# Asynchronous Multi-Agent Reinforcement Learning with General Function Approximation

Anonymous Author(s) Affiliation Address email

## Abstract

1	We study multi-agent reinforcement learning (RL) where agents cooperate through
2	asynchronous communications with a central server to learn a shared environ-
3	ment. Our first focus is on the case of multi-agent contextual bandits with general
4	function approximation, for which we introduce the Async-NLin-UCB algorithm.
5	This algorithm is proven to achieve a regret of $\widetilde{O}(\sqrt{T \dim_E(\mathcal{F}) \log N(\mathcal{F})})$ and a
6	communication complexity of $\widetilde{O}(M^2 \dim_E(\mathcal{F}))$ , where M is the total number of
7	agents and T is the number of rounds, while $\dim_E(\mathcal{F})$ and $N(\mathcal{F})$ are the Eluder
8	dimension and the covering number of function space $\mathcal{F}$ respectively. We then
9	progress to the more intricate setting of multi-agent RL with general function ap-
10	proximation, and present the Async-NLSVI-UCB algorithm. This algorithm enjoys
11	a regret of $\widetilde{O}(H^2\sqrt{K\dim_E(\mathcal{F})\log N(\mathcal{F})})$ and a communication complexity of
12	$\widetilde{O}(HM^2 \dim_E(\mathcal{F}))$ , where H is the horizon length and K the number of episodes.
13	Our findings showcase the provable efficiency of both algorithms for collaborative
14	learning within nonlinear environments and minimal communication overhead.

## 15 **1** Introduction

Multi-agent reinforcement learning (RL) is an important paradigm in RL, and has been successfully 16 applied to real-world tasks such as robotics [Williams et al., 2016, Liu et al., 2019, Ding et al., 2020, 17 Liu et al., 2020, Na et al., 2022], games [Vinyals et al., 2017, Berner et al., 2019, Jaderberg et al., 18 2019, Ye et al., 2020], and control systems [Bazzan, 2009, Yu et al., 2014, 2020, Min et al., 2022, Xu 19 et al., 2023]. By learning cooperatively, agents benefit from sharing learning experiences, enabling 20 them to collectively enhance their decision-making capabilities. This collaborative process is usually 21 accomplished through the utilization of a central server, whose task is to aggregate local data and 22 deliver feedback for the agents. 23

There has been an excellent line of work establishing provably efficient algorithms for multi-agent 24 25 bandits and RL. However, most existing works are restricted to the synchronous setting, where communications between all agents and the server must happen simultaneously. This is impractical since 26 27 in many scenarios the availability of agents may vary and be unpredictable. Ideally, communication should be allowed to happen asynchronously to offer the agents more flexibility. He et al. [2022] and 28 Min et al. [2023] studied this setting respectively for linear contextual bandits and linear Markov 29 Decision Processes (MDPs), both of which assumes linearity in the environment, and introduced 30 algorithms with low regret and communication cost. Yet the linear function class is quite limited, and 31 does not encompass practical reinforcement learning scenarios where nonlinearity is prevalent. 32 To address the aforementioned drawback, in this work, we tackle environments with general function 33 approximation, broadening the applicability of the algorithm to more realistic and complex scenarios. 34 We first delve into multi-agent contextual bandits with general function approximation, where multiple 35 agents interact with homogeneous environments in parallel to solve a common objective. Notably, 36 37 the communication protocol is designed to be flexible and asynchronous, allowing agents to initiate communication with the server and acquire new policy functions whenever the need arises. The 38

39 primary objective is to minimize total regret while reducing communication cost as much as possible. Submitted to 38th Conference on Neural Information Processing Systems (NeurIPS 2024). Do not distribute.

- 40 We propose an algorithm Async-NLin-UCB, which adapts a fully asynchronous communication
- 41 protocol, and leverages various methods for tackling nonlinear function approximation. Despite the
- 42 flexibility of communication, our algorithm performs almost as well as a single agent, in terms of a
- <sup>43</sup> regret that is mostly independent of the number of agents and a low communication cost.
- 44 We then progress to multi-agent RL with general function approximation under similar requirements
- and objectives. We propose an algorithm named Async-NLSVI-UCB based on Least-Squares Value
   Iteration (LSVI) to learn the underlying Markov decision processes (MDPs), which demonstrates
- Iteration (LSVI) to learn the underlying Markov decision processes (MDPs
   similar advantages with provably low regret and communication cost.
- 47 similar advantages with provably low registrand communication
   48 Our main contributions are summarized in the following:
- 49 For asynchronous multi-agent nonlinear contextual bandits, we propose the algorithm
- Async-NLin-UCB, which enjoys an  $\widetilde{O}(\sqrt{T \dim_E(\mathcal{F}) \log N(\mathcal{F})} + \dim_E(\mathcal{F}))$  regret and an
- 51  $\widetilde{O}(M^2 \dim_E(\mathcal{F}))$  communication complexity, where  $\dim_E(\mathcal{F})$  and  $N(\mathcal{F})$  are respectively the
- Eluder dimension and the covering number of function space  $\mathcal{F}$ .
- For asynchronous multi-agent nonlinear MDPs, we propose the algorithm Async-NLSVI-UCB, which enjoys an  $\widetilde{O}(H^2\sqrt{K\dim_E(\mathcal{F})\log N(\mathcal{F})} + H^2\dim_E(\mathcal{F}))$  regret and a communication
- complexity of  $\widetilde{O}(HM^2 \dim_E(\mathcal{F}))$ .
- At the core of our algorithm, we design a *communication criterion* in order to tackles the challenges
   posed by both asynchronous communication and the nonlinearity of function approximation. To
   guarantee a low communication cost, we propose a low switching communication criterion that
   allows the agent to trigger communication rounds.

• We carefully design our *download content* from server to local agents, which consist only of decision and bonus functions, with no mention of any specific historical data. This effectively protects user data against exposure by disallowing local users from obtaining the data of others.

Notation. We use lower case letters to denote scalars. We denote by [n] the set  $\{1, \ldots, n\}$ . For two positive sequences  $\{a_n\}$  and  $\{b_n\}$  with  $n = 1, 2, \ldots$ , we write  $a_n = O(b_n)$  if there exists an absolute constant C > 0 such that  $a_n \le Cb_n$  holds for all  $n \ge 1$ . We use  $\widetilde{O}(\cdot)$  to further hide the polylogarithmic factors. For two non-negative integers a, b satisfying a < b and a sequence  $\{s_i\}$ indexed by integers i, we use  $s_{[a:b]}$  to denote the subsequence  $\{s_a, s_{a+1}, \cdots, s_b\}$ .

## 68 2 Related Work

### 69 2.1 Multi-Agent Bandits

First, there is a multitude of previous work on distributed or federated multi-armed bandits and 70 stochastic linear bandits [Liu and Zhao, 2010, Szorenyi et al., 2013, Landgren et al., 2016, Chakraborty 71 et al., 2017, Landgren et al., 2018, Martínez-Rubio et al., 2019, Sankararaman et al., 2019, Wang et al., 72 2020a,c, Zhu et al., 2021, Huang et al., 2021]. For the more realistic setting of contextual bandits, most 73 previous work are within the scope of linear contextual bandits with synchronized communication. 74 Korda et al. [2016] introduced two novel distributed confidence ball (DCB) algorithms for linear 75 bandit problems in peer-to-peer networks. Wang et al. [2020c] considered both P2P and star-shaped 76 communication, achieving near-optimal regret and low communication cost that is largely independent 77 of the time horizon in their algorithm DisLinUCB. Dubey and Pentland [2020] proposed FedUCB, 78 an algorithm focusing on differential-privacy. 79

Li and Wang [2022] first considered an asynchronous communication protocol and proposed the 80 algorithm Async-LinUCB with near-optimal regret, yet the algorithm contains a download step 81 for all agents triggered by the central server. Their results are flexible and contains a parameter to 82 control the trade-off between regret and communication cost. He et al. [2022] improved the setting 83 to a fully asynchronous communication, proposing the algorithm FedLinUCB with near-optimal 84 regret of  $\tilde{O}(d\sqrt{T})$  and low communication cost of  $\tilde{O}(dm^2)$ , comparable to the benchmark in single-85 agent contextual linear bandits [Abbasi-Yadkori et al., 2011]. We consider the same communication 86 protocol in our results. A summary of these results along with ours can be found in the first four rows 87 of Table 1. 88

## 89 2.2 Multi-Agent RL

Multi-agent reinforcement learning is decidedly more challenging than contextual bandits. There is
 also a vast literature on this setting, with many works discussing different aspects of multi-agent RL

Algorithm	Regret	Communication	Fully asynchrnous
DisLinUCB [Wang et al., 2020c]	$d\sqrt{MT}\log^2 T$	$d^3 M^{3/2}$	×
Async-LinUCB [Li and Wang, 2022]	$dM^{(1-\gamma)/2}\sqrt{T}\log T$	$dM^{1+\gamma}\log T$	×
FedLinUCB [He et al., 2022]	$d\sqrt{T}\log T$	$dM^2 \log T$	1
Async-NLin-UCB (ours)	$\sqrt{\dim_E \log NT} \log T$	$\dim_E M^2 \log^2 T$	1
Coop-LSVI [Dubey and Pentland, 2021]	$d^{3/2}H^2\sqrt{MK}\log K$	$dHM^3$	×
Async-Coop-LSVI-UCB [Min et al., 2023]	$d^{3/2}H^2\sqrt{K\log K}$	$dHM^2 \log K$	1
Async-NLSVI-UCB (ours)	$\sqrt{\dim_E \log N} H^2 \sqrt{K} \log K$	$\dim_E HM^2 \log^2 K$	✓

Table 1: Comparison of our result against baseline methods for multi-agent contextual bandits and MDPs. Note that the first four rows are for contextual bandits, and the last three are for reinforcement learning. Only our algorithms are in the general function approximation setting. We abbreviate  $\dim_E = \dim_E(\mathcal{F})$  and  $N = N(\mathcal{F})$ , and hide logarithmic factors. For algorithms with synchronized communication, each communication round actually corresponds to M rounds in asynchronous settings, which explains the extra M terms.

<sup>92</sup> than ours. For example, there are works focusing on convergence guarantees [Zhang et al., 2018b,a,

Wai et al., 2018], non-stationary or heterogeneous environments [Lowe et al., 2017, Yu et al., 2021,

Dubey and Pentland, 2021, Kuba et al., 2022, Liu et al., 2022, Jin et al., 2022], and deep federated RL
[Clemente et al., 2017, Espeholt et al., 2018, Horgan et al., 2018, Nair et al., 2015, Zhuo et al., 2019],

[Clemente et al., 2017, Espeholt et al., 2018, Horgan et al., 2018, Nair et al., 2015, Zhuo et al., 2019],
 to name a few. We refer to a recent survey on federated reinforcement learning Qi et al. [2021] for a

<sup>97</sup> more comprehensive summary.

98 Narrowing it down to multi-agent RL with function approximation, the benchmark is the LSVI-UCB

algorithm in the single-agent setting [Jin et al., 2020], with an  $\widetilde{O}(d^{3/2}H^2\sqrt{K})$  regret. Dubey and

Pentland [2021] proposed CoopLSVI for multi-agent linear MDPs, which requires a synchronized communication through central server, and proves a regret of  $\tilde{O}(d^{3/2}H^2\sqrt{MK})$ . They also extended their result to the heterogeneous setting. Min et al. [2023] considered the fully asynchronous setting and introduced the Async-Coop-LSVI-UCB algorithm, with a  $\tilde{O}(d^{3/2}H^2\sqrt{K})$  regret not dependent on the number of agents M, as well as a low communication cost. A summary of these results along

with ours can be found in the last three rows of Table 1.

### **106 2.3 General function approximation**

Reinforcement learning with general function approximation extends the well-studied case of linear 107 MDPs to more general classes of MDPs, and has gained a lot of traction in recent years [Wang et al., 108 2020b, Jin et al., 2021, Foster et al., 2023, Du et al., 2021, Agarwal and Zhang, 2022, Agarwal 109 et al., 2023]. Previous works focus on different measures of complexity for the function classes, for 110 example the Bellman rank proposed by Jiang et al. [2017], the Bellman Eluder dimension introduced 111 in Jin et al. [2021], the Decision-Estimation Coefficient in Foster et al. [2023], and generalized Eluder 112 dimension in Agarwal et al. [2023]. Our work considers the Eluder dimension with the introduction 113 of uncertainty estimators  $D^2$ , which has been widely utilized to establish results in RL with general 114 function approximation [Agarwal et al., 2023, Zhao et al., 2023, Ye et al., 2023, Di et al., 2023]. 115

## **116 3 Preliminaries**

In this section, we introduce the formal definition of both multi-agent nonlinear contextual bandits
 and MDPs and some related concepts, and discuss the asynchronous communication protocol.

## 119 3.1 Multi-Agent Contextual Bandits with General Function Approximation

We assume a global action set  $\mathcal{A}$  that is known to all agents. At each round  $t \in [T]$ , a single arbitrary agent  $m_t \in [M]$  is chosen to participate. The agent receives a contextual decision set  $\mathcal{A}_t \subseteq \mathcal{A}$  and chooses from the set an action  $a_t \in \mathcal{A}_t$  to perform, and subsequently receives a random reward  $r_t$ . <sup>123</sup> The assumption of general function approximation is that the reward is generated according to

$$r_t = f^*(a_t) + \eta_t,\tag{1}$$

- where  $f^*$  is the ground truth objective function, and  $\eta_t$  is a random noise variable. We assume the
- the objective function lies within a known function class  $\mathcal{F}$ . In addition, we also make the following
- assumptions regarding the function class and noise variables, which are standard assumptions for
- contextual bandits [Abbasi-Yadkori et al., 2011, He et al., 2022]:
- Assumption 3.1. Suppose the following conditions hold for the contextual bandits environment:
- For any  $f \in \mathcal{F}$  and  $a \in \mathcal{A}$ ,  $|f(a)| \leq 1$ ;
- $\eta_t$  is *R*-sub-Gaussian conditioned on data history:  $\mathbb{E}\left[e^{\lambda\eta_t} | a_{1:t}, m_{1:t}, r_{1:t-1}\right] \leq \exp(R^2\lambda^2/2), \forall \lambda.$
- 131 Learning Objective. The primary goal of contextual bandits is to minimize the cumulative regret

$$\operatorname{Reg}(T) = \sum_{t=1}^{T} [f^*(a_t) - \max_{a \in \mathcal{A}_t} f^*(a)].$$

Notice that this summation is across all time steps does not depend on agent participation order, as should be the case for the resulting regret bound. To achieve this goal, agents are allowed to communicate with the server to upload their interaction history and update their policy. The secondary learning objective is to reduce communication overhead. We will explain the communication protocol further in Section 3.4.

## 137 3.2 Multi-Agent Episodic MDPs with General Function Approximation

We consider episodic MDPs, which are a classic family of models in reinforcement learning [Sutton 138 and Barto, 2018]. It is characterized by the following elements, which we assume to be homogeneous 139 across all agents: a state space S, an action space A, the horizon length H, transition probability functions  $\mathbb{P} = \{\mathbb{P}_h(\cdot|\cdot,\cdot)\}_{h=1}^H$  and reward functions  $\{r_h(\cdot,\cdot)\}_{h=1}^H$ ). Similar to the bandit case, for each episode  $k = 1, \dots, K$ , a single agent  $m = m_k$  is chosen to participate. An episode 140 141 142 k begins with an initial state  $s_1^k$ , which is drawn from an unknown fixed distribution. Then for 143 steps  $h = 1, \dots, H$ , the participating agent m selects an action  $a_h^k$  based on the observed state 144  $s_h^k$ . After each action, the agent receives a reward  $r_h^k = r_h(s_h^k, a_h^k)$ , where  $r_h : S \times A \to \mathbb{R}$  is the reward function at step h. Here for the sake of convenience, we assume the reward function to be 145 146 deterministic, but it is not difficult to generalize our result to stochastic rewards. We also assume 147  $r_h(s, a) \in [0, 1]$  for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$  without loss of generality. The environment then transitions 148 to the next state according to  $s_{h+1}^k \sim \mathbb{P}_h(\cdot | s_h^k, a_h^k)$ , where  $\mathbb{P}_h$  is the transition probability at step h. The episode terminates when  $r_H$  is observed. 149 150

The strategy an agent employs to interact with the environment is called the agent's *policy*, which can be described by a set of decision functions  $\pi = {\pi_h}_{h=1}^H$ , where  $\pi_h : S \to A$  is the decision function at level *h*, mapping the current state to an action to select.

154 Value Functions. For any policy  $\pi = {\pi_h}$ , we define Q-value functions and V-value functions:

$$Q_{h}^{\pi}(s_{h},a_{h}) := \mathbb{E}\left[\sum_{h'=h}^{H} r_{h'}(s_{h'},a_{h'}) \middle| s_{h},a_{h}\right], \quad V_{h}^{\pi}(s_{h}) := \mathbb{E}\left[\sum_{h'=h}^{H} r_{h'}(s_{h'},a_{h'}) \middle| s_{h}\right], \quad (2)$$

where the expectation is taken over the trajectory  $(s_1, a_1, \dots, s_h, a_h)$ , determined by the transition probability functions  $\mathbb{P}$  and policy  $\pi$ . The optimal strategy  $\pi^*$  is the maximizer of the value functions:

$$\pi^* := \operatorname{argmax}_{\pi} V_1^{\pi}(s_1), \forall s_1.$$

We also have optimal value functions  $Q_h^* := Q_h^{\pi^*}$  and  $V_h^* := V_h^{\pi^*}$ , which satisfy Bellman equations

$$Q_{h}^{*}(s_{h}, a_{h}) = r_{h}(s_{h}, a_{h}) + \mathbb{E}\left[V_{h+1}^{*}(s_{h+1})\big|s_{h}, a_{h}\right], \quad V_{h}^{*}(s_{h}) = \max_{a \in \mathcal{A}} Q_{h}^{*}(s_{h}, a).$$
(3)

**Function Approximation.** We approximate Q-value functions with function classes  $\{\mathcal{F}_h\}_{h=1}^H$ , which contain real value functions with domain  $S \times A$ . One basic assumption is that  $Q_h^* \in \mathcal{F}_h$  for all steps  $h \in [H]$ . Now with the convention that functions at level H + 1 are uniformly zero, i.e.,  $f_{H+1} = 0$ , we define the Bellman operator  $\mathcal{T}_h$ :

$$(\mathcal{T}_h f_{h+1})(s_h, a_h) := \mathbb{E} \left| r_h(s_h, a_h) + f_{h+1}(s_{h+1}) \right| s_h, a_h \right|_{\mathcal{T}_h}$$

- and we expect  $\mathcal{T}_h$  to map any function in  $\mathcal{F}_{h+1}$  to a function in  $\mathcal{F}_h$ , i.e.,  $\mathcal{T}_h \mathcal{F}_{h+1} \subseteq \mathcal{F}_h$ . This is called the completeness assumption, which is a fundamental assumption in RL with general function
- approximation [Wang et al., 2020b, Jin et al., 2021].

Learning Objective. The primary goal in multi-agent MDPs is to minimize the cumulative regret over K episodes

$$\operatorname{Reg}(K) = \sum_{k=1}^{K} \left[ V_1^*(s_1^k) - V_1^{\pi_{m,k}}(s_1^k) \right],$$

where  $\pi_{m,k}$  is the policy of agent  $m = m_k$  at round k, while the secondary objective is to minimize 159 the communication cost. 160

#### 3.3 Eluder Dimension and Covering Number 161

To measure the complexity of the learning objective, Russo and Van Roy [2013] first proposed the 162 concept of Eluder dimension, which we define below. 163

**Definition 3.2** ( $\epsilon$ -dependence). For a function class  $\mathcal{F}$  on domain  $\mathcal{D}$ , a point  $z \in \mathcal{D}$  is  $\epsilon$ -dependent 164 on  $\mathcal{Z} \subseteq \mathcal{D}$  if, for any  $f_1, f_2 \in \mathcal{F}$  satisfying  $\sqrt{\sum_{z' \in \mathcal{Z}} (f_1(z') - f_2(z'))^2} \leq \epsilon$ , it must hold that  $|f_1(z) - f_2(z)| \leq \epsilon$ . Accordingly, z is  $\epsilon$ -independent of  $\mathcal{Z}$  if it is not  $\epsilon$ -dependent on  $\mathcal{Z}$ . 165

166

**Definition 3.3** (Eluder dimension). The  $\epsilon$ -Eluder dimension dim<sub>E</sub>( $\mathcal{F}, \epsilon$ ) is the length of the longest 167 sequence of elements in  $\mathcal{D}$  satisfying that, for some  $\epsilon_0 > \epsilon$ , each element is  $\epsilon_0$ -independent of the set 168 consisting of its predecessors. 169

It has been demonstrated that the Eluder dimension roughly corresponds to regular dimension 170 concepts in linear and quadratic cases [Russo and Van Roy, 2013], and that the Eluder family is 171 strictly larger than the generalized linear class [Li et al., 2022]. Note that our Eluder definition can be 172

applied to either the contextual bandit case with  $\mathcal{D} = \mathcal{A}$  or the MDPs case with  $\mathcal{D} = \mathcal{S} \times \mathcal{A}$ . 173

We also introduce covering number for function classes [Wainwright, 2019] in the following: 174

**Definition 3.4** (Covering number). An  $\epsilon$ -cover of  $\mathcal{F}$  is any subset  $\mathcal{F}_{\epsilon} \subseteq \mathcal{F}$  such that for any  $f \in \mathcal{F}$ , 175 there exists  $f' \in \mathcal{F}_{\epsilon}$  that  $||f - f'||_{\infty} \leq \epsilon$ . The *covering number* of  $\mathcal{F}$ , denoted by  $N(\mathcal{F}, \epsilon)$ , is the 176 minimal cardinality of its  $\epsilon$ -cover. 177

3.4 Communication Protocol 178

We consider a star-shaped communication model [He et al., 2022, Min et al., 2023], where the agents 179 communicate through a central server to collaborate. To ensure asynchronous communication, we 180 mandate that all communications must be initiated by a participating agent. Specifically, at the end of 181 a time step / episode, the agent will decide whether or not to trigger a communication round. If so, 182 the agent uploads its local data history and receives some global data for future decision making. The 183 *communication cost* is the total number of communication rounds initiated by the agents. 184

One variability is the form of global data that the communicating agent downloads from server. It 185 may be tempting to have the server send all its stored trajectories to the agent for future decision 186 making, but this will unnecessarily expose other agents' data to the current participating agent. We 187 will come back to this issue and our solution in Section 4.2. 188

#### 4 **Multi-Agent Contextual Bandits** 189

In this section, we introduce the Asynchronous Nonlinear UCB (Async-NLin-UCB) algorithm 190 designed for multi-agent contextual bandits with general function approximation, and provide a 191 theoretical result for its regret and communication cost. 192

#### 4.1 Algorithm: Async-NLin-UCB 193

Algorithm 1 takes as input the total number of time steps T, regularization parameter  $\lambda$ , communica-194 tion parameter  $\alpha$  and exploration radii  $\{\beta_t\}_{t=1}^T$ . 195

In the algorithm, there are some variables that go through different versions as t progresses through 196

 $1, \dots, T$ . For clarity, here we give them an extra subscript t to denote the version of that variable 197 198 before (not included) the least squares calculation on Line 12 at round t.

