

Sample Complexity of Asynchronous Q-Learning: Sharper Analysis and Variance Reduction

Gen Li, *Student Member, IEEE*, Yuting Wei[✉], *Member, IEEE*, Yuejie Chi[✉], *Senior Member, IEEE*,
Yuantao Gu[✉], *Senior Member, IEEE*, and Yuxin Chen[✉], *Senior Member, IEEE*

Abstract—Asynchronous Q-learning aims to learn the optimal action-value function (or Q-function) of a Markov decision process (MDP), based on a single trajectory of Markovian samples induced by a behavior policy. Focusing on a γ -discounted MDP with state space \mathcal{S} and action space \mathcal{A} , we demonstrate that the ℓ_∞ -based sample complexity of classical asynchronous Q-learning — namely, the number of samples needed to yield an entrywise ε -accurate estimate of the Q-function — is at most on the order of $\frac{1}{\mu_{\min}(1-\gamma)^5\varepsilon^2} + \frac{t_{\text{mix}}}{\mu_{\min}(1-\gamma)}$ up to some logarithmic factor, provided that a proper constant learning rate is adopted. Here, t_{mix} and μ_{\min} denote respectively the mixing time and the minimum state-action occupancy probability of the sample trajectory. The first term of this bound matches the sample complexity in the synchronous case with independent samples drawn from the stationary distribution of the trajectory. The second term reflects the cost taken for the empirical distribution of the Markovian trajectory to reach a steady state, which is incurred at the very beginning and becomes amortized as the algorithm runs. Encouragingly, the above bound improves upon the state-of-the-art result by a factor of at least $|\mathcal{S}||\mathcal{A}|$ for all scenarios, and by a factor of at least $t_{\text{mix}}|\mathcal{S}||\mathcal{A}|$ for any sufficiently small accuracy level ε . Further, we demonstrate that the scaling on the effective horizon $\frac{1}{1-\gamma}$ can be improved by means of variance reduction.

Index Terms—Model-free reinforcement learning, asynchronous Q-learning, Markovian samples, variance reduction, TD learning, mixing time.

I. INTRODUCTION

MODEL-FREE algorithms such as Q-learning [3] play a central role in recent breakthroughs of reinforcement learning (RL) [4]. In contrast to model-based algorithms that decouple model estimation and planning, model-free algorithms attempt to directly interact with the environment — in the form of a policy that selects actions based on perceived states of the environment — from the collected data samples, without modeling the environment explicitly. Therefore, model-free algorithms are able to process data in an online fashion and are often memory-efficient. Understanding and improving the sample efficiency of model-free algorithms lie at the core of recent research activity [5], whose importance is particularly evident for the class of RL applications in which data collection is costly and time-consuming (such as clinical trials, online advertisements, and so on).

The current paper concentrates on Q-learning, an off-policy model-free algorithm that seeks to learn the optimal action-value function by observing what happens under a behavior policy. The off-policy feature makes it appealing in various RL applications where it is infeasible to change the policy under evaluation on the fly. There are two basic update models in Q-learning. The first one is termed a *synchronous* setting, which hypothesizes on the existence of a simulator (also called a generative model); at each time, the simulator generates an independent sample for every state-action pair, and the estimates are updated simultaneously across all state-action pairs. The second model concerns an *asynchronous* setting, where only a single sample trajectory following a behavior policy is accessible; at each time, the algorithm updates its estimate of a single state-action pair using one state transition from the trajectory. Obviously, understanding the asynchronous setting is considerably more challenging than the synchronous model, due to the Markovian (and hence non-i.i.d.) nature of its sampling process.

Focusing on an infinite-horizon Markov decision process (MDP) with state space \mathcal{S} and action space \mathcal{A} , this work investigates asynchronous Q-learning on a single *Markovian trajectory* induced by a behavior policy. We ask a fundamental question:

Manuscript received September 27, 2020; revised August 3, 2021; accepted October 6, 2021. Date of publication October 14, 2021; date of current version December 23, 2021. The work of Yuting Wei was supported in part by NSF under Grant CCF-2007911, Grant DMS-2015447, and Grant CCF-2106778. The work of Yuejie Chi was supported in part by the Office of Naval Research (ONR) under Grant N00014-18-1-2142 and Grant N00014-19-1-2404; in part by the Army Research Office (ARO) under Grant W911NF-18-1-0303; and in part by NSF under Grant CCF-1806154, Grant CCF-2007911, and Grant CCF-2106778. The work of Yuantao Gu was supported by NSFC under Grant 61971266. The work of Yuxin Chen was supported in part by the Air Force Office of Scientific Research (AFOSR) Young Investigator Program (YIP) Award under Grant FA9550-19-1-0030; in part by ONR under Grant N00014-19-1-2120; in part by the ARO YIP Award under Grant W911NF-20-1-0097 and Grant W911NF-18-1-0303; and in part by NSF under Grant CCF-2106739, Grant CCF-1907661, Grant IIS-1900140, and Grant IIS-2100158. (Corresponding author: Yuting Wei.)

Gen Li and Yuxin Chen are with the Department of Electrical and Computer Engineering, Princeton University, Princeton, NJ 08544 USA (e-mail: genl@princeton.edu; yuxin.chen@princeton.edu).

Yuting Wei is with the Department of Statistics and Data Science, Wharton School, University of Pennsylvania, Philadelphia, PA 19104 USA (e-mail: ytwei@wharton.upenn.edu).

Yuejie Chi is with the Department of Electrical and Computer Engineering, Carnegie Mellon University, Pittsburgh, PA 15213 USA.

Yuantao Gu is with the Department of Electronic Engineering, Tsinghua University, Beijing 100084, China.

Communicated by R. Venkataramanan, Associate Editor for Machine Learning.

Color versions of one or more figures in this article are available at <https://doi.org/10.1109/TIT.2021.3120096>.

Digital Object Identifier 10.1109/TIT.2021.3120096

How many samples are needed for asynchronous Q-learning to learn the optimal Q-function?

Despite a considerable number of prior works analyzing this algorithm (ranging from the classical works [6], [7] to the very recent paper [2]), it remains unclear whether existing sample complexity analysis of asynchronous Q-learning is tight. As we shall elucidate momentarily, there exists a large gap — at least as large as $|\mathcal{S}||\mathcal{A}|$ — between the state-of-the-art sample complexity bound for asynchronous Q-learning [2] and the one derived for the synchronous counterpart [8]. This raises a natural desire to examine whether there is any bottleneck intrinsic to the asynchronous setting that significantly limits its performance.

A. Main Contributions

This paper develops a refined analysis framework that sharpens our understanding about the sample efficiency of classical asynchronous Q-learning on a single sample trajectory. Setting the stage, consider an infinite-horizon MDP with state space \mathcal{S} , action space \mathcal{A} , and a discount factor $\gamma \in (0, 1)$. What we have access to is a sample trajectory of the MDP induced by a stationary behavior policy. In contrast to the synchronous setting with i.i.d. samples, we single out two parameters intrinsic to the Markovian sample trajectory: (i) the mixing time t_{mix} , which characterizes how fast the trajectory disentangles itself from the initial state; (ii) the smallest state-action occupancy probability μ_{\min} of the stationary distribution of the trajectory, which captures how frequent each state-action pair has been at least visited.

With these parameters in place, our findings unveil that: the sample complexity required for asynchronous Q-learning to yield an ε -optimal Q-function estimate — in a strong ℓ_∞ sense — is at most¹

$$\tilde{O}\left(\frac{1}{\mu_{\min}(1-\gamma)^5\varepsilon^2} + \frac{t_{\text{mix}}}{\mu_{\min}(1-\gamma)}\right). \quad (1)$$

The first component of (1) is consistent with the sample complexity derived for the setting with independent samples drawn from the stationary distribution of the trajectory [8]. In comparison, the second term of (1) — which is unaffected by the accuracy level ε — is intrinsic to the Markovian nature of the trajectory; in essence, this term reflects the cost taken for the empirical distribution of the sample trajectory to converge to a steady state, and becomes amortized as the algorithm runs. In other words, the behavior of asynchronous Q-learning would resemble what happens in the setting with independent samples, as long as the algorithm has been run for reasonably long. In addition, our analysis framework readily yields another sample complexity bound

$$\tilde{O}\left(\frac{t_{\text{cover}}}{(1-\gamma)^5\varepsilon^2}\right), \quad (2)$$

where t_{cover} stands for the cover time — namely, the time taken for the trajectory to visit all state-action pairs at least

once. This facilitates comparisons with several prior results based on the cover time.

Furthermore, we leverage the idea of variance reduction to improve the scaling with the discount complexity $\frac{1}{1-\gamma}$. We demonstrate that a variance-reduced variant of asynchronous Q-learning attains ε -accuracy using at most

$$\tilde{O}\left(\frac{1}{\mu_{\min}(1-\gamma)^3\min\{1, \varepsilon^2\}} + \frac{t_{\text{mix}}}{\mu_{\min}(1-\gamma)}\right) \quad (3)$$

samples, matching the complexity of its synchronous counterpart if $\varepsilon \leq \min\left\{1, \frac{1}{(1-\gamma)\sqrt{t_{\text{mix}}}}\right\}$ [12]. Moreover, by taking the action space to be a singleton set, the aforementioned results immediately lead to ℓ_∞ -based sample complexity guarantees for temporal difference (TD) learning [13] on Markovian samples.

Comparisons with past results. A large fraction of the classical literature focused on asymptotic convergence analysis of asynchronous Q-learning (e.g. [6], [7], [14]); these results, however, did not lead to non-asymptotic sample complexity bounds. The state-of-the-art sample complexity analysis was due to the recent work [2], which derived a sample complexity bound $\tilde{O}\left(\frac{t_{\text{mix}}}{\mu_{\min}^2(1-\gamma)^5\varepsilon^2}\right)$. Given the obvious lower bound $1/\mu_{\min} \geq |\mathcal{S}||\mathcal{A}|$, our result (1) improves upon that of [2] by a factor at least on the order of $|\mathcal{S}||\mathcal{A}|\min\left\{t_{\text{mix}}, \frac{1}{(1-\gamma)^4\varepsilon^2}\right\}$. In particular, for sufficiently small accuracy level ε , our improvement exceeds a factor of at least

$$t_{\text{mix}}|\mathcal{S}||\mathcal{A}|.$$

In addition, we note that several prior works [9], [10] developed sample complexity bounds in terms of the cover time t_{cover} of the sample trajectory; our result strengthens these bounds by a factor of at least

$$t_{\text{cover}}^2|\mathcal{S}||\mathcal{A}| \geq |\mathcal{S}|^3|\mathcal{A}|^3.$$

The interested reader is referred to Table I for more precise comparisons, and to Section V for a discussion of further related works.

B. Paper Organization, Notation, and Basic Concept

The remainder of the paper is organized as follows. Section II formulates the problem and introduces some basic quantities and assumptions. Section III presents the asynchronous Q-learning algorithm along with its theoretical guarantees, whereas Section IV accommodates the extension: asynchronous variance-reduced Q-learning. A more detailed account of related works is given in Section V. The analyses of our main theorems are described in Sections VI-IX. We conclude this paper with a summary of our results and a list of future directions in Section X. Several preliminary facts about Markov chains and the proofs of technical lemmas are postponed to the appendix.

Next, we introduce a set of notation that will be used throughout the paper. Denote by $\Delta(\mathcal{S})$ (resp. $\Delta(\mathcal{A})$) the probability simplex over the set \mathcal{S} (resp. \mathcal{A}). For any vector $\mathbf{z} = [z_i]_{1 \leq i \leq n} \in \mathbb{R}^n$, we overload the notation $\sqrt{\cdot}$ and $|\cdot|$ to denote entry-wise operations, such that $\sqrt{\mathbf{z}} := [\sqrt{z_i}]_{1 \leq i \leq n}$ and $|\mathbf{z}| := [|z_i|]_{1 \leq i \leq n}$. For any vectors $\mathbf{z} = [a_i]_{1 \leq i \leq n}$ and

¹Let $\mathcal{X} := (|\mathcal{S}|, |\mathcal{A}|, \frac{1}{1-\gamma}, \frac{1}{\varepsilon})$. The notation $f(\mathcal{X}) = O(g(\mathcal{X}))$ means there exists a universal constant $C_1 > 0$ such that $f \leq C_1 g$. The notation $\tilde{O}(\cdot)$ is defined analogously except that it hides any logarithmic factor.

TABLE I

SAMPLE COMPLEXITY OF ASYNCHRONOUS Q-LEARNING AND ITS VARIANTS TO COMPUTE AN ε -OPTIMAL Q-FUNCTION IN THE ℓ_∞ NORM, WHERE WE HIDE ALL LOGARITHMIC FACTORS. WITH REGARDS TO THE MARKOVIAN TRAJECTORY INDUCED BY THE BEHAVIOR POLICY, WE DENOTE BY t_{cover} , t_{mix} , AND μ_{\min} THE COVER TIME, MIXING TIME, AND MINIMUM STATE-ACTION OCCUPANCY PROBABILITY OF THE ASSOCIATED STATIONARY DISTRIBUTION, RESPECTIVELY

Algorithm	Sample complexity	Learning rate
Asynchronous Q-learning Even-Dar and Mansour, 2003 [9]	$\frac{(t_{\text{cover}})^{\frac{1}{1-\gamma}}}{(1-\gamma)^4 \varepsilon^2}$	linear: $\frac{1}{t}$
Asynchronous Q-learning Even-Dar and Mansour, 2003 [9]	$\left(\frac{t_{\text{cover}}^{1+3\omega}}{(1-\gamma)^4 \varepsilon^2}\right)^{\frac{1}{\omega}} + \left(\frac{t_{\text{cover}}}{1-\gamma}\right)^{\frac{1}{1-\omega}}$	polynomial: $\frac{1}{t^\omega}$, $\omega \in (\frac{1}{2}, 1)$
Asynchronous Q-learning Beck and Srikant, 2012 [10]	$\frac{t_{\text{cover}}^3 S \mathcal{A} }{(1-\gamma)^5 \varepsilon^2}$	constant: $\frac{(1-\gamma)^4 \varepsilon^2}{ S \mathcal{A} t_{\text{cover}}^2}$
Asynchronous Q-learning Qu and Wierman, 2020 [2]	$\frac{t_{\text{mix}}}{\mu_{\min}^2 (1-\gamma)^5 \varepsilon^2}$	rescaled linear: $\frac{\frac{1}{\mu_{\min} (1-\gamma)}}{t + \max\{\frac{1}{\mu_{\min} (1-\gamma)}, t_{\text{mix}}\}}$
Speedy Q-learning Azar et al., 2011 [11]	$\frac{t_{\text{cover}}}{(1-\gamma)^4 \varepsilon^2}$	rescaled linear: $\frac{1}{t+1}$
Asynchronous Q-learning This work (Theorem 1)	$\frac{1}{\mu_{\min} (1-\gamma)^5 \varepsilon^2} + \frac{t_{\text{mix}}}{\mu_{\min} (1-\gamma)}$	constant: $\min\left\{\frac{(1-\gamma)^4 \varepsilon^2}{\gamma^2}, \frac{1}{t_{\text{mix}}}\right\}$
Asynchronous Q-learning This work (Theorem 2)	$\frac{t_{\text{cover}}}{(1-\gamma)^5 \varepsilon^2}$	constant: $\min\left\{\frac{(1-\gamma)^4 \varepsilon^2}{\gamma^2}, 1\right\}$
Asynchronous Q-learning This work (Theorem 3)	$\frac{1}{\mu_{\min} (1-\gamma)^5 \varepsilon^2} + \frac{t_{\text{mix}}}{\mu_{\min} (1-\gamma)}$	piecewise constant rescaled linear: (23)
Variance-reduced Q-learning This work (Theorem 4)	$\frac{1}{\mu_{\min} (1-\gamma)^3 \varepsilon^2} + \frac{t_{\text{mix}}}{\mu_{\min} (1-\gamma)}$	constant: $\min\left\{\frac{(1-\gamma)^2}{\gamma^2}, \frac{1}{t_{\text{mix}}}\right\}$

$\mathbf{w} = [w_i]_{1 \leq i \leq n}$, the notation $\mathbf{z} \geq \mathbf{w}$ (resp. $\mathbf{z} \leq \mathbf{w}$) means $z_i \geq w_i$ (resp. $z_i \leq w_i$) for all $1 \leq i \leq n$. Additionally, we denote by $\mathbf{1}$ the all-one vector, \mathbf{I} the identity matrix, and $\mathbb{1}\{\cdot\}$ the indicator function. For any matrix $\mathbf{P} = [P_{ij}]$, we denote $\|\mathbf{P}\|_1 := \max_i \sum_j |P_{ij}|$. Throughout this paper, we use c, c_0, c_1, \dots to denote universal constants that do not depend either on the parameters of the MDP or the target levels (ε, δ) , and their exact values may change from line to line.

Finally, let us introduce the concept of uniform ergodicity for Markov chains. Consider any Markov chain (X_0, X_1, X_2, \dots) with transition kernel P , finite state space \mathcal{X} and stationary distribution μ , and denote by $P^t(\cdot | x)$ the distribution of X_t conditioned on $X_0 = x \in \mathcal{X}$. This Markov chain is said to be *uniformly ergodic* if, for some $\rho < 1$ and $M < \infty$, one has

$$\sup_{x \in \mathcal{X}} d_{\text{TV}}(\mu, P^t(\cdot | x)) \leq M \rho^t, \quad (4)$$

where $d_{\text{TV}}(\mu, \nu)$ stands for the total variation distance between two distributions μ and ν [15]:

$$d_{\text{TV}}(\mu, \nu) := \frac{1}{2} \sum_{x \in \mathcal{X}} |\mu(x) - \nu(x)| = \sup_{A \subseteq \mathcal{X}} |\mu(A) - \nu(A)|. \quad (5)$$

II. MODELS AND BACKGROUND

This paper studies an infinite-horizon MDP with discounted rewards, as represented by a quintuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, r, \gamma)$. Here, \mathcal{S} and \mathcal{A} denote respectively the (finite) state space and action space, whereas $\gamma \in (0, 1)$ indicates the discount

factor. Particular emphasis is placed on the scenario with large state/action space and long effective horizon, namely, $|\mathcal{S}|$, $|\mathcal{A}|$ and the effective horizon $\frac{1}{1-\gamma}$ can all be quite large. We use $P : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ to represent the probability transition kernel of the MDP, where for each state-action pair $(s, a) \in \mathcal{S} \times \mathcal{A}$, $P(s' | s, a)$ denotes the probability of transiting to state s' from state s when action a is executed. The reward function is represented by $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$, such that $r(s, a)$ denotes the immediate reward from state s when action a is taken; for simplicity, we assume throughout that all rewards lie within $[0, 1]$. We focus on the tabular setting which, despite its basic form, has not yet been well understood. See [16] for an in-depth introduction of this model.

A. Q-Function and Bellman Operator

An action selection rule is termed a *policy* and represented by a mapping $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$, which maps a state to a distribution over the set of actions. A policy is said to be *stationary* if it is time-invariant. We denote by $\{s_t, a_t, r_t\}_{t=0}^\infty$ a sample trajectory, where s_t (resp. a_t) denotes the state (resp. the action taken) at time t , and $r_t = r(s_t, a_t)$ denotes the reward received at time t . It is assumed throughout that the rewards are deterministic and depend solely upon the current state-action pair. We denote by $V^\pi : \mathcal{S} \rightarrow \mathbb{R}$ the value function of a policy π , namely,

$$\forall s \in \mathcal{S} : \quad V^\pi(s) := \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s \right],$$

which is the expected discounted cumulative reward received when (i) the initial state is $s_0 = s$, (ii) the actions are taken

based on the policy π (namely, $a_t \sim \pi(s_t)$ for all $t \geq 0$) and the trajectory is generated based on the transition kernel (namely, $s_{t+1} \sim P(\cdot | s_t, a_t)$). It can be easily verified that $0 \leq V^\pi(s) \leq \frac{1}{1-\gamma}$ for any π . The action-value function (also Q-function) $Q^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ of a policy π is defined for all $(s, a) \in \mathcal{S} \times \mathcal{A}$ by

$$Q^\pi(s, a) := \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, a_0 = a \right],$$

where the actions are taken according to the policy π except the initial action (i.e. $a_t \sim \pi(s_t)$ for all $t \geq 1$). As is well-known, there exists an optimal policy — denoted by π^* — that simultaneously maximizes $V^\pi(s)$ and $Q^\pi(s, a)$ uniformly over all state-action pairs $(s, a) \in (\mathcal{S} \times \mathcal{A})$. Here and throughout, we shall denote by $V^* := V^{\pi^*}$ and $Q^* := Q^{\pi^*}$ the optimal value function and the optimal Q-function, respectively.

In addition, the Bellman operator \mathcal{T} , which is a mapping from $\mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|}$ to itself, is defined such that the (s, a) -th entry of $\mathcal{T}(Q)$ is given by

$$\mathcal{T}(Q)(s, a) := r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[\max_{a' \in \mathcal{A}} Q(s', a') \right]. \quad (6)$$

It is well known that the optimal Q-function Q^* is the unique fixed point of the Bellman operator.

B. Sample Trajectory and Behavior Policy

Imagine we have access to a sample trajectory $\{s_t, a_t, r_t\}_{t=0}^{\infty}$ generated by the MDP \mathcal{M} under a given stationary policy π_b — called a *behavior policy*. The behavior policy is deployed to help one learn the “behavior” of the MDP under consideration, which often differs from the optimal policy being sought. Given the stationarity of π_b , the sample trajectory can be viewed as a sample path of a time-homogeneous Markov chain over the set of state-action pairs $\{(s, a) \mid s \in \mathcal{S}, a \in \mathcal{A}\}$. Throughout this paper, we impose the following uniform ergodicity assumption [17] (see the definition of uniform ergodicity in Section I-B).

Assumption 1: The Markov chain induced by the stationary behavior policy π_b is uniformly ergodic.

There are several properties concerning the behavior policy and its resulting Markov chain that play a crucial role in learning the optimal Q-function. Specifically, denote by μ_{π_b} the stationary distribution (over all state-action pairs) of the aforementioned behavior Markov chain, and define

$$\mu_{\min} := \min_{(s, a) \in \mathcal{S} \times \mathcal{A}} \mu_{\pi_b}(s, a). \quad (7)$$

Intuitively, μ_{\min} reflects an information bottleneck; that is, the smaller μ_{\min} is, the more samples are needed in order to ensure all state-action pairs are visited sufficiently many times. In addition, we define the associated mixing time of the chain as

$$t_{\text{mix}} := \min \left\{ t \mid \max_{(s_0, a_0) \in \mathcal{S} \times \mathcal{A}} d_{\text{TV}}(P^t(\cdot | s_0, a_0), \mu_{\pi_b}) \leq \frac{1}{4} \right\}, \quad (8)$$

where $P^t(\cdot | s_0, a_0)$ denotes the distribution of (s_t, a_t) conditional on the initial state-action pair (s_0, a_0) , and $d_{\text{TV}}(\mu, \nu)$ is

the total variation distance between μ and ν (see (5)). In words, the mixing time t_{mix} captures how fast the sample trajectory decorrelates from its initial state. Moreover, we define the cover time associated with this Markov chain as follows

$$t_{\text{cover}} := \min \left\{ t \mid \min_{(s_0, a_0) \in \mathcal{S} \times \mathcal{A}} \mathbb{P}(\mathcal{B}_t | s_0, a_0) \geq \frac{1}{2} \right\}, \quad (9)$$

where \mathcal{B}_t denotes the event such that all $(s, a) \in \mathcal{S} \times \mathcal{A}$ have been visited at least once between time 0 and time t , and $\mathbb{P}(\mathcal{B}_t | s_0, a_0)$ denotes the probability of \mathcal{B}_t conditional on the initial state (s_0, a_0) .

Remark 1: It is known that for a finite-state Markov chain, having a finite mixing time t_{mix} implies uniform ergodicity of the chain [17, Page 4]. Thus, our uniform ergodicity assumption is equivalent to the assumption imposed in [2] (which assumes ergodicity in addition to a finite t_{mix}).

C. Goal

Given a single sample trajectory $\{s_t, a_t, r_t\}_{t=0}^{\infty}$ generated by the behavior policy π_b , we aim to compute/approximate the optimal Q-function Q^* in an ℓ_∞ sense. This setting — in which a state-action pair can be updated only when the Markovian trajectory reaches it — is commonly referred to as *asynchronous* Q-learning [2], [6] in tabular RL. The current paper focuses on characterizing, in a non-asymptotic manner, the sample efficiency of classical Q-learning and its variance-reduced variant.

