

On the Convergence Rates of Federated Q-Learning across Heterogeneous Environments

Leo (Muxing) Wang
Northeastern University

wang.muxin@northeastern.edu

Pengkun Yang
Tsinghua University

yangpengkun@tsinghua.edu.cn

Lili Su
Northeastern University

l.su@northeastern.edu

Reviewed on OpenReview: <https://openreview.net/forum?id=EkLAG3gt3g>

Abstract

Large-scale multi-agent systems are often deployed across wide geographic areas, where agents interact with heterogeneous environments. There is an emerging interest in understanding the role of heterogeneity in the performance of the federated versions of classic reinforcement learning algorithms. In this paper, we study synchronous federated Q-learning, which aims to learn an optimal Q-function by having K agents average their local Q-estimates per E iterations. We provide a characterization of the error evolution, which decays to zero as the number of iterations T increases. We show that when $K(E - 1)$ is below a certain threshold, similar to the homogeneous environment settings, there is a linear speed-up concerning K . In sharp contrast, when $K(E - 1)$ is above the threshold, heterogeneous environments lead to significant performance degradation. In particular, as E increases, the convergence rate deteriorates. The slow convergence of having $E > 1$ turns out to be fundamental rather than an artifact of our analysis. We prove that, for a wide range of stepsizes, the ℓ_∞ norm of the error cannot decay faster than $\Theta_R(E/((1 - \gamma)T))$, where Θ_R only hides numerical constants and the specific choice of reward values. In addition, our experiments demonstrate that the convergence exhibits an interesting two-phase phenomenon. For any given stepsize, there is a sharp phase transition of the convergence: the error decays rapidly in the beginning yet later bounces up and stabilizes.

1 Introduction

Advancements in unmanned capabilities are rapidly transforming industries and national security by enabling fast-paced and versatile operations across domains such as advanced manufacturing (Park et al., 2019), autonomous driving (Kiran et al., 2021), and battlefields (Möhlenhof et al., 2021). Reinforcement learning (RL) – a cornerstone for unmanned capabilities – is a powerful machine learning method that aims to enable an agent to learn an optimal policy via interacting with its operating environment to solve sequential decision-making problems (Bertsekas & Tsitsiklis, 1996; Bertsekas, 2019). However, the ever-increasing complexity of the environment results in a high-dimensional state-action space, often imposing overwhelmingly high sample collection requirements on individual agents. This limited-data challenge becomes a significant hurdle that must be addressed to realize the potential of reinforcement learning.

In this paper, we study reinforcement learning within a federated learning framework (also known as Federated Reinforcement Learning (Qi et al., 2021; Jin et al., 2022; Woo et al., 2023)), wherein multiple agents independently collect samples and collaboratively train a common policy under the orchestration of a parameter server without disclosing the local data trajectories. A simple illustration can be found in Fig. 1.

When the environments of all agents are homogeneous, it has been shown that the federated version of classic reinforcement learning algorithms can significantly alleviate the data collection burden on individual agents (Woo et al., 2023; Khodadadian et al., 2022) – error bounds derived therein exhibit a linear speedup in terms of the number of agents. Moreover, by tuning the synchronization period E (i.e., the number of iterations between agent synchronization), the communication cost can be significantly reduced compared with $E = 1$ yet without significant performance degradation.

However, many large-scale multi-agent systems are often deployed across wide geographic areas, resulting in agents interacting with heterogeneous environments. For instance, connected and autonomous vehicles (CAVs) operating in various regions of a metropolitan area encounter diverse conditions such as varying traffic patterns, road infrastructure, and local regulations. Intuitively, the environmental heterogeneity may lead to misaligned learning signals across agents, potentially hinder the convergence, and degrade the generalization performance of the learned policies. Hence, the clients’ federation must be managed in a way that ensures the learned policy is robust to environmental heterogeneity.

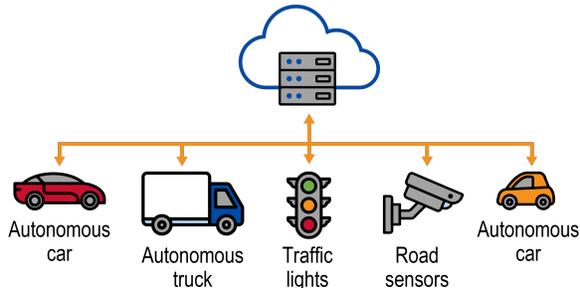


Figure 1: An illustration of a federated learning system.

There is an emerging interest in mathematically understanding the role of heterogeneity in the performance of the federated versions of classic reinforcement learning algorithms (Jin et al., 2022; Woo et al., 2023; Doan et al., 2019; Wang et al., 2023; Xie & Song, 2023) such as Q-learning, policy gradient methods, and temporal difference (TD) methods. In this paper, we study synchronous federated Q-learning (FQL) in the presence of environmental heterogeneity, which aims to learn an optimal Q-function by averaging local Q-estimates per E (where $E \geq 1$) update iterations on their local data. We leave the exploration of asynchronous Q-learning for future work. Federated Q-learning is a natural integration of FedAvg and Q-learning (Jin et al., 2022; Woo et al., 2023). The former is the most widely adopted classic federated learning algorithm (Kairouz et al., 2021; McMahan et al., 2017), and the latter is one of the most fundamental model-free reinforcement learning algorithms (Watkins & Dayan, 1992). Despite intensive study, the tight sample complexity of Q-learning in the single-agent setting was open until recently (Li et al., 2024). Similarly, the understanding of FedAvg is far from complete; a detailed discussion can be found in Section 2. A concise comparison of our work to the related work can be found in Table 1.

Contributions. Our contributions can be summarized as follows. All the asymptotic notations, e.g., \mathcal{O} and $\tilde{\mathcal{O}}$, unless otherwise specified, hide only numerical constants.

- We characterize the error evolution of synchronous Federated Q-learning, showing that it decays to zero as the number of iterations T increases. When $K(E - 1)$ is below a threshold of $\tilde{\mathcal{O}}((\kappa\epsilon)^{-1}(1 - \gamma)^{-2})$, similar to the homogeneous environment settings, there is a linear speed-up concerning K and the sample complexity is $\tilde{\mathcal{O}}\left(\frac{|\mathcal{S}||\mathcal{A}|}{K(1-\gamma)^5\epsilon^2}\right)$, matching the homogenous setting (Woo et al., 2023). Here, \mathcal{S} and \mathcal{A} are the state and action spaces, $\gamma \in (0, 1)$ denotes the discount factor, and κ is a scalar characterizing the environment heterogeneity. In sharp contrast, when $K(E - 1)$ is above the threshold, heterogeneous environments lead to significant performance degradation and results in a unique sample complexity of $\tilde{\mathcal{O}}\left(\frac{|\mathcal{S}||\mathcal{A}|\kappa E}{(1-\gamma)^3\epsilon}\right)$. Note that E is also ϵ -dependent. Hence, the sample complexity $\tilde{\mathcal{O}}\left(\frac{|\mathcal{S}||\mathcal{A}|\kappa E}{(1-\gamma)^3\epsilon}\right)$ does not contradict the lower bound of $\Omega(\epsilon^{-2})$ for the single-agent case.
- We prove that the convergence slowing down for $E > 1$ is fundamental. We show that the ℓ_∞ norm of the error cannot decay faster than $\Theta_R\left(\frac{E}{(1-\gamma)^T}\right)$, where Θ_R only hides numerical constants and the specific choice of reward values. A practical implication of this impossibility result is that, eventually, having multiple local updates (i.e., $E > 1$) ends up consuming more samples (i.e., $E \times$ more) than using $E = 1$ to reach a target accuracy.
- Our numerical results illustrate that when the environments are heterogeneous and $E > 1$, there exists a sharp phase-transition of the error convergence for not too small stepsizes: The error decays rapidly in the beginning yet later bounces up and stabilizes. In addition, provided that the phase-transition time

can be estimated, choosing different stepsizes for the two phases can lead to faster overall convergence for both constant and time-decaying stepsizes. We conjecture that this is because the error in phase 1 is mainly controlled by the initial error with impacting factor decay exponentially in time, and the error in phase 2 is dominated by the collective perturbation caused by environment heterogeneity and multiple local updates (i.e., $E > 1$).

2 Related Work

Federated Learning. Federated learning is a communication-efficient distributed machine learning approach that enables training global models without sharing raw local data (McMahan et al., 2017; Kairouz et al., 2021). Federated learning has been adopted in commercial applications that involve diverse edge devices such as autonomous vehicles (Du et al., 2020; Chen et al., 2021; Zeng et al., 2022; Posner et al., 2021; Peng et al., 2023), internet of things (Nguyen et al., 2019; Yu et al., 2020), industrial automation (Liu et al., 2020), healthcare (Yan et al., 2021; Sheller et al., 2019), and natural language processing (Yang et al., 2018; Ramaswamy et al., 2019). Multiple open-source frameworks and libraries are available such as FATE, Flower, OpenMinded-PySyft, OpenFL, TensorFlow Federated, and NVIDIA Clara.

FedAvg was proposed in the seminal work (McMahan et al., 2017), and has been one of the most widely implemented federated learning algorithms. It also has inspired many follow-up algorithms such as FedProx (Li et al., 2020b), FedNova (Wang et al., 2020), SCAFFOLD (Karimireddy et al., 2020), and adaptive federated methods (Deng et al., 2020). Despite intensive efforts, the theoretical understanding of FedAvg is far from complete. Most existing theoretical work on FedAvg overlooks the underlying data statistics at the agents, which often leads to misalignment of the pessimistic theoretical predictions and empirical success (Su et al., 2023; Pathak & Wainwright, 2020; Wang et al., 2022a;b). This theory and practice gap was studied in a recent work (Su et al., 2023) in the context of solving general non-parametric regression problems.

Reinforcement Learning. For the single-agent setup, there has been extensive research on the convergence guarantees of reinforcement learning algorithms. A recent surge of work studied non-asymptotic convergence and the corresponding sample complexity. Bhandari et al. (2018) analyzed non-asymptotic Temporal Difference (TD) learning with linear function approximation under a variety of noise conditions, including noiseless, independent noise, and Markovian noise. The results were extended to TD(λ) and Q-learning. Li et al. (2020a) investigated the sample complexity of asynchronous Q-learning with different families of learning rates. They also provided an extension of using variance reduction methods inspired by the seminal SVRG algorithm. Li et al. (2024) shows the sample complexity of Q-learning. Recall that \mathcal{A} is the set of actions. When $|\mathcal{A}| = 1$, the sample complexity of synchronous Q-learning is sharp and minimax optimal; however, when $|\mathcal{A}| \geq 2$, it was shown that synchronous Q-learning has a lower bound which is not minimax optimal.

Multi-Agent RL. Yu et al. (2022) tested multi-agent Proximal Policy Optimization in four multi-agent testbeds wherein agents fully share the parameters, and showed its competitive performance. Christiansos et al. (2021) proposed a selective parameter sharing technique, which automatically partitions agents so that they can benefit from the parameter sharing. Zhong et al. (2024) further proposed provably correct heterogeneous-agent algorithms, which allow agents to have different policy functions. The algorithms showed superior effectiveness and stability in various challenging benchmarks.

Federated Reinforcement Learning. Woo et al. (2023) provided sample complexity guarantees for both synchronous and asynchronous distributed Q-learning. They revealed that, under the same transition probability (i.e., homogeneous environment) for all agents, the convergence speed in learning the optimal Q-function can be accelerated linearly in the number of agents. They also uncovered the blessing of heterogeneity in terms of state-action exploration – a completely different notion of heterogeneity from our focus. Salgia & Chi (2025) explored the frontier of sample and communication complexities under homogeneous environments. Via variance reduction and communication quantization, they designed an algorithm that achieves order-optimal sample and communication complexities. Doan et al. (2019) investigated the distributed TD(0) with linear function approximation for a setting where multiple agents act in a shared environment and each agent has its own reward function. Khodadadian et al. (2022) studied on-policy federated TD learning, off-policy federated TD learning, and federated Q-learning of homogeneous environments and reward with Markovian noises. The sample complexity derived exhibits linear speedup with respect to the number of agents.

Work	RL Algorithm	Heterogeneity	Optimality	Lower bound	Sampling	Finite-time	Task
Wang et al. (2023)	TD(0)	✓	✗	✗	✓	✓	Pred
Xie & Song (2023)	Policy Gradient	✓	✗	✗	✓	✗	Pred, Plan
Zhang et al. (2024)	SARSA	✓	✗	✗	✓	✓	Pred, Plan
Khodadadian et al. (2022)	TD, Q-Learning	✗	✓	✗	✓	✓	Pred, Plan
Jin et al. (2022)	Q-Learning, Policy Gradient	✓	✓	✗	✗	✓	Pred, Plan
Woo et al. (2023)	Q-Learning	✗	✓	✗	✓	✓	Pred, Plan
Zheng et al. (2023)	Q-Learning	✗	✓	✗	✓	✓	Pred, Plan
Our work	Q-Learning	✓	✓	✓	✓	✓	Pred, Plan

Table 1: Comparison of various works in the context of FRL. Pred and Plan stand for prediction (policy evaluation) and planning (policy optimization), respectively.

Heterogeneous environments were considered in Jin et al. (2022); Wang et al. (2023); Xie & Song (2023); Zhang et al. (2023). Jin et al. (2022) studied federated Q-learning and policy gradient methods assuming known transition probabilities. To address heterogeneity in both environments and rewards, Wang et al. (2023) proposed FedTD(0) with linear function approximation. They proved that, in a low-heterogeneity regime, there is a linear convergence speedup in the number of agents. Xie & Song (2023) used KL-divergence to penalize the deviation of local update from the global policy, and proved that the local update is beneficial for global convergence. Zhang et al. (2024) proposed FedSARSA using the classic on-policy RL algorithm SARSA with linear function approximation. They theoretically proved that the algorithm can converge to a near-optimal solution. Neither Xie & Song (2023) nor Zhang et al. (2024) characterized sample complexity.

Technical comparisons with Woo et al. (2023); Zhang et al. (2024); Wang et al. (2023).

Zhang et al. (2024) and Wang et al. (2023) examined federated versions of TD learning and SARSA, whereas our paper studied federated Q-learning, which offers distinct theoretical and practical advantages for optimal policy learning. Specifically, the upper bound in Zhang et al. (2024) and Wang et al. (2023) do not indicate how fundamentally the convergence rates are impacted by the heterogeneity κ and synchronization period E . In addition, their upper bounds do not decay to 0 as $T \rightarrow \infty$. Our upper bound converges to 0 as $T \rightarrow \infty$. Furthermore, we derived a lower bound on the convergence rates, showing the fundamental limitation of multiple local updates (i.e., $E > 1$) in the presence of environmental heterogeneity. To the best of our knowledge, this is the first result of its kind.

While our analysis of Theorem 1 builds upon the roadmap established by Woo et al. (2023), adapting their analysis to our setting introduces significant challenges. In Woo et al. (2023), agents operate in homogeneous environments, i.e., each of the K agents shares the same transition distributions. This homogeneity allows the concentration bound on the difference between the true transition distribution and sampled estimates to become arbitrarily small as the number of samples increases. However, in our setting, each agent has its own environment with a distinct transition distribution. This heterogeneity introduces a perturbation term in the error upper bound that does not decrease with additional samples. Additionally, when $E > 1$ and $\kappa > 0$, the term involving $\kappa(E - 1)$ in the upper bound becomes the dominant term, resulting in a unique sample complexity. Further technical details and implications of these adjustments are provided in Corollary 2.

3 Preliminary on Q-Learning

Markov decision process. A Markov decision process (MDP) is defined by the tuple $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \gamma, R \rangle$, where \mathcal{S} represents the set of states, \mathcal{A} represents the set of actions, the transition probability $\mathcal{P} : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ provides the probability distribution over the next states given a current state s and action a , the reward function $R : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ assigns a reward value to each state-action pair, and the discount factor $\gamma \in (0, 1)$ models

the preference for immediate rewards over future rewards. It is worth noting that $\mathcal{P} = \{P(\cdot | s, a)\}_{s \in \mathcal{S}, a \in \mathcal{A}}$ is a collection of $|\mathcal{S}| \times |\mathcal{A}|$ probability distributions over \mathcal{S} , one for each state-action pair (s, a) .

Policy, value function, Q-function, and optimality. A policy π specifies the action-selection strategy and is defined by the mapping $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$, where $\pi(a | s)$ denotes the probability of choosing action a when in state s . For a given policy π , the value function $V^\pi : \mathcal{S} \rightarrow \mathbb{R}$ measures the expected total discounted reward starting from state s :

$$V^\pi(s) = \mathbb{E}_{a_t \sim \pi(\cdot | s_t), s_{t+1} \sim P(\cdot | s_t, a_t)} \left[\sum_t \gamma^t R(s_t, a_t) \mid s_0 = s \right], \quad \forall s \in \mathcal{S}.$$

The Q-function (a.k.a. state-action value function), $Q^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, evaluates the expected total discounted reward from taking action a in state s and then following policy π :

$$Q^\pi(s, a) = R(s, a) + \mathbb{E}_{a_t \sim \pi(\cdot | s_t), s_{t+1} \sim P(\cdot | s_t, a_t)} \left[\sum_t \gamma^t R(s_t, a_t) \mid s_0 = s, a_0 = a \right], \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}.$$

An optimal policy π^* is one that maximizes the value function for every state, that is $\forall s \in \mathcal{S}, V^{\pi^*}(s) \geq V^\pi(s)$ for any other $\pi \neq \pi^*$. Such a policy ensures the highest possible cumulative reward. The optimal value function V^* (shorthand for V^{π^*}) and the optimal Q-function Q^* (shorthand for Q^{π^*}) are defined under the optimal policy π^* .