Throughout the learning process, the server maintains a global history set  $Z_t^{ser}$  that stores action-199

reward pairs  $(a, r) \in \mathcal{A} \times [0, 1]$ , initialized on Line 2 and updated only during communication rounds. 200

Each local agent m maintains a decision function  $f_{m,t}$  for taking action, a bonus function  $b_{m,t}$  for 201

checking communication criterion, and a local data history set  $Z_{m,t}^{\text{loc}}$ , all initialized on Line 3. Each 202 step of Algorithm 1 contains two parts: local exploration and server updates. 203

**Part I: Local Exploration.** At step t a single agent  $m = m_t$  is active (Line 5). It receives a decision 204

set, finds the greedy action according to its decision function  $f_{m,t}$ , receives a reward, and updates its 205

Algorithm 1 Async-NLin-UCB

- 1: **Input:** total number of rounds T, parameters  $\lambda$ ,  $\alpha$ ,  $\beta_t$  for t = 1, ..., T.
- 2: Server init: Set  $Z^{ser} = \emptyset$ .
- 3: Local init: For all  $m \in [M]$ , set  $f_m = 1$ ,  $b_m = \mathcal{B}_{\mathcal{A}}(\emptyset, \mathcal{F}; \lambda, \beta_0)$  and  $Z_m^{\text{loc}} = \emptyset$ .
- 4: for t = 1, ..., T do
- Agent  $m = m_t \in [M]$  is active. 5:
- Receive decision set  $A_t \subseteq A$  and take action  $a_t \in \operatorname{argmax}_{a \in D_t} f_m(a)$  and receive reward  $r_t$ . 6:
- Update local history  $Z_m^{\text{loc}} = Z_m^{\text{loc}} \cup \{(a_t, r_t)\}.$ if switch condition (4) is met **then** 7:
- 8:
- Send new data  $Z_m^{\text{loc}}$  to server. 9:
- on server: 10:
- Update  $Z^{\text{ser}} = Z^{\text{ser}} \cup Z_m^{\text{loc}}$ . 11:
- Calculate  $\hat{f}$  according to (5) and the bonus function  $b = \mathcal{B}_{\mathcal{A}}(Z^{\text{ser}}, \mathcal{F}; \lambda, \beta_t)$ . 12:
- Send  $\hat{f} + b$  and b to agent m. 13:
- 14: end of server
- Agent m receives decision and bonus functions  $f_m = \hat{f} + b$ ,  $b_m = b$ , then set  $Z_m^{\text{loc}} = \emptyset$ . 15:
- 16: end if
- 17: end for

After exploration, the agent checks if the switch condition is true using its bonus function: 207

$$\sum_{(a,r)\in Z_{m,t}^{\text{loc}}} b_{m,t}^2(a) / \left(\beta_{t'}^2 + \lambda\right) \ge \alpha,\tag{4}$$

where t' is the last time step when agent m communicated with the server. If so, the agent initiates a 208 communication round and uploads its local data (Line 9), prompting the server to begin global policy 209 updates. We will discuss the reasons behind this switch condition in Section 4.2. 210

Part II: Server Updates. After receiving a new local data history from an agent, the server merges 211 the data into its global dataset  $Z_t^{\text{ser}}$  (Line 11), and calculate a function  $\hat{f}_{t+1} \in \mathcal{F}$  which minimizes 212 the sum of squares error according to the current dataset  $Z_t^{\text{ser}}$  (Line 12): 213

$$\widehat{f}_{t+1} = \operatorname{argmin}_{f \in \mathcal{F}} \sum_{(a,r) \in Z_t^{\operatorname{ser}}} (f(a) - r)^2.$$
(5)

The next step is to obtain a bonus function  $b_{t+1}$  from the oracle  $\mathcal{B}_{\mathcal{A}}$  from Definition 4.1 (Line ??). 214 We discuss the specifics of this construction in detail up next in Section 4.2. Finally, the server sends 215 the optimistic value function  $\hat{f}_{t+1} + b_{t+1}$  and the bonus function  $b_{t+1}$  back to agent m for future exploration and updates; agent m also resets its local data history to an empty set (Lines 13 and 15). 216 217

#### 4.2 Uncertainty Estimators and Bonus Functions 218

In this section, we introduce uncertainty estimators and bonus functions, and give a detailed explana-219 tion for our communication criterion (4). Most of these apply to the MDPs setting as well. 220

**Uncertainty Estimators.** First we define the uncertainty estimator of new data *a* against data history 221 Z, which is considered in many works on bandits and RL with general function approximation 222

[Gentile et al., 2022, Agarwal et al., 2023]: 223

$$D_{\lambda,\mathcal{F}}(a;Z) = \sup_{f_1,f_2\in\mathcal{F}} |f_1(a) - f_2(a)| / \sqrt{\lambda + \sum_{(a',r)\in Z} |f_1(a') - f_2(a')|^2},$$
(6)

here  $\lambda$  is the regularization parameter,  $\mathcal{F}$  is a function class. Intuitively, the uncertainty estimator 224 measures the difference between functions on new data a against the difference on historical data Z. 225 Switch Condition Based On Uncertainty Estimators. The determinant-based criterion is a 226 common technique used in contextual bandits and RL with linear function approximation to reduce 227 policy switching or communication cost [Abbasi-Yadkori et al., 2011]. For nonlinear function 228 approximation, one can use uncertainty estimators to formulate a new form of switch condition: 229

$$\sum_{(a,r)\in Z_t^{\text{new}}} D^2_{\lambda,\mathcal{F}}(a; Z_t^{\text{old}}) \ge \alpha.$$
(7)

where we use  $Z_t^{\text{new}}$  and  $Z_t^{\text{old}}$  to denote newly accumulated data and old historical data. This criterion 230 has a similar function as the determinant-based criterion in linear settings. Parameter  $\alpha$  controls 231 communication frequency: smaller  $\alpha$  indicates more frequent communication, more accurate decision 232 functions and smaller regret, thus implying a trade-off between regret and communication cost. 233

Next, we introduce bonus functions obtained through oracles that **Bonus Function Oracle.** 234 approximate the uncertainty estimators. 235

**Definition 4.1** (Bonus Function Oracle  $\mathcal{B}_{\mathcal{D}}$ ). Given domain  $\mathcal{D}$ , the oracle  $\mathcal{B}_{\mathcal{D}}(Z, \mathcal{F}; \lambda, \beta)$  takes the following as inputs: a dataset Z consisting of a series of data points (z, e), where  $z \in \mathcal{D}$  and e is some additional data content; function class  $\mathcal{F}$  with functions  $f : \mathcal{D} \to \mathbb{R}_{\geq 0}$ ; regularization parameter  $\lambda$ and exploration radius  $\beta$ . It returns a function  $b \in \mathcal{W}_{\mathcal{D}} : \mathcal{D} \to \mathbb{R}_{>0}$  satisfying for any  $z \in \mathcal{D}$  that

• 
$$b(z) \ge \max\left\{ \left| f_1(z) - f_2(z) \right| : f_1, f_2 \in \mathcal{F}, \sum_{(z,e) \in Z} \left( f_1(z) - f_2(z) \right)^2 \le \beta^2 \right\};$$

• 
$$D_{\lambda,\mathcal{F}}(z;Z) \leq b(z)/\sqrt{\beta^2 + \lambda} \leq C_{\mathcal{B}} D_{\lambda,\mathcal{F}}(z;Z),$$

where  $C_{\mathcal{B}}$  is an absolute constant.

Remark 4.2. Similar bonus function oracles have been proposed in previous works (Definition 3 243 244 in Agarwal et al. [2023]). The accessibility of these oracles is also supported by previous works that proposed methods to compute bonus functions [Kong et al., 2023, Wang et al., 2020b]. In this 245 definition, we leave the domain and data format to be variable so the oracle can be applied to both 246 contextual bandits and MDPs. For bandits, the domain is A, and the data format has z = a and e = r. 247 The first property of the bonus function guarantees the optimism of decision functions  $f_{t+1} + b_{t+1}$ 248 (see Lemma 6.1 for MDPs or Lemma A.2 for bandits), while the second property links bonuses to 249 uncertainty estimators. 250

**Switch Condition Based On Bonus Functions.** If we try to adapt the switch condition (7) in our setting, a local agent will require access to historical data  $Z_t^{\text{old}}$  to calculate uncertainty estimators  $D_{\lambda,\mathcal{F}}^2(a; Z_t^{\text{old}})$ . For multi-agent learning, this dataset consists of the collective data from all agents, and giving local agent access is a clear violation of data privacy. Our solution is to let local agents download bonus functions and set communication criterion to (4), using bonus functions instead of uncertainty estimators.

**Decision Functions Based On Bonus Functions.** Another benefit of introducing the bonus function is evident from our exploration method in line 6. A common practice for nonlinear RL algorithms is to construct *confidence sets* of functions during policy update, and find the optimal function within the confidence sets during exploration [Agarwal et al., 2023, Ye et al., 2023]. However, in a multi-agent setting, this would involve the download of confidence sets, which is impractical due to the complex nature of function classes. With the bonus function, local agents need only download the *decision function* from the server for future exploration, which for contextual bandits is simply  $\hat{f}_{t+1} + b_{t+1}$ .

#### 264 4.3 Theoretical Results

Our main results for Algorithm 1 are summarized in the following theorem, which provides a regret upper bound and communication complexity order.

**Theorem 4.3.** By taking  $\gamma = O(1/T)$ ,  $\beta_t = C_{\beta,1}(\sqrt{\lambda} + RC(M,\alpha)\log(3MN(\mathcal{F},\gamma)/\delta))$  and  $C(M,\alpha) = \sqrt{1 + M\alpha}(\sqrt{1 + M\alpha} + M\sqrt{\alpha})$ , the regret of Algorithm 1 within T rounds is

$$O\left(\sqrt{T}\widetilde{\beta}_1\sqrt{(1+M\alpha)\dim_E}\log(T/\min\{1,\lambda\}) + (1+M\alpha)\dim_E\log^2(T/\min\{1,\lambda\})\right),$$

where we abbreviate  $\dim_E := \dim_E(\mathcal{F}, \lambda/T)$ ; the total communication complexity is

$$O((1+M\alpha)^2/\alpha \dim_E \log^2(T/\min\{1,\lambda\}))$$

*Remark* 4.4. When reduced to linear contextual bandits, where  $\dim_E(\mathcal{F}, \lambda/T) = \widetilde{O}(d)$  and log  $N(\mathcal{F}, \gamma) = \widetilde{O}(d)$ , our result on regret correspond exactly to Theorem 5.1 of He et al. [2022], except for an extra  $1 + M\alpha$  term in the communication cost, an unimportant term when taking  $\alpha = 1/M^2$  that comes from the complication of communication cost analysis in nonlinear settings.

## 274 **5** Multi-Agent Reinforcement Learning

In this section, we introduce the Asynchronous Nonlinear Least Squares Value Iteration UCB (Async-NLin-UCB) algorithm for multi-agent MDPs with general function approximation, and a corresponding theoretical result.

### 278 5.1 Algorithm: Async-NLSVI-UCB

To better represent the elements in the datasets, we sometimes use  $o_h$  to represent the tuple ( $s_h, a_h, r_h, s_{h+1}$ ) and  $z_h$  to represent ( $s_h, a_h$ ) when there is no confusion. Similar to the bandit case, we give some variables an extra subscript k here for clarity, which denotes the version of the variable before (not included) Line 14 at episode k. Algorithm 2 Federated Nonlinear MDPs

1: **Input:** total number of rounds K, parameters  $\lambda$ ,  $\alpha$ ,  $\beta_{k,h}$  for k = [K] and  $h \in [H]$ 2: Server init: Set  $Z_h^{\text{ser}} = \emptyset$  for all  $h \in [H]$ . 3: Local init:  $\forall m \in [M]$  and  $h \in [H]$ , set  $Q_{m,h} = 1$ ,  $b_{m,h} = \mathcal{B}(\emptyset, \mathcal{F}_h; \lambda, \beta_{0,h})$ ,  $Z_{m,h}^{\text{loc}} = \emptyset$ . 4: for k = 1, ..., K do Agent  $m = m_k \in [M]$  is active and receives initial state  $s_1^k \in S$ . 5: for h = 1, ..., H do 6: Take action  $a_h^k = \operatorname{argmax}_{a \in \mathcal{A}} Q_{m,h}(s_h^k, a)$ , receive reward  $r_h^k$  and next state  $s_{h+1}^k$ . Update  $Z_{m,h}^{\text{loc}} = Z_{m,h}^{\text{loc}} \cup \{(s_h^k, a_h^k, r_h^k, s_{h+1}^k)\}$ . 7: 8: 9: end for 10: if switch condition (8) is met then Send new data  $\{Z_{m,h}^{\text{loc}}\}_{h\in[H]}$  to server. 11: 12: on server: Update  $Z_h^{\text{ser}} = Z_h^{\text{ser}} \cup Z_{m,h}^{\text{loc}}$ . Initialize  $Q_{H+1} = V_{H+1} = 0$ . for  $h = H, H - 1, \dots, 1$  do 13: 14: 15: Calculate  $\widehat{f}_h$  according to (9) and bonus function  $b_h = \mathcal{B}_{S \times \mathcal{A}}(Z_h^{\text{ser}}, \mathcal{F}_h; \lambda, \beta_{k,h}).$ 16: Calculate  $Q_h$  and  $V_h$  according to (11). 17: 18: end for Send  $\{Q_h\}_{h=1}^H$  and  $\{b_h\}_{h=1}^H$  to agent m. 19: 20: end of server Agent m receives  $Q_{m,h} = Q_h$ ,  $b_{m,h} = b_h$  and resets  $Z_{m,h}^{\text{loc}} = \emptyset$  for all  $h \in [H]$ . 21: 22: end if 23: end for

The server maintains global historical datasets  $Z_{k,h}^{\text{ser}}$  containing sequences of tuples  $(s_h, a_h, r_h, s_{h+1})$ , initialized in Line 2. Each local agent *m* maintains optimistic value functions  $\{Q_{m,k,h}\}_{h=1}^{H}$ , bonus functions  $\{b_{m,k,h}\}_{h=1}^{H}$ , and local datasets  $\{Z_{m,k,h}^{\text{loc}}\}_{h=1}^{H}$ , all initialized in Line 3.

Each episode k of Algorithm 2 also consists of the two parts local exploration and server updates.

**Part I: Local Exploration.** At step k an agent  $m = m_k$  is active (Line 5). It interacts with the environment by executing the greedy policy according to  $\{Q_{m,k,h}\}_{h=1}^{H}$ , obtaining a trajectory  $\{(s_h^k, a_h^k, r_h^k, s_{h+1}^k)\}_{h=1}^{H}$ , which is then stored into the local historical datasets  $Z_{m,k,h}^{\text{loc}}$  (lines 6 - 9).

After exploration, the agent checks for the following switch condition: there exists  $h \in [H]$  so that

$$\sum_{o_h \in \mathbb{Z}_{m,k,h}^{\mathrm{loc}}} b_{m,k,h}^2(s_h, a_h) / \left(\beta_{k',h}^2 + \lambda\right) \ge \alpha,\tag{8}$$

where k' is the last communication round for m. If so, the agent triggers communication (Line 11).

Part II: Server Updates. After receiving new data, the server merges it with its global datasets  $Z_{k,h}^{ser}$ 

(Line 13) and calculates value function estimates  $\{Q_{k+1,h}\}_{h=1}^{H}$  and  $\{V_{k+1,h}\}_{h=1}^{H}$  using LSVI.

Suppose we already have Q- and V-value function estimates  $Q_{k+1,h+1}$  and  $V_{k+1,h+1}$  at level h + 1. We solve the least squares problem for  $\hat{f}_h$  to minimize the Bellman error (Line 16):

$$\hat{f}_{k+1,h} = \operatorname{argmin}_{f_h \in \mathcal{F}_h} \sum_{o_h \in Z_{k,h}^{\operatorname{ser}}} \left( f_h(z_h) - r_h - V_{k+1,h+1}(s_{h+1}) \right)^2.$$
(9)

We now also define the uncertainty estimator of a new pair of data z = (s, a) against data history Z with normalization parameter  $\lambda$  and function class  $\mathcal{F}$  as

$$D_{\lambda,\mathcal{F}}(z;Z) = \sup_{f_1, f_2 \in \mathcal{F}} |f_1(z) - f_2(z)| / \sqrt{\lambda + \sum_{o' \in Z} |f_1(z') - f_2(z')|^2}.$$
 (10)

Similar to the bandits setting, the uncertainty can be approximated with the bonus function acquired from an oracle  $\mathcal{B}_{S \times A}$  in Definition 4.1. In this case, the domain  $\mathcal{D} = S \times A$ , and the data format corresponds to z = (s, a) and e = (r, s'). Despite these definitions not depending on the step h, we expect the parameters  $z, Z, \mathcal{F}$  to always come from same step h. Finally, we allow the bonus function classes  $\mathcal{W}_h = \mathcal{W}_{h,S \times A}$  to vary between different levels.

After calling oracle for  $b_{k+1,h}$  (Line 16), we can obtain value function estimates (Line 17):

$$Q_{k+1,h}(s,a) = f_{k+1,h}(s,a) + b_{k+1,h}(s,a), \quad V_{k+1,h}(s) = \sup_{a \in \mathcal{A}} Q_{k+1,h}(s,a).$$
 (11)  
Iterating through  $h = H, \dots, 1$ , the server calculates a set of updated Q-value functions  
 $\{Q_{k+1,h}\}_{h=1}^{H}$  and bonus functions  $\{b_{k+1,h}\}_{h=1}^{H}$ , and send them back to agent  $m$  for future ex-  
ploration and updates (lines 19 and 21).

### 307 5.2 Theoretical Results

- <sup>308</sup> We summarize the regret and communication cost of Algorithm 2 in the following theorem:
- **Theorem 5.1.** Taking  $\gamma = O(1/HK)$ ,  $\beta_{h,k} = C_{\beta,2} \left| \sqrt{\lambda} + HC(M,\alpha) \sqrt{\log(3HMN(\gamma)/\delta)} \right|$  and
- 310  $N(\gamma) := \max_h N(\mathcal{F}_h, \gamma) N(\mathcal{F}_{h+1}, \gamma) N(\mathcal{W}_{h+1}, \gamma)$ , the regret within K rounds is bounded by

 $O\left(H\widetilde{\beta}_2\sqrt{(1+M\alpha)\dim_E K}\log(K/\min\{1,\lambda\}) + H^2(1+M\alpha)\dim_E\log^2(K/\min\{1,\lambda\})\right).$ 

where we abbreviate  $\dim_E := \dim_E(\mathcal{F}, \lambda/K)$ ; the total communication complexity is

 $O(H(1+M\alpha)^2\alpha \dim_E(\mathcal{F},\lambda/K)\log^2(K/\min\{1,\lambda\})).$ 

*Remark* 5.2. This result when reduced to linear MDPs correspond well to Theorem 5.1 in Min et al. [2023]. Taking  $\alpha = 1/M^2$ , we get a regret of  $\widetilde{O}(H^2\sqrt{K \dim_E \log N} + H^2 \dim_E)$  and a communication cost of  $\widetilde{O}(HM^2 \dim_E)$ , where  $N = \max_h \{N(\mathcal{F}_h, \gamma), N(\mathcal{W}_h, \gamma)\}$ .

