

# Supplementary Materials for “Fully Decentralized Multi-Agent Reinforcement Learning with Networked Agents”

## A. Pseudocode of Algorithm 1 and Algorithm 2

In this section, we present the pseudocode of Algorithm 1 and Algorithm 2. Algorithm 1 follows the updates (3.3)-(3.5), which is based on the estimate of the local advantage function  $A_\theta^i$ . This can be achieved by maintaining a consensual approximation of the global action-value function  $Q_\theta$  at each agent.

---

**Algorithm 1** The networked actor-critic algorithm based on action-value function approximation

---

**Input:** Initial values of the parameters  $\mu_0^i, \omega_0^i, \tilde{\omega}_0^i, \theta_0^i, \forall i \in \mathcal{N}$ ; the initial state  $s_0$  of the MDP, and stepsizes  $\{\beta_{\omega,t}\}_{t \geq 0}$  and  $\{\beta_{\theta,t}\}_{t \geq 0}$ .

Each agent  $i \in \mathcal{N}$  executes action  $a_0^i \sim \pi_{\theta_0^i}^i(s_0, \cdot)$  and **observes joint actions**  $a_0 = (a_0^1, \dots, a_0^N)$ .

Initialize the iteration counter  $t \leftarrow 0$ .

**Repeat:**

**for all**  $i \in \mathcal{N}$  **do**

    Observe state  $s_{t+1}$ , and reward  $r_{t+1}^i$ .

    Update  $\mu_{t+1}^i \leftarrow (1 - \beta_{\omega,t}) \cdot \mu_t^i + \beta_{\omega,t} \cdot r_{t+1}^i$ .

    Select and execute action  $a_{t+1}^i \sim \pi_{\theta_t^i}^i(s_{t+1}, \cdot)$ .

**end for**

  Observe joint actions  $a_{t+1} = (a_{t+1}^1, \dots, a_{t+1}^N)$ .

**for all**  $i \in \mathcal{N}$  **do**

    Update  $\delta_t^i \leftarrow r_{t+1}^i - \mu_t^i + Q_{t+1}(\omega_t^i) - Q_t(\omega_t^i)$ .

**Critic step:**  $\tilde{\omega}_t^i \leftarrow \omega_t^i + \beta_{\omega,t} \cdot \delta_t^i \cdot \nabla_{\omega} Q_t(\omega_t^i)$ .

    Update  $A_t^i \leftarrow Q_t(\omega_t^i) - \sum_{a^i \in \mathcal{A}^i} \pi_{\theta_t^i}^i(s_t, a^i) \cdot Q(s_t, a^i, a^{-i}; \omega_t^i)$ ,  $\psi_t^i \leftarrow \nabla_{\theta^i} \log \pi_{\theta_t^i}^i(s_t, a_t^i)$ .

**Actor step:**  $\theta_{t+1}^i \leftarrow \theta_t^i + \beta_{\theta,t} \cdot A_t^i \cdot \psi_t^i$ .

    Send  $\tilde{\omega}_t^i$  to the neighbors  $\{j \in \mathcal{N} : (i, j) \in \mathcal{E}_t\}$  over the communication network  $\mathcal{G}_t$ .

**end for**

**for all**  $i \in \mathcal{N}$  **do**

**Consensus step:**  $\omega_{t+1}^i \leftarrow \sum_{j \in \mathcal{N}} c_t(i, j) \cdot \tilde{\omega}_t^j$ .

**end for**

  Update the iteration counter  $t \leftarrow t + 1$ .

**Until Convergence**

---

In addition, Algorithm 2 follows the updates (3.8), (3.11), and (3.12), which is based on the estimate of the global advantage function  $A_\theta$ . This can be achieved by maintaining consensual approximation of both the global state-value function  $V_\theta$  and the globally averaged reward function  $\bar{R}$ , at each agent.

---

**Algorithm 2** The networked actor-critic algorithm based on state-value function approximation

---

**Input:** Initial values of  $\mu_0^i, \tilde{\mu}_0^i, v_0^i, \tilde{v}_0^i, \lambda_0^i, \tilde{\lambda}_0^i, \theta_0^i, \forall i \in \mathcal{N}$ ; the initial state  $s_0$  of the MDP, and stepsizes  $\{\beta_{v,t}\}_{t \geq 0}$  and  $\{\beta_{\theta,t}\}_{t \geq 0}$ .

Each agent  $i$  implements  $a_0^i \sim \pi_{\theta_0^i}(s_0, \cdot)$ .

Initialize the step counter  $t \leftarrow 0$ .

**Repeat:**

**for all**  $i \in \mathcal{N}$  **do**

    Observe state  $s_{t+1}$ , and reward  $r_{t+1}^i$ .

    Update  $\tilde{\mu}_t^i \leftarrow (1 - \beta_{v,t}) \cdot \mu_t^i + \beta_{v,t} \cdot r_{t+1}^i, \quad \tilde{\lambda}_t^i \leftarrow \lambda_t^i + \beta_{v,t} \cdot [r_{t+1}^i - \bar{R}_t(\lambda_t^i)] \cdot \nabla_{\lambda} \bar{R}_t(\lambda_t^i)$ .

    Update  $\delta_t^i \leftarrow r_{t+1}^i - \mu_t^i + V_{t+1}(v_t^i) - V_t(v_t^i)$

**Critic step:**  $\tilde{v}_t^i \leftarrow v_t^i + \beta_{v,t} \cdot \delta_t^i \cdot \nabla_v V_t(v_t^i)$ .

    Update  $\tilde{\delta}_t^i \leftarrow \bar{R}_t(\lambda_t^i) - \mu_t^i + V_{t+1}(v_t^i) - V_t(v_t^i), \quad \psi_t^i \leftarrow \nabla_{\theta^i} \log \pi_{\theta_t^i}(s_t, a_t^i)$ .

**Actor step:**  $\theta_{t+1}^i = \theta_t^i + \beta_{\theta,t} \cdot \tilde{\delta}_t^i \cdot \psi_t^i$ .

    Send  $\tilde{\mu}_t^i, \tilde{\lambda}_t^i, \tilde{v}_t^i$  to the neighbors over  $\mathcal{G}_t$ .

**end for**

**for all**  $i \in \mathcal{N}$  **do**

**Consensus step:**  $\mu_{t+1}^i \leftarrow \sum_{j \in \mathcal{N}} c_t(i, j) \cdot \tilde{\mu}_t^j, \lambda_{t+1}^i \leftarrow \sum_{j \in \mathcal{N}} c_t(i, j) \cdot \tilde{\lambda}_t^j, v_{t+1}^i \leftarrow \sum_{j \in \mathcal{N}} c_t(i, j) \cdot \tilde{v}_t^j$ .

**end for**

  Update the iteration counter  $t \leftarrow t + 1$ .

**Until Convergence**

---

## B. Proofs of the Main Results

In this section, we provide the proofs of the theoretical results presented in the paper. We first prove the policy gradient theorem for MARL, and then provide proofs for the convergence results in §4. Detailed proofs for the convergence of Algorithm 1 are given, followed by succinct proofs for the convergence of Algorithm 2 that draw parallels with the first ones.

**Notation.** For any vector  $x \in \mathbb{R}^n$  and matrix  $Y \in \mathbb{R}^{m \times n}$ , we use  $\|x\|$  and  $\|Y\|$  to denote the Euclidean norm of  $x$  and the induced 2-norm of  $Y$ , respectively. We also use  $\|x\|_\infty$  and  $\|Y\|_\infty$  to denote the infinite norm of  $x$  and the induced infinite-norm of  $Y$ , respectively.

### B.1. Proof of Theorem 3.1

The proof of this theorem follows the proof of the policy gradient theorem in single-agent reinforcement learning (Sutton et al., 2000), which shows that

$$\begin{aligned} \nabla_\theta J(\theta) &= \mathbb{E}_{s \sim d_\theta, a \sim \pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) \cdot Q_\theta(s, a)] \\ &= \sum_{s \in \mathcal{S}, a \in \mathcal{A}} d_\theta(s) \pi_\theta(s, a) \left[ \nabla_\theta \sum_{i \in \mathcal{N}} \log \pi_{\theta^i}^i(s, a^i) \right] \cdot Q_\theta(s, a), \end{aligned} \quad (\text{B.1})$$

where  $Q_\theta$  is the action-value function defined in (2.3), and  $d_\theta$  denotes the stationary distribution of the Markov chain induced by policy  $\pi_\theta$ . Here the second equality in (B.1) holds because  $\pi_\theta$  is the product of local policy functions. Hence, the gradient with respect to the parameter  $\theta^i$  becomes

$$\nabla_{\theta^i} J(\theta) = \sum_{s \in \mathcal{S}, a \in \mathcal{A}} d_\theta(s) \pi_\theta(s, a) \cdot \nabla_{\theta^i} \log \pi_{\theta^i}^i(s, a^i) \cdot Q_\theta(s, a), \quad (\text{B.2})$$

which proves the first equality in (3.2). Moreover, since  $\sum_{a^i \in \mathcal{A}^i} \pi_{\theta^i}^i(s, a^i) = 1$ , we have  $\nabla_{\theta^i} \left[ \sum_{a^i \in \mathcal{A}^i} \pi_{\theta^i}^i(s, a^i) \right] = 0$ . To simplify the notation, for each  $i \in \mathcal{N}$ , we define  $a^{-i}$  as the joint actions of all agents except  $i$ , and let  $\mathcal{A}^{-i} = \prod_{j \neq i} \mathcal{A}^j$ . Thus, for any function  $F: \mathcal{S} \times \mathcal{A}^{-i} \rightarrow \mathbb{R}$  which does not rely on  $a^i$ , we have for any  $s \in \mathcal{S}$  that

$$\begin{aligned} &\sum_{a \in \mathcal{A}} \pi_\theta(s, a) \cdot [\nabla_{\theta^i} \log \pi_{\theta^i}^i(s, a^i)] \cdot F(s, a^{-i}) \\ &= \sum_{a^{-i} \in \mathcal{A}^{-i}} F(s, a^{-i}) \cdot \left[ \prod_{j \in \mathcal{N}, j \neq i} \pi_{\theta^j}^j(s, a^j) \right] \cdot \left[ \sum_{a^i \in \mathcal{A}^i} \nabla_{\theta^i} \pi_{\theta^i}^i(s, a^i) \right] = 0. \end{aligned} \quad (\text{B.3})$$

Thus, replacing  $F$  in (B.3) by the value function  $V_\theta$  and function  $\tilde{V}_\theta^i$  defined in (3.1) and combined with (B.2), we establish (3.2), which concludes the proof.  $\square$

### B.2. Proof of Theorem 4.6

To proceed with the proof, we first establish the stability of the update  $\{\omega_t\}$ . This stability condition serves as an assumption in the original two-time-scale SA analysis (Borkar, 2008, Chapter 6.1). It is usually verified separately using some other sufficient conditions (Borkar & Meyn, 2000; Andrieu et al., 2005). We will directly use the lemma in the convergence analysis to follow and defer its proof to Appendix §C.

**Lemma B.1.** Under Assumptions 2.2, and 4.2-4.5, the sequence  $\{\omega_t^i\}$  generated from (3.3) is bounded almost surely, i.e.,  $\sup_t \|\omega_t^i\| < \infty$  a.s., for any  $i \in \mathcal{N}$ .

As in the classical two-time-scale SA analysis (Borkar, 2008), we let the policy parameter  $\theta_t$  be fixed as  $\theta_t \equiv \theta$  when analyzing the convergence of the critic step. This allows us to show that  $\omega_t$  will converge to some  $\omega_\theta$  depending on  $\theta$ , which can be further utilized to simplify the proof of convergence for the slower time scale. In fact, with linear function approximation, one can rewrite the actor step (3.5) for Algorithm 1 as

$$\theta_{t+1}^i = \Gamma^i \left( \theta_t^i + \beta_{\omega,t} \cdot \frac{\beta_{\theta,t}}{\beta_{\omega,t}} \cdot A_t^i \cdot \psi_t^i \right), \quad (\text{B.4})$$

where the projection  $\Gamma^i$  follows from Assumption 4.1. Note that  $A_t^i$  is also bounded a.s., since the parameter  $\omega_t^i$  is bounded by Lemma B.1, and the feature  $\phi_t$  is bounded by Assumption 4.5. Moreover,  $\psi_t^i$  is bounded a.s. by Assumption 2.2 since it is a continuous function over a bounded set  $\Theta^i$ . Therefore,  $\sup_t \|A_t^i \cdot \psi_t^i\| < \infty$  a.s. Now since  $\beta_{\theta,t} \cdot \beta_{\omega,t}^{-1} \rightarrow 0$  by Assumption 4.3, it follows that  $\beta_{\theta,t} \cdot \beta_{\omega,t}^{-1} \cdot A_t^i \cdot \psi_t^i \rightarrow 0$  as  $t \rightarrow \infty$ . Now the update in (B.4) can be viewed to track the ordinary differential equation (ODE)  $\dot{\theta}^i(t) = 0$ . Hence, one may let  $\theta_t$  be a constant when analyzing the faster update of  $\omega_t^i$ . For notational simplicity, we eliminate the notations associated with  $\theta$  unless otherwise noted.

Let  $\{\mathcal{F}_{t,1}\}$  be the filtration with  $\mathcal{F}_{t,1} = \sigma(r_\tau, \mu_\tau, \omega_\tau, s_\tau, a_\tau, C_{\tau-1}, \tau \leq t)$ , which is an increasing  $\sigma$ -algebra over time  $t$ . For notational convenience, let  $r_t = (r_t^1, \dots, r_t^N)^\top$ ,  $\mu_t = (\mu_t^1, \dots, \mu_t^N)^\top$ ,  $\omega_t = [(\omega_t^1)^\top, \dots, (\omega_t^N)^\top]^\top$ , and  $\delta_t = [(\delta_t^1)^\top, \dots, (\delta_t^N)^\top]^\top$ . The update of  $\omega_t$  in (3.3) can be rewritten in a compact form as

$$\omega_{t+1} = (C_t \otimes I)(\omega_t + \beta_{\omega,t} \cdot y_{t+1}), \quad (\text{B.5})$$

where  $\otimes$  denotes the Kronecker product,  $I \in \mathbb{R}^{K \times K}$  is the identity matrix, and  $y_{t+1} = (\delta_t^1 \phi_t^\top, \dots, \delta_t^N \phi_t^\top)^\top \in \mathbb{R}^{KN}$ . Define the operator  $\langle \cdot \rangle : \mathbb{R}^{KN} \rightarrow \mathbb{R}^K$  by letting

$$\langle \omega \rangle = \frac{1}{N} (\mathbb{1}^\top \otimes I) \omega = \frac{1}{N} \sum_{i \in \mathcal{N}} \omega^i \quad (\text{B.6})$$

for any  $\omega = [(\omega^1)^\top, \dots, (\omega^N)^\top]^\top \in \mathbb{R}^{KN}$  and  $\omega^i \in \mathbb{R}^K$  with  $i \in \mathcal{N}$ . That is,  $\langle \omega \rangle \in \mathbb{R}^K$  represents the average of the vectors in  $\{\omega^1, \dots, \omega^N\}$ , which are local to individual agents. Let  $\mathcal{J} = (1/N \cdot \mathbb{1} \mathbb{1}^\top) \otimes I$  be the projection operator that projects the vector into the *consensus* subspace  $\{\mathbb{1} \otimes u : u \in \mathbb{R}^K\}$ . Thus we have  $\mathcal{J}\omega = \mathbb{1} \otimes \langle \omega \rangle$ . Moreover, we define  $\mathcal{J}_\perp$  as the operator that projects the vector to the *disagreement* subspace, i.e.,  $\mathcal{J}_\perp = I - \mathcal{J}$ . Thus the disagreement vector  $\omega_\perp = \mathcal{J}_\perp \omega$  is written as

$$\omega_\perp = \mathcal{J}_\perp \omega = \omega - \mathbb{1} \otimes \langle \omega \rangle. \quad (\text{B.7})$$

The proof of Theorem 4.6 then consists of two steps. In particular, we separate the iteration  $\omega_t$  as the sum of a vector in this *consensus* space and a vector in the *disagreement* space, i.e.,  $\omega_t = \omega_{\perp,t} + \mathbb{1} \otimes \langle \omega_t \rangle$ . We first show the a.s. convergence of the disagreement vector sequence  $\{\omega_{\perp,t}\}$  to zero. Then, we prove that the consensus vector sequence  $\{\mathbb{1} \otimes \langle \omega_t \rangle\}$  converges to the equilibrium such that  $\langle \omega_t \rangle$  satisfies (4.2).