The Bellman optimality equation for the value function and state-value function are:

$$\begin{aligned} V^*(s) &= \max_a [R(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s' | s, a) V^*(s')] \\ Q^*(s, a) &= R(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s' | s, a) \max_{a' \in \mathcal{A}} Q^*(s', a'). \end{aligned}$$

Q-learning. Q-learning (Watkins & Dayan, 1992) is a model-free reinforcement learning algorithm that aims to learn the value of actions of all states by updating Q-values through iterative exploration of the environment, ultimately converging to the optimal state-action function. Based on the Bellman optimality equation for the state-action function, the update rule for Q-Learning is formulated as:

$$Q_{t+1}(s, a) = (1 - \lambda)Q_t(s, a) + \lambda[R(s, a) + \gamma \max_{a' \in \mathcal{A}} Q_t(s', a')], \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A},$$

where $s' \sim P(\cdot | s, a)$, and λ is the stepsize.

4 Federated Q-learning

The federated learning system consists of one parameter server (PS) and K agents. The K agents are deployed in possibly heterogeneous yet independent environments. The K agents are modeled as MDPs with $\mathcal{M}_k = \langle \mathcal{S}, \mathcal{A}, \mathcal{P}^k, \gamma, R \rangle$ for $k = 1, \dots, K$, where $\mathcal{P}^k = \{P^k(\cdot | s, a)\}_{s \in \mathcal{S}, a \in \mathcal{A}}$ is the collection of probability distributions. In the synchronous setting, each agent k has access to a generative model. At each iteration t , it generates a new state sample $s_t^k(s, a) \sim P^k(\cdot | s, a)$ for each (s, a) , i.e., $\mathbb{P}\{s_t^k(s, a) = s'\} = P^k(s' | s, a)$ for all $s' \in \mathcal{S}$, independently across state-action pairs (s, a) . For each (s, a) , the global environment $\bar{P}(\cdot | s, a)$ (Jin et al., 2022) is defined as

$$\bar{P}(s' | s, a) = \frac{1}{K} \sum_{k=1}^K P^k(s' | s, a), \quad \forall s' \tag{1}$$

with the corresponding global MDP defined as $\mathcal{M}_g = \langle \mathcal{S}, \mathcal{A}, \bar{\mathcal{P}}, \gamma, R \rangle$. Define transition heterogeneity κ as

$$\sup_{k, s, a} \|\bar{P}(\cdot | s, a) - P^k(\cdot | s, a)\|_\infty := \kappa. \tag{2}$$

Let Q^* denote the optimal Q-function of the global MDP. By the Bellman optimality equation, we have,

$$Q^*(s, a) = R(s, a) + \gamma \sum_{s' \in \mathcal{S}} \bar{P}(s' | s, a) V^*(s'), \quad \forall (s, a) \quad (3)$$

where $V^*(s) = \max_{a \in \mathcal{A}} Q^*(s, a)$ is the optimal value function.

The goal of federated Q-learning is to have the K agents collaboratively learn Q^* . We consider synchronous federated Q-learning, which is a natural integration of FedAvg and Q-learning (Woo et al., 2023; Jin et al., 2022) – described in Algorithm 1. Every agent initializes its local Q^k estimate as Q_0 and performs standard synchronous Q-learning based on the locally collected samples $s_t^k(s, a)$. Whenever $t + 1 \bmod E = 0$, through the parameter server, the K agents average their local estimate of Q ; that is, all agents report their $Q_{t+\frac{1}{2}}^k$ to the parameter server, which computes the average and sends back to agents.

5 Main Results

With a little abuse of notation, let the matrix $P^k \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}| \times |\mathcal{S}|}$ represent the transition kernel of the MDP of agent k with the (s, a) -th row being $P^k(\cdot | s, a) \in \mathbb{R}^{|\mathcal{S}|}$ – the transition probability of the state-action pair (s, a) . For ease of exposition, we write $P^k(\cdot | s, a) = P^k(s, a)$ as the state transition probability at the state-action pair (s, a) when its meaning is clear from the context.

5.1 Main Convergence Results.

Let $\tilde{P}_t^k \in \{0, 1\}^{|\mathcal{S}| \times |\mathcal{A}| \times |\mathcal{S}|}$ denote the local empirical transition matrix at the t -th iteration, defined as

$$\tilde{P}_t^k(s' | s, a) = \mathbf{1}\{s' = s_t^k(s, a)\}.$$

Denoting $\tilde{P}_i^k V^* \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}| \times 1}$ with the (s, a) -th entry as $\tilde{P}_i^k(s, a) V^* = \sum_{s' \in \mathcal{S}} \tilde{P}_i^k(s' | s, a) V^*(s')$. Let $\bar{Q}_{t+1} := \frac{1}{K} \sum_{k=1}^K Q_{t+\frac{1}{2}}^k$. From lines 6, 8, and 10 of Algorithm 1, it follows that

$$\bar{Q}_{t+1} = \frac{1}{K} \sum_{k=1}^K \left((1 - \lambda) Q_t^k + \lambda (R + \gamma \tilde{P}_t^k V_t^k) \right), \quad (4)$$

where $V_t^k(s) := \max_{a \in \mathcal{A}} Q_t^k(s, a)$ for all $s \in \mathcal{S}$. Define

$$\Delta_{t+1} := Q^* - \bar{Q}_{t+1}, \quad \text{and} \quad \Delta_0 := Q^* - Q_0. \quad (5)$$

The error iteration Δ_t is captured in the following lemma.

Lemma 1 (Error iteration). *For any $t \geq 0$,*

$$\begin{aligned} \Delta_{t+1} &= (1 - \lambda)^{t+1} \Delta_0 + \gamma \lambda \sum_{i=0}^t (1 - \lambda)^{t-i} \frac{1}{K} \sum_{k=1}^K (\bar{P} - \tilde{P}_i^k) V^* \\ &\quad + \gamma \lambda \sum_{i=0}^t (1 - \lambda)^{t-i} \frac{1}{K} \sum_{k=1}^K \tilde{P}_i^k (V^* - V_i^k). \end{aligned} \quad (6)$$

Algorithm 1 Synchronous Federated Q-Learning

Inputs: discount factor γ , E , total iteration T , stepsize λ , initial estimate Q_0

```

1: for  $k \in [K]$  do
2:    $Q_0^k = Q_0$ 
3: end for
4: for  $t = 0$  to  $T - 1$  do
5:   for  $k \in [K]$  and  $(s, a) \in \mathcal{S} \times \mathcal{A}$  do
6:      $Q_{t+\frac{1}{2}}^k(s, a) = (1 - \lambda) Q_t^k(s, a) +$ 
        $\lambda (R(s, a) + \gamma \max_{a' \in \mathcal{A}} Q_t^k(s_t^k(s, a), a'))$ 
7:     if  $(t + 1) \bmod E = 0$  then
8:        $Q_{t+1}^k = \frac{1}{K} \sum_{k=1}^K Q_{t+\frac{1}{2}}^k$ 
9:     else
10:       $Q_{t+1}^k = Q_{t+\frac{1}{2}}^k$ 
11:    end if
12:  end for
13: end for
14: return  $Q_T = \frac{1}{K} \sum_{k=1}^K Q_T^k$ 

```

To show the convergence of $\|\Delta_{t+1}\|_\infty$, we bound each of the three terms in the right-hand side of (6). The following lemma is a coarse error upper bound.

Lemma 2. *Choosing $R(s, a) \in [0, 1]$ for each state-action pair (s, a) , and choose $0 \leq Q_0(s, a) \leq \frac{1}{1-\gamma}$ for any $(s, a) \in \mathcal{S} \times \mathcal{A}$, then $0 \leq Q_t^k(s, a) \leq \frac{1}{1-\gamma}$, $0 \leq Q^*(s, a) \leq \frac{1}{1-\gamma}$,*

$$\|Q^* - Q_t^k\|_\infty \leq \frac{1}{1-\gamma}, \quad \text{and} \quad \|V^* - V_t^k\|_\infty \leq \frac{1}{1-\gamma}, \quad \forall t \geq 0, \text{ and } k \in [K]. \quad (7)$$

With the choice of Q_0 in Lemma 2, the first term in (6) can be bounded as $\|(1-\lambda)^{t+1}\Delta_0\|_\infty \leq (1-\lambda)^{t+1}\frac{1}{1-\gamma}$. In addition, as detailed in the proof of Lemma 4 and Theorem 1, the boundedness in Lemma 2 enables us to bound the second term in (6) via invoking the Hoeffding's inequality. It remains to bound the third term in (6), for which we follow the analysis roadmap of Woo et al. (2023) by a two-step procedure that is described in Lemma 3 and Lemma 4. Let

$$\Delta_t^k = Q^* - Q_t^k, \quad \text{and} \quad \chi(t) = t - (t \bmod E), \quad (8)$$

i.e., Δ_t^k is the local error of agent k , and $\chi(t)$ is the most recent synchronization iteration of t .

Lemma 3. *If $t \bmod E = 0$, then $\left\| \frac{1}{K} \sum_{k=1}^K \tilde{P}_t^k(V^* - V_t) \right\|_\infty \leq \|\Delta_t\|_\infty$. Otherwise,*

$$\begin{aligned} \left\| \frac{1}{K} \sum_{k=1}^K \tilde{P}_t^k(V^* - V_t^k) \right\|_\infty &\leq \|\Delta_{\chi(t)}\|_\infty + 2\lambda \frac{1}{K} \sum_{k=1}^K \sum_{t'=\chi(t)}^{t-1} \|\Delta_{t'}^k\|_\infty \\ &\quad + \gamma\lambda \frac{1}{K} \sum_{k=1}^K \max_{s,a} \left| \sum_{t'=\chi(t)}^{t-1} \left(\tilde{P}_{t'}^k(s, a) - \bar{P}(s, a) \right) V^* \right|. \end{aligned}$$

where we use the convention that $\sum_{t'=\chi(t)}^{\chi(t)-1} \|\Delta_{t'}^k\|_\infty = 0$.

Lemma 4. *Choose $\lambda \leq \frac{1}{E}$. For any $\delta \in (0, 1)$, with probability at least $(1-\delta)$,*

$$\|\Delta_i^k\|_\infty \leq \|\Delta_{\chi(i)}\|_\infty + \frac{3\gamma}{1-\gamma} \lambda(E-1)\kappa + \frac{3\gamma}{1-\gamma} \sqrt{\lambda \log \frac{|\mathcal{S}||\mathcal{A}|KT}{\delta}}, \quad \forall i \leq T, k \in [K]. \quad (9)$$

To bound the ℓ_∞ norm of the third term in (6), we first invoke Lemma 3, followed by Lemma 4. It is worth noting that directly applying Lemma 4 can also lead to a valid error bound yet the resulting bound will not decay as T increases with proper choice of stepsizes.

Both Lemma 3 and Lemma 4 are non-trivial adaptations of the approach in Woo et al. (2023) due to the absence of a common optimal action for any given state in heterogeneous environments. Moreover, in the homogeneous setting, each agent draws samples from the same true transition distribution, allowing concentration inequalities to bound the discrepancy between the true distribution and sampled estimates. However, this line of reasoning does not go through in the presence of environmental heterogeneity. When $\kappa > 0$, each of the K agents has its own transition distribution, and the discrepancy is captured by the environmental heterogeneity parameter κ .

Theorem 1 (Convergence). *Choose $E-1 \leq \frac{1}{\lambda} \min\{\frac{1-\gamma}{4\gamma}, \frac{1}{K}\}$ and $\lambda \leq \frac{1}{E}$. For any $\delta \in (0, \frac{1}{3})$, with probability at least $1-3\delta$, it holds that*

$$\|\Delta_T\|_\infty \leq \frac{4}{(1-\gamma)^2} \exp\left\{-\frac{1}{2}\sqrt{(1-\gamma)\lambda T}\right\} + \frac{14\gamma^2}{(1-\gamma)^2} \lambda(E-1)\kappa + \frac{16}{(1-\gamma)^2} \sqrt{\frac{\lambda}{K} \log \frac{|\mathcal{S}||\mathcal{A}|KT}{\delta}}.$$

The first term of Theorem 1 is the standard error bound in the absence of environmental heterogeneity and sampling noises. The second term arises from environmental heterogeneity. It is clear that when $E=1$, the environmental heterogeneity does not negatively impact the convergence. The last term results from the randomness in sampling.

Remark 1 (Eventual zero error). It is common to choose the stepsize λ based on the time horizon T . Let $\lambda = g(T)$ be a non-increasing function of T , and other parameters be fixed with respect to T . As long as $\lambda = g(T)$ decay in T , terms 2 and 3 in Theorem 1 will go to 0 as T increases. In addition, when $\lambda = \omega(1/T)$, the first term will decay to 0. Conversely, the convergence bounds in Zhang et al. (2024) and Wang et al. (2023) do not decay to 0.

There is a tradeoff in the convergence rates of the first term and the remaining terms – the slower λ decay in T leads to faster decay in the first term but slower in the remaining terms. Forcing these terms to decay around the same speed leads to slow overall convergence. Corollary 1 follows immediately from Theorem 1 via carefully choosing λ to balance the decay rates of different terms.

Corollary 1. Choose $(E - 1) \leq \min \frac{1}{\lambda} \{ \frac{1-\gamma}{4\gamma}, \frac{1}{K} \}$, and $\lambda = \frac{4 \log^2(TK)}{T(1-\gamma)}$. Let $T \geq E$. For any $\delta \in (0, \frac{1}{3})$, with probability at least $1 - 3\delta$,

$$\|\Delta_T\|_\infty \leq \frac{4}{(1-\gamma)^2 TK} + \frac{32}{(1-\gamma)^{2.5}} \frac{\log(TK)}{\sqrt{TK}} \sqrt{\log \frac{|\mathcal{S}||\mathcal{A}|TK}{\delta}} + \frac{56 \log^2(TK)}{(1-\gamma)^3} \frac{E-1}{T} \kappa.$$

Remark 2. Intuitively, both terms 1 and 2 decay as if there are TK iterations. In fact, the decay rate of the sampling noises in Corollary 1, with respect to TK , is the minimax optimal up to polylog factors (Vershynin, 2018). The decay of the third term is controlled by environmental heterogeneity when $E > 1$. In sharp contrast to the homogeneous settings, larger E significantly slows down the convergence of this term. We show in the next subsection that this slow convergence is fundamental.

Remark 3 (Communication cost and convergence). From Corollary 1, by choosing $E = \tilde{\Theta}(\sqrt{T})$, and other parameters are fixed with respect to T , we can reach the same error bound of $\tilde{\mathcal{O}}(1/\sqrt{T})$ with communication cost of $\tilde{\mathcal{O}}(\sqrt{T})$, which is better than $\tilde{\mathcal{O}}(T)$.

Corollary 2. Choose $E - 1 \leq \frac{1}{\lambda} \min \{ \frac{1-\gamma}{4\gamma}, \frac{1}{K} \}$ and $\lambda \leq \frac{1}{E}$, and define $x_1 = \frac{4096 \log \frac{|\mathcal{S}||\mathcal{A}|K}{\delta} \log^2(\frac{(1-\gamma)^2 \epsilon}{8})}{K(1-\gamma)^5 \epsilon^2}$, $x_2 = \frac{168\kappa(E-1)\gamma^2}{(1-\gamma)^3 \epsilon} \log^2(\frac{(1-\gamma)^2 \epsilon}{12})$, and $x_3 = \frac{9216}{K(1-\gamma)^5 \epsilon^2} \log \frac{|\mathcal{S}||\mathcal{A}|K}{\delta} \log^2(\frac{(1-\gamma)^2 \epsilon}{12})$,

- When $\kappa = 0$ or $E = 1$, for any $\delta \in (0, \frac{1}{3})$, with probability at least $1 - 3\delta$, it holds that

$$\|\Delta_T\|_\infty \leq \epsilon,$$

when $T \geq \exp\{-W_{-1}(-\frac{1}{x_1})\}$, and $\lambda = \frac{\epsilon^2(1-\gamma)^4 K}{2304 \log \frac{|\mathcal{S}||\mathcal{A}|KT}{\delta}}$, where W_{-1} is the Lambert W function. The resulting sample complexity is $\tilde{\mathcal{O}}\left(\frac{|\mathcal{S}||\mathcal{A}|}{K(1-\gamma)^5 \epsilon^2}\right)$.

- When $\kappa > 0$ and $E > 1$,

- If $K(E - 1) \geq \frac{55 \log \frac{|\mathcal{S}||\mathcal{A}|K}{\delta}}{\kappa\gamma^2 \epsilon(1-\gamma)^2}$, for any $\delta \in (0, \frac{1}{3})$, with probability at least $1 - 3\delta$, it holds that

$$\|\Delta_T\|_\infty \leq \epsilon,$$

when $T \geq \exp\{-W_{-1}(-\frac{1}{x_2})\}$ and $\lambda = \frac{\epsilon(1-\gamma)^2}{42\kappa(E-1)\gamma^2}$. The sample complexity is $\tilde{\mathcal{O}}\left(\frac{\kappa|\mathcal{S}||\mathcal{A}|E}{(1-\gamma)^3 \epsilon}\right)$.

- If $K(E - 1) \leq \frac{55 \log \frac{|\mathcal{S}||\mathcal{A}|K}{\delta}}{\kappa\gamma^2 \epsilon(1-\gamma)^2}$, for any $\delta \in (0, \frac{1}{3})$, with probability at least $1 - 3\delta$, it holds that

$$\|\Delta_T\|_\infty \leq \epsilon,$$

when $T \geq \exp\{-W_{-1}(-\frac{1}{x_3})\}$ and $\lambda = \frac{\epsilon^2(1-\gamma)^4 K}{2304 \log \frac{|\mathcal{S}||\mathcal{A}|KT}{\delta}}$. The sample complexity is $\tilde{\mathcal{O}}\left(\frac{|\mathcal{S}||\mathcal{A}|}{K(1-\gamma)^5 \epsilon^2}\right)$.