## 314 6 Proof Sketch

In this section, we provide an outline for the proof of Theorem 5.1, while a more detailed proof can be found in Appendix B, and the full versions of the following lemmas are in Appendix B.1.

### 317 6.1 Regret Upper Bound

- <sup>318</sup> For the regret upper bound, the first lemma establishes optimism of value function estimates.
- **Lemma 6.1.** Taking  $\beta_{k,h}$  as in Theorem 5.1, with probability at least  $1 \delta$ , for all  $k, z \in S \times A$  and  $h \in [H], |\mathcal{T}_h Q_{k+1,h+1}(z) - \widehat{f}_{k+1,h}(z)| \leq b_{k+1,h}(z).$
- This allows us to decompose regret into a sum of bonuses:

 $\operatorname{Reg}(K) = \sum_{k=1}^{K} \left[ V_1^*(s_1^k) - V_1^{\pi_{m,k}}(s_1^k) \right]$ 

$$\leq \sum_{k=1}^{K} \sum_{h=1}^{H} \mathbb{E}_{\pi_{m,k}} \left[ Q_{m,k,h} - \mathcal{T}_h Q_{m,k,h+1} \right] (s_h^k, a_h^k) \leq \sum_{k=1}^{K} \sum_{h=1}^{H} 2b_{m,k,h} (s_h^k, a_h^k).$$
(12)

322 The sum of bonuses is equal to the sum of uncertainty up to a constant, which we bound in the

- following lemma corresponding to the elliptical potential lemma [Abbasi-Yadkori et al., 2011].
- Lemma 6.2. Define universal datasets as  $Z_{k,h}^{all} = \{o_h^{k'}\}_{k' \in [k]}$ . Then we have for any  $h \in [H]$ :

$$\sum_{k=1}^{K} D_{\lambda,\mathcal{F}}^2(z_h^k; Z_{k-1,h}^{all}) = O\left(\dim_E(\mathcal{F}, \lambda/K) \log^2(K/\min\{1,\lambda\})\right)$$

- Careful examination exposes a problem: the uncertainty  $D_{\lambda,\mathcal{F}}(z; Z_{k,h}^{ser})$  corresponding to bonuses are
- based on server data  $Z_{k,h}^{\text{ser}}$  instead of universal data  $Z_{k,h}^{\text{all}}$ . The next lemma bridges this gap:

**Lemma 6.3.** For any 
$$z \in \mathcal{S} \times \mathcal{A}$$
,  $k \in [K]$ ,  $h \in [H]$ ,  $D^2_{\lambda,\mathcal{F}}(z; Z^{ser}_{k,h}) \leq (1 + M\alpha) D^2_{\lambda,\mathcal{F}}(z; Z^{all}_{k,h})$ .

With these, we can deduce the regret bound from (12).

## 329 6.2 Communication Cost

For communication cost, we employ an *epoch segmentation* scheme, which defines N epochs segmented by episodes  $\{k_i\}_{i=1}^N$ , with  $k_i$  being the smallest episode satisfying

$$\sum_{o_h \in Z_{k_i,h}^{\text{ser}} \setminus Z_{k_{i-1},h}^{\text{ser}}} \sum_{h=1}^{H} D_{\lambda,\mathcal{F}_h}^2(z_h; Z_{k_{i-1},h}^{\text{ser}}) \ge 1.$$
(13)

This is a generalization of epoch segmentation based on doubling determinants in linear settings, yet the lack of determinant in the nonlinear case dramatically increases its complexity. Intuitively, switch condition (8) suggests an agent must gather a substantial amount of data to trigger communication, yet a careful analysis according to (13) yields a maximum of  $M + C/\alpha$  communication rounds within one epoch. With this we only need an upper bound for the number of epochs N. This is derived by summing (13) over all epochs, then using Lemma 6.1 and Lemma 6.3 to bound the left hand side.

## 338 7 Conclusions

We propose the algorithms Async-NLin-UCB and Async-NLSVI-UCB to tackle multi-agent nonlinear contextual bandits and MDPs with asynchronous communication. We prove that our algorithms enjoy low regret and communication cost, which are comparable to previous results.

iow regret and communication cost, which are comparable to previous results.

Our algorithms employ a communication criterion that allows the agents to trigger communication rounds, effectively controlling communication cost while promoting the asynchronous protocol.

<sup>344</sup> Moreover, we carefully design the contents of server download to guard against data exposure.

## 345 **References**

- Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits.
   *Advances in neural information processing systems*, 24:2312–2320, 2011.
- 348 Alekh Agarwal and Tong Zhang. Model-based RL with optimistic posterior sampling: Structural conditions
- and sample complexity. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors,
   *Advances in Neural Information Processing Systems*, 2022.
- Alekh Agarwal, Yujia Jin, and Tong Zhang. Voql: Towards optimal regret in model-free rl with nonlinear function
   approximation. In Gergely Neu and Lorenzo Rosasco, editors, *Proceedings of Thirty Sixth Conference on Learning Theory*, volume 195 of *Proceedings of Machine Learning Research*, pages 987–1063. PMLR, 12–15
   Jul 2023.
- Ana LC Bazzan. Opportunities for multiagent systems and multiagent reinforcement learning in traffic control.
   *Autonomous Agents and Multi-Agent Systems*, 18(3):342–375, 2009.
- Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Dębiak, Christy Dennison,
   David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, Rafal Józefowicz, Scott Gray, Catherine Olsson,
   Jakub Pachocki, Michael Petrov, Henrique P. d. O. Pinto, Jonathan Raiman, Tim Salimans, Jeremy Schlatter,
   Jonas Schneider, Szymon Sidor, Ilya Sutskever, Jie Tang, Filip Wolski, and Susan Zhang. Dota 2 with large
- scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*, 2019.
- Mithun Chakraborty, Kee Yuan Peh Chua, Sanmay Das, and Brendan Juba. Coordinated versus decentralized exploration in multi-agent multi-armed bandits. In *IJCAI*, 2017.
- Alfredo V Clemente, Humberto N Castejón, and Arjun Chandra. Efficient parallel methods for deep reinforcement learning. *arXiv preprint arXiv:1705.04862*, 2017.
- Qiwei Di, Heyang Zhao, Jiafan He, and Quanquan Gu. Pessimistic nonlinear least-squares value iteration for
   offline reinforcement learning. *arXiv preprint arXiv:2310.01380*, 2023.
- Guohui Ding, Joewie J Koh, Kelly Merckaert, Bram Vanderborght, Marco M Nicotra, Christoffer Heckman,
   Alessandro Roncone, and Lijun Chen. Distributed reinforcement learning for cooperative multi-robot object
   manipulation. In *Proceedings of the 19th International Conference on Autonomous Agents and MultiAgent*
- *Systems*, pages 1831–1833, 2020.
- Simon Du, Sham Kakade, Jason Lee, Shachar Lovett, Gaurav Mahajan, Wen Sun, and Ruosong Wang. Bilinear
   classes: A structural framework for provable generalization in rl. In Marina Meila and Tong Zhang, editors,
   *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of*
- 375 *Machine Learning Research*, pages 2826–2836. PMLR, 18–24 Jul 2021.
- Abhimanyu Dubey and Alex Pentland. Provably efficient cooperative multi-agent reinforcement learning with function approximation. *arXiv preprint arXiv:2103.04972*, 2021.
- Abhimanyu Dubey and AlexSandy' Pentland. Differentially-private federated linear bandits. *Advances in Neural Information Processing Systems*, 33:6003–6014, 2020.
- Lasse Espeholt, Hubert Soyer, Remi Munos, Karen Simonyan, Vlad Mnih, Tom Ward, Yotam Doron, Vlad
   Firoiu, Tim Harley, Iain Dunning, et al. Impala: Scalable distributed deep-rl with importance weighted
   actor-learner architectures. In *International conference on machine learning*, pages 1407–1416. PMLR, 2018.
- Dylan J. Foster, Sham M. Kakade, Jian Qian, and Alexander Rakhlin. The statistical complexity of interactive
   decision making, 2023.
- 285 Claudio Gentile, Zhilei Wang, and Tong Zhang. Fast rates in pool-based batch active learning, 2022.
- Jiafan He, Tianhao Wang, Yifei Min, and Quanquan Gu. A simple and provably efficient algorithm for
   asynchronous federated contextual linear bandits. In *Advances in Neural Information Processing Systems*,
   2022.
- Dan Horgan, John Quan, David Budden, Gabriel Barth-Maron, Matteo Hessel, Hado van Hasselt, and David
   Silver. Distributed prioritized experience replay. In *International Conference on Learning Representations*,
   2018.
- Ruiquan Huang, Weiqiang Wu, Jing Yang, and Cong Shen. Federated linear contextual bandits. In A. Beygelz imer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021.

- Max Jaderberg, Wojciech M Czarnecki, Iain Dunning, Luke Marris, Guy Lever, Antonio Garcia Castaneda,
   Charles Beattie, Neil C Rabinowitz, Ari S Morcos, Avraham Ruderman, et al. Human-level performance in
- 397 3d multiplayer games with population-based reinforcement learning. *Science*, 364(6443):859–865, 2019.

Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, John Langford, and Robert E. Schapire. Contextual decision
 processes with low Bellman rank are PAC-learnable. In Doina Precup and Yee Whye Teh, editors, *Proceedings* of the 34th International Conference on Machine Learning, volume 70 of Proceedings of Machine Learning
 *Research*, pages 1704–1713. PMLR, 06–11 Aug 2017.

- Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement learning with
   linear function approximation. In *Conference on Learning Theory*, pages 2137–2143. PMLR, 2020.
- Chi Jin, Qinghua Liu, and Sobhan Miryoosefi. Bellman eluder dimension: New rich classes of RL problems,
   and sample-efficient algorithms. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors,
   *Advances in Neural Information Processing Systems*, 2021.
- 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*, pages 18–37.
   PMLR, 2022.
- Dingwen Kong, Ruslan Salakhutdinov, Ruosong Wang, and Lin F. Yang. Online sub-sampling for reinforcement
   learning with general function approximation, 2023.
- Nathan Korda, Balázs Szörényi, and Shuai Li. Distributed clustering of linear bandits in peer to peer networks.
   In *Proceedings of the 33rd International Conference on International Conference on Machine Learning -* Volume 48, ICML'16, page 1301–1309. JMLR.org, 2016.
- Jakub Grudzien Kuba, Ruiqing Chen, Muning Wen, Ying Wen, Fanglei Sun, Jun Wang, and Yaodong Yang. Trust
   region policy optimisation in multi-agent reinforcement learning. In *International Conference on Learning Representations*, 2022.
- Peter Landgren, Vaibhav Srivastava, and Naomi Ehrich Leonard. On distributed cooperative decision-making in
   multiarmed bandits. In 2016 European Control Conference (ECC), pages 243–248, 2016.
- Peter Landgren, Vaibhav Srivastava, and Naomi Ehrich Leonard. Social imitation in cooperative multiarmed
   bandits: Partition-based algorithms with strictly local information. In 2018 IEEE Conference on Decision and
   Control (CDC), 2018.
- Chuanhao Li and Hongning Wang. Asynchronous upper confidence bound algorithms for federated linear
   bandits. In *International Conference on Artificial Intelligence and Statistics*, pages 6529–6553. PMLR, 2022.
- Gene Li, Pritish Kamath, Dylan J Foster, and Nati Srebro. Understanding the eluder dimension. In S. Koyejo,
   S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 23737–23750. Curran Associates, Inc., 2022.
- Boyi Liu, Lujia Wang, and Ming Liu. Lifelong federated reinforcement learning: a learning architecture for navigation in cloud robotic systems. *IEEE Robotics and Automation Letters*, 4(4):4555–4562, 2019.
- Dianbo Liu, Vedant Shah, Oussama Boussif, Cristian Meo, Anirudh Goyal, Tianmin Shu, Michael Mozer,
   Nicolas Heess, and Yoshua Bengio. Stateful active facilitator: Coordination and environmental heterogeneity
   in cooperative multi-agent reinforcement learning. *arXiv preprint arXiv:2210.03022*, 2022.
- Dongfang Liu, Yiming Cui, Zhiwen Cao, and Yingjie Chen. Indoor navigation for mobile agents: A multimodal
   vision fusion model. In 2020 International Joint Conference on Neural Networks (IJCNN), pages 1–8. IEEE,
   2020.
- Keqin Liu and Qing Zhao. Distributed learning in multi-armed bandit with multiple players. *IEEE Transactions* on Signal Processing, 58(11):5667–5681, 2010. doi: 10.1109/TSP.2010.2062509.
- Ryan Lowe, Yi I Wu, Aviv Tamar, Jean Harb, OpenAI Pieter Abbeel, and Igor Mordatch. Multi-agent actor-critic
   for mixed cooperative-competitive environments. *Advances in neural information processing systems*, 30, 2017.
- David Martínez-Rubio, Varun Kanade, and Patrick Rebeschini. Decentralized cooperative stochastic bandits.
   In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*. Curran Associates, Inc., 2019.
- Yifei Min, Tianhao Wang, Ruitu Xu, Zhaoran Wang, Michael Jordan, and Zhuoran Yang. Learn to match with no
   regret: Reinforcement learning in markov matching markets. In *Advances in Neural Information Processing*

- 447 Yifei Min, Jiafan He, Tianhao Wang, and Quanquan Gu. Cooperative multi-agent reinforcement learning:
   448 Asynchronous communication and linear function approximation. In Andreas Krause, Emma Brunskill,
- Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett, editors, *Proceedings of the 40th*
- International Conference on Machine Learning, volume 202 of Proceedings of Machine Learning Research,

451 pages 24785–24811. PMLR, 23–29 Jul 2023.

- 452 Seongin Na, Tomáš Rouček, Jiří Ulrich, Jan Pikman, Tomáš Krajník, Barry Lennox, and Farshad Arvin.
   453 Federated reinforcement learning for collective navigation of robotic swarms, 2022.
- Arun Nair, Praveen Srinivasan, Sam Blackwell, Cagdas Alcicek, Rory Fearon, Alessandro De Maria, Vedavyas
   Panneershelvam, Mustafa Suleyman, Charles Beattie, Stig Petersen, et al. Massively parallel methods for
   deep reinforcement learning. *arXiv preprint arXiv:1507.04296*, 2015.
- Jiaju Qi, Qihao Zhou, Lei Lei, and Kan Zheng. Federated reinforcement learning: techniques, applications, and
   open challenges. *arXiv preprint arXiv:2108.11887*, 2021.
- Daniel Russo and Benjamin Van Roy. Eluder dimension and the sample complexity of optimistic exploration.
   In C.J. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc., 2013.
- Abishek Sankararaman, Ayalvadi Ganesh, and Sanjay Shakkottai. Social learning in multi agent multi armed
   bandits. *Proc. ACM Meas. Anal. Comput. Syst.*, 3(3), 2019.
- 464 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction.* MIT press, 2018.
- Balazs Szorenyi, Robert Busa-Fekete, Istvan Hegedus, Robert Ormandi, Mark Jelasity, and Balazs Kegl. Gossip based distributed stochastic bandit algorithms. In *Proceedings of the 30th International Conference on Machine Learning*, 2013.
- Oriol Vinyals, Timo Ewalds, Sergey Bartunov, Petko Georgiev, Alexander Sasha Vezhnevets, Michelle Yeo,
   Alireza Makhzani, Heinrich Küttler, John Agapiou, Julian Schrittwieser, et al. Starcraft ii: A new challenge
   for reinforcement learning. *arXiv preprint arXiv:1708.04782*, 2017.
- Hoi-To Wai, Zhuoran Yang, Zhaoran Wang, and Mingyi Hong. Multi-agent reinforcement learning via double
   averaging primal-dual optimization. *Advances in Neural Information Processing Systems*, 31, 2018.
- 473 Martin Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. 02 2019. ISBN 9781108498029.
- Po-An Wang, Alexandre Proutiere, Kaito Ariu, Yassir Jedra, and Alessio Russo. Optimal algorithms for
   multiplayer multi-armed bandits. In Silvia Chiappa and Roberto Calandra, editors, *Proceedings of the Twenty*
- Third International Conference on Artificial Intelligence and Statistics, volume 108 of Proceedings of Machine
   Learning Research, pages 4120–4129. PMLR, 26–28 Aug 2020a.
- Ruosong Wang, Russ R Salakhutdinov, and Lin Yang. Reinforcement learning with general value function
  approximation: Provably efficient approach via bounded eluder dimension. In H. Larochelle, M. Ranzato,
  R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33,
- 481 pages 6123–6135. Curran Associates, Inc., 2020b.
- Yuanhao Wang, Jiachen Hu, Xiaoyu Chen, and Liwei Wang. Distributed bandit learning: Near-optimal regret
   with efficient communication. In *International Conference on Learning Representations*, 2020c.
- 484 Grady Williams, Paul Drews, Brian Goldfain, James M Rehg, and Evangelos A Theodorou. Aggressive
   driving with model predictive path integral control. In 2016 IEEE International Conference on Robotics and
   486 Automation (ICRA), pages 1433–1440. IEEE, 2016.
- Ruitu Xu, Yifei Min, Tianhao Wang, Michael I Jordan, Zhaoran Wang, and Zhuoran Yang. Finding regular ized competitive equilibria of heterogeneous agent macroeconomic models via reinforcement learning. In
   *International Conference on Artificial Intelligence and Statistics*, pages 375–407. PMLR, 2023.
- Chenlu Ye, Wei Xiong, Quanquan Gu, and Tong Zhang. Corruption-robust algorithms with uncertainty
   weighting for nonlinear contextual bandits and Markov decision processes. In Andreas Krause, Emma
   Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 39834–39863. PMLR, 23–29 Jul 2023.
- Deheng Ye, Guibin Chen, Wen Zhang, Sheng Chen, Bo Yuan, Bo Liu, Jia Chen, Zhao Liu, Fuhao Qiu, Hongsheng
   Yu, et al. Towards playing full moba games with deep reinforcement learning. *Advances in Neural Information Processing Systems*, 33:621–632, 2020.

- 498 Chao Yu, Akash Velu, Eugene Vinitsky, Yu Wang, Alexandre Bayen, and Yi Wu. The surprising effectiveness of 499 ppo in cooperative, multi-agent games. *arXiv preprint arXiv:2103.01955*, 2021.
- 500 Shuai Yu, Xu Chen, Zhi Zhou, Xiaowen Gong, and Di Wu. When deep reinforcement learning meets federated
- learning: Intelligent multitimescale resource management for multiaccess edge computing in 5g ultradense
   network. *IEEE Internet of Things Journal*, 8(4):2238–2251, 2020.
- Tao Yu, HZ Wang, Bin Zhou, Ka Wing Chan, and J Tang. Multi-agent correlated equilibrium q ( $\lambda$ ) learning for coordinated smart generation control of interconnected power grids. *IEEE transactions on power systems*, 30 (4):1669–1679, 2014.
- Kaiqing Zhang, Zhuoran Yang, and Tamer Basar. Networked multi-agent reinforcement learning in continuous
   spaces. In 2018 IEEE conference on decision and control (CDC), pages 2771–2776. IEEE, 2018a.
- Kaiqing Zhang, Zhuoran Yang, Han Liu, Tong Zhang, and Tamer Basar. Fully decentralized multi-agent
   reinforcement learning with networked agents. In *International Conference on Machine Learning*, pages
   5872–5881. PMLR, 2018b.
- Heyang Zhao, Jiafan He, and Quanquan Gu. A nearly optimal and low-switching algorithm for reinforcement
   learning with general function approximation, 2023.
- Zhaowei Zhu, Jingxuan Zhu, Ji Liu, and Yang Liu. Federated bandit: A gossiping approach. *Proc. ACM Meas. Anal. Comput. Syst.*, 5(1), 2021.
- Hankz Hankui Zhuo, Wenfeng Feng, Yufeng Lin, Qian Xu, and Qiang Yang. Federated deep reinforcement
   learning. *arXiv preprint arXiv:1901.08277*, 2019.

#### **Impact Statement** 517

Our work has the potential to enhance cooperative learning systems across diverse fields. By 518

519 introducing algorithms that enable efficient collaboration among agents with minimal communication

overhead, our research paves the way for advancements in distributed systems, including robotics, 520

traffic management, and distributed sensor networks. This could lead to more adaptive, efficient, 521

and scalable systems capable of tackling complex problems in dynamic environments, ultimately 522 contributing to technological progress and societal well-being.

523

As far as we can tell, there is hardly any negative social impact from our work, mainly because we do 524

525 not include experiments apart from our theoretical analysis.

## 526 A The Bandit Case: Proof of Theorem 4.3

Before we begin the analysis of Algorithm 1, we reiterate and add some notations for clarity and convenience. Define the data collected by agent m that has already been uploaded to the server by round t as  $Z_{m,t}^{up}$ , and the universal data at round t as  $Z_t^{all}$ . Apart from these we also have from the algorithm the datasets  $Z_{m,t}^{loc}$  and  $Z_t^{ser}$ . It is not difficult to check that they satisfy the following relation:

$$Z_t^{\text{all}} = \bigcup_{m=1}^M \left( Z_{m,t}^{\text{up}} \cup Z_{m,t}^{\text{loc}} \right).$$

Furthermore, when t is not a communication round, we also have

$$Z_t^{\rm ser} = \bigcup_{m=1}^M Z_{m,t}^{\rm up}$$

and when it is a communication round that

$$Z^{\mathrm{ser}}_t = \left[\bigcup_{m=1}^M Z^{\mathrm{up}}_{m,t}\right] \cup Z^{\mathrm{loc}}_{m_t,t},$$

<sup>527</sup> which will be useful in our proof of Lemma A.1 and B.1 in Section C.1.