**Step 1.** In this step, we establish that  $\lim_t \omega_{\perp,t} = 0$  a.s. To this end, we first have the following lemma on the boundedness of the sequence  $\{\mu_t^i\}$  for any  $i \in \mathcal{N}$ .

**Lemma B.2.** Under Assumptions 2.2 and 4.2, the sequence  $\{\mu_t^i\}$  generated as in (3.3) is bounded almost surely, i.e.,  $\sup_t |\mu_t^i| < \infty$  a.s., for any  $i \in \mathcal{N}$ .

*Proof.* The local update in (3.3) forms a stochastic approximation iteration, whose asymptotic behavior can be captured by the ODE

$$\dot{\mu}^i = -\mu^i + \sum_{s \in \mathcal{S}} d_\theta(s) \sum_{a \in \mathcal{A}} \pi_\theta(s, a) R^i(s, a). \quad (\text{B.8})$$

Let  $f(\mu^i)$  denote the right hand side (RHS) of (B.8), which is Lipschitz continuous in  $\mu^i$ . Moreover, define  $f_c(\mu^i) = f(c\mu^i) \cdot c^{-1}$ , then  $f_\infty(\mu^i) = \lim_c f(c\mu^i) \cdot c^{-1} = -\mu^i$  exists. Therefore, the ODE  $\dot{\mu}^i = f_\infty(\mu^i)$  has origin as the unique asymptotically stable equilibrium. In addition, since  $r_t^i$  is uniformly bounded, we have

$$\mathbb{E} \left[ |r_{t+1}^i - \mathbb{E}(r_{t+1}^i | \mathcal{F}_{t,1})|^2 | \mathcal{F}_{t,1} \right] \leq K_0 \cdot (1 + |\mu_t^i|^2)$$

for some  $K_0 < \infty$ . Therefore, the conditions (a.1) and (a.4) in Assumption D.1 are satisfied. (See Appendix §D.1 for details.) By Assumption 2.2, (a.2) in Assumption D.1 also holds. We thus conclude that  $\sup_t |\mu_t^i| < \infty$  from Theorem D.3 (see also Theorem 9 on page 74-75 in Borkar (2008)).  $\square$

Let  $z_t^i = [\mu_t^i, (\omega_t^i)^\top]^\top$  and  $z_t = [(z_t^1)^\top, \dots, (z_t^N)^\top]^\top$ . By Lemma B.1, we have  $\mathbb{P}(\sup_t \|z_t\| < \infty) = 1$ , which means that  $\mathbb{P}(\bigcup_{M \in \mathbb{Z}^+} \{\sup_t \|z_t\| \leq M\}) = 1$ , with  $\mathbb{Z}^+$  denoting the set of positive integers. Hence, it suffices to show that  $\lim_t \omega_{\perp,t} \mathbb{I}_{\{\sup_t \|z_t\| \leq M\}} = 0$ , for any  $M \in \mathbb{Z}^+$ , where  $\mathbb{I}_{\{\cdot\}}$  is the indicator function. We then establish that  $\mathbb{E}(\|\beta_{\omega,t}^{-1} \omega_{\perp,t}\|^2)$  is bounded on  $\{\sup_t \|z_t\| \leq M\}$ , for any  $M > 0$ .

**Lemma B.3.** Under Assumptions 4.2-4.5, for any  $M > 0$ , we have

$$\sup_t \mathbb{E}(\|\beta_{\omega,t}^{-1} \omega_{\perp,t}\|^2 \mathbb{I}_{\{\sup_{\tau \leq t} \|z_\tau\| \leq M\}}) < \infty.$$

*Proof.* First note that by (a.1) in Assumption 4.4 and the fact that  $(A \otimes B)(C \otimes D) = (AC) \otimes (BD)$ , we have

$$(C_t \otimes I)(\mathbb{1} \otimes \langle \omega \rangle) = (C_t \mathbb{1}) \otimes \langle \omega \rangle = \mathbb{1} \otimes \langle \omega \rangle.$$

Hence,  $\omega_{\perp,t+1}$  has the form  $\omega_{\perp,t+1} = \mathcal{J}_\perp[(C_t \otimes I)(\omega_t + \beta_{\omega,t} y_{t+1})] = \mathcal{J}_\perp[(C_t \otimes I)(\omega_{\perp,t} + \beta_{\omega,t} y_{t+1})]$ , since  $\mathcal{J}_\perp(\mathbb{1} \otimes \langle \omega \rangle)$  is zero. Thus, by the definition of  $\mathcal{J}_\perp$  in (B.7), the vector  $\omega_{\perp,t+1}$  satisfies

$$\omega_{\perp,t+1} = [(I - \mathbb{1}\mathbb{1}^\top/N) \otimes I](C_t \otimes I)(\omega_{\perp,t} + \beta_{\omega,t} y_{t+1}) = [(I - \mathbb{1}\mathbb{1}^\top/N)C_t \otimes I](\omega_{\perp,t} + \beta_{\omega,t} y_{t+1}). \quad (\text{B.9})$$

Thus, we have

$$\begin{aligned} & \mathbb{E}(\|\beta_{\omega,t+1}^{-1} \omega_{\perp,t+1}\|^2 \mid \mathcal{F}_{t,1}) \\ &= \frac{\beta_{\omega,t}^2}{\beta_{\omega,t+1}^2} \cdot \mathbb{E}\{(\beta_{\omega,t}^{-1} \omega_{\perp,t} + y_{t+1})^\top [C_t^\top (I - \mathbb{1}\mathbb{1}^\top/N)C_t \otimes I](\beta_{\omega,t}^{-1} \omega_{\perp,t} + y_{t+1}) \mid \mathcal{F}_{t,1}\} \\ &\leq \frac{\beta_{\omega,t}^2}{\beta_{\omega,t+1}^2} \cdot \rho \cdot \mathbb{E}[(\beta_{\omega,t}^{-1} \omega_{\perp,t} + y_{t+1})^\top (\beta_{\omega,t}^{-1} \omega_{\perp,t} + y_{t+1}) \mid \mathcal{F}_{t,1}] \\ &\leq \frac{\beta_{\omega,t}^2}{\beta_{\omega,t+1}^2} \cdot \rho \cdot \left\{ \|\beta_{\omega,t}^{-1} \omega_{\perp,t}\|^2 + 2 \cdot \|\beta_{\omega,t}^{-1} \omega_{\perp,t}\| \cdot [\mathbb{E}(\|y_{t+1}\|^2 \mid \mathcal{F}_{t,1})]^\frac{1}{2} + \mathbb{E}(\|y_{t+1}\|^2 \mid \mathcal{F}_{t,1}) \right\}, \quad (\text{B.10}) \end{aligned}$$

where  $\rho$  represents the spectral norm of  $\mathbb{E}[C_t^\top (I - \mathbb{1}\mathbb{1}^\top/N)C_t]$ . By (a.2) in Assumption 4.4, we have  $\rho \in [0, 1)$ . The first inequality in (B.10) is due to the conditional independence of  $C_t$  and  $r_{t+1}^i$  for all  $i \in \mathcal{N}$ , and thus  $y_{t+1}$ , by (a.3) in Assumption 4.4, and the second inequality is due to the Cauchy-Schwarz inequality. Moreover, by the definition of  $y_{t+1}$ , we have

$$\begin{aligned} \mathbb{E}(\|y_{t+1}\|^2 \mid \mathcal{F}_{t,1}) &= \mathbb{E}\left(\sum_{i \in \mathcal{N}} \|\delta_t^i \phi_t\|^2 \mid \mathcal{F}_{t,1}\right) = \mathbb{E}\left[\sum_{i \in \mathcal{N}} \|(r_{t+1}^i - \mu_t^i + \phi_{t+1}^\top \omega_t^i - \phi_t^\top \omega_t^i) \phi_t\|^2 \mid \mathcal{F}_{t,1}\right] \\ &\leq 3 \cdot \mathbb{E}\left[\sum_{i \in \mathcal{N}} \|r_{t+1}^i \phi_t\|^2 + \|\mu_t^i \phi_t\|^2 + \|\phi_t(\phi_{t+1}^\top - \phi_t^\top)\|^2 \|\omega_t^i\|^2 \mid \mathcal{F}_{t,1}\right], \quad (\text{B.11}) \end{aligned}$$

where the inequality follows from that by Assumption 4.5,  $\mathbb{E}[\|\phi_t(\phi_{t+1}^\top - \phi_t^\top)\|^2 \mid \mathcal{F}_{t,1}]$  and  $\mathbb{E}(\|\phi_t\|^2 \mid \mathcal{F}_{t,1})$  are both uniformly bounded for any  $s_t \in \mathcal{S}$  and  $a_t \in \mathcal{A}$ . Moreover, by Assumption 4.2, we have  $\mathbb{E}(|r_{t+1}^i|^2 \mid \mathcal{F}_{t,1}) = \mathbb{E}(|r_{t+1}^i|^2 \mid s_t, a_t)$  also uniformly bounded. Thus the RHS of (B.11) is bounded on the set  $\{\sup_{\tau \leq t} \|z_\tau\| \leq M\}$  for any  $M > 0$  as follows, i.e., there exists  $K_1 < \infty$ , such that

$$\mathbb{E}(\|y_{t+1}\|^2 \mathbb{I}_{\{\sup_{\tau \leq t} \|z_\tau\| \leq M\}} \mid \mathcal{F}_{t,1}) \leq K_1 \cdot \left[1 + \sum_{i \in \mathcal{N}} \mathbb{E}(|r_{t+1}^i|^2 \mid \mathcal{F}_{t,1})\right]. \quad (\text{B.12})$$

Let  $\eta_t = \|\beta_{\omega,t}^{-1} \omega_{\perp,t}\|^2 \mathbb{I}_{\{\sup_{\tau \leq t} \|z_\tau\| \leq M\}}$  and note that  $\mathbb{I}_{\{\sup_{\tau \leq t+1} \|z_\tau\| \leq M\}} \leq \mathbb{I}_{\{\sup_{\tau \leq t} \|z_\tau\| \leq M\}}$ . Then by taking expectation over both sides of (B.10), we obtain that there exists  $K_2 = K_1 \cdot [1 + \mathbb{E}(\sum_{i \in \mathcal{N}} |r_{t+1}^i|^2)] < \infty$  such that

$$\mathbb{E}(\eta_{t+1}) \leq \frac{\beta_{\omega,t}^2}{\beta_{\omega,t+1}^2} \cdot \rho \cdot [\mathbb{E}(\eta_t) + 2\sqrt{\mathbb{E}(\eta_t)} \cdot \sqrt{K_2} + K_2]. \quad (\text{B.13})$$

Since  $\lim_t \beta_{\omega,t}^2 \cdot \beta_{\omega,t+1}^{-2} = 1$  and  $\rho < 1$ , for any  $\delta > 0$ , there exists a large enough  $t_0$  such that for any  $t > t_0$ ,  $\beta_{\omega,t}^2 \cdot \beta_{\omega,t+1}^{-2} \cdot \rho \leq 1 - \delta$ . Hence, there exist positive constants  $K_3$  and  $b$  such that for any  $t \geq t_0$ ,

$$\mathbb{E}(\eta_{t+1}) \leq (1 - \delta) \cdot [\mathbb{E}(\eta_t) + 2\sqrt{\mathbb{E}(\eta_t)} \cdot \sqrt{K_2} + K_2] \leq (1 - \delta/2) \cdot \mathbb{E}(\eta_t) + b \cdot \mathbb{I}_{\{\mathbb{E}(\eta_t) \leq K_3\}}.$$

By induction, we obtain that  $\mathbb{E}(\eta_t) \leq (1 - \delta/2)^{t-t_0} \mathbb{E}(\eta_{t_0}) + 2b/\delta$ . Hence, we have  $\sup_t \mathbb{E}(\eta_t) < \infty$ . In addition, since  $\mathbb{I}_{\{\sup_t \|z_t\| \leq M\}} \leq \mathbb{I}_{\{\sup_{\tau \leq t} \|z_\tau\| \leq M\}}$ , we further obtain

$$\sup_t \mathbb{E}(\|\beta_{\omega,t}^{-1} \omega_{\perp,t}\|^2 \mathbb{I}_{\{\sup_t \|z_t\| \leq M\}}) < \infty,$$

which concludes the proof.  $\square$

Therefore, by Lemma B.3, we obtain that for any  $M > 0$ , there exists a constant  $K_4 < \infty$ , such that for any  $t \geq 0$ ,  $\mathbb{E}(\|\omega_{\perp,t}\|^2) \leq K_4 \cdot \beta_{\omega,t}^2$  on the set  $\{\sup_t \|z_t\| \leq M\}$ . Since  $\sum_t \beta_{\omega,t}^2 < \infty$  by Assumption 4.3, we have that  $\sum_t \mathbb{E}(\|\omega_{\perp,t}\|^2 \mathbb{I}_{\{\sup_t \|z_t\| \leq M\}})$  is finite by Fubini's theorem. This shows that  $\sum_t \|\omega_{\perp,t}\|^2 \mathbb{I}_{\{\sup_t \|z_t\| \leq M\}} < \infty$  a.s., which further yields  $\lim_t \omega_{\perp,t} \mathbb{I}_{\{\sup_t \|z_t\| \leq M\}} = 0$  a.s. By Lemmas B.1 and B.2,  $\{\sup_t \|z_t\| < \infty\}$  holds with probability 1. This shows that  $\lim_t \omega_{\perp,t} = 0$  a.s., and thus concludes Step 1.