III. ASYNCHRONOUS Q-LEARNING ON A SINGLE MARKOVIAN TRAJECTORY

A. Algorithm

The Q-learning algorithm [3] is arguably one of the most famous off-policy algorithms aimed at learning the optimal Q-function. Given the Markovian trajectory $\{s_t, a_t, r_t\}_{t=0}^{\infty}$ generated by the behavior policy π_b , the asynchronous Q-learning algorithm maintains a Q-function estimate $Q_t : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ at each time t and adopts the following iterative update rule

$$\begin{aligned} Q_t(s_{t-1}, a_{t-1}) &= (1 - \eta_t) Q_{t-1}(s_{t-1}, a_{t-1}) \\ &\quad + \eta_t \mathcal{T}_t(Q_{t-1})(s_{t-1}, a_{t-1}) \\ Q_t(s, a) &= Q_{t-1}(s, a), \quad \forall (s, a) \neq (s_{t-1}, a_{t-1}) \end{aligned} \quad (10)$$

for any $t \geq 0$, whereas η_t denotes the learning rate or the stepsize. Here, \mathcal{T}_t denotes the empirical Bellman operator w.r.t. the t -th sample, that is,

$$\mathcal{T}_t(Q)(s_{t-1}, a_{t-1}) := r(s_{t-1}, a_{t-1}) + \gamma \max_{a' \in \mathcal{A}} Q(s_{t-1}, a'). \quad (11)$$

It is worth emphasizing that at each time t , only a single entry — the one corresponding to the sampled state-action pair (s_{t-1}, a_{t-1}) — is updated, with all remaining entries unaltered. While the estimate Q_0 can be initialized to arbitrary values, we shall set $Q_0(s, a) = 0$ for all (s, a) unless otherwise noted. The corresponding value function estimate $V_t : \mathcal{S} \rightarrow \mathbb{R}$ at time t is thus given by

$$\forall s \in \mathcal{S} : \quad V_t(s) := \max_{a \in \mathcal{A}} Q_t(s, a). \quad (12)$$

The complete algorithm is described in Algorithm 1.

Algorithm 1: Asynchronous Q-Learning

-
- 1 **input parameters:** learning rates $\{\eta_t\}$, number of iterations T .
 - 2 **initialization:** $Q_0 = 0$.
 - 3 **for** $t = 1, 2, \dots, T$ **do**
 - 4 Draw action $a_{t-1} \sim \pi_b(s_{t-1})$, observe reward $r(s_{t-1}, a_{t-1})$, and draw next state $s_t \sim P(\cdot | s_{t-1}, a_{t-1})$.
 - 5 Update Q_t according to (10).
-

B. Theoretical Guarantees for Asynchronous Q-Learning

We are in a position to present our main theory regarding the non-asymptotic sample complexity of asynchronous Q-learning, for which the key parameters μ_{\min} and t_{mix} defined respectively in (7) and (8) play a vital role. The proof of this result is provided in Section VI.

Theorem 1 (Asynchronous Q-Learning): For the asynchronous Q-learning algorithm detailed in Algorithm 1, there exist some universal constants $c_0, c_1 > 0$ such that for any $0 < \delta < 1$ and $0 < \varepsilon \leq \frac{1}{1-\gamma}$, one has

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : |Q_T(s, a) - Q^*(s, a)| \leq \varepsilon$$

with probability at least $1 - \delta$, provided that the iteration number T and the learning rates $\eta_t \equiv \eta$ obey

$$T \geq \frac{c_0}{\mu_{\min}} \left\{ \frac{1}{(1-\gamma)^5 \varepsilon^2} + \frac{t_{\text{mix}}}{1-\gamma} \right\} \cdot \log \left(\frac{|\mathcal{S}||\mathcal{A}|T}{\delta} \right) \log \left(\frac{1}{(1-\gamma)^2 \varepsilon} \right), \quad (13a)$$

$$\eta = \frac{c_1}{\log \left(\frac{|\mathcal{S}||\mathcal{A}|T}{\delta} \right)} \min \left\{ \frac{(1-\gamma)^4 \varepsilon^2}{\gamma^2}, \frac{1}{t_{\text{mix}}} \right\}. \quad (13b)$$

Remark 2: The careful reader might immediately remark that the learning rate η studied in Theorem 1 relies on prior knowledge of ε , δ and T . This is more stringent than the learning rates in [2], which do not require pre-determining these parameters. To address this issue, we will explore a more adaptive learning rate schedule shortly in Section III-D, which achieves the same sample complexity without the need of knowing these parameters *a priori*.

Theorem 1 delivers a finite-sample/finite-time analysis of asynchronous Q-learning, given that a fixed learning rate is adopted and chosen appropriately. The ℓ_∞ -based sample complexity required for Algorithm 1 to attain ε accuracy is at most

$$\tilde{O} \left(\frac{1}{\mu_{\min}(1-\gamma)^5 \varepsilon^2} + \frac{t_{\text{mix}}}{\mu_{\min}(1-\gamma)} \right). \quad (14)$$

A few implications are in order.

1) *Dependency on the Minimum State-Action Occupancy Probability μ_{\min} :* Our sample complexity bound (14) scales linearly in $1/\mu_{\min}$, which is in general unimprovable. Consider, for instance, the ideal scenario where state-action occupancy is nearly uniform across all state-action pairs, in which case $1/\mu_{\min}$ is on the order of $|\mathcal{S}||\mathcal{A}|$. In such a “near-uniform” case, the sample complexity scales linearly with

$|\mathcal{S}||\mathcal{A}|$, and this dependency matches the known minimax lower bound [18] derived for the setting with independent samples. In comparison, [2, Theorem 7] depends at least quadratically on $1/\mu_{\min}$, which is at least $|\mathcal{S}||\mathcal{A}|$ times larger than our result (14).

2) *Dependency on the Effective Horizon $\frac{1}{1-\gamma}$:* The sample size bound (14) scales as $\frac{1}{(1-\gamma)^5 \varepsilon^2}$, which coincides with both [8], [19] (for the synchronous setting) and [2], [10] (for the asynchronous setting) with either a rescaled linear learning rate or a constant learning rate. This turns out to be the sharpest scaling known to date for the classical form of Q-learning.

3) *Dependency on the Mixing Time t_{mix} :* The second additive term of our sample complexity (14) depends linearly on the mixing time t_{mix} and is (almost) independent of the target accuracy ε . The influence of this mixing term is a consequence of the expense taken for the Markovian trajectory to reach a steady state, which is a one-time cost that can be amortized over later iterations if the algorithm is run for reasonably long. Put another way, if the behavior chain mixes not too slowly with respect to ε (in the sense that $t_{\text{mix}} \leq \frac{1}{(1-\gamma)^4 \varepsilon^2}$), then the algorithm behaves as if the samples were independently drawn from the stationary distribution of the trajectory. In comparison, the influences of t_{mix} and $\frac{1}{(1-\gamma)^5 \varepsilon^2}$ in [2] (cf. Table I) are multiplicative regardless of the value of ε , thus resulting in a much higher sample complexity. For instance, if $\varepsilon = O(\frac{1}{(1-\gamma)^2 \sqrt{t_{\text{mix}}}})$, then the sample complexity result therein is at least

$$\frac{t_{\text{mix}}}{\mu_{\min}} \geq t_{\text{mix}} |\mathcal{S}||\mathcal{A}|$$

times larger than our result (modulo some log factor).

4) *Schedule of Learning Rates:* An interesting aspect of our analysis lies in the adoption of a time-invariant learning rate, under which the ℓ_∞ error decays linearly — down to some error floor whose value is dictated by the learning rate. Therefore, a desired statistical accuracy can be achieved by properly setting the learning rate based on the target accuracy level ε and then determining the sample complexity accordingly. In comparison, classical analyses typically adopted a (rescaled) linear or a polynomial learning rule [2], [9]. While the work [10] studied Q-learning with a constant learning rate, their bounds were conservative and fell short of revealing the optimal scaling. Furthermore, we note that adopting time-invariant learning rates is not the only option that enables the advertised sample complexity; as we shall elucidate in Section III-D, one can also adopt carefully designed diminishing learning rates to achieve the same performance guarantees.

5) *Mean Estimation Error:* The high-probability bound in Theorem 1 readily translates to a mean estimation error guarantee. To see this, let us first make note of the following basic crude bound (see e.g. [10], [20])

$$|Q_t(s, a)| \leq \frac{1}{1-\gamma}, \quad |Q_t(s, a) - Q^*(s, a)| \leq \frac{1}{1-\gamma} \quad (15)$$

for all $t \geq 0$ and all $(s, a) \in \mathcal{S} \times \mathcal{A}$. By taking $\delta = \varepsilon(1 - \gamma)$ in Theorem 1, we immediately reach

$$\mathbb{E} \left[\max_{s,a} |Q_T(s, a) - Q^*(s, a)| \right] \leq \varepsilon(1 - \delta) + \delta \frac{1}{1 - \gamma} \leq 2\varepsilon, \quad (16)$$

provided that T obeys (13a). As a result, the sample complexity remains unchanged (up to some logarithmic factor) when the goal is to achieve the mean error bound $\mathbb{E} [\max_{s,a} |Q_T(s, a) - Q^*(s, a)|] \leq 2\varepsilon$.

In addition, our analysis framework immediately leads to another sample complexity guarantee stated in terms of the cover time t_{cover} (cf. (9)), which facilitates comparisons with several past work [9], [10]. The proof follows essentially that of Theorem 1, with a sketch provided in Section VII.

Theorem 2: For the asynchronous Q-learning algorithm detailed in Algorithm 1, there exist some universal constants $c_0, c_1 > 0$ such that for any $0 < \delta < 1$ and $0 < \varepsilon \leq \frac{1}{1-\gamma}$, one has

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : |Q_T(s, a) - Q^*(s, a)| \leq \varepsilon$$

with probability at least $1 - \delta$, provided that the iteration number T and the learning rates $\eta_t \equiv \eta$ obey

$$T \geq \frac{c_0 t_{\text{cover}}}{(1 - \gamma)^5 \varepsilon^2} \log^2 \left(\frac{|\mathcal{S}| |\mathcal{A}| T}{\delta} \right) \log \left(\frac{1}{(1 - \gamma)^2 \varepsilon} \right), \quad (17a)$$

$$\eta = \frac{c_1}{\log \left(\frac{|\mathcal{S}| |\mathcal{A}| T}{\delta} \right)} \min \left\{ \frac{(1 - \gamma)^4 \varepsilon^2}{\gamma^2}, 1 \right\}. \quad (17b)$$

Remark 3: The main difference between the cover-time-based analysis and the mixing-time-based analysis lies in the number of visits to each state-action pair (s, a) in every time frame. Owing to the measure concentration of Markov chains, we can see that the number of visits to each (s, a) concentrates around its expected value in each time frame, which in turn ensures that all state-action pairs have been visited at least once as long as the time frame is sufficiently long. This important property allows one to establish an intimate connection between the analysis of Theorem 1 and that of Theorem 2.

In a nutshell, this theorem tells us that the ℓ_∞ -based sample complexity of classical asynchronous Q-learning is bounded above by

$$\tilde{O} \left(\frac{t_{\text{cover}}}{(1 - \gamma)^5 \varepsilon^2} \right), \quad (18)$$

which scales linearly with the cover time. This improves upon the prior result [9] (resp. [10]) by an order of at least

$$t_{\text{cover}}^{3.29} \geq |\mathcal{S}|^{3.29} |\mathcal{A}|^{3.29} \quad (\text{resp. } t_{\text{cover}}^2 |\mathcal{S}| |\mathcal{A}| \geq |\mathcal{S}|^3 |\mathcal{A}|^3).$$

See Table I for detailed comparisons. We shall further make note of some connections between t_{cover} and $t_{\text{mix}}/\mu_{\min}$ to help compare Theorem 1 and Theorem 2: (i) in general, $t_{\text{cover}} = \tilde{O}(t_{\text{mix}}/\mu_{\min})$ for uniformly ergodic chains; (ii) one can find some cases where $t_{\text{mix}}/\mu_{\min} = \tilde{O}(t_{\text{cover}})$. Consequently, while Theorem 1 does not strictly dominate Theorem 2 in all instances, the aforementioned connections reveal that Theorem 1 is tighter for the worst-case scenarios. The interested reader is referred to Section B for details.

C. A Special Case: TD Learning

In the special circumstance that the set of allowable actions \mathcal{A} is a singleton, the corresponding MDP reduces to a Markov reward process (MRP), where the state transition kernel $P : \mathcal{S} \rightarrow \Delta(\mathcal{S})$ describes the probability of transitioning between different states, and $r : \mathcal{S} \rightarrow [0, 1]$ denotes the reward function (so that $r(s)$ is the immediate reward in state s). The goal is to estimate the value function $V : \mathcal{S} \rightarrow \mathbb{R}$ from the trajectory $\{s_t, r_t\}_{t=0}^\infty$, which arises commonly in the task of policy evaluation for a given deterministic policy.

The Q-learning procedure in this special setting reduces to the well-known TD learning algorithm, which maintains an estimate $V_t : \mathcal{S} \rightarrow \mathbb{R}$ at each time t and proceeds according to the following iterative update²

$$\begin{aligned} V_t(s_{t-1}) &= (1 - \eta_t) V_{t-1}(s_{t-1}) + \eta_t (r(s_{t-1}) + \gamma V_{t-1}(s_t)), \\ V_t(s) &= V_{t-1}(s), \quad \forall s \neq s_{t-1}. \end{aligned} \quad (19)$$

As usual, η_t denotes the learning rate at time t , and V_0 is taken to be 0. Consequently, our analysis for asynchronous Q-learning with a Markovian trajectory immediately leads to non-asymptotic ℓ_∞ guarantees for TD learning, stated below as a corollary of Theorem 1. A similar result can be stated in terms of the cover time as a corollary to Theorem 2, which we omit for brevity.

Corollary 1 (Asynchronous TD learning): Consider the TD learning algorithm (19). There exist some universal constants $c_0, c_1 > 0$ such that for any $0 < \delta < 1$ and $0 < \varepsilon \leq \frac{1}{1-\gamma}$, one has

$$\forall s \in \mathcal{S} : |V_T(s) - V(s)| \leq \varepsilon$$

with probability at least $1 - \delta$, provided that the iteration number T and the learning rates $\eta_t \equiv \eta$ obey

$$\begin{aligned} T &\geq \frac{c_0}{\mu_{\min}} \left\{ \frac{1}{(1 - \gamma)^5 \varepsilon^2} + \frac{t_{\text{mix}}}{1 - \gamma} \right\} \\ &\quad \cdot \log \left(\frac{|\mathcal{S}| T}{\delta} \right) \log \left(\frac{1}{(1 - \gamma)^2 \varepsilon} \right), \end{aligned} \quad (20a)$$

$$\eta = \frac{c_1}{\log \left(\frac{|\mathcal{S}| T}{\delta} \right)} \min \left\{ \frac{(1 - \gamma)^4 \varepsilon^2}{\gamma^2}, \frac{1}{t_{\text{mix}}} \right\}. \quad (20b)$$

The above result reveals that the ℓ_∞ -sample complexity for TD learning is at most

$$\tilde{O} \left(\frac{1}{\mu_{\min} (1 - \gamma)^5 \varepsilon^2} + \frac{t_{\text{mix}}}{\mu_{\min} (1 - \gamma)} \right), \quad (21)$$

provided that an appropriate constant learning rate is adopted. We note that prior finite-sample analysis on asynchronous TD learning typically focused on (weighted) ℓ_2 estimation errors with linear function approximation [21], [22], and it is hence difficult to make fair comparisons. The recent paper [23] developed ℓ_∞ guarantees for TD learning, focusing on the synchronous settings with i.i.d. samples rather than Markovian samples.

²When $\mathcal{A} = \{a\}$ is a singleton, the Q-learning update rule (10) reduces to the TD update rule (19) by relating $Q(s, a) = V(s)$.

D. Adaptive and Implementable Learning Rates

As alluded to previously, the learning rates recommended in (13b) depend on the mixing time t_{mix} , a parameter that might be either *a priori* unknown or difficult to estimate. Fortunately, it is feasible to adopt a more adaptive learning rate schedule, which does not rely on prior knowledge of t_{mix} while still being capable of achieving the performance advertised in Theorem 1.

1) *Learning Rates*: In order to describe our new learning rate schedule, we need to keep track of the following quantities for all $(s, a) \in \mathcal{S} \times \mathcal{A}$:

- $K_t(s, a)$: the number of times that the sample trajectory visits (s, a) during the first t iterations.

In addition, we maintain an estimate $\hat{\mu}_{\min, t}$ of μ_{\min} , computed recursively as follows

$$\hat{\mu}_{\min, t} = \begin{cases} \frac{1}{|\mathcal{S}||\mathcal{A}|}, & \text{if } \min_{s,a} K_t(s, a) = 0; \\ \hat{\mu}_{\min, t-1}, & \text{if } \frac{1}{3} < \frac{\min_{s,a} K_t(s, a)/t}{\hat{\mu}_{\min, t-1}} < 3; \\ \min_{s,a} K_t(s, a)/t, & \text{otherwise.} \end{cases} \quad (22)$$

With the above quantities in place, we propose the following learning rate schedule:

$$\eta_t = \min \left\{ 1, c_\eta \exp \left(\left\lfloor \log \frac{\log t}{\hat{\mu}_{\min, t}(1-\gamma)\gamma^2 t} \right\rfloor \right) \right\}, \quad (23)$$

where $c_\eta > 0$ is some universal constant independent of any MDP parameter³ and $\lfloor x \rfloor$ denotes the nearest integer less than or equal to x . If $\hat{\mu}_{\min, t}$ forms a reliable estimate of μ_{\min} , then one can view (23) as a sort of “piecewise constant approximation” of the rescaled linear stepsizes $\frac{c_\eta \log t}{\mu_{\min}(1-\gamma)\gamma^2 t}$; in fact, this can be viewed as a sort of “doubling trick” — reducing the learning rate by a constant factor every once a while — to approximate rescaled linear learning rates. Theorem 1 can then be readily applied to analyze the performance for each constant segment of this learning rate schedule (23). Noteworthy, such learning rates are fully data-driven and do not rely on any prior knowledge about the Markov chain (like t_{mix} and μ_{\min}) or the target accuracy level ε .

2) *Performance Guarantees*: Encouragingly, our theoretical framework can be readily extended without difficulty to accommodate this adaptive learning rate choice. Specifically, for the Q-function estimates

$$\hat{Q}_t = \begin{cases} Q_t, & \text{if } \eta_{t+1} \neq \eta_t, \\ \hat{Q}_{t-1}, & \text{otherwise,} \end{cases} \quad (24)$$

where Q_t is provided by the Q-learning iterations (cf. (10)). We can then establish the following theoretical guarantees, whose proof is deferred to Section VIII.

Theorem 3: Consider asynchronous Q-learning with learning rates (23) and the output (24). There exists some universal constant $C > 0$ such that: for any $0 < \delta < 1$ and $0 < \varepsilon \leq \frac{1}{1-\gamma}$, one has

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad |\hat{Q}_T(s, a) - Q^*(s, a)| \leq \varepsilon \quad (25)$$

³More precisely, $c_\eta > 0$ can be any universal constant obeying $c_\eta \geq 74c_0c_1$ and $c_\eta > 11$, with c_0 and c_1 being the universal constants stated in Theorem 1.

with probability at least $1 - \delta$, provided that

$$T \geq \frac{C}{\gamma^2} \max \left\{ \frac{1}{\mu_{\min}(1-\gamma)^5 \varepsilon^2}, \frac{t_{\text{mix}}}{\mu_{\min}(1-\gamma)} \right\} \cdot \log \left(\frac{|\mathcal{S}||\mathcal{A}|T}{\delta} \right) \log \left(\frac{T}{(1-\gamma)^2 \varepsilon} \right). \quad (26)$$

Remark 4: The interested reader might wonder whether our sample complexity guarantees continue to hold under the linear learning rate $\eta_t = \frac{1}{K_t(s_t, a_t)}$ — a learning rate schedule that has been previously studied in [6], [9]. Nevertheless, as discussed in [8, Section 3.3.1], this linear learning rate can lead to a sample complexity that scales exponentially in the effective horizon $\frac{1}{1-\gamma}$, which is clearly outperformed by a properly rescaled linear learning rate.

IV. EXTENSION: ASYNCHRONOUS VARIANCE-REDUCED Q-LEARNING

As pointed out in prior literature, the classical form of Q-learning (10) often suffers from sub-optimal dependence on the effective horizon $\frac{1}{1-\gamma}$. For instance, in the synchronous setting, the minimax lower bound is proportional to $\frac{1}{(1-\gamma)^3}$ (see, [18]), while the sharpest known upper bound for vanilla Q-learning scales as $\frac{1}{(1-\gamma)^5}$; see detailed discussions in [8]. To remedy this issue, recent work proposed to leverage the idea of variance reduction to develop accelerated RL algorithms in the synchronous setting [12], [24], as inspired by the seminal SVRG algorithm [25] that originates from the stochastic optimization literature. In this section, we adapt this idea to asynchronous Q-learning and characterize its sample efficiency.

A. Algorithm

In order to accelerate the convergence, it is instrumental to reduce the variability of the empirical Bellman operator \mathcal{T}_t employed in the update rule (10) of classical Q-learning. This can be achieved via the following means. Simply put, assuming we have access to (i) a reference Q-function estimate, denoted by \bar{Q} , and (ii) an estimate of $\mathcal{T}(\bar{Q})$, denoted by $\tilde{\mathcal{T}}(\bar{Q})$, the variance-reduced Q-learning update rule is given by

$$\begin{aligned} Q_t(s_{t-1}, a_{t-1}) &= (1 - \eta_t)Q_{t-1}(s_{t-1}, a_{t-1}) \\ &\quad + \eta_t \left(\mathcal{T}_t(Q_{t-1}) - \mathcal{T}_t(\bar{Q}) + \tilde{\mathcal{T}}(\bar{Q}) \right)(s_{t-1}, a_{t-1}), \\ Q_t(s, a) &= Q_{t-1}(s, a), \quad \forall (s, a) \neq (s_{t-1}, a_{t-1}), \end{aligned} \quad (27)$$

where \mathcal{T}_t denotes the empirical Bellman operator at time t (cf. (11)). The empirical estimate $\tilde{\mathcal{T}}(\bar{Q})$ can be computed using a set of samples; more specifically, by drawing N consecutive sample transitions $\{(s_i, a_i, s_{i+1})\}_{0 \leq i < N}$ from the observed trajectory, we compute

$$\begin{aligned} \tilde{\mathcal{T}}(\bar{Q})(s, a) &= r(s, a) \\ &\quad + \frac{\gamma \sum_{i=0}^{N-1} \mathbb{1}\{(s_i, a_i) = (s, a)\} \max_{a'} \bar{Q}(s_{i+1}, a')}{\sum_{i=0}^{N-1} \mathbb{1}\{(s_i, a_i) = (s, a)\}}. \end{aligned} \quad (28)$$

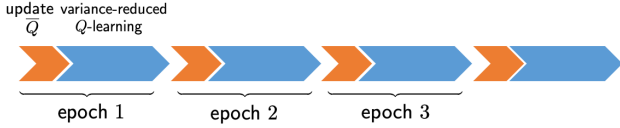


Fig. 1. A pictorial illustration of variance-reduced Q-learning.

Compared with the classical form (10), the original update term $\mathcal{T}_t(Q_{t-1})$ has been replaced by $\mathcal{T}_t(Q_{t-1}) - \mathcal{T}_t(\bar{Q}) + \tilde{\mathcal{T}}(\bar{Q})$, in the hope of achieving reduced variance as long as \bar{Q} (which serves as a proxy to Q^*) is chosen properly.

We now take a moment to elucidate the rationale behind the variance-reduced update rule (27). In the vanilla Q-learning update rule (10), the variability in each iteration (conditional on the past) comes primarily from the stochastic term $\mathcal{T}_t(Q_{t-1})$. In order to accelerate convergence, it is advisable to reduce the variability of this term. Suppose now that we have access to a reference point \bar{Q} that is close to Q_{t-1} . By replacing $\mathcal{T}_t(Q_{t-1})$ with

$$\{\mathcal{T}_t(Q_{t-1}) - \mathcal{T}_t(\bar{Q})\} + \tilde{\mathcal{T}}(\bar{Q}),$$

we see that the variability of the first term $\mathcal{T}_t(Q_{t-1}) - \mathcal{T}_t(\bar{Q})$ can be small if $Q_{t-1} \approx \bar{Q}$, while the uncertainty of the second term $\tilde{\mathcal{T}}(\bar{Q})$ can also be well controlled via the use of batch data. Motivated by this simple idea, the variance-reduced Q-learning rule attempts to operate in an epoch-based manner, computing $\tilde{\mathcal{T}}(\bar{Q})$ once every epoch (so as not to increase the overall sampling burden) and leveraging it to help reduce variability.

For convenience of presentation, we introduce the following notation

$$Q = \text{VR-Q-RUN-EPOCH}(\bar{Q}, N, t_{\text{epoch}}) \quad (29)$$

to represent the above-mentioned update rule, which starts with a reference point \bar{Q} and operates upon a total number of $N + t_{\text{epoch}}$ consecutive sample transitions. The first N samples are employed to construct $\tilde{\mathcal{T}}(\bar{Q})$ via (28), with the remaining samples employed in t_{epoch} iterative updates (27); see Algorithm 3. To achieve the desired acceleration, the proxy \bar{Q} needs to be periodically updated so as to better approximate the truth Q^* and hence reduce the bias. It is thus natural to run the algorithm in a multi-epoch manner. Specifically, we divide the samples into contiguous subsets called epochs, each containing t_{epoch} iterations and using $N + t_{\text{epoch}}$ samples. We then proceed as follows

$$Q_m^{\text{epoch}} = \text{VR-Q-RUN-EPOCH}(Q_{m-1}^{\text{epoch}}, N, t_{\text{epoch}}) \quad (30)$$

for $m = 1, \dots, M$, where M is the total number of epochs, and Q_m^{epoch} denotes the output of the m -th epoch. The whole procedure is summarized in Algorithm 2. Clearly, the total number of samples used in this algorithm is given by $M(N + t_{\text{epoch}})$. We remark that the idea of performing variance reduction in RL is certainly not new, and has been explored in a number of recent works [12], [23], [24], [26]–[28].