Remark 4 (Sample complexity on K and E and conditional linear speedup.). From Corollary 2, we can conclude that when the setting is homogeneous (i.e., $\kappa = 0$) or the agents communicate every step (i.e., $E = 1$), the sample complexity $\tilde{\mathcal{O}}\left(\frac{|\mathcal{S}||\mathcal{A}|}{K(1-\gamma)^5 \epsilon^2}\right)$ matches the one in Woo et al. (2023). On the other hand, when the setting is heterogeneous (i.e., $\kappa > 0$) and $E > 1$, it is evident that if the total computation steps

per synchronization are sufficiently small, i.e., $K(E-1) \leq \tilde{O}((\kappa\epsilon)^{-1}(1-\gamma)^{-2})$, the sample complexity also matches the one in the homogeneous setting, where there is a linear speedup. Otherwise, the sample complexity $\tilde{O}\left(\frac{\kappa|\mathcal{S}||\mathcal{A}|E}{(1-\gamma)^3\epsilon}\right)$ increases with E , meaning that multiple local rounds only consume more samples (i.e., E -times more) samples without achieving linear speedup.

5.2 On the Fundamentals of Convergence Slowdown for $E > 1$ in Heterogeneous Environments.

Theorem 2. *Let $Q_0 = \mathbf{0}$. For any even $K \geq 2$, there exists a collection of $\{(\mathcal{S}, \mathcal{A}, \mathcal{P}^k, R, \gamma) : k \in [K]\}$ where $|\mathcal{S}| = 2$, $|\mathcal{A}| = 1$, and $R := \begin{bmatrix} r_1 \\ r_2 \end{bmatrix}$ fixed for all MDPs in the collection, such that, for $E \geq 2$ and time-invariant stepsize $\lambda \leq \frac{1}{1+\gamma}$,*

$$\|\Delta_T\|_\infty \geq c_R \frac{E}{(1-\gamma)T},$$

when $T/E \in \mathbb{N}$ and $T \geq E \cdot \max\left\{\exp\left\{\frac{4E(\gamma+2)}{(1+\gamma)\gamma^2(E-1)}\right\}, \exp\left\{-W_{-1}\left(-\frac{1-\gamma}{2(1+\gamma)}\right)\right\}\right\}$, where W_{-1} is the Lambert W function, $c_R = \frac{1}{2} \min\left\{\frac{|r_1+r_2|}{e}, |r_1-r_2|\right\}$ when $r_1 \neq r_2$ and $c_R = \frac{|r_1+r_2|}{2e}$ otherwise.

Proof Sketch. Below we provide the proof sketch of Theorem 2. The full proof is deferred to Appendix F.

The eventual slow rate convergence is due to the heterogeneous environments \mathcal{P}^k regardless of the cardinality of the action space. In particular, we prove the slow rate when the action space is a singleton, in which case the Q-function coincides with the V-function. The process is also known as the Markov reward process. According to Algorithm 1, when $(t+1) \bmod E \neq 0$, we have

$$Q_{t+1}^k = ((1-\lambda)I + \lambda\gamma P^k) Q_t^k + \lambda R.$$

Following Algorithm 1, we let z denote the z -th synchronization round, and obtain the following recursion between two synchronization rounds:

$$\Delta_{(z+1)E} = \bar{A}^{(E)} \Delta_{zE} + \left((I - \bar{A}^{(E)}) - (I + \bar{A}^{(1)} + \dots + \bar{A}^{(E-1)}) (I - \bar{A}^{(1)}) \right) Q^*, \quad (10)$$

where $\bar{A}^{(\ell)} \triangleq \frac{1}{K} \sum_{k=1}^K (A^k)^\ell$ and $A^k \triangleq (1-\lambda)I + \lambda\gamma P^k$. While the first term on the right-hand side of (10) decays rapidly to zero, the second term is non-vanishing due to environment heterogeneity for $E \geq 2$. Specifically, to ensure the rapid decay of the first term, it is necessary to select a stepsize $\lambda = \tilde{\Omega}\left(\frac{1}{zE}\right)$. However, this choice results in the dominating residual error from the second term, which increases linearly with $\lambda E = \tilde{\Omega}(1/z)$.

Next, we instantiate the analyses by constructing the set \mathcal{P}^k over two states and an even number of clients with

$$P^{2k-1} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad P^{2k} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, \quad \text{for } k \in \mathbb{N}. \quad (11)$$

Applying the formula of $\bar{A}^{(\ell)}$ yields the following eigen-decomposition:

$$\bar{A}^{(\ell)} = \alpha_\ell (I - \bar{P}) + \beta_\ell \bar{P},$$

where $\bar{P} = \frac{1}{2} \mathbf{1}\mathbf{1}^\top$, $\alpha_\ell \triangleq \frac{1}{2}(\nu_1^\ell + \nu_2^\ell)$, $\beta_\ell \triangleq \nu_2^\ell$, $\nu_1 \triangleq 1 - (1+\gamma)\lambda$, and $\nu_2 \triangleq 1 - (1-\gamma)\lambda$. For this instance of \mathcal{P}_k , the error evolution (10) reduces to $\Delta_{(z+1)E} = (\alpha_E (I - \bar{P}) + \beta_E \bar{P}) \Delta_{zE} + \kappa_E (I - \bar{P}) Q^*$ with $\kappa_E \triangleq -\frac{\gamma}{2} \left(\frac{1-\nu_2^E}{1-\gamma} - \frac{1-\nu_1^E}{1+\gamma} \right)$, which further yields the following full error recursion:

$$\Delta_{zE} = (\alpha_E^z (I - \bar{P}) + \beta_E^z \bar{P}) \Delta_0 + \frac{1 - \alpha_E^z}{1 - \alpha_E} \kappa_E (I - \bar{P}) Q^*.$$

Starting from $Q_0 = \mathbf{0}$, the error can be decomposed into

$$\Delta_{zE} = \beta_E^z \bar{P} Q^* + \left(\alpha_E^z + \frac{1 - \alpha_E^z}{1 - \alpha_E} \kappa_E \right) (I - \bar{P}) Q^*. \quad (12)$$

The two terms of the error are orthogonal and both non-vanishing. Therefore, it remains to lower bound the maximum magnitude of two coefficients irrespective of the stepsize λ .

To this end, we analyze two regimes of λ separated by a threshold $\lambda_0 \triangleq \frac{\log(T/E)}{(1-\gamma)T}$:

- Slow rate due to small stepsize when $\lambda \leq \lambda_0$. Since β_E^z decreases as λ increases,

$$\beta_E^z \geq (1 - (1 - \gamma)\lambda_0)^{zE} = \left(1 - \frac{\log z}{zE}\right)^{zE} \geq \frac{E}{eT}.$$

- Slow rate due to environment heterogeneity when $\lambda \geq \lambda_0$. We show that

$$\left| \alpha_E^z + \frac{1 - \alpha_E^z \kappa_E}{1 - \alpha_E} \right| \geq \frac{E}{(1 - \gamma)T}.$$

We conclude that at least one component of the error in (12) must be slower than the rate $\Omega(E/T)$.

Remark 5. The explicit calculations are based on a set \mathcal{P}^k over a pair of states. Nevertheless, the evolution (10) is generally applicable. Similar analyses can be extended to scenarios involving more than two states, provided that the sequence of matrices $\bar{A}^{(\ell)}$ is simultaneously diagonalizable. For instance, the construction of the transition kernels in (11) can be readily extended to multiple states if the set \mathcal{S} can be partitioned into two different classes. The key insight is the non-vanishing residual on the right-hand side of (10) when $E \geq 2$ due to the environment heterogeneity.

5.3 Discussion on Time-varying Stepsize

Although using time-varying stepsize is common and simple when implementing the algorithm, it is not easy to transfer from current time-invariant stepsize analysis to time-varying stepsize analysis. This is because in the time-invariant stepsize analysis we are dealing with a function of one variable, however, in the time-varying case, we are dealing with a function of T variables.

For example, in our lower bound analysis, we picked a threshold λ_0 and showed that no matter λ is greater or smaller than λ_0 , the convergence rate is greater than $\Theta_R(E/((1-\gamma)T))$, and we can claim we have covered all the cases. However, for time-varying stepsize, the number of stepsizes is T , and it is not easy to generalize a similar result by just considering several cases because each stepsize gives an additional dimension. Even if we know the sequence is decaying, without specifying a particular family of stepsizes, it is not possible to divide it into several cases as we did for time-invariant stepsize.

We conjecture that both approaches lead to comparable residual error levels over extended training. For example, the stepsizes used in Figure 4a and Figure 4b are $\frac{1}{\sqrt{T}}$ and $\frac{1}{\sqrt{t+1}}$, respectively. While the time-decaying stepsize appears to have faster initial convergence due to its larger values, we observe that as t increases, the convergence rates of the two strategies seem to align, suggesting a similar asymptotic behavior.

6 Experiments

Description of the setup. In our experiments, we consider $K = 20$ agents (Jin et al., 2022), each interacting with an independently and randomly generated 5×5 maze environment $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}^k, R, \gamma \rangle$ for $k \in \{1, 2, \dots, 20\}$. The state set \mathcal{S} contains 25 cells that the agent is currently in. The action set contains 4 actions $\mathcal{A} = \{\text{left, up, right, down}\}$. Thus, $|\mathcal{S}| \times |\mathcal{A}| = 100$. We choose $\gamma = 0.99$. For ease of verifying our theory, each entry of the reward $R \in \mathbb{R}^{100}$ is sampled from $\text{Bern}(p = 0.05)$, which slightly departs from a typical maze environment wherein only two state-action pairs have nonzero rewards. We choose this reward so that $\|\Delta_0\|_\infty \approx 100 = \frac{1}{1-\gamma}$, which is the coarse upper bound of $\|\Delta_t\|_\infty$ for all t . For each agent k , its state transition probability vectors \mathcal{P}^k are constructed on top of standard state transition probability vectors of maze environments incorporated with a drifting probability 0.1 in each non-intentional action as in *WindyCliff* (Jin et al., 2022; Paul et al., 2019). In this way, the environment heterogeneity lies not only in the differences of the non-zero probability values (Jin et al., 2022; Paul et al., 2019) but also in the probability supports (i.e., the locations of non-zero entries). Our construction is more challenging: The

environment heterogeneity κ as per (2) of our environment construction was calculated to be 1.2. Yet, the largest environment heterogeneity of the WindyCliff construction in Jin et al. (2022) is about 0.31.

We choose $Q_0 = \mathbf{0} \in \mathbb{R}^{100}$. All numerical results are based on 5 independent runs to capture the variability. The dark lines represent the mean of the runs, while the shaded areas around each line illustrate the range obtained by adding and subtracting one standard deviation from the mean. The maximum time duration is $T = 20,000$ in our experiment since it is sufficient to capture the characteristics of the training process.

Convergence behavior and two-phase phenomenon. We demonstrate through numerical simulations that our analysis aligns with the observed behaviors. For algorithms with a time-invariant stepsize, convergence requires sufficiently small stepsizes and a sufficiently large number of iterations T .

To explore the impact of stepsizes on convergence, we use $\lambda \in \{0.9, 0.5, 0.2, 0.1, 0.05\}$, spanning a range within $(0, 1)$. As shown in Figure 2a, these stepsizes are not sufficiently small, leading to a two-phase phenomenon: the ℓ_∞ -norm of $\Delta_t = Q^* - \bar{Q}_t$ has a rapid decay in the first phase followed by a bounce back in the second phase. This phenomenon is distinctive to heterogeneous settings. In contrast, Figure 2b indicates that in homogeneous environments, no drastic bounce occurs, irrespective of the stepsize. Note that if multiple plots on the same page share the same legend, we display the legend only once for clarity.

Figure 4a (light blue curve) demonstrates that with a sufficiently small stepsize, such as $\lambda = \frac{1}{\sqrt{T}}$, the error continuously decreases, reaching approximately 24 by iteration 20,000.

A useful practice implication of our results is that: While constant stepsizes are often used in reinforcement learning problems because of the great performance in applications as described in Sutton & Barto (2018), they suffer significant performance degradation in the presence of environmental heterogeneity.

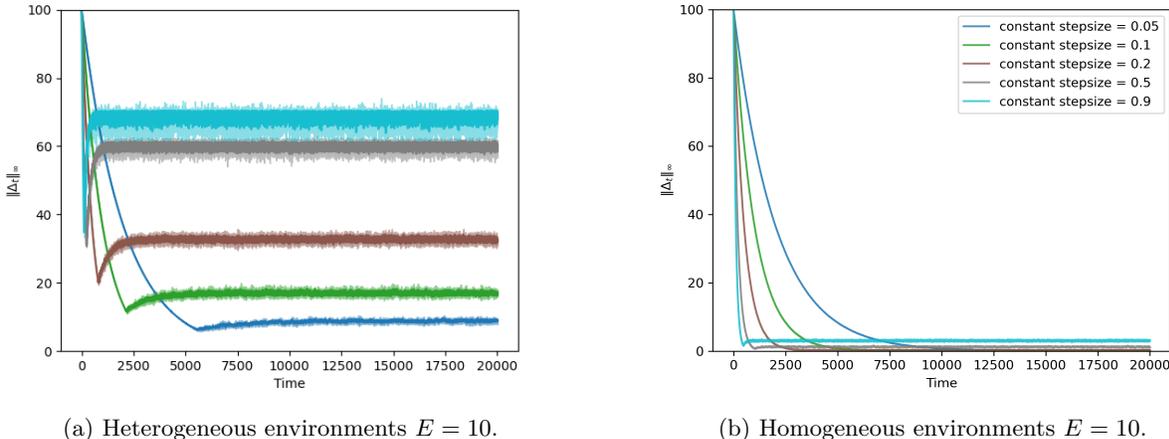


Figure 2: The ℓ_∞ error of different constant stepsizes under the heterogeneous and homogenous settings.

Impacts of the synchronization period E . In homogeneous settings (refer to Figure 5 in Appendix G.1), the synchronization period E has negligible impact, consistent with prior findings in the literature (Woo et al., 2023; Khodadadian et al., 2022). However, under heterogeneous conditions, larger E values lead to increased final error across the five constant stepsizes, as depicted in Figure 3 and Figure 2a. This degradation persists even with time-decaying stepsizes $\lambda_t = \frac{1}{\sqrt{t+1}}$, as shown in Figure 6. We hypothesize that larger E values require either smaller or more rapidly decaying stepsizes to mitigate the degradation caused by increased synchronization periods.

Potential utilization of the two-phase phenomenon. As shown in Figures 2a and 3, in the presence of environmental heterogeneity, the smaller the stepsizes, the smaller error $\|\Delta_t\|_\infty$ can reach and less significant of the error bouncing in the second phase. In our preliminary experiments, we tested small stepsizes $\lambda = 1/T^\alpha$

for $\alpha \in \{0.4, 0.5, \dots, 1\}$, which eventually lead to small errors yet at the cost of being extremely slow. Among these choices, $\lambda = 1/\sqrt{T}$ has the fastest convergence performance yet is still ≈ 24 up to iteration 20,000.

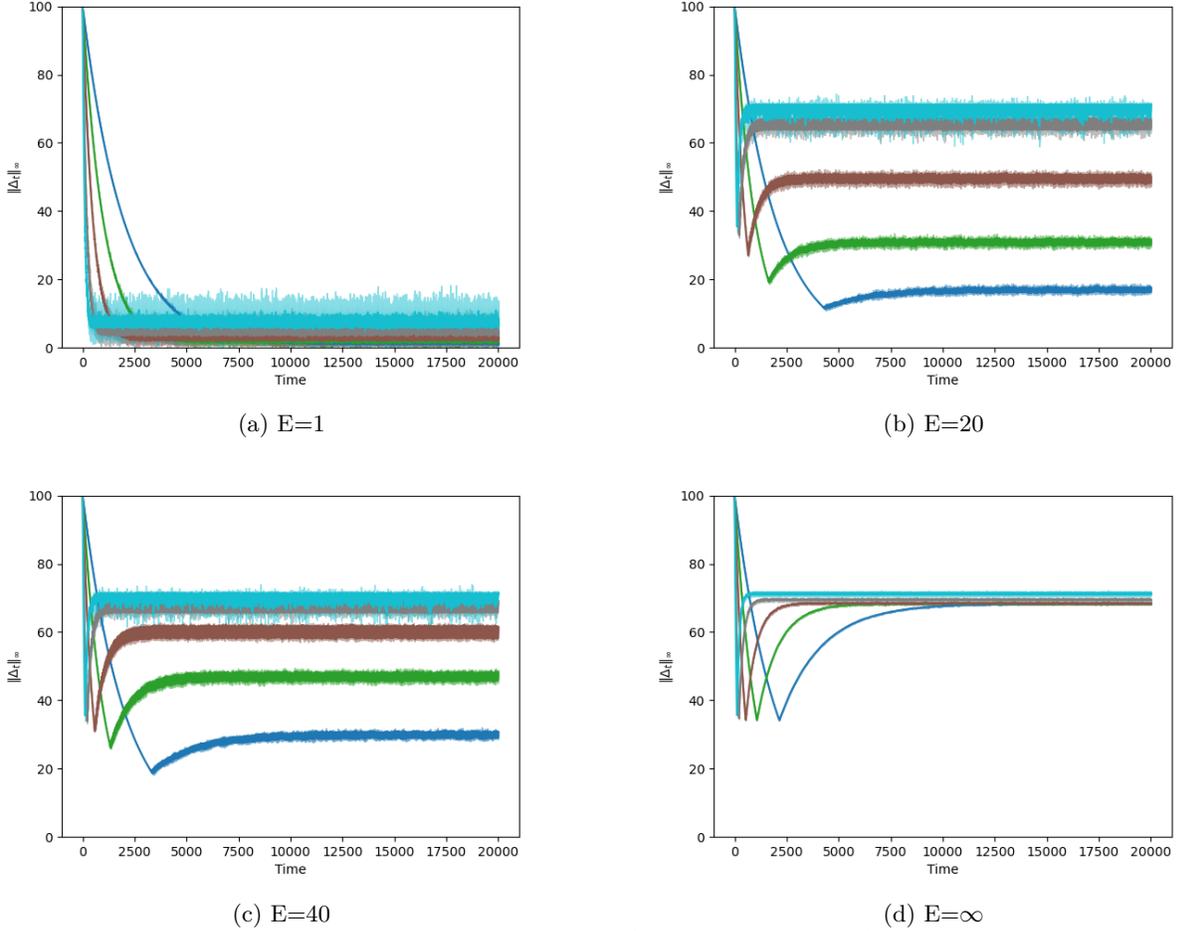


Figure 3: Convergence behavior for constant stepsizes (0.05, 0.1, 0.2, 0.5, 0.9) under various synchronization intervals E (1, 20, 40, ∞). In heterogeneous settings, higher E and larger λ lead to higher residual errors.