Next, we assume that at rounds  $0 = t_0 < t_1 < \cdots < t_L < t_{L+1} = T + 1$ , the participating agent communicates with the server, where  $t_0$  and  $t_{K+1}$  are dummy rounds. The subscripts will be denoted as  $l = 1, \cdots, L$  in the future.

We now describe a participant reordering trick for our asynchronous multi-agent setting, which we will 531 use multiple times in the proof. The basic idea is that, as long as the *communication order* remains the 532 same, and for any given agent, the *number of rounds* between two consecutive communication rounds 533 remains the same, one can switch the episodes around and change the order of agent participation to 534 a certain degree. For example, we may assume that  $m_t = m_{t_l}$  for all  $t \in (t_{l-1}, t_l]$  by reordering the 535 participants, which means all participation of any given agent happens immediately before a certain 536 communication round; as another example, we may assume  $m_t = m_{t_{l-1}}$  for all  $t \in [t_{l-1}, t_l)$ , which 537 means all participation happen immediately after communication rounds. It should be noted that one 538 needs to be careful when utilizing this argument, since switching the participation order changes the 539 values of  $t_l$  and many associated elements, so applying this trick twice in succession would lead to 540 contradictions. 541

For a dataset Z, we define the Z-norm on function set  $\mathcal{F}$  as  $||f||_Z^2 := \sum_{(a,r)\in Z} f^2(a)$  for any  $f \in \mathcal{F}$ . Then we have the shortened notation

$$D_{\lambda,\mathcal{F}}(a;Z) = \sup_{f_1,f_2 \in \mathcal{F}} \frac{|f_1(a) - f_2(a)|}{\sqrt{\lambda + \|f_1 - f_2\|_Z^2}}$$

Finally, we define the confidence set of functions at round t + 1 as:

$$\mathcal{F}_{t+1} = \left\{ f \in \mathcal{F} : \sum_{(a,r)\in Z_t^{\text{ser}}} \left( f(a) - \widehat{f}_{t+1}(a) \right)^2 \le \beta_t^2 \right\},\tag{14}$$

<sup>543</sup> which is a common construction in reinforcement learning.

#### 544 A.1 Auxiliary Lemmas

<sup>545</sup> In this section we present some auxiliary lemmas that will be used in the proof of Theorem 4.3. Note these lemmas correspond well to the lemmas presented in 6, only that these are for the contextual

547 bandit case. The proofs for these lemmas can be found in Section C.

**Lemma A.1.** For any  $t \in [T]$ ,  $m \in [M]$  and  $f_1, f_2 \in \mathcal{F}$ , as long as agent m does not communicate with the server at time step t, we have

$$\lambda + \sum_{m' \in [M]} \|f_1 - f_2\|_{Z_{m',t}^{up}}^2 \ge \frac{1}{\alpha} \|f_1 - f_2\|_{Z_{m,t}^{loc}}^2.$$

Furthermore, for any  $t \in [T]$  and  $f_1, f_2 \in \mathcal{F}$ ,

$$\lambda + \|f_1 - f_2\|_{Z_t^{ser}}^2 \ge \frac{1}{1 + M\alpha} \left(\lambda + \|f_1 - f_2\|_{Z_t^{all}}^2\right),$$

and as a corollary, for any  $a \in A$ ,

$$D^2_{\lambda,\mathcal{F}}(a; Z_t^{ser}) \leq (1 + M\alpha) D^2_{\lambda,\mathcal{F}}(a; Z_t^{all})$$

This lemma describes the discrepancy between different datasets. Crucially, it provides a worst case ratio between uncertainty measured on the server dataset and universal dataset. This is an important tool for bridging between the different uncertainty estimators in the following proofs. The proof can be found in Section C.1.

**Lemma A.2.** By taking  $\gamma = O(1/T)$  and

$$\beta_t = \widetilde{\beta}_1 := C_{\beta,1} \left[ \sqrt{\lambda} + \sqrt{(\gamma^2 + \gamma R)T} + RC(M, \alpha) \log(3MN(\mathcal{F}, \gamma)/\delta) \right],$$

with  $C_{\beta,1} = 6$ , where  $C(M, \alpha) := \sqrt{1 + M\alpha} + M\sqrt{\alpha}$ , we have  $f^* \in \mathcal{F}_{t+1}$  for all  $t \in \{t_l\}_{l=1}^L$  with probability at least  $1 - \delta$ . As a corollary, we also have  $|f_*(a) - \hat{f}_{t+1}(a)| \le b_{t+1}(a)$  for any  $a \in \mathcal{A}_t$ and  $t \in \{t_l\}_{l=1}^L$ .

555 This is the central optimism lemma present in all provably efficient reinforcement learning literature.

It states that the confidence function set contains the ground truth function  $f^*$  with high probability, and in our case, that the decision function  $\hat{f}_t + b_t$  is optimistic. With this, we define the good event  $\mathcal{E}_T = \{f^* \in \mathcal{F}_{t+1}, \forall t \in \{t_l\}_{l=1}^L\}$ . Then according to A.2,  $\mathbb{P}(\mathcal{E}_T) \ge 1 - \delta$ . The proof can be found in Section C.2.

**Lemma A.3.** The sum of squared uncertainty estimators of new data over all historical data can be bounded as follows with some absolute constant  $C_D$ :

$$\sum_{t=1}^{T} D_{\lambda,\mathcal{F}}^2(a_t; Z_{t-1}^{all}) \le C_D \dim_E(\mathcal{F}, \lambda/T) \log^2(T/\min\{1, \lambda\})$$

This lemma corresponds to the elliptical potential argument from the linear setting [Abbasi-Yadkori et al., 2011]. In the nonlinear setting, this lemma essentially reveals the relationship between the sum of Eluder-like confidence quantities and the Eluder dimension. The proof can be found in Section C.3.

#### 564 A.2 The Epoch Segmentation Scheme

In this section, we introduce an epoch segmentation scheme, which is needed for both the regret and communication cost proofs presented in the next two sections. It is a generalization of the epoch segmentation scheme based on doubling determinant in the linear bandits / MDPs setting [He et al., 2022, Min et al., 2023], but the lack of a Gram matrix (used for linear regression) in the nonlinear

case complicates matters significantly.

We segment the entire run of  $t = 1, \dots, T$  into N epochs as follows. Define iteratively  $0 = l_0 < l_1 < \dots < l_N \leq L$  as

$$l_{i} = \min\left\{l > l_{i-1} : \sum_{l'=l_{i-1}+1}^{l} \sum_{(a,r) \in Z_{m,t_{l'}}^{\text{loc}}} D_{\lambda,\mathcal{F}}^{2}(a; Z_{t_{l_{i-1}}}^{\text{ser}}) \ge 1\right\},\$$

where for a given l' in the summation,  $m = m_{t_{l'}}$  is the participating agent at  $t_{l'}$ . In the iterative process, if the above minimum does not exist, simply define N = i - 1 and end the process there.

<sup>572</sup> Correspondingly, the *i*-th epoch is defined by the time steps  $[t_{l_{i-1}}, t_{l_i})$ .

The full is the full point is defined by the time steps  $[v_{l_{i-1}}, v_{l_i}]$ .

The following sections will make use of this epoch scheme as befit their needs, but here we shall give an upper bound for the total number of epochs N. Based on the definition of  $l_i$ , we have for any 575  $l_{i-1} \le l < l_i$  that

$$1 \geq \sum_{l'=l_{i-1}+1}^{l} \sum_{(a,r)\in Z_{m_{t_{l'}},t_{l'}}} D_{\lambda,\mathcal{F}}^2(a; Z_{t_{l_{i-1}}}^{\operatorname{ser}})$$

$$= \sum_{l'=l_{i-1}+1}^{l} \sum_{(a,r)\in Z_{t_{l'}}^{\operatorname{ser}} \setminus Z_{t_{l'-1}}^{\operatorname{ser}}} \sup_{f_{1},f_{2}\in\mathcal{F}} \frac{[f_{1}(a) - f_{2}(a)]^2}{\lambda + \|f_{1} - f_{2}\|_{Z_{t_{l-1}}}^{2\operatorname{ser}}}$$

$$\geq \sup_{f_{1},f_{2}\in\mathcal{F}} \frac{\sum_{(a,r)\in Z_{t_{l}}^{\operatorname{ser}} \setminus Z_{t_{l_{i-1}}}^{\operatorname{ser}}}{1 + \|f_{1} - f_{2}\|_{Z_{t_{i-1}}}^{2\operatorname{ser}}}}{\lambda + \|f_{1} - f_{2}\|_{Z_{t_{i-1}}}^{2\operatorname{ser}}} - 1,$$

576 which gives  $\lambda + \|f_1 - f_2\|_{Z_{t_l}^{\text{ser}}}^2 \le 2\left(\lambda + \|f_1 - f_2\|_{Z_{t_{l_{i-1}}}^{\text{ser}}}^2\right)$  for any  $f_1, f_2 \in \mathcal{F}$ . Then we have  $D_{\lambda,\mathcal{F}}^2(a; Z_{t_{l_{i-1}}}^{\text{ser}}) \le 2D_{\lambda,\mathcal{F}}^2(a; Z_{t_l}^{\text{ser}})$  (15)

577 for any a, and so

$$\begin{split} 1 &\leq \sum_{(a,r)\in Z_{t_{l_{i}}}^{\mathrm{ser}}\setminus Z_{t_{l_{i-1}}}^{\mathrm{ser}}} D_{\lambda,\mathcal{F}}^{2}(a;Z_{t_{l_{i-1}}}^{\mathrm{ser}}) \\ &= \sum_{l=l_{i-1}+1}^{l_{i}} \sum_{(a,r)\in Z_{t_{l}}^{\mathrm{ser}}\setminus Z_{t_{l-1}}^{\mathrm{ser}}} D_{\lambda,\mathcal{F}}^{2}(a;Z_{t_{l_{i-1}}}^{\mathrm{ser}}) \\ &\leq 2\sum_{l=l_{i-1}+1}^{l_{i}} \sum_{(a,r)\in Z_{t_{l}}^{\mathrm{ser}}\setminus Z_{t_{l-1}}^{\mathrm{ser}}} D_{\lambda,\mathcal{F}}^{2}(a;Z_{t_{l-1}}^{\mathrm{ser}}), \end{split}$$

and summing over  $i = 1, \cdots, N-1$  that:

$$N-1 \le 2\sum_{l=1}^{L} \sum_{(a,r)\in Z_{t_l}^{\text{ser}} \setminus Z_{t_{l-1}}^{\text{ser}}} D_{\lambda,\mathcal{F}}^2(a; Z_{t_{l-1}}^{\text{ser}}).$$

If we apply the participant reordering trick and let  $m_t = m_{t_l}$  for all  $t \in (t_{l-1}, t_l]$  and  $l \in [L]$ , we get  $Z_{t_l}^{\text{ser}} \setminus Z_{t_{l-1}}^{\text{ser}} = \{(a_t, r_t)\}_{t=t_{l-1}+1}^{t_l}$ , and so applying Lemma A.1 and Lemma A.3, we get

$$\begin{split} N-1 &\leq 2\sum_{l=1}^{L}\sum_{t=t_{l-1}+1}^{t_l} D_{\lambda,\mathcal{F}}^2(a_t;Z_{t_{l-1}}^{\mathrm{ser}}) \\ &\leq 2(1+M\alpha)\sum_{l=1}^{L}\sum_{t=t_{l-1}+1}^{t_l} D_{\lambda,\mathcal{F}}^2(a_t;Z_{t-1}^{\mathrm{all}}) \\ &\leq 2(1+M\alpha)\sum_{t=1}^{T} D_{\lambda,\mathcal{F}}^2(a_t;Z_{t-1}^{\mathrm{all}}) \\ &\leq C(1+M\alpha)\dim_E(\mathcal{F},\lambda/T)\log(T/\lambda)\log T, \end{split}$$

<sup>581</sup> which gives the order of total number of epochs:

$$N = O\bigg((1 + M\alpha) \dim_E(\mathcal{F}, \lambda/T) \log^2(T/\min\{1, \lambda\})\bigg).$$
(16)

Notice that the participant reordering trick is only used to bound the number of epochs, which itself does not depend on the specific order of participation. This is crucial since it suggests this reordering does not change anything essential, and is in fact not necessary for the proof - it just made the proof easier to read. Therefore we can still reorder participants as we see fit in other parts of our proof.

## 586 A.3 Proof of Regret Upper Bound

Now we are ready to prove the first part of Theorem 4.3 concerning the regret upper bound. We begin by applying the participation reordering trick to assume, without loss of generality, that the same agent is active within the rounds  $[t_l, t_{l+1} - 1]$ , i.e.  $m_{t_l} = m_{t_l+1} = \cdots = m_{t_{l+1}-1}$ . Under this assumption, we have  $t_1 = 1$ .

Let  $a_t^* := \operatorname{argmax}_{a \in D_t} f_*(a)$  be the best arm at time t. Then by Lemma A.2,  $f_*(a_t^*) \leq (\widehat{f}_{m_t,t} + b_{m_t,t})(a_t^*) \leq (\widehat{f}_{m_t,t} + b_{m_t,t})(a_t)$ , where the second inequality is due to the choice of  $a_t$  at round t. Hence we get

$$\operatorname{Reg}(T) = \sum_{t=1}^{T} \left[ f_*(a_t^*) - f_*(a_t) \right]$$
  

$$\leq \min\left\{ \sum_{t=1}^{T} \left( \widehat{f}_{m_t,t} + b_{m_t,t} - f^* \right)(a_t), 4 \right\}$$
  

$$\leq \sum_{t=1}^{T} \min\{2b_{m_t,t}(a_t), 4\}$$
  

$$= 2\sum_{l=1}^{L} \sum_{t=t_l+1}^{t_{l+1}-1} b_{t_l}(a_t) + 2\sum_{l=1}^{L} \min\{b_{m_{t_l},t_l}(a_{t_l}), 2\},$$
(17)

- where the first inequality is due to  $|f| \le 1$  from Assumption 3.1, and the second inequality again
- <sup>595</sup> uses Lemma A.2. We first bound the second term here using the epoch scheme in Section A.2. We <sup>596</sup> start by converting the bonus term to uncertainty:

$$b_{m_{t_l},t_l}(a_{t_l}) = b_{t_{l-1}}(a_{t_l}) \leq C_{\mathcal{B}} \sqrt{\beta_{t_l-1}^2 + \lambda} \cdot D_{\lambda,\mathcal{F}}(a_{t_l}; Z_{t_{l-1}}^{\mathrm{ser}}).$$
(18)

Now consider the episodes in an epoch *i*, specifically  $\{t_{l_{i-1}}, t_{l_{i-1}+1}, \cdots, t_{l_i}\}$ . For any  $l_{i-1} < l < l_i$ , since  $Z_{t_{l_{i-1}}}^{\text{ser}} \subseteq Z_{t_{l-1}}^{\text{ser}}$ , we can deduce that

$$D^{2}_{\lambda,\mathcal{F}}(z_{t_{l}}; Z^{\text{ser}}_{t_{l-1}}) \leq D^{2}_{\lambda,\mathcal{F}}(z_{t_{l}}; Z^{\text{ser}}_{t_{l_{i-1}}}) \leq 2D^{2}_{\lambda,\mathcal{F}}(z_{t_{l}}; Z^{\text{ser}}_{t_{l}}),$$

where the second inequality is borrowed from (15) from Section A.2. Therefore continuing from (18),

$$\sum_{l=1}^{L} \min\{b_{m_{t_l},t_l}(z_{t_l}), 2\} \leq \sum_{l \notin \{l_i\}_{i=1}^{N}} \left[ \sqrt{2}C_{\mathcal{B}} \sqrt{\beta_{t_l-1}^2 + \lambda} \cdot D_{\lambda,\mathcal{F}}(z_{t_l}; Z_{t_l}^{\text{ser}}) \right] + \sum_{i=1}^{N} 2$$
$$\leq \sqrt{2}C_{\mathcal{B}} \sum_{l=1}^{L} D_{\lambda,\mathcal{F}}(z_{t_l}; Z_{t_l}^{\text{ser}}) \sqrt{\beta_h^2 + \lambda} + 2N.$$
(19)

Now combine this result with the first term in (17) and use again (18), we get

$$\operatorname{Reg}(T) \leq 2C_{\mathcal{B}} \sum_{l=1}^{L} \sum_{t=t_{l}+1}^{t_{l+1}-1} D_{\lambda,\mathcal{F}}(a_{t}; Z_{t_{l}}^{\operatorname{ser}}) \sqrt{\beta_{t_{l}}^{2} + \lambda} + 2\sqrt{2}C_{\mathcal{B}} \sum_{l=1}^{L} D_{\lambda,\mathcal{F}}(z_{t_{l}}; Z_{t_{l}}^{\operatorname{ser}}) \sqrt{\beta_{h}^{2} + \lambda} + 4N$$
  
$$\leq 2\sqrt{2}C_{\mathcal{B}} \sum_{l=1}^{L} \sum_{t=t_{l}}^{t_{l+1}-1} D_{\lambda,\mathcal{F}}(a_{t}; Z_{t_{l}}^{\operatorname{ser}}) \sqrt{\beta_{t_{l}}^{2} + \lambda} + 4N$$
  
$$\leq 2\sqrt{2}C_{\mathcal{B}} \left[ \sum_{l=1}^{L} \sum_{t=t_{l}}^{t_{l+1}-1} D_{\lambda,\mathcal{F}}^{2}(a_{t}; Z_{t_{l}}^{\operatorname{ser}}) \right]^{1/2} \left[ \sum_{t=1}^{T} \left( \widetilde{\beta}_{1}^{2} + \lambda \right) \right]^{1/2} + 4N$$

where

$$\widetilde{\beta}_1 = C_\beta \bigg[ \sqrt{\lambda} + RC(M, \alpha) \sqrt{\log(3N(M/\delta))} \bigg].$$

602 According to Lemma A.1 and Lemma A.3, the term

$$\sum_{l=1}^{L} \sum_{t=t_{l}}^{t_{l+1}-1} D_{\lambda,\mathcal{F}}^{2}(a_{t}; Z_{t_{l}}^{\text{ser}}) \leq (1+M\alpha) \sum_{l=1}^{L} \sum_{t=t_{l}}^{t_{l+1}-1} D_{\lambda,\mathcal{F}}^{2}(a_{t}; Z_{t-1}^{\text{all}})$$
$$= (1+M\alpha) \sum_{t=1}^{T} D_{\lambda,\mathcal{F}}^{2}(a_{t}; Z_{t-1}^{\text{all}})$$
$$\leq C(1+M\alpha) \dim_{E}(\mathcal{F}, \lambda/T) \log^{2} \left(T/\min\{1,\lambda\}\right).$$

combining this with (16), we get

$$\begin{split} \operatorname{Reg}(T) &\leq C \big[ (1 + M\alpha) \dim_E(\mathcal{F}, \lambda/T) \log^2 \big( T/\min\{1, \lambda\} \big) \big]^{1/2} \bigg[ \sum_{t=1}^T (\beta_t^2 + \lambda) \bigg]^{1/2} + 4N \\ &= O \bigg( \sqrt{T} \widetilde{\beta}_1 \sqrt{(1 + M\alpha) \dim_E(\mathcal{F}, \lambda/T)} \log(T/\min\{1, \lambda\}) \\ &+ (1 + M\alpha) \dim_E(\mathcal{F}, \lambda/T) \log^2(T/\min\{1, \lambda\}) \bigg). \end{split}$$

## 604 A.4 Proof of Communication Cost

In this section we prove the second part of Theorem 4.3, by calculating the communication complexity. First, for each communication round  $t_l$ , assume the last time before  $t_l$  when the agent  $m_{t_l}$ communicated with the server was  $t_{l'}$ , then

$$\sum_{(a,r)\in Z_{m,t_l}^{\mathrm{loc}}} D^2_{\lambda,\mathcal{F}}(a; Z_{m,t_l}^{\mathrm{up}}) \geq \sum_{(a,r)\in Z_{m,t_l}^{\mathrm{loc}}} \frac{\left[b_{t_{l'}}(a)/C\right]^2}{\beta_{t_{l'}}^2 + \lambda} \geq \frac{\alpha}{C^2},$$

Now employing the epoch segmentation scheme from section A.2, for the *i*-th epoch consisting of the time steps  $[t_{l_{i-1}}, t_{l_i})$ , we have the inequality

$$\begin{split} 1 &\geq \sum_{l=l_{i-1}+1}^{l_{i}-1} \sum_{(a,r)\in Z_{m,t_{l}}^{\mathrm{loc}}} D_{\lambda,\mathcal{F}}^{2}(a;Z_{t_{l_{i-1}}}^{\mathrm{ser}}) \\ &\geq \sum_{l=l_{i-1}+1}^{l_{i}-1} \sum_{(a,r)\in Z_{m,t_{l}}^{\mathrm{loc}}} D_{\lambda,\mathcal{F}}^{2}\left(a;Z_{m,t_{l}}^{\mathrm{up}}\cup Z_{t_{l_{i-1}}}^{\mathrm{ser}}\right). \end{split}$$

