# PROVABLE LEARNING FOR DEC-POMDPS: FACTORED MODELS AND MEMORYLESS AGENTS

Anonymous authors

Paper under double-blind review

# Abstract

This paper studies cooperative Multi-Agent Reinforcement Learning (MARL) under the mathematical model of Decentralized Partially Observable Markov Decision Process (DEC-POMDP). Despite the empirical success of cooperative MARL, its theoretical foundation, particularly in the realm of provable learning of DEC-POMDPs, remains limited. In this paper, we first present a hardness result in theory demonstrating that, without additional structural assumptions, learning DEC-POMDPs requires several samples that grows exponentially with the number of agents in the worst case, which is also known as the curse of multiagency. This motivates us to explore important subclasses of DEC-POMDPs for which efficient solutions can be found. Specifically, we propose new algorithms and establish sample-efficiency guarantees that break the curse of multiagency, for finding both local and global optima in two important scenarios: (1) when agents employ memoryless policies, selecting actions based solely on their current observations; and (2) when a factored structure is present, which enables key properties similar to value decomposition in VDN or Qmix.

024 025 026

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

## 1 INTRODUCTION

027 028

029 Cooperative multi-agent reinforcement learning (MARL) has gained significant attention due to its wide range of applications in real-world problems, such as autonomous vehicles and robotic swarms. However, one of the primary challenges in cooperative MARL is the exponential growth of the ac-031 tion space as the number of agents increases. Consequently, the number of samples required to learn an optimal policy grows exponentially with the number of agents. To address this issue, Oliehoek 033 et al. (2008a); Kraemer and Banerjee (2016) proposed the method of centralized training with de-034 centralized execution (CTDE). In CTDE, agents' policies are trained using global information, but the resulting policies are decentralized, allowing each agent to execute them using only local information during execution. Rashid et al. (2018) introduced the decentralized partially observable 037 Markov decision process (DEC-POMDP) as a model for fully cooperative multi-agent reinforce-038 ment learning. The DEC-POMDP model is similar to a traditional partially observable Markov game, with the key distinction that in a DEC-POMDP, the objective is to maximize the cumulative rewards across all agents. Although cooperative MARL has achieved empirical success across a 040 variety of real-world problems, and numerous significant algorithms have been developed—such as 041 VDN (Sunehag et al., 2017), Q-MIX (Rashid et al., 2018), and Q-PLEX (Wang et al., 2021)-the 042 theoretical foundations of MARL remain underdeveloped. Therefore, the objective of this paper is 043 to develop a sample-efficient algorithm for DEC-POMDPs, with an emphasis on theoretical rigor 044 and comprehensive guarantees.

This paper considers two types of optimality: (1) **global optimality**, where agents seek policies that maximize the total cumulative reward across all agents; and (2) **local optimality**, where each agent aims to find a policy that maximizes total cumulative reward, assuming the policies of all other agents are fixed. The latter concept is commonly referred to as the Nash equilibrium (Kreps, 1989) in game theory.

The challenge of learning in DEC-POMDPs arises from two primary factors. First, agents have
 access only to their individual trajectory histories, making it difficult to identify either global or
 local optima, both of which require cooperation among agents. Second, the joint action and observation space grows exponentially with the number of agents, creating significant obstacles in

designing algorithms with polynomial sample complexity, particularly as the number of agents in creases. We illustrate these challenges by first proving statistical hardness results for learning
 DEC-POMDPs without additional structural assumptions (Theorem 5.2). This motivates us to focus
 on important subclasses of DEC-POMDPs that are rich enough to encompass practical applications
 but constrained enough to admit sample-efficient algorithmic solutions. Specifically, we consider
 the following two subclasses:

060 Factored Structure Model: Motivated by the value decomposition approach, which has been 061 widely employed in practical algorithms such as VDN and Qmix, we consider DEC-POMDPs with 062 a factored structure that enables key properties similar to value decomposition. We demonstrate that 063 a factored value decomposition property holds for our model and develop algorithms with provable 064 efficiency to achieve both local and global optimality. Specifically, the value decomposition-like conditions in our model allow the total action-value function to be decomposed into distinct com-065 ponents, each dependent only on the trajectory of one or a small subset of agents. As a result, 066 each agent only needs to focus on a limited number of components when making decisions, which 067 mitigates the exponential scaling of complexity as the number of agents increases. 068

DEC-POMDP with Memoryless Policy: We first investigate a model frequently studied in partially
 observable settings (Kara and Yuksel, 2022; Kara and Yüksel, 2023), where agents determine their
 actions based solely on their current observations. We propose a sample-efficient algorithm tailored
 to achieve local optimality in this context.

Our Contribution: Our contributions are centered around the development of a sample-efficient algorithm under reasonable conditions. We summarize our key contributions and results as follows:

- 1. We demonstrate that, without further assumptions, designing a sample-efficient algorithm to achieve global optimality in DEC-POMDPs is infeasible, regardless of whether agents employ general policies or are restricted to memoryless policies.
  2. In the setting where agents use memoryless policies.
  - 2. In the setting where agents use memoryless policies, we propose a sample-efficient algorithm that achieves local optimality.
    - 3. For the factored structure model, we prove that the value function can be decomposed into components that depend only on the trajectories of a few agents rather than all agents. We also introduce a sample-efficient algorithm that can achieve both local and global optimality.
    - 4. Our analysis of the factored structure model provides a partial theoretical explanation for empirical algorithms such as VDN, as we establish a sufficient condition for value decomposition and offer a sample-efficient guarantee under this condition.
- 085 086 087

088

079

080

081

082

083

084

# 2 RELATED WORK

Due to space limits, we briefly present a few previous works closely related to this paper and leave the comprehensive discussion on additional related work in the appendix.

091 Learning POMDPs planning in POMDPs is known to be PSPACE-complete (Papadimitriou and 092 Tsitsiklis, 1987; Littman, 1994; Burago et al., 1996; Lusena et al., 2001). Uehara et al. (2022) 093 impose a deterministic latent transition assumption on the model and design computationally effi-094 cient algorithms. Jin et al. (2020) design the observable operator model with the upper confidence 095 bound algorithm for weakly revealing POMDPs, while Liu et al. (2022a) propose the optimistic 096 maximum likelihood estimation (OMLE) algorithm for learning weakly revealing POMDPs. Chen 097 et al. (2022) derive a unified analysis for OMLE with a sharper sample complexity. Furthermore, 098 Liu et al. (2023) provides a generic framework for applying OMLE to a wide range of partially observable problems, including low-rank sequential decision-making problems and general sequential 099 decision-making problems under the SAIL condition. Since OMLE learns the near-optimal poli-100 cies of an enormously rich class of sequential decision-making problems in a polynomial number of 101 samples, we also build our work upon the generic framework of OMLE. 102

Learning DEC-POMDPs Rashid et al. (2018) introduced the decentralized-partially observable
 Markov decision process (DEC-POMDP) as a fully cooperative multiagent reinforcement learning
 task. Empirical algorithms for solving DEC-POMDP include VDN (Sunehag et al., 2017), Q-MIX
 (Rashid et al., 2018), and Q-PLEX (Wang et al., 2021). Recent works on DEC-POMDP, such as
 those by Hu and Foerster (2019), Foerster et al. (2019), and Lerer et al. (2020), adopt ideas similar to
 the common information approach, leading to breakthroughs in challenging DEC-POMDP problems

like Hanabi. The common information approach of Dec-POMDP has been studied theoretically in (Zhang et al., 2019; Mao et al., 2023; Liu and Zhang, 2023). In particular, (Liu and Zhang, 2023) establish a sample quasi-efficient algorithm with quasi-polynomial sample complexity. In comparison, by imposing additional structural assumptions, our algorithms have polynomial sample complexity. Moreover, our lower bound shows that polynomial sample complexity is impossible without additional assumptions.