B. Theoretical Guarantees for Variance-Reduced Q-Learning

This subsection develops a non-asymptotic sample complexity bound for asynchronous variance-reduced Q-learning on a single trajectory. Before presenting our theoretical guarantees, there are several algorithmic parameters that we shall specify; for given target levels (ε, δ) , choose

$$\eta_t \equiv \eta = \frac{c_0}{\log\left(\frac{|\mathcal{S}||\mathcal{A}|t_{\text{epoch}}}{\delta}\right)} \min\left\{\frac{(1-\gamma)^2}{\gamma^2}, \frac{1}{t_{\text{mix}}}\right\}, \quad (31a)$$

$$N \geq \frac{c_1}{\mu_{\min}} \left(\frac{1}{(1-\gamma)^3 \min\{1, \varepsilon^2\}} + t_{\text{mix}} \right) \cdot \log\left(\frac{|\mathcal{S}||\mathcal{A}|t_{\text{epoch}}}{\delta}\right), \quad (31b)$$

$$t_{\text{epoch}} \geq \frac{c_2}{\mu_{\min}} \left(\frac{1}{(1-\gamma)^3} + \frac{t_{\text{mix}}}{1-\gamma} \right) \cdot \log\left(\frac{1}{(1-\gamma)^2 \varepsilon}\right) \log\left(\frac{|\mathcal{S}||\mathcal{A}|t_{\text{epoch}}}{\delta}\right), \quad (31c)$$

where $c_0 > 0$ is some sufficiently small constant, $c_1, c_2 > 0$ are some sufficiently large constants, and we recall the definitions of μ_{\min} and t_{mix} in (7) and (8), respectively. Note that the learning rate (31a) chosen here could be larger than the choice (13b) for the classical form by a factor of $O(\frac{1}{(1-\gamma)^2})$ (which happens if t_{mix} is not too large), allowing the algorithm to progress more aggressively.

Theorem 4 (Asynchronous Variance-Reduced Q-Learning): Let Q_M^{epoch} be the output of Algorithm 2 with parameters chosen according to (31). There exists some constant $c_3 > 0$ such that for any $0 < \delta < 1$ and $0 < \varepsilon \leq \frac{1}{1-\gamma}$, one has

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : \quad |Q_M^{\text{epoch}}(s, a) - Q^*(s, a)| \leq \varepsilon$$

with probability at least $1 - \delta$, provided that the total number of epochs exceeds

$$M \geq c_3 \log \frac{1}{\varepsilon(1-\gamma)^2}. \quad (32)$$

The proof of this result is postponed to Section IX.

In view of Theorem 4, the ℓ_∞ -based sample complexity for variance-reduced Q-learning to yield ε accuracy — which is characterized by $M(N + t_{\text{epoch}})$ — can be as low as

$$\tilde{O}\left(\frac{1}{\mu_{\min}(1-\gamma)^3 \min\{1, \varepsilon^2\}} + \frac{t_{\text{mix}}}{\mu_{\min}(1-\gamma)}\right). \quad (33)$$

Except for the second term that depends on the mixing time, the first term matches the result of [12] derived for the synchronous settings with independent samples. In the range $\varepsilon \in (0, \min\{1, \frac{1}{(1-\gamma)\sqrt{t_{\text{mix}}}}\}]$, the sample complexity reduce to $\tilde{O}(\frac{1}{\mu_{\min}(1-\gamma)^3 \varepsilon^2})$; the scaling $\frac{1}{(1-\gamma)^3}$ matches the minimax lower bound derived in [18] for the synchronous setting.

Once again, we can immediately deduce guarantees for asynchronous variance-reduced TD learning by reducing the action space to a singleton set (akin to Section III-C), which extends the analysis [23] to Markovian noise. In addition, similar to Section III-D, we can also employ adaptive learning rates in variance-reduced Q-learning — which do not require prior knowledge of t_{mix} and μ_{\min} — without compromising

Algorithm 2: Asynchronous Variance-Reduced Q-Learning

```

1 input parameters: number of epochs  $M$ , epoch length
   $t_{\text{epoch}}$ , recentering length  $N$ , learning rate  $\eta$ .
2 initialization: set  $Q_0^{\text{epoch}} \leftarrow 0$ .
3 for each epoch  $m = 1, \dots, M$  do
  | /* Call Algorithm 3.                                     */
4   $Q_m^{\text{epoch}} = \text{VR-Q-RUN-EPOCH}(Q_{m-1}^{\text{epoch}}, N, t_{\text{epoch}})$ .
  
```

Algorithm 3: function $Q = \text{VR-Q-RUN-EPOCH}(\bar{Q}, N, t_{\text{epoch}})$

```

1 Draw  $N$  new consecutive samples from the sample
  trajectory; compute  $\tilde{T}(\bar{Q})$  according to (28).
2 Set  $s_0 \leftarrow$  current state, and  $Q_0 \leftarrow \bar{Q}$ .
3 for  $t = 1, 2, \dots, t_{\text{epoch}}$  do
4  | Draw action  $a_{t-1} \sim \pi_b(s_{t-1})$ , observe reward
    |  $r(s_{t-1}, a_{t-1})$ , and draw next state
    |  $s_t \sim P(\cdot | s_{t-1}, a_{t-1})$ .
5  | Update  $Q_t$  according to (27).
6 return:  $Q \leftarrow Q_{t_{\text{epoch}}}$ .
  
```

the sample complexity. For the sake of brevity, we omit these extensions in the current paper.

V. RELATED WORKS

In this section, we review several recent lines of works and compare our results with them.

A. The Q-Learning Algorithm and Its Variants

The Q-learning algorithm, originally proposed in [29], has been analyzed in the asymptotic regime by [6], [7], [14], [30] since more than two decades ago. Additionally, finite-time performance of Q-learning and its variants have been analyzed by [2], [8]–[10], [19], [31]–[34] in the tabular setting, by [21], [35]–[43] in the context of function approximations, and by [44] with nonparametric regression. In addition, [11], [12], [24], [45]–[47] studied modified Q-learning algorithms that might potentially improve sample complexities and accelerate convergence. Another line of work studied Q-learning with sophisticated exploration strategies such as UCB exploration (e.g. [48]–[51]), which is beyond the scope of the current work.

B. Finite-Sample ℓ_∞ Guarantees for Q-Learning

We now expand on non-asymptotic ℓ_∞ guarantees available in prior literature, which are the most relevant to the current work. An interesting aspect that we shall highlight is the importance of learning rates. For instance, when a linear learning rate (i.e. $\eta_t = 1/t$) is adopted, the sample complexity results derived in past works [9], [14] exhibit an exponential blow-up in $\frac{1}{1-\gamma}$, which is clearly undesirable. In the synchronous setting, [8]–[10], [19] studied the finite-sample complexity of Q-learning under various learning rate rules; the best

sample complexity known to date is $\tilde{O}(\frac{|S||A|}{(1-\gamma)^3 \varepsilon^2})$, achieved via either a rescaled linear learning rate [8], [19] or a constant learning rate [19]. When it comes to asynchronous Q-learning (in its classical form), our work provides the first analysis that achieves linear scaling with $1/\mu_{\min}$ or t_{cover} ; see Table I for detailed comparisons. Going beyond classical Q-learning, the speedy Q-learning algorithm, which adds a momentum term in the update by using previous Q-function estimates, provably achieves a sample complexity of $\tilde{O}(\frac{t_{\text{cover}}}{(1-\gamma)^4 \varepsilon^2})$ [11] in the asynchronous setting, whose update rule takes twice the storage of classical Q-learning. However, the proof idea adopted in the speedy Q-learning paper relies heavily on the specific update rules of speedy Q-learning, which cannot be readily used here to help improve the sample complexity of asynchronous Q-learning in terms of its dependency on $\frac{1}{1-\gamma}$. In comparison, our analysis of the variance-reduced Q-learning algorithm achieves a sample complexity of $\tilde{O}(\frac{1}{\mu_{\min}(1-\gamma)^3 \varepsilon^2} + \frac{t_{\text{mix}}}{\mu_{\min}(1-\gamma)})$ when $\varepsilon < 1$.

C. Finite-Sample Guarantees for Model-Free Algorithms

Convergence properties of several model-free RL algorithms have been studied recently in the presence of Markovian data, including but not limited to TD learning and its variants [21], [22], [28], [52]–[60], Q-learning [35], [36], and SARSA [61]. However, these recent papers typically focused on the (weighted) ℓ_2 error rather than the ℓ_∞ risk, where the latter is often more relevant in the context of RL. In addition, [23] investigated the ℓ_∞ bounds of (variance-reduced) TD learning, although they did not account for Markovian noise.

D. Finite-Sample Guarantees for Model-Based Algorithms

Another contrasting approach for learning the optimal Q-function is the class of model-based algorithms, which has been shown to enjoy minimax-optimal sample complexity in the synchronous setting. More precisely, it is known that by planning over an empirical MDP constructed from $\tilde{O}(\frac{|S||A|}{(1-\gamma)^3 \varepsilon^2})$ samples, we are guaranteed to find not only an ε -optimal Q-function but also an ε -optimal policy [18], [62], [63]. It is worth emphasizing that the minimax optimality of model-based approach has been shown to hold for the entire ε -range; in comparison, the sample optimality of the model-free approach has only been shown for a smaller range of accuracy level ε in the synchronous setting. We also remark that existing sample complexity analysis for model-based approaches might be generalizable to Markovian data.

VI. ANALYSIS OF ASYNCHRONOUS Q-LEARNING

This section is devoted to establishing Theorem 1. Before proceeding, we find it convenient to introduce some matrix notation. Let $\mathbf{\Lambda}_t \in \mathbb{R}^{|S||A| \times |S||A|}$ be a diagonal matrix obeying

$$\mathbf{\Lambda}_t((s, a), (s, a)) := \begin{cases} \eta, & \text{if } (s, a) = (s_{t-1}, a_{t-1}), \\ 0, & \text{otherwise,} \end{cases} \quad (34)$$

where $\eta > 0$ is the learning rate. In addition, we use the vector $\mathbf{Q}_t \in \mathbb{R}^{|S||A|}$ (resp. $\mathbf{V}_t \in \mathbb{R}^{|S|}$) to represent our estimate Q_t

(resp. V_t) in the t -th iteration, so that the (s, a) -th (resp. s th) entry of \mathbf{Q}_t (resp. \mathbf{V}_t) is given by $Q_t(s, a)$ (resp. $V_t(s)$). Similarly, let the vectors $\mathbf{Q}^* \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ and $\mathbf{V}^* \in \mathbb{R}^{|\mathcal{S}|}$ represent the optimal Q-function Q^* and the optimal value function V^* , respectively. We also let the vector $\mathbf{r} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ stand for the reward function r , so that the (s, a) -th entry of \mathbf{r} is given by $r(s, a)$. In addition, we define the matrix $\mathbf{P}_t \in \{0, 1\}^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}|}$ such that

$$\mathbf{P}_t((s, a), s') := \begin{cases} 1, & \text{if } (s, a, s') = (s_{t-1}, a_{t-1}, s_t), \\ 0, & \text{otherwise.} \end{cases} \quad (35)$$

Clearly, this set of notation allows us to express the Q-learning update rule (10) in the following matrix form

$$\mathbf{Q}_t = (\mathbf{I} - \mathbf{\Lambda}_t)\mathbf{Q}_{t-1} + \mathbf{\Lambda}_t(\mathbf{r} + \gamma\mathbf{P}_t\mathbf{V}_{t-1}). \quad (36)$$

A. Error Decay in the Presence of Constant Learning Rates

The main step of the analysis is to establish the following result concerning the dynamics of asynchronous Q-learning. In order to state it formally, we find it convenient to introduce several auxiliary quantities

$$t_{\text{frame}} := \frac{443t_{\text{mix}}}{\mu_{\min}} \log\left(\frac{4|\mathcal{S}||\mathcal{A}|T}{\delta}\right), \quad (37a)$$

$$t_{\text{th}} := \max\left\{\frac{2\log\frac{1}{(1-\gamma)^2\varepsilon}}{\eta\mu_{\min}}, t_{\text{frame}}\right\}, \quad (37b)$$

$$\mu_{\text{frame}} := \frac{1}{2}\mu_{\min}t_{\text{frame}}, \quad (37c)$$

$$\rho := (1-\gamma)(1-(1-\eta)^{\mu_{\text{frame}}}). \quad (37d)$$

With these quantities in mind, we have the following result.

Theorem 5: Consider the asynchronous Q-learning algorithm in Algorithm 1 with $\eta_t \equiv \eta$. For any $\delta \in (0, 1)$ and any $\varepsilon \in (0, \frac{1}{1-\gamma}]$, there exists a universal constant $c > 0$ such that with probability at least $1 - 6\delta$, the following relation holds uniformly for all $t \leq T$ (defined in (13a))

$$\begin{aligned} \|\mathbf{Q}_t - \mathbf{Q}^*\|_{\infty} &\leq (1-\rho)^k \frac{\|\mathbf{Q}_0 - \mathbf{Q}^*\|_{\infty}}{1-\gamma} \\ &\quad + \frac{c\gamma}{1-\gamma} \|\mathbf{V}^*\|_{\infty} \sqrt{\eta \log\left(\frac{|\mathcal{S}||\mathcal{A}|T}{\delta}\right)} + \varepsilon, \end{aligned} \quad (38)$$

provided that $0 < \eta \log\left(\frac{|\mathcal{S}||\mathcal{A}|T}{\delta}\right) < 1$. Here, we define $k := \max\left\{0, \left\lfloor \frac{t-t_{\text{th}}}{t_{\text{frame}}} \right\rfloor\right\}$.

In words, Theorem 5 asserts that the ℓ_{∞} estimation error decays linearly — in a blockwise manner — to some error floor that scales with $\sqrt{\eta}$. This result suggests how to set the learning rate based on the target accuracy level, which in turn allows us to pin down the sample complexity under consideration. In what follows, we shall first establish Theorem 5, and then return to prove Theorem 1 using this result.

Before embarking on the proof of Theorem 5, we would like to point out a few key technical ingredients: (i) an epoch-based analysis that focuses on macroscopic dynamics as opposed to per-iteration dynamics, (ii) measure concentration of Markov chains (see Section A) that helps reveal the similarity between epoch-based dynamics and the synchronous

counterpart, and (iii) careful analysis of recursive relations. These key ingredients taken collectively lead to a sample complexity bound that improves upon prior analysis in [2].

B. Proof of Theorem 5

We are now positioned to outline the proof of Theorem 5. We remind the reader that for any two vectors $\mathbf{z} = [z_i]$ and $\mathbf{w} = [w_i]$, the notation $\mathbf{z} \leq \mathbf{w}$ (resp. $\mathbf{z} \geq \mathbf{w}$) denotes entrywise comparison (cf. Section I), meaning that $z_i \leq w_i$ (resp. $z_i \geq w_i$) holds for all i . As a result, for any non-negative matrix \mathbf{A} , one has $\mathbf{Az} \leq \mathbf{Aw}$ as long as $\mathbf{z} \leq \mathbf{w}$.

1) *Key Decomposition and a Recursive Formula:* The starting point of our proof is the following elementary decomposition

$$\begin{aligned} \Delta_t &:= \mathbf{Q}_t - \mathbf{Q}^* \\ &= (\mathbf{I} - \mathbf{\Lambda}_t)\mathbf{Q}_{t-1} + \mathbf{\Lambda}_t(\mathbf{r} + \gamma\mathbf{P}_t\mathbf{V}_{t-1}) - \mathbf{Q}^* \\ &= (\mathbf{I} - \mathbf{\Lambda}_t)(\mathbf{Q}_{t-1} - \mathbf{Q}^*) + \mathbf{\Lambda}_t(\mathbf{r} + \gamma\mathbf{P}_t\mathbf{V}_{t-1} - \mathbf{Q}^*) \\ &= (\mathbf{I} - \mathbf{\Lambda}_t)(\mathbf{Q}_{t-1} - \mathbf{Q}^*) + \gamma\mathbf{\Lambda}_t(\mathbf{P}_t\mathbf{V}_{t-1} - \mathbf{P}\mathbf{V}^*) \\ &= (\mathbf{I} - \mathbf{\Lambda}_t)\Delta_{t-1} + \gamma\mathbf{\Lambda}_t(\mathbf{P}_t - \mathbf{P})\mathbf{V}^* + \gamma\mathbf{\Lambda}_t\mathbf{P}_t(\mathbf{V}_{t-1} - \mathbf{V}^*) \end{aligned} \quad (39)$$

for any $t > 0$, where the first line results from the update rule (36), and the penultimate line follows from the Bellman equation $\mathbf{Q}^* = \mathbf{r} + \gamma\mathbf{P}\mathbf{V}^*$ (see [16]). Applying this relation recursively gives

$$\begin{aligned} \Delta_t &= \gamma \underbrace{\sum_{i=1}^t \prod_{j=i+1}^t (\mathbf{I} - \mathbf{\Lambda}_j) \mathbf{\Lambda}_i (\mathbf{P}_i - \mathbf{P}) \mathbf{V}^*}_{=: \beta_{1,t}} \\ &\quad + \gamma \underbrace{\sum_{i=1}^t \prod_{j=i+1}^t (\mathbf{I} - \mathbf{\Lambda}_j) \mathbf{\Lambda}_i \mathbf{P}_i (\mathbf{V}_{i-1} - \mathbf{V}^*)}_{=: \beta_{2,t}} \\ &\quad + \underbrace{\prod_{j=1}^t (\mathbf{I} - \mathbf{\Lambda}_j) \Delta_0}_{=: \beta_{3,t}}. \end{aligned} \quad (40)$$

Applying the triangle inequality, we obtain

$$|\Delta_t| \leq |\beta_{1,t}| + |\beta_{2,t}| + |\beta_{3,t}|, \quad (41)$$

where we recall the notation $|\mathbf{z}| := [z_i]_{1 \leq i \leq n}$ for any vector $\mathbf{z} = [z_i]_{1 \leq i \leq n}$. In what follows, we shall look at these terms separately.

- First of all, given that $\mathbf{I} - \mathbf{\Lambda}_j$ and $\mathbf{\Lambda}_j$ are both non-negative diagonal matrices and that

$$\begin{aligned} \|\mathbf{P}_i(\mathbf{V}_{i-1} - \mathbf{V}^*)\|_{\infty} &\leq \|\mathbf{P}_i\|_1 \|\mathbf{V}_{i-1} - \mathbf{V}^*\|_{\infty} \\ &= \|\mathbf{V}_{i-1} - \mathbf{V}^*\|_{\infty} \\ &\leq \|\mathbf{Q}_{i-1} - \mathbf{Q}^*\|_{\infty} = \|\Delta_{i-1}\|_{\infty}, \end{aligned}$$

we can easily see that

$$|\beta_{2,t}| \leq \gamma \sum_{i=1}^t \|\Delta_{i-1}\|_{\infty} \prod_{j=i+1}^t (\mathbf{I} - \mathbf{\Lambda}_j) \mathbf{\Lambda}_i \mathbf{1}. \quad (42)$$

- Next, the term $\beta_{1,t}$ can be controlled by exploiting some sort of statistical independence across different transitions and applying the Bernstein inequality. This is summarized in the following lemma, with the proof deferred to Section C.

Lemma 1: Consider any fixed vector $\mathbf{V}^* \in \mathbb{R}^{|S|}$. There exists some universal constant $c > 0$ such that for any $0 < \delta < 1$, one has for all $1 \leq t \leq T$

$$\left| \gamma \sum_{i=1}^t \prod_{j=i+1}^t (\mathbf{I} - \Lambda_j) \Lambda_i (\mathbf{P}_i - \mathbf{P}) \mathbf{V}^* \right| \leq \tau_1 \|\mathbf{V}^*\|_\infty \mathbf{1} \quad (43)$$

with probability at least $1 - \delta$, provided that $0 < \eta \log \left(\frac{|S||\mathcal{A}|T}{\delta} \right) < 1$. Here, we define

$$\tau_1 := c\gamma \sqrt{\eta \log \left(\frac{|S||\mathcal{A}|T}{\delta} \right)}. \quad (44)$$

- Additionally, we develop an upper bound on the term $\beta_{3,t}$, which follows directly from the concentration of the empirical distribution of the Markov chain (see Lemma 8). The proof is deferred to Section D.

Lemma 2: For any $\delta > 0$, recall the definition of t_{frame} in (37a). Suppose that $T > t_{\text{frame}}$ and $0 < \eta < 1$. Then with probability exceeding $1 - \delta$ one has

$$\left| \prod_{j=1}^t (\mathbf{I} - \Lambda_j) \Delta_0 \right| \leq (1 - \eta)^{\frac{1}{2} t \mu_{\min}} \|\Delta_0\|_\infty \mathbf{1} \quad (45)$$

uniformly over all t obeying $T \geq t \geq t_{\text{frame}}$ and all vector $\Delta_0 \in \mathbb{R}^{|S||\mathcal{A}|}$.

Moreover, in the case where $t < t_{\text{frame}}$, we make note of the straightforward bound

$$\left| \prod_{j=1}^t (\mathbf{I} - \Lambda_j) \Delta_0 \right| \leq \|\Delta_0\|_\infty \mathbf{1}, \quad (46)$$

given that $\mathbf{I} - \Lambda_j$ is a diagonal non-negative matrix whose entries are bounded by $1 - \eta < 1$.

Substituting the preceding bounds into (41), we arrive at

$$|\Delta_t| \leq \begin{cases} \gamma \sum_{i=1}^t \|\Delta_{i-1}\|_\infty \prod_{j=i+1}^t (\mathbf{I} - \Lambda_j) \Lambda_i \mathbf{1} \\ \quad + \tau_1 \|\mathbf{V}^*\|_\infty \mathbf{1} + \|\Delta_0\|_\infty \mathbf{1}, & t < t_{\text{frame}} \\ \gamma \sum_{i=1}^t \|\Delta_{i-1}\|_\infty \prod_{j=i+1}^t (\mathbf{I} - \Lambda_j) \Lambda_i \mathbf{1} \\ \quad + \tau_1 \|\mathbf{V}^*\|_\infty \mathbf{1} + (1 - \eta)^{\frac{1}{2} t \mu_{\min}} \|\Delta_0\|_\infty \mathbf{1}, & t_{\text{frame}} \leq t \leq T \end{cases} \quad (47)$$

with probability at least $1 - 2\delta$, where t_{frame} is defined in (37a). The rest of the proof is thus dedicated to bounding $|\Delta_t|$ based on the above recursive formula (47).

2) *Recursive Analysis:* We shall start by presenting a crude bound, followed by more refined analysis.

a) *A crude bound:* We start by observing the following recursive relation

$$|\Delta_t| \leq \gamma \sum_{i=1}^t \|\Delta_{i-1}\|_\infty \prod_{j=i+1}^t (\mathbf{I} - \Lambda_j) \Lambda_i \mathbf{1} \\ + \tau_1 \|\mathbf{V}^*\|_\infty \mathbf{1} + \|\Delta_0\|_\infty \mathbf{1}, \quad 1 \leq t \leq T, \quad (48)$$

which is a direct consequence of (47). In the sequel, we invoke mathematical induction to establish, for all $1 \leq t \leq T$, the following crude upper bound

$$\|\Delta_t\|_\infty \leq \frac{\tau_1 \|\mathbf{V}^*\|_\infty + \|\Delta_0\|_\infty}{1 - \gamma}, \quad (49)$$

which implies the stability of the asynchronous Q-learning updates.

