

Important Instructions: The appendix has an independent reference list.

APPENDIX A

PSEUDO-CODES FOR CENTRALIZED DQL, MULTI-AGENT DQL, AND KNOWLEDGE DISTILLATION

Algorithm 2 shows the pseudocode of centralized DQL, Algorithm 3 shows the pseudocode of multi-agent DQL, and Algorithm 4 shows the pseudocode of knowledge distillation.

Algorithm 2 Centralized deep Q-learning

- 1: Initialize replay memory \mathcal{V} .
 - 2: Initialize Q-network at the central controller with random weights θ
 - 3: Initialize target Q-network at the central controller with weights $\theta^- = \theta$
 - 4: **for** episode = 1, 2, ... **do**
 - 5: Initialize the state of edge computing system $s_1 \in \mathcal{S}$
 - 6: **for** $t = 1, 2, \dots$ in the episode **do**
 - 7: With probability ϵ_p select a random system rental decision $\mathbf{a}_t \in \mathcal{A}$
 - 8: Otherwise select $\mathbf{a}_t = \arg \max_{\mathbf{a}} Q(s_t, \mathbf{a}; \theta)$
 - 9: Central controller sends resource rental decisions to edge servers and the edge servers execute the decisions.
 - 10: Observe the system reward r_t and state s_{t+1} .
 - 11: Store the experience $(s_t, \mathbf{a}_t, r_t, s_{t+1})$ in \mathcal{V}
 - 12: Sample random minibatch $(s_j, \mathbf{a}_j, r_j, s_{j+1})$ from \mathcal{V}
 - 13: Set $y_j = \begin{cases} r_j, & \text{if episode ends at } t+1 \\ r_j + \gamma \max_{\mathbf{a}} Q(s_{t+1}, \mathbf{a}; \theta^-), & \text{otherwise} \end{cases}$
 - 14: LSM n Perform a gradient descent step on $(y_j - Q(s_j, \mathbf{a}_j; \theta))^2$ with respect to the network parameters θ
 - 15: Every C steps set $\theta^- = \theta$
 - 16: **end for**
 - 17: **end for**
-

Algorithm 3 Multi-agent deep Q-learning

- 1: Initialize replay memory \mathcal{V}_n for each LSM n
 - 2: Initialize Q-network at each LSM n with random weights θ_n
 - 3: Initialize target Q-network at each LSM n with weights $\theta_n^- = \theta_n$
 - 4: **for** episode = 1, 2, ... **do**
 - 5: Initialize the edge computing system and each LSM obtains its local observation o_n
 - 6: **for** $t = 1, 2, \dots$ in the episode **do**
 - 7: With probability ϵ_p , LSM n selects a random resource rental decision $a_{n,t} \in \mathcal{A}_n$
 - 8: Otherwise LSM n determines its rental decision $a_{n,t} = \pi_n(o_n; \theta_n)$ based on $\alpha_{n,t} = \arg \max_{\alpha_n} Q_n(s_t, \alpha_n; \theta_n)$
 - 9: Each LSM n configures the computing resource according to the rental decision \mathbf{a}_n .
 - 10: Each LSM n observes the rental decisions of nearby edge servers $\alpha_{n,t}$, reward $r_{n,t}$, and new observation $o_{n,t+1}$
 - 11: Store the experience $(o_{n,t}, \alpha_{n,t}, r_{n,t}, o_{n,t+1})$ of LSM n in \mathcal{V}_n
 - 12: Each LSM n samples random mini-batch $(o_{n,j}, \alpha_{n,j}, r_{n,j}, o_{n,j+1})$ from \mathcal{V}_n
 - 13: Set $y_{n,j} = \begin{cases} r_{n,j}, & \text{if episode ends at } t+1 \\ r_{n,j} + \gamma \max_{\alpha_n} Q_n(o_{t+1}, \alpha_n; \theta_n^-), & \text{otherwise} \end{cases}$
 - 14: Each LSM n performs gradient descent on $(y_{n,j} - Q_n(o_{n,j}, \alpha_{n,j}; \theta_n))^2$ with respect to θ_n
 - 15: Each LSM n sets $\theta_n^- = \theta_n$ every C steps
 - 16: **end for**
 - 17: **end for**
-

APPENDIX B

PERFORMANCE OF MULTI-AGENT DQL AND STANDARD DQL

Multi-agent DQL runs in a distributed fashion and the learning process locally on each edge server resembles the standard Deep Q-Learning (DQL) except that each DNN outputs a localized decision and observes the service provision decisions taken by one-hop neighbors as part of experiences. Therefore, on each edge server, the sample complexity, computational complexity, and stability analysis MA-DQL method are similar to that of standard DQL.

Algorithm 4 Knowledge Distillation

```

1: Initialize replay memory  $\mathcal{V}_n$  for each LSM  $n$ 
2: Initialize actor network at each LSM  $n$  with random weights  $\theta_n^\mu$ 
3: for  $t = 1, 2, \dots$  do
4:   for each LSM  $n \in \mathcal{N}$  do
5:     Get local observation  $o_{n,t}$  and run N2O cooperatively to obtain  $\tilde{a}_{n,t}$ 
6:     Store experience  $(o_{n,t}, \tilde{a}_{n,t})$  of time slot  $t$  in replay memory  $\mathcal{V}_n$ 
7:     if update actor network then
8:       Randomly sample a mini-batch  $\mathcal{H}$  from  $\mathcal{V}$ 
9:       Calculate loss  $L(\theta_n^\mu) = \sum_{(o_n, \tilde{a}_n) \in \mathcal{H}} (\tilde{a}_n - \mu_n(o_n, \theta_n^\mu))^2$ 
10:      Update actor parameter  $\theta_n^\mu = \theta_n^\mu + \delta \nabla_{\theta_n^\mu} L(\theta_n^\mu)$ 
11:    end if
12:  end for
13: end for

```

Sample complexity. The high sample complexity is a common issue in the deep learning community. The bound of sample complexity for DQL is still an open problem. Two very recent works [1], [2] have provided some empirical and theoretical results on the sample complexity of DQL. [1] shows via experiments that the sample complexity of DQL varies significantly based on the environment. [2] gives a bound of sample complexity for a simplified version of DQL with several assumptions on the reward function and MDP, and the bound is also related to the difficulty of the target problem. There also exist works that improve the sample efficiency of DQL, interested readers are recommended to refer to [3] and references therein. These learning techniques are also compatible with multi-agent DQL.