114 115

116

# **3** PROBLEM FORMULATION

# 117 3.1 PRELIMINARY

Partially Observable Model: We study the decentralized partially observable Markov deci-119 sion process (DEC-POMDP)(Rashid et al., 2018) with n agents, which is denoted by a tuple 120  $(\mathcal{S}, \mathcal{O}, \mathcal{A}; H, \mu_1, \mathbb{T}, \mathbb{O}; r)$ . Here  $\mathcal{S}$  is the set of all possible states, where the states are not observ-121 able by the agents. For each agent  $m \in [n]$ , we let  $\mathcal{A}_m$  and  $\mathcal{O}_m$  denote the action and observation 122 space of agent m respectively. We define the joint action space and the joint observation space by 123  $\mathcal{A} := \mathcal{A}_1 \times \cdots \times \mathcal{A}_n$  and  $\mathcal{O} := \mathcal{O}_1 \times \cdots \times \mathcal{O}_n$ , respectively. Besides, H is the episode length, 124  $\mu_1$  is the distribution of the initial state  $s_1$ ,  $\mathbb{T} = \{\mathbb{T}_{h,\mathbf{a}}\}_{(h,\mathbf{a})\in[H-1]\times\mathcal{A}}$  is the Markov state transition 125 kernel,  $\mathbb{O} = \{\mathbb{O}_{h,m}\}_{h \in [H], m \in [n]}$  is the observation emission kernel, and  $r = \{r_{h,m}\}_{h \in [H], m \in [n]}$ 126 denotes the reward function. In particular, starting from  $s_1 \sim \mu_1$ , at each step  $h \in [H]$ , each agent 127 m observes an observation  $o_{h,m} \in \mathcal{O}_m$  according to distribution  $\mathbb{O}_{h,m}(\cdot | s_h)$ , takes an action  $a_{h,m}$ , 128 and receives a reward  $r_{h,m}(o_{h,m}) \in [0,1]$ , which is a function of  $o_{h,m}$ . We consider this type 129 of reward function to prevent information about latent states from leaking through rewards beyond 130 what is provided by the observations. Let  $\mathbf{a}_h = (a_{h,1}, a_{h,2}, \dots, a_{h,n})$  denote the joint action at step 131 h. The next state  $s_{h+1}$  is sampled from distribution  $\mathbb{T}_{h,\mathbf{a}_h}(\cdot | s_h)$ . Such a process terminates when  $s_{H+1}$  is reached. For each agent m, it collects data  $\{o_{h,m}, a_{h,m}, r_{h,m}\}_{h \in [H]}$ . Notably, each agent 132 has access only to their own observation and reward, and therefore does not know the total payoff. 133 The goal of the agents is to maximize the social welfare, i.e. the summation of cumulative rewards 134 obtained by all the *n* agents. Moreover, to simplify the notation, we let S, A, and O denote |S|, 135  $\max_{m \in [n]} |\mathcal{A}_m|$ , and  $\max_{m \in [n]} |\mathcal{O}_m|$ , respectively. 136

**DEC-POMDP:** In this work, we study the case where agents adopt decentralized policies. Namely, each agent only selects actions based on the history of her own trajectories. The joint policy  $\pi_h$ represents the decentralized policy product of the *n* agents. We formally denote the policy class as  $\pi = \{\{\bigotimes_{m=1}^{n} \pi_{h,m}\}_{h \in [H]} \mid \pi_{h,m} : (\mathcal{O}_m \times \mathcal{A}_m)^{h-1} \times \mathcal{O}_m \to \Delta_m\}$ . When considering a product policy  $\pi$ , we denote its value function as  $V^{\pi}$ , defined as the expected total reward received by all agents under policy  $\pi$ :

143

156

161

$$V^{\boldsymbol{\pi}} = \mathbb{E}_{\boldsymbol{\pi}} \Big[ \sum_{h=1}^{H} \sum_{m=1}^{n} r_{h,m}(o_{h,m}) \Big].$$

