

Momentum-Based Federated Reinforcement Learning with Interaction and Communication Efficiency

Sheng Yue¹, Xingyuan Hua², Lili Chen¹, Ju Ren^{1,3*}

¹Department of Computer Science and Technology, BNRist, Tsinghua University, Beijing, China

²School of Computer Science and Technology, Beijing Institute of Technology, Beijing, China

³Zhongguancun Laboratory, Beijing, China

{shengyue, lilichen, renju}@tsinghua.edu.cn, xingyuanhua@bit.edu.cn

Abstract—Federated Reinforcement Learning (FRL) has garnered increasing attention recently. However, due to the intrinsic spatio-temporal non-stationarity of data distributions, the current approaches typically suffer from high interaction and communication costs. In this paper, we introduce a new FRL algorithm, named MFPO, that utilizes momentum, importance sampling, and additional server-side adjustment to control the shift of stochastic policy gradients and enhance the efficiency of data utilization. We prove that by proper selection of momentum parameters and interaction frequency, MFPO can achieve $\tilde{O}(HN^{-1}\epsilon^{-3/2})$ and $\tilde{O}(\epsilon^{-1})$ interaction and communication complexities (N represents the number of agents), where the interaction complexity achieves linear speedup with the number of agents, and the communication complexity aligns the best achievable of existing first-order FL algorithms. Extensive experiments corroborate the substantial performance gains of MFPO over existing methods on a suite of complex and high-dimensional benchmarks.

I. INTRODUCTION

With the rapid proliferation of Artificial Internet of Things (AIoT) applications and the increasing significance of data security, Federated Learning (FL) has emerged as a key enabler in the era of edge intelligence [1]–[3]. Recently, to reconcile FL with ever-growing intelligent decision-making applications, there has been a surge of interest towards *Federated Reinforcement Learning* (FRL), whereby distributed agents collaborate to build a decision policy with no need to share their raw trajectories [4]–[8]. FRL has been deemed as a practically appealing approach to address the data hungry of Reinforcement Learning (RL) [8], and demonstrated remarkable potential in a wide range of real-world systems, including robotics [4], autonomous driving [9], resource management in networking [10], and control of IoT devices [11].

However, the majority of current studies in FRL heuristically repurpose well-established supervised FL methods for the RL setting, e.g., directly combining FedAvg with classical PG or Q-learning [5], [8], [11], neglecting a unique challenge embedded therein: *the spatio-temporal non-stationarity of data distributions*. That is, in contrast to supervised FL operating on fixed datasets, FRL’s intrinsic trial-and-error learning process typically necessitates each agent to explore the environment and sample new data using the current policy in each local update, causing continually varying data distributions *across participating agents and training rounds*. As a result, it

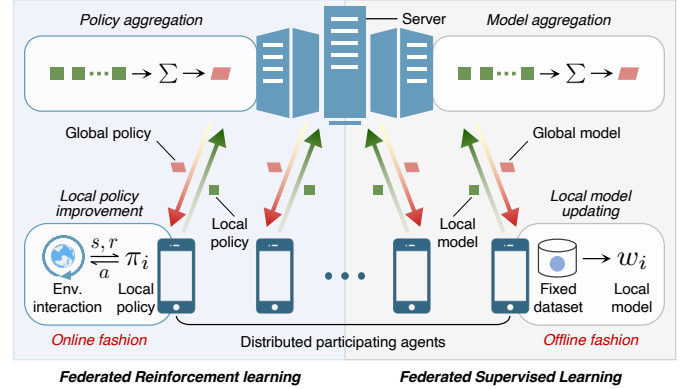


Fig. 1. Federated Reinforcement Learning vs Federated Supervised Learning.

would inflict two major issues on the existing methods: on one hand, since interaction with real systems can be slow, expensive, or fragile, the simplistic combinations easily suffer from *excessive interaction/sampling cost* during the continual environmental exploration; on the other hand, the dynamic data distributions can lead to substantially increased inter- and intra-agent shifts of stochastic gradients, inducing *brittle convergence properties and high communication complexity*. For instance, drawing upon the analysis for FedAvg [12], we can show that the direct combination of FedAvg and PG [13] requires $\tilde{O}(H\sigma_g^4\epsilon^{-2})$ environmental interactions and $\tilde{O}(\sigma_g^4\epsilon^{-2})$ communication rounds to reach an ϵ -stationary solution (H and σ_g^2 represent the trajectory length and the potentially large variance of policy gradient, respectively), which can be quite expensive for resource-sensitive edge users. Of note, despite the existence of variance reduction techniques in supervised FL [14]–[16], the specific spatio-temporal non-stationarity of data distributions renders them inapplicable to FRL settings [7].

Contributions. To overcome these challenges, this paper proposes *Momentum-assisted Federated Policy Optimization* (MFPO), capable of jointly optimizing both the interaction and communication complexities. Specifically, we introduce a new FRL framework that utilizes momentum, importance sampling, and extra server-side adjustment to control the variates of stochastic policy gradients and improve the efficiency of data utilization. Building on this, we rigorously quantify the impacts of the inter- and intra-agent gradient errors on the performance. We prove that by proper selection of momentum parameters

*Corresponding author

and interaction frequency, MFPO can effectively counteract the gradient shifts and achieve $\tilde{\mathcal{O}}(H\sigma_g^2N^{-1}\epsilon^{-3/2})$ interaction complexity along with $\tilde{\mathcal{O}}(\sigma^2\epsilon^{-1})$ communication complexity (N represents the number of agents). Notably, the interaction complexity achieves linear speedup with the number of agents, and the communication complexity recovers the best achievable across existing first-order stochastic FL algorithms. Finally, we evaluate MFPO on a suite of complex and high-dimensional RL benchmarks, including Classic Control, MuJoCo, and image-based Atari games. The results demonstrate that MFPO surpasses the state-of-the-art baseline methods by a significant margin, in terms of the performance and the efficiency of communication and interaction.

II. RELATED WORK

In recent years, several FRL algorithms have been proposed to address data-sharing constraints and facilitate safe co-training of policies [17]–[19]. Nadiger et al. [5] propose an FRL approach that combines DQN [20] and FedAvg [1] to obtain personalized policies for individual players in the Pong game by employing the smoothing average technique. Lim et al. [11] propose an FRL algorithm that combines Proximal Policy Optimization (PPO) [21] with FedAvg. Utilizing transfer learning, Liang et al. [9] adapt Deep Deterministic Policy Gradient (DDPG) [22] for FedAvg to operate in autonomous driving scenarios. Cha et al. [23] introduce a privacy-preserving variant of policy distillation, where a pre-arranged set of states and time-averaged policies is exchanged instead of raw data during training. Zhuo et al. [6] propose a two-agent FRL framework in the discrete state-action spaces built on Q-learning [24], where agents share local encrypted Q-values and alternately update the global Q-network using multilayer perceptron (MLP). Anwar and Raychowdhury [13] study the adversarial attack issue in FRL via combining the policy gradient method with FedAvg, with the goal of training a unified policy for individual tasks. However, these works heuristically repurpose popular supervised FL methods for RL settings, not equipped with rigorous communication/interaction assurances. It remains a critical drawback due to the (potentially excessive) communication/interaction cost in real systems [25].

More recently, analogous to [23], Khodadadian et al. [8] introduce an FRL algorithm by combining FedAvg with classical Q-learning and provide corresponding convergence guarantees, whereas they mainly concentrate on discrete cases (the state and action spaces are finite). Instead, this work operates in high-dimensional and continuous spaces. In a separate development, Fan et al. [7] develop a fault-tolerant FRL algorithm, namely FedPG-BR, where a certain percentage (denoted by $\alpha \leq 0.5$) of the participating agents are subject to random system failures or adversarial attacks. FedPG-BR requires $\mathcal{O}(HN^{-2/3}\epsilon^{-5/3} + H\alpha^{4/3}\epsilon^{-5/3})$ interaction steps for each agent, under the assumption that the server can continually interact with the environment. In contrast, our algorithm offers a more favorable interaction complexity of $\tilde{\mathcal{O}}(HN^{-1}\epsilon^{-3/2})$ with the linear speedup and does not require any interaction between the server and the environment.

III. FEDERATED REINFORCEMENT LEARNING

Reinforcement Learning (RL) is typically modeled as a *Markov Decision Process* (MDP) $\mathcal{M} \doteq \langle \mathcal{S}, \mathcal{A}, T, H, r, \mu, \gamma \rangle$, consisting of state space \mathcal{S} , action space \mathcal{A} , transition dynamics $T : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{P}(\mathcal{S})$, episode horizon H , reward function $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, R_{\max}]$, initial state distribution $\mu : \mathcal{S} \rightarrow [0, 1]$, and discount factor $\gamma \in (0, 1]$ [26]. A (stationary stochastic) policy, $\pi_\theta(a|s)$, defines the probability of taking action a at state s , which is parameterized by $\theta \in \mathbb{R}^d$. A trajectory, denoted as $\tau = \{s_1, a_1, \dots, s_H, a_H\}$, is a sequence of state-action pairs when rolling out π with M . The objective of RL is to find a policy that can maximize the expected cumulative reward over the traversed trajectories:

$$\max_{\theta} \mathbb{E} \left[\sum_{h=1}^H \gamma^{h-1} r(s_h, a_h) \mid T, \mu, \pi_\theta \right], \quad (1)$$

where the expectation is taken w.r.t. $s_1 \sim \mu$, $a_h \sim \pi_\theta(\cdot|s_h)$, and $s_{h+1} \sim T(\cdot|s_h, a_h)$. Due to the intrinsic complexity (e.g., high-dimensional state-action spaces and delayed reward feedback), RL algorithms are often time-consuming and sample-inefficient.

Federated Reinforcement Learning (FRL) solves Problem (1) in a distributed fashion, where N distributed agents federatively build a policy under the orchestration of a central server, without sharing their raw trajectories (as illustrated in Fig. 1). The goal is to speed up policy search and improve sampling efficiency via collaboration among agents while complying with the requirement of information privacy or data confidentiality.