Computational complexity. The computational complexity of our method lies in training and running DNNs, and therefore we give the computational complexity of DNN inference (forward propagation) and DNN training (backward propagation). The computational complexity depends on the structure of DNNs. In this paper, we use the Multi-Layer Perceptron (MLP) network which consists of multiple fully connected layers.

Computational complexity of DNN Inference (Forward Propagation): For a standard MLP, the computational complexity is dominated by the matrix multiplication operations. For example, an MLP with n inputs, H hidden layers, where i -th hidden layer contains h_i hidden nodes, and k

output nodes will perform $nh_1 + h_H k + \sum_{i=1}^{H-1} h_i h_{i+1}$ multiplications. Therefore, the computational complexity of DNN inference is $O(nh_1 + h_H k + \sum_{i=1}^{H-1} h_i h_{i+1})$.

Computational complexity of DNN training (Backward Propagation): In general, using the automatic differentiation [4], the backward propagation of MLP is at most a constant factor slower than the forward propagation to the output. Therefore, the computational complexity of backward propagation is $O(\sum_{i=1}^{H-1} h_i h_{i+1})$. We use mini-batches for updating the DNN weights and let B be the batch size the computational complexity of DNN training is $O(B \sum_{i=1}^{H-1} h_i h_{i+1})$.

Stability: Based on the classic result in [5], DQL is able to achieve stability by utilizing the experience replay technique. Note that the learning process of our method on each edge server follows standard DQL, therefore the stability of our method is also guaranteed.

APPENDIX C

PROOF OF THEOREM 1

We begin by defining auxiliary variables and establishing lemmas useful in the proof. Using the techniques in [6], we first define two auxiliary sequences:

$$\bar{\mathbf{z}}(\tau) := \frac{1}{N} \sum_{n=1}^N \mathbf{z}_n(\tau) \quad \text{and} \quad \mathbf{y}(\tau) = \Pi_{\mathcal{A}}^{\psi}(\bar{\mathbf{z}}(\tau), \beta(\tau-1)) \quad (7)$$

The sequence $\bar{\mathbf{z}}(\tau)$ evolves as:

$$\begin{aligned} \bar{\mathbf{z}}(\tau+1) &= \frac{1}{N} \sum \bar{\mathbf{z}}_n(\tau+1) \\ &= \frac{1}{N} \sum_{n=1}^N \left(\sum_{m=1}^N W_{m,n} \mathbf{z}_m(\tau) + g_n(\tau) \right) \\ &\stackrel{\dagger}{=} \frac{1}{N} \sum_{m=1}^N \mathbf{z}_m(\tau) + \frac{1}{N} \sum_{n=1}^N g_n(\tau) \\ &= \bar{\mathbf{z}}(\tau) + \frac{1}{N} \sum_{n=1}^N g_n(\tau) \end{aligned}$$

where the equality $\stackrel{\dagger}{=}$ in above equation follows from double-stochasticity of matrix W . Next, we state a few useful results regarding the converge of the standard dual averaging. Let us begin with a result about Lipschitz continuity of the projection mapping $\Pi_{\mathcal{A}}^{\psi_n}$.

Lemma 1. For a LSM $n \in \mathcal{N}$ and an arbitrary pair $\mathbf{z}, \mathbf{z}' \in \mathbb{R}^N$, we have $\|\Pi_{\mathcal{A}}^{\psi_n}(\mathbf{z}, \beta) - \Pi_{\mathcal{A}}^{\psi_n}(\mathbf{z}', \beta)\| \leq \beta \|\mathbf{z} - \mathbf{z}'\|_*$, where $\|\cdot\|_*$ is dual norm to $\|\cdot\|$.

Lemma 1 is a standard result in convex analysis ([7], Lemma 1). We next give the convergence guarantee for the standard dual averaging.

Lemma 2. Consider an arbitrary sequence of vectors $\{g(\tau)\}_{\tau=1}^{\infty}$ and the sequence given by $\mathbf{a}(\tau+1) = \Pi_{\mathcal{A}}^{\psi}(\sum_{l=1}^{\tau} g(l), \beta(\tau)) := \arg \min_{\mathbf{a} \in \mathcal{A}} \left\{ \sum_{l=1}^{\tau} \langle g(l), \mathbf{a} \rangle + \frac{1}{\beta(\tau)} \psi(\mathbf{a}) \right\}$. Then, for any non-increasing sequence $\{\beta(\tau)\}_{\tau=0}^{\infty}$ of positive stepsizes and for any $\mathbf{a}^* \in \mathcal{A}$, we have:

$$\sum_{\tau=1}^T \langle g(\tau), \mathbf{a}(\tau) - \mathbf{a}^* \rangle \leq \frac{1}{2} \sum_{\tau=1}^T \beta(\tau-1) \|g(\tau)\|_*^2 + \frac{1}{\beta(T)} \psi(\mathbf{a}^*)$$

Proof. This lemma is a consequence of Theorem 2 and Equation (3.3) in [7]. \square

With the above definitions and lemma, we now give the proof for Theorem 1. Since the local observations o_n and Q-network parameters θ_n does not change during execution of N₂O, we let $Q_n(\mathbf{a})$ denote $Q_n(o_n, \alpha_n, \theta_n)$ for ease of exposition with $\mathbf{a} = \{a_1, a_2, \dots, a_N\}$ and $\alpha_n = \{a_n \cup \{a_i\}_{i \in \mathcal{B}_n}\}$. Our proof is based on analyzing the sequence $\mathbf{y}(\tau)_{\tau=1}^{\infty}$. Given an arbitrary $\mathbf{a}^* \in \mathcal{A}$, we have

$$\begin{aligned} & \frac{1}{N} \sum_{\tau=1}^T \sum_{n=1}^N Q_n(\mathbf{a}^*) - Q_n(\mathbf{y}(\tau)) \\ &= \frac{1}{N} \sum_{\tau=1}^T \sum_{n=1}^N (Q_n(\mathbf{a}^*) - Q_n(\mathbf{a}_n(\tau))) + \frac{1}{N} \sum_{\tau=1}^T \sum_{n=1}^N (Q_n(\mathbf{a}(\tau)) - Q_n(\mathbf{y}(\tau))) \\ &\stackrel{\dagger}{\leq} \frac{1}{N} \sum_{\tau=1}^T \sum_{n=1}^N (Q_n(\mathbf{a}^*) - Q_n(\mathbf{a}_n(\tau)) + L \|\mathbf{a}_n(\tau) - \mathbf{y}(\tau)\|) \end{aligned}$$