144 145  $\forall \boldsymbol{\tau}_h = (\mathbf{o}_1, \mathbf{a}_1, \dots, \mathbf{a}_h), \text{ we define the } Q\text{-function as } Q^{\boldsymbol{\pi}}(\boldsymbol{\tau}_h) = \mathbb{E}_{\boldsymbol{\pi}}[\sum_{j=h}^{H} \sum_{m=1}^{n} r_{j,m}(o_{j,m}) \mid \boldsymbol{\tau}_h].$ 146 For each agent  $m \in [n]$  and step  $h \in [H]$ , we define the trajectory notation of the  $m^{\text{th}}$  agent as 147  $\tau_{h,m} = (o_{1,m}, a_{1,m}, \dots, o_{h,m}, a_{h,m}).$ 

Learning Target: We define two types of optimality as learning objectives for DEC-POMDPs:global optimality and local optimality. Their formal definitions are provided as follows:

**Definition 3.1** (Local Optimality). A policy  $\pi = \pi_1 \times \pi_2 \times \cdots \times \pi_n$  is considered a local optimal policy for agents  $1, 2, \ldots, n$  if it satisfies  $V^{\pi} = \max_{i \in [n]} [\max_{\pi'_i} V^{\pi'_i, \pi_{-i}}].$ 

Our aim is to minimize the number of samples required to obtain an  $\epsilon$ -approximate local optimal policy. We define an  $\epsilon$ -approximate local optimal policy as a policy  $\pi = \pi_1 \times \pi_2 \times \cdots \times \pi_n$  that satisfies

$$V^{\boldsymbol{\pi}} \ge \max_{i \in [n]} \left[ \max_{\pi'_i} V^{\pi'_i, \pi_{-i}} \right] - \epsilon.$$

**Definition 3.2** (Global Optimality). A policy  $\pi = \pi_1 \times \pi_2 \times \cdots \times \pi_n$  is deemed a global optimal policy for agents  $1, 2, \ldots, n$  if it satisfies  $V^{\pi} = \max_{\pi'_1, \pi'_2, \ldots, \pi'_n} V^{\pi'_1 \times \pi'_2 \times \cdots \times \pi'_n}$ .

160 Our goal is to minimize the regret, which is defined as

$$\operatorname{Regret}(T) = \sum_{k=1}^{T} (V^{\pi^*} - V^{\pi^k})$$

where  $\pi^* = \arg \max_{\pi} V^{\pi}$ . We assume that agents interact with DEC-POMDPs for T episodes, and in the k-th iteration for any  $k \in [T]$ , they follow the policy  $\pi^k = \pi_1^k \times \pi_2^k \times \cdots \times \pi_n^k$ . Similarly, we can define an  $\epsilon$ -approximate global optimal policy in the local context, where a uniform mixture of  $\pi^1, \ldots, \pi^T$  satisfies the definition when sub-linear regret is achieved.

Weakly Revealing Condition: Liu et al. (2022a) demonstrated that without any assumption on the model, there exist hard instances such that the number of samples required to learn an  $\epsilon$ -approximate optimal policy in single-agent POMDPs is exponential in the horizon length *H*. Given the difficulty of learning POMDPs without assumptions on the model, even in single-agent settings, we consider the weakly revealing condition. This assumption is commonly adopted in previous works on partially observable contextual settings (Jin et al., 2020; Liu et al., 2022a; Chen et al., 2022).

**Assumption 3.1.** We define  $\mathbb{O}_h^i$  as  $\mathbb{O}_h^i(o_{h,i} | s_h) = \sum_{\{o_{h,j}\}_{j \in [n]/\{i\}}} \mathbb{O}_h(\mathbf{o}_h | s_h)$ . There exist  $\alpha > 0$ such that  $\min_{h,i} \sigma_S(\mathbb{O}_h^i) \ge \alpha$ , where for matrix  $\mathbf{A}$  we use  $\sigma_S(\mathbf{A})$  to denote the  $S^{th}$  singular value of emission matrix  $\mathbf{A}$ .

This condition guarantees that, with enough samples, the observations provide adequate information
 to differentiate between any two combinations of states.

