

1 We thank the reviewers for very helpful comments. This letter addresses the **major questions** raised by the reviewers.
 2 **Learning rates.** To address the reviewers’ comments on learning rates, we will add results with *easy-to-implement*
 3 *learning rates*, without compromising sample complexities. Specifically, for some constant $c > 0$ let

$$\eta_t = \min \left\{ 1, c \exp \left(\left\lfloor \log \frac{\log t}{\widehat{\mu}_{\min,t}(1-\gamma)\gamma^2 t} \right\rfloor \right) \right\} \quad (\text{an epoch-based choice}) \quad (1)$$

4 which can be viewed as a “piecewise approximation” of the rescaled linear stepsizes $\eta_t = \min \left\{ 1, \frac{c \log t}{\widehat{\mu}_{\min,t}(1-\gamma)\gamma^2 t} \right\}$.
 5 Here, $\widehat{\mu}_{\min,t}$ is the minimum entry of certain empirical state-action visitation probability vector.¹. Clearly, this choice
 6 does *not* rely on the mixing time t_{mix} , minimum state-action occupancy probability μ_{\min} , and target accuracy ε .
 7 Encouragingly, our current theory can be easily extended to cover this easier-to-implement learning rate choice:

8 **Theorem 5** (ℓ_∞ sample complexity for achieving ε accuracy). *Consider asynchronous Q-learning with learning*
 9 *rates (1). There exists some universal constant $C > 0$ such that: for any $0 < \delta < 1$ and $0 < \varepsilon \leq \frac{1}{1-\gamma}$, one has*
 10 $\|Q_T - Q^*\|_\infty \leq \varepsilon$ with probability at least $1 - \delta$, provided that the sample size (or number of iterations) T obeys

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

11 Similarly, our theory for variance-reduced Q-learning can also be extended to a stepsize that does not depend on t_{mix} .
 12 More specifically, this requires two changes: (1) the epoch length needs to keep increasing (i.e. at the end of every
 13 epoch, run $t_{\text{epoch}} \leftarrow 2t_{\text{epoch}}$); (2) set $\eta_t = \frac{c \log t_{\text{epoch}}}{\widehat{\mu}_{\min,t}(1-\gamma)t_{\text{epoch}}}$. This can be analyzed via a similar argument.

14 *Proof of Theorem 5.* We sketch the proof for the piecewise choice (1), which follows easily from our Theorem 1.

15 1) Set $T_0 = T/2$. Given that $\eta_t \leq 1$, it is easily seen that $\|Q_{T_0} - Q^*\|_\infty \lesssim T_0$. To simplify presentation, we assume
 16 here that T is the point where η_t undergoes a change (we can easily cover general cases via epoch-based analysis).

17 2) Choose $\tilde{\varepsilon}$ s.t. $\frac{\mu_{\min}(1-\gamma)T/2}{\log(\frac{|\mathcal{S}||\mathcal{A}|T/2}{\delta}) \log \frac{T}{2}} = \frac{C}{(1-\gamma)^4 \tilde{\varepsilon}^2}$, which obeys $\tilde{\varepsilon} \leq \varepsilon$ under Condition (2). Combining the piecewise
 18 choice (1) and Condition (2) implies: $\eta_t \equiv \frac{c'}{\log(\frac{|\mathcal{S}||\mathcal{A}|T}{\delta})} \min \left\{ (1-\gamma)^4 \tilde{\varepsilon}^2, \frac{1}{t_{\text{mix}}} \right\}, \forall t \in [T_0, T]$, where c' is some constant.

19 3) With the above learning rate condition in mind, invoking Theorem 1 of our paper with initialization Q_{T_0} ensures that
 20 $\|Q_T - Q^*\|_\infty \leq \tilde{\varepsilon} \leq \varepsilon$ with probability at least $1 - \delta$, provided that the sample size condition (2) holds. \square

21 **Specific questions by Reviewer 1:** 1. “*Implementable learning rates*”: See our response above on “learning rates”.
 22 While the constant c in η_t can also be specified explicitly (by using specific constants in Bernstein inequality, etc),
 23 we caution that such a theoretical constant might be overly conservative in practice, given that our theory focuses on
 24 orderwise sample complexity bounds and does not strive to sharpen the constant. We will clarify this in the revision.

25 2. “*Main contributions*”: In comparison to the state-of-the-art [33] which unveiled tight scaling w.r.t. the important
 26 factors $\frac{1}{1-\gamma}$ and $\frac{1}{\varepsilon}$, our main focus is towards sharpening the dependency on the problem dimension $|\mathcal{S}||\mathcal{A}|$ through
 27 improving the dependency on μ_{\min} . Specifically, we improve prior sample complexity bound by a factor of $1/\mu_{\min} \geq$
 28 $|\mathcal{S}||\mathcal{A}|$. Given that $|\mathcal{S}||\mathcal{A}|$ is often enormous in practice, our theory potentially leads to a notable improvement.

29 **Specific questions by Reviewer 2:** 1. “*Dependency of stepsizes on t_{mix}* ”: See the response above on “learning rates”.

30 2. “*Bounds in expectation*”. A bound on expectation can also be extracted by (1) using the boundedness nature of the
 31 Q-update and (2) choosing δ to be sufficiently small. We will add this in the revision.

32 **Specific questions by Reviewer 3:** “*Asynchronous Q-learning vs. A3C*”: We’d like to clarify a possible source of
 33 confusion due to the different use of terminology in two different topics. The word “asynchronous” in Q-learning
 34 was often used in classical Q-learning literature (e.g. Tsitsiklis [41]) to indicate that: at every time only a single state-
 35 action pair in the Q-function is updated. This is in stark contrast to another line of recent literature on asynchronous
 36 optimization, which studies asynchronous updates of multiple CPU threads in parallel/multi-agent optimization (for
 37 instance, the A3C paper uses asynchronous SGD to simultaneously deploy/coordinate multiple CPU threads). Hence,
 38 the two settings are indeed quite different, although the Q-learning algorithm studied here might be used in any single
 39 thread to perform some component of an RL algorithm. We will clarify this in the revision to avoid confusion.

40 **Specific questions by Reviewer 5:** 1. “*Dependency of stepsizes on t_{mix}* ”: See the response above on “learning rates”.

41 2. “*Advantages of model-free methods*”: Thanks for raising concerns about our statements on model-based vs. model-
 42 free RL. We will rephrase these statements in the revision based on the reviewer’s suggestion.

¹Using concentration bounds in the supplement, we can ensure $\widehat{\mu}_{\min,t} = (1 + o(1))\mu_{\min}$ for all $t \gtrsim t_{\text{mix}}/\mu_{\min}$ (up to log factor).