IV. MFPO: MOMENTUM-BASED FEDERATED POLICY OPTIMIZATION

Denote the probability of $\tau = \{s_1, a_1, \dots, s_H, a_H\}$ under policy π_θ as $p(\tau|\theta) \doteq \mu(s_1) \prod_{h=1}^H T(s_{h+1}|s_h, a_h) \pi_\theta(a_h|s_h)$, and the objective function as $J(\theta) \doteq -\mathbb{E}_{\tau \sim p(\cdot|\theta)} [r(\tau)]$ with $r(\tau) = \sum_{h=1}^H \gamma^{h-1} r(s_h, a_h)$ being the cumulative reward of trajectory τ . Using the log-gradient trick and substituting the expression of $p(\tau|\theta)$, we can obtain the gradient of $J(\theta)$:

$$\nabla J(\theta) = -\mathbb{E}_{\tau \sim p(\cdot|\theta)} \left[\left(\sum_{h=1}^H \nabla_{\theta} \pi_{\theta}(a_h|s_h) \right) r(\tau) \right]. \quad (2)$$

Define $g(\theta; \tau)$ as an unbiased gradient estimator, which can be selected as widely used REINFORCE [27] or GPOMDP [28]. For example, the REINFORCE estimator can be expressed as $g(\theta; \tau) = (b - \sum_{h=1}^H r(s_h, a_h)) \sum_{h=1}^H \nabla_{\theta} \log \pi_{\theta}(a_h|s_h)$ with b being the baseline reward.

A natural FRL solution is to integrate Policy Gradient (PG) directly into current supervised FL frameworks. Yet, due to the stochasticity exponential in H , the gradient estimates inevitably suffer from substantially high variance [29]. Besides, since the local policy updates will alter the agent-side distribution, $p(\cdot|\theta)$, on which $\nabla J(\theta)$ depends, the trajectory distributions across agents are spatio-temporally non-stationary. As a result, this sort of combination may cause pronounced inter- and intra-agent gradient shifts, significantly impeding learning performance.

Motivated by the recent advance of momentum-based distributed optimization [16], [30], we next introduce a novel

momentum-assisted FRL algorithm, exploiting the techniques of momentum, importance sampling, and server-side adjustment, to tackle the above-mentioned problem. To be specific, the algorithm begins by initializing the policy parameters as $\theta_i^{(1)} = \bar{\theta}^{(1)}$ and then computes the corresponding initial directions as $\tilde{u}_i^{(1)} = (1/\tilde{D}) \sum_{j=1}^{\tilde{D}} g(\theta_i^{(1)}; \tau_{i,j}^{(1)})$, with \tilde{D} the number of trajectories generated from the initial policy. Subsequently, it alternates between local and global phases as follows.

1) Local phase: In step t , each agent samples D trajectories (denoted as $\tau_{i,j}^{(t)}$) via interacting with the environment using policy $\theta_i^{(t)}$. Then, it locally computes its updating direction by

$$\begin{aligned} \tilde{u}_i^{(t)} = & \nu^{(t)} \left(\tilde{u}_i^{(t-1)} - \frac{1}{D} \sum_{j=1}^D w_{i,j}^{(t)} \cdot g(\theta_i^{(t-1)}; \tau_{i,j}^{(t)}) \right) \\ & + \frac{1}{D} \sum_{j=1}^D g(\theta_i^{(t)}; \tau_{i,j}^{(t)}), \end{aligned} \quad (3)$$

with $\nu^{(t)} \in [0, 1]$ being the momentum parameter and $w_{i,j}^{(t)}$ the importance weight, computed as

$$w_{i,j}^{(t)} = w(\theta_i^{(t)}, \theta_i^{(t-1)}; \tau_{i,j}^{(t)}) = \frac{\prod_{h=1}^H \pi_i^{(t-1)}(a_{h,i} | s_{h,i})}{\prod_{h=1}^H \pi_i^{(t)}(a_{h,i} | s_{h,i})}. \quad (4)$$

If $t \bmod K \neq 0$, agents update their policy parameters locally:

$$\theta_i^{(t+1)} = \theta_i^{(t)} - \alpha^{(t)} \tilde{u}_i^{(t)}, \quad \forall i \in [N], \quad (5)$$

with K the number of local steps and $\alpha^{(t)}$ the stepsize.

2) Global phase: If $t \bmod K = 0$, agents upload their local parameters and directions to the server for aggregation:

$$\bar{u}^{(t)} = \frac{1}{N} \sum_{i=1}^N \tilde{u}_i^{(t)}, \quad \bar{\theta}^{(t)} = \frac{1}{N} \sum_{i=1}^N \theta_i^{(t)}. \quad (6)$$

The server carries out server-side adjustment as follows:

$$\bar{\theta}^{(t+1)} = \bar{\theta}^{(t)} - \alpha^{(t)} \bar{u}^{(t)}. \quad (7)$$

Then, it distributes the parameter and direction to all agents:

$$\theta_i^{(t+1)} = \bar{\theta}^{(t+1)}, \quad \tilde{u}_i^{(t)} = \bar{u}^{(t)}, \quad (8)$$

and starts the next round.

We term our algorithm *Momentum-assisted Federated Policy Optimization* (MFPO) and outline the pseudocode in Alg. 1. Intuitively, the momentum term along with importance sampling can keep track of past gradients in an off-policy manner, capable of improving sample efficiency while reducing the effect of fluctuations in the intra-agent gradient estimates [30]. On the other hand, the server-side adjustment enables the global policy to continue moving along the dominant dimension and hence alleviating the inter-agent gradient shift. Next, we show how to select momentum parameters and interaction frequency to jointly optimize interaction and communication complexities.

V. THEORETICAL ANALYSIS

In this section, we analyze the performance of MFPO. We first introduce the necessary assumptions and then present our main result, followed by detailed proofs.

Algorithm 1: MFPO

```

1 Initialize policy parameters and updating directions;
2 for  $t = 1$  to  $T$  do
3   for  $i = 1$  to  $N$  do
4     // Local phase
5     Agent  $i$  rolls out local policy with environment,
6     generates  $D$  trajectories, and computes local
7     direction by Eq. (3);
8     if  $t \bmod K \neq 0$  then
9       Agent  $i$  updates local policy by Eq. (5);
10  if  $t \bmod K = 0$  then
11    // Global phase
12    Agents upload local policies and directions;
13    Server updates global policy by Eqs. (6) and (7);
14    Server distributes global policy to all agents;

```

A. Notations and Assumptions

For $\tau \bmod K \neq 0$, we define auxiliary variables as follows:

$$\bar{u}^{(\tau)} \doteq \frac{1}{N} \sum_{i=1}^N \tilde{u}_i^{(\tau)}, \quad \bar{\theta}^{(\tau)} \doteq \frac{1}{N} \sum_{i=1}^N \theta_i^{(\tau)}, \quad (9)$$

and for each $i \in [N]$, we define $\tilde{u}_i^{(0)} \doteq 0$. When clear from the context, we use α_t and ν_t instead of $\alpha^{(t)}$ and $\nu^{(t)}$ for conciseness. We denote $t_q \doteq qK$ with $q \in \{0, \dots, M\}$ the index of communication rounds, and denote $[n] = \{1, 2, \dots, n\}$ for any $n \in \mathbb{N}$. In addition, we represent $\tilde{L} \doteq \max\{\tilde{L}_g, L\}$ where L and \tilde{L}_g are defined in Lemmas 1 and 3 respectively.

Due to the non-concavity of Problem (1) [29], it is generally not feasible to measure the optimality by function values. Instead, the convergence of non-convex problems is typically characterized via finding an ϵ -first-order stationary point (ϵ -FOSP), defined as follows.

Definition 1. A solution $\theta \in \mathbb{R}^d$ is called an ϵ -first-order stationary point (ϵ -FOSP) of Problem (1), if $\|\nabla J(\theta)\|^2 \leq \epsilon$.

We impose two commonly used assumptions in analyzing policy gradient methods as follows [7], [31], [32].

Assumption 1. There exist $\beta_1, \beta_2 > 0$ such that for all $s \in \mathcal{S}$, $a \in \mathcal{A}$, the log-density of the policy function satisfies

$$\|\nabla_{\theta} \log \pi_{\theta}(a|s)\| \leq \beta_1, \quad \|\nabla_{\theta}^2 \log \pi_{\theta}(a|s)\| \leq \beta_2. \quad (10)$$

Assumption 2. There exists $\sigma_g, \sigma_w > 0$ such that the following fact holds:

$$\mathbb{E}_{\tau \sim p(\cdot|\theta)} [\|g(\theta; \tau) - \nabla J(\theta)\|^2] \leq \sigma_g^2, \quad \forall \theta \in \mathbb{R}^d, \quad (11)$$

$$\text{Var}(w(\theta, \theta'; \tau)) \leq \sigma_w^2, \quad \forall \theta, \theta' \in \mathbb{R}^d, \tau \sim p(\cdot|\theta), \quad (12)$$

where $w(\theta, \theta'; \tau) = p(\tau|\theta')/p(\tau|\theta)$ is the importance weight used in the algorithm.

Assumptions 1 and 2 bound the gradient of the policy log-density and the variance of the gradient estimator, respectively. We suppose Assumptions 1 and 2 and $T = MK$ ($M \in \mathbb{N}$) hold throughout this section.

B. Main Results

We define the *communication complexity* as the total number of communication rounds necessary for the algorithm to reach an ϵ -FOSP, and the *interaction complexity* as the total number of actions that each agent requires taking in the environment to achieve the ϵ -FOSP. We define c_α , c_ν , and c_t as follows:

$$c_t \doteq \max \left\{ \frac{c_\nu^3 c_\alpha^3}{2^{12} K^3 \tilde{L}^3}, 2^{12} K^3 D^2 N^2 \sigma_g^2 - \sigma_g^2 t, 2\sigma_g^2 \right\}$$

$$c_\nu \doteq \frac{\tilde{L}^2}{24K(DN)^2} + \frac{64\tilde{L}^2}{DN}, \quad c_\alpha \doteq \frac{(DN\sigma_g)^{2/3}}{\tilde{L}}, \quad (13)$$

where positive integers N , K , and D represent the numbers of agents, local updates, and trajectories required in each local update, respectively. Our main result is presented below.