**Step 2.** We now proceed to show the convergence of the consensus vector  $\mathbb{1} \otimes \langle \omega_t \rangle$ . According to the update in (B.5) and definition (B.6), the iteration of  $\langle \omega_t \rangle$  has the form

$$\langle \omega_{t+1} \rangle = \frac{1}{N} (\mathbb{1}^\top \otimes \mathbf{I}) (\mathbf{C}_t \otimes \mathbf{I}) (\mathbb{1} \otimes \langle \omega_t \rangle + \omega_{t,\perp} + \beta_{\omega,t} y_{t+1}) = \langle \omega_t \rangle + \beta_{\omega,t} \langle (\mathbf{C}_t \otimes \mathbf{I}) (y_{t+1} + \beta_{\omega,t}^{-1} \omega_{\perp,t}) \rangle.$$

Hence, we write the updates for  $\langle \omega_t \rangle$  and  $\langle \mu_t \rangle$  as

$$\langle \mu_{t+1} \rangle = \langle \mu_t \rangle + \beta_{\omega,t} \cdot \mathbb{E}(\bar{r}_{t+1} - \langle \mu_t \rangle \mid \mathcal{F}_{t,1}) + \beta_{\omega,t} \cdot \xi_{t+1,1}, \quad (\text{B.14})$$

$$\langle \omega_{t+1} \rangle = \langle \omega_t \rangle + \beta_{\omega,t} \cdot \mathbb{E}(\langle \delta_t \rangle \phi_t \mid \mathcal{F}_{t,1}) + \beta_{\omega,t} \cdot \xi_{t+1,2}, \quad (\text{B.15})$$

where  $\xi_{t+1,1} = \bar{r}_{t+1} - \mathbb{E}(\bar{r}_{t+1} \mid \mathcal{F}_{t,1})$  and  $\xi_{t+1,2}$  is

$$\xi_{t+1,2} = \langle (\mathbf{C}_t \otimes \mathbf{I}) (y_{t+1} + \beta_{\omega,t}^{-1} \omega_{\perp,t}) \rangle - \mathbb{E}(\langle \delta_t \rangle \phi_t \mid \mathcal{F}_{t,1}).$$

Note that  $\mathbb{E}(\bar{r}_{t+1} - \langle \mu_t \rangle \mid \mathcal{F}_{t,1})$  is Lipschitz continuous in  $\langle \mu_t \rangle$ . Recall that  $\langle \delta_t \rangle$  has the form

$$\langle \delta_t \rangle = \frac{1}{N} \sum_{i \in \mathcal{N}} r_{t+1}^i - \mu_t^i + \phi_{t+1}^\top \omega_t^i - \phi_t^\top \omega_t^i = \bar{r}_{t+1} - \langle \mu_t \rangle + \phi_{t+1}^\top \langle \omega_t \rangle - \phi_t^\top \langle \omega_t \rangle.$$

Hence,  $\mathbb{E}(\langle \delta_t \rangle \phi_t \mid \mathcal{F}_{t,1})$  is Lipschitz continuous in both  $\langle \omega_t \rangle$  and  $\langle \mu_t \rangle$ , and thus the condition (a.1) in Assumption D.1 (See Appendix §D.1) is satisfied.

Note that  $\xi_{t,1}$  is a martingale difference sequence and satisfies

$$\mathbb{E}(\|\xi_{t+1,1}\|^2 \mid \mathcal{F}_{t,1}) \leq K_5 \cdot (1 + \|\langle \omega_t \rangle\|^2 + \|\langle \mu_t \rangle\|^2), \quad (\text{B.16})$$

for some  $K_5 < \infty$ , since  $\bar{r}_{t+1}$  is uniformly bounded. In addition, the term  $\xi_{t,2}$  is also a martingale difference sequence, since

$$\mathbb{E}[\langle (\mathbf{C}_t \otimes \mathbf{I}) (y_{t+1} + \beta_{\omega,t}^{-1} \omega_{\perp,t}) \rangle \mid \mathcal{F}_{t,1}] = \mathbb{E}[\langle (\mathbf{C}_t \otimes \mathbf{I}) y_{t+1} \rangle \mid \mathcal{F}_{t,1}] = \mathbb{E}(\langle y_{t+1} \rangle \mid \mathcal{F}_{t,1}) = \mathbb{E}(\langle \delta_t \rangle \phi_t \mid \mathcal{F}_{t,1}),$$

which results from the facts that  $\langle \omega_{\perp,t} \rangle = 0$  and that  $\mathbb{1}^\top \mathbb{E}(\mathbf{C}_t) = \mathbb{1}^\top$ . Moreover, we have

$$\mathbb{E}(\|\xi_{t+1,2}\|^2 \mid \mathcal{F}_{t,1}) \leq 2 \cdot \mathbb{E}(\|y_{t+1} + \beta_{\omega,t}^{-1} \omega_{\perp,t}\|_{\mathbf{G}_t}^2 \mid \mathcal{F}_{t,1}) + 2 \cdot \|\mathbb{E}(\langle \delta_t \rangle \phi_t \mid \mathcal{F}_{t,1})\|^2, \quad (\text{B.17})$$

where  $\mathbf{G}_t = \mathbf{C}_t^\top \mathbb{1} \mathbb{1}^\top \mathbf{C}_t \otimes \mathbf{I} \cdot N^{-2}$ . Note that  $\mathbf{G}_t$  has bounded spectral norm since  $\mathbf{C}_t$  is a stochastic matrix. Thus the first term in (B.17) can be further bounded over the set  $\{\sup_t \|z_t\| \leq M\}$ , for any  $M > 0$ . Notably, there exist  $K_6, K_7 < \infty$  such that

$$\begin{aligned} & \mathbb{E}(\|y_{t+1} + \beta_{\omega,t}^{-1} \omega_{\perp,t}\|_{\mathbf{G}_t}^2 \mid \mathcal{F}_{t,1}) \cdot \mathbb{I}_{\{\sup_t \|z_t\| \leq M\}} \\ & \leq K_6 \cdot \mathbb{E}(\|y_{t+1}\|^2 + \|\beta_{\omega,t}^{-1} \omega_{\perp,t}\|^2 \mid \mathcal{F}_{t,1}) \cdot \mathbb{I}_{\{\sup_t \|z_t\| \leq M\}} < K_7, \end{aligned}$$

where the second inequality follows from (B.12) and Lemma B.3. Moreover, the second term in (B.17) can be bounded by  $\|\mathbb{E}(\langle \delta_t \rangle \phi_t | \mathcal{F}_{t,1})\|^2 \leq K_8 \cdot (1 + \|\langle \omega_t \rangle\|^2 + \|\langle \mu_t \rangle\|^2)$  with some  $K_8 < \infty$ , due to the boundedness of  $r_{t+1}^i$  and  $\phi_t$  (from Assumptions 4.2 and 4.5). Hence, for any  $M > 0$ , it follows that

$$\mathbb{E}(\|\xi_{t+1,1}\|^2 | \mathcal{F}_{t,1}) \leq K_9 \cdot (1 + \|\langle \omega_t \rangle\|^2 + \|\langle \mu_t \rangle\|^2), \quad (\text{B.18})$$

over the set  $\{\sup_t \|z_t\| \leq M\}$  for some  $K_9 < \infty$ . This verifies that on the set  $\{\sup_t \|z_t\| \leq M\}$  for any  $M > 0$ , the condition (a.4) in Assumption D.1 is satisfied.

Now consider the following ODE that captures the asymptotic behavior of (B.14) and (B.15)

$$\dot{\langle z \rangle} = \begin{pmatrix} \dot{\langle \mu \rangle} \\ \dot{\langle \omega \rangle} \end{pmatrix} = \begin{pmatrix} -1 & 0 \\ -\Phi^\top D_\theta^{s,a} \mathbf{1} & \Phi^\top D_\theta^{s,a} (P^\theta - \mathbf{I}) \Phi \end{pmatrix} \begin{pmatrix} \langle \mu \rangle \\ \langle \omega \rangle \end{pmatrix} + \begin{pmatrix} J(\theta) \\ \Phi^\top D_\theta^{s,a} \bar{R} \end{pmatrix}. \quad (\text{B.19})$$

Recall that  $D_\theta^{s,a} = \text{diag}[d_\theta(s) \cdot \pi_\theta(s, a), s \in \mathcal{S}, a \in \mathcal{A}]$ . Let the RHS of the ODE (B.19) be  $h(\langle z \rangle)$ , then  $h(\langle z \rangle)$  is Lipschitz continuous in  $\langle z \rangle$ , which satisfies the condition (a.1) in Assumption D.1. By the Perron-Frobenius theorem and Assumption 2.2, the stochastic matrix  $P^\theta$  has a simple eigenvalue of 1, and its remaining eigenvalues have real parts less than 1. Hence,  $(P^\theta - \mathbf{I})$  has all eigenvalues with negative real parts but one zero, so does the matrix  $\Phi^\top D_\theta^{s,a} (P^\theta - \mathbf{I}) \Phi$ , since  $\Phi$  is full column rank by Assumption 4.5. The simple eigenvalue of zero has eigen-vector  $\nu$  that satisfies  $\Phi \nu = \alpha \mathbf{1}$  for some  $\alpha \neq 0$ , since  $\alpha \mathbf{1}$  lies in the eigen-space of  $D_\theta^{s,a} (P^\theta - \mathbf{I})$  associated with zero. By Assumption 4.5, however, this will not happen with any choice of  $\Phi$  since  $\Phi \nu \neq \alpha \mathbf{1}$  for any  $\nu \in \mathbb{R}^K$ . Hence, the ODE (B.19) is globally asymptotically stable and has its equilibrium satisfying

$$-\langle \mu \rangle = J(\theta), \quad \Phi^\top D_\theta^{s,a} [\bar{R} - \langle \mu \rangle \mathbf{1} + P^\theta \Phi \langle \omega \rangle - \Phi \langle \omega \rangle] = 0. \quad (\text{B.20})$$

Note that the corresponding solution for  $\langle \mu \rangle$  at equilibrium is  $J(\theta)$ , whereas the solution for  $\langle \omega \rangle$  has the form  $\omega_\theta + \alpha \nu$  with any  $\alpha \in \mathbb{R}$  and  $\nu \in \mathbb{R}^K$  such that  $\Phi \nu = \mathbf{1}$ . While by Assumption 4.5,  $\Phi \nu \neq \mathbf{1}$ , thus the term  $\omega_\theta$  is unique, and it follows that  $\Phi^\top D_\theta^{s,a} [T_\theta^Q(\Phi \omega_\theta) - \Phi \omega_\theta] = 0$  with  $T_\theta^Q$  as defined in (4.1).

Recall from Lemmas B.1 and B.2 that  $\{z_t\}$  is bounded a.s., so is the sequence  $\{\langle z_t \rangle\}$ . Hence all conditions for Theorem D.2 to hold are satisfied. For the concatenated vector  $\langle z_t \rangle = (\langle \mu_t \rangle, \langle \omega_t \rangle^\top)^\top$ , we thus have  $\lim_t \langle \mu_t \rangle = J(\theta)$  and  $\lim_t \langle \omega_t \rangle = \omega_\theta$  over the set  $\{\sup_t \|z_t\| \leq M\}$  for any  $M > 0$ . By Lemmas B.1 and B.2, this holds with probability 1, which concludes Step 2. Combined with Step 1, we arrive at the conclusion that  $\lim_t \omega_t^i = \omega_\theta$  for any  $i \in \mathcal{N}$ , which completes the proof for Theorem 4.6.  $\square$

### B.3. Proof of Theorem 4.7

Let  $\mathcal{F}_{t,2} = \sigma(\theta_\tau, \tau \leq t)$  be the  $\sigma$ -field generated by  $\{\theta_\tau, \tau \leq t\}$ . In addition, we define

$$\zeta_{t+1,1}^i = A_t^i \cdot \psi_t^i - \mathbb{E}_{s_t \sim d_{\theta_t}, a_t \sim \pi_{\theta_t}} (A_t^i \cdot \psi_t^i | \mathcal{F}_{t,2}), \quad \zeta_{t+1,2}^i = \mathbb{E}_{s_t \sim d_{\theta_t}, a_t \sim \pi_{\theta_t}} [(A_t^i - A_{t,\theta_t}^i) \cdot \psi_t^i | \mathcal{F}_{t,2}],$$

where  $A_{t,\theta_t}^i$  is as defined in (4.3) with  $\theta = \theta_t$ . Then the actor update in (3.5) with a local projection becomes

$$\theta_{t+1}^i = \Gamma^i \left[ \theta_t^i + \beta_{\theta,t} \mathbb{E}_{s_t \sim d_{\theta_t}, a_t \sim \pi_{\theta_t}} (A_{t,\theta_t}^i \cdot \psi_t^i | \mathcal{F}_{t,2}) + \beta_{\theta,t} \zeta_{t+1,1}^i + \beta_{\theta,t} \zeta_{t+1,2}^i \right]. \quad (\text{B.21})$$

Note that  $\zeta_{t+1,2}^i = o(1)$  since the critic converges, i.e.,  $A_t^i \rightarrow A_{t,\theta_t}^i$ , at the faster time scale. Moreover, letting  $M_t^i = \sum_{\tau=0}^t \beta_{\theta,\tau} \zeta_{\tau+1,1}^i$ ,  $\{M_t^i\}$  is a martingale sequence. Since the sequences  $\{\omega_t^i\}$ ,  $\{\psi_t^i\}$ , and  $\{\phi_t\}$  are all bounded, the sequence  $\{\zeta_{t,1}^i\}$  is also bounded. Hence, by Assumption 4.3, we have

$$\sum_t \mathbb{E}(\|M_{t+1}^i - M_t^i\|^2 | \mathcal{F}_{t,2}) = \sum_{t \geq 1} \|\beta_{\theta,t} \zeta_{t+1,1}^i\|^2 < \infty \text{ a.s.}$$

By the martingale convergence theorem (Proposition VII-2-3(c) on page 149 of Neveu (1975)), the martingale sequence  $\{M_t^i\}$  converges a.s. Thus, for any  $\epsilon > 0$ , we have

$$\lim_t \mathbb{P} \left( \sup_{n \geq t} \left\| \sum_{\tau=t}^n \beta_{\theta,\tau} \zeta_{\tau+1,1}^i \right\| \geq \epsilon \right) = 0.$$

In addition, let

$$g^i(\theta_t) = \mathbb{E}_{s_t \sim d_{\theta_t}, a_t \sim \pi_{\theta_t}} (A_{t,\theta_t}^i \cdot \psi_{t,\theta_t}^i | \mathcal{F}_{t,2}) = \sum_{s_t \in \mathcal{S}, a_t \in \mathcal{A}} d_{\theta_t}(s_t) \cdot \pi_{\theta_t}(s_t, a_t) \cdot A_{t,\theta_t}^i \cdot \psi_{t,\theta_t}^i,$$

then we show that  $g^i(\theta_t)$  is continuous in  $\theta_t^i$  as follows. First,  $\psi_{t,\theta_t}^i$  is continuous by Assumption 2.2. Also, the term  $d_{\theta_t}(s_t) \cdot \pi_{\theta_t}(s_t, a_t)$  is continuous in  $\theta_t^i$  since it is the stationary distribution and thus is the solution to  $d_{\theta_t}(s) \cdot \pi_{\theta_t}(s, a) = \sum_{s' \in \mathcal{S}, a' \in \mathcal{A}} P^{\theta_t}(s', a' | s, a) \cdot d_{\theta_t}(s') \cdot \pi_{\theta_t}(s', a')$  and  $\sum_{s \in \mathcal{S}, a \in \mathcal{A}} d_{\theta_t}(s) \cdot \pi_{\theta_t}(s, a) = 1$ , where  $P^{\theta_t}(s', a' | s, a) = P(s' | s, a) \cdot \pi_{\theta_t}(s', a')$ . The unique solution to this set of linear equations can be verified to be continuous in  $\theta_t$ , noting that  $\pi_{\theta_t}(s, a) > 0$  by Assumption 2.2. In addition,  $A_{t,\theta_t}^i$  is continuous in  $\theta_t^i$  since  $\omega_{\theta_t}$  is the unique solution to the linear equation  $\Phi^\top D_{\theta}^{s,a} [T_{\theta}^Q(\Phi \omega_{\theta}) - \Phi \omega_{\theta}] = 0$  and can also be verified to be continuous in  $\theta_t$ . Therefore, by Kushner-Clark lemma (Kushner & Clark, 1978, page 191-196) (see also Theorem D.5 in Appendix §D.2), the update in (B.21) converges a.s. to the set of asymptotically stable equilibria of the ODE (4.5) for each  $i \in \mathcal{N}$ , which concludes the proof.  $\square$

The proof for the convergence of Algorithm 2 is similar to the proofs in B.2 and B.3. To avoid duplication, we leave out some of the details in the proofs to follow.