Let t_0 be the iteration at which the error trajectory $\|\Delta_t\|_\infty$ switches from phase 1 to phase 2. Provided that t_0 can be estimated, choosing different stepsizes for the two phases can lead to faster overall convergence, compared with using the same stepsize throughout.

Figure 4a illustrates two-phase training with different phase 1 stepsizes and phase 2 stepsize $\lambda = 1/\sqrt{T}$ compared with using $\lambda = 1/\sqrt{T}$ throughout. Overall, using $\lambda = 1/\sqrt{T}$ throughout leads to the slowest convergence, highlighting the benefits of the two-phase training strategy. Among all two-phase stepsize choices, the stepsize of 0.05 in the first phase results in a longer phase 1 duration ($t_0 = 5550$) but the lowest final error (2.75327), suggesting a better convergence. We further test the convergence performance with respect to different target error levels, details can be found in Appendix G.3.

We also evaluated the two-phase training strategy using various time-decaying step sizes, including $\frac{1}{\sqrt{t+1}}$, $\frac{c+1}{t+c}$, $\frac{1}{t+1}$, and $\frac{1}{(t+1)^{0.7}}$. In all cases, Figure 4 shows the two-phase training has an advantage.

We leave the estimation and characterization of t_0 for future work.

Acknowledgments

We thank He Cheng for his valuable assistance on the experiments. P. Yang is supported in part by the National Key R&D Program of China 2024YFA1015800, and Tsinghua University Dushi Program 2025Z11DSZ001. L. Su is supported in part by NSF CAREER CIF-2340482.

References

- Dimitri Bertsekas. *Reinforcement learning and optimal control*, volume 1. Athena Scientific, 2019.
- Dimitri Bertsekas and John N Tsitsiklis. *Neuro-dynamic programming*. Athena Scientific, 1996.
- Jalaj Bhandari, Daniel Russo, and Raghav Singal. A finite time analysis of temporal difference learning with linear function approximation. In *Conference on learning theory*, pp. 1691–1692. PMLR, 2018.
- Jin-Hua Chen, Min-Rong Chen, Guo-Qiang Zeng, and Jia-Si Weng. Bdf: a byzantine-fault-tolerance decentralized federated learning method for autonomous vehicle. *IEEE Transactions on Vehicular Technology*, 70(9):8639–8652, 2021.
- Filippos Christianos, Georgios Papoudakis, Muhammad A Rahman, and Stefano V Albrecht. Scaling multi-agent reinforcement learning with selective parameter sharing. In *International Conference on Machine Learning*, pp. 1989–1998. PMLR, 2021.
- Yuyang Deng, Mohammad Mahdi Kamani, and Mehrdad Mahdavi. Adaptive personalized federated learning. *arXiv preprint arXiv:2003.13461*, 2020.
- Thinh Doan, Siva Maguluri, and Justin Romberg. Finite-time analysis of distributed td (0) with linear function approximation on multi-agent reinforcement learning. In *International Conference on Machine Learning*, pp. 1626–1635. PMLR, 2019.
- Zhaoyang Du, Celimuge Wu, Tsutomu Yoshinaga, Kok-Lim Alvin Yau, Yusheng Ji, and Jie Li. Federated learning for vehicular internet of things: Recent advances and open issues. *IEEE Open Journal of the Computer Society*, 1:45–61, 2020.
- Hao Jin, Yang Peng, Wenhao Yang, Shusen Wang, and Zhihua Zhang. Federated reinforcement learning with environment heterogeneity. In *International Conference on Artificial Intelligence and Statistics*, pp. 18–37. PMLR, 2022.
- Peter Kairouz, H. Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, Rafael G. L. D’Oliveira, Hubert Eichner, Salim El Rouayheb, David Evans, Josh Gardner, Zachary Garrett, Adrià Gascón, Badi Ghazi, Phillip B. Gibbons, Marco Gruteser, Zaid Harchaoui, Chaoyang He, Lie He, Zhouyuan Huo, Ben Hutchinson, Justin Hsu, Martin Jaggi, Tara Javidi, Gauri Joshi, Mikhail Khodak, Jakub Konečný, Aleksandra Korolova, Farinaz Koushanfar, Sanmi Koyejo, Tancrede Lepoint, Yang Liu, Prateek Mittal, Mehryar Mohri, Richard Nock, Ayfer Özgür, Rasmus Pagh, Hang Qi, Daniel Ramage, Ramesh Raskar, Mariana Raykova, Dawn Song, Weikang Song, Sebastian U. Stich, Ziteng Sun, Ananda Theertha Suresh, Florian Tramèr, Praneeth Vepakomma, Jianyu Wang, Li Xiong, Zheng Xu, Qiang Yang, Felix X. Yu, Han Yu, and Sen Zhao. Advances and open problems in federated learning. *Foundations and Trends® in Machine Learning*, 14(1–2):1–210, 2021.
- Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In *International Conference on Machine Learning*, pp. 5132–5143. PMLR, 2020.
- Sajad Khodadadian, Pranay Sharma, Gauri Joshi, and Siva Theja Maguluri. Federated reinforcement learning: Linear speedup under markovian sampling. In *International Conference on Machine Learning*, pp. 10997–11057. PMLR, 2022.

- B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A Al Sallab, Senthil Yogamani, and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. *IEEE Transactions on Intelligent Transportation Systems*, 23(6):4909–4926, 2021.
- Gen Li, Yuting Wei, Yuejie Chi, Yuantao Gu, and Yuxin Chen. Sample complexity of asynchronous q-learning: Sharper analysis and variance reduction. *Advances in neural information processing systems*, 33:7031–7043, 2020a.
- Gen Li, Changxiao Cai, Yuxin Chen, Yuting Wei, and Yuejie Chi. Is q-learning minimax optimal? a tight sample complexity analysis. *Operations Research*, 72(1):222–236, 2024.
- Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. *Proceedings of Machine Learning and Systems*, 2:429–450, 2020b.
- Yang Liu, Anbu Huang, Yun Luo, He Huang, Youzhi Liu, Yuanyuan Chen, Lican Feng, Tianjian Chen, Han Yu, and Qiang Yang. Fedvision: An online visual object detection platform powered by federated learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 13172–13179, 2020.
- Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pp. 1273–1282. PMLR, 2017.
- Thies Möhlenhof, Norman Jansen, and Wiam Rachid. Reinforcement learning environment for tactical networks. In *2021 International Conference on Military Communication and Information Systems (ICMCIS)*, pp. 1–8. IEEE, 2021.
- Thien Duc Nguyen, Samuel Marchal, Markus Miettinen, Hossein Fereidooni, N Asokan, and Ahmad-Reza Sadeghi. Diot: A federated self-learning anomaly detection system for iot. In *2019 IEEE 39th International conference on distributed computing systems (ICDCS)*, pp. 756–767. IEEE, 2019.
- In-Beom Park, Jaeseok Huh, Joongkyun Kim, and Jonghun Park. A reinforcement learning approach to robust scheduling of semiconductor manufacturing facilities. *IEEE Transactions on Automation Science and Engineering*, 17(3):1420–1431, 2019.
- Reese Pathak and Martin J Wainwright. Fedsplit: An algorithmic framework for fast federated optimization. *Advances in neural information processing systems*, 33:7057–7066, 2020.
- Supratik Paul, Michael A Osborne, and Shimon Whiteson. Fingerprint policy optimisation for robust reinforcement learning. In *International Conference on Machine Learning*, pp. 5082–5091. PMLR, 2019.
- Muzi Peng, Jiangwei Wang, Dongjin Song, Fei Miao, and Lili Su. Privacy-preserving and uncertainty-aware federated trajectory prediction for connected autonomous vehicles. In *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 11141–11147, 2023. doi: 10.1109/IROS55552.2023.10341638.
- Jason Posner, Lewis Tseng, Moayad Aloqaily, and Yaser Jararweh. Federated learning in vehicular networks: Opportunities and solutions. *IEEE Network*, 35(2):152–159, 2021.
- Jiaju Qi, Qihao Zhou, Lei Lei, and Kan Zheng. Federated reinforcement learning: techniques, applications, and open challenges. *Intelligence & Robotics*, 2021. doi: 10.20517/ir.2021.02. URL <https://doi.org/10.20517/2Fir.2021.02>.
- Swaroop Ramaswamy, Rajiv Mathews, Kanishka Rao, and Françoise Beaufays. Federated learning for emoji prediction in a mobile keyboard. *arXiv preprint arXiv:1906.04329*, 2019.
- Sudeep Salgia and Yuejie Chi. The sample-communication complexity trade-off in federated q-learning. *Advances in Neural Information Processing Systems*, 37:39694–39747, 2025.

- Micah J Sheller, G Anthony Reina, Brandon Edwards, Jason Martin, and Spyridon Bakas. Multi-institutional deep learning modeling without sharing patient data: A feasibility study on brain tumor segmentation. In *Brainlesion: Glioma, Multiple Sclerosis, Stroke and Traumatic Brain Injuries: 4th International Workshop, BrainLes 2018, Held in Conjunction with MICCAI 2018, Granada, Spain, September 16, 2018, Revised Selected Papers, Part I 4*, pp. 92–104. Springer, 2019.
- Lili Su, Jiaming Xu, and Pengkun Yang. A non-parametric view of fedavg and fedprox: Beyond stationary points. *Journal of Machine Learning Research*, 24(203):1–48, 2023.
- Richard S. Sutton and Andrew G. Barto. *Chapter 2.5 Tracking a Nonstationary Problem, Reinforcement Learning: An Introduction*, chapter 8, pp. 33. The MIT Press, 2018.
- Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.
- Chunnan Wang, Xiang Chen, Junzhe Wang, and Hongzhi Wang. Atpfl: Automatic trajectory prediction model design under federated learning framework. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 6563–6572, June 2022a.
- Han Wang, Aritra Mitra, Hamed Hassani, George J Pappas, and James Anderson. Federated temporal difference learning with linear function approximation under environmental heterogeneity. *arXiv preprint arXiv:2302.02212*, 2023.
- Jianyu Wang, Qinghua Liu, Hao Liang, Gauri Joshi, and H Vincent Poor. Tackling the objective inconsistency problem in heterogeneous federated optimization. *Advances in neural information processing systems*, 33: 7611–7623, 2020.
- Jianyu Wang, Rudrajit Das, Gauri Joshi, Satyen Kale, Zheng Xu, and Tong Zhang. On the unreasonable effectiveness of federated averaging with heterogeneous data. *arXiv preprint arXiv:2206.04723*, 2022b.
- Christopher Watkins and Peter Dayan. Q-learning. *Machine Learning*, 8:279–292, 1992. URL <https://api.semanticscholar.org/CorpusID:208910339>.
- Jiin Woo, Gauri Joshi, and Yuejie Chi. The blessing of heterogeneity in federated q-learning: Linear speedup and beyond. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 37157–37216. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/woo23a.html>.
- Zhijie Xie and Shenghui Song. FedKL: Tackling Data Heterogeneity in Federated Reinforcement Learning by Penalizing KL Divergence. *IEEE Journal on Selected Areas in Communications*, 41(4):1227–1242, April 2023. ISSN 1558-0008. doi: 10.1109/JSAC.2023.3242734. URL https://ieeexplore.ieee.org/abstract/document/10038492?casa_token=yGyMDlnL_FsAAAAA:hqNvzWEb6yVKwTZVdHKLvNorDg07AWx4uuuDjDsLLTTY_7unjr1ew8Yv4_UUAWfCz3X1b9wHNYSP8.
- Bingjie Yan, Jun Wang, Jieren Cheng, Yize Zhou, Yixian Zhang, Yifan Yang, Li Liu, Haojiang Zhao, Chunjuan Wang, and Boyi Liu. Experiments of federated learning for covid-19 chest x-ray images. In *Advances in Artificial Intelligence and Security: 7th International Conference, ICAIS 2021, Dublin, Ireland, July 19-23, 2021, Proceedings, Part II 7*, pp. 41–53. Springer, 2021.
- Timothy Yang, Galen Andrew, Hubert Eichner, Haicheng Sun, Wei Li, Nicholas Kong, Daniel Ramage, and Françoise Beaufays. Applied federated learning: Improving google keyboard query suggestions. *arXiv preprint arXiv:1812.02903*, 2018.
- Chao Yu, Akash Velu, Eugene Vinitzky, Jiaxuan Gao, Yu Wang, Alexandre Bayen, and Yi Wu. The surprising effectiveness of ppo in cooperative multi-agent games. *Advances in neural information processing systems*, 35:24611–24624, 2022.

- Tianlong Yu, Tian Li, Yuqiong Sun, Susanta Nanda, Virginia Smith, Vyas Sekar, and Srinivasan Seshan. Learning context-aware policies from multiple smart homes via federated multi-task learning. In *2020 IEEE/ACM Fifth International Conference on Internet-of-Things Design and Implementation (IoTDI)*, pp. 104–115. IEEE, 2020.
- Tengchan Zeng, Omid Semiari, Mingzhe Chen, Walid Saad, and Mehdi Bennis. Federated learning on the road autonomous controller design for connected and autonomous vehicles. *IEEE Transactions on Wireless Communications*, 21(12):10407–10423, 2022.
- Chenyu Zhang, Han Wang, Aritra Mitra, and James Anderson. Finite-time analysis of on-policy heterogeneous federated reinforcement learning. *arXiv preprint arXiv:2401.15273*, 2024.
- Shangdong Zhang, Remi Tachet Des Combes, and Romain Laroche. On the convergence of sarsa with linear function approximation. In *International Conference on Machine Learning*, pp. 41613–41646. PMLR, 2023.
- Zhong Zheng, Fengyu Gao, Lingzhou Xue, and Jing Yang. Federated q-learning: Linear regret speedup with low communication cost. *arXiv preprint arXiv:2312.15023*, 2023.
- Yifan Zhong, Jakub Grudzien Kuba, Xidong Feng, Siyi Hu, Jiaming Ji, and Yaodong Yang. Heterogeneous-agent reinforcement learning. *Journal of Machine Learning Research*, 25(32):1–67, 2024. URL <http://jmlr.org/papers/v25/23-0488.html>.

Appendices

A Proof of Lemma 1

The evolution of Δ_{t+1} can be decomposed as follows:

$$\begin{aligned}
\Delta_{t+1} &= Q^* - \bar{Q}_{t+1} \\
&\stackrel{(a)}{=} \frac{1}{K} \sum_{k=1}^K (Q^* - ((1-\lambda)Q_t^k + \lambda(R + \gamma\tilde{P}_t^k V_t^k))) \\
&= \frac{1}{K} \sum_{k=1}^K ((1-\lambda)(Q^* - Q_t^k) + \lambda(Q^* - R - \gamma\tilde{P}_t^k V_t^k)) \\
&\stackrel{(b)}{=} (1-\lambda)\Delta_t + \gamma\lambda \frac{1}{K} \sum_{k=1}^K (\bar{P}V^* - \tilde{P}_t^k V_t^k) \\
&= (1-\lambda)\Delta_t + \frac{\gamma\lambda}{K} \sum_{k=1}^K (\bar{P} - \tilde{P}_t^k)V^* + \frac{\gamma\lambda}{K} \sum_{k=1}^K \tilde{P}_t^k (V^* - V_t^k) \\
&= (1-\lambda)^{t+1}\Delta_0 + \gamma\lambda \sum_{i=0}^t (1-\lambda)^{t-i} \frac{1}{K} \sum_{k=1}^K (\bar{P} - \tilde{P}_i^k)V^* \\
&\quad + \gamma\lambda \sum_{i=0}^t (1-\lambda)^{t-i} \frac{1}{K} \sum_{k=1}^K \tilde{P}_i^k (V^* - V_i^k),
\end{aligned}$$

where equality (a) follows from Eq. (4), equality (b) follows from the Bellman optimality equation in Eq. (3), and the last equality follows from unrolling the updates $t+1$ times and using the fact that $\Delta_0 = Q^* - Q_0$.

B Proof of Lemma 2

We first show $0 \leq Q_t^k(s, a) \leq \frac{1}{1-\gamma}$ by inducting on t . When $t=0$, this is true by the choice of Q_0 . Suppose that $0 \leq Q_{t-1}^k(s, a) \leq \frac{1}{1-\gamma}$ for any state-action pair (s, a) and any client k . Let's focus on time t . When t is not a synchronization iteration (i.e., $t+1 \bmod E \neq 0$), we have

$$\begin{aligned}
Q_t^k(s, a) &= (1-\lambda)Q_{t-1}^k(s, a) + \lambda(R(s, a) + \gamma\tilde{P}_t^k(s, a)V_{t-1}^k) \\
&\leq \frac{1-\lambda}{1-\gamma} + \lambda(R(s, a) + \gamma\tilde{P}_t^k(s, a)V_{t-1}^k) \\
&\stackrel{(a)}{\leq} \frac{1-\lambda}{1-\gamma} + \lambda\left(1 + \frac{\gamma}{1-\gamma}\right) \\
&\leq \frac{1}{1-\gamma} - \frac{\lambda}{1-\gamma} + \frac{\lambda}{1-\gamma} \\
&= \frac{1}{1-\gamma},
\end{aligned}$$

where inequality (a) holds because for any $s, V_{t-1}^k(s) = \max_{a \in \mathcal{A}} Q_{t-1}^k(s, a) \leq \frac{1}{1-\gamma}$ by the inductive hypothesis, and each element of $\tilde{P}_t^k(s, a) \in [0, 1]$. Then $\tilde{P}_t^k(s, a)V_{t-1}^k \leq \|\tilde{P}_t^k(s, a)\|_1 \|V_{t-1}^k\|_\infty \leq \frac{1}{1-\gamma}$ by Hölder's inequality. Similarly, we can show the case when t is a synchronization iteration.