Theorem 1. *Suppose the stepsizes and momentum parameters are selected as $\alpha_t = c_\alpha / (c_t + \sigma_g^2 t)^{1/3}$ and $\nu_{t+1} = 1 - c_\nu \alpha_t^2$. For any $\lambda \in [0, 1]$, if $D = \mathcal{O}((T/N^2)^{1/2-\lambda/2})$, $\tilde{D} = DK$, and $K = \mathcal{O}((T/N^2)^{\lambda/3})$, MFPO finds an ϵ -FOSP after at most $\tilde{\mathcal{O}}(\epsilon^{-1})$ communication rounds and $\tilde{\mathcal{O}}(HN^{-1}\epsilon^{-3/2})$ environmental interactions.*

Proof. The result can be obtained by substituting the expressions of K , D and \tilde{D} in Eq. (46). We omit it for brevity. \square

Remarks. Theorem 1 indicates that MFPO achieves $\tilde{\mathcal{O}}(\epsilon^{-1})$ and $\tilde{\mathcal{O}}(HN^{-1}\epsilon^{-3/2})$ communication and interaction complexities by appropriate selection of the momentum parameters and the interaction frequency. The communication complexity recovers the best achievable of existing first-order FL algorithms [16], [33]. The interaction complexity exhibits linear speedup with the number of agents, making it superior to current FRL methods [7]. It implies that MFPO can effectively cope with the spatio-temporal non-stationary data distributions. In addition, Theorem 1 reveals a tradeoff between the local updates, K , and the required trajectories per step, D , characterized by $\lambda \in [0, 1]$. This means with a large number of local updates, the required trajectories per step can be set relatively small, and vice versa. In practical terms, this flexibility in adjusting K and D allows MFPO to adapt to different scenarios and requirements.

C. Detailed Proofs

In this subsection, we detail the proof for Theorem 1. We begin by bounding the *successive difference* of the objective function in the following lemma.

Lemma 1. *For $t \in (t_q, t_{q+1}]$, the following fact holds true:*

$$\begin{aligned} \mathbb{E}[J(\bar{\theta}^{(t+1)})] &\leq \mathbb{E}[J(\bar{\theta}^{(t)})] - \frac{\alpha_t}{2} \mathbb{E}[\|\nabla J(\bar{\theta}^{(t)})\|^2] \\ &\quad - \frac{\alpha_t - \alpha_t^2 L}{2} \mathbb{E}[\|\bar{u}^{(t)}\|^2] + \alpha_t \mathbb{E}[\|\bar{\varepsilon}^{(t)}\|^2] \\ &\quad + \frac{\alpha_t L^2}{N} \sum_{i=1}^N \mathbb{E}[\|\theta_i^{(t)} - \bar{\theta}^{(t)}\|^2], \end{aligned} \quad (14)$$

with $\bar{\varepsilon}^{(t)} \doteq \bar{u}^{(t)} - (1/N) \sum_{i=1}^N \nabla J(\theta_i^{(t)})$ being the gradient error and $L \doteq HR_{\max}(H\beta_1^2 + \beta_2)/(1 - \gamma)$.

Proof. Built upon Assumption 1 and [32, Proposition 5.2], $J(\theta)$ is L -smooth, which implies

$$J(\theta_1) \leq J(\theta_2) + \nabla J(\theta_2)^\top (\theta_1 - \theta_2) + \frac{L}{2} \|\theta_1 - \theta_2\|^2. \quad (15)$$

Based on the L -smooth and Eqs. (5) and (9), we have

$$\begin{aligned} J(\bar{\theta}^{(t+1)}) &= J(\bar{\theta}^{(t)}) - \alpha_t \nabla J(\bar{\theta}^{(t)})^\top \bar{u}^{(t)} + \frac{\alpha_t^2 L}{2} \|\bar{u}^{(t)}\|^2 \\ &= J(\bar{\theta}^{(t)}) - \frac{\alpha_t}{2} \mathbb{E}[\|\nabla J(\bar{\theta}^{(t)})\|^2] - \frac{\alpha_t - \alpha_t^2 L}{2} \\ &\quad \cdot \underbrace{\mathbb{E}[\|\bar{u}^{(t)}\|^2]}_{(a)} + \frac{\alpha_t}{2} \underbrace{\|\bar{u}^{(t)} - \nabla J(\bar{\theta}^{(t)})\|^2}_{(a)}, \end{aligned} \quad (16)$$

where second equality is derived by adding and subtracting $\alpha_t \|\bar{u}^{(t)}\|^2$ and utilizing $2\theta_1^\top \theta_2 = \|\theta_1\|^2 + \|\theta_2\|^2 - \|\theta_1 - \theta_2\|^2$. Regarding (a), we have

$$\begin{aligned} (a) &= \left\| \bar{u}^{(t)} - \sum_{i=1}^N \frac{\nabla J(\theta_i^{(t)})}{N} + \sum_{i=1}^N \frac{\nabla J(\theta_i^{(t)})}{N} - \nabla J(\bar{\theta}^{(t)}) \right\|^2 \\ &\leq 2\|\bar{\varepsilon}^{(t)}\|^2 + \frac{2}{N} \sum_{i=1}^N \|\nabla J(\theta_i^{(t)}) - \nabla J(\bar{\theta}^{(t)})\|^2 \\ &\leq 2\|\bar{\varepsilon}^{(t)}\|^2 + \frac{2L^2}{N} \sum_{i=1}^N \|\theta_i^{(t)} - \bar{\theta}^{(t)}\|^2, \quad (\text{the smoothness}) \end{aligned} \quad (17)$$

where derived Eq. (17) using the following relationship:

$$\|\theta_1 + \theta_2 + \dots + \theta_n\|^2 \leq n\|\theta_1\|^2 + \dots + n\|\theta_n\|^2. \quad (18)$$

Plugging (a) in Eq. (16) and taking expectations on both sides yield the result. \square

Next, we bound the last term of Eq. (14) in Lemma 2.

Lemma 2. *For $t \in (t_q, t_{q+1}]$ and $i \in [N]$, we have*

$$\mathbb{E}[\|\theta_i^{(t)} - \bar{\theta}^{(t)}\|^2] \leq (K-1) \sum_{\tau=t_q+1}^t \alpha_\tau^2 \cdot \mathbb{E}[\|\tilde{u}_i^{(\tau)} - \bar{u}^{(\tau)}\|^2]. \quad (19)$$

Proof. When $t = t_q + 1$, $\mathbb{E}[\|\theta_i^{(t_q+1)} - \bar{\theta}^{(t_q+1)}\|^2] = 0$ holds due to Eq. (8). When $t \in (t_q + 1, t_{q+1}]$, it follows

$$\begin{aligned} \mathbb{E}[\|\theta_i^{(t)} - \bar{\theta}^{(t)}\|^2] &= \mathbb{E} \left[\left\| \sum_{\tau=t_q+1}^{t-1} \alpha_\tau (\tilde{u}_i^{(\tau)} - \bar{u}^{(\tau)}) \right\|^2 \right] \\ &\quad (\text{from Eq. (5) and } \theta_i^{(t_q+1)} = \bar{\theta}^{(t_q+1)}) \\ &\leq (K-1) \sum_{\tau=t_q+1}^{t-1} \alpha_\tau^2 \cdot \mathbb{E}[\|\tilde{u}_i^{(\tau)} - \bar{u}^{(\tau)}\|^2], \end{aligned} \quad (20)$$

where the last inequality is derived via $t_{q+1} - t_q = K$ and Eq. (18). Thus, we complete the proof. \square

Recall the definition of $\bar{\varepsilon}^{(t)}$ in Lemma 1. Lemma 2 characterizes the error accumulation in the iterates of Alg. 1. Substituting Eq. (19) in Eq. (14), we obtain

$$\mathbb{E}[J(\bar{\theta}^{(t+1)})] \leq \mathbb{E}[J(\bar{\theta}^{(t)})] - \frac{\alpha_t}{2} \mathbb{E}[\|\nabla J(\bar{\theta}^{(t)})\|^2]$$

$$\begin{aligned}
& - \frac{\alpha_t - \alpha_t^2 L}{2} \mathbb{E}[\|\bar{u}^{(t)}\|^2] + \underbrace{\alpha_t \mathbb{E}[\|\bar{\varepsilon}^{(t)}\|^2]}_{\text{Gradient error}} \\
& + \frac{(K-1)L^2 \alpha_t}{N} \sum_{\tau=t_q+1}^{t-1} \alpha_\tau^2 \\
& \cdot \underbrace{\sum_{i=1}^N \mathbb{E}[\|\bar{u}_i^{(\tau)} - \bar{u}^{(\tau)}\|^2]}_{\text{Gradient consensus error}}. \tag{21}
\end{aligned}$$

It suggests that the expected descent of $J(\cdot)$ relies on both the expected *gradient error* and the expected *gradient consensus error*. For conciseness, we denote the gradient consensus error as $\delta_t \doteq \sum_{i=1}^N \mathbb{E}[\|\bar{u}_i^{(t)} - \bar{u}^{(t)}\|^2]$. In what follows, we bound the two errors respectively.

Regarding the gradient error, we introduce Lemma 3 to show how it contracts over time.

Lemma 3. *Denote constants $L_g \doteq H\beta_2(R_{\max} + b)/(1 - \gamma)$, $G_g \doteq H\beta_1(R_{\max} + b)/(1 - \gamma)$, $c_w \doteq H(2H\beta_1^2 + \beta_2)(\sigma_w^2 + 1)$ and $\tilde{L}_g \doteq \sqrt{2(L_g^2 + G_g^2 c_w)}$ where b is the baseline reward. Then, for $t \in [T]$, the following fact holds:*

$$\begin{aligned}
\mathbb{E}[\|\bar{\varepsilon}^{(t)}\|^2] & \leq \nu_t^2 \cdot \mathbb{E}[\|\bar{\varepsilon}^{(t-1)}\|^2] + \frac{4\tilde{L}_g^2 \nu_t^2 \alpha_{t-1}^2}{DN} \mathbb{E}[\|\bar{u}^{(t-1)}\|^2] \\
& + \frac{2\sigma_g^2(1 - \nu_t)^2}{DN} + \frac{8(K-1)\tilde{L}_g^2 \nu_t^2 \alpha_{t-1}^2 \delta_{t-1}}{KDN^2} \tag{22}
\end{aligned}$$

with $\bar{\varepsilon}^{(t)}$ being the gradient error defined in Lemma 1.