#### B.4. Proof of Theorem 4.9

Let  $z_t^i = [\mu_t^i, (\lambda_t^i)^\top, (v_t^i)^\top]^\top \in \mathbb{R}^{1+M+L}$ . We first have the following lemma on the stability of the updates of  $\{z_t^i\}$ , as in Lemma B.1, the proof of which is provided in Appendix §C.

**Lemma B.4.** Under Assumptions 2.2, 4.2-4.4, and 4.8, the sequence  $\{z_t^i\}$  generated from (3.8) and (3.11) satisfies  $\sup_t \|z_t^i\| < \infty$  a.s., for any  $i \in \mathcal{N}$ .

Note that  $\tilde{\delta}_t^i \cdot \psi_t^i$  is bounded by Assumptions 2.2, 4.8, and Lemma B.4. Thus the actor step (3.12) can be viewed to track ODE  $\dot{\theta}^i = 0$  when analyzing the faster time scale update. Thus by the same argument as in §B.2, we fix the value of  $\theta_t$  as a constant  $\theta$ . For notational convenience, let  $v_t = [(v_t^1)^\top, \dots, (v_t^N)^\top]^\top$ ,  $\delta_t = [(\delta_t^1)^\top, \dots, (\delta_t^N)^\top]^\top$ , and  $z_t = [(z_t^1)^\top, \dots, (z_t^N)^\top]^\top$ . By little abuse of notation, we here use  $\{\mathcal{F}_{t,1}\}$  to denote the filtration with  $\mathcal{F}_{t,1} = \sigma(r_\tau, z_\tau, s_\tau, C_{\tau-1}, \tau \leq t)$ , an increasing  $\sigma$ -algebra. Then the updates of  $z_t$  in (3.8) and (3.11) have the following compact form

$$z_{t+1} = (C_t \otimes I)(z_t + \beta_{v,t} \cdot y_{t+1}), \quad (\text{B.22})$$

where  $y_t = [(y_t^1)^\top, \dots, (y_t^N)^\top]^\top \in \mathbb{R}^{(1+M+L)N}$ . Here we denote  $[r_{t+1}^i - \mu_t^i, (r_{t+1}^i - f_t^\top \lambda_t^i) f_t^\top, \delta_t^i \varphi_t^\top]^\top$  by  $y_{t+1}^i$ . Recall that  $f_t = f(s_t, a_t)$  and  $\varphi_t = \varphi(s_t)$ . With the same definitions for  $\langle \cdot \rangle$ ,  $\mathcal{J}$ , and  $\mathcal{J}_\perp$ , we can also separate the iteration of  $z_t$  as the sum of the consensus vector and the disagreement vector, i.e.,  $z_t = z_{\perp,t} + \mathbb{1} \otimes \langle z_t \rangle$ , with  $z_{\perp,t} = \mathcal{J}_\perp z_t$ . Then the proof proceeds again in two steps as follows.

**Step 1.** We first establish that  $\lim_t z_{\perp,t} = 0$  a.s. By Lemma B.4, it suffices to show that for any  $M \in \mathbb{Z}^+$ ,  $\lim_t z_{\perp,t} \mathbb{I}_{\{\sup_t \|z_t\| \leq M\}} = 0$ . We first establish the boundedness of  $\mathbb{E}(\|\beta_{v,t}^{-1} z_{\perp,t}\|^2)$  over the set  $\{\sup_t \|z_t\| \leq M\}$ , for any  $M > 0$ .

**Lemma B.5.** Under Assumptions 4.2- 4.4, and 4.8, for any  $M > 0$ , we have

$$\sup_t \mathbb{E}(\|\beta_{v,t}^{-1} z_{\perp,t}\|^2 \mathbb{I}_{\{\sup_t \|z_t\| \leq M\}}) < \infty.$$

*Proof.* Following the derivation of (B.9), we obtain the iteration of  $z_{\perp,t}$  as

$$z_{\perp,t+1} = [(I - \mathbb{1}\mathbb{1}^\top/N) \otimes I](C_t \otimes I)(z_{\perp,t} + \beta_{v,t} y_{t+1}) = [(I - \mathbb{1}\mathbb{1}^\top/N) C_t \otimes I](z_{\perp,t} + \beta_{v,t} y_{t+1}). \quad (\text{B.23})$$

Thus, similar to the derivation of (B.10), we obtain

$$\begin{aligned} \mathbb{E}(\|\beta_{v,t+1}^{-1} z_{\perp,t+1}\|^2 | \mathcal{F}_{t,1}) &\leq \rho \cdot \frac{\beta_{v,t}^2}{\beta_{v,t+1}^2} \cdot \left\{ \|\beta_{v,t}^{-1} z_{\perp,t}\|^2 + 2\|\beta_{v,t}^{-1} z_{\perp,t}\| \right. \\ &\quad \cdot [\mathbb{E}(\|y_{t+1}\|^2 | \mathcal{F}_{t,1})]^\frac{1}{2} + \mathbb{E}(\|y_{t+1}\|^2 | \mathcal{F}_{t,1}) \Big\}, \end{aligned} \quad (\text{B.24})$$



where  $\rho$  represents the spectral norm of  $\mathbb{E}[C_t^\top (\mathbf{I} - \mathbb{1}\mathbb{1}^\top / N) C_t]$  and  $\rho \in [0, 1)$ . Then we have

$$\mathbb{E}(\|y_{t+1}\|^2 \mid \mathcal{F}_{t,1}) = \mathbb{E} \left[ \sum_{i \in \mathcal{N}} |r_{t+1}^i - \mu_t^i|^2 + \|(r_{t+1}^i - f_t^\top \lambda_t^i) f_t\|^2 + \|\delta_t^i \varphi_t\|^2 \mid \mathcal{F}_{t,1} \right]. \quad (\text{B.25})$$

By Assumption 4.8, we have  $\mathbb{E}[\|\varphi_t(\varphi_{t+1}^\top - \varphi_t^\top)\|^2 \mid \mathcal{F}_{t,1}]$  and  $\mathbb{E}(\|\varphi_t\|^2 \mid \mathcal{F}_{t,1})$  uniformly bounded for any  $s_t \in \mathcal{S}$ . Moreover, by Assumption 4.2, we have  $\mathbb{E}(|r_{t+1}^i|^2 \mid \mathcal{F}_{t,1}) = \mathbb{E}(|r_{t+1}^i|^2 \mid s_t, a_t)$  also uniformly bounded for any  $s_t \in \mathcal{S}, a_t \in \mathcal{A}$ . Thus for any  $M > 0$ , there exist  $K_1, K_2 < \infty$ <sup>1</sup> such that (B.25) is further bounded as

$$\mathbb{E}(\|y_{t+1}\|^2 \mathbb{I}_{\{\sup_{\tau \leq t} \|z_\tau\| \leq M\}} \mid \mathcal{F}_{t,1}) \leq K_1 \cdot \left[ 1 + \sum_{i \in \mathcal{N}} \mathbb{E}(|r_{t+1}^i|^2 \mid \mathcal{F}_{t,1}) \right] < K_2. \quad (\text{B.26})$$

Let  $\eta_t = \|\beta_{v,t}^{-1} z_{\perp,t}\|^2 \mathbb{I}_{\{\sup_{\tau \leq t} \|z_\tau\| \leq M\}}$ . By taking expectation on both sides of (B.24), we obtain

$$\mathbb{E}(\eta_{t+1}) \leq \frac{\beta_{v,t}^2}{\beta_{v,t+1}^2} \cdot \rho \cdot [\mathbb{E}(\eta_t) + 2\sqrt{\mathbb{E}(\eta_t)} \cdot \sqrt{K_2} + K_2].$$

Following the same argument as in the proof of Lemma B.3, we obtain  $\sup_t \mathbb{E}(\eta_t) < \infty$  and thus

$$\sup_t \mathbb{E}(\|\beta_{v,t}^{-1} z_{\perp,t}\|^2 \mathbb{I}_{\{\sup_{\tau \leq t} \|z_\tau\| \leq M\}}) < \infty,$$

which concludes the proof.  $\square$

Therefore, by Lemma B.5 and Assumption 4.3, we arrive at  $\lim_t z_{\perp,t} \mathbb{I}_{\{\sup_{\tau \leq t} \|z_\tau\| \leq M\}} = 0$  a.s. for any  $M > 0$ . By Lemma B.4,  $\{\sup_t \|z_t\| < \infty\}$  holds with probability 1. This shows that  $\lim_t z_{\perp,t} = 0$  a.s., and thus concludes **Step 1**.

**Step 2.** We now proceed to show the convergence of the consensus vector  $\mathbb{1} \otimes \langle z_t \rangle$ . The iteration of  $\langle z_t \rangle$  has the form

$$\langle z_{t+1} \rangle = \frac{1}{N} (\mathbb{1}^\top \otimes \mathbf{I}) (C_t \otimes \mathbf{I}) (\mathbb{1} \otimes \langle z_t \rangle + z_{t,\perp} + \beta_{v,t} y_{t+1}) = \langle z_t \rangle + \beta_{v,t} \langle (C_t \otimes \mathbf{I}) (y_{t+1} + \beta_{v,t}^{-1} z_{\perp,t}) \rangle.$$

Hence, the update for  $\langle z_t \rangle$  becomes

$$\langle z_{t+1} \rangle = \langle z_t \rangle + \beta_{v,t} \cdot \mathbb{E}(\langle y_{t+1} \rangle \mid \mathcal{F}_{t,1}) + \beta_{v,t} \cdot \xi_{t+1}, \quad (\text{B.27})$$

where  $\xi_{t+1}$  and  $\langle y_{t+1} \rangle$  have the form

$$\begin{aligned} \xi_{t+1} &= \langle (C_t \otimes \mathbf{I}) (y_{t+1} + \beta_{v,t}^{-1} z_{\perp,t}) \rangle - \mathbb{E}(\langle y_{t+1} \rangle \mid \mathcal{F}_{t,1}), \\ \langle y_{t+1} \rangle &= [\bar{r}_{t+1} - \langle \mu_t \rangle, (\bar{r}_{t+1} - f_t^\top \langle \lambda_t \rangle) f_t^\top, \langle \delta_t \rangle \varphi_t^\top]^\top, \end{aligned}$$

respectively. Recall that  $\langle \delta_t \rangle = \bar{r}_{t+1} - \langle \mu_t \rangle + \varphi_{t+1}^\top \langle v_t \rangle - \varphi_t^\top \langle v_t \rangle$ . Note that  $\mathbb{E}(\langle y_{t+1} \rangle \mid \mathcal{F}_{t,1})$  is Lipschitz continuous in  $\langle z_t \rangle = (\langle \mu_t \rangle, \langle \lambda_t \rangle^\top, \langle v_t \rangle^\top)^\top$ , and thus the condition (a.1) in Assumption D.1 is satisfied. Moreover, one can verify that the term  $\xi_t$  is a martingale difference sequence. The conditional second moment of  $\xi_t$  can be bounded as

$$\mathbb{E}(\|\xi_{t+1}\|^2 \mid \mathcal{F}_{t,1}) \leq 2 \cdot \mathbb{E}(\|y_{t+1} + \beta_{v,t}^{-1} z_{\perp,t}\|_{G_t}^2 \mid \mathcal{F}_{t,1}) + 2 \cdot \|\mathbb{E}(\langle y_{t+1} \rangle \mid \mathcal{F}_{t,1})\|^2, \quad (\text{B.28})$$

where  $G_t = C_t^\top \mathbb{1}\mathbb{1}^\top C_t \otimes \mathbf{I} \cdot N^{-2}$  has bounded spectral norm. Thus the first term in (B.28) is bounded over the set  $\{\sup_t \|z_t\| \leq M\}$ , for any  $M > 0$ , i.e., there exist  $K_3 < \infty$  such that

$$\mathbb{E}(\|y_{t+1} + \beta_{v,t}^{-1} z_{\perp,t}\|_{G_t}^2 \mid \mathcal{F}_{t,1}) \cdot \mathbb{I}_{\{\sup_{\tau \leq t} \|z_\tau\| \leq M\}} \leq K_3 \cdot \mathbb{E}(\|y_{t+1}\|^2 + \|\beta_{v,t}^{-1} z_{\perp,t}\|^2 \mid \mathcal{F}_{t,1}) \cdot \mathbb{I}_{\{\sup_{\tau \leq t} \|z_\tau\| \leq M\}}.$$

<sup>1</sup>We note that  $K_1$  and  $K_2$  here are absolute constant values, with slight abuse of notation, we use the same notation as in the proof of Theorem 4.6. The same abuse applies to other constants with notation  $K_j$ , for any  $j \in \mathbb{N}$ .

From (B.26) and Lemma B.5, we obtain that the RHS can be further bounded by some  $K_4 < \infty$ . Moreover, the second term in (B.28) can be bounded by

$$\|\mathbb{E}(\langle y_{t+1} \rangle \mid \mathcal{F}_{t,1})\|^2 \leq \mathbb{E}(\|\langle y_{t+1} \rangle\|^2 \mid \mathcal{F}_{t,1}) \leq K_5 \cdot (1 + \|\langle \mu_t \rangle\|^2 + \|\langle \lambda_t \rangle\|^2 + \|\langle v_t \rangle\|^2) = K_5 \cdot (1 + \|\langle z_t \rangle\|^2)$$

with some  $K_5 < \infty$  due to the boundedness of  $r_{t+1}^i$ ,  $f_t$ , and  $\varphi_t$ . Hence, for any  $M > 0$ , it follows that

$$\mathbb{E}(\|\xi_{t+1}\|^2 \mid \mathcal{F}_{t,1}) \leq K_6 \cdot (1 + \|\langle z_t \rangle\|^2), \quad (\text{B.29})$$

over the set  $\{\sup_t \|z_t\| \leq M\}$  for some  $K_6 < \infty$ . This verifies the condition (a.4) in Assumption D.1.