With the above argument, we can also show that $0 \leq Q^*(s, a) \leq \frac{1}{1-\gamma}$ for any state-action pair (s, a) . Therefore, we have that $\|Q^* - Q_t^k\|_\infty \leq \frac{1}{1-\gamma}$.

Next, we show that bound on $\|V^* - V_t^k\|_\infty$.

$$\begin{aligned}
\|V^* - V_t^k\|_\infty &= \max_{s \in \mathcal{S}} |V^*(s) - V_t^k(s)| \\
&= \max_{s \in \mathcal{S}} \left| \max_{a \in \mathcal{A}} Q^*(s, a) - \max_{a' \in \mathcal{A}} Q_t^k(s, a') \right| \\
&\leq \max_{s \in \mathcal{S}, a \in \mathcal{A}} |Q^*(s, a) - Q_t^k(s, a)| \\
&= \|Q^* - Q_t^k\|_\infty \\
&\leq \frac{1}{1 - \gamma}.
\end{aligned}$$

C Proof of Lemma 3

When $t \bmod E = 0$, i.e., i is a synchronization round, $Q_t^k = Q_t^{k'}$ for any pair of agents $k, k' \in [K]$. Hence,

$$\begin{aligned}
\frac{1}{K} \sum_{k=1}^K \tilde{P}_t^k(s, a)(V^* - V_t^k) &= \left(\frac{1}{K} \sum_{k=1}^K \tilde{P}_t^k(s, a) \right) (V^* - \bar{V}_t) \\
&\leq \left\| \frac{1}{K} \sum_{k=1}^K \tilde{P}_t^k(s, a) \right\|_1 \|V^* - \bar{V}_t\|_\infty \\
&\leq \|V^* - \bar{V}_t\|_\infty \\
&\leq \|Q^* - \bar{Q}_t\|_\infty \\
&= \|\Delta_t\|_\infty.
\end{aligned} \tag{13}$$

For general t , we anchor the error term to that of the synchronization round as follows:

$$\begin{aligned}
\left\| \frac{1}{K} \sum_{k=1}^K \tilde{P}_t^k(V^* - V_t^k) \right\|_\infty &= \left\| \frac{1}{K} \sum_{k=1}^K \tilde{P}_t^k(V^* - V_{\chi(t)}^k + V_{\chi(t)}^k - V_i^k) \right\|_\infty \\
&\leq \left\| \frac{1}{K} \sum_{k=1}^K \tilde{P}_t^k(V^* - V_{\chi(t)}^k) \right\|_\infty + \left\| \frac{1}{K} \sum_{k=1}^K \tilde{P}_t^k(V_{\chi(t)}^k - V_t^k) \right\|_\infty \\
&\stackrel{(a)}{\leq} \|\Delta_{\chi(t)}\|_\infty + \left\| \frac{1}{K} \sum_{k=1}^K \tilde{P}_t^k(V_{\chi(t)}^k - V_t^k) \right\|_\infty \\
&\leq \|\Delta_{\chi(t)}\|_\infty + \frac{1}{K} \sum_{k=1}^K \|V_{\chi(t)}^k - V_t^k\|_\infty.
\end{aligned} \tag{14}$$

where inequality (a) follows from Eq. (13). For any state $s \in \mathcal{S}$, let $a_t^k(s) \in \arg \max_{a \in \mathcal{A}} Q_t^k(s, a)$ for all t and k . We have

$$\begin{aligned}
&V_t^k(s) - V_{\chi(t)}^k(s) \\
&= Q_t^k(s, a_t^k(s)) - Q_{\chi(t)}^k(s, a_{\chi(t)}^k(s)) \\
&\stackrel{(a)}{\leq} Q_t^k(s, a_t^k(s)) - Q_{\chi(t)}^k(s, a_t^k(s)) \\
&= Q_t^k(s, a_t^k(s)) - Q_{t-1}^k(s, a_t^k(s)) + Q_{t-1}^k(s, a_t^k(s)) - Q_{t-2}^k(s, a_t^k(s)) \\
&\quad + \cdots + Q_{\chi(t)+1}^k(s, a_t^k(s)) - Q_{\chi(t)}^k(s, a_t^k(s)).
\end{aligned} \tag{15}$$

where inequality (a) holds because $Q_{\chi(t)}^k(s, a_t^k(s)) \leq Q_{\chi(t)}^k(s, a_{\chi(t)}^k(s))$. For each t' such that $\chi(t) \leq t' \leq t$, it holds that,

$$\begin{aligned}
& Q_{t'+1}^k(s, a_t^k(s)) - Q_{t'}^k(s, a_t^k(s)) \\
&= (1 - \lambda)Q_{t'}^k(s, a_t^k(s)) + \lambda(R(s, a_t^k(s)) + \gamma\tilde{P}_{t'}^k(s, a_t^k(s))V_{t'}^k) - Q_{t'}^k(s, a_t^k(s)) \\
&\stackrel{(a)}{=} -\lambda Q_{t'}^k(s, a_t^k(s)) + \lambda \left(Q^*(s, a_t^k(s)) - R(s, a_t^k(s)) - \gamma\bar{P}(s, a_t^k(s))V^* + R(s, a_t^k(s)) + \gamma\tilde{P}_{t'}^k(s, a_t^k(s))V_{t'}^k \right) \\
&= \lambda\Delta_{t'}^k(s, a_t^k(s)) + \gamma\lambda \left(\tilde{P}_{t'}^k(s, a_t^k(s))V_{t'}^k - \bar{P}(s, a_t^k(s))V^* \right) \\
&= \lambda\Delta_{t'}^k(s, a_t^k(s)) + \gamma\lambda \left((\tilde{P}_{t'}^k(s, a_t^k(s)) - \bar{P}(s, a_t^k(s)))V^* + \tilde{P}_{t'}^k(s, a_t^k(s))(V_{t'}^k - V^*) \right) \\
&\leq 2\lambda \|\Delta_{t'}^k\|_\infty + \gamma\lambda \left(\tilde{P}_{t'}^k(s, a_t^k(s)) - \bar{P}(s, a_t^k(s)) \right) V^*,
\end{aligned}$$

where equality (a) follows from the Bellman equation Eq. (3), and the last inequality follows from Eq.(13) and that

$$\gamma\lambda\tilde{P}_{t'}^k(s, a_t^k(s))(V_{t'}^k - V^*) \leq \gamma\lambda\|\tilde{P}_{t'}^k(s, a_t^k(s))\|_1\|V_{t'}^k - V^*\|_\infty \leq \gamma\|V_{t'}^k - V^*\|_\infty.$$

Thus, $V_t^k(s) - V_{\chi(t)}^k(s)$ can be upper bounded as

$$\begin{aligned}
V_t^k(s) - V_{\chi(t)}^k(s) &\leq \sum_{t'=\chi(t)}^{t-1} Q_{t'+1}^k(s, a_t^k(s)) - Q_{t'}^k(s, a_t^k(s)) \\
&= 2\lambda \sum_{t'=\chi(t)}^{t-1} \|\Delta_{t'}^k\|_\infty + \gamma\lambda \sum_{t'=\chi(t)}^{t-1} \left(\tilde{P}_{t'}^k(s, a_t^k(s)) - \bar{P}(s, a_t^k(s)) \right) V^*. \tag{16}
\end{aligned}$$

Similarly, we have the following lower bound:

$$\begin{aligned}
V_t^k(s) - V_{\chi(t)}^k(s) &\geq \sum_{t'=\chi(t)}^{t-1} Q_{t'+1}^k(s, a_{\chi(t)}^k(s)) - Q_{t'}^k(s, a_{\chi(t)}^k(s)) \\
&\geq -2\lambda \sum_{t'=\chi(t)}^{t-1} \|\Delta_{t'}^k\|_\infty + \gamma\lambda \sum_{t'=\chi(t)}^{t-1} \left(\tilde{P}_{t'}^k(s, a_{\chi(t)}^k(s)) - \bar{P}(s, a_{\chi(t)}^k(s)) \right) V^*. \tag{17}
\end{aligned}$$

Plugging the bounds in Eq. (16) and in Eq. (17) back into Eq. (14), we get

$$\begin{aligned}
\left\| \frac{1}{K} \sum_{k=1}^K \tilde{P}_t^k(V^* - V_t^k) \right\|_\infty &\leq \|\Delta_{\chi(t)}\|_\infty + \frac{1}{K} \sum_{k=1}^K \|V_{\chi(t)}^k - V_t^k\|_\infty \\
&\leq \|\Delta_{\chi(t)}\|_\infty + 2\lambda \frac{1}{K} \sum_{k=1}^K \sum_{t'=\chi(t)}^{t-1} \|\Delta_{t'}^k\|_\infty \\
&\quad + \gamma\lambda \frac{1}{K} \sum_{k=1}^K \max_{s,a} \left| \sum_{t'=\chi(t)}^{t-1} \left(\tilde{P}_{t'}^k(s, a) - \bar{P}(s, a) \right) V^* \right|,
\end{aligned}$$

proving the lemma.

D Proof of Lemma 4

When $i \bmod E = 0$, then $\Delta_i^k = \Delta_{\chi(i)}$. When $i \bmod E \neq 0$, we have

$$\begin{aligned}
Q_i^k &= (1 - \lambda)Q_{i-1}^k + \lambda \left(R + \gamma\tilde{P}_{i-1}^k V_{i-1}^k \right) \\
&= (1 - \lambda)Q_{i-1}^k + \lambda \left(Q^* - R - \gamma\bar{P}V^* + R + \gamma\tilde{P}_{i-1}^k V_{i-1}^k \right).
\end{aligned}$$

So,

$$\begin{aligned}
\Delta_i^k &= (1-\lambda)\Delta_{i-1}^k + \lambda\gamma\left(\bar{P}V^* - \tilde{P}_{i-1}^k V_{i-1}^k\right) \\
&= (1-\lambda)\Delta_{i-1}^k + \lambda\gamma(\bar{P} - \tilde{P}_{i-1}^k)V^* + \lambda\gamma\tilde{P}_{i-1}^k(V^* - V_{i-1}^k) \\
&\leq (1-\lambda)^{i-\chi(i)}\Delta_{\chi(i)} + \gamma\lambda\sum_{j=\chi(i)}^{i-1}(1-\lambda)^{i-j-1}(\bar{P} - \tilde{P}_j^k)V^* \\
&\quad + \gamma\lambda\sum_{j=\chi(i)}^{i-1}(1-\lambda)^{i-j-1}\tilde{P}_j^k(V^* - V_j^k). \tag{18}
\end{aligned}$$

For any state-action pair (s, a) ,

$$|(1-\lambda)^{i-\chi(i)}\Delta_{\chi(i)}(s, a)| \leq (1-\lambda)^{i-\chi(i)}\|\Delta_{\chi(i)}\|_\infty. \tag{19}$$

For the second term, we have

$$\left\|\gamma\lambda\sum_{j=\chi(i)}^{i-1}(1-\lambda)^{i-j-1}(\bar{P} - \tilde{P}_j^k)V^*\right\|_\infty \tag{20}$$

$$\leq \left\|\gamma\lambda\sum_{j=\chi(i)}^{i-1}(1-\lambda)^{i-j-1}(\bar{P} - P_j^k)V^*\right\|_\infty + \left\|\gamma\lambda\sum_{j=\chi(i)}^{i-1}(1-\lambda)^{i-j-1}(P_j^k - \tilde{P}_j^k)V^*\right\|_\infty \tag{21}$$

$$\leq \frac{\gamma}{1-\gamma}\lambda\sum_{j=\chi(i)}^{i-1}(1-\lambda)^{i-1-j}\kappa + \frac{\gamma}{1-\gamma}\sqrt{\lambda\log\frac{|\mathcal{S}||\mathcal{A}|KT}{\delta}} \tag{22}$$

$$\leq \frac{\gamma}{1-\gamma}\lambda(E-1)\kappa + \frac{\gamma}{1-\gamma}\sqrt{\lambda\log\frac{|\mathcal{S}||\mathcal{A}|KT}{\delta}}, \tag{23}$$

for all $(s, a) \in \mathcal{S} \times \mathcal{A}, i \in [T], k \in [K]$. From Eq. (21) to Eq. (22), since for each timestep, each agent independently samples the next state for all state-action pairs to form the sample transition matrix, we can use Hoeffding's inequality by treating $\lambda(1-\lambda)^{i-j-1}(P_j^k - \tilde{P}_j^k)V^*$ as the independent random variables with their absolute values bounded by $\lambda(1-\lambda)^{i-j-1}\|V^*\|_\infty$.

In addition, we have

$$\left\|\gamma\lambda\sum_{j=\chi(i)}^{i-1}(1-\lambda)^{i-j-1}\tilde{P}_j^k(V^* - V_j^k)\right\|_\infty \leq \gamma\lambda\sum_{j=\chi(i)}^{i-1}(1-\lambda)^{i-j-1}\|\Delta_j^k\|_\infty. \tag{24}$$

Combining the bounds in Eq. (19), Eq. (20), and Eq. (24), we get

$$\begin{aligned}
\|\Delta_i^k\|_\infty &\leq (1-\lambda)^{i-\chi(i)}\|\Delta_{\chi(i)}\|_\infty + \frac{\gamma}{1-\gamma}\lambda(E-1)\kappa + \frac{\gamma}{1-\gamma}\sqrt{\lambda\log\frac{|\mathcal{S}||\mathcal{A}|KT}{\delta}} \\
&\quad + \gamma\lambda\sum_{j=\chi(i)}^{i-1}(1-\lambda)^{i-j-1}\|\Delta_j^k\|_\infty \\
&\leq (1-(1-\gamma)\lambda)^{i-\chi(i)}\|\Delta_{\chi(i)}\|_\infty \\
&\quad + (1+\gamma\lambda)^{i-\chi(i)}\left(\frac{\gamma}{1-\gamma}\lambda(E-1)\kappa + \frac{\gamma}{1-\gamma}\sqrt{\lambda\log\frac{|\mathcal{S}||\mathcal{A}|KT}{\delta}}\right), \tag{25}
\end{aligned}$$

where the last inequality can be shown via inducting on $i - \chi(i) \in \{0, \dots, E-1\}$. When $\lambda \leq \frac{1}{E}$,

$$(1+\gamma\lambda)^{i-\chi(i)} \leq (1+\lambda)^E \leq (1+1/E)^E \leq e \leq 3.$$

We get

$$\|\Delta_i^k\|_\infty \leq \|\Delta_{\chi(i)}\|_\infty + 3\frac{\gamma}{1-\gamma}\lambda(E-1)\kappa + 3\frac{\gamma}{1-\gamma}\sqrt{\lambda \log \frac{|\mathcal{S}||\mathcal{A}|KT}{\delta}}.$$

E Proof of Theorem 1

By Lemma 1,

$$\Delta_{t+1} = (1-\lambda)^{t+1}\Delta_0 + \sum_{i=0}^t (1-\lambda)^i \frac{\gamma\lambda}{K} \sum_{k=1}^K (\bar{P} - \tilde{P}_{t-i}^k) V^* + \sum_{i=0}^t (1-\lambda)^i \frac{\gamma\lambda}{K} \sum_{k=1}^K \tilde{P}_{t-i}^k (V^* - V_{t-i}^k).$$

Taking the ℓ_∞ norm on both sides, we get

$$\begin{aligned} \|\Delta_{t+1}\|_\infty &\leq \underbrace{(1-\lambda)^{t+1} \|\Delta_0\|_\infty}_{(I.1)} + \underbrace{\left\| \sum_{i=0}^t (1-\lambda)^i \lambda \gamma \frac{1}{K} \sum_{k=1}^K (\bar{P} - \tilde{P}_{t-i}^k) V^* \right\|_\infty}_{(I.2)} \\ &\quad + \underbrace{\left\| \sum_{i=0}^t (1-\lambda)^i \lambda \gamma \frac{1}{K} \sum_{k=1}^K \tilde{P}_{t-i}^k (V^* - V_{t-i}^k) \right\|_\infty}_{(I.3)}. \end{aligned} \quad (26)$$

We bound the three terms in the right-hand-side of the above-displayed equation separately.