Proof. From the definition of $\bar{\varepsilon}^{(t)}$ and Eq. (9), we have

$$\begin{aligned}
\mathbb{E}[\|\bar{\varepsilon}^{(t)}\|^2] & = \frac{1}{D^2 N^2} \sum_{i=1}^N \sum_{j=1}^D \mathbb{E} \left[\left\| g(\theta_i^{(t)}; \tau_{i,j}^{(t)}) - \nabla J(\theta_i^{(t)}) \right. \right. \\
& \quad \left. \left. - \nu_t \left(w_{i,j}^{(t)} g(\theta_i^{(t-1)}; \tau_{i,j}^{(t)}) - \nabla J(\theta_i^{(t-1)}) \right) \right\|^2 \right] \\
& + \nu_t^2 \mathbb{E}[\|\bar{\varepsilon}^{(t-1)}\|^2], \tag{23}
\end{aligned}$$

where we add and subtract $(1/N) \sum_{i=1}^N \nu_t \nabla J(\theta_i^{(t-1)})$, expand the norms, and use the fact that the corresponding cross terms are zero, which can be easily verified via the tower rule and the unbiasedness of the importance-weighted gradient estimator $w_{i,j}^{(t)} g(\theta_i^{(t-1)}; \tau_{i,j}^{(t)})$. That is, for any $\theta, \theta' \in \mathbb{R}^d$, the following holds:

$$\begin{aligned}
\nabla J(\theta) & = \mathbb{E}_{\tau \sim p(\cdot|\theta)} [g(\theta; \tau)] \quad (\text{the unbiasedness of } g(\theta; \tau)) \\
& = \int p(\tau|\theta') \cdot \frac{p(\tau|\theta)}{p(\tau|\theta')} \cdot g(\theta; \tau) d\tau \\
& = \mathbb{E}_{\tau \sim p(\cdot|\theta')} [w(\theta', \theta; \tau) g(\theta; \tau)]. \tag{24}
\end{aligned}$$

For the first term in Eq. (23), we have

$$\begin{aligned}
& \mathbb{E} \left[\left\| g(\theta_i^{(t)}; \tau_{i,j}^{(t)}) - \nabla J(\theta_i^{(t)}) \right. \right. \\
& \quad \left. \left. - \nu_t \left(w_{i,j}^{(t)} g(\theta_i^{(t-1)}; \tau_{i,j}^{(t)}) - \nabla J(\theta_i^{(t-1)}) \right) \right\|^2 \right] \\
& = 2\nu_t^2 \mathbb{E} \left[\left\| g(\theta_i^{(t)}; \tau_{i,j}^{(t)}) - w_{i,j}^{(t)} g(\theta_i^{(t-1)}; \tau_{i,j}^{(t)}) \right. \right. \\
& \quad \left. \left. - \left(\nabla J(\theta_i^{(t)}) - \nabla J(\theta_i^{(t-1)}) \right) \right\|^2 \right] + 2(1 - \nu_t)^2
\end{aligned}$$

$$\begin{aligned}
& \cdot \mathbb{E} \left[\left\| g(\theta_i^{(t)}; \tau_{i,j}^{(t)}) - \nabla J(\theta_i^{(t)}) \right\|^2 \right] \quad (\text{from Eq. (18)}) \\
& \leq 2\nu_t^2 \mathbb{E} \left[\left\| g(\theta_i^{(t)}; \tau_{i,j}^{(t)}) - w_{i,j}^{(t)} g(\theta_i^{(t-1)}; \tau_{i,j}^{(t)}) \right\|^2 \right] + 2\sigma_g^2 \\
& \quad \cdot (1 - \nu_t)^2 \quad (\text{Assumption 2, mean variance inequality}) \\
& \leq 4\nu_t^2 \mathbb{E} \left[\left\| g(\theta_i^{(t)}; \tau_{i,j}^{(t)}) - g(\theta_i^{(t-1)}; \tau_{i,j}^{(t)}) \right\|^2 \right] \\
& \quad + 4\nu_t^2 \mathbb{E} \left[\left\| (1 - w_{i,j}^{(t)}) g(\theta_i^{(t-1)}; \tau_{i,j}^{(t)}) \right\|^2 \right] + 2(1 - \nu_t)^2 \sigma_g^2 \\
& (\text{adding and subtracting } g(\theta_i^{(t-1)}; \tau_{i,j}^{(t)}) \text{ and using Eq. (18)}) \\
& \stackrel{(a)}{\leq} 4\nu_t^2 L_g^2 \mathbb{E}[\|\theta_i^{(t)} - \theta_i^{(t-1)}\|^2] + 4\nu_t^2 G_g^2 \mathbb{E}[(1 - w_{i,j}^{(t)})^2] \\
& \quad + 2(1 - \nu_t)^2 \sigma_g^2 \quad (\text{from Eq. (27)}) \\
& \stackrel{(b)}{\leq} 4\nu_t^2 L_g^2 \mathbb{E}[\|\theta_i^{(t)} - \theta_i^{(t-1)}\|^2] + 4\nu_t^2 G_g^2 c_w \mathbb{E}[\|\theta_i^{(t)} \\
& \quad - \theta_i^{(t-1)}\|^2] + 2(1 - \nu_t)^2 \sigma_g^2 \quad (\text{from Eq. (28)}) \\
& = 2\nu_t^2 \tilde{L}_g^2 \mathbb{E}[\|\theta_i^{(t)} - \theta_i^{(t-1)}\|^2] + 2(1 - \nu_t)^2 \sigma_g^2 \tag{25} \\
& \leq 2\alpha_{t-1}^2 \nu_t^2 \tilde{L}_g^2 \mathbb{E}[\|\bar{u}_i^{(t-1)}\|^2] + 2(1 - \nu_t)^2 \sigma_g^2 \quad (\text{from Eq. (5)}) \\
& \leq 2(1 - \nu_t)^2 \sigma_g^2 + 4\alpha_{t-1}^2 \nu_t^2 \tilde{L}_g^2 \mathbb{E}[\|\bar{u}_i^{(t-1)} - \bar{u}^{(t-1)}\|^2] \\
& \quad + 4\alpha_{t-1}^2 \nu_t^2 \tilde{L}_g^2 \mathbb{E}[\|\bar{u}^{(t-1)}\|^2] \quad (\text{from Eq. (18)}) \\
& \leq \frac{8(K-1)\tilde{L}_g^2 \nu_t^2 \alpha_{t-1}^2}{K} \cdot \mathbb{E}[\|\bar{u}_i^{(t-1)} - \bar{u}^{(t-1)}\|^2] \\
& \quad + 4\alpha_{t-1}^2 \nu_t^2 \tilde{L}_g^2 \mathbb{E}[\|\bar{u}^{(t-1)}\|^2] + 2(1 - \nu_t)^2 \sigma_g^2, \tag{26}
\end{aligned}$$

where the last inequality follows from the fact: (i) when $K = 1$, $\bar{u}_i^{(t-1)} = \bar{u}^{(t-1)}$, and when $K \geq 2$, $K - 1/K \geq 1/2$. Built on Assumptions 1 and 2, inequality (a) holds due to [32, Proposition 5.2]:

$$\|g(\theta_1; \tau) - g(\theta_2; \tau)\| \leq L_g \|\theta_1 - \theta_2\|, \quad \|g(\theta; \tau)\| \leq G_g, \tag{27}$$

and inequality (b) follows [34, Lemma 1] and [32, Lemma 6.1]: for $\tau \sim p(\cdot|\theta)$, we have

$$\mathbb{E}[w(\theta, \theta'; \tau)] = 1, \quad \text{Var}(w(\theta, \theta'; \tau)) \leq c_w \|\theta - \theta'\|^2. \tag{28}$$

Plugging Eq. (26) in Eq. (23) completes the proof. \square

Next, we bound the gradient consensus error in Lemma 4.

Lemma 4. *For $t \in (t_q, t_{q+1}]$, the following fact holds:*

$$\begin{aligned}
\delta_t & \leq \nu_t^2 \left(1 + \frac{1}{K} + 8\alpha_{t-1}^2 \tilde{L}_g^2 K \right) \delta_{t-1} + 32(1 - \nu_t)^2 L^2 K^2 \\
& \quad \cdot \sum_{\tau=t_q+1}^t \alpha_\tau^2 \delta_\tau + 8\alpha_{t-1}^2 \tilde{L}_g^2 \nu_t^2 N K \mathbb{E}[\|\bar{u}^{(t-1)}\|^2] \\
& \quad + \frac{8(1 - \nu_t)^2 \sigma_g^2 N K}{D}, \tag{29}
\end{aligned}$$

where L and \tilde{L}_g are defined in Lemmas 1 and 3, respectively.