Then the ODE associated with (B.27) has the form

$$\begin{pmatrix} \dot{\langle \mu \rangle} \\ \dot{\langle \lambda \rangle} \\ \dot{\langle v \rangle} \end{pmatrix} = \begin{pmatrix} -1 & 0 & 0 \\ 0 & -F^\top D_\theta^{s,a} F & 0 \\ -\Phi^\top D_\theta^s \mathbb{1} & 0 & \Phi^\top D_\theta^s (P^\theta - I) \Phi \end{pmatrix} \begin{pmatrix} \langle \mu \rangle \\ \langle \lambda \rangle \\ \langle v \rangle \end{pmatrix} + \begin{pmatrix} J(\theta) \\ F^\top D_\theta^{s,a} \bar{R} \\ \Phi^\top D_\theta^s \bar{R}_\theta \end{pmatrix}. \quad (\text{B.30})$$

Letting the RHS of the ODE (B.30) be  $h(\langle z \rangle)$ , we have  $h(\langle z \rangle)$  Lipschitz continuous in  $\langle z \rangle$ . Similar to the proof in **Step 2** of §B.2, one can verify that the ODE has a unique globally asymptotically stable equilibrium  $[J(\theta), \lambda_\theta^\top, v_\theta^\top]^\top$ , by Assumption 4.8 on the feature matrices  $F$  and  $\Phi$ . Here  $\lambda_\theta$  and  $v_\theta$  are the unique solutions to  $F^\top D_\theta^{s,a} (\bar{R} - F \lambda_\theta) = 0$  and  $\Phi^\top D_\theta^s [T_\theta^V (\Phi v_\theta) - \Phi v_\theta] = 0$ , respectively. Recall the operator  $T_\theta^V$  defined in (4.6). Moreover, the sequence  $\{z_t\}$  is bounded almost surely by Assumption B.4. Hence all conditions for Theorem D.2 to hold are satisfied. We thus have  $\lim_t \langle \mu_t \rangle = J(\theta)$ ,  $\lim_t \langle \lambda_t \rangle = \lambda_\theta$ , and  $\lim_t \langle v_t \rangle = v_\theta$  over the set  $\{\sup_t \|z_t\| \leq M\}$  for any  $M > 0$ . By Lemma B.4 and the results from **Step 1**, we obtain that  $\lim_t \mu_t^i = J(\theta)$ ,  $\lim_t \lambda_t^i = \lambda_\theta$ , and  $\lim_t v_t^i = v_\theta$  for any  $i \in \mathcal{N}$  a.s., which completes the proof.  $\square$

### B.5. Proof of Theorem 4.10

Let  $\mathcal{F}_{t,2} = \sigma(\theta_\tau, \tau \leq t)$  be the  $\sigma$ -field generated by  $\theta_\tau, \tau \leq t$ . Let

$$\zeta_{t+1,1}^i = \tilde{\delta}_t^i \cdot \psi_t^i - \mathbb{E}_{s_t \sim d_{\theta_t}, a_t \sim \pi_{\theta_t}} (\tilde{\delta}_t^i \cdot \psi_t^i \mid \mathcal{F}_{t,2}), \quad \zeta_{t+1,2}^i = \mathbb{E}_{s_t \sim d_{\theta_t}, a_t \sim \pi_{\theta_t}} [(\tilde{\delta}_t^i - \tilde{\delta}_{t,\theta_t}^i) \cdot \psi_t^i \mid \mathcal{F}_{t,2}],$$

where  $\tilde{\delta}_{t,\theta_t}^i$  is as defined in (4.9) with  $\theta = \theta_t$ . Then the actor update in (3.12) with a local projection becomes

$$\theta_{t+1}^i = \Gamma^i \left[ \theta_t^i + \beta_{\theta,t} \mathbb{E}_{s_t \sim d_{\theta_t}, a_t \sim \pi_{\theta_t}} (\tilde{\delta}_t^i \cdot \psi_t^i \mid \mathcal{F}_{t,2}) + \beta_{\theta,t} \zeta_{t+1,1}^i + \beta_{\theta,t} \zeta_{t+1,2}^i \right]. \quad (\text{B.31})$$

Note that  $\zeta_{t+1,2}^i = o(1)$  since the critic converges, i.e.,  $\tilde{\delta}_t^i \rightarrow \tilde{\delta}_{t,\theta_t}^i$ , at the faster time scale. Moreover, letting  $M_t^i = \sum_{\tau=0}^t \beta_{\theta,\tau} \zeta_{\tau+1,1}^i$ , we have  $\{M_t^i\}$  a martingale sequence. Note that the sequences  $\{z_t^i\}$ ,  $\{\psi_t^i\}$ , and  $\{\phi_t\}$  are all bounded, and so is the sequence  $\{\zeta_{t,1}^i\}$ . Hence, we have  $\sum_t \mathbb{E}(\|M_{t+1}^i - M_t^i\|^2 \mid \mathcal{F}_{t,2}) < \infty$  a.s., and further obtain that the martingale sequence  $\{M_t^i\}$  converges a.s. (Neveu, 1975, page 149). Thus the condition (a.4) in Assumption D.4 is satisfied. (See Appendix §D.2 for details.) One can also verify that  $\mathbb{E}_{s_t \sim d_{\theta_t}, a_t \sim \pi_{\theta_t}} (\tilde{\delta}_t^i \cdot \psi_t^i \mid \mathcal{F}_{t,2})$  is continuous in  $\theta_t^i$ , similar as the argument in §B.3. Therefore, we can apply the Kushner-Clark lemma (Theorem D.5 in Appendix §D.2) to show that the update in (B.31) converges a.s. to the set of asymptotically stable equilibria of the ODE (4.10), for each  $i \in \mathcal{N}$ , which concludes the proof.  $\square$

### C. Proofs of the Stability of Consensus Updates

As mentioned before, the stability of the updates in stochastic approximation is usually proved separately. Here we provide the proof for the stability of the slower time scale update  $\{\omega_t^i\}$  in Algorithm 1 and  $\{z_t^i\}$  in Algorithm 2, i.e., Lemmas B.1 and B.4. In particular, we provide a sufficient condition for the stability of the consensus-based SA updates following the spirit of the results in Borkar & Meyn (2000) and Mathkar & Borkar (2016). We first state a main theorem for stability and then verify that Lemmas B.1 and B.4 are two special cases of it.

Let  $\mathcal{N}$  be the set of agents with  $|\mathcal{N}| = N$  and  $x^i \in \mathbb{R}^d$  for any  $i \in \mathcal{N}$ . Consider the consensus update for  $x_n^i \in \mathbb{R}^d$  as<sup>2</sup>

$$x_{n+1}^i = \sum_{j \in \mathcal{N}} c_n(i, j) \{x_n^j + \gamma_n [h^j(x_n, Y_n) + M_{n+1}^j]\}, \text{ for any } i \in \mathcal{N}, \quad (\text{C.1})$$

where  $\{Y_n\}_{n \geq 0}$  is an irreducible and aperiodic Markov chain over the finite set  $A$ . Let  $\eta$  denote the stationary distribution of  $\{Y_n\}$  and  $\bar{h}^i(x_n) = \mathbb{E}_{Y_n \sim \eta} [h^i(x_n, Y_n)]$  denote the expectation of  $h^i(x_n, Y_n)$  over  $\eta$ . Let  $x_n = [(x_n^1)^\top, \dots, (x_n^N)^\top]^\top \in \mathbb{R}^{dN}$ ,  $C_n = [c_n(i, j)]_{N \times N}$ ,  $h = [(h^1)^\top, \dots, (h^N)^\top]^\top \in \mathbb{R}^{dN}$ ,  $\bar{h}(x) = [(\bar{h}^1)^\top, \dots, (\bar{h}^N)^\top]^\top \in \mathbb{R}^{dN}$ , and  $M_n = [(M_n^1)^\top, \dots, (M_n^N)^\top]^\top \in \mathbb{R}^{dN}$ . Let  $\{\mathcal{F}_n\}$  be the filtration with  $\mathcal{F}_n = \sigma(x_m, M_m, Y_m, C_{m-1}, m \leq n)$ .

**Assumption C.1.** We make the following assumptions:

- (a.1) The consensus weight matrices  $\{C_n\}$  satisfy Assumption 4.4;
- (a.2)  $h^i : \mathbb{R}^n \times A \rightarrow \mathbb{R}^n$  is Lipschitz continuous in its first argument for any  $i \in \mathcal{N}$ ;
- (a.3)  $\{M_n\}$  is a martingale difference sequence satisfying

$$\mathbb{E}(\|M_{n+1}\|^2 | \mathcal{F}_n) \leq K \cdot (1 + \|x_n\|^2),$$

for some  $K > 0$ ;

- (a.4) The difference  $\zeta_{n+1} = \bar{h}(x_n) - h(x_n, Y_n)$  satisfies

$$\|\zeta_{n+1}\|^2 \leq K' \cdot (1 + \|x_n\|^2) \text{ a.s.,}$$

for some  $K' > 0$ ;

- (a.5) The stepsize sequence  $\{\gamma_n\}$  satisfies  $\sum_n \gamma_n = \infty$  and  $\sum_n \gamma_n^2 < \infty$ ;
- (a.6) Define  $h_c : \mathbb{R}^{dN} \rightarrow \mathbb{R}^{dN}$  as  $h_c(x) = \bar{h}(cx) \cdot c^{-1}$  with some  $c > 0$ , and  $\tilde{h}_c(y) : \mathbb{R}^d \rightarrow \mathbb{R}^d$  be  $\tilde{h}_c(y) = \langle h_c(\mathbf{1} \otimes y) \rangle$ . Then  $\tilde{h}_c(y) \rightarrow h_\infty(y)$  as  $c \rightarrow \infty$  uniformly on compact sets for some  $h_\infty(y) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ . Also, for some  $\epsilon < N^{-1/2}$ ,  $B^\epsilon = \{y \mid \|y\| < \epsilon\}$  contains a globally asymptotically stable attractor of the ODE  $\dot{y} = h_\infty(y)$ .

Note that the definitions of  $h_c$  and  $h_\infty$  in (a.6) in Assumption C.1 are different from those in Mathkar & Borkar (2016). Here we consider the averaged ODE in the consensus subspace for each agent, while that reference considers the overall ODE associated with (C.1), i.e., define  $\tilde{h} : \mathbb{R}^{dN} \rightarrow \mathbb{R}^{dN}$  as  $\tilde{h}(x) = C_* \bar{h}(x)$  and let  $h_\infty(x) = \lim_c \tilde{h}(cx) \cdot c^{-1}$ , where  $C_* = \lim_n \prod_{m=1}^n C_m$ . In fact, from Nedic & Ozdaglar (2009); Nedic et al. (2016), the limit  $C_*$  exists and has identical rows and rank one, provided the sequence  $\{C_t\}$  satisfies Assumption 4.4. Therefore, the globally asymptotical stability of the ODE  $\dot{x} = h_\infty(x)$  (see Assumption (A5) in Mathkar & Borkar (2016)) does not hold for the linear ODE we consider in the convergence proof of the critic steps in both algorithms. In contrast, we can verify our condition (a.6) in Assumption C.1 later in the proof of Lemmas B.1 and B.4. We then have the following theorem on the stability of the sequence  $\{x_n\}$ .

**Theorem C.2.** Under Assumption C.1, the sequence  $\{x_n\}$  generated from (C.1) is bounded almost surely, i.e.,  $\sup_n \|x_n^i\| < \infty$  a.s. for any  $i \in \mathcal{N}$ .

<sup>2</sup>To avoid possible confusion with the notation of continuous time  $t$  needed in the stability analysis, we use subscript  $n$  to denote the iteration index in the proof of Theorem C.2.

### C.1. Proof of Theorem C.2

Let  $\vartheta_c(y, t)$  denote the solution to the ODE

$$\dot{y} = \tilde{h}_c(y), \quad y(0) = y \quad (\text{C.2})$$

We first have the following lemma, which is similar to Lemma 5 in [Mathkar & Borkar \(2016\)](#), and thus we leave out its proof here.

**Lemma C.3.** There exist constants  $c_0 > 0$  and  $T > 0$  such that for all initial conditions  $y$  within the sphere  $\{y \mid \|y\| \leq N^{-1/2}\}$  and all  $c \geq c_0$ , we have  $\|\vartheta_c(y, t)\| \leq (1 - \epsilon') \cdot N^{-1/2}$  for  $t \in [T, T + 1]$ , for some  $0 < \epsilon' < 1$ , where  $N$  is the number of agents.

Now, before stating the next set of lemmas, we introduce some notations and terminology. First, by the convention adopted in [Borkar & Meyn \(2000\)](#), we define  $t_0 = 0$  and  $t_n = \sum_{i=0}^n \gamma_i, n \geq 0$ . Then we define  $\bar{x}(t), t \geq 0$  as  $\bar{x}(t_n) = x_n, n \geq 0$  with linear interpolation on each interval  $[t_n, t_{n+1}]$ . Moreover, we let  $T_0 = 0$  and  $T_{n+1} = \min\{t_m : t_m \geq T_n + T\}$  for any  $n \geq 0$ . Then  $T_{n+1} \in [T_n + T, T_n + T + \sup_n \gamma_n]$ . Let  $m(n)$  be such that  $T_n = t_{m(n)}$ , for  $n \geq 0$ . Define the piecewise continuous trajectory  $\hat{x}(t) = \bar{x}(t) \cdot r_n^{-1}$  for  $t \in [T_n, T_{n+1})$ , where  $r_n = \max_{i \in \mathcal{N}} \{\|\bar{x}(T_n)\|, 1\}$ . This implies that  $\|\hat{x}(T_n)\| \leq 1$  for any  $n \geq 0$ . We also define  $\hat{x}(T_{n+1}) = \bar{x}(T_{n+1}) \cdot r_n^{-1}$ ,  $\hat{M}_{k+1} = M_{k+1} \cdot r_n^{-1}$ , and  $\hat{\zeta}_{k+1} = \zeta_{k+1} \cdot r_n^{-1}$  for  $k \in [m(n), m(n+1))$ .

Note that  $\{\hat{M}_k\}$  is also a martingale difference sequence as  $\{M_k\}$ . We first establish boundedness of  $\mathbb{E}[\|\hat{x}(t)\|^2]$  as follows.