The inequality $\stackrel{\dagger}{\leq}$ in the above following by the L -Lipschitz condition of $Q_n(\cdot)$. Now let $g_n(\tau) = -\partial Q_n(\mathbf{a}_n(\tau)) / \partial \mathbf{a}_n(\tau)$ be the negative gradient of $Q_n(\cdot)$ at $\mathbf{a}_n(\tau)$. Using the convexity of $-Q_n(\mathbf{a}_n) + \frac{1}{\beta(\tau)} \psi_n(\mathbf{a}_n)$, we have the following inequality:

$$\sum_{n=1}^N \left(-Q_n(\mathbf{a}_n(\tau)) + \frac{\psi_n(\mathbf{a}_n(\tau))}{\beta(\tau)} - \left(-Q_n(\mathbf{a}^*) + \frac{\psi_n(\mathbf{a}^*)}{\beta(\tau)} \right) \right) \leq \sum_{n=1}^N \left\langle g_n(\tau) + \frac{\partial \psi_n(\mathbf{a}_n(\tau))}{\beta(\tau)}, \mathbf{a}_n(\tau) - \mathbf{a}^* \right\rangle$$

Rearranging the above inequality yields:

$$\sum_{n=1}^N (Q_n(\mathbf{a}^*) - Q_n(\mathbf{a}_n(\tau))) \leq \sum_{n=1}^N \left(\langle g_n(\tau), \mathbf{a}_n(\tau) - \mathbf{a}^* \rangle + \left\langle \frac{\partial \psi_n(\mathbf{a}_n(\tau))}{\beta(\tau)}, \mathbf{a}_n(\tau) - \mathbf{a}^* \right\rangle + \frac{\psi_n(\mathbf{a}^*)}{\beta(\tau)} \right) \quad (8)$$

We next bound the terms on the right-hand side of (8) separately. Notice that the first term can be decompose into two parts:

$$\sum_{n=1}^N \langle g_n(\tau), \mathbf{a}_n(\tau) - \mathbf{a}^* \rangle = \sum_{n=1}^N \langle g_n(\tau), \mathbf{y}(\tau) - \mathbf{a}^* \rangle + \sum_{n=1}^N \langle g_n(\tau), \mathbf{a}_n(\tau) - \mathbf{y}(\tau) \rangle \quad (9)$$

Recalling the definition of $\bar{\mathbf{z}}(\tau)$ and $\mathbf{y}(\tau)$ in (7), we can write the first term in the decomposition (9) in the similar way as the bound in Lemma 2:

$$\begin{aligned} \frac{1}{N} \sum_{\tau=1}^T \left\langle \sum_{n=1}^N g_n(\tau), \mathbf{y}(\tau) - \mathbf{a}^* \right\rangle &= \left\langle \sum_{\tau=1}^T \left(\frac{1}{N} \sum_{n=1}^N g_n(\tau) \right), \mathbf{y}(\tau) - \mathbf{a}^* \right\rangle \\ &\leq \frac{1}{2} \sum_{\tau=1}^T \beta(\tau-1) \left\| \frac{1}{N} \sum_{n=1}^N g_n(\tau) \right\|_*^2 + \frac{1}{\beta(T)} \psi(\mathbf{a}^*) \\ &\leq \frac{L^2}{2} \sum_{\tau=1}^T \beta(\tau-1) + \frac{1}{\beta(T)} \psi(\mathbf{a}^*) \end{aligned}$$

For the second term in decomposition (9), we have:

$$\frac{1}{N} \sum_{\tau=1}^T \sum_{n=1}^N \langle g_n(\tau), \mathbf{a}_n(\tau) - \mathbf{y}(\tau) \rangle \leq \frac{L}{N} \sum_{\tau=1}^T \sum_{n=1}^N \|\mathbf{a}_n(\tau) - \mathbf{y}(\tau)\|.$$

The inequality follows from $\|g_n(\tau)\|_* \leq L$. Combining the above results we have:

$$\begin{aligned} \frac{1}{N} \sum_{\tau=1}^T \sum_{n=1}^N Q_n(\mathbf{a}^*) - Q_n(\mathbf{y}(\tau)) &\leq \frac{L^2}{2} \sum_{\tau=1}^T \beta(\tau-1) + \frac{1}{\beta(T)} \psi(\mathbf{a}^*) + \frac{2L}{N} \sum_{\tau=1}^T \sum_{n=1}^N \|\mathbf{a}_n(\tau) - \mathbf{y}(\tau)\| \\ &\quad + \frac{1}{N} \sum_{\tau=1}^T \sum_{n=1}^N \left(\left\langle \frac{\partial \psi_n(\mathbf{a}(\tau))}{\beta(\tau)}, \mathbf{a}_n(\tau) - \mathbf{a}^* \right\rangle + \frac{\psi_n(\mathbf{a}^*)}{\beta(\tau)} \right) \end{aligned} \quad (10)$$

For an arbitrary action $\mathbf{a}_i(\tau)$, $i \in \mathcal{N}$, considering the L -Lipschitz continuity of $Q_n(\cdot)$, we have

$$\begin{aligned} \frac{1}{TN} \sum_{\tau=1}^T \sum_{n=1}^N (Q_n(\mathbf{a}^*) - Q_n(\mathbf{a}_i(\tau))) &= \frac{1}{TN} \sum_{\tau=1}^T \sum_{n=1}^N (Q_n(\mathbf{a}^*) - Q_n(\mathbf{y}(\tau)) + Q_n(\mathbf{y}(\tau)) - Q_n(\mathbf{a}_i(\tau))) \\ &\leq \frac{1}{TN} \sum_{\tau=1}^T \sum_{n=1}^N (Q_n(\mathbf{a}^*) - Q_n(\mathbf{y}(\tau))) + \frac{L}{T} \sum_{\tau=1}^T \|\mathbf{a}_i(\tau) - \mathbf{y}(\tau)\| \end{aligned}$$