Towards this, we first observe that (49) holds trivially for the base case (namely, $t = 0$). Now suppose that the inequality (49) holds for all iterations up to $t - 1$. In view of (48) and the induction hypotheses,

$$|\Delta_t| \leq \frac{\gamma(\tau_1 \|\mathbf{V}^*\|_\infty + \|\Delta_0\|_\infty)}{1 - \gamma} \sum_{i=1}^t \prod_{j=i+1}^t (\mathbf{I} - \Lambda_j) \Lambda_i \mathbf{1} \\ + \tau_1 \|\mathbf{V}^*\|_\infty \mathbf{1} + \|\Delta_0\|_\infty \mathbf{1}, \quad (50)$$

where we invoke the fact that the vector $\prod_{j=i+1}^t (\mathbf{I} - \Lambda_j) \Lambda_i \mathbf{1}$ is non-negative. Next, define the diagonal matrix $\mathbf{M}_i := \prod_{j=i+1}^t (\mathbf{I} - \Lambda_j) \Lambda_i$, and denote by $N_i^j(s, a)$ the number of visits to the state-action pair (s, a) between the i -th and the j -th iterations (including i and j). Then the diagonal entries of \mathbf{M}_i satisfy

$$\mathbf{M}_i((s, a), (s, a)) \\ = \begin{cases} \eta(1 - \eta)^{N_{i+1}^t(s, a)}, & \text{if } (s, a) = (s_{i-1}, a_{i-1}), \\ 0, & \text{if } (s, a) \neq (s_{i-1}, a_{i-1}). \end{cases}$$

Letting $\mathbf{e}_{(s, a)} \in \mathbb{R}^{|S||\mathcal{A}|}$ be a standard basis vector whose only nonzero entry is the (s, a) -th entry, we can easily verify that

$$\prod_{j=i+1}^t (\mathbf{I} - \Lambda_j) \Lambda_i \mathbf{1} = \mathbf{M}_i \mathbf{1} = \mathbf{M}_i \mathbf{e}_{(s_{i-1}, a_{i-1})} \\ = \eta(1 - \eta)^{N_{i+1}^t(s_{i-1}, a_{i-1})} \mathbf{e}_{(s_{i-1}, a_{i-1})} \quad (51a)$$

and

$$\sum_{i=1}^t \prod_{j=i+1}^t (\mathbf{I} - \Lambda_j) \Lambda_i \mathbf{1} \\ = \sum_{i=1}^t \eta(1 - \eta)^{N_{i+1}^t(s_{i-1}, a_{i-1})} \mathbf{e}_{(s_{i-1}, a_{i-1})} \\ = \sum_{(s, a) \in \mathcal{S} \times \mathcal{A}} \mathbf{e}_{(s, a)} \\ \cdot \left\{ \sum_{i=1}^t \eta(1 - \eta)^{N_{i+1}^t(s, a)} \mathbb{1} \{ (s_{i-1}, a_{i-1}) = (s, a) \} \right\} \\ \leq \sum_{(s, a) \in \mathcal{S} \times \mathcal{A}} \sum_{j=0}^{\infty} \eta(1 - \eta)^j \mathbf{e}_{(s, a)} = \sum_{j=0}^{\infty} \eta(1 - \eta)^j \mathbf{1} = \mathbf{1}. \quad (51b)$$

Combining the above relations with the inequality (50), one deduces that

$$\begin{aligned}\|\Delta_t\|_\infty &\leq \frac{\gamma(\tau_1\|\mathbf{V}^*\|_\infty + \|\Delta_0\|_\infty)}{1-\gamma} + \tau_1\|\mathbf{V}^*\|_\infty + \|\Delta_0\|_\infty \\ &= \frac{\tau_1\|\mathbf{V}^*\|_\infty + \|\Delta_0\|_\infty}{1-\gamma},\end{aligned}$$

thus establishing (49) for the t -th iteration. This induction analysis thus validates (49) for all $1 \leq t \leq T$.

b) Refined analysis: Now, we strengthen the bound (49) by means of a recursive argument. To begin with, it is easily seen that the term $(1-\eta)^{\frac{1}{2}t\mu_{\min}}\|\Delta_0\|_\infty$ is bounded above by $(1-\gamma)\varepsilon$ for any $t > t_{\text{th}}$, where we remind the reader of the definition of t_{th} in (37b) and the fact that $\|\Delta_0\|_\infty = \|\mathbf{Q}^*\|_\infty \leq \frac{1}{1-\gamma}$. It is assumed that $T > t_{\text{th}}$. To facilitate our argument, we introduce a collection of auxiliary quantities u_t as follows

$$u_0 = \frac{\|\Delta_0\|_\infty}{1-\gamma}, \quad (52a)$$

$$u_t = \|\mathbf{v}_t\|_\infty, \quad \text{with}$$

$$\mathbf{v}_t = \begin{cases} \gamma \sum_{i=1}^t \prod_{j=i+1}^t (\mathbf{I} - \Lambda_j) \Lambda_i \mathbf{1} u_{i-1} + \|\Delta_0\|_\infty \mathbf{1}, & \text{for } 1 \leq t \leq t_{\text{th}}, \\ \gamma \sum_{i=1}^t \prod_{j=i+1}^t (\mathbf{I} - \Lambda_j) \Lambda_i \mathbf{1} u_{i-1}, & \text{for } t > t_{\text{th}}. \end{cases} \quad (52b)$$

These auxiliary quantities are useful as they provide upper bounds on $\|\Delta_t\|_\infty$, as asserted by the following lemma. The proof is deferred to Section E.

Lemma 3: Recall the definition (44) of τ_1 in Lemma 1. With probability at least $1 - 2\delta$, the quantities $\{u_t\}$ defined in (52) satisfy

$$\|\Delta_t\|_\infty \leq \frac{\tau_1\|\mathbf{V}^*\|_\infty}{1-\gamma} + u_t + \varepsilon. \quad (53)$$

The preceding result motivates us to turn attention to bounding the quantities $\{u_t\}$. Towards this end, we resort to a frame-based analysis by dividing the iterations $[1, t]$ into contiguous frames each comprising t_{frame} (cf. (37a)) iterations. Further, we define another auxiliary sequence:

$$w_k := (1-\rho)^k \frac{\|\Delta_0\|_\infty}{1-\gamma} = (1-\rho)^k \frac{\|\mathbf{Q}_0 - \mathbf{Q}^*\|_\infty}{1-\gamma}, \quad (54)$$

where we remind the reader of the definition of ρ in (37d). The connection between $\{w_k\}$ and $\{u_t\}$ is made precise as follows, whose proof is postponed to Section F.

Lemma 4: For any $\delta \in (0, \frac{1}{2})$, with probability at least $1 - 2\delta$, one has

$$u_t \leq w_k, \quad \text{with } k = \max \left\{ 0, \left\lfloor \frac{t - t_{\text{th}}}{t_{\text{frame}}} \right\rfloor \right\}. \quad (55)$$

Combining Lemmas 3-4, we arrive at

$$\begin{aligned}\|\mathbf{Q}_t - \mathbf{Q}^*\|_\infty &= \|\Delta_t\|_\infty \leq \frac{\tau_1\|\mathbf{V}^*\|_\infty}{1-\gamma} + w_k + \varepsilon \\ &\leq \frac{(1-\rho)^k \|\mathbf{Q}_0 - \mathbf{Q}^*\|_\infty}{1-\gamma} + \frac{\tau_1\|\mathbf{V}^*\|_\infty}{1-\gamma} + \varepsilon,\end{aligned}$$

which finishes the proof of Theorem 5.

C. Proof of Theorem 1

Now we return to complete the proof of Theorem 1. To control $\|\Delta_t\|_\infty$ to the desired level, we first claim that the first term of (38) obeys

$$(1-\rho)^k \frac{\|\Delta_0\|_\infty}{1-\gamma} \leq \varepsilon \quad (56)$$

whenever

$$t \geq t_{\text{th}} + t_{\text{frame}} + \frac{4}{(1-\gamma)\eta\mu_{\min}} \log \left(\frac{\|\Delta_0\|_\infty}{\varepsilon(1-\gamma)} \right), \quad (57)$$

provided that $\eta < 1/\mu_{\text{frame}}$. Furthermore, by taking the learning rate as

$$\eta = \min \left\{ \frac{(1-\gamma)^4 \varepsilon^2}{c^2 \gamma^2 \log \frac{|S||\mathcal{A}|T}{\delta}}, \frac{1}{\mu_{\text{frame}}} \right\}, \quad (58)$$

one can easily verify that the second term of (38) satisfies

$$\frac{c\gamma}{1-\gamma} \|\mathbf{V}^*\|_\infty \sqrt{\eta \log \left(\frac{|S||\mathcal{A}|T}{\delta} \right)} \leq \varepsilon, \quad (59)$$

where the last step follows since $\|\mathbf{V}^*\|_\infty \leq \frac{1}{1-\gamma}$. Putting the above bounds together ensures $\|\Delta_t\|_\infty \leq 3\varepsilon$. By replacing ε with $\varepsilon/3$, we can readily conclude the proof, as long as the claim (56) can be justified.

Proof of the Inequality (56): Observe that

$$(1-\rho)^k \frac{\|\Delta_0\|_\infty}{1-\gamma} \leq \exp(-\rho k) \frac{\|\Delta_0\|_\infty}{1-\gamma} \leq \varepsilon$$

holds true whenever $k \geq \frac{\log \left(\frac{\|\Delta_0\|_\infty}{\varepsilon(1-\gamma)} \right)}{\rho}$, which would hold as long as (according to the definition (55) of k)

$$t \geq t_{\text{th}} + t_{\text{frame}} + \frac{t_{\text{frame}}}{\rho} \log \left(\frac{\|\Delta_0\|_\infty}{\varepsilon(1-\gamma)} \right). \quad (60)$$

In addition, if $\eta < 1/\mu_{\text{frame}}$, then one has $(1-\eta)^{\mu_{\text{frame}}} \leq 1 - \eta\mu_{\text{frame}}/2$, thus guaranteeing that

$$\begin{aligned}\rho &= (1-\gamma) \left(1 - (1-\eta)^{\mu_{\text{frame}}} \right) \geq (1-\gamma) \left(1 - 1 + \frac{\eta\mu_{\text{frame}}}{2} \right) \\ &= \frac{1}{2} (1-\gamma) \eta \mu_{\text{frame}}.\end{aligned}$$

This taken collectively with (60) demonstrates that $(1-\rho)^k \frac{\|\Delta_0\|_\infty}{1-\gamma} \leq \varepsilon$ holds as long as

$$\begin{aligned}t &\geq t_{\text{th}} + t_{\text{frame}} + \frac{2t_{\text{frame}}}{(1-\gamma)\eta\mu_{\text{frame}}} \log \left(\frac{\|\Delta_0\|_\infty}{\varepsilon(1-\gamma)} \right) \\ &= t_{\text{th}} + t_{\text{frame}} + \frac{4}{(1-\gamma)\eta\mu_{\min}} \log \left(\frac{\|\Delta_0\|_\infty}{\varepsilon(1-\gamma)} \right),\end{aligned} \quad (61)$$

where we have made use of the definition of μ_{frame} (cf. (37c)). \square

VII. COVER-TIME-BASED ANALYSIS OF ASYNCHRONOUS Q-LEARNING

In this section, we sketch the proof of Theorem 2. Before continuing, we recall the definition of t_{cover} in (9), and further introduce a quantity

$$t_{\text{cover,all}} := t_{\text{cover}} \log \frac{T}{\delta}. \quad (62)$$

There are two useful facts regarding $t_{\text{cover,all}}$ that play an important role in the analysis.

Lemma 5: Define the event

$$\mathcal{K}_l := \left\{ \exists (s, a) \in \mathcal{S} \times \mathcal{A} \text{ s.t. it is not visited within iterations } (lt_{\text{cover,all}}, (l+1)t_{\text{cover,all}}] \right\},$$

and set $L := \lfloor \frac{T}{t_{\text{cover,all}}} \rfloor$. Then one has $\mathbb{P} \left\{ \bigcup_{l=0}^L \mathcal{K}_l \right\} \leq \delta$.

Proof: See Section H. \square

In other words, Lemma 5 tells us that with high probability, all state-action pairs are visited at least once in every time frame $(lt_{\text{cover,all}}, (l+1)t_{\text{cover,all}}]$ with $0 \leq l \leq \lfloor T/t_{\text{cover,all}} \rfloor$. The next result is a consequence of Lemma 5 as well as the analysis of Lemma 2; the proof can be found in Section D.

Lemma 6: For any $\delta > 0$, recall the definition of $t_{\text{cover,all}}$ in (62). Suppose that $T > t_{\text{cover,all}}$ and $0 < \eta < 1$. Then with probability exceeding $1 - \delta$ one has

$$\left| \prod_{j=1}^t (I - \Lambda_j) \Delta_0 \right| \leq (1 - \eta)^{\frac{t}{2t_{\text{cover,all}}}} \|\Delta_0\|_\infty \mathbf{1} \quad (63)$$

uniformly over all t obeying $T \geq t \geq t_{\text{cover,all}}$ and all vector $\Delta_0 \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$.

With the above two lemmas in mind, we are now positioned to prove Theorem 2. Repeating the analysis of (47) (except that Lemma 2 is replaced by Lemma 6) yields

$$|\Delta_t| \leq \begin{cases} \gamma \sum_{i=1}^t \|\Delta_{i-1}\|_\infty \prod_{j=i+1}^t (I - \Lambda_j) \Lambda_i \mathbf{1} \\ + \tau_1 \|\mathbf{V}^*\|_\infty \mathbf{1} + \|\Delta_0\|_\infty \mathbf{1}, & t < t_{\text{cover,all}} \\ \gamma \sum_{i=1}^t \|\Delta_{i-1}\|_\infty \prod_{j=i+1}^t (I - \Lambda_j) \Lambda_i \mathbf{1} \\ + \tau_1 \|\mathbf{V}^*\|_\infty \mathbf{1} + (1 - \eta)^{\frac{t}{2t_{\text{cover,all}}}} \|\Delta_0\|_\infty \mathbf{1}, & t_{\text{cover,all}} \leq t \leq T \end{cases}$$

with probability at least $1 - 2\delta$. This observation resembles (47), except that t_{frame} (resp. μ_{\min}) is replaced by $t_{\text{cover,all}}$ (resp. $\frac{1}{t_{\text{cover,all}}}$). As a consequence, we can immediately use the recursive analysis carried out in Section VI-B.2 to establish a convergence guarantee based on the cover time. More specifically, define

$$\tilde{\rho} := (1 - \gamma) \left(1 - (1 - \eta)^{\frac{t_{\text{cover,all}}}{2t_{\text{cover,all}}}} \right) = (1 - \gamma) \left(1 - (1 - \eta)^{\frac{1}{2}} \right). \quad (64)$$

Replacing ρ by $\tilde{\rho}$ in Theorem 5 reveals that with probability at least $1 - 6\delta$,

$$\begin{aligned} \|\mathbf{Q}_t - \mathbf{Q}^*\|_\infty &\leq (1 - \tilde{\rho})^k \frac{\|\mathbf{Q}_0 - \mathbf{Q}^*\|_\infty}{1 - \gamma} \\ &\quad + \frac{c\gamma}{1 - \gamma} \|\mathbf{V}^*\|_\infty \sqrt{\eta \log \left(\frac{|\mathcal{S}||\mathcal{A}|T}{\delta} \right)} + \varepsilon \end{aligned} \quad (65)$$

holds for all $t \leq T$, where $k := \max \left\{ 0, \left\lfloor \frac{t - t_{\text{th,cover}}}{t_{\text{cover,all}}} \right\rfloor \right\}$ and we abuse notation to define

$$t_{\text{th,cover}} := 2t_{\text{cover,all}} \log \frac{1}{(1 - \gamma)^2 \varepsilon}.$$

Repeating the proof of the inequality (56) yields

$$(1 - \tilde{\rho})^k \frac{\|\Delta_0\|_\infty}{1 - \gamma} \leq \varepsilon,$$

whenever $t \geq t_{\text{th,cover}} + t_{\text{cover,all}} + \frac{2t_{\text{cover,all}}}{(1 - \gamma)\eta} \log \left(\frac{1}{\varepsilon(1 - \gamma)^2} \right)$, with the proviso that $\eta < 1/2$. In addition, setting $\eta = \frac{(1 - \gamma)^4}{c^2 \gamma^2 \varepsilon^2 \log \left(\frac{|\mathcal{S}||\mathcal{A}|T}{\delta} \right)}$ guarantees that

$$\begin{aligned} \frac{c\gamma}{1 - \gamma} \|\mathbf{V}^*\|_\infty \sqrt{\eta \log \left(\frac{|\mathcal{S}||\mathcal{A}|T}{\delta} \right)} \\ \leq \frac{c\gamma}{(1 - \gamma)^2} \sqrt{\eta \log \left(\frac{|\mathcal{S}||\mathcal{A}|T}{\delta} \right)} \leq \varepsilon. \end{aligned}$$

In conclusion, we have $\|\mathbf{Q}_t - \mathbf{Q}^*\|_\infty \leq 3\varepsilon$ as long as

$$t \geq \frac{c't_{\text{cover,all}}}{(1 - \gamma)^5 \varepsilon^2} \log \left(\frac{|\mathcal{S}||\mathcal{A}|T}{\delta} \right) \log \left(\frac{1}{\varepsilon(1 - \gamma)^2} \right),$$

for some sufficiently large constant $c' > 0$. This together with the definition (62) completes the proof.

VIII. ANALYSIS UNDER ADAPTIVE LEARNING RATES (PROOF OF THEOREM 3)

A. Useful Preliminary Facts About η_t

To begin with, we make note of several useful properties about η_t .

- Invoking the concentration result in Lemma 8, one can easily show that with probability at least $1 - \delta$,

$$\frac{1}{2} \mu_{\min} < \min_{s,a} \frac{K_t(s, a)}{t} < \frac{3}{2} \mu_{\min} \quad (66)$$

holds simultaneously for all t obeying $T \geq t \geq \frac{443t_{\text{mix}} \log \left(\frac{4|\mathcal{S}||\mathcal{A}|t}{\delta} \right)}{\mu_{\min}}$. In addition, this concentration result taken collectively with the update rule (22) of $\hat{\mu}_{\min,t}$ — in particular, the second case of (22) — implies that $\hat{\mu}_{\min,t}$ “stabilizes” as t grows; to be precise, there exists some quantity $c' \in [1/6, 9/2]$ such that

$$\hat{\mu}_{\min,t} \equiv c' \mu_{\min} \quad (67)$$

holds simultaneously for all t obeying $T \geq t \geq \frac{443t_{\text{mix}} \log \left(\frac{4|\mathcal{S}||\mathcal{A}|t}{\delta} \right)}{\mu_{\min}}$.

- For any t obeying $t \geq \frac{6c_\eta t_{\text{mix}} \log \left(\frac{2|\mathcal{S}||\mathcal{A}|t}{\delta} \right)}{\mu_{\min}(1 - \gamma)^2}$ (so that $\frac{\log t}{\hat{\mu}_{\min,t}(1 - \gamma)^2} \leq \frac{1}{c_\eta}$ and $t \geq \frac{443t_{\text{mix}} \log \left(\frac{2|\mathcal{S}||\mathcal{A}|t}{\delta} \right)}{\mu_{\min}}$ for $c_\eta \geq 11$), the learning rate (23) simplifies to

$$\eta_t = c_\eta \exp \left(\left\lfloor \log \frac{\log t}{c' \mu_{\min}(1 - \gamma)^2 t} \right\rfloor \right). \quad (68)$$

Clearly, there exists a sequence of endpoints $t_1 < t_2 < t_3 < \dots$ with $t_1 \leq \frac{6ec_\eta t_{\text{mix}} \log \left(\frac{2|\mathcal{S}||\mathcal{A}|t_1}{\delta} \right)}{\mu_{\min}(1 - \gamma)^2}$ such that:

$$2t_k < t_{k+1} < 3t_k \quad \text{and} \quad (69)$$

$$\eta_t = \eta_{(k)} := \frac{\alpha_k \log t_{k+1}}{\mu_{\min}(1 - \gamma)^2 t_{k+1}}, \quad \forall t_k < t \leq t_{k+1} \quad (70)$$

for some positive constant $\alpha_k \in [\frac{2c_\eta}{9e}, 6c_\eta]$; in words, (70) provides a concrete expression/bound for the piecewise constant learning rate, where the t_k 's form the change points.

Combining (70) with the definition of \hat{Q}_t (cf. (22)), one can easily check that for $t > t_1$,

$$\hat{Q}_t = Q_{t_k}, \quad \forall t_k < t \leq t_{k+1}, \quad (71)$$

meaning that \hat{Q}_t remains fixed within each time segment $(t_k, t_{k+1}]$. With this property in mind, we only need to analyze Q_{t_k} in the sequel, which can be easily accomplished by invoking Theorem 1.

B. A Crude Bound

Given that $0 < \eta_t \leq 1$ and $0 \leq r(s, a) \leq 1$, the update rule (10) of Q_t implies that

$$\begin{aligned} & \|Q_t\|_\infty \\ & \leq \max \left\{ (1 - \eta_t) \|Q_{t-1}\|_\infty + \eta_t (1 + \gamma \|Q_{t-1}\|_\infty), \|Q_{t-1}\|_\infty \right\} \\ & \leq \|Q_{t-1}\|_\infty + \gamma, \end{aligned}$$

thus leading to the following crude bound that for any $t > \frac{1}{1-\gamma}$

$$\|Q_t - Q^*\|_\infty \leq t + \|Q_0\|_\infty + \|Q^*\|_\infty \leq t + \frac{2}{1-\gamma} \leq 3t. \quad (72)$$

Remark 5: As we shall see momentarily, this crude bound allows one to control — in a coarse manner — the error at the beginning of each time interval $[t_{k-1}, t_k]$, which is needed when invoking Theorem 1.

C. Refined Analysis

Let us define

$$\varepsilon_k := \sqrt{\frac{c_{k,0} \log\left(\frac{|S||A|t_k}{\delta}\right) \log t_k}{\mu_{\min}(1-\gamma)^5 \gamma^2 t_k}}, \quad (73)$$

where the constant $c_{k,0}$ is chosen to be $c_{k,0} = \alpha_{k-1}/c_1 > 0$, with $c_1 > 0$ the universal constant stated in Theorem 1. The property (70) of η_t together with the definition (73) implies that

$$\eta_t = \frac{c_1(1-\gamma)^4 \varepsilon_k^2}{\log\left(\frac{|S||A|t_k}{\delta}\right)} = \frac{c_1}{\log\left(\frac{|S||A|t_k}{\delta}\right)} \min \left\{ (1-\gamma)^4 \varepsilon_k^2, \frac{1}{t_{\text{mix}}} \right\}$$

for any $t \in (t_{k-1}, t_k]$, as long as $(1-\gamma)^4 \varepsilon_k^2 \leq 1/t_{\text{mix}}$, or more explicitly, when

$$t_k \geq \frac{c_{k,0} t_{\text{mix}} \log\left(\frac{|S||A|t_k}{\delta}\right) \log t_k}{\mu_{\min}(1-\gamma) \gamma^2}. \quad (74)$$

In addition, the condition (69) and the definition (73) further tell us that

$$t_k - t_{k-1} > t_{k-1} > \frac{1}{3} t_k = \frac{c_{k,0} \log\left(\frac{|S||A|t_k}{\delta}\right) \log t_k}{3\mu_{\min}(1-\gamma)^5 \gamma^2 \varepsilon_k^2}.$$

Invoking Theorem 1 with an initialization $Q_{t_{k-1}}$ (which clearly satisfies the crude bound (72)) ensures that

$$\|Q_{t_k} - Q^*\|_\infty \leq \varepsilon_k \quad (75)$$

with probability at least $1 - \delta$, with the proviso that

$$\begin{aligned} \frac{1}{3} t_k & \geq \frac{c_0}{\mu_{\min}} \left\{ \frac{1}{(1-\gamma)^5 \varepsilon_k^2} + \frac{t_{\text{mix}}}{1-\gamma} \right\} \\ & \cdot \log\left(\frac{|S||A|t_k}{\delta}\right) \log\left(\frac{t_k}{(1-\gamma)^2 \varepsilon_k}\right) \end{aligned} \quad (76)$$

with $c_0 > 0$ the universal constant stated in Theorem 1. Under the sample size condition (74), this requirement (76) can be guaranteed by adjusting the constant c_η in (23) to satisfy the following inequality:

$$c_{k,0} = \frac{\alpha_{k-1}}{c_1} \geq \frac{2c_\eta}{9ec_1} > 6c_0.$$

Finally, taking $t_{k_{\max}}$ to be the largest change point that does not exceed T , we see from (69) that $\frac{1}{3}T \leq t_{k_{\max}} \leq T$. Then one has

$$\begin{aligned} \|Q_T - Q^*\|_\infty & = \|Q_{t_{k_{\max}}} - Q^*\|_\infty \\ & \leq \varepsilon_{k_{\max}} = \sqrt{\frac{c_{k,0} \log\left(\frac{|S||A|t_{k_{\max}}}{\delta}\right) \log t_{k_{\max}}}{\mu_{\min}(1-\gamma)^5 \gamma^2 t_{k_{\max}}}} \\ & \leq \sqrt{\frac{3c_{k,0} \log\left(\frac{|S||A|T}{\delta}\right) \log T}{\mu_{\min}(1-\gamma)^5 \gamma^2 T}} \end{aligned} \quad (77)$$

These immediately conclude the proof of the theorem under the sample size condition (26), provided that

$$C > \frac{18c_\eta}{c_1} > \frac{3\alpha_{k-1}}{c_1} = 3c_{k,0}.$$

IX. ANALYSIS OF ASYNCHRONOUS VARIANCE-REDUCED Q-LEARNING

This section aims to establish Theorem 4. We carry out an epoch-based analysis, that is, we first quantify the progress made over each epoch, and then demonstrate how many epochs are sufficient to attain the desired accuracy. In what follows, we shall overload the notation by defining

$$t_{\text{frame}} := \frac{443t_{\text{mix}}}{\mu_{\min}} \log\left(\frac{4|S||A|t_{\text{epoch}}}{\delta}\right), \quad (78a)$$

$$t_{\text{th}} := \max \left\{ \frac{2 \log\left(\frac{1}{(1-\gamma)^2 \varepsilon}\right)}{\eta \mu_{\min}}, t_{\text{frame}} \right\}, \quad (78b)$$

$$\rho := (1-\gamma)(1 - (1-\eta)^{\mu_{\text{frame}}}), \quad (78c)$$

$$\mu_{\text{frame}} := \frac{1}{2} \mu_{\min} t_{\text{frame}}. \quad (78d)$$

A. Per-Epoch Analysis

We start by analyzing the progress made over each epoch. Before proceeding, we denote by $\hat{P} \in [0, 1]^{|S||A| \times |S|}$ a matrix corresponding to the empirical probability transition kernel used in (28) from N new sample transitions. Further, we use the vector $\bar{Q} \in \mathbb{R}^{|S||A|}$ to represent the reference Q-function, and introduce the vector $\bar{V} \in \mathbb{R}^{|S|}$ to represent the corresponding value function so that $\bar{V}(s) := \max_a \bar{Q}(s, a)$ for all $s \in S$.