**Lemma C.4.** Under Assumption C.1,  $\sup_t \mathbb{E}[\|\hat{x}(t)\|^2] < \infty$ .

*Proof.* It suffices to show that  $\sup_{m(n) \leq k < m(n+1)} \mathbb{E}[\|\hat{x}(t_k)\|^2] < M$  for some  $M > 0$  independent of  $n$ . We first write the update of  $\hat{x}(t_k)$  for  $k \in [m(n), m(n+1))$  in a compact form as

$$\hat{x}(t_{k+1}) = (C_k \otimes I)(\hat{x}(t_k) + \gamma_k \{h_{r_n}[\hat{x}(t_k)] + \hat{M}_{k+1} + \hat{\zeta}_{k+1}\}). \quad (\text{C.3})$$

Note that the additional term  $\hat{\zeta}_{k+1}$  also satisfies  $\mathbb{E}(\|\hat{\zeta}_{k+1}\|^2 \mid \mathcal{F}_k) \leq K' \cdot (1 + \|\hat{x}(t_k)\|^2)$  by condition (a.4) in Assumption C.1 since  $r_n \geq 1$ . Moreover, since  $C_k \otimes I$  has bounded norm, it follows similarly as in the proof for Lemma 4 on page 25 in [Borkar \(2008\)](#) that

$$\mathbb{E}[\|\hat{x}(t_{k+1})\|^2]^{1/2} \leq \mathbb{E}[\|\hat{x}(t_k)\|^2]^{1/2} (1 + \gamma_k K_1) + \gamma_k K_2,$$

for some  $K_1, K_2 > 0$  and  $k \in [m(n), m(n+1))$ . Then, by Grönwall inequality, we have the desired boundedness of  $\mathbb{E}[\|\hat{x}(t)\|^2]$ .  $\square$

By Lemma C.4, we immediately have the following result.

**Lemma C.5** (Lemma 5 on page 25 in [Borkar \(2008\)](#)). The sequence  $\{\sum_{k=0}^{n-1} \gamma_k \hat{M}_{k+1}\}$  converges almost surely.

We thus obtain the almost sure boundedness of the trajectory  $\{\hat{x}(t)\}$ .

**Lemma C.6.** Under Assumption C.1,  $\sup_t \|\hat{x}(t)\| < \infty$  a.s.

*Proof.* Recall the update in (C.3) and note that  $\|(C_k \otimes I)x\|_\infty \leq \|x\|_\infty$  since  $C_k$  is a row stochastic matrix, where  $\|\cdot\|_\infty$  denotes the infinite norm of a vector. Thus we have

$$\|\hat{x}(t_{k+1})\|_\infty \leq \|\hat{x}(t_k)\|_\infty + \gamma_k \|h_{r_n}[\hat{x}(t_k)] + \hat{M}_{k+1} + \hat{\zeta}_{k+1}\|_\infty. \quad (\text{C.4})$$

By iterating (C.4), we obtain

$$\begin{aligned} \|\hat{x}(t_{k+1})\|_\infty &\leq \|\hat{x}(t_{m(n)})\|_\infty + \sum_{l=0}^{k-m(n)} \gamma_{m(n)+l} (\|h_{r_n}[\hat{x}(t_{m(n)+l})]\|_\infty + \|\hat{M}_{m(n)+l+1}\|_\infty + \|\hat{\zeta}_{m(n)+l+1}\|_\infty) \\ &\leq \|\hat{x}(t_{m(n)})\|_\infty + \sum_{l=0}^{k-m(n)} \gamma_{m(n)+l} \cdot K_3 [1 + \|\hat{x}(t_{m(n)+l})\|_\infty] + \sum_{l=0}^{k-m(n)} \gamma_{m(n)+l} \|\hat{M}_{m(n)+l+1}\|_\infty, \quad (\text{C.5}) \end{aligned}$$

for some  $K_3 > 0$ . The second inequality is due to the Lipschitz continuity of  $h_{r_n}$  and condition (a.4) on  $\zeta_{n+1}$ , by the equivalence of vector norms. Moreover, by Lemma C.5, the third term on the RHS of (C.5) is bounded a.s. since  $\{\sum_{k=0}^{n-1} \gamma_k \hat{M}_{k+1}\}$  converges a.s. Recall that  $\sum_{l=0}^{k-m(n)} \gamma_{m(n)+l} \leq T + \sup_n \gamma_n < \infty$  by definition of  $m(n)$  and  $T_n$ . Thus, there exist  $K_4, K_5 > 0$  such that

$$\|\hat{x}(t_{k+1})\|_\infty \leq K_4 + K_5 \sum_{l=0}^{k-m(n)} \gamma_{m(n)+l} \|\hat{x}(t_{m(n)+l})\|_\infty.$$

By discrete-time Grönwall inequality, we have

$$\sup_{m(n) \leq k < m(n+1)} \|\hat{x}(t_k)\|_\infty \leq K_4 \cdot \exp[K_5 \cdot (T + \sup_n \gamma_n)], \quad (\text{C.6})$$

where the RHS of (C.6) is a (random) constant independent of  $n$ . Hence by equivalence of vector norms, we further obtain  $\sup_t \|\hat{x}(t)\| < \infty$ , which concludes the proof.  $\square$

The stability of  $\|\hat{x}(t)\|$  is essential in showing the convergence of the consensus update in (C.3). For  $n \geq 0$ , let  $y^n(t)$  denote the trajectory of  $\dot{y} = \tilde{h}_c(y)$  with  $c = r_n$  and  $y^n(T_n) = \langle \hat{x}(T_n) \rangle$ , for  $t \in [T_n, T_{n+1})$ . Then we have the following lemma.

**Lemma C.7.** Under Assumption C.1,  $\lim_n \sup_{t \in [T_n, T_{n+1})} \|\hat{x}(t) - \mathbf{1} \otimes y^n(t)\| = 0$ .

*Proof.* Since  $\hat{x}(t)$  is bounded a.s. on  $[T_n, T_{n+1})$ , we can mimic our proofs for Theorems 4.6 and 4.9 to show the convergence of  $\hat{x}(t_k)$  for  $k \in [T_n, T_{n+1})$  as  $n \rightarrow \infty$ . We will provide here only a sketch. One can first show that over the set  $\{\sup_k \|\hat{x}(t_k)\| \leq M\}$  for any  $M > 0$ ,  $\lim_k \|\mathcal{J}_\perp \hat{x}(t_k)\| = 0$ . The iteration of  $\mathcal{J}_\perp \hat{x}(t_k)$  has the form (similar to (B.9))

$$\mathcal{J}_\perp \hat{x}(t_{k+1}) = [(\mathbf{I} - \mathbf{1}\mathbf{1}^\top/N)C_k \otimes \mathbf{I}][\mathcal{J}_\perp \hat{x}(t_k) + \gamma_k y_{k+1}], \quad (\text{C.7})$$

where  $y_{k+1} = h_{r_n}[\hat{x}(t_k)] + \hat{M}_{k+1} + \hat{\zeta}_{k+1}$  here. One can easily verify that  $\mathbb{E}(\|y_{k+1}\|^2 \mathbb{I}_{\{\sup_k \|\hat{x}(t_k)\| \leq M\}} | \mathcal{F}_k) < K_6$  for some  $K_6 > 0$ , due to Lipschitz continuity of  $h_{r_n}$  and the conditions (a.3) and (a.4) in Assumption C.1. Hence, by similar arguments as in the proof of Lemma B.3, we obtain  $\lim_k \|\mathcal{J}_\perp \hat{x}(t_k)\| = 0$  almost surely, i.e., the vector  $\hat{x}(t_k)$  reaches consensus as  $k \rightarrow \infty$ . Then we proceed to show the convergence of the sequence  $\{\langle \hat{x}(t_k) \rangle\}$ . Define  $\hat{h}_c : \mathbb{R}^d \times A \rightarrow \mathbb{R}^d$  as  $\hat{h}_c(y, Y_k) = \langle h(c \cdot \mathbf{1} \otimes y) \cdot c^{-1} \rangle$ ; then the iteration can be written as follows

$$\begin{aligned} \langle \hat{x}(t_{k+1}) \rangle &= \langle \hat{x}(t_k) \rangle + \gamma_k \cdot \mathbb{E}(\langle y_{k+1} \rangle | \mathcal{F}_k) + \gamma_k \cdot \xi_{k+1} \\ &= \langle \hat{x}(t_k) \rangle + \gamma_k \cdot \hat{h}_{r_n}[\langle \hat{x}(t_k) \rangle, Y_k] + \gamma_k \cdot \xi_{k+1} + \gamma_k \cdot \beta_{k+1}, \end{aligned}$$

where  $\xi_{k+1} = \langle (C_k \otimes \mathbf{I})(y_{k+1} + \gamma_k^{-1} \mathcal{J}_\perp \hat{x}(t_k)) - \mathbb{E}(\langle y_{k+1} \rangle | \mathcal{F}_k) \rangle$ ,  $\beta_{k+1} = \mathbb{E}(\langle y_{k+1} \rangle | \mathcal{F}_k) - \hat{h}_{r_n}[\langle \hat{x}(t_k) \rangle, Y_k]$ , and  $\langle y_{k+1} \rangle = \langle h_{r_n}[\hat{x}(t_k)] \rangle + \hat{M}_{k+1} + \langle \hat{\zeta}_{k+1} \rangle$ . One can verify that  $\{\xi_{k+1}\}$  is a martingale difference sequence satisfying  $\mathbb{E}(\|\xi_{t+1}\|^2 | \mathcal{F}_{t,1}) < K_7 \cdot (1 + \|\langle \hat{x}(t_k) \rangle\|^2)$  for some  $K_7 < \infty$  over the set  $\{\sup_k \|\hat{x}(t_k)\| \leq M\}$ . In addition, note that  $\mathbb{E}(\hat{M}_{k+1} | \mathcal{F}_k) = 0$  and thus  $\mathbb{E}(\langle y_{k+1} \rangle | \mathcal{F}_k) = \langle h[\hat{x}(t_k), Y_k] \rangle \cdot r_n^{-1}$ . Thus we have  $\|\beta_{k+1}\| \leq L \cdot \|\mathcal{J}_\perp \hat{x}(t_k)\| \cdot r_n^{-1}$  for some  $L < \infty$  due to the Lipschitz continuity of  $h$ . Hence  $\beta_k \rightarrow 0$  a.s. since  $\|\mathcal{J}_\perp \hat{x}(t_k)\| \rightarrow 0$  a.s. and  $r_n \geq 1$ . Moreover,  $\hat{h}_{r_n}[\langle \hat{x}(t_k) \rangle, Y_k]$  is Lipschitz continuous in  $\langle \hat{x}(t_k) \rangle$ . Therefore, by Theorem D.2, we obtain that  $\langle \hat{x}(t_k) \rangle \rightarrow y^n(t)$  as  $n \rightarrow \infty$ , namely  $k \rightarrow \infty$ . Further we obtain that  $\hat{x}^i(t_k) \rightarrow y^n(t)$  for any  $i \in \mathcal{N}$ , which concludes the proof following Theorem 2 in Chapter 2 of Borkar (2008).  $\square$

Now suppose that  $\|\bar{x}_n^i\| \rightarrow \infty$  for some  $i \in \mathcal{N}$ ; then there exists a subsequence of  $\{n_q\}$  such that  $\|\bar{x}^i(T_{n_q})\| \rightarrow \infty$  and thus  $\|\bar{x}(T_{n_q})\| \rightarrow \infty$ . Hence  $r_{n_q} \rightarrow \infty$ . If  $r_n > c_0 \geq 1$ , then  $\|\hat{x}(T_n)\| = 1$ , and thus  $\|y^n(T_n)\| = \|\langle \hat{x}(T_n) \rangle\| \leq N^{-1/2}$ . By Lemma C.3, we have  $\|\mathbf{1} \otimes y^n(T_n^-)\| = N^{1/2} \cdot \|y^n(T_n^-)\| \leq 1 - \epsilon'$ . Thus by Lemma C.7, we have  $\|\hat{x}(T_{n+1}^-)\| < 1 - \epsilon''$  for some  $0 < \epsilon'' < \epsilon'$ . Hence for  $r_n > c_0$  and sufficiently large  $n$ ,

$$\frac{\|\bar{x}(T_{n+1})\|}{\|\bar{x}(T_n)\|} = \frac{\|\hat{x}(T_{n+1}^-)\|}{\|\hat{x}(T_n)\|} < 1 - \epsilon''.$$

It thus follows that if  $\|\bar{x}(T_n)\| > 1$ ,  $\|\bar{x}(T_k)\|$  falls back to the unit ball at an exponential rate for  $k \geq n$ . The rest of the argument follows directly from the proof of Theorem 2 in [Mathkar & Borkar \(2016\)](#), which concludes the proof.  $\square$

Now we are ready to apply Theorem C.2 to prove Lemmas B.1 and B.4. We will return to the notations in §4.

### C.2. Proof of Lemma B.1

The proof follows by verifying the conditions for Theorem C.2 to hold. Recall that the critic step in (3.3) has the form

$$\omega_{t+1} = (C_t \otimes I)(\omega_t + \beta_{\omega,t} \cdot y_{t+1}),$$

with  $y_{t+1} = (\delta_t^1 \phi_t^\top, \dots, \delta_t^N \phi_t^\top)^\top \in \mathbb{R}^{KN}$ . Thus the terms corresponding to (C.1) are

$$h^i(\omega_t^i, \mu_t^i, s_t, a_t) = \mathbb{E}(\delta_t^i \phi_t^\top | \mathcal{F}_{t,1}), \quad M_{t+1}^i = \delta_t^i \phi_t^\top - \mathbb{E}(\delta_t^i \phi_t^\top | \mathcal{F}_{t,1}). \quad (\text{C.8})$$

Since the Markov chain  $\{(s_t, a_t)\}_{t \geq 0}$  is irreducible and aperiodic given policy  $\pi_\theta$ , we have  $\bar{h}^i(\omega_t^i, \mu_t^i) = \Phi^\top D_\theta^{s,a} [R^i - \mu_t^i \mathbf{1} + P^\theta \Phi \omega_t^i - \Phi \omega_t^i]$ . By Lemma B.2, it is established that  $\{\mu_t^i\}$  is bounded a.s. Hence, over the set  $\{\sup_t \|\mu_t\| \leq M\}$  for any  $M > 0$ , there exists some  $K' > 0$  such that  $\|\bar{h}(\omega_t, \mu_t) - h(\omega_t, \mu_t, s_t, a_t)\|^2 \leq K' \cdot (1 + \|\omega_t\|^2)$ , since the Markov chain is finite. This verifies the condition (a.4) in Assumption C.1. Moreover, since  $r_{t+1}^i$  and  $\|\phi_t\|$  are uniformly bounded,  $\mathbb{E}(\|M_{t+1}\|^2 | \mathcal{F}_{t,1}) \leq K \cdot (1 + \|\omega_t\|^2)$  is also verified for some  $K > 0$ . More importantly, over the set  $\{\sup_t \|\mu_t\| \leq M\}$ ,  $h_\infty(y)$  exists and has the form

$$h_\infty(y) = \lim_c \tilde{h}_c(y) = \Phi^\top D_\theta^{s,a} (P^\theta - I) \Phi y. \quad (\text{C.9})$$

Clearly  $\dot{y} = h_\infty(y)$  has origin as its globally asymptotically stable attractor (see the proof of Theorem 4.6 in §B.2). Hence we apply Theorem C.2 to conclude the proof.  $\square$