Proof. We denote $\tilde{d}_i^{(t-1)} \doteq (1/D) \sum_{j=1}^D w_{i,j}^{(t-1)} \cdot g(\theta_i^{(t-1)}; \tau_{i,j}^{(t-1)})$ and $d_i^{(t)} \doteq (1/D) \sum_{j=1}^D g(\theta_i^{(t)}; \tau_{i,j}^{(t)})$. For any $y > 0$, we have

$$\begin{aligned}
\delta_t & \leq (1 + y) \nu_t^2 \delta_{t-1} + (1 + \frac{1}{y}) \mathbb{E} \left[\sum_{i=1}^N \left\| d_i^{(t)} - \frac{1}{N} \sum_{j=1}^N d_j^{(t)} \right\|^2 \right. \\
& \quad \left. - \nu_t \left(\tilde{d}_i^{(t-1)} - \frac{1}{N} \sum_{j=1}^N \tilde{d}_j^{(t-1)} \right) \right], \tag{30}
\end{aligned}$$

which is derived by substituting the expressions of $\tilde{u}_i^{(t)}$ and $\bar{u}^{(t)}$, extending the norm, and using $2\theta_1^\top\theta_2 \leq q\|\theta_1\|^2 + (1/q)\|\theta_2\|^2$. For the second term of Eq. (30), we can write

$$\begin{aligned} & \mathbb{E} \left[\sum_{i=1}^N \left\| d_i^{(t)} - \frac{1}{N} \sum_{j=1}^N d_j^{(t)} - \nu_t \left(\tilde{d}_i^{(t-1)} - \frac{1}{N} \sum_{j=1}^N \tilde{d}_j^{(t-1)} \right) \right\|^2 \right] \\ & \leq 2(1 - \nu_t)^2 \underbrace{\mathbb{E} \left[\sum_{i=1}^N \left\| d_i^{(t)} - \frac{1}{N} \sum_{j=1}^N d_j^{(t)} \right\|^2 \right]}_{(a)} \\ & + 2\nu_t^2 \underbrace{\sum_{i=1}^N \mathbb{E} \left[\left\| d_i^{(t)} - \tilde{d}_i^{(t-1)} \right\|^2 \right]}_{(b)}, \end{aligned} \quad (31)$$

where the inequality is derived from Eq. (18) and the fact: for any $\theta_1, \theta_2, \dots, \theta_n \in \mathbb{R}^d$ and $\bar{\theta} = (1/n) \sum_{i=1}^n \theta_i$, it follows

$$\sum_{i=1}^n \|\theta_i - \bar{\theta}\|^2 \leq \sum_{i=1}^n \|\theta_i\|^2. \quad (32)$$

For (b), analogous to Eq. (25), we have

$$(b) \leq \tilde{L}_g^2 \mathbb{E} [\|\theta_i^{(t)} - \theta_i^{(t-1)}\|^2], \quad (33)$$

For (a), we have

$$\begin{aligned} (a) & \leq \mathbb{E} \left[2 \sum_{i=1}^N \left\| d_i^{(t)} - \nabla J(\theta_i^{(t)}) - \left(\frac{1}{N} \sum_{j=1}^N d_j^{(t)} - \frac{1}{N} \sum_{j=1}^N \nabla J(\theta_j^{(t)}) \right) \right\|^2 + 2 \sum_{i=1}^N \left\| \nabla J(\theta_i^{(t)}) - \frac{1}{N} \sum_{j=1}^N \nabla J(\theta_j^{(t)}) \right\|^2 \right] \quad (\text{from Eq. (18)}) \\ & = \frac{2}{D^2} \sum_{i=1}^N \mathbb{E} \left[\underbrace{\left\| \sum_{j=1}^D g(\theta_i^{(t)}; \tau_{i,j}^{(t)}) - \nabla J(\theta_i^{(t)}) \right\|^2}_{(c)} \right] \\ & + 2 \sum_{i=1}^N \mathbb{E} \left[\underbrace{\left\| \nabla J(\theta_i^{(t)}) - \frac{1}{N} \sum_{j=1}^N \nabla J(\theta_j^{(t)}) \right\|^2}_{(d)} \right] \\ & \quad (\text{using Eq. (32) and rearranging terms}) \\ & = \frac{2N\sigma_g^2}{D} + 8L^2 \sum_{i=1}^N \mathbb{E} [\|\theta_i^{(t)} - \bar{\theta}^{(t)}\|^2], \end{aligned} \quad (34)$$

Term (c) can be bounded via expanding the norm, eliminating zero expected cross terms, and using Assumption 2 as follows:

$$(c) = \sum_{j=1}^D \mathbb{E} \left[\left\| g(\theta_i^{(t)}; \tau_{i,j}^{(t)}) - \nabla J(\theta_i^{(t)}) \right\|^2 \right] \leq D\sigma_g^2. \quad (35)$$

From Eqs. (15) and (18), term (d) is bounded by

$$\begin{aligned} (d) & = \left\| \nabla J(\theta_i^{(t)}) - \nabla J(\bar{\theta}^{(t)}) + \nabla J(\bar{\theta}^{(t)}) - \sum_{j=1}^N \frac{\nabla J(\theta_j^{(t)})}{N} \right\|^2 \\ & \leq 2L^2 \|\theta_i^{(t)} - \bar{\theta}^{(t)}\|^2 + \frac{2L^2}{N} \sum_{j=1}^N \|\theta_j^{(t)} - \bar{\theta}^{(t)}\|^2. \end{aligned} \quad (36)$$

Substituting (a), (b) in Eq. (31) and rearranging terms yield

$$\begin{aligned} & \mathbb{E} \left[\sum_{i=1}^N \left\| d_i^{(t)} - \frac{1}{N} \sum_{j=1}^N d_j^{(t)} - \nu_t \left(\tilde{d}_i^{(t-1)} - \frac{1}{N} \sum_{j=1}^N \tilde{d}_j^{(t-1)} \right) \right\|^2 \right] \\ & \leq \frac{4(1 - \nu_t)^2 \sigma_g^2 N}{D} + 2\tilde{L}_g^2 \nu_t^2 \sum_{i=1}^N \mathbb{E} [\|\theta_i^{(t)} - \theta_i^{(t-1)}\|^2] \\ & + 16(1 - \nu_t)^2 L^2 \sum_{i=1}^N \mathbb{E} [\|\theta_i^{(t)} - \bar{\theta}^{(t)}\|^2]. \end{aligned} \quad (37)$$

Plugging Eq. (37) into Eq. (30), for $t \in (t_q, t_{q+1}]$, we obtain

$$\begin{aligned} \delta_t & \leq \nu_t^2 \left(1 + y + 4\alpha_{t-1}^2 \tilde{L}_g^2 (1 + \frac{1}{y}) \right) \delta_{t-1} \\ & + 16L^2 (1 - \nu_t)^2 (K - 1) (1 + \frac{1}{y}) \sum_{\tau=t_q+1}^t \alpha_\tau^2 \delta_\tau \\ & + 4\tilde{L}_g^2 \nu_t^2 \alpha_{t-1}^2 N (1 + \frac{1}{y}) \cdot \mathbb{E} [\|\bar{u}^{(t-1)}\|^2] \\ & + (1 + \frac{1}{y}) \cdot \frac{4(1 - \nu_t)^2 \sigma_g^2 N}{D}, \end{aligned} \quad (38)$$

where the inequality is derived via using Lemma 2, adding and subtracting $\bar{u}^{(t-1)}$ in $\|\tilde{u}_i^{(t-1)}\|^2$, and then applying Eq. (18). Letting $y = 1/K$ and using $1 + K \leq 2K$ yield the result. \square

Lemmas 3 and 4 bound the expected gradient error and the expected gradient shift while quantifying the impacts of learning rates, momentum parameters, local steps and interaction frequency on the convergence. We proceed to show how to select these parameters correctly to optimize the communication and interaction complexity.

Lemma 5. For $t \in [T]$, if $\nu_{t+1}/\alpha_t + 64\tilde{L}^2\alpha_t/(DN) \leq 1/\alpha_{t-1}$ and $\nu_{t+1} = 1 - c_\nu\alpha_t^2$ hold with $\tilde{L} = \max\{\tilde{L}_g, L\}$, we have

$$\begin{aligned} \alpha_t \mathbb{E} [\|\bar{\varepsilon}^{(t)}\|^2] & \leq \frac{DN}{64\tilde{L}^2} \left(\frac{\mathbb{E} [\|\bar{\varepsilon}^{(t)}\|^2]}{\alpha_{t-1}} - \frac{\mathbb{E} [\|\bar{\varepsilon}^{(t+1)}\|^2]}{\alpha_t} \right) + \alpha_t \delta_t \\ & \cdot \frac{K-1}{8NK} + \frac{\alpha_t}{16} \mathbb{E} [\|\bar{u}^{(t)}\|^2] + \frac{c_\nu^2 \sigma_g^2 \alpha_t^3}{32\tilde{L}^2}. \end{aligned} \quad (39)$$

Proof. From Lemma 3 and $\nu_{t+1}^2 \leq \nu_{t+1} \leq 1$, for all $t \in [T]$, we can write

$$\begin{aligned} & \frac{\mathbb{E} [\|\bar{\varepsilon}^{(t+1)}\|^2]}{\alpha_t} - \frac{\mathbb{E} [\|\bar{\varepsilon}^{(t)}\|^2]}{\alpha_{t-1}} \\ & \leq \left(\frac{\nu_{t+1}}{\alpha_t} - \frac{1}{\alpha_{t-1}} \right) \mathbb{E} [\|\bar{\varepsilon}^{(t)}\|^2] + \frac{4\tilde{L}_g^2 \alpha_t}{DN} \mathbb{E} [\|\bar{u}^{(t)}\|^2] \\ & + \frac{8(K-1)\tilde{L}_g^2 \alpha_t \delta_t}{KDN^2} + \frac{2\sigma_g^2 (1 - \nu_{t+1})^2}{DN\alpha_t}. \end{aligned} \quad (40)$$

Using the contions and rearranging terms yields the result. \square

Lemma 6. For each $t \in (t_q, t_{q+1}]$, if $c_\nu \leq 128\sqrt{2}L^2/(DN)$, $\alpha_t \leq 1/(16\tilde{L}K)$, and $\nu_t = 1 - c_\nu\alpha_{t-1}^2$ hold, then we have

$$\frac{K-1}{4NK} \sum_{\tau=t_q+1}^t \alpha_\tau \delta_\tau \leq \sum_{\tau=t_q}^{t-1} \frac{\alpha_\tau}{64} \mathbb{E} [\|\bar{u}^{(\tau)}\|^2] + \frac{c_\nu^2 \sigma_g^2 \alpha_\tau^3}{64D\tilde{L}^2}. \quad (41)$$

Proof. Due to $\nu_t^2 \leq 1$, $\alpha_t \leq 1/(16K\tilde{L}_g)$ and Lemma 4, for each $t \in (t_q, t_{q+1}]$, we have

$$\begin{aligned} \delta_t &\leq \left(1 + \frac{33}{32K}\right) \delta_{t-1} + 32K^2L^2(1-\nu_t)^2 \sum_{\tau=t_q+1}^t \alpha_\tau^2 \delta_\tau \\ &\quad + \frac{N\tilde{L}_g\alpha_{t-1} \mathbb{E}[\|\bar{u}^{(t-1)}\|^2]}{2} + \frac{8NK\sigma_g^2(1-\nu_t)^2}{D}. \end{aligned} \quad (42)$$