Additional Notations: Throughout this paper, we adopt the following notation for sets of elements with subscripts: Let  $\mathcal{R} = \{x_i\}_{i \in S}$ , where S denotes the set of subscripts of the elements in  $\mathcal{R}$ . For simplicity, we represent  $\mathcal{R}$  as  $\mathcal{R} = x_S$ .

181 182 183

### 3.2 HARDNESS RESULT FOR GENERAL DEC-POMDP

The following theorem demonstrates that, in the absence of specific assumptions on the model,
 achieving global optimality in DEC-POMDPs is not possible with a sample complexity that is not
 exponential in the number of agents.

**Theorem 3.1.** For any randomized or deterministic algorithms, there exists an instance of DEC-MDP wherein the regret scales at least as  $\Omega(\sqrt{A^nT})$ .

This result highlights the limitations of achieving sample efficiency in algorithms for DEC-POMDPs
without making assumptions about the transition model. Consequently, in the following sections,
we aim to develop a sample-efficient algorithm for DEC-POMDPs under reasonable assumptions
about the model.

194 195

196 197

203

# 4 LEARNING DEC-POMDP WITH FACTORED STRUCTURE MODEL

4.1 FACTORED STRUCTURE MODEL

In this section, we consider a factored structure model, where the state space is decomposed as the Cartesian product of *n* individual spaces,  $S = S_1 \times S_2 \times \cdots \times S_n$ . For all  $\mathbf{s}' = (s'_1, \ldots, s'_n) \in S$ ,  $\mathbf{s} = (s_1, \ldots, s_n) \in S$ ,  $\mathbf{a} \in \mathcal{A}$ ,  $\mathbf{o} \in \mathcal{O}$ , and  $h \in [H]$ , the observation distribution and transition probability are factorized as:

$$\mathbb{O}_h(\mathbf{o} \mid \mathbf{s}) = \prod_{m=1}^n \mathbb{O}_{h,m}(o_m \mid s_m), \quad \mathbb{T}_h(\mathbf{s}' \mid \mathbf{s}, \mathbf{a}) = \prod_{m=1}^n \mathbb{T}_{h,m}(s'_m \mid s_m, a_m, a_{\mathsf{pa}(m)}),$$

where  $pa(m) \subset [n]$  represents the set of agents whose actions influence the transition of agent m. We further define  $\overline{pa}(m) = pa(m) \cup \{m\}$ . We assume that the local state transition of each individual agent depends solely on the actions of other agents, with no dependency between the states of different agents.

209 To represent the correlation between different agents, we introduce the following influence graph:

**210 Definition 4.1.** We define a directed graph  $G = (\mathcal{V}, \mathcal{E})$ , where  $\mathcal{V} = \{1, 2, ..., n\}$ , and there is 211 a directed edge from vertex *i* to vertex *j* for distinct vertices  $i, j \in \mathcal{V}$  if and only if  $i \in pa(j)$ . 212 Additionally, we assume that the maximum indegree of the graph G is d.

- Additionally, for clarity in presentation, we introduce several notations from graph theory:
- **Definition 4.2.** For each agent  $m \in [n]$ , we define two sets: the children set ch(m) and the ancestor set an(m). Specifically, for a vertex  $i \in [n]$ , if there exists a directed path  $i = j_1, j_2, ..., j_l = m$

222

223 224

225 226 227

228

238 239

240 241

242

243

244

245

246

247

248

249

250

253

255

256

257

259 260

261 262

264 265

266

267

216 in the influence graph G, where there is a directed edge from  $j_r$  to  $j_{r+1}$  for all  $r \in [l-1]$ , then 217  $i \in an(m)$  and  $m \in ch(i)$ . Moreover, for all  $m \in [n]$ , we define  $ch(m) = \{m\} \cup ch(m)$  and 218  $\overline{pa}(m) = \{m\} \cup pa(m)$ . We also define the complement set of ch(m) as  $nch(m) = [n] \setminus ch(m)$ . 219