For convenience, this subsection abuses notation to assume that an epoch starts with an estimate $\mathbf{Q}_0 = \bar{\mathbf{Q}}$, and consists of the subsequent

$$t_{\text{epoch}} := t_{\text{frame}} + t_{\text{th}} + \frac{8 \log \frac{2}{1-\gamma}}{(1-\gamma)\eta\mu_{\min}} \quad (79)$$

iterations of variance-reduced Q-learning updates, where t_{frame} and t_{th} are defined in (78a) and (78b), respectively. In the sequel, we divide all epochs into two phases, depending on the quality of the initial estimate $\bar{\mathbf{Q}}$ in each epoch.

1) *Phase 1: When $\|\bar{\mathbf{Q}} - \mathbf{Q}^*\|_\infty > 1/\sqrt{1-\gamma}$:* Recalling the matrix notation of $\mathbf{\Lambda}_t$ and \mathbf{P}_t in (34) and (35), respectively, we can rewrite (27) as follows

$$\mathbf{Q}_t = (\mathbf{I} - \mathbf{\Lambda}_t)\mathbf{Q}_{t-1} + \mathbf{\Lambda}_t \left(\mathbf{r} + \gamma \mathbf{P}_t(\mathbf{V}_{t-1} - \bar{\mathbf{V}}) + \gamma \tilde{\mathbf{P}}\bar{\mathbf{V}} \right). \quad (80)$$

Following similar steps as in the expression (39), we arrive at the following error decomposition

$$\begin{aligned} \mathbf{\Theta}_t &:= \mathbf{Q}_t - \mathbf{Q}^* = (\mathbf{I} - \mathbf{\Lambda}_t)\mathbf{Q}_{t-1} \\ &\quad + \mathbf{\Lambda}_t \left(\mathbf{r} + \gamma \mathbf{P}_t(\mathbf{V}_{t-1} - \bar{\mathbf{V}}) + \gamma \tilde{\mathbf{P}}\bar{\mathbf{V}} \right) - \mathbf{Q}^* \\ &= (\mathbf{I} - \mathbf{\Lambda}_t)(\mathbf{Q}_{t-1} - \mathbf{Q}^*) \\ &\quad + \mathbf{\Lambda}_t \left(\mathbf{r} + \gamma \mathbf{P}_t(\mathbf{V}_{t-1} - \bar{\mathbf{V}}) + \gamma \tilde{\mathbf{P}}\bar{\mathbf{V}} - \mathbf{Q}^* \right) \\ &= (\mathbf{I} - \mathbf{\Lambda}_t)(\mathbf{Q}_{t-1} - \mathbf{Q}^*) \\ &\quad + \gamma \mathbf{\Lambda}_t \left(\mathbf{P}_t(\mathbf{V}_{t-1} - \bar{\mathbf{V}}) + \tilde{\mathbf{P}}\bar{\mathbf{V}} - \mathbf{P}\mathbf{V}^* \right) \\ &= (\mathbf{I} - \mathbf{\Lambda}_t)\mathbf{\Theta}_{t-1} + \gamma \mathbf{\Lambda}_t(\tilde{\mathbf{P}} - \mathbf{P})\bar{\mathbf{V}} \\ &\quad + \gamma \mathbf{\Lambda}_t(\mathbf{P}_t - \mathbf{P})(\mathbf{V}^* - \bar{\mathbf{V}}) + \gamma \mathbf{\Lambda}_t \mathbf{P}_t(\mathbf{V}_{t-1} - \mathbf{V}^*), \end{aligned} \quad (81)$$

which once again leads to a recursive relation

$$\begin{aligned} \mathbf{\Theta}_t &= \gamma \sum_{i=1}^t \prod_{j=i+1}^t (\mathbf{I} - \mathbf{\Lambda}_j) \underbrace{\mathbf{\Lambda}_i(\tilde{\mathbf{P}} - \mathbf{P})\bar{\mathbf{V}}}_{=: \mathbf{h}_{0,t}} \\ &\quad + \gamma \sum_{i=1}^t \prod_{j=i+1}^t (\mathbf{I} - \mathbf{\Lambda}_j) \underbrace{\mathbf{\Lambda}_i(\mathbf{P}_i - \mathbf{P})(\mathbf{V}^* - \bar{\mathbf{V}})}_{=: \mathbf{h}_{1,t}} \\ &\quad + \gamma \sum_{i=1}^t \prod_{j=i+1}^t (\mathbf{I} - \mathbf{\Lambda}_j) \underbrace{\mathbf{\Lambda}_i \mathbf{P}_i(\mathbf{V}_{i-1} - \mathbf{V}^*)}_{=: \mathbf{h}_{2,t}} \\ &\quad + \underbrace{\prod_{j=1}^t (\mathbf{I} - \mathbf{\Lambda}_j) \mathbf{\Theta}_0}_{=: \mathbf{h}_{3,t}}. \end{aligned} \quad (82)$$

This identity takes a very similar form as (40) except for the additional term $\mathbf{h}_{0,t}$.

Let us begin by controlling the first term, towards which we have the following lemma. The proof is postponed to Section G.

Lemma 7: Suppose that $\tilde{\mathbf{P}}$ is constructed using N consecutive sample transitions. If $N > t_{\text{frame}}$, then with probability greater than $1 - \delta$, one has

$$\begin{aligned} \|\mathbf{h}_{0,t}\|_\infty &\leq \gamma \sqrt{\frac{4 \log \left(\frac{6N|\mathcal{S}||\mathcal{A}|}{\delta} \right)}{N\mu_{\min}}} \|\bar{\mathbf{V}} - \mathbf{V}^*\|_\infty \\ &\quad + \frac{\gamma}{1-\gamma} \sqrt{\frac{4 \log \left(\frac{6N|\mathcal{S}||\mathcal{A}|}{\delta} \right)}{N\mu_{\min}}}. \end{aligned} \quad (83)$$

If $t < t_{\text{frame}}$, then it is straightforwardly seen that

$$\|\mathbf{h}_{3,t}\|_\infty \leq \|\mathbf{\Theta}_0\|_\infty \mathbf{1}.$$

Taking this together with the results from Lemma 1 and Lemma 2, we are guaranteed that

$$\begin{aligned} \|\mathbf{h}_{1,t}\|_\infty &\leq \tau_2 \|\mathbf{V}^* - \bar{\mathbf{V}}\|_\infty \mathbf{1} \\ \|\mathbf{h}_{3,t}\|_\infty &\leq \begin{cases} (1-\gamma)^{\frac{1}{2}t\mu_{\min}} \|\mathbf{\Theta}_0\|_\infty \mathbf{1}, & \text{if } t_{\text{frame}} \leq t \leq t_{\text{epoch}} \\ \|\mathbf{\Theta}_0\|_\infty \mathbf{1}, & \text{if } t < t_{\text{frame}} \end{cases} \end{aligned}$$

with probability at least $1 - 2\delta$, where

$$\tau_2 := c' \gamma \sqrt{\eta \log \left(\frac{|\mathcal{S}||\mathcal{A}|t_{\text{epoch}}}{\delta} \right)}$$

for some constant $c' > 0$ (similar to (44)). In addition, the term $\mathbf{h}_{2,t}$ can be bounded in the same way as $\beta_{2,t}$ in (42). Therefore, repeating the same argument as for Theorem 5 and taking $\xi = \frac{1}{16\sqrt{1-\gamma}}$, we conclude that with probability at least $1 - \delta$,

$$\begin{aligned} \|\mathbf{\Theta}_t\|_\infty &\leq (1-\rho)^k \frac{\|\mathbf{\Theta}_0\|_\infty}{1-\gamma} + \tilde{\tau} + \xi \\ &= (1-\rho)^k \frac{\|\bar{\mathbf{Q}} - \mathbf{Q}^*\|_\infty}{1-\gamma} + \tilde{\tau} + \xi \end{aligned} \quad (84)$$

holds simultaneously for all $0 < t \leq t_{\text{epoch}}$, where $k = \max\{0, \lfloor \frac{t-t_{\text{th},\xi}}{t_{\text{frame}}} \rfloor\}$, and

$$\begin{aligned} \tilde{\tau} &:= \frac{c\gamma}{1-\gamma} \left\{ \sqrt{\frac{\log \frac{N|\mathcal{S}||\mathcal{A}|}{\delta}}{(1-\gamma)^2 N\mu_{\min}}} + \|\mathbf{V}^* - \bar{\mathbf{V}}\|_\infty \right. \\ &\quad \cdot \left(\sqrt{\eta \log \left(\frac{|\mathcal{S}||\mathcal{A}|t_{\text{epoch}}}{\delta} \right)} + \sqrt{\frac{4 \log \left(\frac{6N|\mathcal{S}||\mathcal{A}|}{\delta} \right)}{N\mu_{\min}}} \right) \Big\}, \\ t_{\text{th},\xi} &:= \max \left\{ \frac{2 \log \frac{1}{(1-\gamma)^2 \xi}}{\eta\mu_{\min}}, t_{\text{frame}} \right\} \end{aligned}$$

for some constant $c > 0$.

Let $C > 0$ be some sufficient large constant. Setting $\eta_t \equiv \eta = \min \left\{ \frac{(1-\gamma)^2}{C\gamma^2 \log \frac{|\mathcal{S}||\mathcal{A}|t_{\text{epoch}}}{\delta}}, \frac{1}{\mu_{\text{frame}}} \right\}$, and ensuring $N \geq \max\{t_{\text{frame}}, C \frac{\log \frac{N|\mathcal{S}||\mathcal{A}|}{\delta}}{(1-\gamma)^3 \mu_{\min}}\}$, we can easily demonstrate that

$$\begin{aligned} \|\mathbf{\Theta}_t\|_\infty &\leq (1-\rho)^k \frac{\|\bar{\mathbf{Q}} - \mathbf{Q}^*\|_\infty}{1-\gamma} \\ &\quad + \frac{1}{8\sqrt{1-\gamma}} + \frac{1}{4} \|\mathbf{V}^* - \bar{\mathbf{V}}\|_\infty. \end{aligned}$$

As a consequence, if $t_{\text{epoch}} \geq t_{\text{frame}} + t_{\text{th},\xi} + \frac{8 \log \frac{2}{1-\gamma}}{(1-\gamma)\eta\mu_{\min}}$, one has

$$(1 - \rho)^k \leq \frac{1}{8}(1 - \gamma),$$

which in turn implies that

$$\begin{aligned} \|\Theta_{t_{\text{epoch}}}\|_{\infty} &\leq \frac{1}{8}\|\bar{Q} - Q^*\|_{\infty} + \frac{1}{8\sqrt{1-\gamma}} + \frac{1}{4}\|V^* - \bar{V}\|_{\infty} \\ &\leq \frac{1}{2} \max \left\{ \frac{1}{\sqrt{1-\gamma}}, \|\bar{Q} - Q^*\|_{\infty} \right\}, \end{aligned} \quad (85)$$

where the last step invokes the simple relation $\|V^* - \bar{V}\|_{\infty} \leq \|\bar{Q} - Q^*\|_{\infty}$. Thus, we conclude that

$$\|Q_{t_{\text{epoch}}} - Q^*\|_{\infty} \leq \frac{1}{2} \max \left\{ \frac{1}{\sqrt{1-\gamma}}, \|\bar{Q} - Q^*\|_{\infty} \right\}. \quad (86)$$

2) *Phase 2: When $\|\bar{Q} - Q^*\|_{\infty} \leq 1/\sqrt{1-\gamma}$:* The analysis of Phase 2 follows by straightforwardly combining the analysis of Phase 1 and that of the synchronous counterpart in [12]. For the sake of brevity, we only sketch the main steps.

Following the proof idea of [12, Section B.2], we introduce an auxiliary vector \hat{Q} which is the unique fix point to the following equation, which can be regarded as a population-level Bellman equation with proper reward perturbation, namely,

$$\hat{Q} = r + \gamma P(\hat{V} - \bar{V}) + \gamma \tilde{P}\bar{V}. \quad (87)$$

Here, as usual, $\hat{V} \in \mathbb{R}^{|\mathcal{S}|}$ represents the value function corresponding to \hat{Q} . This can be viewed as a Bellman equation when the reward vector r is replaced by $\tilde{r} := r + \gamma(\tilde{P} - P)\bar{V}$. Repeating the arguments in the proof of [12, Lemma 4] (except that we need to apply the measure concentration of \tilde{P} in the manner performed in the proof of Lemma 7 due to Markovian data), we reach

$$\|\hat{Q} - Q^*\|_{\infty} \leq c' \sqrt{\frac{\log \frac{N|\mathcal{S}||\mathcal{A}|}{\delta}}{(1-\gamma)^3 N \mu_{\min}}} \leq \varepsilon \quad (88)$$

with probability at least $1 - \delta$ for some constant $c' > 0$, provided that $N \geq (c')^2 \frac{\log \frac{N|\mathcal{S}||\mathcal{A}|}{\delta}}{(1-\gamma)^3 \varepsilon^2}$ and that $\|\bar{Q} - Q^*\|_{\infty} \leq 1/\sqrt{1-\gamma}$. It is worth noting that \hat{Q} only serves as a helper in the proof and is never explicitly constructed in the algorithm, as we don't have access to the probability transition matrix P .

In addition, we claim that

$$\|Q_{t_{\text{epoch}}} - \hat{Q}\|_{\infty} \leq \frac{\|\hat{Q} - Q^*\|_{\infty}}{8} + \frac{\|\bar{Q} - Q^*\|_{\infty}}{8} + \varepsilon. \quad (89)$$

Under this claim, the triangle inequality yields

$$\begin{aligned} \|Q_{t_{\text{epoch}}} - Q^*\|_{\infty} &\leq \|Q_{t_{\text{epoch}}} - \hat{Q}\|_{\infty} + \|\hat{Q} - Q^*\|_{\infty} \\ &\leq \frac{1}{8}\|\bar{Q} - Q^*\|_{\infty} + \frac{9}{8}\|\hat{Q} - Q^*\|_{\infty} + \varepsilon \\ &\leq \frac{1}{8}\|\bar{Q} - Q^*\|_{\infty} + \frac{17}{8}\varepsilon, \end{aligned} \quad (90)$$

where the last inequality follows from (88).

a) *Proof of the inequality (88):* Suppose that

$$\begin{aligned} |\tilde{r} - r| &= \gamma |(\tilde{P} - P)\bar{V}| \\ &\leq c \left\{ \frac{1}{\sqrt{1-\gamma}} \mathbf{1} + \sqrt{\text{Var}_P(V^*)} \right\} \sqrt{\frac{\log \frac{N|\mathcal{S}||\mathcal{A}|}{\delta}}{N\mu_{\min}}}, \end{aligned} \quad (91)$$

holds for some constant $c > 0$. By replacing Lemma 5 in the proof of [12, Lemma 4] with this bound, we can arrive at (88) immediately. In what follows, we demonstrate how to prove the bound (91), which follows a similar argument as in the proof of Lemma 7.

Let us begin with the following triangle inequality:

$$|(\tilde{P} - P)\bar{V}| \leq |(\tilde{P} - P)(\bar{V} - V^*)| + |(\tilde{P} - P)V^*|, \quad (92)$$

leaving us with two terms to control.

- Similar to (141), by applying the Hoeffding inequality and taking the union bound over all $(s, a) \in \mathcal{S} \times \mathcal{A}$, we can control the first term on the right-hand side of (92) as follows:

$$\begin{aligned} \|(\tilde{P} - P)(\bar{V} - V^*)\|_{\infty} &\leq \max_{(s,a) \in \mathcal{S} \times \mathcal{A}} \sqrt{\frac{2 \log \left(\frac{2N|\mathcal{S}||\mathcal{A}|}{\delta} \right)}{K_N(s,a)}} \|\bar{V} - V^*\|_{\infty} \\ &\leq \sqrt{\frac{4 \log \left(\frac{2N|\mathcal{S}||\mathcal{A}|}{\delta} \right)}{N\mu_{\min}(1-\gamma)}} \end{aligned} \quad (93)$$

with probability at least $1 - \delta$. Here, we have made use of the following property of this phase that

$$\|\bar{V} - V^*\|_{\infty} \leq \|\bar{Q} - Q^*\|_{\infty} \leq 1/\sqrt{1-\gamma}$$

and $K_N(s, a) \geq N\mu_{\min}/2$ for all (s, a) (see Lemma 8).

- Next, we turn attention to the second term on the right-hand side of (92), towards which we resort to the Bernstein inequality. Note that the (s, a) -th entry of $|(\tilde{P} - P)V^*|$ is given by

$$\left| \frac{1}{K_N(s,a)} \sum_{i=1}^{K_N(s,a)} (P_{t_i+1}(s,a) - P(s,a)) V^* \right|, \quad (94)$$

where $K_N(s, a)$ denotes the total number of visits to (s, a) during the first N time instances (see also (113)). In addition, let $t_i := t_i(s, a)$ denote the time stamp when the trajectory visits (s, a) for the i -th time (see also (112)). In view of our derivation for (117), the state transitions happening at times t_1, t_2, \dots, t_k (which are random) are independent for any given integer $k > 0$. It can be calculated that

$$|(P_{t_i+1}(s, a) - P(s, a)) V^*| \leq \frac{1}{1-\gamma}; \quad (95a)$$

$$\begin{aligned} \text{Var} \left(\frac{1}{k} \sum_{i=1}^k (P_{t_i+1}(s, a) - P(s, a)) V^* \right) \\ = \frac{1}{k} \text{Var}_{P(s,a)}(V^*). \end{aligned} \quad (95b)$$

Consequently, invoking the Bernstein inequality implies that with probability at least $1 - \frac{\delta}{|\mathcal{S}||\mathcal{A}|}$,

$$\left| \frac{1}{k} \sum_{i=1}^k (P_{t_{i+1}}(s, a) - P(s, a)) V^* \right| \leq \sqrt{\frac{4 \log \left(\frac{2N|\mathcal{S}||\mathcal{A}|}{\delta} \right)}{k} \text{Var}_{P(s,a)}(V^*)} + \frac{4 \log \left(\frac{2N|\mathcal{S}||\mathcal{A}|}{\delta} \right)}{3(1-\gamma)k}$$

holds simultaneously for all $1 \leq k \leq N$. Recognizing the bound $\frac{1}{2}N\mu_{\min} \leq K_N(s, a) \leq N$ and applying the union bound over all $(s, a) \in \mathcal{S} \times \mathcal{A}$ yield

$$|(\tilde{P} - P)V^*| \leq \sqrt{\frac{2 \log \left(\frac{2N|\mathcal{S}||\mathcal{A}|}{\delta} \right)}{N\mu_{\min}} \text{Var}_P(V^*)} + \frac{8 \log \left(\frac{2N|\mathcal{S}||\mathcal{A}|}{\delta} \right)}{3(1-\gamma)N\mu_{\min}}. \quad (96)$$

- Finally, combining (93) and (96) immediately establishes the claim (91).

b) *Proof of the Inequality (89)*: Recalling the variance-reduced update rule (80) and using the Bellman-type equation (87), we obtain

$$\begin{aligned} \hat{\Theta}_t &:= Q_t - \hat{Q} = (I - \Lambda_t)(Q_{t-1} - \hat{Q}) + \Lambda_t(r + \gamma P_t(V_{t-1} - \bar{V}) + \gamma \tilde{P}\bar{V} - r - \gamma P(\hat{V} - \bar{V}) - \gamma \tilde{P}\bar{V}) \\ &= (I - \Lambda_t)(Q_{t-1} - \hat{Q}) \\ &\quad + \Lambda_t(\gamma P_t(V_{t-1} - \bar{V}) - \gamma P(\hat{V} - \bar{V})) \\ &= (I - \Lambda_t)\hat{\Theta}_{t-1} \\ &\quad + \gamma \Lambda_t((P_t - P)(\hat{V} - \bar{V}) + P_t(V_{t-1} - \hat{V})). \end{aligned} \quad (97)$$

Adopting the same expansion as before (see (40)), we arrive at

$$\begin{aligned} \hat{\Theta}_t &= \gamma \underbrace{\sum_{i=1}^t \prod_{j=i+1}^t (I - \Lambda_j) \Lambda_i (P_i - P)(\hat{V} - \bar{V})}_{=: \vartheta_{1,t}} \\ &\quad + \gamma \underbrace{\sum_{i=1}^t \prod_{j=i+1}^t (I - \Lambda_j) \Lambda_i P_i (V_{i-1} - \hat{V})}_{=: \vartheta_{2,t}} \\ &\quad + \underbrace{\prod_{j=1}^t (I - \Lambda_j) \hat{\Theta}_0}_{=: \vartheta_{3,t}}. \end{aligned}$$

Inheriting the results in Lemma 1 and Lemma 2, we can demonstrate that, with probability at least $1 - 2\delta$,

$$\begin{aligned} |\vartheta_{1,t}| &\leq c\gamma \|\hat{V} - \bar{V}\|_{\infty} \sqrt{\eta \log \left(\frac{|\mathcal{S}||\mathcal{A}|t_{\text{epoch}}}{\delta} \right)} \mathbf{1}; \\ |\vartheta_{3,t}| &\leq \begin{cases} (1-\eta)^{\frac{1}{2}t\mu_{\min}} \|\hat{\Theta}_0\|_{\infty} \mathbf{1}, & \text{if } t_{\text{frame}} \leq t \leq t_{\text{epoch}}, \\ \|\hat{\Theta}_0\|_{\infty} \mathbf{1}, & \text{if } t < t_{\text{frame}}. \end{cases} \end{aligned}$$

Repeating the same argument as for Theorem 5, we reach

$$\begin{aligned} \|\hat{\Theta}_t\|_{\infty} &\leq (1-\rho)^k \frac{\|\hat{Q} - \bar{Q}\|_{\infty}}{1-\gamma} \\ &\quad + \frac{c\gamma}{1-\gamma} \|\hat{V} - \bar{V}\|_{\infty} \sqrt{\eta \log \left(\frac{|\mathcal{S}||\mathcal{A}|t_{\text{epoch}}}{\delta} \right)} + \varepsilon \end{aligned}$$

for some constant $c > 0$, where $k = \max\{0, \lfloor \frac{t-t_{\text{th}}}{t_{\text{frame}}} \rfloor\}$ with t_{th} defined in (78b).

By taking $\eta = c_5 \min \left\{ \frac{(1-\gamma)^2}{\gamma^2 \log \left(\frac{|\mathcal{S}||\mathcal{A}|t_{\text{epoch}}}{\delta} \right)}, \frac{1}{\mu_{\text{frame}}} \right\}$ for some sufficiently small constant $c_5 > 0$ and ensuring that

$$t_{\text{epoch}} \geq t_{\text{th}} + t_{\text{frame}} + \frac{c_6}{(1-\gamma)\eta\mu_{\min}} \log \frac{1}{(1-\gamma)^2}$$

for some large constant $c_6 > 0$, we obtain

$$\begin{aligned} \|\hat{\Theta}_{t_{\text{epoch}}}\|_{\infty} &\leq \frac{\|\hat{Q} - \bar{Q}\|_{\infty}}{8} + \varepsilon \\ &\leq \frac{\|\hat{Q} - Q^*\|_{\infty}}{8} + \frac{\|\bar{Q} - Q^*\|_{\infty}}{8} + \varepsilon, \end{aligned}$$

where the last line follows by the triangle inequality.

B. How Many Epochs Are Needed?

We are now ready to pin down how many epochs are needed to achieve ε -accuracy.

- In Phase 1, the contraction result (86) indicates that, if the algorithm is initialized with $Q_0 = \mathbf{0}$ at the very beginning, then it takes at most

$$\begin{aligned} \log_2 \left(\frac{\|Q^*\|_{\infty}}{\max \left\{ \varepsilon, \frac{1}{\sqrt{1-\gamma}} \right\}} \right) &\leq \log_2 \left(\frac{1}{\sqrt{1-\gamma}} \right) \\ &\quad + \log_2 \left(\frac{1}{\varepsilon(1-\gamma)} \right) \end{aligned}$$

epochs to yield $\|\bar{Q} - Q^*\|_{\infty} \leq \max \left\{ \frac{1}{\sqrt{1-\gamma}}, \varepsilon \right\}$ (so as to enter Phase 2). Clearly, if the target accuracy level $\varepsilon > \frac{1}{\sqrt{1-\gamma}}$, then the algorithm terminates in this phase.

