on*
376 *Machine Learning*, pages 2681–2690. PMLR, 2017.
- 377 [21] S. Qiu, Z. Yang, J. Ye, and Z. Wang. On the finite-time convergence of actor-critic algorithm.
378 In *Optimization Foundations for Reinforcement Learning Workshop at Advances in Neural*
379 *Information Processing Systems (NeurIPS)*, 2019.
- 380 [22] T. Rashid, M. Samvelyan, C. Schroeder, G. Farquhar, J. Foerster, and S. Whiteson. Qmix: Mono-
381 tonic value function factorisation for deep multi-agent reinforcement learning. In *International*
382 *Conference on Machine Learning*, pages 4295–4304. PMLR, 2018.
- 383 [23] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization
384 algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 385 [24] H. Shen, K. Zhang, M. Hong, and T. Chen. Asynchronous advantage actor critic: Non-
386 asymptotic analysis and linear speedup. *ArXiv:2012.15511*, 2020.
- 387 [25] D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker,
388 M. Lai, A. Bolton, et al. Mastering the game of go without human knowledge. *nature*,
389 550(7676):354–359, 2017.
- 390 [26] K. Son, D. Kim, W. J. Kang, D. E. Hostallero, and Y. Yi. Qtran: Learning to factorize with
391 transformation for cooperative multi-agent reinforcement learning. In *International Conference*
392 *on Machine Learning*, pages 5887–5896. PMLR, 2019.
- 393 [27] J. Sun, G. Wang, G. B. Giannakis, Q. Yang, and Z. Yang. Finite-sample analysis of decentral-
394 ized temporal-difference learning with linear function approximation. In *Proc. International*
395 *Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 4485–4495, 2020.
- 396 [28] R. S. Sutton, D. A. McAllester, S. P. Singh, and Y. Mansour. Policy gradient methods for
397 reinforcement learning with function approximation. In *Proc. Advances in Neural Information*
398 *Processing Systems (NIPS)*, pages 1057–1063, 2000.
- 399 [29] J. N. Tsitsiklis and B. Van Roy. Analysis of temporal-difference learning with function
400 approximation. In *Advances in neural information processing systems (NIPS)*, pages 1075–
401 1081, 1997.
- 402 [30] O. Vinyals, I. Babuschkin, W. M. Czarnecki, M. Mathieu, A. Dudzik, J. Chung, D. H. Choi,
403 R. Powell, T. Ewalds, P. Georgiev, et al. Grandmaster level in starcraft ii using multi-agent
404 reinforcement learning. *Nature*, 575(7782):350–354, 2019.
- 405 [31] Y. Wang, S. Zou, and Y. Zhou. Non-asymptotic analysis for two time-scale TDC with general
406 smooth function approximation. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and
407 J. W. Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages
408 9747–9758. Curran Associates, Inc., 2021.
- 409 [32] Y. F. Wu, W. Zhang, P. Xu, and Q. Gu. A finite-time analysis of two time-scale actor-critic
410 methods. *Advances in Neural Information Processing Systems*, 33:17617–17628, 2020.
- 411 [33] P. Xu, F. Gao, and Q. Gu. An improved convergence analysis of stochastic variance-reduced
412 policy gradient. In *Proc. International Conference on Uncertainty in Artificial Intelligence*
413 *(UAI)*, 2019.
- 414 [34] T. Xu and Y. Liang. Sample complexity bounds for two timescale value-based reinforcement
415 learning algorithms. In *International Conference on Artificial Intelligence and Statistics*, pages
416 811–819. PMLR, 2021.

- 417 [35] T. Xu, Z. Wang, and Y. Liang. Improving sample complexity bounds for (natural) actor-critic
418 algorithms. In *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, volume 33,
419 2020.
- 420 [36] T. Xu, Z. Yang, Z. Wang, and Y. Liang. Doubly robust off-policy actor-critic: Convergence and
421 optimality. *ArXiv:2102.11866*, 2021.
- 422 [37] C. Yu, X. Wang, X. Xu, M. Zhang, H. Ge, J. Ren, L. Sun, B. Chen, and G. Tan. Distributed
423 multiagent coordinated learning for autonomous driving in highways based on dynamic co-
424 ordination graphs. *IEEE Transactions on Intelligent Transportation Systems*, 21(2):735–748,
425 2019.
- 426 [38] S. Zeng, T. Chen, A. Garcia, and M. Hong. Learning to coordinate in multi-agent systems: A
427 coordinated actor-critic algorithm and finite-time guarantees. *arXiv preprint arXiv:2110.05597*,
428 2021.
- 429 [39] H. Zhang, W. Chen, Z. Huang, M. Li, Y. Yang, W. Zhang, and J. Wang. Bi-level actor-critic
430 for multi-agent coordination. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
431 volume 34, pages 7325–7332, 2020.
- 432 [40] K. Zhang, A. Koppel, H. Zhu, and T. Başar. Global convergence of policy gradient methods to
433 (almost) locally optimal policies. *arXiv preprint arXiv:1906.08383*, 2019.
- 434 [41] K. Zhang, Z. Yang, and T. Başar. Multi-agent reinforcement learning: A selective overview of
435 theories and algorithms. *Handbook of Reinforcement Learning and Control*, pages 321–384,
436 2021.
- 437 [42] K. Zhang, Z. Yang, H. Liu, T. Zhang, and T. Basar. Fully decentralized multi-agent reinforcement
438 learning with networked agents. In *International Conference on Machine Learning*, pages
439 5872–5881. PMLR, 2018.

440 Checklist

- 441 1. For all authors...
- 442 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
443 contributions and scope? [Yes]
- 444 (b) Did you describe the limitations of your work? [Yes] The limitation is written in an
445 equivalent form as future works in the conclusion section; see Section 6.
- 446 (c) Did you discuss any potential negative societal impacts of your work? [N/A] We
447 conduct research about the design and analysis of the fundamental actor-critic algorithm,
448 which should not bring any negative societal impact.
- 449 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
450 them? [Yes]
- 451 2. If you are including theoretical results...
- 452 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 453 (b) Did you include complete proofs of all theoretical results? [Yes]
- 454 3. If you ran experiments...
- 455 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
456 mental results (either in the supplemental material or as a URL)? [Yes]
- 457 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
458 were chosen)? [Yes]
- 459 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
460 ments multiple times)? [Yes]
- 461 (d) Did you include the total amount of compute and the type of resources used (e.g., type
462 of GPUs, internal cluster, or cloud provider)? [Yes]
- 463 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 464 (a) If your work uses existing assets, did you cite the creators? [Yes]

- 465 (b) Did you mention the license of the assets? [Yes]
466 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
467
468 (d) Did you discuss whether and how consent was obtained from people whose data you're
469 using/curating? [Yes]
470 (e) Did you discuss whether the data you are using/curating contains personally identifiable
471 information or offensive content? [N/A]
472 5. If you used crowdsourcing or conducted research with human subjects...
473 (a) Did you include the full text of instructions given to participants and screenshots, if
474 applicable? [N/A]
475 (b) Did you describe any potential participant risks, with links to Institutional Review
476 Board (IRB) approvals, if applicable? [N/A]
477 (c) Did you include the estimated hourly wage paid to participants and the total amount
478 spent on participant compensation? [N/A]