For  $m \in [M]$ , assume the agent m communicated with the server a total of  $n_m$  times within  $[t_{l_{i-1}}, t_{l_i})$ . Then except for the first of these communication rounds, for each  $l \in [l_{i-1}+1, l_i-1]$  with  $m_{t_l} = m$ , there exists  $l' \in [l_{i-1}, l)$  with  $m_{t_{l'}} = m$ , thus we have  $Z_{m,t_l}^{up} \supset Z_{m,t_{l'}+1}^{up} = Z_{t_{l'}}^{ser} \supset Z_{t_{l_{i-1}}}^{ser}$ . With this we have the corresponding term

$$\sum_{(a,r)\in Z_{m,t_l}^{\mathrm{loc}}} D^2_{\lambda,\mathcal{F}}(a; Z_{m,t_l}^{\mathrm{up}} \cup Z_{t_{l_{i-1}}}^{\mathrm{ser}}) = \sum_{(a,r)\in Z_{m,t_l}^{\mathrm{loc}}} D^2_{\lambda,\mathcal{F}}(a; Z_{m,t_l}^{\mathrm{up}}) \ge \frac{\alpha}{C_{\mathcal{B}}^2},$$

610 therefore

$$1 \ge \sum_{m=1}^{M} (n_m - 1) \cdot \frac{\alpha}{4C^2} \Rightarrow \sum_{m=1}^{M} n_m \le M + \frac{C_{\mathcal{B}}^2}{\alpha}$$

Notice that  $\sum_{m=1}^{M} n_m = l_i - l_{i-1}$  is the number of communication rounds within  $[t_{l_{i-1}}, t_{l_i})$ , hence summing over *i* the total number of communication rounds is upper bounded by  $N(M + C_{\mathcal{B}}^2/\alpha)$ . Combine this with (16), we have the total number of communication rounds throughout the algorithm is

$$O\left(\frac{(1+M\alpha)^2}{\alpha}\dim_E(\mathcal{F},\lambda/T)\log^2\left(T/\min\{1,\lambda\}\right)\right).$$

#### B The MDPs Case: Proof of Theorem 5.1 611

Similar to the bandit case, we define  $Z_{m,k,h}^{\text{loc}}$ ,  $Z_{m,k,h}^{\text{up}}$ ,  $Z_{k,h}^{\text{ser}}$ , and  $Z_{k,h}^{\text{all}}$  to be the local, uploaded, server and universal data, with corresponding subscripts of agent  $m \in [M]$ , episode  $k \in [K]$ ,  $h \in [H]$ . 612

613

Suppose at rounds  $0 = k_0 < k_1 < \cdots < k_L < k_{L+1} = T + 1$ , the participating agent communicates 614 with the server, where  $k_0$  and  $k_{L+1}$  are dummy rounds. 615

For a dataset  $Z_h$  in the MDPs setting, we again define the  $Z_h$ -norm on function set  $\mathcal{F}_h$  as  $||f||_Z^2 := \sum_{o_h \in Z} f^2(z_h)$  for any  $f \in \mathcal{F}$ . As a reminder, the tuples  $o_h = (s_h, a_h, r_h, s_{h+1})$  and  $z_h = (s_h, a_h)$ . Then we have the shortened notation

$$D_{\lambda,\mathcal{F}_h}(z_h; Z_h) = \sup_{f_1, f_2 \in \mathcal{F}_h} \frac{|f_1(z_h) - f_2(z_h)|}{\sqrt{\lambda + \|f_1 - f_2\|_{Z_h}^2}}$$

Finally, we define the confidence set of functions at round k + 1 and step h as: 616

$$\mathcal{F}_{k+1,h} = \left\{ f \in \mathcal{F}_h : \sum_{o_h \in \mathbb{Z}_{k,h}^{\mathrm{ser}}} \left( f(z_h) - \widehat{f}_{k+1,h}(z_h) \right)^2 \le (\beta_{k,h})^2 \right\}.$$
 (20)

#### **B.1** Auxiliary Lemmas 617

In this section we present some auxiliary lemmas that will be used in the proof of Theorem 5.1. These 618

lemmas are generalizations / restatements to the lemmas presented in 6, and their detailed proofs can 619 620 be found in Section C.

**Lemma B.1** (Restatement of Lemma 6.3). For any  $k \in [K]$ ,  $m \in [M]$ ,  $h \in [H]$  and  $f_1, f_2 \in \mathcal{F}$ , as long as agent m does not communicate with the server at episode k, we have

$$\lambda + \sum_{m' \in [M]} \|f_1 - f_2\|_{Z_{m',k,h}^{up}}^2 \ge \frac{1}{\alpha} \|f_1 - f_2\|_{Z_{m,k,h}^{loc}}^2.$$

Furthermore, we have for any  $k \in [K]$  and  $f_1, f_2 \in \mathcal{F}$ ,

$$\lambda + \|f_1 - f_2\|_{Z_{k,h}^{\text{ser}}}^2 \ge \frac{1}{1 + M\alpha} \left(\lambda + \|f_1 - f_2\|_{Z_{k,h}^{\text{all}}}^2\right),$$

and as a corollary, for any  $z = (s, a) \in \mathcal{S} \times \mathcal{A}$ ,

$$D^2_{\lambda,\mathcal{F}}(z; Z^{ser}_{k,h}) \le (1 + M\alpha) D^2_{\lambda,\mathcal{F}}(z; Z^{all}_{k,h})$$

- Similar to Lemma A.1, this lemma provides a worst case ratio between uncertainty measured on the 621 server dataset and universal dataset. The proof can be found in Section C.1. 622
- **Lemma B.2** (Restatement of Lemma 6.1). By taking  $\gamma = 1/(C_{\gamma}KH)$  with  $C_{\gamma} \ge 20$ , as well as 623

$$\beta_{k,h} = \widetilde{\beta}_2 := C_{\beta,2} \left\lfloor \sqrt{\lambda} + HC(M,\alpha) \sqrt{\log(3HMN_h(\gamma)/\delta)} \right\rfloor$$

with  $C_{\beta,2} = 12$  for all  $k \in [K]$  and  $h \in [H]$ , where  $N_h(\gamma) = N(\mathcal{F}_h, \gamma) \cdot N(\mathcal{F}_{h+1}, \gamma) \cdot N(\mathcal{W}_{h+1}, \gamma)$ , 624

we have with probability at least  $1-\delta$  that  $\mathcal{T}_h Q_{k+1,h+1} \in \mathcal{F}_{k+1,h}$  for all  $k \in \{k_l\}_{l=1}^L$  with probability 625

at least  $1 - \delta$ . As a corollary, we also have  $|\mathcal{T}_h Q_{k+1,h+1}(s,a) - \widehat{f}_{k+1,h}(s,a)| \leq b_{k+1,h}(s,a)$  for 626 any  $(s, a) \in \mathcal{S} \times \mathcal{A}$ ,  $k \in \{k_l\}_{l=1}^L$  and  $h \in [H]$ . 627

This is the central optimism lemma. It states that the Bellman operator of Q-value function at level 628 h + 1 is within the confidences set at level h. The conclusion immediately gives the optimism 629 inequality  $\mathcal{T}_h Q_{k+1,h+1}(s,a) \leq Q_{k+1,h}(s,a)$ , which we will use at the start of the regret upper 630 bound prove. The proof of the lemma can be found in Section C.2. 631

With this, we define the good event  $\mathcal{E}_T = \{\mathcal{T}_h Q_{k+1,h+1} \in \mathcal{F}_{k+1,h}, \forall k \in \{k_l\}_{l=1}^L, h \in [H]\}$ . Then 632 according to Lemma B.2,  $\mathbb{P}(\mathcal{E}_T) \geq 1 - \delta$ . 633

**Lemma B.3.** For some absolute constant  $C_D$ , the following holds for all level  $h \in [H]$ :

$$\sum_{k=1}^{K} D_{\lambda,\mathcal{F}}^2(z_h^k; Z_{k-1,h}^{all}) \le C_D \dim_E(\mathcal{F}, \lambda/T) \log^2(T/\min\{1,\lambda\})$$

- This lemma is essentially the same as Lemma A.3. It reveals the relationship between the sum of 634
- Eluder-like confidence quantities and the Eluder dimension. The proof can be found in Section C.3. 635

### 636 B.2 The Epoch Segmentation Scheme

In this section, we introduce the epoch segmentation scheme for MDPs, which is again needed for both the regret and communication cost proofs presented in the next two sections. All of this is quite similar to the bandit case in Section A.2, but the introduction of multiple levels  $h \in [H]$  does

640 complicate things a bit.

We segment the entire run of episodes  $k = 1, \dots, K$  into N epochs as follows. Define iteratively  $0 = l_0 < l_1 < \dots < l_N \leq L$  as

$$l_{i} = \min\left\{l > l_{i-1} : \sum_{l'=l_{i-1}+1}^{l} \sum_{h=1}^{H} \sum_{o_{h} \in Z_{m,k_{l'},h}^{\text{loc}}} D_{\lambda,\mathcal{F}_{h}}^{2}(z_{h}; Z_{k_{l_{i-1}},h}^{\text{ser}}) \ge 1\right\},\$$

where for a given l' in the summation,  $m = m_{k_{l'}}$  is the participating agent at  $k_{l'}$ . In the iterative process, if the above minimum does not exist, simply define N = i - 1 and end the process there. Correspondingly, the *i*-th epoch is defined by the episodes  $[k_{l_{i-1}}, k_{l_i}]$ .

The following sections will make use of this epoch scheme as befit their needs, but here we shall give an upper bound for the total number of epochs N. Based on the definition of  $l_i$ , we have for any  $l_{i-1} \leq l < l_i$  that

$$\begin{split} 1 &\geq \sum_{l'=l_{i-1}+1}^{l} \sum_{h=1}^{H} \sum_{o_h \in Z_{m,k_{l'},h}^{\text{sec}}} D_{\lambda,\mathcal{F}}^2(z_h; Z_{k_{l_{i-1}},h}^{\text{ser}}) \\ &= \sum_{l'=l_{i-1}+1}^{l} \sum_{h=1}^{H} \sum_{o_h \in Z_{k_{l'},h}^{\text{ser}} \setminus Z_{k_{l'-1},h}^{\text{ser}}} \sup_{f_1, f_2 \in \mathcal{F}} \frac{[f_1(z_h) - f_2(z_h)]^2}{\lambda + \|f_1 - f_2\|_{Z_{k_{l_{i-1}},h}}^{\text{ser}}} \\ &\geq \sum_{h=1}^{H} \sup_{f_1, f_2 \in \mathcal{F}_h} \frac{\sum_{o_h \in Z_{k_{l'},h}^{\text{ser}} \setminus Z_{k_{l_{i-1},h}}^{\text{ser}} [f_1(z_h) - f_2(z_h)]^2}{\lambda + \|f_1 - f_2\|_{Z_{k_{l_{i-1},h}}}^{\text{ser}}} \\ &= \sum_{h=1}^{H} \bigg[ \sup_{f_1, f_2 \in \mathcal{F}_h} \frac{\lambda + \|f_1 - f_2\|_{Z_{k_{l,h}}}^2}{\lambda + \|f_1 - f_2\|_{Z_{k_{l-1},h}}^2} - 1 \bigg], \end{split}$$

647 which gives  $\lambda + \|f_1 - f_2\|_{Z^{\text{ser}}_{k_l,h}}^2 \le 2(\lambda + \|f_1 - f_2\|_{Z^{\text{ser}}_{k_{l_{i-1}},h}}^2)$  for any  $h \in [H]$  and  $f_1, f_2 \in \mathcal{F}_h$ . Then 648 we have

$$D^{2}_{\lambda,\mathcal{F}}(z_{h}; Z^{\text{ser}}_{k_{l_{i-1}},h}) \leq 2D^{2}_{\lambda,\mathcal{F}}(z_{h}; Z^{\text{ser}}_{k_{l},h})$$

$$\tag{21}$$

649 for any  $h \in [H]$  and  $z_h \in \mathcal{S} \times \mathcal{A}$ , and so

$$1 \leq \sum_{l=l_{i-1}+1}^{l_i} \sum_{h=1}^{H} \sum_{o_h \in Z_{m,k_l,h}^{\text{loc}}} D_{\lambda,\mathcal{F}_h}^2(z_h; Z_{k_{l_{i-1}},h}^{\text{ser}})$$
$$\leq 2 \sum_{l=l_{i-1}+1}^{l_i} \sum_{h=1}^{H} \sum_{o_h \in Z_{k_l,h}^{\text{ser}} \setminus Z_{k_{l-1},h}^{\text{ser}}} D_{\lambda,\mathcal{F}_h}^2(z_h; Z_{k_{l-1},h}^{\text{ser}})$$

and summing over  $i = 1, \dots, N-1$  that:

$$N-1 \le 2\sum_{h=1}^{H} \sum_{l=1}^{L} \sum_{o_h \in Z_{k_l,h}^{\text{ser}} \setminus Z_{k_{l-1},h}^{\text{ser}}} D_{\lambda,\mathcal{F}_h}^2(z_h; Z_{k_{l-1},h}^{\text{ser}}).$$

If we apply the participant reordering trick and let  $m_k = m_{k_l}$  for all  $k \in (k_{l-1}, k_l]$  and  $l \in [L]$ , we get  $Z_{k_l,h}^{\text{ser}} \setminus Z_{k_{l-1},h}^{\text{ser}} = \{o_h^k\}_{k=k_{l-1}+1}^{k_l}$ , and so applying Lemma 6.3 and Lemma 6.2, we get

$$\begin{split} N-1 &\leq 2 \sum_{h=1}^{H} \sum_{l=1}^{L} \sum_{k=k_{l-1}+1}^{k_{l}} D_{\lambda,\mathcal{F}_{h}}^{2}(z_{h}^{k};Z_{k_{l-1},h}^{\mathrm{ser}}) \\ &\leq 2(1+M\alpha) \sum_{h=1}^{H} \sum_{l=1}^{L} \sum_{k=k_{l-1}+1}^{k_{l}} D_{\lambda,\mathcal{F}_{h}}^{2}(z_{h}^{k};Z_{k-1,h}^{\mathrm{all}}) \\ &\leq 2(1+M\alpha) \sum_{h=1}^{H} \sum_{k=1}^{K} D_{\lambda,\mathcal{F}_{h}}^{2}(z_{h}^{k};Z_{k-1,h}^{\mathrm{all}}) \\ &\leq CH(1+M\alpha) \dim_{E}(\mathcal{F},\lambda/T) \log(T/\lambda) \log T, \end{split}$$

<sup>653</sup> which gives the order of total number of epochs:

$$N = O\bigg(H(1+M\alpha)\dim_E(\mathcal{F},\lambda/T)\log^2(T/\min\{1,\lambda\})\bigg).$$
(22)

## 654 B.3 Proof of Regret Upper Bound

In this section, we prove the first half of Theorem 5.1, which gives an upper bound for the cumulative regret of Algorithm 2.

- Using the participant reordering trick, assume without loss of generality that the same agent is active
- within the rounds  $[k_l, k_{l+1} 1]$ , i.e.  $m_{k_l} = m_{k_l+1} = \cdots = m_{k_{l+1}-1}$ . Under this assumption, we have  $k_1 = 1$ .

We first prove via induction that  $Q_h^* \leq Q_{m,k,h}$  for any  $m \in [M], k \in [K]$  and  $h \in [H+1]$ . This holds true for h = H + 1 trivially since both value functions at H + 1 are uniformly 0. Suppose we already have  $Q_{h+1}^* \leq Q_{m,k,h+1}$ , we have from Lemma B.2 that for the last communication round k' for agent m, the server functions satisfy  $\mathcal{T}_h Q_{k'+1,h+1}(s,a) \leq \widehat{f}_{k'+1,h}(s,a) + b_{k'+1,h}(s,a) = Q_{k'+1,h}(s,a)$ . Couple this with the fact that  $Q_{m,k,h} = Q_{k'+1,h}$ , we can prove that

$$Q_h^* = \mathcal{T}_h Q_{h+1}^* \le \mathcal{T}_h Q_{m,k,h+1} \le Q_{m,k,h},$$

<sup>660</sup> which finishes the induction process.

Now let 
$$a_h^{k*} := \operatorname{argmax}_{a \in \mathcal{A}} Q_h^*(s_h^k, a)$$
 be the best action at time  $t$ , then  $V_h^*(s_h^k) = Q_h^*(s_h^k, a_h^{k*}) \leq Q_{m,k,h}(s_h^k, a_h^{k*}) \leq Q_{m,k,h}(s_h^k, a_h^k)$ , where the second inequality is due to the choice of  $a_h^k$  at round  $k$ . Hence we get

$$\operatorname{Reg}(K) = \sum_{k=1}^{K} \left[ V_{1}^{*}(s_{1}^{k}) - V_{1}^{\pi_{k}}(s_{1}^{k}) \right]$$

$$\leq \sum_{k=1}^{K} \min \left\{ V_{m,k,1}(s_{1}^{k}) - V_{1}^{\pi_{k}}(s_{1}^{k}), 2H \right\}$$

$$= \sum_{k=1}^{K} \sum_{h=1}^{H} \min \left\{ \mathbb{E}_{\pi_{k}} \left[ Q_{m,k,h}(s_{h}^{k}, a_{h}^{k}) - \mathcal{T}_{h} Q_{m,k,h+1}(s_{h}^{k}, a_{h}^{k}) \right], 2 \right\}$$

$$= \sum_{k=1}^{K} \sum_{h=1}^{H} \min \left\{ \mathbb{E}_{\pi_{k}} \left[ \widehat{f}_{k'+1,h}(s_{h}^{k}, a_{h}^{k}) + b_{k'+1,h}(s_{h}^{k}, a_{h}^{k}) - \mathcal{T}_{h} Q_{m,k,h+1}(s_{h}^{k}, a_{h}^{k}) \right], 2 \right\}$$

$$\leq \sum_{k=1}^{K} \sum_{h=1}^{H} \min \left\{ 2b_{k',h}(s_{h}^{k}, a_{h}^{k}), 2 \right\}$$

$$= 2\sum_{l=1}^{L} \sum_{k=k_{l}+1}^{k_{l+1}-1} \sum_{h=1}^{H} b_{k_{l}+1,h}(z_{h}^{k}) + 2\sum_{l=1}^{L} \sum_{h=1}^{H} \min \{ b_{m_{k_{l}},k_{l},h}(z_{h}^{k_{l}}), 1 \}.$$
(23)

where the second equality uses the Value-decomposition Lemma from Jiang et al. [2017], the second inequality uses again Lemma B.2, and from the third inequality onward we let k' be the last time agent m communicated with the server. We now bound the second term here using the epoch scheme in Section B.2. We start by converting the bonus term to uncertainty:

$$b_{m_{k_{l}},k_{l},h}(z_{h}^{k_{l}}) = b_{k_{l-1},h}(z_{h}^{k_{l}})$$
  
$$\leq C_{\mathcal{B}}\sqrt{\beta_{k_{l}-1,h}^{2} + \lambda} \cdot D_{\lambda,\mathcal{F}_{h}}(z_{h}^{k_{l}}; Z_{k_{l-1},h}^{\text{ser}}).$$
(24)

Now consider the episodes in an epoch *i*, specifically  $\{k_{l_{i-1}}, k_{l_{i-1}+1}, \cdots, k_{l_i}\}$ . For any  $l_{i-1} < l < l_i$ , since  $Z_{k_{l_{i-1}},h}^{\text{ser}} \subseteq Z_{k_{l_{i-1}},h}^{\text{ser}}$ , we can deduce that

$$D^{2}_{\lambda,\mathcal{F}_{h}}(z_{h}^{k_{l}}; Z_{k_{l-1},h}^{\text{ser}}) \leq D^{2}_{\lambda,\mathcal{F}_{h}}(z_{h}^{k_{l}}; Z_{k_{l_{l-1}},h}^{\text{ser}}) \leq 2D^{2}_{\lambda,\mathcal{F}_{h}}(z_{h}^{k_{l}}; Z_{k_{l,h}}^{\text{ser}}),$$

where the second inequality is borrowed from (21) from Section B.2. Therefore continuing from (24),

$$\sum_{l=1}^{L} \sum_{h=1}^{H} \min\{b_{m_{k_{l}},k_{l},h}(z_{h}^{k_{l}}),1\} \leq \sum_{l \notin \{l_{i}\}_{i=1}^{N}} \sum_{h=1}^{H} \left[\sqrt{2}C_{\mathcal{B}}\sqrt{\beta_{k_{l}-1,h}^{2} + \lambda} \cdot D_{\lambda,\mathcal{F}_{h}}(z_{h}^{k_{l}};Z_{k_{l},h}^{\mathrm{ser}})\right] + \sum_{i=1}^{N} \sum_{h=1}^{H} 1 \sum_{h=1}^{H} 1 \sum_{h=1}^{L} \sum_{h=1}^{H} D_{\lambda,\mathcal{F}_{h}}(z_{h}^{k_{l}};Z_{k_{l},h}^{\mathrm{ser}})\sqrt{\beta_{h}^{2} + \lambda} + NH.$$
(25)