Bounding (I.1). Since $0 \leq Q_0(s, a) \leq \frac{1}{1-\gamma}$, the first term can be bounded as

$$(I.1) = (1-\lambda)^{t+1} \|\Delta_0\|_\infty \leq (1-\lambda)^{t+1} \frac{1}{1-\gamma}. \quad (27)$$

Bounding (I.2). To bound the second term (I.2) in Eq. (26), we have

$$\begin{aligned} \sum_{i=0}^t (1-\lambda)^i \lambda \gamma \frac{1}{K} \sum_{k=1}^K (\bar{P} - \tilde{P}_{t-i}^k) V^* &= \sum_{i=0}^t (1-\lambda)^i \lambda \gamma \frac{1}{K} \sum_{k=1}^K (P^k - \tilde{P}_{t-i}^k) V^* \\ &= \frac{1}{K} \sum_{k=1}^K \sum_{i=0}^t (1-\lambda)^i \lambda \gamma (P^k - \tilde{P}_{t-i}^k) V^*. \end{aligned}$$

Let $X_{i,k} = \frac{1}{K} \lambda \gamma (1-\lambda)^i (P^k - \tilde{P}_{t-i}^k) V^*$. It is easy to see that $\mathbb{E}[X_{i,k}(s, a)] = 0$ for all (s, a) . By Lemma 2, we have $|X_{i,k}(s, a)| \leq \frac{2}{K(1-\gamma)} \lambda \gamma (1-\lambda)^i$ for all (s, a) . Since the sampling across clients and across iterations are independent, via invoking Hoeffding's inequality, for any given $\delta \in (0, 1)$, with probability at least $1 - \delta$,

$$(I.2) = \left\| \sum_{i=0}^t (1-\lambda)^i \lambda \gamma \frac{1}{K} \sum_{k=1}^K (\bar{P} - \tilde{P}_{t-i}^k) V^* \right\|_\infty \leq \frac{\gamma}{1-\gamma} \sqrt{\frac{1}{K} \lambda \log \frac{|\mathcal{S}||\mathcal{A}|TK}{\delta}}. \quad (28)$$

Bounding (I.3). To bound the third term (I.3) in Eq. (26), following the roadmap of Woo et al. (2023), we divide the summation into two parts as follows. For any $\beta E \leq t \leq T$, we have

$$\begin{aligned}
(\text{I.3}) &= \sum_{i=0}^t (1-\lambda)^i \lambda \gamma \left\| \frac{1}{K} \sum_{k=1}^K \tilde{P}_{t-i}^k (V^* - V_{t-i}^k) \right\|_{\infty} \\
&= \sum_{i=0}^t (1-\lambda)^{t-i} \lambda \gamma \left\| \frac{1}{K} \sum_{k=1}^K \tilde{P}_i^k (V^* - V_i^k) \right\|_{\infty} \\
&= \sum_{i=0}^{\chi(t)-\beta E} (1-\lambda)^{t-i} \lambda \gamma \left\| \frac{1}{K} \sum_{k=1}^K \tilde{P}_i^k (V^* - V_i^k) \right\|_{\infty} + \sum_{i=\chi(t)-\beta E+1}^t (1-\lambda)^{t-i} \lambda \gamma \left\| \frac{1}{K} \sum_{k=1}^K \tilde{P}_i^k (V^* - V_i^k) \right\|_{\infty} \\
&\leq \frac{\gamma}{1-\gamma} (1-\lambda)^{t-\chi(t)+\beta E} + \sum_{i=\chi(t)-\beta E+1}^t (1-\lambda)^{t-i} \lambda \gamma \left\| \frac{1}{K} \sum_{k=1}^K \tilde{P}_i^k (V^* - V_i^k) \right\|_{\infty}.
\end{aligned}$$

By Lemma 3,

$$\begin{aligned}
&\sum_{i=\chi(t)-\beta E+1}^t (1-\lambda)^{t-i} \lambda \gamma \left\| \frac{1}{K} \sum_{k=1}^K \tilde{P}_i^k (V^* - V_i^k) \right\|_{\infty} \\
&\leq \sum_{i=\chi(t)-\beta E+1}^t (1-\lambda)^{t-i} \lambda \gamma \left(\|\Delta_{\chi(i)}\|_{\infty} + 2\lambda \frac{1}{K} \sum_{k=1}^K \sum_{j=\chi(i)}^{i-1} \|\Delta_{t'}^k\|_{\infty} \right. \\
&\quad \left. + \gamma \lambda \frac{1}{K} \sum_{k=1}^K \max_{s,a} \left| \sum_{j=\chi(i)}^{i-1} \left(\tilde{P}_j^k(s,a) - \bar{P}(s,a) \right) V^* \right| \right).
\end{aligned}$$

Since $\tilde{P}_j^k(s,a)$'s are independent across time j and across state action pair (s,a) , and $|\tilde{P}_j^k(s,a) - \bar{P}(s,a)V^*| \leq \frac{1}{1-\gamma}$ (from Lemma 2), with Hoeffding's inequality and union bound, we get for any $\delta \in (0,1)$, with probability at least $1 - \delta$,

$$\left| \sum_{j=\chi(i)}^{i-1} \left(\tilde{P}_j^k(s,a) - \bar{P}(s,a) \right) V^* \right| \leq (E-1) \frac{1}{1-\gamma} \kappa + \frac{1}{1-\gamma} \sqrt{(E-1) \log \frac{|\mathcal{S}||\mathcal{A}|KT}{\delta}} \quad (29)$$

for all $(s,a) \in \mathcal{S} \times \mathcal{A}$, $k \in K$, and i . By Lemma 4, with probability at least $(1 - \delta)$, we have

$$\begin{aligned}
&\sum_{i=\chi(t)-\beta E+1}^t (1-\lambda)^{t-i} \lambda \gamma 2\lambda \frac{1}{K} \sum_{k=1}^K \sum_{j=\chi(i)}^{i-1} \|\Delta_j^k\|_{\infty} \\
&\leq 2\lambda^2 \gamma \sum_{i=\chi(t)-\beta E+1}^t (1-\lambda)^{t-i} \frac{1}{K} \sum_{k=1}^K \sum_{j=\chi(i)}^{i-1} \left(\|\Delta_{\chi(i)}\|_{\infty} + 3 \frac{\gamma}{1-\gamma} \lambda (E-1) \kappa + 3 \frac{\gamma}{1-\gamma} \sqrt{\lambda \log \frac{|\mathcal{S}||\mathcal{A}|KT}{\delta}} \right) \\
&\leq 2\lambda \gamma (E-1) \max_{\chi(t)-\beta E \leq i \leq t} \|\Delta_{\chi(i)}\|_{\infty} + \frac{6\gamma^2 \lambda^2}{1-\gamma} (E-1)^2 \kappa + \frac{6\gamma^2 \lambda}{1-\gamma} (E-1) \sqrt{\lambda \log \frac{|\mathcal{S}||\mathcal{A}|KT}{\delta}}.
\end{aligned}$$

Thus, by applying the union bound, we get with probability at least $(1 - 2\delta)$,

$$\begin{aligned}
& \sum_{i=\chi(t)-\beta E+1}^t (1-\lambda)^{t-i} \lambda \gamma \left\| \frac{1}{K} \sum_{k=1}^K \tilde{P}_i^k (V^* - V_i^k) \right\|_{\infty} \\
& \leq \gamma \max_{\chi(t)-\beta E \leq i \leq t} \|\Delta_{\chi(i)}\|_{\infty} + 2\lambda\gamma(E-1) \max_{\chi(t)-\beta E \leq i \leq t} \|\Delta_{\chi(i)}\|_{\infty} + \frac{6\gamma^2\lambda^2}{1-\gamma}(E-1)^2\kappa \\
& \quad + \frac{6\gamma^2\lambda}{1-\gamma}(E-1) \sqrt{\lambda \log \frac{|\mathcal{S}||\mathcal{A}|KT}{\delta}} \\
& \quad + \sum_{i=\chi(t)-\beta E+1}^t (1-\lambda)^{t-i} \lambda \gamma \left(\frac{\gamma\lambda}{1-\gamma}(E-1)\kappa + \frac{\gamma\lambda}{1-\gamma} \sqrt{(E-1) \log \frac{|\mathcal{S}||\mathcal{A}|KT}{\delta}} \right) \\
& = \gamma(1+2\lambda(E-1)) \max_{\chi(t)-\beta E \leq i \leq t} \|\Delta_{\chi(i)}\|_{\infty} + \frac{\gamma^2}{1-\gamma}(6\lambda^2(E-1)^2 + \lambda(E-1))\kappa \\
& \quad + \frac{\gamma^2\lambda}{1-\gamma} \sqrt{(E-1) \log \frac{|\mathcal{S}||\mathcal{A}|KT}{\delta}} + \frac{6\gamma^2\lambda}{1-\gamma}(E-1) \sqrt{\lambda \log \frac{|\mathcal{S}||\mathcal{A}|KT}{\delta}}.
\end{aligned}$$

Hence, the third term (I.3) in Eq. (26) can be bounded as

$$\begin{aligned}
\text{(I.3)} & = \sum_{i=0}^t (1-\lambda)^i \lambda \gamma \left\| \frac{1}{K} \sum_{k=1}^K \tilde{P}_i^k (V^* - V_i^k) \right\|_{\infty} \\
& \leq \frac{\gamma}{1-\gamma} (1-\lambda)^{t-\chi(t)+\beta E} + \gamma(1+2\lambda(E-1)) \max_{\chi(t)-\beta E \leq i \leq t} \|\Delta_{\chi(i)}\|_{\infty} + \frac{\gamma^2}{1-\gamma}(6\lambda^2(E-1)^2 + \lambda(E-1))\kappa \\
& \quad + \frac{\gamma^2\lambda}{1-\gamma} \sqrt{(E-1) \log \frac{|\mathcal{S}||\mathcal{A}|KT}{\delta}} + \frac{6\gamma^2\lambda}{1-\gamma}(E-1) \sqrt{\lambda \log \frac{|\mathcal{S}||\mathcal{A}|KT}{\delta}}. \tag{30}
\end{aligned}$$

Combing the bounds of (I.1), (I.2), and (I.3) in Eq. (26). Combing the bounds for terms (27), (28), and (30), we get the following recursion holds for all rounds T with probability at least $(1 - 3\delta)$:

$$\begin{aligned}
\|\Delta_{t+1}\|_{\infty} & \leq (1-\lambda)^{t+1} \frac{1}{1-\gamma} + \frac{\gamma}{1-\gamma} \sqrt{\frac{1}{K} \lambda \log \frac{|\mathcal{S}||\mathcal{A}|TK}{\delta}} + \frac{\gamma}{1-\gamma} (1-\lambda)^{t-\chi(t)+\beta E} \\
& \quad + \gamma(1+2\lambda(E-1)) \max_{\chi(t)-\beta E \leq i \leq t} \|\Delta_{\chi(i)}\|_{\infty} + \frac{\gamma^2}{1-\gamma}(6\lambda^2(E-1)^2 + \lambda(E-1))\kappa \\
& \quad + \frac{\gamma^2\lambda}{1-\gamma} \sqrt{(E-1) \log \frac{|\mathcal{S}||\mathcal{A}|KT}{\delta}} + \frac{6\gamma^2\lambda}{1-\gamma}(E-1) \sqrt{\lambda \log \frac{|\mathcal{S}||\mathcal{A}|KT}{\delta}} \\
& \leq \gamma(1+2\lambda(E-1)) \max_{\chi(t)-\beta E \leq i \leq t} \|\Delta_{\chi(i)}\|_{\infty} + \frac{2}{1-\gamma} (1-\lambda)^{\beta E} + \frac{\gamma^2}{1-\gamma}(6\lambda^2(E-1)^2 + \lambda(E-1))\kappa \\
& \quad + \frac{\gamma^2\lambda}{1-\gamma} \sqrt{(E-1) \log \frac{|\mathcal{S}||\mathcal{A}|KT}{\delta}} + \frac{6\gamma^2\lambda}{1-\gamma}(E-1) \sqrt{\lambda \log \frac{|\mathcal{S}||\mathcal{A}|KT}{\delta}} \\
& \quad + \frac{\gamma}{1-\gamma} \sqrt{\frac{1}{K} \lambda \log \frac{|\mathcal{S}||\mathcal{A}|TK}{\delta}}.
\end{aligned}$$

Let

$$\begin{aligned}
\rho & := \frac{2}{1-\gamma} (1-\lambda)^{\beta E} + \frac{\gamma^2}{1-\gamma}(6\lambda^2(E-1)^2 + \lambda(E-1))\kappa \\
& \quad + \frac{\gamma^2\lambda}{1-\gamma} \sqrt{(E-1) \log \frac{|\mathcal{S}||\mathcal{A}|KT}{\delta}} + \frac{6\gamma^2\lambda}{1-\gamma}(E-1) \sqrt{\lambda \log \frac{|\mathcal{S}||\mathcal{A}|KT}{\delta}} \\
& \quad + \frac{\gamma}{1-\gamma} \sqrt{\frac{1}{K} \lambda \log \frac{|\mathcal{S}||\mathcal{A}|TK}{\delta}}. \tag{31}
\end{aligned}$$

With the assumption that $\lambda \leq \frac{1-\gamma}{4\gamma(E-1)}$, the above recursion can be written as

$$\|\Delta_{t+1}\|_\infty \leq \frac{1+\gamma}{2} \max_{\chi(t)-\beta E \leq i \leq t} \|\Delta_{\chi(i)}\|_\infty + \rho.$$

Unrolling the above recursion L times where $L\beta E \leq t \leq T$, we obtain that

$$\begin{aligned} \|\Delta_{t+1}\|_\infty &\leq \left(\frac{1+\gamma}{2}\right)^L \max_{\chi(t)-L\beta E \leq i \leq t} \|\Delta_{\chi(i)}\|_\infty + \sum_{i=0}^{L-1} \left(\frac{1+\gamma}{2}\right)^i \rho \\ &\leq \left(\frac{1+\gamma}{2}\right)^L \frac{1}{1-\gamma} + \frac{2}{1-\gamma} \rho. \end{aligned}$$

Choosing $\beta = \left\lfloor \frac{1}{E} \sqrt{\frac{(1-\gamma)T}{2\lambda}} \right\rfloor$, $L = \left\lceil \sqrt{\frac{\lambda T}{1-\gamma}} \right\rceil$, $t+1 = T$, we get

$$\begin{aligned} \|\Delta_T\|_\infty &\leq \frac{1}{1-\gamma} \left(\frac{1+\gamma}{2}\right)^{\sqrt{\frac{\lambda T}{1-\gamma}}} + \frac{2}{1-\gamma} \left(\frac{2}{1-\gamma} (1-\lambda)^{\beta E} + \frac{\gamma^2}{1-\gamma} (6\lambda^2(E-1)^2 + \lambda(E-1))\kappa \right. \\ &\quad \left. + \left(\frac{6\gamma^2\lambda}{1-\gamma} \sqrt{E-1} + \frac{\gamma^2\sqrt{\lambda}}{1-\gamma} \right) \sqrt{\lambda(E-1) \log \frac{|\mathcal{S}||\mathcal{A}|KT}{\delta}} + \frac{\gamma}{1-\gamma} \sqrt{\frac{1}{K} \lambda \log \frac{|\mathcal{S}||\mathcal{A}|TK}{\delta}} \right) \\ &\leq \frac{1}{1-\gamma} \exp \left\{ -\frac{1}{2} \sqrt{(1-\gamma)\lambda T} \right\} + \frac{4}{(1-\gamma)^2} \exp \left\{ -\frac{1}{2} \sqrt{(1-\gamma)\lambda T} \right\} \\ &\quad + \frac{2\gamma^2}{(1-\gamma)^2} (6\lambda^2(E-1)^2 + \lambda(E-1))\kappa \\ &\quad + \left(\frac{12\gamma^2\lambda}{(1-\gamma)^2} \sqrt{E-1} + \frac{2\gamma^2\sqrt{\lambda}}{(1-\gamma)^2} \right) \sqrt{\lambda(E-1) \log \frac{|\mathcal{S}||\mathcal{A}|KT}{\delta}} + \frac{2\gamma}{(1-\gamma)^2} \sqrt{\frac{1}{K} \lambda \log \frac{|\mathcal{S}||\mathcal{A}|TK}{\delta}} \\ &\leq \frac{4}{(1-\gamma)^2} \exp \left\{ -\frac{1}{2} \sqrt{(1-\gamma)\lambda T} \right\} + \frac{2\gamma^2}{(1-\gamma)^2} (6\lambda^2(E-1)^2 + \lambda(E-1))\kappa \\ &\quad + \left(\frac{14\gamma^2\lambda}{(1-\gamma)^2} \sqrt{E-1} \right) \sqrt{\log \frac{|\mathcal{S}||\mathcal{A}|KT}{\delta}} + \frac{2\gamma}{(1-\gamma)^2} \sqrt{\frac{1}{K} \lambda \log \frac{|\mathcal{S}||\mathcal{A}|TK}{\delta}}, \end{aligned}$$

where the second inequality follows from

$$\begin{aligned} \left(\frac{1+\gamma}{2}\right)^{\sqrt{\frac{\lambda T}{1-\gamma}}} &= \left(1 - \frac{1-\gamma}{2}\right)^{\sqrt{\frac{\lambda T}{1-\gamma}}} \leq \exp \left\{ -\frac{1}{2} \sqrt{(1-\gamma)\lambda T} \right\}, \\ (1-\lambda)^{\beta E} &\leq \exp \left\{ -\lambda \sqrt{\frac{(1-\gamma)T}{2\lambda}} \right\} \leq \exp \left\{ -\frac{1}{2} \sqrt{(1-\gamma)\lambda T} \right\}. \end{aligned}$$

By the assumption that $(E-1) \leq \frac{1}{K\lambda}$, the above can be further simplified as

$$\|\Delta_T\|_\infty \leq \frac{4}{(1-\gamma)^2} \exp \left\{ -\frac{1}{2} \sqrt{(1-\gamma)\lambda T} \right\} + \frac{14\gamma^2}{(1-\gamma)^2} \lambda(E-1)\kappa + \frac{16}{(1-\gamma)^2} \sqrt{\frac{\lambda}{K} \log \frac{|\mathcal{S}||\mathcal{A}|KT}{\delta}}.$$

F Proof of Theorem 2

Let $|\mathcal{A}| = 1$, in which case Q -function coincides with the V -function. According to Algorithm 1, when $(t+1) \bmod E \neq 0$, we have

$$Q_{t+1}^k = ((1-\lambda)I + \lambda\gamma P^k) Q_t^k + \lambda R.$$

Define $A^k \triangleq (1-\lambda)I + \lambda\gamma P^k$. We obtain the following recursion between two synchronization rounds:

$$Q_{(z+1)E}^k = (A^k)^E Q_{zE}^k + ((A^k)^0 + \dots + (A^k)^{E-1}) \lambda R.$$

Define

$$\bar{A}^{(\ell)} \triangleq \frac{1}{K} \sum_{k=1}^K (A^k)^\ell. \quad (32)$$

Note that Q^* is the fixed point under the transition kernel \bar{P} , we have $\lambda R = \lambda(I - \gamma\bar{P})Q^* = (I - \bar{A}^{(1)})Q^*$ since $\bar{A}^{(1)} = I - \lambda(I - \gamma\bar{P})$. Furthermore, since $Q_{tE}^1, \dots, Q_{tE}^K$ are identical due to synchronization, we get

$$\bar{Q}_{(z+1)E} = \bar{A}^{(E)} \bar{Q}_{zE} + \left(I + \bar{A}^{(1)} + \dots + \bar{A}^{(E-1)} \right) (I - \bar{A}^{(1)}) Q^*.$$