- Suppose now that the target accuracy level $\varepsilon \leq \frac{1}{\sqrt{1-\gamma}}$. Once the algorithm enters Phase 2, the dynamics can be characterized by (90). Given that \bar{Q} is also the last iterate of the preceding epoch, the property (90) provides a recursive relation across epochs. Standard recursive analysis thus reveals that: within at most

$$c_7 \log \left(\frac{1}{\varepsilon \sqrt{1-\gamma}} \right) \leq c_7 \log \left(\frac{1}{\varepsilon(1-\gamma)} \right)$$

epochs (with $c_7 > 0$ some constant), we are guaranteed to attain an ℓ_{∞} estimation error at most 3ε .

To summarize, a total number of $O(\log \frac{1}{\varepsilon(1-\gamma)} + \log \frac{1}{1-\gamma})$ epochs are sufficient for our purpose. This concludes the proof.

X. DISCUSSION

This work develops a sharper finite-sample analysis of the classical asynchronous Q-learning algorithm, highlighting and refining its dependency on intrinsic features of the Markovian trajectory induced by the behavior policy. Our sample complexity bound strengthens the state-of-the-art result by an order

of at least $|\mathcal{S}||\mathcal{A}|$. A variance-reduced variant of asynchronous Q-learning is also analyzed, exhibiting improved scaling with the effective horizon $\frac{1}{1-\gamma}$.

Our findings and the analysis framework developed herein suggest a couple of directions for future investigation. For instance, our improved sample complexity of asynchronous Q-learning has a dependence of $\frac{1}{(1-\gamma)^5}$ on the effective horizon, which is inferior to its model-based counterpart. In the synchronous setting, [32], [34] recently demonstrated Q-learning has a dependence of $\frac{1}{(1-\gamma)^4}$, which is tight up to logarithmic factors. In light of this development, it would be important to determine the exact scaling for the asynchronous setting, which is left as future work. In addition, it would be interesting to see whether the techniques developed herein can be exploited towards understanding model-free algorithms with more sophisticated exploration schemes [64]. Finally, asynchronous Q-learning on a single Markovian trajectory is closely related to coordinate descent with coordinates selected according to a Markov chain; one would naturally ask whether our analysis framework can yield improved convergence guarantees for general Markov-chain-based optimization algorithms [65], [66].

APPENDIX

In this section, we gather some basic facts about Markov chains. Before proceeding, we remind the readers of some notation. For any two probability distributions μ and ν , denote by $d_{\text{TV}}(\mu, \nu)$ the total variation distance between μ and ν (cf. (5)). Recall the definition of uniform ergodicity in Section I-B. For any *time-homogeneous* and *uniformly ergodic* Markov chain (X_0, X_1, X_2, \dots) with transition kernel P , finite state space \mathcal{X} and stationary distribution μ , we let $P^t(\cdot | x)$ denote the distribution of X_t conditioned on $X_0 = x$. Then the mixing time t_{mix} of this Markov chain is defined by

$$t_{\text{mix}}(\epsilon) := \min \left\{ t \mid \max_{x \in \mathcal{X}} d_{\text{TV}}(P^t(\cdot | x), \mu) \leq \epsilon \right\}; \quad (98a)$$

$$t_{\text{mix}} := t_{\text{mix}}(1/4). \quad (98b)$$

A. Concentration of Empirical Distributions of Markov Chains

We first record a result concerning the concentration of measure of the empirical distribution of a uniformly ergodic Markov chain, which makes clear the role of the mixing time.

Lemma 8: Consider the above-mentioned Markov chain. For any $0 < \delta < 1$, if $t \geq \frac{443t_{\text{mix}}}{\mu_{\min}} \log \frac{4|\mathcal{X}|}{\delta}$, then for any $y \in \mathcal{X}$, one has

$$\mathbb{P}_{X_1=y} \left\{ \exists x \in \mathcal{X} : \left| \sum_{i=1}^t \mathbb{1}\{X_i = x\} - t\mu(x) \right| \geq \frac{1}{2}t\mu(x) \right\} \leq \delta. \quad (99)$$

Proof: To begin with, consider the scenario when $X_1 \sim \mu$, namely, when X_1 follows the stationary distribution of the chain. Then [17, Theorem 3.4] tells us that: for any given

$x \in \mathcal{X}$ and any $\tau \geq 0$,

$$\begin{aligned} \mathbb{P}_{X_1 \sim \mu} \left\{ \left| \sum_{i=1}^t \mathbb{1}\{X_i = x\} - t\mu(x) \right| \geq \tau \right\} \\ \leq 2 \exp \left(-\frac{\tau^2 \gamma_{\text{ps}}}{8(t + 1/\gamma_{\text{ps}})\mu(x) + 20\tau} \right) \\ \leq 2 \exp \left(-\frac{\tau^2/t_{\text{mix}}}{16(t + 2t_{\text{mix}})\mu(x) + 40\tau} \right), \end{aligned} \quad (100)$$

where γ_{ps} stands for the so-called *pseudo spectral gap* as defined in [17, Section 3.1]. Here, the first inequality relies on the fact $\text{Var}_{X_i \sim \mu}[\mathbb{1}\{X_i = x\}] = \mu(x)(1 - \mu(x)) \leq \mu(x)$, while the last inequality results from the fact $\gamma_{\text{ps}} \geq 1/(2t_{\text{mix}})$ that holds for uniformly ergodic chains (cf. [17, Proposition 3.4]). Consequently, for any $t \geq t_{\text{mix}}$ and any $\tau \geq 0$, one can continue the bound (100) to obtain

$$\begin{aligned} (100) &\leq 2 \exp \left(-\frac{\tau^2}{48t\mu(x)t_{\text{mix}} + 40\tau t_{\text{mix}}} \right) \\ &\leq 2 \max \left\{ \exp \left(-\frac{\tau^2}{96t\mu(x)t_{\text{mix}}} \right), \exp \left(-\frac{\tau}{80t_{\text{mix}}} \right) \right\} \\ &\leq \frac{\delta}{|\mathcal{X}|}, \end{aligned}$$

provided that

$$\tau \geq \max \left\{ 10\sqrt{t\mu(x)t_{\text{mix}} \log \frac{2|\mathcal{X}|}{\delta}}, 80t_{\text{mix}} \log \frac{2|\mathcal{X}|}{\delta} \right\}.$$

As a result, by taking $\tau = \frac{10}{21}t\mu(x)$ and applying the union bound, we reach

$$\begin{aligned} \mathbb{P}_{X_1 \sim \mu} \left\{ \exists x \in \mathcal{X} : \left| \sum_{i=1}^t \mathbb{1}\{X_i = x\} - t\mu(x) \right| \geq \frac{10}{21}t\mu(x) \right\} \\ \leq \sum_{x \in \mathcal{X}} \mathbb{P}_{X_1 \sim \mu} \left\{ \left| \sum_{i=1}^t \mathbb{1}\{X_i = x\} - t\mu(x) \right| \geq \frac{10}{21}t\mu(x) \right\} \leq \delta, \end{aligned} \quad (101)$$

as long as

$$\frac{10}{21}t\mu(x) \geq \max \left\{ 10\sqrt{t\mu(x)t_{\text{mix}} \log \frac{2|\mathcal{X}|}{\delta}}, 80t_{\text{mix}} \log \frac{2|\mathcal{X}|}{\delta} \right\}$$

for all $x \in \mathcal{X}$, or equivalently, when

$$t \geq \frac{441t_{\text{mix}}}{\mu_{\min}} \log \frac{2|\mathcal{X}|}{\delta} \quad \text{with } \mu_{\min} := \min_{x \in \mathcal{X}} \mu(x).$$

Next, we seek to extend the above result to the more general case when X_1 takes an arbitrary state $y \in \mathcal{X}$. From the definition of $t_{\text{mix}}(\cdot)$ (cf. (98a)), we know that

$$d_{\text{TV}} \left(\sup_{y \in \mathcal{X}} P^{t_{\text{mix}}(\delta)}(\cdot | y), \mu \right) \leq \delta. \quad (102)$$

This taken together with the definition of d_{TV} (cf. (5)) reveals that: for any event \mathcal{B} belonging to the σ -algebra generated by $\{X_\tau\}_{\tau \geq t_{\text{mix}}(\delta)}$, one has

$$\begin{aligned} &|\mathbb{P}\{\mathcal{B} \mid X_1 = y\} - \mathbb{P}\{\mathcal{B} \mid X_1 \sim \mu\}| \\ &= \left| \sum_{s \in \mathcal{S}} \mathbb{P}\{\mathcal{B} \mid X_{t_{\text{mix}}(\delta)} = s\} \mathbb{P}\{X_{t_{\text{mix}}(\delta)} = s \mid X_1 = y\} \right| \end{aligned}$$

$$\begin{aligned}
& \left| -\sum_{s \in \mathcal{S}} \mathbb{P}\{\mathcal{B} \mid X_{t_{\text{mix}}(\delta)} = s\} \mathbb{P}\{X_{t_{\text{mix}}(\delta)} = s \mid X_1 \sim \mu\} \right| \\
& \leq \max \left\{ \sum_{s \in \mathcal{S}_+} \left[\mathbb{P}\{X_{t_{\text{mix}}(\delta)} = s \mid X_1 = y\} - \mathbb{P}\{X_{t_{\text{mix}}(\delta)} = s \mid X_1 \sim \mu\} \right], \right. \\
& \quad \left. \sum_{s \in \mathcal{S}_-} \left[\mathbb{P}\{X_{t_{\text{mix}}(\delta)} = s \mid X_1 \sim \mu\} - \mathbb{P}\{X_{t_{\text{mix}}(\delta)} = s \mid X_1 = y\} \right] \right\} \\
& \leq \sup_{A \subseteq \mathcal{S}} \left| \mathbb{P}\{X_{t_{\text{mix}}(\delta)} \in A \mid X_1 = y\} - \mathbb{P}\{X_{t_{\text{mix}}(\delta)} \in A \mid X_1 \sim \mu\} \right| \leq \delta,
\end{aligned} \tag{103}$$

where we define

$$\begin{aligned}
\mathcal{S}_+ &:= \{s \in \mathcal{S} : \mathbb{P}\{X_{t_{\text{mix}}(\delta)} = s \mid X_1 = y\} > \mathbb{P}\{X_{t_{\text{mix}}(\delta)} = s \mid X_1 \sim \mu\}\}; \\
\mathcal{S}_- &:= \{s \in \mathcal{S} : \mathbb{P}\{X_{t_{\text{mix}}(\delta)} = s \mid X_1 = y\} < \mathbb{P}\{X_{t_{\text{mix}}(\delta)} = s \mid X_1 \sim \mu\}\}.
\end{aligned}$$

Here, the last inequality in (103) follows from the inequality (102) and the definition (5) of the total-variation distance. As a consequence, one obtains

$$\begin{aligned}
& \sup_{y \in \mathcal{X}} \mathbb{P}_{X_1=y} \left\{ \exists x \in \mathcal{X} : \left| \sum_{i=t_{\text{mix}}(\delta)}^t \mathbb{1}\{X_i = x\} - (t - t_{\text{mix}}(\delta))\mu(x) \right| \right. \\
& \quad \left. \geq \frac{10}{21}(t - t_{\text{mix}}(\delta))\mu(x) \right\} \\
& \leq \mathbb{P}_{X_1 \sim \mu} \left\{ \exists x \in \mathcal{X} : \left| \sum_{i=t_{\text{mix}}(\delta)}^t \mathbb{1}\{X_i = x\} - (t - t_{\text{mix}}(\delta))\mu(x) \right| \right. \\
& \quad \left. \geq \frac{10}{21}(t - t_{\text{mix}}(\delta))\mu(x) \right\} + \delta \leq 2\delta,
\end{aligned} \tag{104}$$

with the proviso that $t \geq t_{\text{mix}}(\delta) + \frac{441t_{\text{mix}}}{\mu_{\min}} \log \frac{2|\mathcal{X}|}{\delta}$.

To finish up, we recall from [17, Section 1.1] that $t_{\text{mix}}(\delta) \leq 2t_{\text{mix}} \log \frac{2}{\delta}$. Consequently, if $t \geq \frac{443t_{\text{mix}}}{\mu_{\min}} \log \frac{2|\mathcal{X}|}{\delta} \geq t_{\text{mix}}(\delta) + \frac{441t_{\text{mix}}}{\mu_{\min}} \log \frac{2|\mathcal{X}|}{\delta}$, then one has

$$\begin{aligned}
& \frac{1}{2}(t - t_{\text{mix}}(\delta))\mu(x) - t_{\text{mix}}(\delta) \geq \frac{441t_{\text{mix}}}{2} \log \frac{2|\mathcal{X}|}{\delta} \geq 100t_{\text{mix}}(\delta), \\
& \implies \frac{1}{2}(t - t_{\text{mix}}(\delta))\mu(x) - t_{\text{mix}}(\delta) \geq \frac{10}{21}(t - t_{\text{mix}}(\delta))\mu(x).
\end{aligned}$$

These taken together lead to (105), shown at the bottom of the next page. Here, the last inequality of (105) results from (104). Replacing δ with $\delta/2$ thus concludes the proof. \square

B. Connection Between the Mixing Time and the Cover Time

Lemma 8 combined with the definition (9) immediately reveals the following upper bound on the cover time:

$$t_{\text{cover}} = O\left(\frac{t_{\text{mix}}}{\mu_{\min}} \log |\mathcal{X}|\right). \tag{106}$$

In addition, while a general matching converse bound (namely, $t_{\text{mix}}/\mu_{\min} = \tilde{O}(t_{\text{cover}})$) is not available, we can come up with some special examples for which the bound (106) is provably tight.

Example 1: Consider a time-homogeneous Markov chain with state space $\mathcal{X} := \{1, \dots, |\mathcal{X}|\}$ and probability transition matrix $\mathbf{P} \in \mathbb{R}^{|\mathcal{X}| \times |\mathcal{X}|}$ as

$$\mathbf{P} = \left(1 - \frac{q(k+1)}{2}\right) \mathbf{I}_{|\mathcal{X}|} + \frac{q}{|\mathcal{X}|} \begin{bmatrix} k\mathbf{1}_{|\mathcal{X}|} \mathbf{1}_{|\mathcal{X}|/2}^\top & \mathbf{1}_{|\mathcal{X}|} \mathbf{1}_{|\mathcal{X}|/2}^\top \end{bmatrix} \tag{107}$$

for some quantities $q > 0$ and $k \geq 1$. Suppose $q(k+1) < 2$ and $|\mathcal{X}| \geq 3$. Then this chain obeys

$$t_{\text{cover}} \geq \frac{t_{\text{mix}}}{(8 \log 2 + 4 \log \frac{1}{\mu_{\min}}) \mu_{\min}}. \tag{108}$$

With the lower bound (108) in place, we conclude that the upper bound (106) is, in general, nearly un-improvable (up to some logarithmic factor).

Remark 6: We shall take a moment to briefly discuss the key design rationale behind Example 1. Let us partition the state space into two halves, denoted respectively by \mathcal{X}_1 and \mathcal{X}_2 . From every state $s \in \mathcal{X}$, it is much easier to transition into the first half \mathcal{X}_1 rather than the second half \mathcal{X}_2 . This leads to two properties: (i) the stationary distribution of any state in \mathcal{X}_2 is much lower than that of a state in \mathcal{X}_1 ; (ii) the cover time also increases as the stationary distribution w.r.t. \mathcal{X}_2 decreases, given that it becomes more difficult to traverse the second half. As a result, we can guarantee that t_{cover} is proportional to μ_{\min} through this type of designs. On the other hand, the example is also constructed in a way such that all states are “lazy”, meaning that they are more inclined to stay unchanged rather than moving to a different state. The level of laziness clearly controls how fast the Markov chain mixes, as well as how long it takes to cover all states. This in turn allows one to ensure that t_{cover} is proportional to t_{mix} . More details can be found in the proof below.

Proof: As can be easily verified, this chain is reversible, whose stationary distribution vector $\boldsymbol{\mu} \in \mathbb{R}^{|\mathcal{X}|}$ obeys

$$\boldsymbol{\mu} = \frac{2}{(k+1)|\mathcal{X}|} \begin{bmatrix} k\mathbf{1}_{|\mathcal{X}|/2} \\ \mathbf{1}_{|\mathcal{X}|/2} \end{bmatrix}.$$

As a result, the minimum state occupancy probability of the stationary distribution is given by

$$\mu_{\min} := \min_{1 \leq x \leq |\mathcal{X}|} \mu_x = \frac{2}{(k+1)|\mathcal{X}|}. \tag{109}$$

In addition, the reversibility of this chain implies that the matrix $\mathbf{P}^d := \mathbf{D}^{\frac{1}{2}} \mathbf{P} \mathbf{D}^{-\frac{1}{2}}$ with $\mathbf{D} := \text{diag}[\boldsymbol{\mu}]$ is symmetric and has the same set of eigenvalues as \mathbf{P} [67]. A little algebra yields

$$\begin{aligned}
\mathbf{P}^d &= \left(1 - \frac{q(k+1)}{2}\right) \mathbf{I}_{|\mathcal{X}|} \\
&+ \frac{q}{|\mathcal{X}|} \begin{bmatrix} k\mathbf{1}_{|\mathcal{X}|/2} \mathbf{1}_{|\mathcal{X}|/2}^\top & \sqrt{k} \mathbf{1}_{|\mathcal{X}|/2} \mathbf{1}_{|\mathcal{X}|/2}^\top \\ \sqrt{k} \mathbf{1}_{|\mathcal{X}|/2} \mathbf{1}_{|\mathcal{X}|/2}^\top & \mathbf{1}_{|\mathcal{X}|/2} \mathbf{1}_{|\mathcal{X}|/2}^\top \end{bmatrix},
\end{aligned}$$

allowing us to determine the eigenvalues $\{\lambda_i\}_{1 \leq i \leq |\mathcal{X}|}$ as follows

$$\lambda_1 = 1 \quad \text{and} \quad \lambda_i = 1 - \frac{q(k+1)}{2} > 0 \quad (i \geq 2).$$

We are now ready to establish the lower bound on the cover time. First of all, the well-known connection between the spectral gap and the mixing time gives [17, Proposition 3.3]

$$t_{\text{mix}} \leq \frac{2 \log 2 + \log \frac{1}{\mu_{\min}}}{2(1 - \lambda_2)} = \frac{2 \log 2 + \log \frac{1}{\mu_{\min}}}{q(k+1)}. \quad (110)$$

In addition, let (x_0, x_1, \dots) be the corresponding Markov chain, and assume that $x_0 \sim \mu$, where μ stands for the stationary distribution. Consider the last state — denoted by $|\mathcal{X}|$, which enjoys the minimum state occupancy probability μ_{\min} . For any integer $t > 0$ one has

$$\begin{aligned} & \mathbb{P}\{x_l \neq |\mathcal{X}|, \forall 0 \leq l \leq t\} \\ & \stackrel{(i)}{=} \mathbb{P}\{x_0 \neq |\mathcal{X}|\} \prod_{l=1}^t \mathbb{P}\{x_l \neq |\mathcal{X}| \mid x_0 \neq |\mathcal{X}|, \dots, x_{l-1} \neq |\mathcal{X}|\} \\ & \stackrel{(ii)}{\geq} \mathbb{P}\{x_0 \neq |\mathcal{X}|\} \prod_{l=1}^t \min_{j: j \neq |\mathcal{X}|} \mathbb{P}\{x_l \neq |\mathcal{X}| \mid x_{l-1} = j\} \\ & \stackrel{(iii)}{=} \left(1 - \frac{2}{(k+1)|\mathcal{X}|}\right) \left(1 - \frac{q}{|\mathcal{X}|}\right)^t \\ & \stackrel{(iv)}{\geq} \left(1 - \frac{2}{(k+1)|\mathcal{X}|}\right) \left(1 - \frac{2qt}{|\mathcal{X}|}\right), \end{aligned}$$

where (i) follows from the chain rule, (ii) relies on the Markovian property, (iii) results from the construction (107), and (iv) holds as long as $\frac{q}{|\mathcal{X}|}t < \frac{1}{2}$. Consequently, if $|\mathcal{X}| \geq 3$ and if $t < \frac{|\mathcal{X}|}{8q}$, then one necessarily has

$$\begin{aligned} & \mathbb{P}\{x_l \neq |\mathcal{X}|, \forall 0 \leq l \leq t\} \\ & \geq \left(1 - \frac{2}{(k+1)|\mathcal{X}|}\right) \left(1 - \frac{2qt}{|\mathcal{X}|}\right) > \frac{1}{2}. \end{aligned}$$

This taken collectively with the definition of t_{cover} (cf. (9)) reveals that

$$t_{\text{cover}} \geq \frac{|\mathcal{X}|}{8q} \geq \frac{t_{\text{mix}}}{(8 \log 2 + 4 \log \frac{1}{\mu_{\min}}) \mu_{\min}},$$

where the last inequality is a direct consequence of (109) and (110). \square

C. Proof of Lemma 1

Fix any state-action pair $(s, a) \in \mathcal{S} \times \mathcal{A}$, and let us look at $\beta_{1,t}(s, a)$, namely, the (s, a) -th entry of

$$\beta_{1,t} = \gamma \sum_{i=1}^t \prod_{j=i+1}^t (I - \Lambda_j) \Lambda_i (P_i - P) V^*.$$

For convenience of presentation, we abuse the notation to let $\Lambda_j(s, a)$ denote the (s, a) -th diagonal entry of the diagonal matrix Λ_j , and $P_t(s, a)$ (resp. $P(s, a)$) the (s, a) -th row of P_t (resp. P). In view of the definition (40), we can write

$$\begin{aligned} \beta_{1,t}(s, a) &= \gamma \sum_{i=1}^t \prod_{j=i+1}^t \left[(1 - \Lambda_j(s, a)) \Lambda_i(s, a) \right. \\ & \quad \left. \cdot (P_i(s, a) - P(s, a)) V^* \right]. \end{aligned} \quad (111)$$

As it turns out, it is convenient to study this expression by defining

$$t_k(s, a) := \text{the time stamp when the trajectory visits } (s, a) \text{ for the } k\text{-th time} \quad (112)$$

and

$$K_t(s, a) := \max\{k \mid t_k(s, a) \leq t\}, \quad (113)$$

namely, the total number of times — during the first t iterations — that the sample trajectory visits (s, a) . With these in place, the special form of Λ_j (cf. (34)) allows us to rewrite (111) as

$$\begin{aligned} \beta_{1,t}(s, a) &= \gamma \sum_{k=1}^{K_t(s, a)} \left[(1 - \eta)^{K_t(s, a) - k} \right. \\ & \quad \left. \cdot (P_{t_k+1}(s, a) - P(s, a)) V^* \right]. \end{aligned} \quad (114)$$

where we suppress the dependency on (s, a) and write $t_k := t_k(s, a)$ to streamline notation. The main step thus boils down to controlling (114).

Towards this, we claim that: with probability at least $1 - \delta$,

$$\begin{aligned} & \left| \sum_{k=1}^K (1 - \eta)^{K-k} \eta (P_{t_k+1}(s, a) - P(s, a)) V^* \right| \\ & \leq \sqrt{\eta \log \left(\frac{|\mathcal{S}||\mathcal{A}|T}{\delta} \right)} \|V^*\|_{\infty} \end{aligned} \quad (115)$$

holds simultaneously for all $(s, a) \in \mathcal{S} \times \mathcal{A}$ and all $1 \leq K \leq T$, provided that $0 < \eta \log \left(\frac{|\mathcal{S}||\mathcal{A}|T}{\delta} \right) < 1$. Recognizing the trivial bound $K_t(s, a) \leq t \leq T$ (by construction (113)) and substituting the claimed bound (115) into the expression (114), we arrive at for all $(s, a) \in \mathcal{S} \times \mathcal{A}$

$$\begin{aligned} |\beta_{1,t}(s, a)| &\leq \gamma \sqrt{\eta \log \left(\frac{|\mathcal{S}||\mathcal{A}|T}{\delta} \right)} \|V^*\|_{\infty} \\ &\leq \sqrt{\eta \log \left(\frac{|\mathcal{S}||\mathcal{A}|T}{\delta} \right)} \|V^*\|_{\infty}, \end{aligned} \quad (116)$$

$$\begin{aligned} & \sup_{y \in \mathcal{X}} \mathbb{P}_{X_1=y} \left\{ \exists x \in \mathcal{X} : \left| \sum_{i=1}^t \mathbb{1}\{X_i = x\} - t\mu(x) \right| \geq \frac{1}{2}t\mu(x) \right\} \\ & \leq \sup_{y \in \mathcal{X}} \mathbb{P}_{X_1=y} \left\{ \exists x \in \mathcal{X} : \left| \sum_{i=t_{\text{mix}}(\delta)}^t \mathbb{1}\{X_i = x\} - (t - t_{\text{mix}}(\delta))\mu(x) \right| \geq \frac{1}{2}(t - t_{\text{mix}}(\delta))\mu(x) - t_{\text{mix}}(\delta) \right\} \\ & \leq \sup_{y \in \mathcal{X}} \mathbb{P}_{X_1=y} \left\{ \exists x \in \mathcal{X} : \left| \sum_{i=t_{\text{mix}}(\delta)}^t \mathbb{1}\{X_i = x\} - (t - t_{\text{mix}}(\delta))\mu(x) \right| \geq \frac{10}{21}(t - t_{\text{mix}}(\delta))\mu(x) \right\} \leq 2\delta. \end{aligned} \quad (105)$$

thus concluding the proof of this lemma. It remains to validate the inequality (115).