By utilizing again the convexity of $-Q_n(\mathbf{a}_n) + \frac{1}{\beta(\tau)} \psi(\mathbf{a}_n)$, we have

$$\begin{aligned} \frac{1}{T} \sum_{\tau=1}^T \left(-Q_n(\mathbf{a}_i(\tau)) + \frac{1}{\beta(\tau)} \psi_n(\mathbf{a}_i(\tau)) \right) &\geq \frac{1}{T} \sum_{\tau=1}^T \left(-Q_n(\mathbf{a}_i(\tau)) + \frac{1}{\beta(0)} \psi_n(\mathbf{a}_i(\tau)) \right) \\ &\geq -Q_n \left(\frac{1}{T} \sum_{\tau=1}^T \mathbf{a}_i(\tau) \right) + \frac{1}{\beta(0)} \psi_n \left(\frac{1}{T} \sum_{\tau=1}^T \mathbf{a}_i(\tau) \right) \\ &= -Q_n(\bar{\mathbf{a}}_i(T)) + \frac{1}{\beta(0)} \psi_n(\bar{\mathbf{a}}_i(T)) \end{aligned}$$

which implies that

$$Q_n(\bar{\mathbf{a}}_i(T)) \geq \frac{1}{T} \sum_{\tau=1}^T \left(Q_n(\mathbf{a}_i(\tau)) - \frac{1}{\beta(\tau)} \psi_n(\mathbf{a}_i(\tau)) \right) + \frac{1}{\beta(0)} \psi_n(\bar{\mathbf{a}}_i(T)) \quad (11)$$

Therefore, we have

$$\begin{aligned} & \frac{1}{N} \sum_{n=1}^N Q_n(\mathbf{a}^*) - Q_n(\bar{\mathbf{a}}_i(T)) \\ & \leq \frac{1}{TN} \sum_{\tau=1}^T \sum_{n=1}^N (Q_n(\mathbf{a}^*) - Q_n(\mathbf{a}_i(\tau))) + \frac{1}{TN} \sum_{\tau=1}^T \sum_{n=1}^N \frac{1}{\beta(\tau)} \psi_n(\mathbf{a}_i(\tau)) - \frac{1}{N} \sum_{n=1}^N \frac{1}{\beta(0)} \psi_n(\bar{\mathbf{a}}_i(T)) \\ & \leq \frac{1}{TN} \sum_{\tau=1}^T \sum_{n=1}^N (Q_n(\mathbf{a}^*) - Q_n(\mathbf{y}(\tau)) + L \|\mathbf{a}_i(\tau) - \mathbf{y}(\tau)\|) + \frac{1}{TN} \sum_{\tau=1}^T \sum_{n=1}^N \frac{1}{\beta(\tau)} \psi_n(\mathbf{a}_i(\tau)) \\ & \leq \frac{L^2}{2T} \sum_{\tau=1}^T \beta(\tau - 1) + \frac{1}{\beta(T)} \psi(\mathbf{a}^*) + \frac{2L}{TN} \sum_{\tau=1}^T \sum_{n=1}^N \|\mathbf{a}_n(\tau) - \mathbf{y}(\tau)\| \\ & \quad + \frac{1}{TN} \sum_{\tau=1}^T \sum_{n=1}^N \left(\left\langle \frac{\partial \psi_n(\mathbf{a}(\tau))}{\beta(\tau)}, \mathbf{a}_n(\tau) - \mathbf{a}^* \right\rangle + \frac{\psi_n(\mathbf{a}^*)}{\beta(\tau)} \right) + \frac{L}{T} \sum_{\tau=1}^T \|\mathbf{a}_i(\tau) - \mathbf{y}(\tau)\| \\ & \quad + \frac{1}{TN} \sum_{\tau=1}^T \sum_{n=1}^N \frac{1}{\beta(\tau)} \psi_n(\mathbf{a}_i(\tau)) \end{aligned}$$

Using Lemma 1, $\psi_n(\mathbf{a}) \leq \psi_n^{\max}$ and $\nabla \psi_n(\mathbf{a}) \leq \psi'^{\max}$, we can easily reaching

$$\begin{aligned} & \frac{1}{N} \sum_{n=1}^N (Q_n(\mathbf{a}^*) - Q_n(\bar{\mathbf{a}}_i(T))) \\ & \leq \frac{L^2}{2T} \sum_{\tau=1}^T \beta(\tau - 1) + \frac{1}{\beta(T)} \psi(\mathbf{a}^*) + \frac{2L}{TN} \sum_{\tau=1}^T \sum_{n=1}^N \beta(\tau) \|\bar{\mathbf{z}}(\tau) - \mathbf{z}_n(\tau)\|_* \\ & \quad + \frac{L}{T} \sum_{\tau=1}^T \beta(\tau) \|\bar{\mathbf{z}}(\tau) - \mathbf{z}_i(\tau)\|_* + \frac{2}{N\beta(T)} \sum_{n=1}^N \psi_n^{\max} + \frac{d^{\max}}{N\beta(T)} \sum_{n=1}^N \psi'_n{}^{\max} \end{aligned}$$

where $d^{\max} = \arg \max_{\mathbf{a}, \mathbf{a}'} \|\mathbf{a} - \mathbf{a}'\|, \forall \mathbf{a}, \mathbf{a}' \in \mathcal{A}$.

APPENDIX D

PROOF OF THEOREM 2

We first introduce the following notational conventions. For an $N \times N$ matrix W , we define its singular values $\sigma_1(W) \geq \sigma_2(W) \geq \dots \geq \sigma_N(W) \geq 0$. For a real symmetric matrix, we use $\lambda_1(W) \geq \lambda_2(W) \geq \dots \geq \lambda_N(W)$ to denote N real eigenvalues of W . Let $\Delta_N = \{x \in \mathbb{R}^N | x \geq 0, \sum_{n=1}^N x_n = 1\}$ denote the N -dimensional probability simplex, and $\mathbb{1}$ denote the vector of all ones. Given these definitions, we introduce the below lemma.