Now combine this result with the first term in (23) and use again (24), we get

$$\begin{aligned} \operatorname{Reg}(K) &\leq C_{\mathcal{B}} \sum_{l=1}^{L} \sum_{k=k_{l}+1}^{k_{l+1}-1} \sum_{h=1}^{H} \left[ D_{\lambda,\mathcal{F}_{h}}(z_{h}^{k}; Z_{k_{l},h}^{\operatorname{ser}}) \sqrt{\beta_{h}^{2} + \lambda} \right] + \sqrt{2} C_{\mathcal{B}} \sum_{l=1}^{L} \sum_{h=1}^{H} \left[ D_{\lambda,\mathcal{F}_{h}}(z_{h}^{k}; Z_{k_{l},h}^{\operatorname{ser}}) \sqrt{\beta_{h}^{2} + \lambda} \right] + NH \\ &\leq \sqrt{2} C_{\mathcal{B}} \sum_{l=1}^{L} \sum_{k=k_{l}}^{k_{l+1}-1} \sum_{h=1}^{H} \left[ D_{\lambda,\mathcal{F}_{h}}(z_{h}^{k}; Z_{k_{l},h}^{\operatorname{ser}}) \sqrt{\beta_{h}^{2} + \lambda} \right] + NH \\ &\leq \sqrt{2} C_{\mathcal{B}} \left[ \sum_{l=1}^{L} \sum_{k=k_{l}}^{k_{l+1}-1} \sum_{h=1}^{H} D_{\lambda,\mathcal{F}_{h}}^{2}(z_{h}^{k}; Z_{k_{l},h}^{\operatorname{ser}}) \right]^{1/2} \left[ \sum_{k=1}^{K} \sum_{h=1}^{H} \left( \beta_{h}^{2} + \lambda \right) \right]^{1/2} + NH. \end{aligned}$$

674 According to Lemma 6.3 and Lemma 6.2, the term

$$\sum_{l=1}^{L} \sum_{k=k_{l}}^{k_{l+1}-1} \sum_{h=1}^{H} D_{\lambda,\mathcal{F}_{h}}^{2}(z_{h}^{k}; Z_{k_{l},h}^{\text{ser}}) \leq (1+M\alpha) \sum_{l=1}^{L} \sum_{k=k_{l}}^{k_{l+1}-1} \sum_{h=1}^{H} D_{\lambda,\mathcal{F}_{h}}^{2}(z_{h}^{k}; Z_{k-1,h}^{\text{all}})$$
$$\leq (1+M\alpha) \sum_{k=1}^{K} \sum_{h=1}^{H} D_{\lambda,\mathcal{F}_{h}}^{2}(z_{h}^{k}; Z_{k-1,h}^{\text{all}})$$
$$\leq H(1+M\alpha) \dim_{E}(\mathcal{F}, \lambda/T) \log(T/\lambda) \log T.$$

Now with  $\gamma = O(1/KH)$ , we have

$$\beta_h = O(1)\beta_{h+1} + C_\beta \left[ \sqrt{\lambda} + H\left(\sqrt{(1+M\alpha)\log(3HN_h(\gamma)/\delta)} + M\sqrt{\alpha\log(3HMN_h(\gamma)/\delta)} \right) \right]$$

therefore, with  $C(M, \alpha) = \sqrt{1 + M\alpha} + M\sqrt{\alpha}$  and the upper bound for number of epochs N in (22), we have

$$\begin{split} &\sum_{l=1}^{L}\sum_{k=k_{l}+1}^{k_{l+1}-1}\sum_{h=1}^{H}b_{k_{l},h}(z_{h}^{k}) \\ &\leq O\bigg(\Big[H(1+M\alpha)\dim_{E}(\mathcal{F},\lambda/K)\log^{2}(K/\min\{1,\lambda\})\Big]^{1/2}\bigg[K\sum_{h=1}^{H}(\beta_{h}^{2}+\lambda)\bigg]^{1/2}+HN\bigg) \\ &= O\bigg(H\sqrt{K}\widetilde{\beta}_{2}\sqrt{(1+M\alpha)\dim_{E}(\mathcal{F},\lambda/K)}\log(K/\min\{1,\lambda\}) \\ &\quad +H^{2}(1+M\alpha)\dim_{E}(\mathcal{F},\lambda/K)\log^{2}(K/\min\{1,\lambda\})\bigg), \end{split}$$

678 where 
$$\tilde{\beta}_2 = C_{\beta,2} \left[ \sqrt{\lambda} + HC(M,\alpha) \log \left( HMN(\mathcal{F},\gamma)N(\mathcal{W},\gamma)/\delta \right) \right]$$
 is the choice of  $\beta_{k,h}$  in the

679 algorithm.

#### 680 B.4 Proof of Communication Cost

Next up, we calculate the communication complexity of Algorithm 2 and prove the second half of Theorem 5.1. For each communication round  $k_l$ , assume the last time before  $k_l$  when the agent  $m = m_{k_l}$  communicated with the server was  $k_{l'}$ , then by the communication rule there exists  $h_l \in [H]$  such that  $\sum_{o_{h_l} \in Z_{m,k_l,h_l}^{\text{loc}}} b_{k_{l'},h_l}^2(z_{h_l})/(\beta_{k_{l'},h_l}^2 + \lambda) \ge \alpha$ ,

$$\sum_{p_{h_l} \in Z_{m,k_l,h_l}^{\text{loc}}} D_{\lambda,\mathcal{F}_{h_l}}^2(z_{h_l}; Z_{m,k_l,h_l}^{\text{up}}) \ge \sum_{o_{h_l} \in Z_{m,k_l,h_l}^{\text{loc}}} \frac{\left[ b_{k_{l'},h_l}(z_{h_l})/C \right]^2}{\beta_{k_{l'},h_l}^2 + \lambda} \ge \frac{\alpha}{C^2},$$

Next we will make use of the epoch segmentation scheme in Section B.2. For the *i*-th epoch consisting of the time steps  $[k_{l_{i-1}}, k_{l_i})$ , we have the inequality

$$1 \ge \sum_{l=l_{i-1}+1}^{l_i-1} \sum_{h=1}^{H} \sum_{o_h \in Z_{m,k_l,h}^{\text{loc}}} D_{\lambda,\mathcal{F}_h}^2(z_h; Z_{k_{l_{i-1}},h}^{\text{ser}})$$
$$\ge \sum_{l=l_{i-1}+1}^{l_i-1} \sum_{h=1}^{H} \sum_{o_h \in Z_{m,k_l,h}^{\text{loc}}} D_{\lambda,\mathcal{F}_h}^2(z_h; Z_{m,k_l,h}^{\text{up}} \cup Z_{k_{l_{i-1}},h}^{\text{ser}})$$

For  $m \in [M]$ , assume the agent m communicated with the server a total of  $n_m$  times within  $[k_{l_{i-1}}, k_{l_i})$ . Then except for the first of these communication rounds, for each  $l \in [l_{i-1} + 1, l_i - 1]$  with  $m_{k_l} = m$ , there exists  $l' \in [l_{i-1}, l)$  with  $m_{k_{l'}} = m$ , thus we have  $Z_{m,k_l,h}^{up} \supset Z_{m,k_{l'}+1,h}^{up} = Z_{k_{l'},h}^{ser} \supset Z_{k_{l-1},h}^{ser}$  for all  $h \in [H]$ . With this we have

$$\sum_{h=1}^{H} \sum_{o_h \in Z_{m,k_l,h}^{\mathrm{loc}}} D_{\lambda,\mathcal{F}_h}^2 \left( z_h; Z_{m,k_l,h}^{\mathrm{up}} \cup Z_{k_{l_{i-1}},h}^{\mathrm{ser}} \right) = \sum_{h=1}^{H} \sum_{o_h \in Z_{m,k_l,h}^{\mathrm{loc}}} D_{\lambda,\mathcal{F}_h}^2 \left( z_h; Z_{m,k_l,h}^{\mathrm{up}} \right) \ge \frac{\alpha}{4C^2},$$

687 therefore

$$1 \ge \sum_{m=1}^{M} (n_m - 1) \cdot \frac{\alpha}{4C^2} \Rightarrow \sum_{m=1}^{M} n_m \le M + \frac{4C^2}{\alpha}$$

Notice that  $\sum_{m=1}^{M} n_m = l_i - l_{i-1}$  is the number of communication rounds within  $[k_{l_{i-1}}, k_{l_i})$ , hence summing over *i* the total number of communication rounds is upper bounded by  $N(M + 4C^2/\alpha)$ . Combine this with the result in (22), we have the total number of communication rounds throughout the algorithm is

$$O\left(H\frac{(1+M\alpha)^2}{\alpha}\dim_E(\mathcal{F},\lambda/K)\log^2(K/\min\{1,\lambda\})\right).$$

### 688 C Proof of Auxiliary Lemmas

In this section we prove all the auxiliary lemmas in Section A.1 and Section *B*.1. Note that some of these lemmas are very similar in nature, for which we will only give the proof for the version for the MDPs case, and briefly remark on the version for the bandit case.

#### 692 C.1 Proof of Lemma A.1 and Lemma B.1

Here we prove Lemma B.1 in detail. The proof for Lemma A.1 is very similar, and so we will only give a short remark on how to apply this to the bandit case.

Proof of Lemma B.1. First, for an episode  $k \in [K]$  and agent  $m \in [M]$  such that m does not communicate with the server at episode k (either m is not participating or k is not a communication <sup>697</sup> round), from the communication criterion we have

$$\begin{split} \alpha &\geq \sum_{o_h \in Z_{m,k,h}^{\text{loc}}} \frac{b_{m,k,h}^2(a)}{\beta_{k',h}^2 + \lambda} \\ &\geq \sum_{o_h \in Z_{m,k,h}^{\text{loc}}} D_{\lambda,\mathcal{F}_h}^2(z_h; Z_{k',h}^{\text{ser}}) \\ &= \sum_{o_h \in Z_{m,k,h}^{\text{loc}}} \sup_{f_1, f_2 \in \mathcal{F}_h} \frac{|f_1(z_h) - f_2(z_h)|^2}{\lambda + \|f_1 - f_2\|_{Z_{k',h}^{\text{ser}}}} \\ &\geq \sup_{f_1, f_2 \in \mathcal{F}_h} \frac{\|f_1 - f_2\|_{Z_{m,k,h}^{\text{loc}}}^2}{\lambda + \|f_1 - f_2\|_{Z_{k',h}^{\text{loc}}}^2}, \end{split}$$

where k' is the last communication round for agent m. This means that for any  $f_1, f_2 \in \mathcal{F}_h$ ,  $(1/\alpha) \|f_1 - f_2\|_{Z^{\text{loc}}_{m,k,h}}^2 \leq \lambda + \|f_1 - f_2\|_{Z^{\text{ser}}_{k',h}}^2$ . Observing that  $Z^{\text{ser}}_{k',h} \subset Z^{\text{ser}}_{k,h} = \bigcup_{m'=1}^M Z^{\text{up}}_{m',k,h}$  proves the first conclusion that

$$\frac{1}{\alpha} \|f_1 - f_2\|_{Z^{\text{loc}}_{m,k,h}}^2 \le \lambda + \sum_{m'=1}^M \|f_1 - f_2\|_{Z^{\text{up}}_{m',k,h}}^2.$$

Second, for any  $f_1, f_2 \in \mathcal{F}_h$ , from the above conclusion we have for any  $k \in [K] \setminus \{k_l\}_{l=1}^L$  that

$$\begin{split} \lambda + \|f_1 - f_2\|_{Z_{k,h}^{\text{ser}}}^2 &= \lambda + \sum_{m=1}^M \|f_1 - f_2\|_{Z_{m,k,h}^{\text{up}}}^2 \\ &\geq \frac{1}{M\alpha} \sum_{m=1}^M \|f_1 - f_2\|_{Z_{m,k,h}^{\text{loc}}}^2 \\ &= \frac{1}{M\alpha} \|f_1 - f_2\|_{Z_{k,h}^{\text{all}} \setminus Z_{k,h}^{\text{ser}}}^2, \end{split}$$

and when  $k = k_l$  for some  $l \in [L]$ , we have alternatively

$$\begin{split} \lambda + \|f_1 - f_2\|_{Z_{k,h}^{\text{ser}}}^2 &= \lambda + \sum_{m' \neq m_t} \|f_1 - f_2\|_{Z_{m',k,h}^{\text{up}}}^2 + \|f_1 - f_2\|_{Z_{m_k,k,h}^{\text{up}} \cup Z_{m_k,k,h}^{\text{loc}}} \\ &\geq \lambda + \sum_{m=1}^M \|f_1 - f_2\|_{Z_{m,k,h}^{\text{up}}}^2 \\ &\geq \frac{1}{(M-1)\alpha} \sum_{m' \neq m_k} \|f_1 - f_2\|_{Z_{m',k,h}^{\text{loc}}}^2 \\ &\geq \frac{1}{M\alpha} \|f_1 - f_2\|_{Z_{k,h}^{\text{ul}} \setminus Z_{k,h}^{\text{ser}}}^2. \end{split}$$

Either way, we can deduce for any  $k \in [K]$  that

$$(1+M\alpha)\left(\lambda+\|f_1-f_2\|_{Z_{k,h}^{\text{ser}}}^2\right) \ge \lambda+\|f_1-f_2\|_{Z_{k,h}^{\text{sel}}}^2$$

<sup>700</sup> Finally, from the above we immediately have

$$\begin{split} D^2_{\lambda,\mathcal{F}}(z_h; Z^{\text{ser}}_{k,h}) &= \sup_{f_1, f_2 \in \mathcal{F}_h} \frac{[f_1(z_h) - f_2(z_h)]^2}{\lambda + \|f_1 - f_2\|^2_{Z^{\text{ser}}_{k,h}}} \\ &\leq (1 + M\alpha) \sup_{f_1, f_2 \in \mathcal{F}} \frac{[f_1(z_h) - f_2(z_h)]^2}{\lambda + \|f_1 - f_2\|^2_{Z^{\text{all}}_{k,h}}} \\ &= (1 + M\alpha) D^2_{\lambda,\mathcal{F}}(a; Z^{\text{all}}_{k,h}). \end{split}$$

701

*Remark* C.1. Notice that this prove does not depend on the multi-level structure of episodic MDPs, but is a direct result of the communication criterion and protocol. This means the proof can be converted to the bandit case of Lemma A.1 without any essential changes: simply change episode k

## into time step t, disregard all mentions of level h, and consider z = a instead of z = (s, a).

#### 706 C.2 Proof of Lemma A.2 and Lemma B.2

<sup>707</sup> We begin with the proof of Lemma A.2, which is an almost direct application of Lemma D.3.

*Proof of Lemma A.2.* We invoke Lemma D.3 with  $\epsilon_0 = 0$ , then with probability at least  $1 - \delta$ , for all  $t \in \{t_l\}_{l=1}^L$ ,

$$\sum_{(a,r)\in Z_t^{\text{ser}}} \left(\widehat{f}_{t+1}(a) - f^*(a)\right)^2 \le C_{\text{ERM}} \left[\lambda + \gamma^2 T + \gamma T R + R^2 (1 + M\alpha) \log(3N/\delta) + R^2 M^2 \alpha \log(3NM/\delta)\right] \le \widetilde{\beta}_1^2,$$

708 if we let  $\gamma = O(1/T)$  be sufficiently small and take  $\tilde{\beta}_1 = C_{\beta,1} \left[ \sqrt{\lambda} + C_{\beta,1} \right]$ 

<sup>709</sup>  $RC(M,\alpha)\log(3MN(\mathcal{F},\gamma)/\delta)$  with  $C_{\beta,1} = \sqrt{C_{\text{ERM}}} = 6$ . Thus taking  $\beta_t = \tilde{\beta}_1$ , accord-<sup>710</sup> ing to the definition of  $\mathcal{F}_{t+1}$ , this directly implies  $f^* \in \mathcal{F}_{t+1}$ .

With this, since the bonus function satisfy

$$b_{t+1}(a) \ge |f_1(a) - f_2(a)|, \quad \forall f_1, f_2 \in \mathcal{F} \quad \text{s.t.} \quad \sum_{(a,r)\in Z_t^{\text{ser}}} \left(f_1(a) - f_2(a)\right)^2 \le \beta_t^2,$$

which is based on the first property of the bonus oracle in Definition 4.1, by taking  $f_1 = \hat{f}_{t+1}$  and

712  $f_2 = f^*$  we get for any  $a \in \mathcal{A}$  that  $b_{t+1}(a) \ge |f_*(a) - \widehat{f}_{t+1}(a)|$ , which finishes the proof.

Next we prove Lemma B.2, which is more challenging and requires an analysis on the least squares
 value iteration method.

<sup>715</sup> Proof of Lemma B.2. Take  $\mathcal{F}_{h+1,\gamma}$  as a  $\gamma$ -cover of  $\mathcal{F}_{h+1}$ , and  $\mathcal{W}_{h+1,\gamma}$  as a  $\gamma$ -cover of  $\mathcal{W}_{h+1}$ . Select <sup>716</sup>  $\overline{f}_{k+1,h+1} \in \mathcal{F}_{h+1,\gamma} \bigoplus \mathcal{W}_{h+1,\gamma}$  so that  $||Q_{k+1,h+1} - \overline{f}_{k+1,h+1}||_{\infty} \leq \overline{\epsilon} := (1 + \beta_{k+1,h+1})\gamma$ . For <sup>717</sup>  $o_h = (s_h, a_h, r_h, s_{h+1})$ , define the corresponding  $y_h = r_h + V_{k+1,h+1}(s_{h+1})$  and  $\overline{y}_h = r_h + \gamma_{18}$ <sup>718</sup>  $\sup_{a \in \mathcal{A}} \overline{f}_{k+1,h+1}(s_{h+1}, a)$ . Let

$$\widetilde{f}_{k+1,h} = \operatorname*{argmin}_{f_h \in \mathcal{F}_h} \sum_{o_h \in Z_{k,h}^{\operatorname{ser}}} \left( f_h(s_h, a_h) - \bar{y}_h \right)^2.$$

719 Then we have

$$\left(\sum_{o_h \in Z_{k,h}^{\operatorname{ser}}} \left(\widehat{f}_{k+1,h}(s_h, a_h) - \bar{y}_h\right)^2\right)^{1/2} \le \left(\sum_{o_h \in Z_{k,h}^{\operatorname{ser}}} \left(\widehat{f}_{k+1,h}(s_h, a_h) - y_h\right)^2\right)^{1/2} + \bar{\epsilon}\sqrt{k}$$
$$\le \left(\sum_{o_h \in Z_{k,h}^{\operatorname{ser}}} \left(\widetilde{f}_{k+1,h}(s_h, a_h) - y_h\right)^2\right)^{1/2} + \bar{\epsilon}\sqrt{k}$$
$$\le \left(\sum_{o_h \in Z_{k,h}^{\operatorname{ser}}} \left(\widetilde{f}_{k+1,h}(s_h, a_h) - \bar{y}_h\right)^2\right)^{1/2} + 2\bar{\epsilon}\sqrt{k}.$$

Now notice that  $\mathbb{E}\bar{y}_h = \mathcal{T}_h \bar{f}_{k+1,h}(s_h, a_h)$ , and the difference  $\bar{y}_h - \mathcal{T}_h \bar{f}_{k+1,h}(s_h, a_h)$  is bounded in [-H, H], hence we may apply Lemma D.3 with  $f^* = \mathcal{T}_h \bar{f}_{k+1,h}$ ,  $r_t = \bar{y}_h$ , R = H,  $\epsilon_0 = 2\bar{\epsilon}$ and  $\delta = \delta/3HN(\mathcal{F}_{h+1}, \gamma) \cdot N(\mathcal{W}_{h+1}, \gamma)$ , taking a union bound over  $\bar{f} \in \mathcal{F}_{h+1,\gamma} \bigoplus \mathcal{W}_{h+1,\gamma}$  and 723  $h \in [H]$ , we have

$$\begin{split} & \left(\sum_{o_h \in Z_{k,h}^{\text{ser}}} \left(\widehat{f}_{k+1,h}(s_h, a_h) - \mathcal{T}_h Q_{k+1,h+1}(s_h, a_h)\right)^2\right)^{1/2} \\ \leq & \left(\sum_{o_h \in Z_{k,h}^{\text{ser}}} \left(\widehat{f}_{k+1,h}(s_h, a_h) - \mathcal{T}_h \overline{f}_{k+1,h+1}(s_h, a_h)\right)^2\right)^{1/2} + \gamma \sqrt{k} \\ \leq & \sqrt{C_{\text{ERM}}} \sqrt{\lambda + (\gamma + 2\overline{\epsilon})^2 K + (\gamma + 2\overline{\epsilon}) K H + H^2 (1 + M\alpha) \log(3HN_h(\gamma)/\delta) + H^2 M^2 \alpha \log(3HMN_h(\gamma)/\delta)} + \gamma \sqrt{k} \\ \leq & \sqrt{C_{\text{ERM}}} \left[\sqrt{\lambda} + \gamma (3 + 2\beta_{k+1,h+1}) \sqrt{K} + \sqrt{\gamma (3 + 2\beta_{k+1,h+1}) K H} + HC(M, \alpha) \sqrt{\log(3HMN_h(\gamma)/\delta)}\right)\right], \end{split}$$

where  $N_h(\gamma) = N(\mathcal{F}_h, \gamma) \cdot N(\mathcal{F}_{h+1}, \gamma) \cdot N(\mathcal{W}_{h+1}, \gamma)$ . By taking  $\gamma = 1/(C_{\gamma}KH)$  with sufficiently large absolute constant  $C_{\gamma}$  (for example,  $C_{\gamma} = 20$ ), the second and third terms within the bracket above are both less than  $(1/2)\beta_{k+1,h+1}$ , and hence we can easily prove via induction on h that the above is no greater than  $\tilde{\beta}_2$ , where

$$\widetilde{\beta}_2 = C_{\beta,2} \left[ \sqrt{\lambda} + HC(M,\alpha) \sqrt{\log(3HMN(\gamma)/\delta)} \right]$$

with  $C_{\beta,2} = 2\sqrt{C_{\text{ERM}}} = 12$  and  $N(\gamma) = \max_{h \in [H]} N_h(\gamma)$ .