Due to $\delta_{t_q} = 0$, applying Eq. (42) recursively for $\tau \in (t_q, t]$, we obtain

$$\begin{aligned} \delta_t &\leq 96K^2L^2c_\nu^2 \sum_{\tau=t_q}^{t-1} \alpha_\tau^4 \sum_{n=t_q+1}^{\tau+1} \alpha_n^2 \delta_n + \frac{24c_\nu^2\sigma_g^2NK}{D} \sum_{\tau=t_q}^{t-1} \alpha_\tau^4 \\ &\quad + \frac{3N\tilde{L}_g}{2} \sum_{\tau=t_q}^{t-1} \alpha_\tau \mathbb{E}[\|\bar{u}^{(\tau)}\|^2] \\ &\quad (\text{from } t-\tau-1 \leq K \text{ and } (1+33/(32K))^K \leq e^{33/22} \leq 3) \\ &\leq 96K^3L^2c_\nu^2 \left(\frac{1}{16LK}\right)^5 \sum_{\tau=t_q+1}^{t-1} \alpha_\tau \delta_\tau + \frac{3c_\nu^2\sigma_g^2N}{2B\tilde{L}} \sum_{\tau=t_q}^{t-1} \alpha_\tau^3 \\ &\quad + \frac{3N\tilde{L}_g}{2} \sum_{\tau=t_q}^{t-1} \alpha_\tau \mathbb{E}[\|\bar{u}^{(\tau)}\|^2], \end{aligned} \quad (43)$$

where the last inequality holds from $\alpha_t \leq 1/(16\tilde{L}_gK)$. Multiplying Eq. (43) by α_t and summing over $t = t_q + 1$ to $\tau \in (t_q, t_{q+1}]$, we get

$$\begin{aligned} \sum_{t=t_q+1}^{\tau} \alpha_t \delta_t &\leq \frac{3N}{32} \sum_{n=t_q}^{\tau-1} \alpha_n \mathbb{E}[\|\bar{u}^{(n)}\|^2] + \frac{3c_\nu^2\sigma_g^2N}{32D\tilde{L}^2} \sum_{n=t_q}^{\tau-1} \alpha_n^3 \\ &\quad + 96K^4L^2c_\nu^2 \left(\frac{1}{16LK}\right)^6 \sum_{n=t_q+1}^{\tau} \alpha_n \delta_n, \end{aligned} \quad (44)$$

where we use $\tau - t_q \leq K$, $\alpha_t \leq 1/(16\tilde{L}_gK)$ and $\tilde{u}_i^{(t_q)} = \bar{u}^{(t_q)}$. Rearranging terms yields

$$\begin{aligned} &\left(1 - 96K^4L^2c_\nu^2 \left(\frac{1}{16LK}\right)^6\right) \sum_{t=t_q+1}^{\tau} \alpha_t \delta_t \\ &\leq \frac{3N}{32} \sum_{n=t_q}^{\tau-1} \alpha_n \mathbb{E}[\|\bar{u}^{(n)}\|^2] + \frac{3Nc_\nu^2\sigma_g^2}{32D\tilde{L}^2} \sum_{n=t_q}^{\tau-1} \alpha_n^3. \end{aligned} \quad (45)$$

Due to $c_\nu \leq 128\sqrt{2}L^2/(DN)$, $1 - 96K^4L^2c_\nu^2/(16LK)^6 \geq 1/4$ holds, thereby completing the proof. \square

Now, we are ready to establish the convergence property.

Theorem 2. *Suppose that the sequences of learning rates and momentum parameters across interaction steps are selected as $\alpha_t = c_\alpha/(c_t + \sigma_g^2t)^{1/3}$ and $\nu_{t+1} = 1 - c_\nu\alpha_t^2$ respectively. Then, the following fact holds true:*

$$\begin{aligned} &\frac{1}{T} \sum_{t=1}^T \mathbb{E}[\|\nabla J(\bar{\theta}^{(t)})\|^2] \\ &\leq \left(\frac{32\tilde{L}K}{T} + \frac{2\tilde{L}}{(DN)^{2/3}T^{2/3}}\right) \left(J(\bar{\theta}^{(1)}) - J^*\right) \\ &\quad + \left(\frac{2^{15}K}{T} + \frac{2^{11}}{(DN)^{2/3}T^{2/3}}\right) \left(1 + \ln(1+T)\right) \sigma_g^2 \end{aligned}$$

$$+ \left(\frac{8DK^2}{\tilde{D}T} + \frac{DK}{2\tilde{D}(DN)^{2/3}T^{2/3}}\right) \sigma_g^2, \quad (46)$$

where $J^* \doteq \min_\theta J(\theta)$, and c_α , c_ν and c_t are defined as

$$\begin{aligned} c_t &\doteq \max \left\{ \frac{c_\nu^3 c_\alpha^3}{2^{12}K^3\tilde{L}^3}, 2^{12}K^3D^2N^2\sigma_g^2 - \sigma_g^2t, 2\sigma_g^2 \right\} \\ c_\nu &\doteq \frac{\tilde{L}^2}{24K(DN)^2} + \frac{64\tilde{L}^2}{DN}, \quad c_\alpha \doteq \frac{(DN\sigma_g)^{2/3}}{\tilde{L}}. \end{aligned}$$

Proof. First, it can be easily verified that $\alpha_t \leq 1/(16\tilde{L}K)$. Due to $c_t \leq c_{t-1}$, the following holds:

$$\begin{aligned} \frac{1}{\alpha_t} - \frac{1}{\alpha_{t-1}} &\leq \frac{(c_t + \sigma_g^2t)^{1/3}}{c_\alpha} - \frac{(c_t + \sigma_g^2(t-1))^{1/3}}{c_\alpha} \quad (47) \\ &\leq \frac{\sigma_g^2}{3c_\alpha(c_t + \sigma_g^2(t-1))^{2/3}} \\ &\quad (\text{from the concavity of } x^{1/3}: (x+y)^{1/3} - x^{1/3} \leq \frac{y}{3x^{2/3}}) \\ &\leq \frac{2^{2/3}\sigma_g^2}{3c_\alpha^3} \cdot \frac{c_\alpha^2}{(c_t + \sigma_g^2t)^{2/3}} \quad (\text{from } c_t \geq 2\sigma_g^2) \\ &\leq \left(c_\nu - \frac{64\tilde{L}^2}{DN}\right) \alpha_t, \end{aligned} \quad (48)$$

where we use $\alpha_t \leq 1/(16\tilde{L}_gK)$ and the definitions of α_t and c_ν . Based on Eq. (48), substituting Eqs. (39) and (41) in Eq. (21), using $\alpha_t \leq 1/(16KL)$ and $D \geq 1$ and summing over $t = t_q + 1$ to t_{q+1} , we obtain

$$\begin{aligned} &\mathbb{E}[J(\bar{\theta}^{(t_{q+1}+1)})] + \frac{DN\mathbb{E}[\|\bar{\varepsilon}^{(t_{q+1}+1)}\|^2]}{64\tilde{L}^2\alpha_{t_{q+1}}} \\ &\leq \mathbb{E}[J(\bar{\theta}^{(t_q+1)})] + \frac{DN\mathbb{E}[\|\bar{\varepsilon}^{(t_q+1)}\|^2]}{64\tilde{L}\alpha_{t_q}} + \frac{\alpha_{t_q}\mathbb{E}[\|\bar{u}^{(t_q)}\|^2]}{64} \\ &\quad - \sum_{t=t_q+1}^{t_{q+1}} \left(\frac{27\alpha_t}{64} - \frac{\alpha_t L}{2}\right) \mathbb{E}[\|\bar{u}^{(t)}\|^2] + \frac{3c_\nu^2\sigma_g^2}{64\tilde{L}^2} \sum_{t=t_q+1}^{t_{q+1}} \alpha_t^3 \\ &\quad + \frac{c_\nu^2\sigma_g^2}{64D\tilde{L}^2} \alpha_{t_q}^3 - \sum_{t=t_q+1}^{t_{q+1}} \frac{\alpha_t}{2} \mathbb{E}[\|\nabla J(\bar{\theta}^{(t)})\|^2] \end{aligned} \quad (49)$$

Based upon $\bar{u}^{(0)} = 0$, $D \geq 1$ and $\alpha_t \leq 1/(16LK)$, summing over $q \in \{0, 1, \dots, M-1\}$ yields

$$\begin{aligned} \mathbb{E}[J(\bar{\theta}^{(t_M+1)})] &\leq J(\bar{\theta}^{(1)}) + \frac{DN\mathbb{E}[\|\bar{\varepsilon}^{(1)}\|^2]}{64\tilde{L}^2\alpha_0} + \frac{c_\nu^2\sigma_g^2}{16\tilde{L}^2} \sum_{t=0}^{t_M} \alpha_t^3 \\ &\quad - \sum_{t=1}^{t_M} \frac{\alpha_t}{2} \mathbb{E}[\|\nabla J(\bar{\theta}^{(t)})\|^2]. \end{aligned} \quad (50)$$

Note that $T = MK$. Rearranging terms in Eq. (50), we have

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\|\nabla J(\bar{\theta}^{(t)})\|^2] &\leq \frac{2(J(\bar{\theta}^{(1)}) - J^*)}{\alpha_T T} + \frac{D\sigma_g^2}{32\tilde{L}^2\tilde{D}\alpha_0\alpha_T T} \\ &\quad + \frac{c_\nu^2\sigma_g^2}{8\tilde{L}^2\alpha_T T} \sum_{t=0}^T \alpha_t^3, \end{aligned} \quad (51)$$

where we use $\mathbb{E}[\|\varepsilon^{(1)}\|^2] \leq \sigma_g^2/(\tilde{D}N)$ (akin to Eq. (23)). For term $\sum_{t=0}^T \alpha_t^3$, due to $c_t \geq 2\sigma_g^2 \geq \sigma_g^2$, we obtain

$$\sum_{t=0}^T \alpha_t^3 = \frac{c_\alpha^3}{\sigma_g^2} \sum_{t=0}^T \frac{1}{1+t} \leq \frac{c_\alpha^3}{\sigma_g^2} (1 + \ln(1+T)), \quad (52)$$

based on the relationship:

$$\sum_{t=1}^T \frac{x_t}{x_0 + \sum_{\tau=1}^t x_\tau} \leq \ln\left(1 + \frac{\sum_{\tau=1}^t x_\tau}{x_0}\right), \quad (53)$$

with $x_0 = 1$ and $x_1, x_2, \dots, x_T = 1$. From the definition of c_α and c_t and the fact, $c_\nu \leq 2^7 \tilde{L}^2/(DN)$, it is clear that

$$c_T \leq \sigma_g^2 \max\left\{\frac{2^9}{DNK^3}, 2^{12}K^3(DN)^2 - T, 2\right\}. \quad (54)$$