Lemma 3. For a stochastic matrix W and $x \in \Delta_N$, the following inequality holds true for any positive integer τ .

$$\|W^\tau x - \mathbb{1}/N\|_1 \leq \sqrt{N} \|W^\tau x - \mathbb{1}/N\|_2 \leq \sigma_2(W)^\tau \sqrt{N}.$$

Proof. The proof can be found in [8] regarding the Perron-Frobenius theory. \square

The key focus is controlling the term $\sum_{n=1}^N \beta(\tau) \|\bar{z}(\tau) - z_n(\tau)\|_*$. Define the matrix $\Phi(\tau, \kappa) = W^{\tau-\kappa+1}$. Let $[\Phi(\tau, \kappa)]_{mn}$ be the m -th entry of the n -th column of $\Phi(\tau, \kappa)$. Then we have:

$$z_n(\tau+1) = \sum_{m=1}^N [\Phi(\tau, \kappa)]_{mn} z_m(\kappa) + \sum_{v=\kappa+1}^{\tau} \left(\sum_{m=1}^N [\Phi(\tau, v)]_{mn} g_m(v-1) \right) + g_n(\tau)$$

The above reduces to the standard update in (3) when $\kappa = \tau$. Recall that $\bar{z}(\tau+1) = \bar{z}(\tau) + \frac{1}{N} \sum_{n=1}^N g_n(\tau)$, we will have

$$\bar{z}(\tau) - z_n(\tau) = \sum_{\kappa=1}^{\tau-1} \sum_{m=1}^N \left(\frac{1}{N} - [\Phi(\tau-1, \kappa)]_{mn} \right) g_m(\kappa-1) + \frac{1}{N} \sum_{m=1}^N (g_m(\tau-1) - g_n(\tau-1))$$

Recall $\|g_n(\tau)\|_* \leq L, \forall n, \tau$. With the definition $\bar{\Phi}(\tau, \kappa) := \mathbb{1}\mathbb{1}^\top/N - \Phi(\tau, \kappa)$, we can reach

$$\|\bar{z}(\tau) - z_n(\tau)\|_* \leq \left\| \sum_{\kappa=1}^{\tau-1} \sum_{m=1}^N [\bar{\Phi}(\tau-1, \kappa)]_{mn} g_m(\kappa-1) \right\|_* + \left\| \frac{1}{N} \sum_{m=1}^N (g_m(\tau-1) - g_n(\tau-1)) \right\|_*.$$

Letting e_n be the n -th standard basis vector, the above is further bounded by

$$\begin{aligned} & \sum_{\kappa=1}^{\tau-1} \sum_{m=1}^N |[\bar{\Phi}(\tau-1, \kappa)]_{mn}| \|g_m(\kappa-1)\|_* + \frac{1}{N} \sum_{m=1}^N \|g_m(\tau-1) - g_n(\tau-1)\|_* \\ & \leq \sum_{\kappa=1}^{\tau-1} L \|\Phi(\tau-1, \kappa) e_n - \mathbb{1}/N\|_1 + 2L \end{aligned}$$

We now break the above sum into two parts separated by a cut off point $\hat{\tau}$:

$$\begin{aligned} & \|\bar{z}(\tau) - z_n(\tau)\|_* \\ & \leq L \sum_{\kappa=\tau-\hat{\tau}}^{\tau-1} \|\Phi(\tau-1, \kappa) e_n - \mathbb{1}/N\|_1 + L \sum_{\kappa=1}^{\tau-1-\hat{\tau}} \|\Phi(\tau-1, \kappa) e_n - \mathbb{1}/N\|_1 + 2L \end{aligned} \quad (12)$$

Note that the indexing on $\Phi(\tau-1, \kappa) = W^{\tau-\kappa+1}$ implies that when κ is small, $\Phi(\tau-1, \kappa)$ is close to uniform. Given $\|\Phi(\tau, \kappa) e_n - \mathbb{1}/N\|_1 \leq \sqrt{N} \sigma_2(W)^{t-s+1}$ in Lemma 3, if we let $\tau - \kappa \geq \frac{\log \epsilon^{-1}}{\log \sigma_2(W)^{-1}} - 1$ then $\|\Phi(\tau, \kappa) e_n - \mathbb{1}/N\|_1 \leq \sqrt{N} \epsilon$. By setting $\epsilon^{-1} = T\sqrt{N}$, for $\tau - \kappa + 1 \geq \log(T\sqrt{N})/\log \sigma_2(W)^{-1}$, we have $\|\Phi(\tau, \kappa) e_n - \mathbb{1}/N\|_1 \leq \frac{1}{T}$. For $\kappa \geq t - \log(T\sqrt{N})/\log \sigma_2(W)^{-1}$,

we simply have $\|\Phi(\tau, \kappa)e_n - \mathbb{1}/N\|_1 \leq 2$. Therefore, if we set $\hat{\tau} = \log(T\sqrt{N})/\log \sigma_2(W)^{-1}$, we will have:

$$\begin{aligned} \|\bar{z}(\tau) - z_n(\tau)\|_* &\leq 2L(\tau - 1 - (\tau - \hat{\tau})) + \frac{L}{T}(\tau - 1 - \hat{\tau} - 1) + 2L \\ &\leq 2L\hat{\tau} + \frac{L\tau}{T} + 2L \\ &\leq 2L \frac{\log(T\sqrt{N})}{\log \sigma_2(W)^{-1}} + 3L \end{aligned}$$

Using the convexity of $\log(\cdot)$, we have $\sigma_2(W)^{-1} \geq 1 - \sigma_2(W)$, which implies $\|\bar{z}(\tau) - z_n(\tau)\|_* \leq 2L \frac{\log(T\sqrt{N})}{1 - \sigma_2(W)} + 3L$. Using $\sum_{\tau}^T \tau^{-1/2} \leq 2\sqrt{T} - 1$ and results in Theorem 1 complete the proof.