## 726 C.3 Proof of Lemma A.3 and Lemma B.3

In this section we prove Lemma B.3 in detail. The proof for Lemma A.3 is very similar, and so we
 will again only give a short remark on how to apply this to the bandit case.

Proof of Lemma B.3. We fix the level  $h \in [H]$  throughout the proof. For an index set  $\mathcal{K}_0 \subseteq [K]$ , we denote  $\mathcal{Z}(\mathcal{K}_0) := \{z_h^k : k \in \mathcal{K}_0\}.$ 

First, let  $n = \lceil \log(K/\lambda) / \log 2 \rceil$ , and we divide the set of episodes  $\mathcal{K} = [K]$  into n + 1 disjoint episode sets as follows. For any  $1 \le l \le L$  and  $k_l \le k < k_{l+1}$ , let

$$(\bar{f}_{k,1}, \bar{f}_{k,2}) = \operatorname*{argmax}_{f_1, f_2 \in \mathcal{F}_h} \frac{\left(f_1(z_h^k) - f_2(z_h^k)\right)^2}{\lambda + \|f_1 - f_2\|_{Z_1^{\mathrm{all}}}^2},$$

and define  $L_k : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$  as  $L_k(z) = (\bar{f}_{k,1}(z) - \bar{f}_{k,2}(z))^2$ . Now we define  $\mathcal{K}^{\iota} := \{k \in \mathcal{K} : L_k(z_h^k) \in (2^{-\iota-1}, 2^{-\iota}]\}$  for  $\iota \in \{0, 1, \cdots, n-1\}$  and  $\mathcal{K}^n := \{k \in \mathcal{K} : L_k(z_h^k) \in [0, 2^{-n}]\}$ . We note that for  $k \in \mathcal{K}^n$ ,  $L_k(z_h^k) \leq \lambda/K$ .

Now define the mapping  $\tau : [K] \to [K]$ , such that for any  $k \in [K]$ ,  $\tau(k)$  is the last episode when agent  $m_k$  communicated with the server (not including k). We will bound  $\sum_{k \in \mathcal{K}^{\iota}} D^2_{\lambda, \mathcal{F}_h}(z_h^k; Z_{h,k-1}^{\text{all}})$ for  $\iota \in \{0, \dots, n-1\}$ .

For a fixed  $\iota \leq n-1$ , we now decompose  $\mathcal{K}^{\iota} = \bigcup_{j=1}^{n^{\iota}+1} \mathcal{K}_{j}^{\iota}$ , where  $n^{\iota} = \left\lceil |\mathcal{K}^{\iota}| / \dim_{E}(\mathcal{F}_{h}, 2^{-\iota-1}) \right\rceil$ . 737 We start off each set  $\mathcal{K}_{i}^{\iota} = \emptyset$ , and fill them up gradually by iterating through  $k \in \mathcal{K}^{\iota}$  one by one 738 in increasing order to decide which subset  $\mathcal{K}_j^\iota$  should k belong to. Specifically, we define j(k) to 739 be the smallest index  $j < n^{\iota}$  such that is  $z_h^k$  is  $2^{-(\iota+1)/2}$ -independent of  $\mathcal{Z}(\mathcal{K}_j^{\iota})$ , and assign k to 740 the set  $\mathcal{K}_{i(k)}^{\iota}$ . If such a j does not exist, we simply let  $j(k) = n^{\iota} + 1$  assign k to  $\mathcal{K}_{n^{\iota}+1}^{\iota}$ . Finally 741 after the assignment process, we define  $\mathcal{K}_{j,k}^{\iota} = \mathcal{K}_{j}^{\iota} \cap [k]$  for any  $k \in [K]$ . Then we have the elements added into  $\mathcal{K}_{j(k)-1,k}^{\iota}$  form a sequence where each data corresponding to a new member 742 743 is  $2^{-(\iota+1)/2}$ -independent of the old members, and so there are no more than  $\dim_E(\mathcal{F}_h, 2^{-\iota-1})$ 744 members within each of them. Moreover, for all  $k \in \mathcal{K}^{\iota}$  that  $z_h^k$  is  $2^{-(\iota+1)/2}$ -dependent on each of 745  $\mathcal{Z}(\mathcal{K}_{1,k}^{\iota}), \cdots, \mathcal{Z}(\mathcal{K}_{j(k)-1,k}^{\iota}).$ 746

Now for any  $k \in \mathcal{K}^{\iota}$  by the definition of  $\mathcal{K}^{\iota}$ , we have  $(\bar{f}_{k,1}(z_h^k) - \bar{f}_{k,2}(z_h^k))^2 \ge 2^{-\iota-1}$ . This combined with the  $2^{-\iota-1}$ -dependencies imply that for each  $j' = 1, \cdots, j(k) - 1, \|\bar{f}_{k,1} - \bar{f}_{k,2}\|_{\mathcal{Z}(\mathcal{K}^{\iota}_{j',k})}^2 \ge 2^{-\iota-1}$ .

Notice that  $\mathcal{Z}(\mathcal{K}_{j',k}^{\iota}) \subset Z_{h,k-1}^{\text{all}}$  for any  $j' \in [j(k) - 1]$ , and that  $\mathcal{Z}(\mathcal{K}_{j',k}^{\iota})$  for  $j' \in [j(k) - 1]$  are disjoint, therefore

$$(j(k)-1)2^{-\iota-1} \le \sum_{j'=1}^{j(k)-1} \|\bar{f}_{k,1} - \bar{f}_{k,2}\|_{\mathcal{Z}(\mathcal{K}_{j',k}^{\iota})}^2 \le \|\bar{f}_{k,1} - \bar{f}_{k,2}\|_{Z_{h,k-1}^{\text{all}}}^2.$$

747 It follows that

$$\begin{split} D^2_{\lambda,\mathcal{F}_h}(z_h^k;Z_{h,k-1}^{\text{all}}) &= \frac{\left(\bar{f}_{k,1}(z_h^k) - \bar{f}_{k,2}(z_h^k)\right)^2}{\lambda + \|\bar{f}_{k,1} - \bar{f}_{k,2}\|_{Z_{h,k-1}^{\text{all}}}^2} \\ &\leq \frac{2^{-\iota}}{\lambda + (j(k) - 1)2^{-\iota-1}} \\ &= \frac{2}{(j(k) - 1) + 2^{\iota+1}\lambda}, \end{split}$$

where the first inequality uses the definition of  $\mathcal{K}^{\iota}$ . Summing over  $k \in \mathcal{K}^{\iota}$ , we have

$$\begin{split} \sum_{k \in \mathcal{K}^{\iota}} D_{\lambda,\mathcal{F}_{h}}^{2}(z_{h}^{k}; Z_{h,k-1}^{\text{all}}) &= \sum_{j=1}^{n^{\iota}+1} \sum_{k \in \mathcal{K}_{j}^{\iota}} D_{\lambda,\mathcal{F}_{h}}^{2}(z_{h}^{k}; Z_{h,k-1}^{\text{all}}) \\ &\leq \sum_{j=1}^{n^{\iota}} \frac{2|\mathcal{K}_{j}^{\iota}|}{(j-1)+2^{\iota+1}\lambda} + \frac{2|\mathcal{K}_{n^{\iota}+1}^{\iota}|}{n^{\iota}+2^{\iota+1}\lambda} \\ &\leq \frac{2\dim_{E}(\mathcal{F}_{h}, 2^{-\iota-1})}{2^{\iota+1}\lambda} + \sum_{j=2}^{n^{\iota}} \frac{2\dim_{E}(\mathcal{F}_{h}, 2^{-\iota-1})}{j-1} + 2|\mathcal{K}^{\iota}| \cdot \frac{\dim_{E}(\mathcal{F}_{h}, 2^{-\iota-1})}{|\mathcal{K}^{\iota}|} \\ &\leq \dim_{E}(\mathcal{F}_{h}, 2^{-\iota-1})(2\log n^{\iota} + 4 + 1/(2^{\iota}\lambda)), \end{split}$$

where we used the relation  $|\mathcal{K}_{j}^{\iota}| \leq \dim_{E}(\mathcal{F}_{h}, 2^{-\iota-1})$  and the definition of  $n^{\iota}$  in the second inequality. Additionally, for  $\iota = n$  we also have

$$\sum_{k \in \mathcal{K}^n} D^2_{\lambda, \mathcal{F}_h}(z_h^k; Z_{h, k-1}^{\text{all}}) \le \sum_{k \in \mathcal{K}^n} \frac{L_k(z_h^k)}{\lambda} \le |\mathcal{K}^n| \cdot \frac{\lambda/K}{\lambda} \le 1,$$

and so finally we sum over  $\iota = 0, \cdots, n$  to get

$$\sum_{k=1}^{K} D_{\lambda,\mathcal{F}_{h}}^{2}(z_{h}^{k}; Z_{h,k-1}^{\text{all}}) \leq \sum_{\iota=0}^{n-1} \dim_{E}(\mathcal{F}_{h}, 2^{-\iota-1}) (2\log n^{\iota} + 4 + 1/(2^{\iota}\lambda)) + 1$$
$$\leq n \dim_{E}(\mathcal{F}_{h}, 2^{-n}) (2\log K + 4 + 1/\lambda) + 1$$
$$\leq C \dim_{E}(\mathcal{F}_{h}, \lambda/K) \log(K/\min\{1,\lambda\}),$$

where the final step makes the assumption that  $\lambda = O(1/\log K)$ , in which case it holds with some absolute constant  $C_D$ .

*Remark* C.2. Again, this prove does not depend on the multi-level structure of episodic MDPs. In fact, it only relies on the Eluder dimensionality of  $\mathcal{F}_h$ . This means the proof can be converted to the bandit case of Lemma A.3 without any essential changes: simply change episode k into time step t, disregard all mentions of level h, and consider z = a instead of z = (s, a).

## 758 **D** Technical Lemmas

In this section, we provide a technical concentration lemma that serves as the core of our results. For
 one, this lemma is based on the following concentration inequality:

**Lemma D.1.** For a sequence of random variables  $\{Z_t\}_{t\in\mathbb{N}}$  adapted to the filtration  $\{S_t\}_{t\in\mathbb{N}}$  and function  $f \in \mathcal{F}$ , for any  $\lambda > 0$ , with probability at least  $1 - \delta$ , for all  $t \in \mathbb{N}$ , we have

$$-\frac{1}{\lambda} \sum_{s=1}^{t} \log \mathbb{E} \left[ \exp[-\lambda f(Z_s)] \middle| \mathcal{S}_{s-1} \right] - \sum_{s=1}^{t} f(Z_s) \le \frac{1}{\lambda \delta}.$$

The proof for this lemma can be found under Lemma 4 of Russo and Van Roy [2013]. Apart from 763 this, we need yet another basic concentration lemma: 764

**Lemma D.2.** Suppose  $\{\eta_t\}_{t=1}^T$  is a sequence of conditional *R*-sub-Gaussian random variables satisfying  $\mathbb{E}[e^{\mu\eta_t}|\mathcal{H}_{t-1}] \leq \exp(R^2\mu^2/2)$ , where  $\mathcal{H}_{t-1}$  denotes all history before time t, with probability  $1 - \delta$ , we have

$$\sum_{t=1}^T \eta_t^2 \le 2T\sigma^2 + 3\sigma^2 \log(1/\delta).$$

A proof of this lemma can be found under Lemma G.2 of Ye et al. [2023]. With this, we can prove 765

the following lemma characterizing the accuracy of least squares solution. Even though we need 766 this lemma for both bandit and RL settings, we will follow the notations presented in multi-agent 767

contextual bandits. Detailed explanation of how this translates to multi-agent MDPs can be found in 768

Section C.2. 769

**Lemma D.3.** Suppose we have a sequence of inputs  $\{(a_t, r_t)\}_{t=1}^T$  that follow the rule  $r_t = f^*(a_t) + \eta_t$  for some ground truth  $f^* \in \mathcal{F}$ , with  $\eta_t$  being conditionally *R*-sub-Gaussian:

$$\mathbb{E}\left[e^{\mu\eta_t} \middle| a_{1:t}, r_{1:t-1}\right] \le \exp(R^2 \mu^2/2), \forall \mu \in \mathbb{R}.$$

We also have server datasets  $Z_t^{ser}$  at different time steps, collected following the communication protocol in our settings. Note that strictly speaking, the conditions under which  $\eta_t$  is sub-Gaussian 770

771

- should also include the former participants  $m_{1:t}$ , but we will omit this dependency for convenience. 772
- Consider  $\hat{f}_{t+1}^{ser}$ , the approximate ERM solution to the least squares problem: 773

$$\left(\sum_{(a,r)\in Z_t^{ser}} \left(\widehat{f}_{t+1}^{ser}(a) - r\right)^2\right)^{1/2} \le \min_{f\in\mathcal{F}_t} \left(\sum_{(a,r)\in Z_t^{ser}} \left(f(a) - r\right)^2\right)^{1/2} + \epsilon_0\sqrt{t}$$

Then abbreviating  $N = N(\mathcal{F}, \gamma)$  and taking  $C_{ERM} = 36$ , with probability at least  $1 - \delta$ , 774

$$\sum_{(a,r)\in Z_t^{ser}} \left(\widehat{f}_{t+1}^{ser}(a) - f^*(a)\right)^2 \le C_{\text{ERM}} \bigg[ \lambda + (\gamma + \epsilon_0)^2 T + (\gamma + \epsilon_0) TR + R^2 (1 + M\alpha) \log(3N/\delta) \\ + R^2 M^2 \alpha \log(3NM/\delta) + R^2 M^2 \alpha \log(3NM/\delta) + R^2 M^2 \alpha \log(3NM/\delta) \bigg]$$

*Proof of Lemma D.3.* Let  $\mathcal{F}_{\gamma}$  be a  $\gamma$ -cover of the function class  $\mathcal{F}$  with respect to the infinity norm  $\|\cdot\|_{\infty}$ . For  $f \in \mathcal{F}$  and  $(a_t, r_t)$  for some  $t \in [T]$ , let

$$\phi(f, a_t, r_t) = -(f(a_t) - r_t)^2 + (f^*(a_t) - r_t)^2,$$

Since  $r_t = f^*(a_t) + \eta_t$ , we can write  $\phi(f, a_t, r_t)$  as 775

$$\phi(f, a_t, r_t) = -(f(a_t) - f^*(a_t) + \eta_t)^2 + \eta_t^2$$
  
= -2(f(a\_t) - f^\*(a\_t))\eta\_t - (f(a\_t) - f^\*(a\_t))^2

Since  $\eta_t$  is R-sub-Gaussian conditional on  $Z_{t-1}^{\text{all}}$ ,  $a_t$ , we have for any positive parameter  $\mu$  that 776

$$\begin{split} \log \mathbb{E} \big[ \exp(\mu \phi(f, a_t, r_t)) \big| Z_{t-1}^{\text{all}}, a_t \big] &\leq 2\mu^2 R^2 (f(a_t) - f^*(a_t))^2 - \mu (f(a_t) - f^*(a_t))^2 \\ &= (2\mu^2 R^2 - \mu) (f(a_t) - f^*(a_t))^2 \end{split}$$

Using Lemma D.1, we have with probability at least  $1 - \delta/3$ , for all  $f \in \mathcal{F}_{\gamma}$  and  $t \in [T]$ , 777

$$\mu_{\text{all}} \sum_{(a,r)\in Z_t^{\text{all}}} \phi(f,a,r) \le (2\mu_{\text{all}}^2 R^2 - \mu_{\text{all}}) \sum_{(a,r)\in Z_t^{\text{all}}} (f(a) - f^*(a))^2 + \log(3N/\delta),$$
(26)

- where  $\mu_{all} > 0$  is a parameter we will determine later. 778
- 779
- On the other hand, if we consider any local agent m, when  $m_t = m$ , we have  $\eta_t$  is R-sub-Gaussian conditional on  $Z_{m,t-1}^{up} \cup Z_{m,t-1}^{loc}$  and  $a_t$ , i.e. all the data agent m has received from the environment up to this point. Thus we have for any  $\mu > 0$  that 780
- 781

$$\log \mathbb{E} \Big[ \exp(-\mu\phi(f, a_t, r_t)) \Big| Z_{m,t-1}^{\text{up}} \cup Z_{m,t-1}^{\text{loc}}, a_t \Big] \le 2\mu^2 R^2 (f(a_t) - f^*(a_t))^2 + \mu (f(a_t) - f^*(a_t))^2 \\ = (2\mu^2 R^2 + \mu) (f(a_t) - f^*(a_t))^2$$

Then again using Lemma D.1 and taking summation on  $Z_{m,t}^{\text{loc}}$ , with probability at least  $1 - \delta/3$ , the following holds for any  $m \in [M]$ :

$$-\mu_{\rm loc} \sum_{(a,r)\in Z_{m,t}^{\rm loc}} \phi(f,a,r) \le (2\mu_{\rm loc}^2 R^2 + \mu_{\rm loc}) \sum_{(a,r)\in Z_{m,t}^{\rm loc}} (f(a) - f^*(a))^2 + \log(3NM/\delta), \quad (27)$$

where  $\mu_{\text{loc}} > 0$  is a parameter we will determine later.

Taking the summation of (27) for all  $m \in [M]$  and combining (26), while observing that  $Z_t^{\text{ser}} = Z_t^{\text{all}} \setminus \bigcup_{m=1}^M Z_{m,t}^{\text{loc}}$ , we get

$$\begin{split} \sum_{(a,r)\in Z_t^{\text{ser}}} \phi(f,a,r) &= \sum_{(a,r)\in Z_t^{\text{all}}} \phi(f,a,r) - \sum_{m=1}^M \sum_{(a,r)\in Z_{m,t}^{\text{loc}}} \phi(f,a,r) \\ &\leq (2\mu_{\text{all}}R^2 - 1) \sum_{(a,r)\in Z_t^{\text{all}}} (f(a) - f^*(a))^2 + \frac{1}{\mu_{\text{all}}} \log(3N/\delta) \\ &\quad + (2\mu_{\text{loc}}R^2 + 1) \sum_{m=1}^M \sum_{(a,r)\in Z_{m,t}^{\text{loc}}} (f(a) - f^*(a))^2 + \frac{1}{\mu_{\text{loc}}} M \log(3NM/\delta) \\ &= 2R^2(\mu_{\text{all}} + \mu_{\text{loc}}) \|f - f^*\|_{Z_t^{\text{all}}}^2 - (2\mu_{\text{loc}}R^2 + 1) \|f - f^*\|_{Z_t^{\text{ser}}}^2 \\ &\quad + \frac{1}{\mu_{\text{all}}} \log(3N/\delta) + \frac{1}{\mu_{\text{loc}}} M \log(3NM/\delta). \end{split}$$

From Lemma A.1, we have  $\lambda + \|f - f^*\|_{Z_t^{all}}^2 \leq (1 + M\alpha) \left(\lambda + \|f - f^*\|_{Z_t^{ser}}^2\right) \Leftrightarrow \|f - f^*\|_{Z_t^{all}}^2 \leq M\alpha\lambda + (1 + M\alpha)\|f - f^*\|_{Z_t^{ser}}^2$ . Plugging this inequality into the above and letting  $\mu_{all} = 1/8R^2(1 + M\alpha)$ and  $\mu_{loc} = 1/8R^2M\alpha$ , we get

$$\sum_{(a,r)\in Z_t^{\text{ser}}} \phi(f,a,r) \le 2R^2 (\mu_{\text{all}} + \mu_{\text{loc}}) M \alpha \lambda - (1 - 2M \alpha \mu_{\text{loc}} R^2 - 2(1 + M \alpha) \mu_{\text{all}} R^2) \|f - f^*\|_{Z_t^{\text{ser}}}^2 + \frac{1}{\mu_{\text{all}}} \log(3N/\delta) + \frac{1}{\mu_{\text{loc}}} M \log(3NM/\delta) \le -\frac{1}{2} \|f - f^*\|_{Z_t^{\text{ser}}}^2 + \frac{1}{2} \lambda + 8R^2 (1 + M \alpha) \log(3N/\delta) + 8R^2 M^2 \alpha \log(3NM/\delta)$$
(28)

Now for  $\hat{f}_{t+1}^{\text{ser}}$ , there exists  $\tilde{f} \in \mathcal{F}_{\gamma}$  such that  $\|\tilde{f} - \hat{f}_{t+1}^{\text{ser}}\|_{\infty} \leq \gamma$ . Using Lemma D.2, this gives us the following with probability at least  $1 - \delta/3$ :

$$\begin{split} -\sum_{(a,r)\in Z_t^{\text{ser}}} \phi(\tilde{f}, a, r) &= \sum_{(a,r)\in Z_t^{\text{ser}}} \left[ \left(\tilde{f}(a) - r\right)^2 - \left(f^*(a) - r\right)^2 \right] \\ &\leq \left( \sqrt{\sum_{(a,r)\in Z_t^{\text{ser}}} \left(\hat{f}_{t+1}^{\text{ser}}(a) - r\right)^2} + \sqrt{t\gamma^2} \right)^2 - \sum_{(a,r)\in Z_t^{\text{ser}}} \left(f^*(a) - r\right)^2 \\ &\leq \left( \sqrt{\sum_{(a,r)\in Z_t^{\text{ser}}} \left(f^*(a) - r\right)^2} + \sqrt{t}(\gamma + \epsilon_0) \right)^2 - \sum_{(a,r)\in Z_t^{\text{ser}}} \left(f^*(a) - r\right)^2 \\ &= (\gamma + \epsilon_0)^2 t + 2(\gamma + \epsilon_0) \sqrt{t} \left(\sum_{s=1}^t \eta_s^2\right)^{1/2} \\ &\leq (\gamma + \epsilon_0)^2 t + 2(\gamma + \epsilon_0) \sqrt{2T^2R^2 + 3TR^2 \log(3/\delta)}, \end{split}$$

where we used the basic inequality  $\sqrt{\sum(a+b)^2} \le \sqrt{\sum a^2} + \sqrt{\sum b^2}$  in the first inequality and used the property of  $\hat{f}_{t+1}^{\text{ser}}$  in the second inequality. Finally, taking a union bound and combining this with (28), we have with probability at least  $1 - \delta$ ,

$$\begin{split} &\sum_{(a,r)\in Z_t^{\text{ser}}} \left( \tilde{f}_{t+1}^{\text{ser}}(a) - f^*(a) \right)^2 \\ &\leq 2\gamma^2 t + 2\sum_{(a,r)\in Z_t^{\text{ser}}} \left( \tilde{f}(a) - f^*(a) \right)^2 \\ &\leq 2\gamma^2 t - 2\sum_{(a,r)\in Z_t^{\text{ser}}} \phi(\tilde{f},a,r) + \lambda + 32R^2(1+M\alpha)\log(3N/\delta) + 32R^2M^2\alpha\log(3NM/\delta) \\ &\leq 2\gamma^2 T + 2(\gamma+\epsilon_0)^2 T + 4(\gamma+\epsilon_0)\sqrt{2T^2R^2 + 3TR^2\log(3/\delta)} + \lambda + 32R^2(1+M\alpha)\log(3N/\delta) + 32R^2M^2\alpha\log(3NM/\delta) \\ &\leq C_{\text{ERM}} \left[ \lambda + (\gamma+\epsilon_0)^2 T + (\gamma+\epsilon_0)TR + R^2(1+M\alpha)\log(3N/\delta) + R^2M^2\alpha\log(3NM/\delta) \right], \end{split}$$

where the first inequality uses again  $\|\tilde{f} - \hat{f}_{t+1}^{\text{ser}}\|_{\infty} \leq \gamma$ , and it can be verified that the last inequality holds when  $C_{\text{ERM}} \geq 36$ .