Accordingly, we have $c_T \leq 2^{12}K^3(DN)^2\sigma_g^2$. Drawing on the definition of α_t and c_α , we can bound term $1/(\alpha_T T)$ by

$$\begin{aligned} \frac{1}{\alpha_T T} &\leq \frac{c_T^{1/3}}{c_\alpha T} + \frac{\sigma_g^{2/3}}{c_\alpha T^{2/3}} \quad (\text{from } (x+y)^{1/3} \leq x^{1/3} + y^{1/3}) \\ &\leq \frac{16K\tilde{L}}{T} + \frac{\tilde{L}}{(DN)^{2/3}T^{2/3}}. \end{aligned} \quad (55)$$

Regarding the second term in Eq. (51), we can write

$$\begin{aligned} \frac{D\sigma_g^2}{32\tilde{L}^2\tilde{D}\alpha_0\alpha_T T} &\leq \left(\frac{16\tilde{L}K}{T} + \frac{\tilde{L}}{(DN)^{2/3}T^{2/3}}\right) \frac{\sigma_g^2 c_0^{1/3} D}{32\tilde{L}^2\tilde{D}c_\alpha} \\ &\quad (\text{from Eq. (55) and the definition of } \alpha_0) \\ &\leq \left(\frac{16\tilde{L}K}{T} + \frac{\tilde{L}}{(DN)^{2/3}T^{2/3}}\right) \frac{16\sigma_g^2 K\tilde{L}D}{32\tilde{L}^2\tilde{D}} \\ &\quad (\text{from } c_0 \leq 2^{12}\tilde{L}^3K^3c_\alpha^3) \\ &\leq \frac{8DK^2\sigma_g^2}{\tilde{D}T} + \frac{DK\sigma_g^2}{2\tilde{D}(DN)^{2/3}T^{2/3}}. \end{aligned} \quad (56)$$

Regarding the third term in Eq. (51), we have

$$\begin{aligned} \frac{c_\nu^2 c_\alpha^3}{8\tilde{L}^2\alpha_T T} &\leq \left(\frac{16\tilde{L}K}{T} + \frac{\tilde{L}}{(DN)^{2/3}T^{2/3}}\right) \left(\frac{2^7\tilde{L}^2}{DN}\right)^2 \frac{(DN)^2\sigma_g^2}{8\tilde{L}^5} \\ &\quad (\text{from } c_\nu \leq \frac{2^7\tilde{L}^2}{DN} \text{ and } c_\alpha = (DN)^{2/3}\sigma_g^{2/3}/\tilde{L}) \\ &\leq \frac{2^{15}K\sigma_g^2}{T} + \frac{2^{11}\sigma_g^2}{(DN)^{2/3}T^{2/3}}. \end{aligned} \quad (57)$$

Finally, substituting the bounds in Eqs. (52) and (55) to (57) in Eq. (51), we complete the proof. \square

VI. EXPERIMENT

In this section, we empirically evaluate the proposed method and answer the following key questions:

- 1) How does MFPO perform on standard benchmarks in comparison to the existing FRL algorithms?
- 2) What is the communication/interaction cost of MFPO?
- 3) How does MFPO speed up with the number of agents?
- 4) What is the impact of interaction frequency?

A. Experimental Setup

1) *Environments*: The experiments are carried out on six challenging gym environments [35], including Classic Control tasks (i.e., Cartpole and Pendulum), continuous control MuJoCo tasks (i.e., Halfcheetah and Hopper), and image-based Atari games (i.e., Pong and Breakout), as shown in Fig. 2.

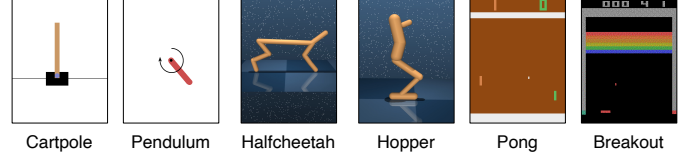


Fig. 2. Benchmark environments.

2) *Baselines*: We compare our proposed algorithm with the following three baseline methods:

- *Federated Policy Gradient with the Byzantine Resilience* (FedPG-BR) [7], a recent fault-tolerant FRL algorithm;
- *Federated Double Q-learning* (Fed-DQN) [8], an FRL algorithm that combines FedAvg with DQN;
- *Federated Soft Actor-Critic* (Fed-SAC), an FRL algorithm combining FedAvg with SAC [36].

3) *Implementation*: The policy is represented as a 2-layer feedforward neural network with 16 hidden units for Classic Control and 256 for the other tasks. It uses ReLU activation functions and Tanh Gaussian outputs. Guided by the analytical results, we set the momentum parameter as $\nu^{(t)} = 1 - 3\alpha^{(t)}$ and the stepsize as $\alpha^{(t)} = 10^{-4} \times 0.99^{-t}$, they are both decrease with updating step t . The discount factor is set to $\gamma = 0.99$. In each round, the number of local updating steps is set as $K = 10$ and $K = 20$ for Classic Control domains and other domains, respectively. We sample $D = 20$ trajectories in each updating step. In addition, we implement the code using Pytorch 1.8.1 framework and run the experiments on Ubuntu 18.04.2 LTS with 8 NVIDIA GeForce RTX A6000 GPUs.

B. Experimental Result

1) *Comparative results*: To answer the first and second questions arised above, we provide comparison results of proposed MFPO with three baselines. As shown in Table I, MFPO yields the best performance with a wide margin among all tasks. Figs. 3 and 4 show that MFPO achieves higher returns while incurring relatively low communication and interaction costs (often by less than 30 rounds and 1000 trajectories), especially in complex and high-dimensional environments. In contrast to the fluctuating performance observed in the baselines, MFPO holds remarkable stability. This demonstrates the efficacy of the momentum-assisted mechanism introduced in controlling policy gradient variates.

2) *Linear speedup*: To answer the third question, we conduct experiments by varying the number of agents from 1 to 50. Fig. 5 reveals a significant improvement in performance as the number of participating agents increases, which aligns well with our theoretical findings. That is, MFPO adeptly controls the inter-agent gradient shift, thereby preventing variance accumulation even with an increasing number of agents involved.

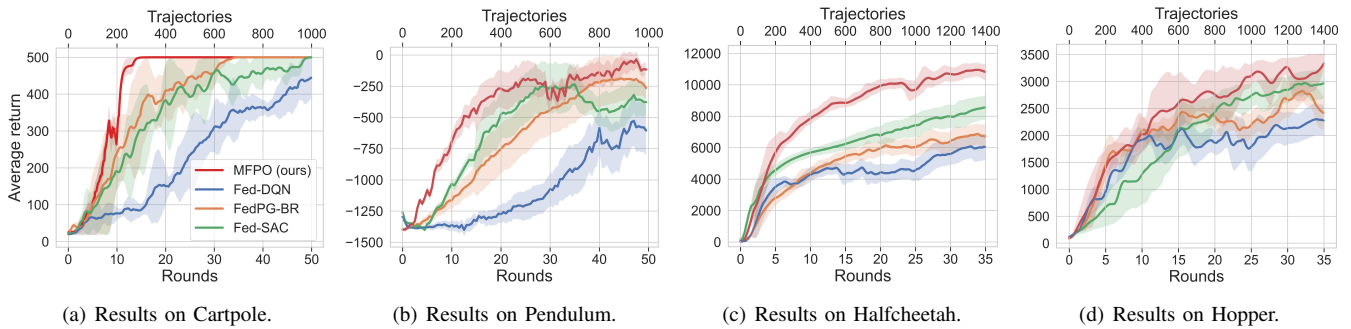


Fig. 3. Convergence results on Classic Control and Mujoco domains.

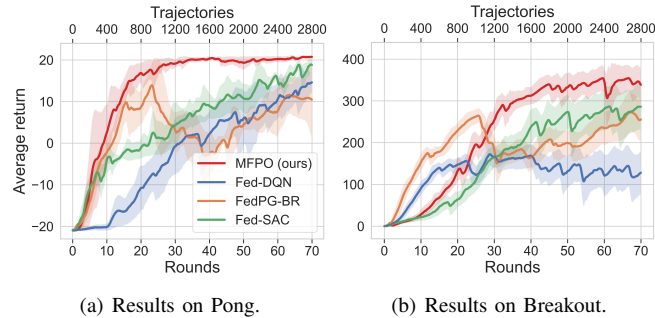


Fig. 4. Convergence results on Atari games.

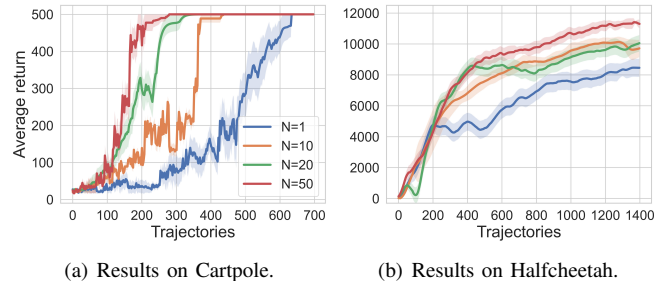


Fig. 5. Performance under different number of agents.

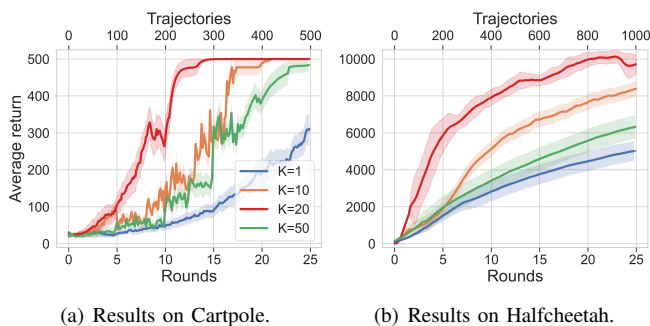


Fig. 6. Impact of local steps on performance.