220 The factored structure of the model leads to the following property of value decomposition.

**Proposition 4.1.** (Value Decomposition) For all  $h \in [H]$  and trajectory  $\tau_h = (\mathbf{o}_1, \mathbf{a}_1, \dots, \mathbf{o}_h, \mathbf{a}_h)$ , the Q-function can be decomposed as follows:

$$Q^{\pi}(\boldsymbol{\tau}_h) = \sum_{m=1}^n Q_m(\boldsymbol{\tau}_{h,\overline{\mathtt{an}}(m)})$$

where we define  $Q_m(\tau_{h,\overline{an}(m)}) = \mathbb{E}_{\pi_{\overline{an}(m)}}[\sum_{j=h}^{H} r_{j,m}(o_{j,m}) \mid \tau_{h,\overline{an}(m)}]$ . In other words, the value function can be expressed as the sum of n terms, where the m-th component depends only on the trajectory of the agents in  $\overline{an}(m)$ .

For each  $i \in [n]$ , we denote  $\theta_i = (\mathbb{T}_i, \mathbb{O}_i, \mu_i)$  as the collection of parameters representing the 229 transition and observation models of the *i*-th agent. We further use  $\Theta_i$  to denote the set of all 230 possible model parameters  $\theta_i$ . According to the factored structure condition, the joint trajectory 231 probability can be rewritten as the product of individual trajectory operators, where the individual 232 operator  $\mathbb{P}_{\theta_i}^{\pi_i}(\tau_{H,i} \mid \tau_{H,\mathsf{pa}(i)})$  is defined as: 233

$$\mathbb{P}_{\theta_{i}}^{\pi_{i}}(\tau_{H,i} \mid \tau_{H,\mathsf{pa}(i)}) = \sum_{s_{[H],i}} \mu(s_{1,i}) \mathbb{O}_{1,i}(o_{1,i} \mid s_{1,i}) \pi_{1,i}(a_{1,i} \mid o_{1,i}) \\
\cdot \left[ \prod_{h=1}^{H-1} \mathbb{T}_{h,i}(s_{h+1,i} \mid s_{h,i}, a_{h,\overline{\mathsf{pa}}(i)}) \mathbb{O}_{h+1,i}(o_{h+1,i} \mid s_{h+1,i}) \pi_{h+1,i}(a_{h+1,i} \mid \tau_{h-1,i}, o_{h,i}) \right].$$
(1)

We further use  $\{\theta_i^*\}_{i \in [n]}$  to denote the model parameters of the true transition model.

### 4.2 ACHIEVING GLOBAL OPTIMALITY WITH FACTORED STRUCTURE

In this section, we introduce a sample-efficient algorithm (outlined in Algorithm 4.2) to achieve global optimality under the factored structure model.

**Algorithm Description:** Algorithm 4.2 consists of three main steps in each episode k:

- Update policy and parameters (Line 3): We construct n distinct confidence intervals, where the *i*-th confidence interval contains only the parameters of the *i*-th agent's transition and observation model. The total value function is considered as a function of the joint product policy of all agents and the model parameters. We select model parameters  $\theta_i^k$  from the *i*-th confidence interval, along with a joint product policy  $\pi^k$ , such that the value function is maximized. After selecting the policy, we iteratively execute the following two steps for each agent  $i \in [n]$  to collect samples and update the confidence intervals.
- Sample trajectories for Agent i (Lines 6-8): We sample trajectories for different agents according to two distinct distributions. At step h, all agents initially select an action according to their policies. For agent  $m \in [i]$ , an observation sample is directly collected from the true model. For 254 the remaining agents  $m \in [n] \setminus [i]$ , we denote  $\mathbb{T}_{h,m}^k$  and  $\mathbb{O}_{h,m}^k$  as the transition and observation models corresponding to the parameter  $\theta_k$ . Given that model  $\mathbb{T}_{h,m}^k$  and  $\mathbb{O}_{h,m}^k$  are known, agent m samples  $s_{h+1,m} \sim \mathbb{T}_{h,m}^k (\cdot \mid s_{h,m})$ . Subsequently, agent m collects an observation  $o_{h+1,m} \sim$  $\mathbb{O}_{h+1,m}^k(\cdot \mid s_{h+1,m})$  and stores a dummy state  $s_{h+1,m}$  for exploration in the next episode. 258
  - Update confidence interval for agent *i* (Line 10): After collecting the trajectories  $\tau_m^k = (o_{1,m}^k, \dots, o_{H,m}^k, a_{H,m}^k)$  for all  $m \in [n]$ , we add the tuple  $(\pi_i^k, \tau_{\overline{pa}(i)}^k)$  to the sample set  $\mathcal{D}_i$ . The confidence set is then updated according to:

$$\mathcal{B}_{i}^{k+1} = \left\{ \hat{\theta}_{i} \in \mathcal{B}_{i}^{1} : \sum_{(\pi_{i}, \tau_{\overline{\mathsf{ps}}(i)}) \in \mathcal{D}_{i}} \log \mathbb{P}_{\hat{\theta}_{i}}^{\pi_{i}}(\tau_{i} \mid \tau_{\mathsf{pa}(i)}) \geq \max_{\theta_{i}' \in \Theta_{i}} \sum_{(\pi_{i}, \tau_{\overline{\mathsf{ps}}(i)}) \in \mathcal{D}_{i}} \log \mathbb{P}_{\theta_{i}'}^{\pi_{i}}(\tau_{i} \mid \tau_{\mathsf{pa}(i)}) - \beta_{i} \right\}$$

$$(2)$$

That is, we include those model parameters  $\theta_m$  for which the total log-likelihood assigned to the data is close to the maximum possible total log-likelihood.

**Technical Challenge and Insights:** The dimensionality of the model grows as  $\Omega(A^n O^n)$ , since 268 the joint action and observation spaces expand exponentially with the number of agents. Conse-269 quently, the sample complexity for estimating the model parameters is susceptible to  $\Omega(A^n O^n)$ . To

289

291

293 294

295

296

297

298

299

300

301

302

307

308

270 address this challenge, we construct separate confidence intervals for estimating the different model 271 parameters. This approach mitigates the exponential sample complexity in n, as the dimension of 272 each parameter  $\theta_m$  does not increase exponentially with n. Additionally, we employ a carefully 273 designed sampling procedure (outlined in lines 6-8) instead of directly sampling from the true tran-274 sition model. This enables us to precisely control the statistical error in estimating the joint trajectory probability by separately managing the statistical error of each individual trajectory operator. 275

Algorithm 1 OMLE for Achieving Global Optimal in Factored Structure Model 277 278 1: Initialize:  $\mathcal{B}_m^1 = \{\hat{\theta}_m \in \Theta_m : \min_h \sigma_S(\mathbb{O}_m(\hat{\theta}_m) \ge \alpha)\}, \mathcal{D}_m = \{\}, \text{ for all agents } m \in [n].$ 2: **for** k = 1 ... T **do** 279  $\begin{array}{l} k = 1 \dots 1 \text{ do} \\ \text{Compute } (\theta_1^k, \theta_2^k, \dots, \theta_n^k, \boldsymbol{\pi}^k) = \arg \max_{\hat{\theta}_1 \in \mathcal{B}_1^k, \hat{\theta}_2 \in \mathcal{B}_2^k, \dots, \hat{\theta}_n \in \mathcal{B}_n^k, \boldsymbol{\pi}} V^{\boldsymbol{\pi}}(\hat{\boldsymbol{\theta}}). \end{array}$ 3: for  $i = 1 \dots, n$  do 281 4: for  $h = 1, \ldots, H$  do 282 5: Selects an action according to  $a_{h,m}^k \sim \pi_{h,m}^k$  ( $\cdot \mid \tau_{h-1,m}, o_{h,m}$ ) for all  $m \in [n]$ . 6: 283 For agent  $m \in [i]$ , collect observation  $o_{h+1,m}^k$  from the environment. 284 7: For agent  $m \in [n] \setminus [i]$ , samples dummy state  $s_{h+1,m} \sim \mathbb{T}_{h,m}^k (\cdot \mid s_{h,m}, a_{h,\overline{\mathsf{pa}}(m)})$ . 8: 286 Collect observation  $o_{h+1,m}^k \sim \mathbb{O}_{h+1,m}^k (\cdot \mid s_{h+1,m})$  for agent  $m \in [n] \setminus [i]$ . 9: 287 Add  $(\pi_i^k, \tau_{\overline{pa}(i)}^k)$  into  $\mathcal{D}_i$ , and then update  $\mathcal{B}_i^{k+1}$  with eq. (2). 10: 288