Consequently,

$$\begin{aligned} \Delta_{(z+1)E} &= Q^* - \bar{Q}_{(z+1)E} \\ &= \bar{A}^{(E)} \Delta_{zE} + \left((I - \bar{A}^{(E)}) - \left(I + \bar{A}^{(1)} + \dots + \bar{A}^{(E-1)} \right) (I - \bar{A}^{(1)}) \right) Q^*. \end{aligned} \quad (33)$$

Next, consider $|\mathcal{S}| = 2$ and even K with

$$P^{2k-1} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad P^{2k} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, \quad \text{for } k \in \mathbb{N}.$$

Then $\bar{P} = \frac{1}{2} \mathbf{1}\mathbf{1}^\top$, where $\mathbf{1}$ denotes the all ones vector. For the above transition kernels, we have

$$\frac{1}{K} \sum_{k=1}^K (P^k)^\ell = \begin{cases} I, & \ell \text{ even,} \\ \bar{P}, & \ell \text{ odd.} \end{cases}$$

Applying the definition of $\bar{A}^{(\ell)}$ in (32) yields that

$$\begin{aligned} \bar{A}^{(\ell)} &= \frac{1}{K} \sum_{k=1}^K (A^k)^\ell \\ &= \frac{1}{K} \sum_{k=1}^K ((1-\lambda)I + \lambda\gamma P^k)^\ell \\ &= \frac{1}{K} \sum_{k=1}^K \sum_{j=0}^{\ell} \binom{\ell}{j} (\lambda\gamma P^k)^j ((1-\lambda)I)^{\ell-j} \\ &= \sum_{j \text{ even}} \binom{\ell}{j} (1-\lambda)^{\ell-j} (\lambda\gamma)^j (I - \bar{P} + \bar{P}) + \sum_{j \text{ odd}} \binom{\ell}{j} (1-\lambda)^{\ell-j} (\lambda\gamma)^j \bar{P} \\ &= \underbrace{\frac{1}{2} ((1-\lambda - \lambda\gamma)^\ell + (1-\lambda + \lambda\gamma)^\ell)}_{\triangleq \alpha_\ell} (I - \bar{P}) + \underbrace{(1-\lambda + \lambda\gamma)^\ell}_{\triangleq \beta_\ell} \bar{P} \\ &= \alpha_\ell (I - \bar{P}) + \beta_\ell \bar{P}, \end{aligned}$$

which is the eigen-decomposition of $\bar{A}^{(\ell)}$. Let

$$\lambda_1 \triangleq (1+\gamma)\lambda, \lambda_2 \triangleq (1-\gamma)\lambda, \quad \nu_1 = 1 - \lambda_1, \nu_2 = 1 - \lambda_2.$$

Then

$$\alpha_\ell = \frac{1}{2}(\nu_1^\ell + \nu_2^\ell), \quad \beta_\ell = \nu_2^\ell. \quad (34)$$

Note that $0 \leq \alpha \leq \beta \leq 1$ and $I - \bar{P}$ and \bar{P} are orthogonal projection matrices satisfying $(I - \bar{P})\bar{P} = 0$. The matrices for the second term of the error on the right-hand side of 33 reduce to

$$\begin{aligned} & \left(I + \bar{A}^{(1)} + \dots + \bar{A}^{(E-1)} \right) \left(I - \bar{A}^{(1)} \right) \\ &= \left(\sum_{\ell=0}^{E-1} \alpha_\ell (I - \bar{P}) + \sum_{\ell=0}^{E-1} \beta_\ell \bar{P} \right) \left((\alpha_0 - \alpha_1)(I - \bar{P}) + (\beta_0 - \beta_1)\bar{P} \right) \\ &= \left((1 - \alpha_1) \sum_{\ell=0}^{E-1} \alpha_\ell (I - \bar{P})^2 + (1 - \beta_1) \sum_{\ell=0}^{E-1} \beta_\ell \bar{P}^2 \right) \text{ since } \alpha_0 = \beta_0 = 1 \\ &= \left((1 - \alpha_1) \sum_{\ell=0}^{E-1} \alpha_\ell (I - \bar{P}) + (1 - \beta_1) \sum_{\ell=0}^{E-1} \beta_\ell \bar{P} \right) \text{ since } (I - \bar{P}) \text{ and } \bar{P} \text{ are idempotent.} \end{aligned}$$

It follow that

$$\begin{aligned} & \left(I - \bar{A}^{(E)} \right) - \left(I + \bar{A}^{(1)} + \dots + \bar{A}^{(E-1)} \right) \left(I - \bar{A}^{(1)} \right) \\ &= \underbrace{\left((1 - \alpha_E) - (1 - \alpha_1) \left(\sum_{i=0}^{E-1} \alpha_i \right) \right)}_{\triangleq \kappa_E} (I - \bar{P}) + \underbrace{\left((1 - \beta_E) - (1 - \beta_1) \left(\sum_{i=0}^{E-1} \beta_i \right) \right)}_{=0} \bar{P}. \end{aligned}$$

Applying (34) yields that

$$\kappa_E = -\frac{\gamma}{2} \left(\frac{1 - \nu_2^E}{1 - \gamma} - \frac{1 - \nu_1^E}{1 + \gamma} \right). \quad (35)$$

It follows from (33) that the error evolves as

$$\Delta_{(z+1)E} = (\alpha_E(I - \bar{P}) + \beta_E \bar{P}) \Delta_{zE} + \kappa_E (I - \bar{P}) Q^*,$$

which further yields the following full recursion of the error:

$$\begin{aligned} \Delta_{zE} &= (\alpha_E(I - \bar{P}) + \beta_E \bar{P})^z \Delta_0 + \sum_{\ell=0}^{z-1} (\alpha_E(I - \bar{P}) + \beta_E \bar{P})^\ell \kappa_E (I - \bar{P}) Q^* \\ &= (\alpha_E^z (I - \bar{P}) + \beta_E^z \bar{P}) \Delta_0 + \sum_{\ell=0}^{z-1} (\alpha_E^\ell (I - \bar{P}) + \beta_E^\ell \bar{P}) \kappa_E (I - \bar{P}) Q^* \\ &\quad \text{since } (\alpha_E(I - \bar{P}) + \beta_E \bar{P})^\ell = \alpha_E^\ell (I - \bar{P}) + \beta_E^\ell \bar{P}, \forall \ell \in \mathbb{N} \\ &= (\alpha_E^z (I - \bar{P}) + \beta_E^z \bar{P}) \Delta_0 + \frac{1 - \alpha_E^z}{1 - \alpha_E} \kappa_E (I - \bar{P}) Q^* \\ &= \left(\alpha_E^z + \frac{1 - \alpha_E^z}{1 - \alpha_E} \kappa_E \right) (I - \bar{P}) Q^* + \beta_E^z \bar{P} Q^*, \end{aligned}$$

where the last equality applied the zero initialization condition.

Note that $(I - \bar{P})Q^*$ and $\bar{P}Q^*$ are orthogonal vectors. Since $|\mathcal{S}| = 2$, we have

$$\|\Delta_{zE}\|_\infty \geq \frac{1}{\sqrt{2}} \|\Delta_{zE}\|_2 \geq \frac{\min\{\|(I - \bar{P})Q^*\|_2, \|\bar{P}Q^*\|_2\}}{\sqrt{2}} \cdot \max\left\{ \left| \alpha_E^z + \frac{1 - \alpha_E^z}{1 - \alpha_E} \kappa_E \right|, \beta_E^z \right\}.$$

Let $R = \begin{bmatrix} r_1 \\ r_2 \end{bmatrix}$, since $Q^* = (I - \gamma\bar{P})^{-1}R = (I - \bar{P})R + \frac{1}{1-\gamma}\bar{P}R$, we obtain that

$$(I - \bar{P})Q^* = (I - \bar{P})R = \frac{1}{2} \begin{bmatrix} r_1 - r_2 \\ r_2 - r_1 \end{bmatrix}, \quad \bar{P}Q^* = \frac{1}{1-\gamma}\bar{P}R = \frac{1}{2(1-\gamma)} \begin{bmatrix} r_1 + r_2 \\ r_1 + r_2 \end{bmatrix}.$$

$$\|(I - \bar{P})Q^*\|_2 = \frac{\sqrt{2}}{2} |r_1 - r_2|, \quad \|\bar{P}Q^*\|_2 = \frac{\sqrt{2}}{2(1-\gamma)} |r_1 + r_2|.$$

When $r_1 = r_2$, the error Δ_{zE} reduces to $\beta_E^z \bar{P}Q^*$, and $\|\Delta_{zE}\|_\infty = \frac{1}{2(1-\gamma)} |r_1 + r_2| |\beta_E^z|$; otherwise, $\min\{\|(I - \bar{P})Q^*\|_2, \|\bar{P}Q^*\|_2\} = \frac{\sqrt{2}}{2} \min\{|r_1 - r_2|, \frac{1}{1-\gamma} |r_1 + r_2|\}$. It remains to analyze the coefficients as functions of λ . To this end, we introduce the following lemma:

Lemma 5. *The following properties hold:*

1. *Negativity:* $\kappa_E < 0$;
2. *Monotonicity:* $\frac{\kappa_E}{1-\alpha_E}$ is monotonically decreasing for $\lambda \in (0, \frac{1}{1+\gamma})$;
3. *Upper bound:* $|\frac{\kappa_E}{1-\alpha_E}| \leq \frac{\gamma^2}{1-\gamma^2}$ for $\lambda \in (0, \frac{1}{1+\gamma})$;
4. *Lower bound:* if $(1+\gamma)\lambda \leq \frac{1}{2E}$, then $|\frac{\kappa_E}{1-\alpha_E}| \geq \frac{\lambda\gamma^2(E-1)}{4}$.

Proof. We prove the properties separately.

1. Note that $\nu_1 < \nu_2$, $1 - \nu_1 = (1 + \gamma)\lambda$, and $1 - \nu_2 = (1 - \gamma)\lambda$. Then it follows from (35) that

$$\kappa_E = -\frac{\lambda\gamma}{2} \sum_{i=1}^{E-1} (\nu_2^i - \nu_1^i) < 0.$$

2. For the monotonicity, it suffices to show that $\frac{d}{d\lambda} \frac{\kappa_E}{1-\alpha_E} \leq 0$. We calculate the derivative as

$$\frac{d}{d\lambda} \frac{\kappa_E}{1-\alpha_E} = \frac{\gamma E (1 - \nu_1^E) (1 - \nu_2^E)}{2(1-\gamma^2)(1-\alpha_E)^2} \left(\frac{(1+\gamma)\nu_1^{E-1}}{1-\nu_1^E} - \frac{(1-\gamma)\nu_2^{E-1}}{1-\nu_2^E} \right).$$

Note that

$$\frac{(1+\gamma)\nu_1^{E-1}}{1-\nu_1^E} - \frac{(1-\gamma)\nu_2^{E-1}}{1-\nu_2^E} = \frac{1}{\lambda} \left(\frac{\nu_1^{E-1}}{1+\nu_1+\dots+\nu_1^{E-1}} - \frac{\nu_2^{E-1}}{1+\nu_2+\dots+\nu_2^{E-1}} \right) \leq 0.$$

3. For the upper bound, it suffices to show the result at $\lambda = \frac{1}{1+\gamma}$ due to the negativity and monotonicity. At $\lambda = \frac{1}{1+\gamma}$, we have

$$\left| \frac{\kappa_E}{1-\alpha_E} \right| = \frac{\gamma}{1-\gamma^2} \left(\gamma - \frac{(\frac{2\gamma}{1+\gamma})^E}{2 - (\frac{2\gamma}{1+\gamma})^E} \right) \leq \frac{\gamma^2}{1-\gamma^2}.$$

4. For the lower bound, the case $E = 1$ trivially holds. Next, consider $E \geq 2$. We have

$$\begin{aligned} \frac{\kappa_E}{1-\alpha_E} &= -\frac{\gamma}{1-\gamma^2} \frac{(1+\gamma)(1-\nu_2^E) - (1-\gamma)(1-\nu_1^E)}{(1-\nu_1^E) + (1-\nu_2^E)} \\ &= -\lambda\gamma \frac{\sum_{\ell=1}^{E-1} (\nu_2^\ell - \nu_1^\ell)}{(1-\nu_1^E) + (1-\nu_2^E)}. \end{aligned}$$

Note that $1 - nx \leq (1 - x)^n \leq 1 - \frac{1}{2}nx$ for $n \geq 1$ and $0 \leq x \leq \frac{1}{n}$. Then, for $(1 + \gamma)\lambda \leq \frac{1}{2E}$, we have

$$\begin{aligned}\nu_1^E &= (1 - (1 + \gamma)\lambda)^E \geq 1 - (1 + \gamma)\lambda E \geq \frac{1}{2}, \\ \nu_2^E &= (1 - (1 - \gamma)\lambda)^E \geq 1 - (1 - \gamma)\lambda E.\end{aligned}$$

Moreover, for all $x \in [\nu_1, \nu_2] \subseteq [0, 1]$ and $\ell - 1 \leq E$, we have

$$x^{\ell-1} \geq x^E \geq \nu_1^E \geq \frac{1}{2}.$$

We obtain that

$$\frac{\sum_{\ell=1}^{E-1} (\nu_2^\ell - \nu_1^\ell)}{(1 - \nu_1^E) + (1 - \nu_2^E)} \geq \frac{\sum_{\ell=1}^{E-1} \int_{\nu_1}^{\nu_2} \ell \cdot x^{\ell-1} dx}{2\lambda E} \geq \frac{\sum_{\ell=1}^{E-1} \ell \frac{1}{2} (\nu_2 - \nu_1)}{2\lambda E} = \frac{1}{4}\gamma(E - 1).$$

The proof is completed. \square

We consider two regimes of the stepsize separated by $\lambda_0 \triangleq \frac{\log z}{(1-\gamma)zE} < \frac{1}{1+\gamma}$, where the dominating error is due to the small stepsize and the environment heterogeneity, respectively:

Slow rate due to small stepsize when $\lambda \leq \lambda_0$. Since β_E^z monotonically decreases as λ increases,

$$\beta_E^z = (1 - (1 - \gamma)\lambda)^{zE} \geq (1 - (1 - \gamma)\lambda_0)^{zE} = \left(1 - \frac{\log z}{zE}\right)^{zE}.$$

Note that $\frac{\log z}{zE} \in (0, \frac{1}{2})$, applying the fact $\log(1 - x) + x \geq -x^2$ for $x \in [0, \frac{1}{2}]$ yields that

$$\log\left(1 - \frac{\log z}{zE}\right) + \frac{\log z}{zE} \geq -\left(\frac{\log z}{zE}\right)^2 \geq -\frac{1}{zE}.$$

Then we get

$$\beta_E^z \geq \left(1 - \frac{\log z}{zE}\right)^{zE} \geq \frac{1}{ez}.$$

Slow rate due to environment heterogeneity when $\lambda \geq \lambda_0$. Recall that $\lambda < \frac{1}{1+\gamma}$. Applying the triangle inequality yields that

$$\left|\alpha_E^z + \frac{1 - \alpha_E^z}{1 - \alpha_E} \kappa_E\right| \geq \left|\frac{\kappa_E}{1 - \alpha_E}\right| - \left(1 + \left|\frac{\kappa_E}{1 - \alpha_E}\right|\right) \alpha_E^z.$$

For the first term, by the negativity and monotonicity in Lemma 5, it suffices to show the lower bound at $\lambda = \lambda_0$. Since $\lambda < \frac{1}{1+\gamma}$, then $\alpha_E = \frac{1}{2}((1 - (1 - \gamma)\lambda)^E + (1 - (1 + \gamma)\lambda)^E)$ decreases as λ increases. For $z \geq \exp\left\{-W_{-1}\left(-\frac{1-\gamma}{2(1+\gamma)}\right)\right\}$, where W_{-1} is the Lambert W function, such that $(1 + \gamma)\lambda_0 \leq \frac{1}{2E}$, we apply the lower bound in Lemma 5 and obtain that

$$\left|\frac{\kappa_E}{1 - \alpha_E}\right| \geq \frac{\lambda_0 \gamma^2 (E - 1)}{4} \geq \frac{\frac{\log z}{(1-\gamma)zE} \gamma^2 (E - 1)}{4} \geq \frac{(E - 1)}{4E} \gamma^2 \frac{\log z}{(1 - \gamma)z}.$$

Additionally, applying the upper bound in Lemma 5 yields

$$\left(1 + \left|\frac{\kappa_E}{1 - \alpha_E}\right|\right) \alpha_E^z \leq \frac{\nu_2^{zE}}{1 - \gamma^2} = \frac{(1 - (1 - \gamma)\lambda)^{zE}}{1 - \gamma^2} \leq \frac{1}{(1 - \gamma^2)z}.$$

Therefore,

$$\begin{aligned}
\left| \alpha_E^z + \frac{1 - \alpha_E^z}{1 - \alpha_E} \kappa_E \right| &\geq \left| \frac{\kappa_E}{1 - \alpha_E} \right| - \left(1 + \left| \frac{\kappa_E}{1 - \alpha_E} \right| \right) \alpha_E^z \\
&\geq \frac{(E-1)}{4E} \gamma^2 \frac{\log z}{(1-\gamma)z} - \frac{1}{(1-\gamma^2)z} \\
&= \frac{1}{(1-\gamma^2)z} \left((1+\gamma)\gamma^2 \log(z)(E-1)/(4E) - 1 \right) \\
&= \frac{1}{(1-\gamma)z} \left(\frac{\gamma^2 \log(z)(E-1) - 4E/(1+\gamma)}{4E} \right).
\end{aligned}$$