APPENDIX E

PROOF OF COROLLARY 1

In order to prove the statement in corollary, we first use graph Laplacian [9] to describe the graph structure. We let $A \in \mathbb{R}^{N \times N}$ be the adjacency matrix of the undirected graph G , satisfying $A_{i,j} = 1$ when $(i, j) \in \mathcal{E}$ and $A_{i,j} = 0$ otherwise. For each node $i \in \mathcal{N}$, we let $\delta_i = |\mathcal{B}_i| = \sum_{j=1}^N A_{ij}$ denote the degree of node i , and we define the diagonal matrix $D = \text{diag}\{\delta_1, \dots, \delta_N\}$. We assume that the graph is connected such that $\delta_i \geq 1$ for all $i \in \mathcal{N}$ and D is invertible. With this notation, the *normalized graph Laplacian* of graph G is

$$\mathcal{L}(G) = I - D^{-1/2}AD^{-1/2}.$$

The graph Laplacian $\mathcal{L} := \mathcal{L}(G)$ is symmetric, positive semi-definite, and satisfies $\mathcal{L}D^{1/2}\mathbb{1} = 0$, where $\mathbb{1}$ is the all ones vector. When the graph is degree-regular, i.e., $\delta_i = \delta, \forall i \in \mathcal{N}$, the standard random walk with self-loops on G given by the matrix $W := I - (\delta/(\delta+1))\mathcal{L}$ is doubly stochastic and valid for our theory. For non-regular graphs, a minor modification is required to obtain a double stochastic matrix: let $\delta_{\max} = \max_{i \in \mathcal{N}} \delta_i$ denote G 's largest degree and define

$$W_N(G) = I - \frac{1}{\delta_{\max} + 1}(D - A) = I - \frac{1}{\delta_{\max} + 1}D^{1/2}\mathcal{L}D^{1/2} \quad (13)$$

This matrix is symmetric by construction and it is also doubly stochastic. Note that if the graph is δ -regular, the $W_N(G)$ is the standard choice mentioned above. Plugging $W_N(G)$ into Theorem 2, we have the convergence rate of N₂O becomes

$$O\left(\frac{L^2}{\sqrt{T}} \frac{\log(T\sqrt{N})}{1 - \sigma_2(W_N(G))}\right).$$

The corollary is based on bounding the spectral gap of $W_N(G)$. We begin with a technical lemma.

Lemma 4. Let $\bar{\delta} = \delta_{\max}$, the matrix $W_N(G)$ satisfies

$$\sigma_2(W_N(G)) \leq \max \left\{ 1 - \frac{\min_i \delta_i}{\bar{\delta} + 1} \lambda_{N-1}(\mathcal{L}), \frac{\bar{\delta}}{\bar{\delta} + 1} \lambda_1(\mathcal{L}) - 1 \right\}$$

where $\lambda_{N-1}(\mathcal{L})$ and $\lambda_1(\mathcal{L})$ is the second smallest eigenvalue and the largest eigenvalue of \mathcal{L} , respectively.

Proof. By a theorem of Ostrowski on congruent matrices (Theorem 4.5.9 in [10]), we have

$$\lambda_k(D^{1/2} \mathcal{L} D^{1/2}) \in \left[\min_i \delta_i \lambda_k(\mathcal{L}), \max_i \delta_i \lambda_k(\mathcal{L}) \right]. \quad (14)$$

Since $\mathcal{L} D^{1/2} \mathbb{1} = 0$, we have $\lambda_N(\mathcal{L}) = 0$ and so it suffice to focus on $\lambda_1(D^{1/2} \mathcal{L} D^{1/2})$ and $\lambda_{N-1}(D^{1/2} \mathcal{L} D^{1/2})$. From the definition of $W_N(G)$ in (13), the eigenvalues pf $W_N(G)$ are of the form $1 - (\delta_{\max} + 1)^{-1} \lambda_k(D^{1/2} \mathcal{L} D^{1/2})$. The bound (14) and the fact that all eigenvalues of \mathcal{L} are non-negative implies that $\sigma_2(W_N(G)) = \max_{k < N} \{|1 - (\delta_{\max} + 1)^{-1} \lambda_k(D^{1/2} \mathcal{L} D^{1/2})|\}$ is upper bounded by the larger of $1 - (\delta_{\min}/(\delta_{\max} + 1)) \lambda_{N-1}(\mathcal{L})$ and $(\delta_{\max}/(\delta_{\max} + 1)) \lambda_1(\mathcal{L}) - 1$.

Computing the upper bound in Lemma 4 requires controlling both $\lambda_{N-1}(\mathcal{L})$ and $\lambda_1(\mathcal{L})$. To circumvent this complication, we use the well-known idea of a lazy random walk [11], in which we replace $W_N(G)$ by $\frac{1}{2}(I + W_N(G))$. The resulting symmetric matrix has the same eigenstructure as $W_N(G)$. Further, $\frac{1}{2}(I + W_N(G))$ is positive semidefinite such that $\sigma_2\left(\frac{1}{2}(I + W_N(G))\right) = \lambda_2\left(\frac{1}{2}(I + W_N(G))\right)$, and hence

$$\begin{aligned} \sigma_2\left(\frac{1}{2}(I + W_N(G))\right) &= \lambda_2\left(I - \frac{1}{2(\delta_{\max} + 1)} D^{1/2} \mathcal{L} D^{1/2}\right) \\ &\leq 1 - \frac{\delta_{\min}}{2(\delta_{\max} + 1)} \lambda_{N-1}(\mathcal{L}). \end{aligned}$$

Consequently, it is sufficient to bound only $\lambda_{N-1}(\mathcal{L})$. The convergence rate implied by the lazy random walk through Theorem D is no worse than twice that of the original walk, which is insignificant for the analysis. We are now equipped to address each of the graph classes covered by Corollary 1.

Regular Grids: Consider a \sqrt{N} -by- \sqrt{N} grid, in particular, a regular k -connected grid in which any node is joined to every node that is fewer than k horizontal or vertical edges away in an axis-aligned direction. In this case, we use results on Cartesian product of graphs [9] to analyze the eigenstructure of the Laplacian. In particular, the \sqrt{N} -by- \sqrt{N} k -connected grid is the Cartesian

product of two regular k -connected paths of \sqrt{N} nodes. The second smallest eigenvalue of a Cartesian product of graphs is half the minimum of second-smallest eigenvalues of the original graphs [9]. Thus, if $k \leq N^{1/4}$, then we have $\lambda_{N-1}(\mathcal{L}) = \Theta(k^2/N)$, and use Lemma 4, it is easy to see

$$1 - \sigma_2(W) = \Theta(k^2/N).$$

The result (a) in Corollary 1 immediately follows.