Proof of the Inequality (115): We first make the observation that: for any fixed integer $K > 0$, the following vectors

$$\{\mathbf{P}_{t_k+1}(s, a) \mid 1 \leq k \leq K\}$$

are identically and independently distributed.⁴ To justify this observation, let us denote by $\mathbb{P}_{s,a}(\cdot)$ the transition probability from state s when action a is taken. For any $i_1, \dots, i_K \in \mathcal{S}$, one obtains

$$\begin{aligned} & \mathbb{P}\{s_{t_k+1} = i_k \ (\forall 1 \leq k \leq K)\} \\ &= \mathbb{P}\{s_{t_k+1} = i_k \ (\forall 1 \leq k \leq K-1) \text{ and } s_{t_K+1} = i_K\} \\ &= \sum_{m>0} \mathbb{P}\{s_{t_k+1} = i_k \ (\forall 1 \leq k \leq K-1) \text{ and } t_K = m \\ & \quad \text{and } s_{m+1} = i_K\} \\ &\stackrel{(i)}{=} \sum_{m>0} \left[\mathbb{P}\{s_{t_k+1} = i_k \ (\forall 1 \leq k \leq K-1) \text{ and } t_K = m\} \right. \\ & \quad \cdot \mathbb{P}\{s_{m+1} = i_K \mid s_m = s, a_m = a\} \left. \right] \\ &= \mathbb{P}_{s,a}(i_K) \\ & \quad \cdot \sum_{m>0} \mathbb{P}\{s_{t_k+1} = i_k \ (\forall 1 \leq k \leq K-1) \text{ and } t_K = m\} \\ &= \mathbb{P}_{s,a}(i_K) \mathbb{P}\{s_{t_k+1} = i_k \ (\forall 1 \leq k \leq K-1)\}, \end{aligned}$$

where (i) holds true from the Markov property as well as the fact that t_K is an iteration in which the trajectory visits state s and takes action a . Invoking the above identity recursively, we arrive at

$$\mathbb{P}\{s_{t_k+1} = i_k \ (\forall 1 \leq k \leq K)\} = \prod_{j=1}^K \mathbb{P}_{s,a}(i_j), \quad (117)$$

meaning that the state transitions happening at times $\{t_1, \dots, t_K\}$ are independent, each following the distribution $\mathbb{P}_{s,a}(\cdot)$. This clearly demonstrates the independence of $\{\mathbf{P}_{t_k+1}(s, a) \mid 1 \leq k \leq K\}$.

With the above observation in mind, we resort to the Hoeffding inequality to bound the quantity of interest (which has zero mean). To begin with, notice the facts that for all $k \geq 1$,

$$0 \leq \mathbf{P}_{t_k+1}(s, a) \mathbf{V}^* \leq \|\mathbf{V}^*\|_\infty, \text{ and } 0 \leq \mathbf{P}(s, a) \mathbf{V}^* \leq \|\mathbf{V}^*\|_\infty, \quad (118)$$

which gives

$$\begin{aligned} & |(1-\eta)^{K-k} \eta (\mathbf{P}_{t_k+1}(s, a) - \mathbf{P}(s, a)) \mathbf{V}^*| \\ & \leq (1-\eta)^{K-k} \eta \|\mathbf{V}^*\|_\infty. \end{aligned}$$

As a consequence, invoking the Hoeffding inequality [68] implies that

$$\left| \sum_{k=1}^K (1-\eta)^{K-k} \eta (\mathbf{P}_{t_k}(s, a) - \mathbf{P}(s, a)) \mathbf{V}^* \right|$$

⁴The Markov chain w.r.t. the sample trajectory should be viewed as being infinitely long, although we only get to observe its first T samples. The random variables $\{t_k\}$ are, in truth, independent of the choice of T .

$$\begin{aligned} & \leq \sqrt{\frac{1}{2} \sum_{k=1}^K \left((1-\eta)^{K-k} \eta \|\mathbf{V}^*\|_\infty \right)^2 \log \left(\frac{2|\mathcal{S}||\mathcal{A}|T}{\delta} \right)} \\ & \leq \sqrt{\eta \log \left(\frac{|\mathcal{S}||\mathcal{A}|T}{\delta} \right)} \|\mathbf{V}^*\|_\infty \end{aligned} \quad (119)$$

with probability exceeding $1 - \frac{\delta}{|\mathcal{S}||\mathcal{A}|T}$, where the last line holds since

$$\sum_{k=1}^K \left((1-\eta)^{K-k} \eta \right)^2 \leq \eta^2 \sum_{j=0}^{\infty} (1-\eta)^j = \frac{\eta^2}{1-(1-\eta)} = \eta.$$

Taking the union bound over all $(s, a) \in \mathcal{S} \times \mathcal{A}$ and all $1 \leq K \leq T$ then reveals that: with probability at least $1 - \delta$, the inequality (119) holds simultaneously over all $(s, a) \in \mathcal{S} \times \mathcal{A}$ and all $1 \leq K \leq T$. This concludes the proof. \square

D. Proof of Lemma 2 and Lemma 6

Proof of Lemma 2: Let $\beta_{3,t} = \prod_{j=1}^t (\mathbf{I} - \Lambda_j) \Delta_0$. Denote by $\beta_{3,t}(s, a)$ (resp. $\Delta_0(s, a)$) the (s, a) -th entry of $\beta_{3,t}$ (resp. Δ_0). From the definition of $\beta_{3,t}$, it is easily seen that

$$|\beta_{3,t}(s, a)| = (1-\eta)^{K_t(s,a)} |\Delta_0(s, a)|, \quad (120)$$

where $K_t(s, a)$ denotes the number of times the sample trajectory visits (s, a) during the iterations $[1, t]$ (cf. (113)). By virtue of Lemma 8 and the union bound, one has, with probability at least $1 - \delta$, that

$$K_t(s, a) \geq t\mu_{\min}/2 \quad (121)$$

simultaneously over all $(s, a) \in \mathcal{S} \times \mathcal{A}$ and all t obeying $\frac{443\tau_{\min}}{\mu_{\min}} \log \frac{4|\mathcal{S}||\mathcal{A}|T}{\delta} \leq t \leq T$. Substitution into the relation (120) establishes that, with probability greater than $1 - \delta$,

$$|\beta_3(s, a)| \leq (1-\eta)^{\frac{1}{2}t\mu_{\min}} |\Delta_0(s, a)|. \quad (122)$$

holds uniformly over all $(s, a) \in \mathcal{S} \times \mathcal{A}$ and all t obeying $\frac{443\tau_{\min}}{\mu_{\min}} \log \frac{4|\mathcal{S}||\mathcal{A}|T}{\delta} \leq t \leq T$, as claimed.

Proof of Lemma 6: The proof of this lemma is essentially the same as that of Lemma 2, except that we use instead the following lower bound on $K_t(s, a)$ (which is an immediate consequence of Lemma 5)

$$K_t(s, a) \geq \left\lfloor \frac{t}{t_{\text{cover,all}}} \right\rfloor \geq \frac{t}{2t_{\text{cover,all}}} \quad (123)$$

for all $t > t_{\text{cover,all}}$. Therefore, replacing $t\mu_{\min}$ with $t/t_{\text{cover,all}}$ in the above analysis, we establish Lemma 6.

E. Proof of Lemma 3

We prove this fact via an inductive argument. The base case with $t = 0$ is a consequence of the crude bound (49). Now, assume that the claim holds for all iterations up to $t-1$, and we would like to justify it for the t -th iteration as well. Towards this, define

$$h(t) := \begin{cases} \|\Delta_0\|_\infty, & \text{if } t \leq t_{\text{th}}, \\ (1-\gamma)\varepsilon, & \text{if } t > t_{\text{th}}. \end{cases} \quad (124)$$

Recall that $(1 - \eta)^{\frac{1}{2}t\mu_{\min}} \leq (1 - \gamma)\varepsilon$ for any $t \geq t_{\text{th}}$. Therefore, combining the inequality (47) with the induction hypotheses indicates that

$$\begin{aligned} |\Delta_t| &\leq \gamma \sum_{i=1}^t \prod_{j=i+1}^t (\mathbf{I} - \Lambda_j) \Lambda_i \mathbf{1} \cdot \left(\frac{\tau_1 \|\mathbf{V}^*\|_{\infty}}{1 - \gamma} + u_{i-1} + \varepsilon \right) \\ &\quad + \tau_1 \|\mathbf{V}^*\|_{\infty} \mathbf{1} + h(t) \mathbf{1} \\ &= \gamma \sum_{i=1}^t \prod_{j=i+1}^t (\mathbf{I} - \Lambda_j) \Lambda_i \mathbf{1} u_{i-1} + \tau_1 \|\mathbf{V}^*\|_{\infty} \mathbf{1} + h(t) \mathbf{1} \\ &\quad + \gamma \sum_{i=1}^t \prod_{j=i+1}^t (\mathbf{I} - \Lambda_j) \Lambda_i \mathbf{1} \left(\frac{\tau_1 \|\mathbf{V}^*\|_{\infty}}{1 - \gamma} + \varepsilon \right). \end{aligned}$$

Taking this together with the inequality (51b) and rearranging terms, we obtain

$$\begin{aligned} |\Delta_t| &\leq \gamma \sum_{i=1}^t \prod_{j=i+1}^t (\mathbf{I} - \Lambda_j) \Lambda_i \mathbf{1} u_{i-1} + \frac{\gamma \tau_1 \|\mathbf{V}^*\|_{\infty}}{1 - \gamma} \mathbf{1} \\ &\quad + \gamma \varepsilon \mathbf{1} + \tau_1 \|\mathbf{V}^*\|_{\infty} \mathbf{1} + h(t) \mathbf{1} \\ &= \frac{\tau_1 \|\mathbf{V}^*\|_{\infty}}{1 - \gamma} \mathbf{1} + \gamma \varepsilon \mathbf{1} + \gamma \sum_{i=1}^t \prod_{j=i+1}^t (\mathbf{I} - \Lambda_j) \Lambda_i \mathbf{1} u_{i-1} \\ &\quad + h(t) \mathbf{1} \\ &= \frac{\tau_1 \|\mathbf{V}^*\|_{\infty}}{1 - \gamma} \mathbf{1} + \gamma \varepsilon \mathbf{1} + \mathbf{v}_t + (1 - \gamma) \varepsilon \mathbf{1} \{t > t_{\text{th}}\} \mathbf{1} \\ &\leq \frac{\tau_1 \|\mathbf{V}^*\|_{\infty}}{1 - \gamma} \mathbf{1} + \varepsilon \mathbf{1} + \mathbf{v}_t, \end{aligned} \quad (125)$$

where we have used the definition of \mathbf{v}_t in (52). This taken collectively with the definition $u_t = \|\mathbf{v}_t\|_{\infty}$ establishes that

$$\|\Delta_t\|_{\infty} \leq \frac{\tau_1 \|\mathbf{V}^*\|_{\infty}}{1 - \gamma} + \varepsilon + u_t$$

as claimed. This concludes the proof.

F. Proof of Lemma 4

We shall prove this result by induction over the index k . To start with, consider the base case where $k = 0$ and $t < t_{\text{th}} + t_{\text{frame}}$. By definition, it is straightforward to see that $u_0 \leq \|\Delta_0\|_{\infty} / (1 - \gamma) = w_0$. In fact, repeating our argument for the crude bound (see Section VI-B.2) immediately reveals that

$$\forall t \geq 0: \quad u_t \leq \frac{\|\Delta_0\|_{\infty}}{1 - \gamma} = w_0, \quad (126)$$

thus indicating that the inequality (55) holds for the base case. In what follows, we assume that the inequality (55) holds up to $k - 1$, and would like to extend it to the case with all t obeying $\lfloor \frac{t - t_{\text{th}}}{t_{\text{frame}}} \rfloor = k$.

Consider any $0 \leq j < t_{\text{frame}}$. In view of the definition of \mathbf{v}_t (cf. (52)) as well as our induction hypotheses, one can arrange terms to derive

$$\begin{aligned} \mathbf{v}_{t_{\text{th}} + kt_{\text{frame}} + j} &= \gamma \sum_{i=1}^{t_{\text{th}} + kt_{\text{frame}} + j} \prod_{n=i+1}^{t_{\text{th}} + kt_{\text{frame}} + j} (\mathbf{I} - \Lambda_n) \Lambda_i \mathbf{1} u_{i-1} \\ &= \gamma \sum_{s=0}^{k-1} \left\{ \sum_{i: \max \left\{ \lfloor \frac{i-j-1-t_{\text{th}}}{t_{\text{frame}}} \rfloor, 0 \right\} = s} \prod_{n=i+1}^{t_{\text{th}} + kt_{\text{frame}} + j} (\mathbf{I} - \Lambda_n) \Lambda_i \mathbf{1} u_{i-1} \right\} \end{aligned}$$

$$\leq \gamma \sum_{s=0}^{k-1} \left\{ \sum_{i: \max \left\{ \lfloor \frac{i-j-1-t_{\text{th}}}{t_{\text{frame}}} \rfloor, 0 \right\} = s} \prod_{n=i+1}^{t_{\text{th}} + kt_{\text{frame}} + j} (\mathbf{I} - \Lambda_n) \Lambda_i \mathbf{1} \right\} w_s, \quad (127)$$

where the last inequality follows from our induction hypotheses, the non-negativity of $(\mathbf{I} - \Lambda_j) \Lambda_i \mathbf{1}$, and the fact that w_s is non-increasing.

Given any state-action pair $(s, a) \in \mathcal{S} \times \mathcal{A}$, let us look at the (s, a) -th entry of $\mathbf{v}_{t_{\text{th}} + kt_{\text{frame}} + j}$ — denoted by $v_{t_{\text{th}} + kt_{\text{frame}} + j}(s, a)$, towards which it is convenient to pause and introduce some notation. Recall that $N_i^n(s, a)$ has been used to denote the number of visits to the state-action pair (s, a) between iteration i and iteration n (including i and n). To help study the behavior in each timeframe, we introduce the following quantities

$$L_h^{k-1} := N_i^n(s, a) \quad (128)$$

with $i = t_{\text{th}} + ht_{\text{frame}} + j + 1$, $n = t_{\text{th}} + kt_{\text{frame}} + j$ for every $h \leq k - 1$. Lemma 8 tells us that, with probability at least $1 - 2\delta$,

$$L_h^{k-1} \geq (k - h) \mu_{\text{frame}} \quad \text{with } \mu_{\text{frame}} = \frac{1}{2} \mu_{\min} t_{\text{frame}}, \quad (129)$$

which holds uniformly over all state-action pairs (s, a) . Armed with this set of notation, it is straightforward to use the expression (127) to verify that

$$\begin{aligned} v_{t_{\text{th}} + kt_{\text{frame}} + j}(s, a) &\leq \gamma \sum_{h=0}^{k-1} \eta \left\{ (1 - \eta)^{L_h^{k-1} - 1} + (1 - \eta)^{L_h^{k-1} - 2} \right. \\ &\quad \left. + \cdots + (1 - \eta)^{L_{h+1}^{k-1}} \right\} w_h \\ &= \gamma \sum_{h=0}^{k-1} \left((1 - \eta)^{L_{h+1}^{k-1}} - (1 - \eta)^{L_h^{k-1}} \right) w_h \\ &=: \gamma \sum_{h=0}^{k-1} (\alpha_{h+1} - \alpha_h) w_h, \end{aligned} \quad (130)$$

where we denote $\alpha_h := (1 - \eta)^{L_h^{k-1}}$ for any $h \leq k - 1$ and $\alpha_k := 1$.

A little algebra further leads to

$$\begin{aligned} \gamma \sum_{h=0}^{k-1} (\alpha_{h+1} - \alpha_h) w_h &= \gamma (\alpha_k w_{k-1} - \alpha_0 w_0) \\ &\quad + \gamma \sum_{h=1}^{k-1} \alpha_h (w_{h-1} - w_h). \end{aligned} \quad (131)$$

Thus, in order to control the quantity $v_{t_{\text{th}} + kt_{\text{frame}} + j}(s, a)$, it suffices to control the right-hand side of (131), for which we start by bounding the last term. Plugging in the definitions of w_h and α_h yields

$$\begin{aligned} \frac{1 - \gamma}{\|\Delta_0\|_{\infty}} \sum_{h=1}^{k-1} \alpha_h (w_{h-1} - w_h) &= \sum_{h=1}^{k-1} (1 - \eta)^{L_h^{k-1}} (1 - \rho)^{h-1} \rho \\ &\leq \rho \sum_{h=1}^{k-1} (1 - \eta)^{(k-h)\mu_{\text{frame}}} (1 - \rho)^{h-1}, \end{aligned}$$

where the last inequality results from the fact (129). Additionally, direct calculation yields

$$\begin{aligned}
& \rho \sum_{h=1}^{k-1} (1-\eta)^{(k-h)\mu_{\text{frame}}} (1-\rho)^{h-1} \\
&= \rho(1-\eta)^{(k-1)\mu_{\text{frame}}} \sum_{h=1}^{k-1} \left(\frac{1-\rho}{(1-\eta)^{\mu_{\text{frame}}}} \right)^{h-1} \\
&= \rho(1-\eta)^{(k-1)\mu_{\text{frame}}} \frac{1 - \left(\frac{1-\rho}{(1-\eta)^{\mu_{\text{frame}}}} \right)^{k-1}}{1 - \frac{1-\rho}{(1-\eta)^{\mu_{\text{frame}}}}} \\
&= \rho(1-\eta)^{\mu_{\text{frame}}} \frac{(1-\rho)^{k-1} - (1-\eta)^{(k-1)\mu_{\text{frame}}}}{(1-\rho) - (1-\eta)^{\mu_{\text{frame}}}} \\
&\leq \rho(1-\eta)^{\mu_{\text{frame}}} \frac{(1-\rho)^{k-1}}{(1-\rho) - (1-\eta)^{\mu_{\text{frame}}}}, \quad (132)
\end{aligned}$$

where the last inequality makes use of the fact that

$$\begin{aligned}
& (1-\rho) - (1-\eta)^{\mu_{\text{frame}}} \\
&= 1 - (1-\gamma)(1 - (1-\eta)^{\mu_{\text{frame}}}) - (1-\eta)^{\mu_{\text{frame}}} \\
&= \gamma \{1 - (1-\eta)^{\mu_{\text{frame}}}\} = \frac{\gamma}{1-\gamma} \rho \geq 0. \quad (133)
\end{aligned}$$

Combining the inequalities (130), (131) and (132) and using the fact $\alpha_0 w_0 \geq 0$ give

$$\begin{aligned}
\mathbf{v}_{t_{\text{th}}+kt_{\text{frame}}+j}(s, a) &\leq \gamma \sum_{h=1}^{k-1} \alpha_h (w_{h-1} - w_h) + \gamma \alpha_k w_{k-1} \\
&\leq \frac{\|\Delta_0\|_{\infty}}{1-\gamma} \left\{ \gamma \rho(1-\eta)^{\mu_{\text{frame}}} \frac{(1-\rho)^{k-1}}{(1-\rho) - (1-\eta)^{\mu_{\text{frame}}}} \right. \\
&\quad \left. + \gamma(1-\rho)^{k-1} \right\}. \quad (134)
\end{aligned}$$

We are now ready to justify that $\mathbf{v}_{t_{\text{th}}+kt_{\text{frame}}+j}(s, a) \leq w_k$. Note that the observation (133) implies

$$\begin{aligned}
\gamma \frac{\rho(1-\eta)^{\mu_{\text{frame}}}}{(1-\rho) - (1-\eta)^{\mu_{\text{frame}}}} &= \gamma \frac{\rho(1-\eta)^{\mu_{\text{frame}}}}{\frac{\gamma}{1-\gamma} \rho} \\
&= (1-\gamma)(1-\eta)^{\mu_{\text{frame}}}.
\end{aligned}$$

This combined with the bound (134) yields

$$\begin{aligned}
& \mathbf{v}_{t_{\text{th}}+kt_{\text{frame}}+j}(s, a) \\
&\leq \frac{\|\Delta_0\|_{\infty}}{1-\gamma} \{ (1-\gamma)(1-\eta)^{\mu_{\text{frame}}} (1-\rho)^{k-1} + \gamma(1-\rho)^{k-1} \} \\
&\leq \frac{\|\Delta_0\|_{\infty}}{1-\gamma} (\gamma + (1-\gamma)(1-\eta)^{\mu_{\text{frame}}}) (1-\rho)^{k-1} \\
&= (1-\rho)^k \frac{\|\Delta_0\|_{\infty}}{1-\gamma} = w_k, \quad (135)
\end{aligned}$$

where the last line follows from the definition of ρ (cf. (37d)). Since the above inequality holds for all state-action pair (s, a) , we conclude that

$$u_{t_{\text{th}}+kt_{\text{frame}}+j} = \|\mathbf{v}_{t_{\text{th}}+kt_{\text{frame}}+j}\|_{\infty} \leq w_k. \quad (136)$$

As a consequence, we have established the inequality (55) for all t obeying $\lfloor \frac{t-t_{\text{th}}}{t_{\text{frame}}} \rfloor = k$, which together with the induction argument completes the proof of this lemma.

G. Proof of Lemma 7

Recalling that $\mathbf{0} \leq \sum_{i=1}^t \prod_{j=i+1}^t (\mathbf{I} - \Lambda_j) \Lambda_i \mathbf{1} \leq \mathbf{1}$ (cf. (51b)), we obtain

$$\begin{aligned}
\|\mathbf{h}_{0,t}\|_{\infty} &\leq \gamma \left\| \sum_{i=1}^t \prod_{j=i+1}^t (\mathbf{I} - \Lambda_j) \Lambda_i \right\|_1 \|\tilde{\mathbf{P}} - \mathbf{P}\|_{\infty} \|\bar{\mathbf{V}}\|_{\infty} \\
&\leq \gamma \|\tilde{\mathbf{P}} - \mathbf{P}\|_{\infty} \|\bar{\mathbf{V}}\|_{\infty}. \quad (137)
\end{aligned}$$

As a result, it remains to upper bound $\|\tilde{\mathbf{P}} - \mathbf{P}\|_{\infty}$.

Suppose that $\tilde{\mathbf{P}}$ is constructed using N consecutive sample transitions. Without loss of generality, assume that these N sample transitions are the transitions between the following $N+1$ samples

$$(s_0, a_0), (s_1, a_1), (s_2, a_2), \dots, (s_N, a_N).$$

Then the (s, a) -th row of $\tilde{\mathbf{P}}$ — denoted by $\tilde{\mathbf{P}}(s, a)$ — is given by

$$\begin{aligned}
\tilde{\mathbf{P}}(s, a) &= \frac{1}{K_N(s, a)} \sum_{i=0}^{N-1} \mathbf{P}_{i+1}(s, a) \bar{\mathbf{V}} \mathbb{1}\{(s_i, a_i) = (s, a)\} \\
&= \frac{1}{K_N(s, a)} \sum_{i=1}^{K_N(s, a)} \mathbf{P}_{t_{i+1}}(s, a) \bar{\mathbf{V}}, \quad (138)
\end{aligned}$$

where \mathbf{P}_i is defined in (35), and $\mathbf{P}_i(s, a)$ denotes its (s, a) -th row. Here, $K_N(s, a)$ denotes the total number of visits to (s, a) during the first N time instances (cf. (113)), and $t_k := t_k(s, a)$ denotes the time stamp when the trajectory visits (s, a) for the k -th time (cf. (112)).

In view of our derivation for (117), the state transitions happening at time t_1, t_2, \dots, t_k are independent for any given integer $k > 0$. This together with the Hoeffding inequality implies that

$$\begin{aligned}
& \mathbb{P} \left\{ \frac{1}{k} \left| \sum_{i=1}^k (\mathbf{P}_{t_{i+1}}(s, a) - \mathbf{P}(s, a)) \bar{\mathbf{V}} \right| \geq \tau \right\} \\
&\leq 2 \exp \left\{ -\frac{k\tau^2}{2\|\bar{\mathbf{V}}\|_{\infty}^2} \right\}. \quad (139)
\end{aligned}$$

Consequently, with probability at least $1 - \frac{\delta}{|S||A|}$ one has

$$\begin{aligned}
& \left| \frac{1}{k} \sum_{i=1}^k (\mathbf{P}_{t_{i+1}}(s, a) - \mathbf{P}(s, a)) \bar{\mathbf{V}} \right| \\
&\leq \sqrt{\frac{2 \log \left(\frac{2N|S||A|}{\delta} \right)}{k}} \|\bar{\mathbf{V}}\|_{\infty}, \quad 1 \leq k \leq N.
\end{aligned}$$

Recognizing the simple bound $K_N(s, a) \leq N$, the above inequality holds for each state-action pair (s, a) when k is replaced by $K_N(s, a)$. Then, applying the union bound over all $(s, a) \in \mathcal{S} \times \mathcal{A}$, we obtain

$$\|\tilde{\mathbf{P}} - \mathbf{P}\|_{\infty} \leq \max_{(s, a) \in \mathcal{S} \times \mathcal{A}} \sqrt{\frac{2 \log \left(\frac{2N|S||A|}{\delta} \right)}{K_N(s, a)}} \|\bar{\mathbf{V}}\|_{\infty} \quad (140)$$

with probability at least $1 - \delta$.