When $r_1 = r_2$,

$$\begin{aligned}
\|\Delta_{zE}\|_\infty &= \frac{|r_1 + r_2|}{2(1-\gamma)} |\beta_E^z| \\
&\geq \frac{|r_1 + r_2|}{2(1-\gamma)} \frac{E}{eT};
\end{aligned}$$

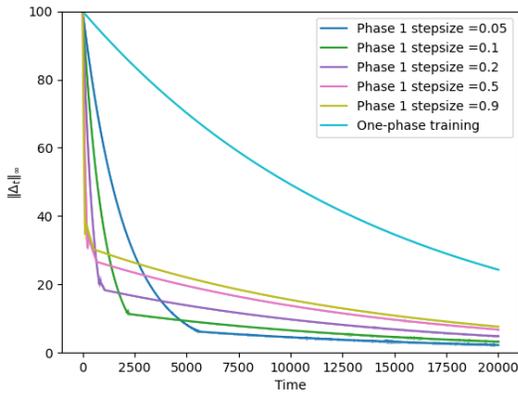
otherwise,

$$\begin{aligned}
\|\Delta_{zE}\|_\infty &\geq \frac{\min\{\|(I - \bar{P})Q^*\|_2, \|\bar{P}Q^*\|_2\}}{\sqrt{2}} \cdot \max \left\{ \left| \alpha_E^z + \frac{1 - \alpha_E^z}{1 - \alpha_E} \kappa_E \right|, \beta_E^z \right\} \\
&\geq \frac{1}{2} \min \left\{ |r_1 - r_2|, \frac{1}{1-\gamma} |r_1 + r_2| \right\} \max \left\{ \left| \alpha_E^z + \frac{1 - \alpha_E^z}{1 - \alpha_E} \kappa_E \right|, \beta_E^z \right\} \\
&\geq \frac{1}{2} \min \left\{ |r_1 - r_2|, \frac{1}{1-\gamma} |r_1 + r_2| \right\} \max \left\{ \frac{1}{(1-\gamma)z} \left(\frac{\gamma^2 \log(z)(E-1) - 4E/(1+\gamma)}{4E} \right), \frac{1}{ez} \right\} \\
&= \frac{1}{2} \min \left\{ |r_1 - r_2|, \frac{1}{1-\gamma} |r_1 + r_2| \right\} \max \left\{ \frac{E}{(1-\gamma)T} \left(\frac{\gamma^2 \log(z)(E-1) - 4E/(1+\gamma)}{4E} \right), \frac{E}{eT} \right\}.
\end{aligned}$$

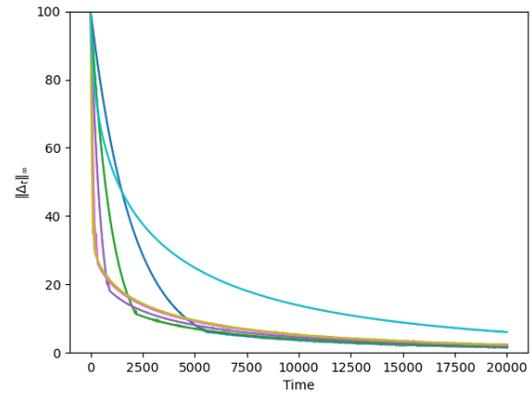
We can choose $\log(z) \geq \frac{4E(\gamma+2)}{(1+\gamma)\gamma^2(E-1)}$, $E \geq 2$ so that $\left(\frac{\gamma^2 \log(z)(E-1) - 4E/(1+\gamma)}{4E} \right) \geq 1$. Then the first term inside the max operator is bigger. Then,

$$\|\Delta_{zE}\|_\infty \geq \frac{1}{2} \min \left\{ |r_1 - r_2|, \frac{1}{1-\gamma} |r_1 + r_2| \right\} \frac{E}{(1-\gamma)T}.$$

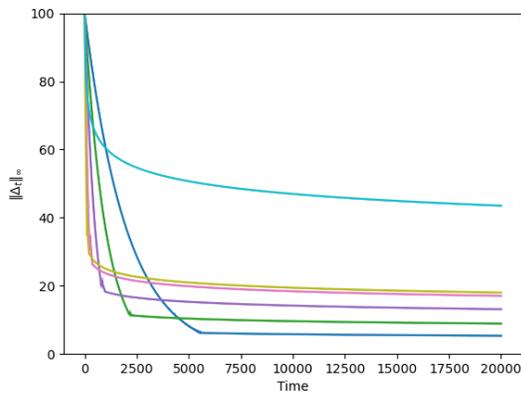
G Additional experiments



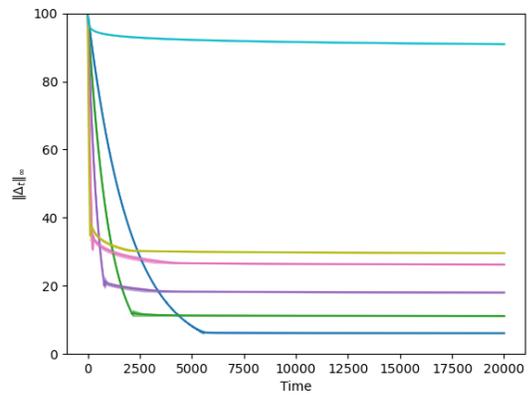
(a) Phase 2 stepsize $\lambda = \frac{1}{\sqrt{T}}$



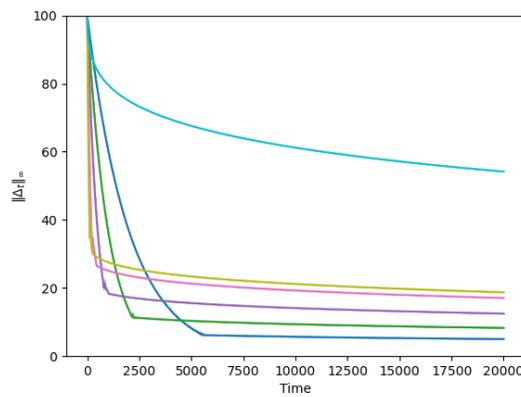
(b) Phase 2 stepsize $\lambda_t = \frac{1}{\sqrt{t+1}}$



(c) Phase 2 stepsize $\lambda_t = \frac{c+1}{t+c}$, where $c = 10$



(d) Phase 2 stepsize $\lambda_t = \frac{1}{t+1}$



(e) Phase 2 stepsize $\lambda_t = \frac{1}{(t+1)^{0.7}}$

Figure 4: Choosing different stepsizes for phases 1 and 2 leads to faster overall convergence. $E = 10$.

G.1 Impacts of E on homogeneous settings.

For the homogeneous settings, in addition to $E = 10$, we also consider $E = \{1, 20, 40, \infty\}$, where $E = \infty$ means no communication among the agents throughout the entire learning process. Similar to Figure 2b, there is no obvious two-phase phenomenon even in the extreme case when $E = \infty$. Also, though there is indeed performance degradation caused by larger E , the overall performance degradation is nearly negligible compared with the heterogeneous settings shown in Figures 2a and 3.

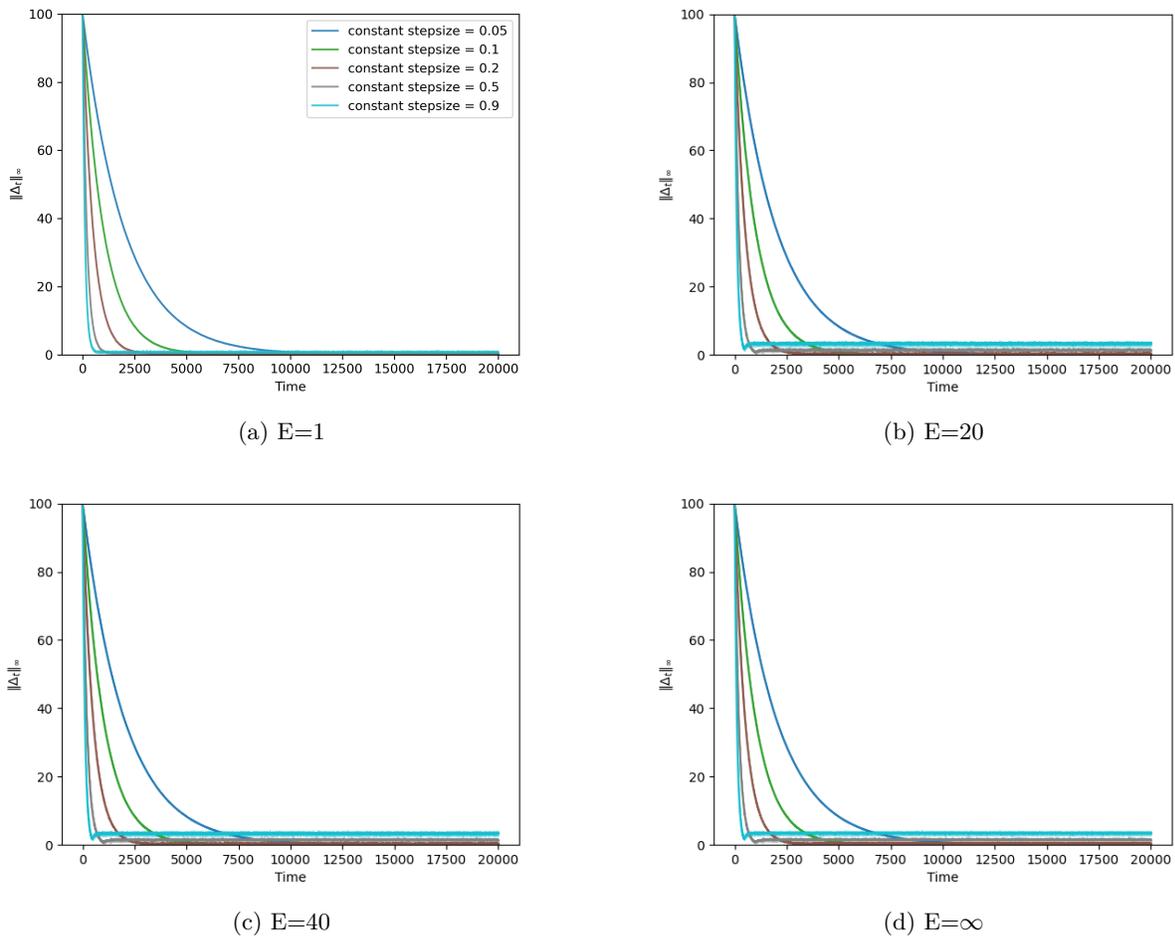


Figure 5: Homogeneous federated Q-learning with varying E .

G.2 Impacts of E on time-decaying stepsize

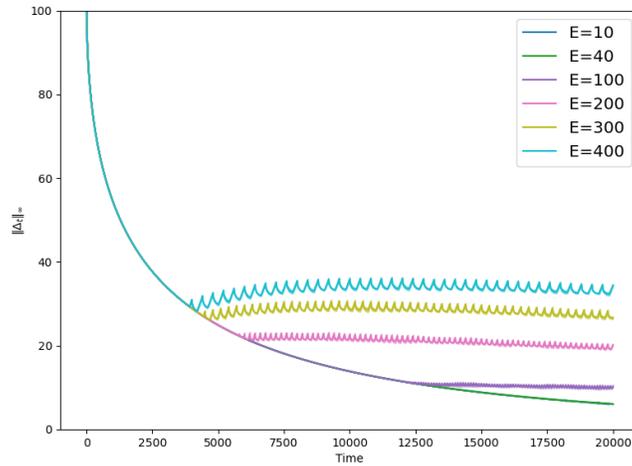


Figure 6: Using time-decaying stepsize $\lambda_t = \frac{1}{\sqrt{t+1}}$, the overall convergence becomes worse as E increases

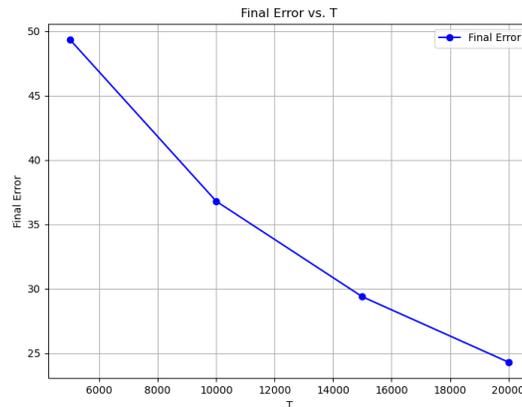
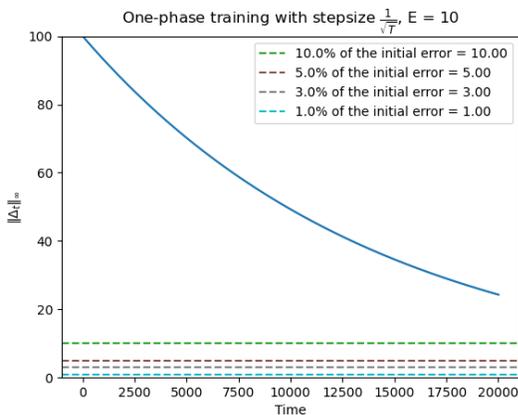


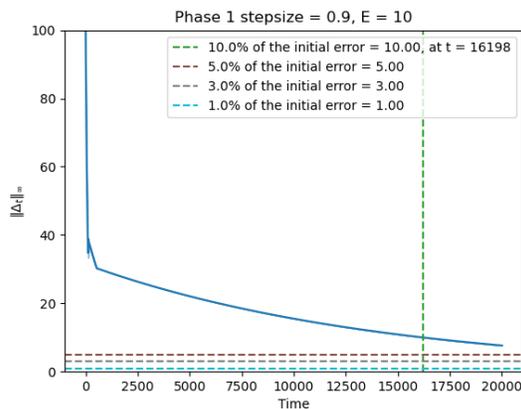
Figure 7: Final error versus T . It is clear that when choosing $\lambda = \frac{1}{\sqrt{T}}$, the final error decays as T increases.

G.3 Different target error levels.

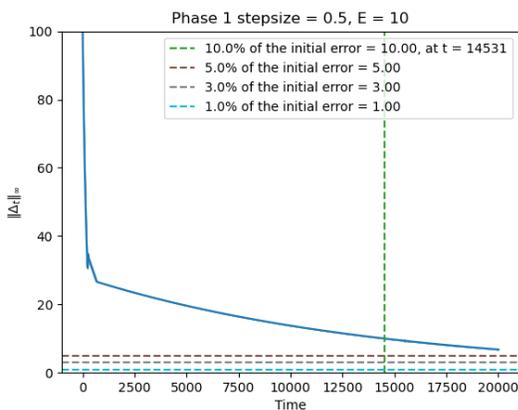
In Figure 8, we show the error levels that these training strategies can achieve within a time horizon $T = 20,000$. The tolerance levels are 10%, 5%, 3%, and 1% of the initial error $\|\Delta_0\|_\infty$, respectively. At a high level, choosing different stepsizes for phases 1 and 2 can speed up convergence.



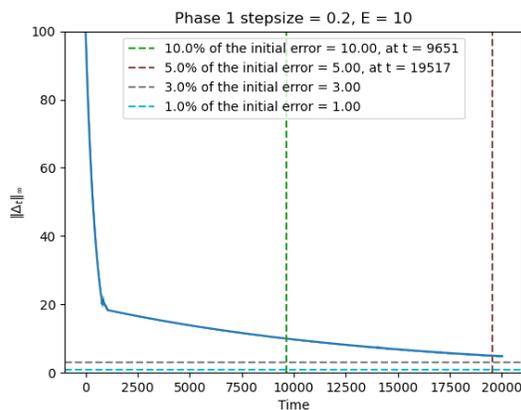
(a) One common $\lambda = \frac{1}{\sqrt{T}}$ throughout. $\|\Delta_t\|_\infty$ does not meet any of the tolerance levels within 20000 iterations



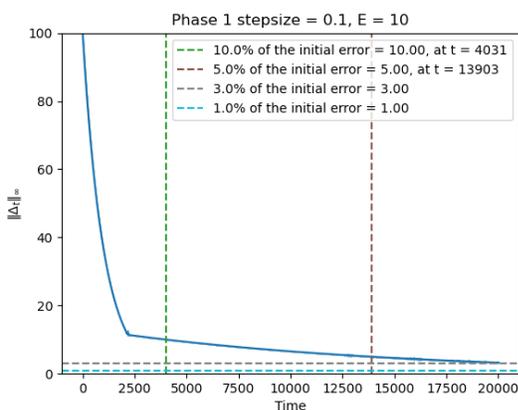
(b) With a phase 1 stepsize of 0.9, it meets the 10% tolerance level at iteration 16198.



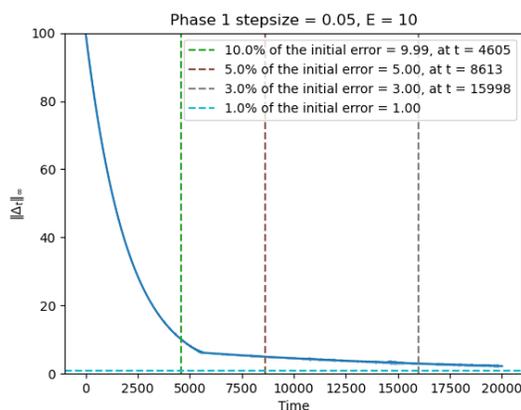
(c) With a phase 1 stepsize of 0.5, it meets the 10% tolerance level at iteration 14531.



(d) With a phase 1 stepsize of 0.2, it meets the 10% and 5% tolerance level at iterations 9651 and 19517, respectively.



(e) With a phase 1 stepsize of 0.1, it meets the 10% and 5% tolerance level at iterations 4031 and 13903, respectively.



(f) With a phase 1 stepsize of 0.05, it meets the 10%, 5%, and 3% tolerance levels at iterations 4605, 8613, and 15998, respectively.

Figure 8: Convergence performance of different tolerance levels of different stepsize choices. The horizontal dashed lines represent the tolerance levels not met, while the vertical dashed lines indicate the iterations at which the training processes meet the corresponding tolerance levels.