Random Geometric Graphs: Using the proof of Lemma 10 in [12], we see that for any ϵ and $c > 0$, if $r = \sqrt{\log^{1+\epsilon} N / (N\pi)}$, then with probability at least $1 - 2/N^{c-1}$

$$\log^{1+\epsilon} N - \sqrt{2}c \log N \leq \delta_i \leq \log^{1+\epsilon} N + \sqrt{2}c \log N \quad (15)$$

for all i . Recent work [13] gives concentration results on the second-smallest eigenvalue of a geometric graph. Theorem 3 in [13] indicates that if $r = \omega\left(\sqrt{\log N/N}\right)$, then with high probability $\lambda_{N-1}(\mathcal{L}) = \Omega(r^2) = \omega\left(\sqrt{\log N/N}\right)$. Using (15), we have for $r = (\log^{1+\epsilon} N/N)^{1/2}$, the ratio $\min_i \delta_i = \Theta(1)$ and $\lambda_{N-1}(\mathcal{L}) = \Omega(\log^{1+\epsilon} N/N)$ with high probability. Therefore, we have

$$1 - \sigma_2(W) = \Omega\left(\frac{\log^{1+\epsilon} N}{N}\right),$$

which gives the result (b) in Corollary 1. \square

APPENDIX F

PROOF OF THEOREM 3

Recall the Theorem 1 involves the sum $\frac{2L}{TN} \sum_{\tau=1}^T \sum_{n=1}^N \beta \tau \|\bar{\mathbf{z}}(\tau) - \mathbf{z}_n(\tau)\|_*$. In the proof of Theorem 2 (Appendix D), we have shown how to control this sum when the communication between agents occurs on a static underlying network structure via a fixed doubly-stochastic matrix W . We now extend the analysis to time-varying $W(\tau)$.

Given $W(\tau)$ at iteration τ , the update policy in (3) becomes:

$$\mathbf{z}_n(\tau+1) = \sum_{m=1}^N W_{m,n}(\tau) \mathbf{z}_m \tau + g_n(\tau), \quad \mathbf{a}_n(\tau+1) = \Pi_{\mathcal{A}}^{\Psi_n}(\mathbf{z}_n(\tau+1), \beta(\tau))$$

We still have the evolution $\bar{\mathbf{z}}(\tau+1) = \bar{\mathbf{z}}(\tau) + \frac{1}{N} \sum_{n=1}^N g_n(\tau)$. Define $\Phi(\tau, \kappa) = W(\kappa)W(\kappa+1) \dots W(\tau)$ with $\kappa \leq \tau$, the following holds

$$\bar{\mathbf{z}}(\tau) - \mathbf{z}_n(\tau) = \sum_{\kappa=1}^{\tau-1} \sum_{m=1}^N \left(\frac{1}{N} - [\Phi(\tau-1, \kappa)]_{mn} \right) g_m(\kappa-1) + \frac{1}{N} \sum_{m=1}^N (g_m(\tau-1) - g_n(\tau-1)).$$

To show the convergence for the random communication model, we must control the convergence of $\Phi(\tau - 1, \kappa)$ to the uniform distribution. We first claim that

$$\Pr\{\|\Phi(\tau, \kappa)\mathbf{e}_n - \mathbb{1}/N\|_2 \geq \epsilon\} \leq \epsilon^{-2} \lambda_2 (\mathbb{E}[W(\tau)^\top W(\tau)])^{\tau-\kappa+1}. \quad (16)$$

This inequality can be established by modifying a few known result in [12]. Let Δ_N denote the N -dimensional probability simplex and $u(0) \in \Delta_N$ be arbitrary. Consider the random sequence $\{u(\tau)\}_{\tau=1}^\infty$ generated by $u(\tau + 1) = W(\tau)u\tau$. Let $v(\tau) := u(\tau) - \mathbb{1}/N$ correspond to the portion of $u(\tau)$ orthogonal to the all one vector. Calculating the second moment of $v(\tau + 1)$:

$$\begin{aligned} \mathbb{E}[\langle v(\tau + 1), v(\tau + 1) \rangle | v(\tau)] &= \mathbb{E}[v(\tau)W(\tau)^\top W(\tau)v(\tau) | v(\tau)] \\ &= v(\tau)^\top \mathbb{E}[W(\tau)^\top W(\tau)] v(\tau) \\ &\leq \|v(\tau)\|_2^2 \lambda_2 (\mathbb{E}[W(\tau)^\top W(\tau)]) \end{aligned}$$

since $\langle v(\tau), \mathbb{1} \rangle = 0$, $v(t)$ is orthogonal to the first eigenvector of $W(\tau)$, and $W(\tau)^\top W(\tau)$ is symmetric and double stochastic. Applying Chebyshev's inequality yields:

$$\Pr\left[\frac{\|u(\tau) - \mathbb{1}/N\|_2}{\|u(0)\|_2} \geq \epsilon\right] \leq \frac{\mathbb{E}[\|v(\tau)\|_2^2]}{\|u(0)\|_2^2 \epsilon^2} \leq \epsilon^{-2} \frac{\|v(0)\|_2^2 \lambda_2 (\mathbb{E}[W(\tau)^\top W(\tau)])^\tau}{\|u(0)\|_2^2}$$

Replacing $u(0)$ with \mathbf{e}_n and noticing that $\|\mathbf{e}_n - \mathbb{1}/N\|_2 \leq 1$ yields the result (16).

We now use the result in (16) to prove Theorem 3. Following similar technique used in the proof of Theorem 2. We begin by choosing a iteration index $\hat{\tau}$ such that for $\tau - \kappa \geq \hat{\tau}$, with high probability, $\Phi(\tau, \kappa)$, is close to the uniform matrix $\mathbb{1}\mathbb{1}^\top/N$. We then break the summation from 1 to T into two terms separated by the cutoff point $\hat{\tau}$. Throughout this derivation, we let λ_2 denote $\lambda_2(\mathbb{E}[W(\tau)^\top W(\tau)])$ to ease notation. Using the probabilistic bound in (16), if $\tau - \kappa \geq (3 \log \epsilon^{-1} / \log \lambda_2^{-1}) - 1$, then $\Pr\{\|\Phi(\tau, \kappa)\mathbf{e}_n - \mathbb{1}/N\|_2 > \epsilon\} \leq \epsilon$. Consequently, the choice

$$\tilde{\tau} = \frac{3 \log(T^2 N)}{\log \lambda_2^{-1}} = \frac{6 \log T + 3 \log N}{\log \lambda_2^{-1}} \leq \frac{6 \log T + 3 \log N}{1 - \lambda_2} \quad (17)$$

guarantees that if $\tau - \kappa \geq \tilde{\tau} - 1$, then

$$\Pr\left[\|\Phi(\tau, \kappa)\mathbf{e}_n - \mathbb{1}/N\|_2 \geq \frac{1}{T^2 N}\right] \leq (T^2 N)^2 \lambda_2^{\frac{3 \log(T^2 N)}{-\log \lambda_2}} = \frac{1}{T^2 N}. \quad (18)$$