### C.3. Proof of Lemma B.4

Recall that the critic step from (3.8) and (3.11) has the compact form

$$z_{t+1} = (C_t \otimes I)(z_t + \beta_{v,t} \cdot y_{t+1}), \quad (\text{C.10})$$

where  $z_t^i = [\mu_t^i, (\lambda_t^i)^\top, (v_t^i)^\top]^\top$  and  $y_t = [(y_t^1)^\top, \dots, (y_t^N)^\top]^\top \in \mathbb{R}^{(1+M+L)N}$ . Here  $y_{t+1}^i$  denotes  $y_{t+1}^i = [r_{t+1}^i - \mu_t^i, (r_{t+1}^i - f_t^\top \lambda_t^i) f_t^\top, \delta_t^i \varphi_t^\top]^\top$ . Thus the terms corresponding to (C.1) are

$$h^i(z_t^i, s_t, a_t) = \mathbb{E}(y_{t+1}^i | \mathcal{F}_{t,1}), \quad M_{t+1}^i = y_{t+1}^i - \mathbb{E}(y_{t+1}^i | \mathcal{F}_{t,1}). \quad (\text{C.11})$$

Furthermore, we have

$$\bar{h}^i(z_t^i) = \begin{pmatrix} -1 & 0 & 0 \\ 0 & -F^\top D_\theta^{s,a} F & 0 \\ -\Phi^\top D_\theta^s \mathbf{1} & 0 & \Phi^\top D_\theta^s (P^\theta - I) \Phi \end{pmatrix} \begin{pmatrix} \mu_t^i \\ \lambda_t^i \\ v_t^i \end{pmatrix} + \begin{pmatrix} J^i(\theta) \\ F^\top D_\theta^{s,a} R^i \\ \Phi^\top D_\theta^s R_\theta^i \end{pmatrix},$$

where  $J^i(\theta) = \sum_{s \in \mathcal{S}, a \in \mathcal{A}} d_\theta(s, a) \cdot R^i(s, a)$  and  $R_\theta^i(s) = \sum_a \pi_\theta(s, a) R^i(s, a)$ . Therefore, one can verify that both conditions (a.3) and (a.4) in Assumption C.1 are satisfied. In addition,  $h_\infty(y)$  exists and has the form

$$h_\infty(y) = \lim_c \tilde{h}_c(y) = \begin{pmatrix} -1 & 0 & 0 \\ 0 & -F^\top D_\theta^{s,a} F & 0 \\ -\Phi^\top D_\theta^s \mathbf{1} & 0 & \Phi^\top D_\theta^s (P^\theta - I) \Phi \end{pmatrix} \cdot y.$$

Clearly  $\dot{y} = h_\infty(y)$  has origin as its globally asymptotically stable attractor (see the proof of Theorem 4.9 in §B.4), which completes the proof.  $\square$

## D. Technical Background

### D.1. A Basic Result of Stochastic Approximation

For the sake of completeness, we reproduce here a key result from [Borkar \(2008\)](#) that has been used repeatedly in our proofs. The results follow by specializing Corollary 8 and Theorem 9 on page 74-75 in [Borkar \(2008\)](#). We note that this is actually an extension of Theorem 2.1 and Theorem 2.2 in [Borkar & Meyn \(2000\)](#) to the case with irreducible Markovian state and diminishing noise in the update. More general conclusions can also be found in [Benaïm \(1999\)](#); [Kushner & Yin \(2003\)](#).

Consider the  $n$ -dimensional stochastic approximation iteration

$$x_{t+1} = x_t + \gamma_t [h(x_t, Y_t) + M_{t+1} + \beta_{t+1}], \quad t \geq 0, \quad (\text{D.1})$$

where  $\gamma_t > 0$  and  $\{Y_t\}_{t \geq 0}$  is a Markov chain on a finite set  $A$ .

**Assumption D.1.** We make the following assumptions:

- (a.1)  $h : \mathbb{R}^n \times A \rightarrow \mathbb{R}^n$  is Lipschitz continuous in its first argument;
- (a.2)  $\{Y_t\}_{t \geq 0}$  is an irreducible Markov chain with stationary distribution  $\pi$ ;
- (a.3) The stepsize sequence  $\{\gamma_t\}$  satisfies  $\sum_t \gamma_t = \infty$  and  $\sum_t \gamma_t^2 < \infty$ ;
- (a.4)  $\{M_t\}$  is a martingale difference sequence, i.e.,  $\mathbb{E}[M_{t+1} | x_\tau, M_\tau, Y_\tau, \tau \leq t] = 0$ , satisfying that for some  $K > 0$  and  $t \geq 0$

$$\mathbb{E}(\|M_{t+1}\|^2 | x_\tau, M_\tau, Y_\tau, \tau \leq t) \leq K \cdot (1 + \|x_t\|^2).$$

- (a.5) The sequence  $\{\beta_t\}$  is a bounded random sequence with  $\beta_t \rightarrow 0$  almost surely as  $t \rightarrow \infty$ .

Then the asymptotic behavior of the iteration (D.1) is related to the behavior of the solution to the ODE

$$\dot{x} = \bar{h}(x) = \sum_i \pi(i) h(x, i). \quad (\text{D.2})$$

Suppose (D.2) has a unique globally asymptotically stable equilibrium  $x^*$ , we then have the following two theorems.

**Theorem D.2.** Under Assumption D.1, if  $\sup_t \|x_t\| < \infty$  a.s., we have  $x_t \rightarrow x^*$ .

**Theorem D.3.** Under Assumption D.1, suppose that

$$\lim_{c \rightarrow \infty} \frac{\bar{h}(cx)}{c} = h_\infty(x)$$

exists uniformly on compact sets for some  $h_\infty \in C(\mathbb{R}^n)$ . If the ODE  $\dot{y} = h_\infty(y)$  has origin as the unique globally asymptotically stable equilibrium, then

$$\sup_t \|x_t\| < \infty \text{ a.s.}$$

### D.2. Kushner-Clark Lemma

We state here the well-known Kushner-Clark Lemma ([Kushner & Clark, 1978](#); [Metivier & Priouret, 1984](#); [Prasad et al., 2014](#)) in the sequel.

Let  $\Gamma$  be an operator that projects a vector onto a compact set  $\mathcal{X} \subseteq \mathbb{R}^N$ . Define a vector  $\hat{\Gamma}(\cdot)$  as

$$\hat{\Gamma}[h(x)] = \lim_{0 < \eta \rightarrow 0} \left\{ \frac{\Gamma[x + \eta h(x)] - x}{\eta} \right\},$$

for any  $x \in \mathcal{X}$  and with  $h : \mathcal{X} \rightarrow \mathbb{R}^N$  continuous. Consider the following recursion in  $N$  dimensions

$$x_{t+1} = \Gamma\{x_t + \gamma_t [h(x_t) + \xi_t + \beta_t]\}. \quad (\text{D.3})$$

The ODE associated with (D.3) is given by

$$\dot{x} = \hat{\Gamma}[h(x)]. \quad (\text{D.4})$$

**Assumption D.4.** We make the following assumptions:

- (a.1)  $h(\cdot)$  is a continuous  $\mathbb{R}^N$ -valued function.
- (a.2) The sequence  $\{\beta_t\}, t \geq 0$  is a bounded random sequence with  $\beta_t \rightarrow 0$  almost surely as  $t \rightarrow \infty$ .
- (a.3) The stepsizes  $\gamma_t, t \geq 0$  satisfy  $\gamma_t \rightarrow 0$  as  $t \rightarrow \infty$  and  $\sum_t \gamma_t = \infty$ .
- (a.4) The sequence  $\xi_t, t \geq 0$  satisfies for any  $\epsilon > 0$

$$\lim_t \mathbb{P} \left( \sup_{n \geq t} \left\| \sum_{\tau=t}^n \gamma_\tau \xi_\tau \right\| \geq \epsilon \right) = 0.$$

Then the Kushner-Clark Lemma says the following.

**Theorem D.5.** Under Assumption D.4, suppose that the ODE (D.4) has a compact set  $\mathcal{K}^*$  as its set of asymptotically stable equilibria. Then  $x_t$  in (D.3) converges almost surely to  $\mathcal{K}^*$  as  $t \rightarrow \infty$ .



## E. Experiment Details

In this section, we provide the details of the simulation settings in §5.

### E.1. Linear Function Approximation

We first consider the setting with linear function approximation to corroborate our theoretical results. Consider in total  $N = 20$  agents, each has a binary-valued action space, i.e.,  $\mathcal{A}^i = \{0, 1\}$ , for all  $i \in \mathcal{N}$ . Thus the cardinality of the set of actions  $\mathcal{A}$  is  $2^{20}$ . In addition, there are in total  $|\mathcal{S}| = 20$  states. The following ways of selecting the model and algorithm parameters, including transition probabilities, rewards, and features, follow from those in Dann et al. (2014). The elements in the transition probability matrix  $P$  are uniformly sampled from the interval  $[0, 1]$  and normalized to be stochastic. We also add a small constant  $10^{-5}$  onto each element in the matrix to ensure ergodicity of the MDP such that Assumption 2.2 is satisfied. For each agent  $i$  and each state-action pair  $(s, a)$ , the mean reward  $R^i(s, a)$  is sampled uniformly from  $[0, 4]$ , which varies among agents. The instantaneous rewards  $r_t^i$  are sampled from the uniform distribution  $[R^i(s, a) - 0.5, R^i(s, a) + 0.5]$ . The policy  $\pi_{\theta^i}^i(s, a^i)$  is parametrized following the Boltzman policies, i.e.,

$$\pi_{\theta^i}^i(s, a^i) = \frac{\exp(q_{s,a^i}^\top \theta^i)}{\sum_{b^i \in \mathcal{A}^i} \exp(q_{s,b^i}^\top \theta^i)}$$

where  $q_{s,b^i} \in \mathbb{R}^{m_i}$  is the feature vector with the same dimension as  $\theta^i$ , for any  $s \in \mathcal{S}$  and  $i \in \mathcal{N}$ . Here we set  $m_1 = m_2 = \dots = m_N = 5$ . The elements of  $q_{s,b^i}$  are also uniformly sampled from  $[0, 1]$ . In particular, the gradient of the score function thus has the form

$$\nabla_{\theta^i} \log \pi_{\theta^i}^i(s, a^i) = q_{s,a^i} - \sum_{b^i \in \mathcal{A}^i} \pi_{\theta^i}^i(s, a^i) q_{s,b^i}.$$

The feature vectors  $\phi \in \mathbb{R}^K$  for the action-value function  $Q(\cdot, \cdot; \omega)$  in Algorithm 1,  $\varphi \in \mathbb{R}^L$  for the state-value function  $V(\cdot; v)$  and  $f \in \mathbb{R}^M$  for the globally averaged reward function  $\bar{R}(\cdot, \cdot; \lambda)$  in Algorithm 2, are all uniformly sampled from  $[0, 1]$ , of dimensions  $K = 10 \ll |\mathcal{S}| \cdot |\mathcal{A}|$ ,  $L = 5 < |\mathcal{S}|$ , and  $M = 10 \ll |\mathcal{S}| \cdot |\mathcal{A}|$ . Moreover, the selected feature matrices  $\Phi$ ,  $\bar{\Phi}$ , and  $F$  are all ensured to have full column rank as required in Assumptions 4.5 and 4.8.

The consensus weight matrix  $C_t$  are chosen independent and identically distributed along time  $t$ , by normalizing the absolute Laplacian matrix of a connected random graph  $\mathcal{G}_t$  over agents  $\mathcal{N}$  to be doubly stochastic<sup>3</sup>. The graph  $\mathcal{G}_t$  is generated by randomly placing communication links among agents such that the connectivity ratio<sup>4</sup> is  $4/N$ . The stepsizes are selected as  $\beta_{\omega,t} = \beta_{v,t} = 1/t^{0.65}$  and  $\beta_{\theta,t} = 1/t^{0.85}$ , which satisfy Assumption 4.3.

The performances of the fully decentralized algorithms are compared with those of the centralized algorithms in which the rewards  $r_t^i$  of all agents are available at a centralized controller and the global policy  $\pi_\theta$  is also updated there. These centralized version algorithms thus reduce to single-agent AC algorithms with linear function approximation. We refer the two centralized AC algorithms for comparison as Central-1 and Central-2, respectively. The algorithm Central-1 has the following critic step, which is based on action-value function approximation as in Algorithm 1

$$\begin{cases} \mu_{t+1} = (1 - \beta_{\omega,t}) \cdot \mu_t + \beta_{\omega,t} \cdot \bar{r}_{t+1}, \\ \delta_t = \bar{r}_{t+1} - \mu_t + Q_{t+1}(\omega_t) - Q_t(\omega_t) \\ \omega_{t+1} = \omega_t + \beta_{\omega,t} \cdot \delta_t \cdot \nabla_\omega Q_t(\omega_t). \end{cases} \quad (\text{E.1})$$

Recall that  $Q_t(\omega) = Q(s_t, a_t; \omega)$  with  $Q(\cdot, \cdot; \omega)$  the estimate for the global action-value function  $Q_\theta$ ,  $\beta_{\omega,t} > 0$  is the stepsize, and  $\bar{r}_t = \sum_{i \in \mathcal{N}} r_t^i \cdot N^{-1}$ . Accordingly, the central controller improves the policy for each agent  $i$ , which results in the actor step

$$\theta_{t+1}^i = \theta_t^i + \beta_{\theta,t} \cdot A_t \cdot \psi_t^i, \quad \forall i \in \mathcal{N}, \quad (\text{E.2})$$

<sup>3</sup>A stochastic matrix  $P$  is doubly stochastic if it is both row and column stochastic.

<sup>4</sup>Note that the connectivity ratio is defined as the ratio between the total degree of the graph and the degree of the complete graph, i.e.,  $2E/[N(N-1)]$ , where  $E$  is the number of edges.

where  $\beta_{\theta,t} > 0$  is the stepsize,  $A_t$  and  $\psi_t^i$  are defined as

$$A_t = Q_t(\omega_t) - \sum_{a \in \mathcal{A}} \pi_{\theta_t}(s_t, a) \cdot Q(s_t, a; \omega_t), \quad \psi_t^i = \nabla_{\theta^i} \log \pi_{\theta_t^i}(s_t, a_t^i). \quad (\text{E.3})$$

The algorithm Central-2 follows the updates of Algorithm 1 in [Bhatnagar et al. \(2009\)](#), based on state-value approximation as in Algorithm 2 here. In particular, it has the following critic step

$$\begin{cases} \mu_{t+1} = (1 - \beta_{v,t}) \cdot \mu_t + \beta_{v,t} \cdot \bar{r}_{t+1}, \\ \delta_t = \bar{r}_{t+1} - \mu_t + V_{t+1}(v_t) - V_t(v_t), \\ v_{t+1} = v_t + \beta_{v,t} \cdot \delta_t \cdot \nabla_v V_t(v_t), \end{cases} \quad (\text{E.4})$$

where we recall that  $V_t(v) = V(s_t; v)$  for any  $v \in \mathbb{R}^L$  and  $\beta_{v,t} > 0$  is the stepsize satisfying Assumption 4.3. Since the rewards of all agents  $\{r_t^i\}_{i \in \mathcal{N}}$  are available to the controller, no estimation for the globally averaged

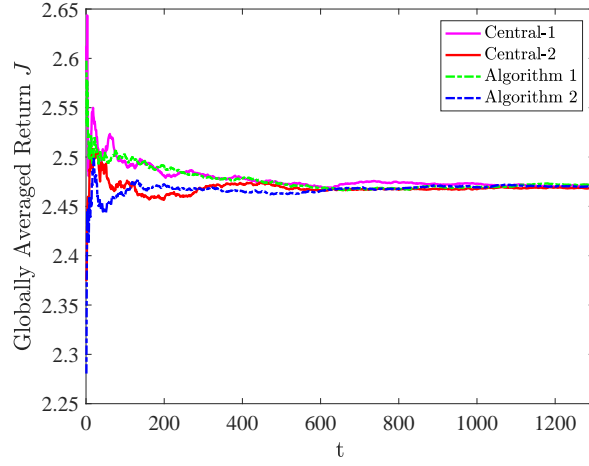


Figure 3. The convergence of globally averaged returns, when linear function approximation is used. We plot the returns achieved by both Algorithm 1 and Algorithm 2, along with their centralized counterparts Central-1 and Central-2.