In addition, for any $N \geq t_{\text{frame}}$, Lemma 8 guarantees that with probability $1 - 2\delta$, each state-action pair (s, a) is visited

at least $N\mu_{\min}/2$ times, namely, $K_N(s, a) \geq \frac{1}{2}N\mu_{\min}$ for all (s, a) . This combined with (141) yields

$$\begin{aligned} & \|(\tilde{\mathbf{P}} - \mathbf{P})\bar{\mathbf{V}}\|_{\infty} \\ & \leq \sqrt{\frac{4 \log\left(\frac{2N|\mathcal{S}||\mathcal{A}|}{\delta}\right)}{N\mu_{\min}}} \|\bar{\mathbf{V}}\|_{\infty} \\ & \leq \sqrt{\frac{4 \log\left(\frac{2N|\mathcal{S}||\mathcal{A}|}{\delta}\right)}{N\mu_{\min}}} (\|\bar{\mathbf{V}} - \mathbf{V}^*\|_{\infty} + \|\mathbf{V}^*\|_{\infty}) \\ & \leq \sqrt{\frac{4 \log\left(\frac{2N|\mathcal{S}||\mathcal{A}|}{\delta}\right)}{N\mu_{\min}}} \|\bar{\mathbf{V}} - \mathbf{V}^*\|_{\infty} \\ & \quad + \frac{1}{1-\gamma} \sqrt{\frac{4 \log\left(\frac{2N|\mathcal{S}||\mathcal{A}|}{\delta}\right)}{N\mu_{\min}}} \end{aligned} \quad (141)$$

with probability at least $1 - 3\delta$, where the second inequality follows from the triangle inequality, and the last inequality follows from $\|\mathbf{V}^*\|_{\infty} \leq \frac{1}{1-\gamma}$. Putting this together with (137) concludes the proof.

H. Proof of Lemma 5

For notational convenience, set $t_l := t_{\text{cover}}l$, and define

$$\mathcal{H}_l := \left\{ \exists (s, a) \in \mathcal{S} \times \mathcal{A} \text{ that is not visited within } (t_l, t_{l+1}] \right\}$$

for any integer $l \geq 0$. In view of the definition of t_{cover} , we see that for any given $(s', a') \in \mathcal{S} \times \mathcal{A}$,

$$\mathbb{P}\{\mathcal{H}_l \mid (s_{t_l}, a_{t_l}) = (s', a')\} \leq \frac{1}{2}. \quad (142)$$

Consequently, for any integer $L > 0$, one can invoke the Markovian property to obtain

$$\begin{aligned} \mathbb{P}\{\mathcal{H}_1 \cap \dots \cap \mathcal{H}_L\} &= \mathbb{P}\{\mathcal{H}_1 \cap \dots \cap \mathcal{H}_{L-1}\} \\ & \quad \cdot \mathbb{P}\{\mathcal{H}_L \mid \mathcal{H}_1 \cap \dots \cap \mathcal{H}_{L-1}\} \\ &= \mathbb{P}\{\mathcal{H}_1 \cap \dots \cap \mathcal{H}_{L-1}\} \sum_{s', a'} \left[\mathbb{P}\{\mathcal{H}_L \mid (s_{t_L}, a_{t_L}) = (s', a')\} \right. \\ & \quad \cdot \mathbb{P}\{(s_{t_L}, a_{t_L}) = (s', a') \mid \mathcal{H}_1 \cap \dots \cap \mathcal{H}_{L-1}\} \left. \right] \\ &\leq \frac{1}{2} \mathbb{P}\{\mathcal{H}_1 \cap \dots \cap \mathcal{H}_{L-1}\} \\ & \quad \cdot \sum_{s', a'} \mathbb{P}\{(s_{t_L}, a_{t_L}) = (s', a') \mid \mathcal{H}_1 \cap \dots \cap \mathcal{H}_{L-1}\} \\ &= \frac{1}{2} \mathbb{P}\{\mathcal{H}_1 \cap \dots \cap \mathcal{H}_{L-1}\}, \end{aligned}$$

where the inequality follows from (142). Repeating this derivation recursively, we deduce that

$$\mathbb{P}\{\mathcal{H}_1 \cap \dots \cap \mathcal{H}_L\} \leq \frac{1}{2^L}.$$

This tells us that

$$\begin{aligned} & \mathbb{P}\{\exists (s, a) \in \mathcal{S} \times \mathcal{A} \text{ that is not visited between } (0, t_{\text{cover}, \text{all}}]\} \\ & \leq \mathbb{P}\left\{\mathcal{H}_1 \cap \dots \cap \mathcal{H}_{\log_2 \frac{T}{\delta}}\right\} \leq \frac{1}{2^{\log_2 \frac{T}{\delta}}} = \frac{\delta}{T}, \end{aligned}$$

which in turn establishes the advertised result by applying the union bound.

ACKNOWLEDGMENT

The authors thank Shicong Cen, Chen Cheng, and Cong Ma for numerous discussions about reinforcement learning.

REFERENCES

- [1] G. Li, Y. Wei, Y. Chi, Y. Gu, and Y. Chen, "Sample complexity of asynchronous Q -learning: Sharper analysis and variance reduction," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 33, 2020, pp. 7031–7043.
- [2] G. Qu and A. Wierman, "Finite-time analysis of asynchronous stochastic approximation and Q -learning," in *Proc. Conf. Learn. Theory*, 2020, pp. 3185–3205.
- [3] C. J. C. H. Watkins and P. Dayan, " Q -learning," *Mach. Learn.*, vol. 8, nos. 3–4, pp. 279–292, 1992.
- [4] V. Mnih *et al.*, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [5] G. Dulac-Arnold, D. Mankowitz, and T. Hester, "Challenges of real-world reinforcement learning," 2019, *arXiv:1904.12901*. [Online]. Available: <http://arxiv.org/abs/1904.12901>
- [6] J. N. Tsitsiklis, "Asynchronous stochastic approximation and Q -learning," *Mach. Learn.*, vol. 16, no. 16, pp. 185–202, Sep. 1994.
- [7] T. Jaakkola, M. I. Jordan, and S. P. Singh, "Convergence of stochastic iterative dynamic programming algorithms," in *Proc. Adv. Neural Inf. Process. Syst.*, 1994, pp. 703–710.
- [8] M. J. Wainwright, "Stochastic approximation with cone-contractive operators: Sharp ℓ_{∞} -bounds for Q -learning," 2019, *arXiv:1905.06265*. [Online]. Available: <https://arxiv.org/abs/1905.06265>
- [9] E. Even-Dar and Y. Mansour, "Learning rates for Q -learning," *J. Mach. Learn. Res.*, vol. 5, pp. 1–25, Dec. 2004.
- [10] C. L. Beck and R. Srikant, "Error bounds for constant step-size Q -learning," *Syst. Control Lett.*, vol. 61, no. 12, pp. 1203–1208, 2012.
- [11] M. G. Azar, R. Munos, M. Ghavamzadeh, and H. Kappen, "Reinforcement learning with a near optimal rate of convergence," INRIA, Tech. Rep. INRIA-00636615v2, 2011.
- [12] M. J. Wainwright, "Variance-reduced Q -learning is minimax optimal," 2019, *arXiv:1906.04697*. [Online]. Available: <https://arxiv.org/abs/1906.04697>
- [13] R. S. Sutton, "Learning to predict by the methods of temporal differences," *Mach. Learn.*, vol. 3, no. 1, pp. 9–44, 1988.
- [14] C. Szepesvári, "The asymptotic convergence-rate of Q -learning," in *Proc. Adv. Neural Inf. Process. Syst.*, 1998, pp. 1064–1070.
- [15] A. B. Tsybakov and V. Zaiats, *Introduction to Nonparametric Estimation*, vol. 11. Springer, 2009.
- [16] D. P. Bertsekas, *Dynamic Programming and Optimal Control*, 4th ed. Athena Scientific, 2017.
- [17] D. Paulin, "Concentration inequalities for Markov chains by Marton couplings and spectral methods," *Electron. J. Probab.*, vol. 20, pp. 1–32, Jan. 2015.
- [18] M. Gheshlaghi Azar, R. Munos, and H. J. Kappen, "Minimax PAC bounds on the sample complexity of reinforcement learning with a generative model," *Mach. Learn.*, vol. 91, no. 3, pp. 325–349, Jun. 2013.
- [19] Z. Chen, S. Theja Maguluri, S. Shakkottai, and K. Shanmugam, "Finite-sample analysis of stochastic approximation using smooth convex envelopes," 2020, *arXiv:2002.00874*. [Online]. Available: <http://arxiv.org/abs/2002.00874>
- [20] A. Gosavi, "Boundedness of iterates in Q -learning," *Syst. Control Lett.*, vol. 55, no. 4, pp. 347–349, 2006.
- [21] J. Bhandari, D. Russo, and R. Singal, "A finite time analysis of temporal difference learning with linear function approximation," in *Proc. Conf. Learn. Theory*, 2018, pp. 1691–1692.
- [22] R. Srikant and L. Ying, "Finite-time error bounds for linear stochastic approximation and TD learning," in *Proc. Conf. Learn. Theory*, 2019, pp. 2803–2830.
- [23] K. Khamaru, A. Pananjady, F. Ruan, M. J. Wainwright, and M. I. Jordan, "Is temporal difference learning optimal? An instance-dependent analysis," 2020, *arXiv:2003.07337*. [Online]. Available: <http://arxiv.org/abs/2003.07337>
- [24] A. Sidford, M. Wang, X. Wu, L. Yang, and Y. Ye, "Near-optimal time and sample complexities for solving Markov decision processes with a generative model," in *Proc. Adv. Neural Inf. Process. Syst.*, 2018, pp. 5186–5196.
- [25] R. Johnson and T. Zhang, "Accelerating stochastic gradient descent using predictive variance reduction," in *Proc. Adv. Neural Inf. Process. Syst.*, 2013, pp. 315–323.

- [26] S. S. Du, J. Chen, L. Li, L. Xiao, and D. Zhou, "Stochastic variance reduction methods for policy evaluation," in *Proc. 34th Int. Conf. Mach. Learn.*, vol. 70, 2017, pp. 1049–1058.
- [27] A. Sidford, M. Wang, X. Wu, and Y. Ye, "Variance reduced value iteration and faster algorithms for solving Markov decision processes," in *Proc. 29th Annu. ACM-SIAM Symp. Discrete Algorithms*. Philadelphia, PA, USA: SIAM, 2018, pp. 770–787.
- [28] T. Xu, Z. Wang, Y. Zhou, and Y. Liang, "Reanalysis of variance reduced temporal difference learning," 2020, *arXiv:2001.01898*. [Online]. Available: <http://arxiv.org/abs/2001.01898>
- [29] C. J. C. H. Watkins, "Learning from delayed rewards," King's College, Cambridge, U.K., Tech. Rep., 1989.
- [30] V. S. Borkar and S. P. Meyn, "The ODE method for convergence of stochastic approximation and reinforcement learning," *SIAM J. Control Optim.*, vol. 38, no. 2, pp. 447–469, 2000.
- [31] M. J. Kearns and S. P. Singh, "Finite-sample convergence rates for Q -learning and indirect algorithms," in *Proc. Adv. Neural Inf. Process. Syst.*, 1999, pp. 996–1002.
- [32] G. Li, C. Cai, Y. Chen, Y. Gu, Y. Wei, and Y. Chi, "Tightening the dependence on horizon in the sample complexity of Q -learning," in *Proc. Int. Conf. Mach. Learn.*, 2021, pp. 6296–6306.
- [33] H. Xiong, L. Zhao, Y. Liang, and W. Zhang, "Finite-time analysis for double Q -learning," 2020, *arXiv:2009.14257*. [Online]. Available: <http://arxiv.org/abs/2009.14257>
- [34] G. Li, C. Cai, Y. Chen, Y. Gu, Y. Wei, and Y. Chi, "Is Q -learning minimax optimal? A tight sample complexity analysis," 2021, *arXiv:2102.06548*. [Online]. Available: <http://arxiv.org/abs/2102.06548>
- [35] Z. Chen, S. Zhang, T. T. Doan, J.-P. Clarke, and S. Theja Maguluri, "Finite-sample analysis of nonlinear stochastic approximation with applications in reinforcement learning," 2019, *arXiv:1905.11425*. [Online]. Available: <http://arxiv.org/abs/1905.11425>
- [36] P. Xu and Q. Gu, "A finite-time analysis of Q -learning with neural network function approximation," in *Proc. Int. Conf. Mach. Learn.*, 2020, pp. 10555–10565.
- [37] J. Fan, Z. Wang, Y. Xie, and Z. Yang, "A theoretical analysis of deep Q -learning," 2019, *arXiv:1901.00137*. [Online]. Available: <http://arxiv.org/abs/1901.00137>
- [38] S. S. Du, Y. Luo, R. Wang, and H. Zhang, "Provably efficient Q -learning with function approximation via distribution shift error checking Oracle," in *Proc. Adv. Neural Inf. Process. Syst.*, 2019, pp. 8058–8068.
- [39] S. S. Du, J. D. Lee, G. Mahajan, and R. Wang, "Agnostic Q -learning with function approximation in deterministic systems: Tight bounds on approximation error and sample complexity," 2020, *arXiv:2002.07125*. [Online]. Available: <http://arxiv.org/abs/2002.07125>
- [40] Q. Cai, Z. Yang, J. D. Lee, and Z. Wang, "Neural temporal-difference and Q -learning converges to global optima," in *Proc. Adv. Neural Inf. Process. Syst.*, 2019, pp. 11312–11322.
- [41] L. Yang and M. Wang, "Sample-optimal parametric Q -learning using linearly additive features," in *Proc. Int. Conf. Mach. Learn.*, 2019, pp. 6995–7004.
- [42] B. Weng, H. Xiong, L. Zhao, Y. Liang, and W. Zhang, "Momentum Q -learning with finite-sample convergence guarantee," 2020, *arXiv:2007.15418*. [Online]. Available: <http://arxiv.org/abs/2007.15418>
- [43] W. Weng, H. Gupta, N. He, L. Ying, and R. Srikant, "Provably-efficient double Q -learning," 2020, *arXiv:2007.05034*. [Online]. Available: <http://arxiv.org/abs/2007.05034>
- [44] D. Shah and Q. Xie, " Q -learning with nearest neighbors," in *Proc. Adv. Neural Inf. Process. Syst.*, 2018, pp. 3111–3121.
- [45] A. L. Strehl, L. Li, E. Wiewiora, J. Langford, and M. L. Littman, "PAC model-free reinforcement learning," in *Proc. 23rd Int. Conf. Mach. Learn. (ICML)*, 2006, pp. 881–888.
- [46] M. Ghavamzadeh, H. J. Kappen, M. G. Azar, and R. Munos, "Speedy Q -learning," in *Proc. Adv. Neural Inf. Process. Syst.*, 2011, pp. 2411–2419.
- [47] A. M. Devraj and S. P. Meyn, " Q -learning with uniformly bounded variance: Large discounting is not a barrier to fast learning," 2020, *arXiv:2002.10301*. [Online]. Available: <http://arxiv.org/abs/2002.10301>
- [48] C. Jin, Z. Allen-Zhu, S. Bubeck, and M. I. Jordan, "Is Q -learning provably efficient?" in *Proc. Adv. Neural Inf. Process. Syst.*, 2018, pp. 4863–4873.
- [49] Y. Wang, K. Dong, X. Chen, and L. Wang, " Q -learning with UCB exploration is sample efficient for infinite-horizon MDP," in *Proc. Int. Conf. Learn. Represent.*, 2020.
- [50] Y. Bai, T. Xie, N. Jiang, and Y.-X. Wang, "Provably efficient Q -learning with low switching cost," in *Proc. Adv. Neural Inf. Process. Syst.*, 2019, pp. 8002–8011.
- [51] G. Li, L. Shi, Y. Chen, Y. Gu, and Y. Chi, "Breaking the sample complexity barrier to regret-optimal model-free reinforcement learning," in *Proc. Neural Inf. Process. Syst. (NeurIPS)*, 2021, pp. 1–60.
- [52] T. Xu, S. Zou, and Y. Liang, "Two time-scale off-policy TD learning: Non-asymptotic analysis over Markovian samples," in *Proc. Adv. Neural Inf. Process. Syst.*, 2019, pp. 10633–10643.
- [53] H. Gupta, R. Srikant, and L. Ying, "Finite-time performance bounds and adaptive learning rate selection for two time-scale reinforcement learning," in *Proc. Adv. Neural Inf. Process. Syst.*, 2019, pp. 4706–4715.
- [54] T. Doan, S. Maguluri, and J. Romberg, "Finite-time analysis of distributed TD(0) with linear function approximation on multi-agent reinforcement learning," in *Proc. Int. Conf. Mach. Learn.*, 2019, pp. 1626–1635.
- [55] M. Kaledin, E. Moulines, A. Naumov, V. Tadic, and H.-T. Wai, "Finite time analysis of linear two-timescale stochastic approximation with Markovian noise," 2020, *arXiv:2002.01268*. [Online]. Available: <http://arxiv.org/abs/2002.01268>
- [56] G. Dalal, G. Thoppe, B. Szörényi, and S. Mannor, "Finite sample analysis of two-timescale stochastic approximation with applications to reinforcement learning," in *Proc. Conf. Learn. Theory*, 2018, pp. 1199–1233.
- [57] G. Dalal, B. Szörényi, G. Thoppe, and S. Mannor, "Finite sample analyses for TD(0) with function approximation," in *Proc. 32nd AAAI Conf. Artif. Intell.*, 2018, pp. 6144–6160.
- [58] D. Lee and N. He, "Target-based temporal difference learning," 2019, *arXiv:1904.10945*. [Online]. Available: <http://arxiv.org/abs/1904.10945>
- [59] Y. Lin, G. Qu, L. Huang, and A. Wierman, "Multi-agent reinforcement learning in time-varying networked systems," 2020, *arXiv:2006.06555*. [Online]. Available: <http://arxiv.org/abs/2006.06555>
- [60] W. Mou, C. Junchi Li, M. J. Wainwright, P. L. Bartlett, and M. I. Jordan, "On linear stochastic approximation: Fine-grained Polyak-ruppert and non-asymptotic concentration," 2020, *arXiv:2004.04719*. [Online]. Available: <http://arxiv.org/abs/2004.04719>
- [61] S. Zou, T. Xu, and Y. Liang, "Finite-sample analysis for SARSA with linear function approximation," in *Proc. Adv. Neural Inf. Process. Syst.*, 2019, pp. 8665–8675.
- [62] A. Agarwal, S. Kakade, and L. F. Yang, "Model-based reinforcement learning with a generative model is minimax optimal," 2019, *arXiv:1906.03804*. [Online]. Available: <http://arxiv.org/abs/1906.03804>
- [63] G. Li, Y. Wei, Y. Chi, Y. Gu, and Y. Chen, "Breaking the sample size barrier in model-based reinforcement learning with a generative model," 2020, *arXiv:2005.12900*. [Online]. Available: <http://arxiv.org/abs/2005.12900>
- [64] C. Dann and E. Brunskill, "Sample complexity of episodic fixed-horizon reinforcement learning," in *Proc. Adv. Neural Inf. Process. Syst.*, 2015, pp. 2818–2826.
- [65] T. Sun, Y. Sun, Y. Xu, and W. Yin, "Markov chain block coordinate descent," *Comput. Optim. Appl.*, vol. 75, pp. 35–61, Oct. 2020.
- [66] T. T. Doan, L. M. Nguyen, N. H. Pham, and J. Romberg, "Convergence rates of accelerated Markov gradient descent with applications in reinforcement learning," 2020, *arXiv:2002.02873*. [Online]. Available: <http://arxiv.org/abs/2002.02873>
- [67] P. Brémaud, *Markov Chains: Gibbs Fields, Monte Carlo Simulation, and Queues*, vol. 31. Springer, 2013.
- [68] S. Boucheron, G. Lugosi, and P. Massart, *Concentration Inequalities: A Nonasymptotic Theory of Independence*. Oxford, U.K.: Oxford Univ. Press, 2013.

Gen Li (Student Member, IEEE) received the B.E. and Ph.D. degrees in electronic engineering from Tsinghua University in 2016 and 2021, respectively, under the advice of Yuantao Gu. He is currently a Post-Doctoral Scholar advised by Yuxin Chen with the Department of Electrical and Computer Engineering, Princeton University. His recent research interests include machine learning, reinforcement learning, high-dimensional statistics, and nonconvex optimization.

Yuting Wei (Member, IEEE) received the Ph.D. degree in statistics from the University of California, Berkeley, advised by Martin Wainwright and Aditya Guntuboyina. She spent two years at Carnegie Mellon University as an Assistant Professor of statistics and one year at Stanford University as a Stein Fellow. She is currently an Assistant Professor with the Department of Statistics and Data Science, Wharton School, University of Pennsylvania. Her research interests include high-dimensional and non-parametric statistics, statistical machine learning, and reinforcement learning. She was a recipient of the 2018 Erich L. Lehmann Citation from the Berkeley Statistics Department for her Ph.D. dissertation in theoretical statistics.

Yuejie Chi (Senior Member, IEEE) received the B.E. (Hons.) degree in electrical engineering from Tsinghua University, Beijing, China, in 2007, and the M.A. and Ph.D. degrees in electrical engineering from Princeton University in 2009 and 2012, respectively. She was with The Ohio State University from 2012 to 2017. Since 2018, she has been with the Department of Electrical and Computer Engineering, Carnegie Mellon University, where she is currently a Professor and held the Inaugural Robert E. Doherty Early Career Development Professorship from 2018 to 2020. Her research interests lie in the theoretical and algorithmic foundations of data science, signal processing, and machine learning and inverse problems, with applications in sensing and societal systems, broadly defined. Among others, she was a recipient of the Presidential Early Career Award for Scientists and Engineers (PECASE), the Inaugural IEEE Signal Processing Society Early Career Technical Achievement Award for contributions to high-dimensional structured signal processing, and named a Goldsmith Lecturer by the IEEE Information Theory Society. She serves as an Associate Editor for IEEE TRANSACTIONS ON INFORMATION THEORY, IEEE TRANSACTIONS ON SIGNAL PROCESSING, and IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE.

Yuantao Gu (Senior Member, IEEE) received the B.E. degree in electronic engineering from Xi'an Jiaotong University in 1998 and the Ph.D. degree (Hons.) in electronic engineering from Tsinghua University in 2003.

He joined the Faculty of Tsinghua University in 2003 where he is currently a Full Professor with the Department of Electronic Engineering. He was a Visiting Scientist at Microsoft Research Asia from 2005 to 2006, the Research Laboratory of Electronics, Massachusetts Institute of Technology, from 2012 to 2013, and the Department of Electrical Engineering and Computer Science, University of Michigan, Ann Arbor, in 2015. His research interests include machine learning, decentralized optimization, graph signal processing, and compressed sensing.

Prof. Gu has been an Elected Member of the IEEE Machine Learning for Signal Processing (MLSP) Technical Committee since 2019 and the IEEE Signal Processing Theory and Methods (SPTM) Technical Committee since 2017. He received the Best Paper Award of IEEE Global Conference on Signal and Information Processing (GlobalSIP) in 2015, the Award for Best Presentation of Journal Paper of IEEE International Conference on Signal and Information Processing (ChinaSIP) in 2015, and the Zhang Si-Ying (CCDC) Outstanding Youth Paper Award (with his student) in 2017. He was one of the most Outstanding Reviewers for IEEE ICASSP in 2019 and received the Outstanding Editorial Board Member Award for the IEEE TRANSACTIONS ON SIGNAL PROCESSING Editorial Board in 2021. He was an Associate Editor for the IEEE TRANSACTIONS ON SIGNAL PROCESSING from 2015 to 2019 and a Handling Editor for Elsevier *Digital Signal Processing* from 2015 to 2017. He has been a Senior Area Editor for the IEEE TRANSACTIONS ON SIGNAL PROCESSING since 2019.

Yuxin Chen (Senior Member, IEEE) received the B.E. degree (Hons.) from Tsinghua University in 2008, the M.S. degree in electrical and computer engineering from The University of Texas at Austin in 2010, and the Ph.D. degree in electrical engineering and the M.S. degree in statistics from Stanford University in January 2015. He is currently an Assistant Professor with the Department of Electrical and Computer Engineering, Princeton University, Princeton, NJ, USA. Prior to joining Princeton, he was a Post-Doctoral Scholar with the Department of Statistics, Stanford University. His research interests include high-dimensional data analysis, convex and nonconvex optimization, reinforcement learning, statistical learning, statistical signal processing, and information theory. He has received the 2020 Princeton Graduate Mentoring Award, the 2019 AFOSR Young Investigator Award, the 2020 ARO Young Investigator Award, the ICCM Best Paper Award (Gold Medal), and the 2021 Princeton SEAS Junior Faculty Award, and was selected as the Finalist for the Best Paper Prize for Young Researchers in Continuous Optimization 2019.