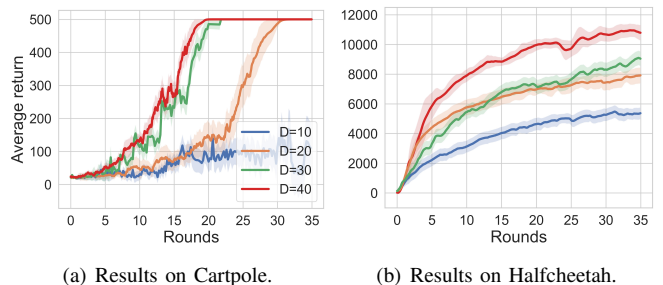


Fig. 7. Performance under different batch sizes.

TABLE I
AVERAGE SCORES ON DIFFERENT ENVIRONMENTS.

Environment	Fed-DQN	FedPG-BR	Fed-SAC	MFPO (ours)
Cartpole	234.3 ± 92.5	446.7 ± 53.3	393.4 ± 55.5	500.0 ± 0.0
Pendulum	-604.1 ± 200.4	-264.2 ± 86.8	-377.9 ± 107.1	-115.5 ± 51.4
Halfcheetah	6055.6 ± 843.1	6719.4 ± 753.8	8561.5 ± 775.3	10794.8 ± 520.2
Hopper	1854.6 ± 294.9	1997.7 ± 297.7	2446.1 ± 376.5	3030.9 ± 68.8
Pong	14.6 ± 2.5	10.5 ± 9.2	18.8 ± 2.2	20.8 ± 0.2
Breakout	128.1 ± 51.6	255.6 ± 47.8	286.2 ± 54.3	336.5 ± 42.1

3) *Impact of local steps*: To validate the impact of the number of local updates, denoted as K , we conduct experiments by varying its value from 1 to 50. The results, displayed in Fig. 6, show that the performance initially improves with an increasing value of K and then starts to decline, consistent with our theoretical results (referring to the last term in Eq. (46)). The reason behind this trend is that a larger value of K exacerbates the gradient shifts across different agents.

4) *Impact of batch sizes*: Fig. 7 shows the impact of the number of collected trajectories in each updating step. As expected, when using a small value of D , the performance degrades dramatically, primarily because of the substantially high variance in the gradient estimator.

VII. CONCLUSION

This paper introduces a new momentum-assisted federated policy optimization algorithm, namely MFPO, to cope with the spatio-temporal non-stationarity of data distributions in FRL. We theoretically show that MFPO offers the best communication and interaction complexities over the existing FRL methods, and provide extensive experiments to corroborate its superior performance over the baselines in continuous and high-dimensional environments. In future work, we will investigate offline/batch FRL approaches that can extract policies only from distributed *static* data with no need to interact with environments. The authors have provided public access to their code at <https://codeocean.com/capsule/1418921/tree/v1>.

ACKNOWLEDGMENTS

This research was supported in part by NSFC under Grant No. 62341201, 62122095, 62072472, 62172445, 62302260, and 62202256, by the National Key R&D Program of China under Grant No. 2022YFF0604502, by CPSF Grant 2023M731956, and by a grant from the Guoqiang Institute, Tsinghua University.

REFERENCES

- [1] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Proc. AISTATS*. PMLR, 2017, pp. 1273–1282.
- [2] D. Xu, T. Li, Y. Li, X. Su, S. Tarkoma, T. Jiang, J. Crowcroft, and P. Hui, "Edge intelligence: Empowering intelligence to the edge of network," *Proc. IEEE*, vol. 109, no. 11, pp. 1778–1837, 2021.
- [3] Z. Zhang, S. Yue, and J. Zhang, "Towards resource-efficient edge ai: From federated learning to semi-supervised model personalization," *IEEE Transactions on Mobile Computing*, 2023.
- [4] B. Liu, L. Wang, and M. Liu, "Lifelong federated reinforcement learning: a learning architecture for navigation in cloud robotic systems," *IEEE Rob. Autom. Lett.*, vol. 4, no. 4, pp. 4555–4562, 2019.
- [5] C. Nadiger, A. Kumar, and S. Abdelhak, "Federated reinforcement learning for fast personalization," in *Proc. AIKE*. IEEE, 2019, pp. 123–127.
- [6] H. H. Zhuo, W. Feng, Y. Lin, Q. Xu, and Q. Yang, "Federated deep reinforcement learning," 2020.
- [7] X. Fan, Y. Ma, Z. Dai, W. Jing, C. Tan, and B. K. H. Low, "Fault-tolerant federated reinforcement learning with theoretical guarantee," in *Proc. NeurIPS*, vol. 34. Curran Associates, Inc., 2021, pp. 1007–1021.
- [8] S. Khodadadian, P. Sharma, G. Joshi, and S. T. Maguluri, "Federated reinforcement learning: Linear speedup under markovian sampling," in *Proc. ICML*, vol. 162. PMLR, 2022, pp. 10997–11057.
- [9] X. Liang, Y. Liu, T. Chen, M. Liu, and Q. Yang, "Federated transfer reinforcement learning for autonomous driving," *arXiv preprint arXiv:1910.06001*, 2019.
- [10] S. Yu, X. Chen, Z. Zhou, X. Gong, and D. Wu, "When deep reinforcement learning meets federated learning: Intelligent multitimescale resource management for multiaccess edge computing in 5g ultradense network," *IEEE Internet Things J.*, vol. 8, no. 4, pp. 2238–2251, 2020.
- [11] H.-K. Lim, J.-B. Kim, J.-S. Heo, and Y.-H. Han, "Federated reinforcement learning for training control policies on multiple iot devices," *Sensors*, vol. 20, no. 5, p. 1359, 2020.
- [12] H. Yang, M. Fang, and J. Liu, "Achieving linear speedup with partial worker participation in non-IID federated learning," in *Proc. ICLR*, 2021.
- [13] A. Anwar and A. Raychowdhury, "Multi-task federated reinforcement learning with adversaries," *arXiv preprint arXiv:2103.06473*, 2021.
- [14] R. Johnson and T. Zhang, "Accelerating stochastic gradient descent using predictive variance reduction," *Proc. NeurIPS*, vol. 26, 2013.
- [15] S. P. Karimireddy, S. Kale, M. Mohri, S. Reddi, S. Stich, and A. T. Suresh, "SCAFFOLD: Stochastic controlled averaging for federated learning," in *Proc. ICML*, vol. 119. PMLR, 2020, pp. 5132–5143.
- [16] P. Khanduri, P. SHARMA, H. Yang, M. Hong, J. Liu, K. Rajawat, and P. Varshney, "Stem: A stochastic two-sided momentum algorithm achieving near-optimal sample and communication complexities for federated learning," in *Proc. NeurIPS*, vol. 34, 2021, pp. 6050–6061.
- [17] Y. Zhang, J. Ren, J. Liu, C. Xu, H. Guo, and Y. Liu, "A survey on emerging computing paradigms for big data," *Chin. J. Electron.*, vol. 26, no. 1, pp. 1–12, 2017.
- [18] Z. Zhou, X. Chen, E. Li, L. Zeng, K. Luo, and J. Zhang, "Edge intelligence: Paving the last mile of artificial intelligence with edge computing," *Proc. IEEE*, vol. 107, no. 8, pp. 1738–1762, 2019.
- [19] Z. Wang, K. Liu, J. Hu, J. Ren, H. Guo, and W. Yuan, "Attrleaks on the edge: Exploiting information leakage from privacy-preserving co-inference," *Chin. J. Electron.*, vol. 32, no. 1, pp. 1–12, 2023.
- [20] H. Hasselt, "Double q-learning," *Proc. NeurIPS*, vol. 23, pp. 2613–2621, 2010.
- [21] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.
- [22] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning," *arXiv preprint arXiv:1509.02971*, 2015.
- [23] H. Cha, J. Park, H. Kim, S.-L. Kim, and M. Bennis, "Federated reinforcement distillation with proxy experience memory," *arXiv preprint arXiv:1907.06536*, 2019.
- [24] C. J. Watkins and P. Dayan, "Q-learning," *Mach. learn.*, vol. 8, pp. 279–292, 1992.
- [25] D. J. Mankowitz, G. Dulac-Arnold, and T. Hester, "Challenges of real-world reinforcement learning," in *ICML Workshop on Real-Life Reinforcement Learn.*, 2019.
- [26] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.
- [27] R. J. Williams, "Simple statistical gradient-following algorithms for connectionist reinforcement learning," *Reinforcement learning*, pp. 5–32, 1992.
- [28] J. Baxter and P. L. Bartlett, "Infinite-horizon policy-gradient estimation," *J. of Artificial Intelligence Res.*, vol. 15, pp. 319–350, 2001.
- [29] A. Agarwal, N. Jiang, S. M. Kakade, and W. Sun, "Reinforcement learning: Theory and algorithms," *Tech. Rep.*, 2019.
- [30] H. Yu, R. Jin, and S. Yang, "On the linear speedup analysis of communication efficient momentum sgd for distributed non-convex optimization," in *Proc. ICML*. PMLR, 2019, pp. 7184–7193.
- [31] M. Papini, D. Binaghi, G. Canonaco, M. Pirodda, and M. Restelli, "Stochastic variance-reduced policy gradient," in *Proc. ICML*. PMLR, 2018, pp. 4026–4035.
- [32] P. Xu, F. Gao, and Q. Gu, "An improved convergence analysis of stochastic variance-reduced policy gradient," in *Uncertainty in Artif. Intell.* PMLR, 2020, pp. 541–551.
- [33] Y. Drori and O. Shamir, "The complexity of finding stationary points with stochastic gradient descent," in *Proc. ICML*. PMLR, 2020, pp. 2658–2667.
- [34] C. Cortes, Y. Mansour, and M. Mohri, "Learning bounds for importance weighting," in *Proc. NeurIPS*, vol. 23. Curran Associates, Inc., 2010.
- [35] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba, "Openai gym," 2016.
- [36] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, "Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor," in *Proc. ICML*. PMLR, 2018, pp. 1861–1870.