# 797 NeurIPS Paper Checklist

849

798	1.	Claims
799		Ouestion: Do the main claims made in the abstract and introduction accurately reflect the
800		paper's contributions and scope?
801		Answer: [Ves]
001		Let' Castien We account he communication and the intermediate in the shot of
802 803		and elaborate further on motivation and main techniques in our introduction.
804		Guidelines:
805		• The answer NA means that the abstract and introduction do not include the claims
806		made in the paper.
807		• The abstract and/or introduction should clearly state the claims made, including the
808		contributions made in the paper and important assumptions and limitations. A No or
809		NA answer to this question will not be perceived well by the reviewers.
810		• The claims made should match theoretical and experimental results, and reflect how
811		much the results can be expected to generalize to other settings.
812		• It is fine to include aspirational goals as motivation as long as it is clear that these goals
813		are not attained by the paper.
814	2.	Limitations
815		Question: Does the paper discuss the limitations of the work performed by the authors?
916		Answer: [Vec]
010		Answei. [105]
817		Justification: All assumptions made for our theoretical analysis are present and stated
818		assumptions are also included
019		
820		Guidelines:
821		• The answer NA means that the paper has no limitation while the answer No means that
822		the paper has limitations, but those are not discussed in the paper.
823		• The authors are encouraged to create a separate "Limitations" section in their paper.
824		• The paper should point out any strong assumptions and how robust the results are to
825		violations of these assumptions (e.g., independence assumptions, noiseless settings,
826		should reflect on how these assumptions might be violated in practice and what the
828		implications would be
820		• The authors should reflect on the scope of the claims made e.g. if the approach was
830		only tested on a few datasets or with a few runs. In general, empirical results often
831		depend on implicit assumptions, which should be articulated.
832		• The authors should reflect on the factors that influence the performance of the approach.
833		For example, a facial recognition algorithm may perform poorly when image resolution
834		is low or images are taken in low lighting. Or a speech-to-text system might not be
835		used reliably to provide closed captions for online lectures because it fails to handle
836		technical jargon.
837		• The authors should discuss the computational efficiency of the proposed algorithms
838		and how they scale with dataset size.
839		• If applicable, the authors should discuss possible limitations of their approach to
840		address problems of privacy and fairness.
841		• While the authors might fear that complete honesty about limitations might be used by
842		reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aron't asknowledged in the paper. The suthers should use their best
843 844		infinations that aren't acknowledged in the paper. The authors should use their best individual actions in favor of transparency play an import
044 845		tant role in developing norms that preserve the integrity of the community Reviewers
846		will be specifically instructed to not penalize honesty concerning limitations.
847	3	Theory Assumptions and Proofs
040	5.	Question: For each theoretical result does the paper provide the full set of accumptions and
040 849		a complete (and correct) proof?
-		

850	Answer: [Yes]
851 852	Justification: All assumptions are listed in the main paper, while a very detailed and sound proof is displayed in the appendices.
853	Guidelines:
854	• The answer NA means that the paper does not include theoretical results.
855	• All the theorems, formulas, and proofs in the paper should be numbered and cross-
856	referenced.
857	• All assumptions should be clearly stated or referenced in the statement of any theorems.
858	• The proofs can either appear in the main paper or the supplemental material but if
859	they appear in the supplemental material, the authors are encouraged to provide a short
860	proof sketch to provide intuition.
861	• Inversely, any informal proof provided in the core of the paper should be complemented
862	by formal proofs provided in appendix or supplemental material.
863	• Theorems and Lemmas that the proof relies upon should be properly referenced.
864 4	. Experimental Result Reproducibility
96E	Question: Does the paper fully disclose all the information needed to reproduce the main ex-
866	perimental results of the paper to the extent that it affects the main claims and/or conclusions
867	of the paper (regardless of whether the code and data are provided or not)?
868	Answer: [NA]
869	Institution: Our theoretical paper does not present any experimental results
970	Guidelines:
070	• The ensure NA means that the paper does not include experiments
8/1	• The answer NA means that the paper does not include experiments.
872	• If the paper includes experiments, a no answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important regardless of
873	whether the code and data are provided or not
875	• If the contribution is a dataset and/or model, the authors should describe the steps taken
876	to make their results reproducible or verifiable.
877	• Depending on the contribution, reproducibility can be accomplished in various ways.
878	For example, if the contribution is a novel architecture, describing the architecture fully
879	might suffice, or if the contribution is a specific model and empirical evaluation, it may
880	be necessary to either make it possible for others to replicate the model with the same
881	dataset, or provide access to the model. In general, releasing code and data is often
882	instructions for how to replicate the results access to a hosted model (e.g. in the case
884	of a large language model) releasing of a model checkpoint or other means that are
885	appropriate to the research performed.
886	• While NeurIPS does not require releasing code, the conference does require all submis-
887	sions to provide some reasonable avenue for reproducibility, which may depend on the
888	nature of the contribution. For example
889	(a) If the contribution is primarily a new algorithm, the paper should make it clear how
890	to reproduce that algorithm.
891	(b) If the contribution is primarily a new model architecture, the paper should describe
892	the architecture clearly and fully.
893	(c) If the contribution is a new model (e.g., a large language model), then there should
894	either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open source dataset or instructions for how to construct
896	the dataset)
897	(d) We recognize that reproducibility may be tricky in some cases, in which case
898	authors are welcome to describe the particular way they provide for reproducibility
899	In the case of closed-source models, it may be that access to the model is limited in
900	some way (e.g., to registered users), but it should be possible for other researchers
901	to have some path to reproducing or verifying the results.
902 5	. Open access to data and code

903 904 905	Question: Does the paper provide open access to the data and code, with sufficient instruc- tions to faithfully reproduce the main experimental results, as described in supplemental material?
906	Answer: [NA]
300	Justification: Our theoretical paper does not present any experimental results
907	
908	Guidelines:
909	<ul> <li>The answer NA means that paper does not include experiments requiring code.</li> </ul>
910 911	• Please see the NeurIPS code and data submission guidelines (https://nips.cc/ public/guides/CodeSubmissionPolicy) for more details.
912	• While we encourage the release of code and data, we understand that this might not be
913	possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not
914 915	including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
916	• The instructions should contain the exact command and environment needed to run to
917	reproduce the results. See the NeurIPS code and data submission guidelines (https: (/ning_es/mublic/muideg/CodeSubmiggierDeligu) for more dataile
918	// https.cc/public/guides/codeSubmissionPolicy) for more details.
919 920	to access the raw data, preprocessed data, intermediate data, and generated data, etc.
921	• The authors should provide scripts to reproduce all experimental results for the new
922	proposed method and baselines. If only a subset of experiments are reproducible, they
923	should state which ones are omitted from the script and why.
924	• At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable)
925	<ul> <li>Providing as much information as possible in supplemental material (appended to the</li> </ul>
927	paper) is recommended, but including URLs to data and code is permitted.
928	6. Experimental Setting/Details
929	Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
930	parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
931	results?
932	Answer: [NA]
933	Justification: Our theoretical paper does not present any experimental results.
934	Guidelines:
935	<ul> <li>The answer NA means that the paper does not include experiments.</li> </ul>
936	• The experimental setting should be presented in the core of the paper to a level of detail
937	that is necessary to appreciate the results and make sense of them.
938 939	• The full details can be provided either with the code, in appendix, or as supplemental material.
940	7. Experiment Statistical Significance
941 942	Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?
943	Answer: [NA]
944	Justification: Our theoretical paper does not present any experimental results.
945	Guidelines:
946	• The answer NA means that the paper does not include experiments.
947	• The authors should answer "Yes" if the results are accompanied by error bars, confi-
948	dence intervals, or statistical significance tests, at least for the experiments that support
949	the main claims of the paper.
950	• The factors of variability that the error bars are capturing should be clearly stated (for
951	example, train/test split, initialization, random drawing of some parameter, or overall
952	• The method for calculating the error bars should be explained (closed form formula)
954	call to a library function, bootstrap, etc.)

955		The assumptions made should be given (e.g., Normany distributed errors).
956		• It should be clear whether the error bar is the standard deviation or the standard error
957		of the mean.
958		• It is OK to report 1-sigma error bars, but one should state it. The authors should
959		of Normality of errors is not verified
900		• For asymmetric distributions, the authors should be careful not to show in tables or
962		figures symmetric error bars that would vield results that are out of range (e.g. negative
963		error rates).
964		• If error bars are reported in tables or plots, The authors should explain in the text how
965		they were calculated and reference the corresponding figures or tables in the text.
966	8.	Experiments Compute Resources
967		Question: For each experiment, does the paper provide sufficient information on the com-
968		puter resources (type of compute workers, memory, time of execution) needed to reproduce
969		
970		Answer: [NA]
971		Justification: Our theoretical paper does not present any experimental results.
972		Guidelines:
973		<ul> <li>The answer NA means that the paper does not include experiments.</li> </ul>
974		• The paper should indicate the type of compute workers CPU or GPU, internal cluster,
975		or cloud provider, including relevant memory and storage.
976 077		• The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute
978		• The paper should disclose whether the full research project required more compute
979		than the experiments reported in the paper (e.g., preliminary or failed experiments that
980		didn't make it into the paper).
981	9.	Code Of Ethics
982 983		Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?
984		Answer: [Yes]
985		Justification: We have thoroughly reviewed the NeurIPS Code of Ethics, and believe our
986		work conforms fully the the Code
986		work conforms fully the the Code.
986 987		<ul><li>work conforms fully the the Code.</li><li>Guidelines:</li><li>The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics</li></ul>
986 987 988		<ul> <li>work conforms fully the the Code.</li> <li>Guidelines:</li> <li>The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.</li> <li>If the authors answer No, they should explain the special circumstances that require a</li> </ul>
986 987 988 989 990		<ul> <li>work conforms fully the the Code.</li> <li>Guidelines:</li> <li>The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.</li> <li>If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.</li> </ul>
986 987 988 989 990 991		<ul> <li>work conforms fully the the Code.</li> <li>Guidelines: <ul> <li>The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.</li> <li>If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.</li> <li>The authors should make sure to preserve anonymity (e.g., if there is a special consid-</li> </ul> </li> </ul>
986 987 988 989 990 991 992		<ul> <li>work conforms fully the the Code.</li> <li>Guidelines: <ul> <li>The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.</li> <li>If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.</li> <li>The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).</li> </ul> </li> </ul>
986 987 988 989 990 991 992 993	10.	<ul> <li>work conforms fully the the Code.</li> <li>Guidelines: <ul> <li>The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.</li> <li>If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.</li> <li>The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).</li> </ul> </li> <li>Broader Impacts</li> </ul>
986 987 988 989 990 991 992 993 994 995	10.	<ul> <li>work conforms fully the the Code.</li> <li>Guidelines: <ul> <li>The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.</li> <li>If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.</li> <li>The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).</li> </ul> </li> <li>Broader Impacts <ul> <li>Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?</li> </ul> </li> </ul>
986 987 988 989 990 991 992 993 993 994 995 996	10.	<ul> <li>work conforms fully the the Code.</li> <li>Guidelines: <ul> <li>The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.</li> <li>If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.</li> <li>The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).</li> </ul> </li> <li>Broader Impacts <ul> <li>Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?</li> </ul> </li> </ul>
986 987 988 990 991 992 993 993 994 995 996	10.	<ul> <li>work conforms fully the the Code.</li> <li>Guidelines: <ul> <li>The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.</li> <li>If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.</li> <li>The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).</li> </ul> </li> <li>Broader Impacts <ul> <li>Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?</li> <li>Answer: [Yes]</li> </ul> </li> <li>Justification: See Impact Statement before appendices</li> </ul>
986 987 988 989 990 991 992 993 993 994 995 996 997 998	10.	<ul> <li>work conforms fully the the Code.</li> <li>Guidelines: <ul> <li>The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.</li> <li>If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.</li> <li>The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).</li> </ul> </li> <li>Broader Impacts <ul> <li>Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?</li> </ul> </li> <li>Answer: [Yes] <ul> <li>Justification: See Impact Statement before appendices</li> <li>Guidelines:</li> </ul> </li> </ul>
986 987 988 989 990 991 992 993 993 994 995 996 997 998	10.	<ul> <li>work conforms fully the the Code.</li> <li>Guidelines: <ul> <li>The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.</li> <li>If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.</li> <li>The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).</li> </ul> </li> <li>Broader Impacts <ul> <li>Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?</li> </ul> </li> <li>Answer: [Yes] Justification: See Impact Statement before appendices Guidelines: <ul> <li>The answer NA means that there is no societal impact of the work performed.</li> </ul> </li> </ul>
986 987 988 989 990 991 992 993 993 994 995 996 997 998 999	10.	<ul> <li>work conforms fully the the Code.</li> <li>Guidelines: <ul> <li>The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.</li> <li>If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.</li> <li>The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).</li> </ul> </li> <li>Broader Impacts <ul> <li>Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?</li> </ul> </li> <li>Answer: [Yes] <ul> <li>Justification: See Impact Statement before appendices</li> <li>Guidelines: <ul> <li>The answer NA means that there is no societal impact of the work performed.</li> <li>If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.</li> </ul> </li> </ul></li></ul>
986 987 988 989 990 991 992 993 993 994 995 996 997 998 999 999 1000 1001	10.	<ul> <li>work conforms fully the the Code.</li> <li>Guidelines: <ul> <li>The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.</li> <li>If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.</li> <li>The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).</li> </ul> </li> <li>Broader Impacts <ul> <li>Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?</li> <li>Answer: [Yes]</li> </ul> </li> <li>Justification: See Impact Statement before appendices</li> <li>Guidelines: <ul> <li>The answer NA means that there is no societal impact of the work performed.</li> <li>If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.</li> <li>Examples of negative societal impacts include potential malicious or unintended uses</li> </ul> </li> </ul>
986 987 988 989 990 991 992 993 993 994 995 996 997 998 997 998 999 1000 1001	10.	<ul> <li>work conforms fully the the Code.</li> <li>Guidelines: <ul> <li>The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.</li> <li>If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.</li> <li>The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).</li> </ul> </li> <li>Broader Impacts <ul> <li>Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?</li> </ul> </li> <li>Answer: [Yes] <ul> <li>Justification: See Impact Statement before appendices</li> <li>Guidelines:</li> <li>The answer NA means that there is no societal impact of the work performed.</li> <li>If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.</li> <li>Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations</li> </ul> </li> </ul>

1020   Inecuback over time, improving the enciency and accessionity of ML).     1021   11. Safeguards	
1021 11. Sateguards	
Ougstion, Dogs the non-an describe sefery and that have been put in place for responsible	10
release of data or models that have a high risk for misuse (e.g. pretrained language mode	s
image generators, or scraped datasets)?	5,
1025 Answer: [NA]	
Institution: Our theoretical paper does not present any experimental results, and thus do	26
not feature such data and models.	/0
1028 Guidelines:	
• The answer NA means that the paper poses no such risks	
• Released models that have a high risk for misuse or dual-use should be released wi	h
necessary safeguards to allow for controlled use of the model, for example by requiri	g
that users adhere to usage guidelines or restrictions to access the model or implementi	g
1033 safety filters.	-
• Datasets that have been scraped from the Internet could pose safety risks. The author	ſS
should describe how they avoided releasing unsafe images.	
• We recognize that providing effective safeguards is challenging, and many papers	0
1037 not require this, but we encourage authors to take this into account and make a be	st
12 Licenses for existing essets	
1039 12. Licenses for existing assets	
1040 Question: Are the creators or original owners of assets (e.g., code, data, models), used	n d
1042 properly respected?	u
1043 Answer: [NA]	
1044 Institution: Our paper does not include any such assets	
1045 Guidelines:	
• The answer NA means that the paper does not use existing assets	
• The authors should gits the original paper that produced the code package or datase	
• The authors should che the original paper that produced the code package of datase	•
1048 • The authors should state which version of the asset is used and, it possible, include	a
• The name of the license (e.g., CC-BY 4.0) should be included for each asset.	
• For scraped data from a particular source (e.g., website), the convright and terms	of
service of that source should be provided.	-
• If assets are released, the license, copyright information, and terms of use in t	e
1054 package should be provided. For popular datasets, paperswithcode.com/dataset	s
has curated licenses for some datasets. Their licensing guide can help determine t	e
1056 license of a dataset.	c
• For existing datasets that are re-packaged, both the original license and the license the derived asset (if it has changed) should be provided.	)t

1059 1060		• If this information is not available online, the authors are encouraged to reach out to the asset's creators.
1061	13.	New Assets
1062 1063		Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?
1064		Answer: [NA]
1065		Justification: Our theoretical paper does not introduce any new assets.
1066		Guidelines:
1067		• The answer NA means that the paper does not release new assets
1069		<ul> <li>Researchers should communicate the details of the dataset/code/model as part of their</li> </ul>
1069		submissions via structured templates. This includes details about training, license,
1070		limitations, etc.
1071 1072		• The paper should discuss whether and how consent was obtained from people whose asset is used.
1073 1074		• At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.
1075	14.	Crowdsourcing and Research with Human Subjects
1076		Question: For crowdsourcing experiments and research with human subjects does the paper
1070		include the full text of instructions given to participants and screenshots, if applicable, as
1078		well as details about compensation (if any)?
1079		Answer: [NA]
1080		Justification: Our theoretical paper does not present any experimental results.
1081		Guidelines:
1082		• The answer NA means that the paper does not involve crowdsourcing nor research with
1083		human subjects.
1084		• Including this information in the supplemental material is fine, but if the main contribu-
1085		tion of the paper involves human subjects, then as much detail as possible should be included in the main paper
1087		<ul> <li>According to the NeurIPS Code of Ethics, workers involved in data collection, curation</li> </ul>
1088		or other labor should be paid at least the minimum wage in the country of the data
1089		collector.
1090	15.	Institutional Review Board (IRB) Approvals or Equivalent for Research with Human
1091		Subjects
1092		Question: Does the paper describe potential risks incurred by study participants, whether
1093		such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
1094		institution) were obtained?
1096		Answer: [NA]
1097		Justification: Our theoretical paper does not present any experimental results
1098		Guidelines:
1099		• The answer NA means that the paper does not involve crowdsourcing nor research with
1100		human subjects.
1101		• Depending on the country in which research is conducted, IRB approval (or equivalent)
1102		may be required for any human subjects research. If you obtained IRB approval, you
1103		should clearly state this in the paper.
1104		• We recognize that the procedures for this may vary significantly between institutions and leastions, and we are at authors to adhere to the NeurIDS Code of Ethics and the
1105		guidelines for their institution.
1107		• For initial submissions, do not include any information that would break approximity (if
1108		applicable), such as the institution conducting the review.