Recalling the bound (12) in the proof of Theorem 2:

$$\|\bar{\mathbf{z}}(\tau) - \mathbf{z}_n(\tau)\|_* \leq L \sum_{\kappa=1}^{\tau-1} \|\Phi(\tau - 1, \kappa)\mathbf{e}_n - \mathbb{1}/N\|_1 + 2L$$

Breaking the above sum into two parts at $\hat{\tau}$ and using $\|\Phi(\tau, \kappa) - \mathbb{1}/N\|_1 \leq 2$ for $\kappa \geq \tau - \hat{\tau}$, we have

$$\begin{aligned} \|\bar{z}(\tau) - z_n(\tau)\|_* &\leq L \sum_{\kappa=\tau-\hat{\tau}}^{\tau-1} \|\Phi(\tau-1, \kappa) \mathbf{e}_n - \mathbb{1}/N\|_1 + L \sum_{\kappa=1}^{\tau-\hat{\tau}-1} \|\Phi(\tau-1, \kappa) \mathbf{e}_n - \mathbb{1}/N\|_1 + 2L \\ &\leq 2L \frac{3 \log(T^2 N)}{1 - \lambda_2} + L\sqrt{N} \sum_{\kappa=1}^{\tau-\hat{\tau}-1} \|\Phi(\tau-1, \kappa) \mathbf{e}_n - \mathbb{1}/N\|_2 + 2L \end{aligned}$$

Now for any $\kappa' \leq \kappa$, since the matrices $W(\tau)$ are doubly stochastic, we have

$$\begin{aligned} \|\Phi(\tau-1, \kappa') \mathbf{e}_n - \mathbb{1}/N\|_2 &= \|\Phi(\kappa-1, \kappa') \Phi(\tau-1, \kappa) \mathbf{e}_n - \mathbb{1}/N\|_2 \\ &\leq \|\Phi(\kappa-1, \kappa')\|_2 \|\Phi(\tau-1, \kappa) \mathbf{e}_n - \mathbb{1}/N\|_2 \\ &\leq \|\Phi(\tau-1, \kappa) \mathbf{e}_n - \mathbb{1}/N\|_2. \end{aligned}$$

where the final inequality uses the bound $\|\Phi(\kappa-1, \kappa')\|_2 \leq 1$. Using the result in (18), we have $\|\Phi(\tau-1, \tau-\hat{\tau}-1) \mathbf{e}_n - \mathbb{1}/N\|_2 \leq 1/(T^2 N)$ with probability at least $1 - 1/(T^2 N)$. Since κ ranges between 1 and $\tau - \hat{\tau}$, we have:

$$L\sqrt{N} \sum_{\kappa=1}^{\tau-\hat{\tau}-1} \|\Phi(\tau-1, \kappa) \mathbf{e}_n - \mathbb{1}/N\|_2 \leq L\sqrt{N} T \frac{1}{T^2 N} = \frac{L}{T\sqrt{N}}$$

Hence we have

$$\|\bar{z}(\tau) - z_n(\tau)\|_* \leq \frac{6L \log(T^2 N)}{1 - \lambda_2} + \frac{L}{T\sqrt{N}} + 2L$$

with probability at least $1 - 1/(T^2 N)$. Applying the union bound over all iterations $\tau = 1, \dots, T$ and nodes $n = 1, \dots, N$.

$$\Pr \left[\max_{\tau, n} \|\bar{z}(\tau) - z_n(\tau)\|_* > \frac{6L \log(T^2 N)}{1 - \lambda_2} + \frac{L}{T\sqrt{N}} + 2L \right] \leq \frac{1}{T}.$$

Recalling the master result in Theorem 1 completes the proof.

REFERENCES

- [1] J. Tyo and Z. Lipton, “How transferable are the representations learned by deep q agents?” *arXiv preprint arXiv:2002.10021*, 2020.
- [2] Z. Yang, Y. Xie, and Z. Wang, “A theoretical analysis of deep q-learning,” in *Learning for Dynamics and Control*. PMLR, 2020, pp. 486–489.
- [3] S. Y. Lee, C. Sungik, and S.-Y. Chung, “Sample-efficient deep reinforcement learning via episodic backward update,” in *Advances in Neural Information Processing Systems*, 2019, pp. 2112–2121.
- [4] A. G. Baydin, B. A. Pearlmutter, A. A. Radul, and J. M. Siskind, “Automatic differentiation in machine learning: a survey,” *The Journal of Machine Learning Research*, vol. 18, no. 1, pp. 5595–5637, 2017.

- [5] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski *et al.*, “Human-level control through deep reinforcement learning,” *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [6] A. Nedic and A. Ozdaglar, “Distributed subgradient methods for multi-agent optimization,” *IEEE Transactions on Automatic Control*, vol. 54, no. 1, pp. 48–61, 2009.
- [7] Y. Nesterov, “Primal-dual subgradient methods for convex problems,” *Mathematical programming*, vol. 120, no. 1, pp. 221–259, 2009.
- [8] R. A. Horn and C. R. Johnson, *Matrix analysis*. Cambridge university press, 2012.
- [9] F. R. Chung and F. C. Graham, *Spectral graph theory*. American Mathematical Soc., 1997, no. 92.
- [10] L. Xiao, “Dual averaging methods for regularized stochastic learning and online optimization,” *Journal of Machine Learning Research*, vol. 11, no. Oct, pp. 2543–2596, 2010.
- [11] D. A. Levin and Y. Peres, *Markov chains and mixing times*. American Mathematical Soc., 2017, vol. 107.
- [12] S. Boyd, A. Ghosh, B. Prabhakar, and D. Shah, “Randomized gossip algorithms,” *IEEE transactions on information theory*, vol. 52, no. 6, pp. 2508–2530, 2006.
- [13] U. von Luxburg, M. Hein, and A. Radl, “Hitting times, commute distances and the spectral gap for large random geometric graphs,” Tech. Rep., 2010.