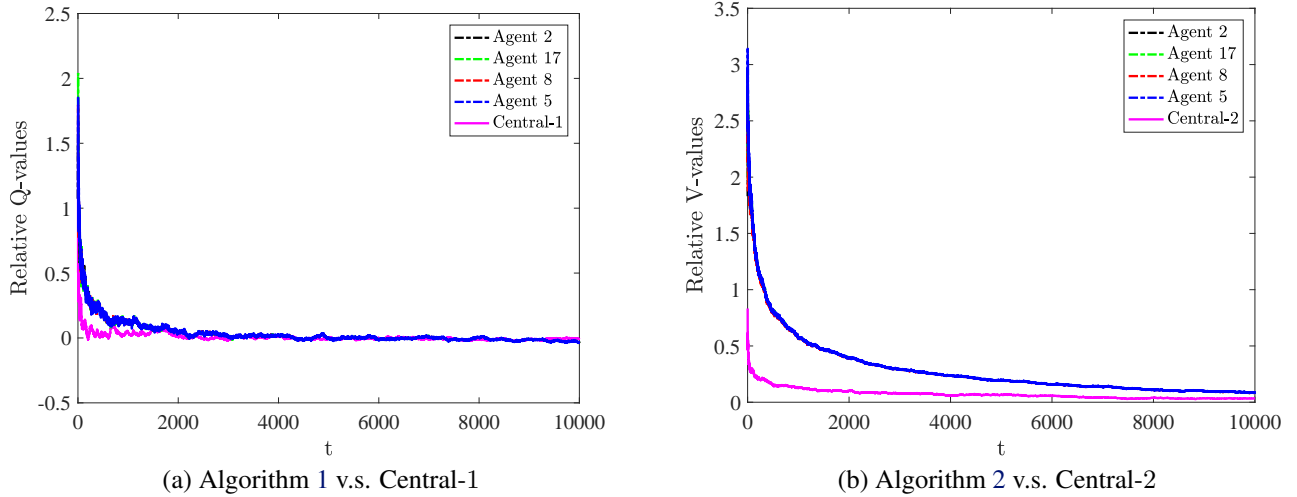


Figure 4. The convergence of relative value functions at four randomly selected agents, when linear function approximation is used. We randomly select the agents 2, 5, 8, and 17. In (a), we plot the convergence curve of the relative action-value at a randomly selected state-action pair, obtained from Central-1 and Algorithm 1. In (b), we plot the convergence curve of the relative state-value at a randomly selected state, obtained from Central-2 and Algorithm 2.

reward function  $\bar{R}$  is needed in the update. Thus, the global state-value TD-error  $\delta_t$  can be computed immediately and used in the actor step as

$$\theta_{t+1}^i = \theta_t^i + \beta_{\theta,t} \cdot \delta_t \cdot \psi_t^i, \quad \forall i \in \mathcal{N}, \quad (\text{E.5})$$

where  $\psi_t^i$  is as defined in (E.3).

Convergence of the globally averaged returns, the relative action-value function, and the relative state-value function are reported in Figure 3 and Figure 4, respectively. Figure 3 shows that both decentralized algorithms converge to the globally averaged return as achieved by the two centralized counterparts. Moreover, Figure 4 illustrates that for each agent, the approximation of global value functions in Algorithms 1 and 2 reach consensus much faster than the AC algorithm itself converges. In addition, the convergence of the value functions is relatively slower for the decentralized algorithms than for the centralized ones, possibly due to the delay of information diffusion across the network.

## E.2. Nonlinear Function Approximation

We also evaluate the performance of Algorithm 1 and Algorithm 2 when nonlinear function approximators, for example, neural networks, are adopted. Although it seems difficult to establish convergence guarantees in this case, we believe that the empirical results are of independent interest, which justify the effectiveness of the proposed fully decentralized algorithms in a more sophisticated environment.

To this end, we consider the simulation environment of the *Cooperative Navigation* task in Lowe et al. (2017). In this environment, agents need to reach a set of  $L$  landmarks through physical movement. Agents are able to observe the position of the landmarks and other agents, and are rewarded based on the proximity of any agent to each landmark (Lowe et al., 2017). To fit in our networked MDP model, we modify the environment there in the following aspects. First, we assume the state is globally observable, i.e., the position of the landmarks and other agents are observable to each agent. Moreover, each agent has a certain target landmark to cover, and the individual reward is determined by the proximity to that certain landmark, as well as the penalty from collision with other agents. In this way, the reward function varies between agents. The reward is further scaled by different positive coefficients, representing the different priority/preferences of different agents. In addition, agents are connected via a time-varying communication network with several other agents nearby. The collaborative goal of the agents is then to maximize the network-wide averaged long-term return. The illustration of the modified Cooperative Navigation environment is provided in Figure 5.

Specifically, we consider  $N = 10$  agents moving in a rectangular region of size  $2 \times 2$ . Each agent has a single target landmark, i.e.,  $L = N = 10$ , which is randomly located in the region. The action set for each agent is the

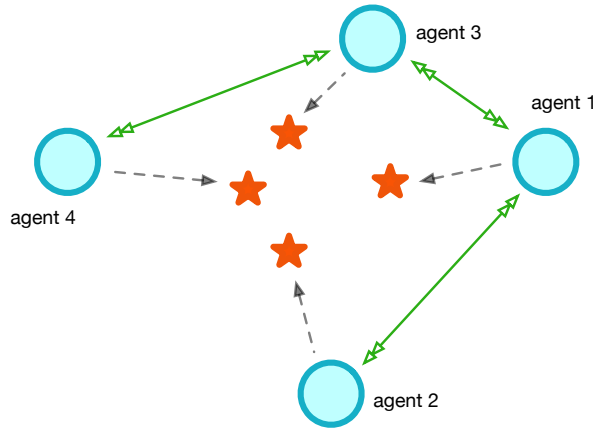


Figure 5. Illustration of the experimental environment for the Cooperative Navigation task we consider, modified from Lowe et al. (2017). In particular, the blue circles represent the agents, the orange stars represent the landmarks, the green arrows represent the communication links between agents, and the gray arrows show the target landmark each agent need to cover.

movement set  $\{left, right, up, down, stay\}$ , and thus  $|\mathcal{A}_i| = 5$  for any  $i \in \mathcal{N}$ . The state  $s$  includes the position of the landmarks and other agents, which thus has a dimension of  $2(N + L) = 40$ . The reward of agent  $i$  is the negative number of the distance to the target landmark, plus  $-1$  if agent  $i$  collides with any other agents. The coefficients that scale the reward of each agent are selected randomly from a uniform distribution over  $[0, 2]$ . Each agent maintains two neural networks for actor and critic, respectively. Both neural networks have one hidden layer containing 24 neural units, which all use ReLU as the activation function. The output layer for the actor network is softmax, and that for the critic network is linear.

The time-varying network  $\mathcal{G}_t$  and consensus matrix  $C_t$  are constructed in the same way as in §E.2. The stepsizes for the actor and critic step are set as constants 0.001 and 0.01, respectively. For each episode, the algorithms terminate either all agents reach the target landmarks or after 1000 iterations, and we run in total 200 episodes in each test run. We report the globally averaged return from 10 test runs in Figure 2, where the algorithms Central-1 and Central-2 follow the update rules in §E.1, but with nonlinear function approximation.

To better illustrate the necessity of cooperation via communication among agents, we compare the results of Algorithms 1 and 2 with the *non-cooperative* counterparts, where each agent performs single-agent RL, with no awareness of the existence of other agents in the environment. These non-cooperative algorithms are equivalent to our proposed consensus-based algorithms, when the communication network is disconnected and contains only a self-loop at each node. Thus, we demonstrate the performance of the non-cooperative counterparts of Algorithms 1 and 2 in Figure 6. Note that the results of Algorithms 1 and 2 follow from those in Figure 2. It is shown that the non-cooperative algorithms are unstable, and achieve much worse long-term return with much larger variance than the cooperative counterparts. These observations showcase the necessity of communication in fully decentralized cooperative MARL.

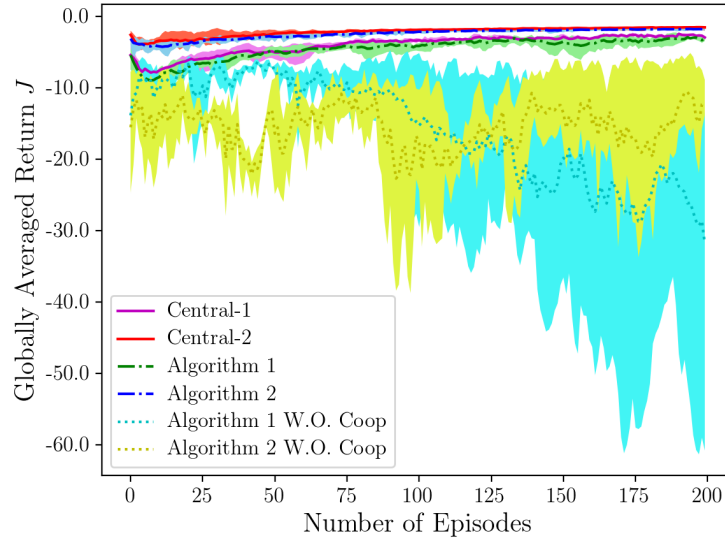


Figure 6. The globally averaged returns for Cooperative Navigation, when neural networks are used for function approximation.

## F. Comparison with Existing Work on Multi-Agent Systems and MARL

In this section, we compare both our *model* and *algorithms* with related work on multi-agent systems and collaborative MARL in details.

Our framework of networked multi-agent systems finds a broad range of applications in distributed cooperative control problems, including formation control of unmanned vehicles (Fax & Murray, 2004), cooperative navigation of robots (Corke et al., 2005), load management in energy networks (Dall’Anese et al., 2013), and flocking of mobile sensor networks (Cortes et al., 2004), etc. Previously, the collective goal of the multi-agent system is to either reach a stable and consensus state for all agents (Fax & Murray, 2004; Corke et al., 2005), or solve a static optimization problem in a distributed fashion (Dall’Anese et al., 2013; Nedic & Ozdaglar, 2009). In the first line of work, including formation control and consensus problems, the objective is not formulated explicitly as an optimization problem, and most of the work focuses on continuous-time dynamic systems. Whereas in the second line of work, the problem is approached in a static setting, in the sense that the optimization objective is deterministic and there is no control input affecting the transition of the system, see recent efforts in Nedic & Ozdaglar (2009); Agarwal & Duchi (2011); Chen & Sayed (2012). In contrast, we here model the interaction of multiple agents and evolution of the system as an MDP, a dynamic setting, and explicitly use the network-wide long-term return as the collaborative goal of all agents. In this regard, our framework is pertinent to the cooperative/distributed optimal control problems (Lewis et al., 2013; Movric & Lewis, 2014), but focuses on the discrete-time setting and falls into the realm of reinforcement learning, where the model of the system may be unknown. One recent work (Macua et al., 2017) for multi-task RL, which is almost concurrent to ours, is also based on the model with networked agents. Nonetheless, the MDP problem solved by different agents are totally decoupled, which excludes the work from the realm of MARL with interactive agents as we consider here.

Our framework also departs from the existing framework on collaborative MARL models in the following aspects. In contrast to the canonical multi-agent MDP (MMDP) model proposed in Boutilier (1996); Lauer & Riedmiller (2000), our model allows the agents to exchange information over a communication network with possibly sparse connectivity at each agent. This improves the *scalability* of the multi-agent model with a high population of agents, which is one of the long-standing challenges in general MARL problems (Shoham et al., 2003). Moreover, we allow heterogeneous agents to have various individual reward functions, while the canonical MMDP assumes a common reward function for all agents. The latter setting greatly simplifies the problem since no information exchange among agents is necessary to approximate the value function for each agent. Our model not only fits in the *multi-task* setting which has gained increasing popularity in MARL (Omidshafiei et al., 2017; Teh et al., 2017), but also applies to the general multi-agent RL setting. One of the few models that also consider heterogeneous reward functions in collaborative MARL is Kar et al. (2013), where the global action is assumed to be actuated by a remote controller, but in our case, the agents are fully decentralized and have local control capabilities. Besides, the models in Guestrin et al. (2002); Kok & Vlassis (2006) also consider heterogeneous rewards, but with a strong assumption that the global Q-function can be factorized as several local Q-functions that depend on the actions of only a subset of agents, which simplified the general setting we considered where the Q-function is affected by the joint action of all agents. It is also worth noting that our model generalizes the team Markov game model for collaborative MARL, see Littman (2001); Wang & Sandholm (2003); Arslan & Yüksel (2017), where all agents have individual action sets but share a common payoff function as in the canonical MMDP.

Moreover, our algorithms designed for networked MMDP are distinct from the existing collaborative MARL algorithms in the following aspects. First, our MARL algorithms belong to the type of actor-critic algorithms, whereas several of the existing MARL algorithms are designed based on Q-learning type (critic-based) algorithms only (Boutilier, 1996; Lauer & Riedmiller, 2000; Guestrin et al., 2002; Kok & Vlassis, 2006; Kar et al., 2013). Moreover, these algorithms assume either the rewards are common to all agents (Boutilier, 1996; Lauer & Riedmiller, 2000), or there exists a remote central controller to take actions for the agents (Kar et al., 2013). The rest of them (Guestrin et al., 2002; Kok & Vlassis, 2006) utilize the factorized structure of the Q-function as mentioned above, which fail to handle our general MARL setting. More recently, some actor-critic type MARL algorithms with *distributed/decentralized* structures have gained increasing attention (Gupta et al., 2017; Lowe et al., 2017; Omidshafiei et al., 2017). They are developed for more complicated settings where both cooperation and competition may appear among agents. However, they all rely on a central controller to perform

the critic step, which are closer to the *hierarchical* structure rather than the *fully decentralized* structure we consider here. Second, our algorithms apply function approximation to handle the setting with massively large state and action spaces, while enjoying theoretical guarantees for convergence as we show in §4. However, the existing collaborative MARL algorithms are either guaranteed to converge only for tabular cases (Hu & Wellman, 2003; Wang & Sandholm, 2003; Kar et al., 2013; Prasad et al., 2014), or only have empirical convergence when function approximation is applied (Foerster et al., 2016; Gupta et al., 2017; Lowe et al., 2017; Omidshafiei et al., 2017). The recent work on multi-task RL with networked agents (Macua et al., 2017) also focuses on empirical results only, with no complete convergence analysis.