**Theorem 4.1.** For all  $m \in [n]$ , we select bonus parameter as  $\beta_m = H^2(S^2A^{|\overline{pa}(m)|} +$ 290 SO)  $\log(TSAOH) + \log(Tn/\delta)$  for some constant c. Then, with probability at least  $1 - \delta$ , Algorithm 4.2 guarantees that the following inequality holds. 292

$$\operatorname{Regret}(k) = \sum_{t=1}^{k} V^{\pi^*} - V^{\pi^t} \le \tilde{\mathcal{O}}\Big(\alpha^{-2}S^2OA^{d+1}\sqrt{k(S^2A^{d+1} + SO)}\Big), \forall k \in [T], \quad (3)$$

where we define  $\pi^* = \arg \max_{\pi} V^{\pi}$ , and recall that d denotes the maximum in-degree of G.

**Remark 4.1.** The term  $\sqrt{S^2 A^{d+1} + SO}$  in 3 arises from the model error, while the additional  $OA^{d+1}$  terms result from the statistical error related to the eluder dimension. The model dimension of the factored model scales exponentially with d. Consequently, the regret also scales exponentially in d, as we incur a model estimation error of  $\mathcal{O}(A^d)$ . Notably, when  $d = \mathcal{O}(1)$ , the regret is bounded by  $poly(S, A^{\mathcal{O}(1)}, O, H, \alpha^{-1}, \log(\delta^{-1}T)).$ 

Theorem 4.2. (Lower Bound) For any randomized or deterministic algorithms, there exists an instance of DEC-POMDP with a factorization structure such that the regret for achieving global optimal is at least  $\Omega(\sqrt{A^{d+1}T})$ . 303

304 The regret scales as  $\Omega(\sqrt{A^{d+1}})$  since the model dimension scales as  $\Omega(\sqrt{A^{d+1}})$ . Therefore, Theo-305 rem 4.2 demonstrates that the dependence on the model's dimension is unavoidable. 306

#### ACHIEVING LOCAL OPTIMALITY WITH FACTORED STRUCTURE 4.3

309 In this section, we derive theoretical guarantees for achieving local optimality within a factored 310 model. Since local optimality is a specific case of global optimality, we can directly apply Algorithm 4.2 with minor modifications to achieve  $\epsilon$ -local optimality with a sample complexity of 311  $K = \tilde{\mathcal{O}}(S^4 A^{2d+2} (S^2 A^{d+1} + SO) \cdot \text{poly}(H) / (\alpha^4 \epsilon^3)).$  However, we demonstrate that a more re-312 fined analysis is possible to further improve the sample complexity. We present algorithm which 313 achieves local optimality with fewer samples compared to the direct application of Algorithm 4.2. 314

315 Algorithm Description: Due to space constraints, we refer the reader to Appendix D.1 for the com-316 plete description of Algorithm D.1. Here, we briefly outline the core idea of the algorithm and illus-317 trate the key sub-routine, emphasizing the novel contributions of our approach. Our approach entails 318 iteratively implementing the following procedure for each agent  $m \in [n]$ : we maintain the policies of agents  $[n] \setminus \{m\}$  (referred to as  $\pi_{-m}$ ) fixed and determine the policy  $\pi_m^* = \arg \max_{\mu_m} V^{\mu_m, \pi_{-m}}$ . 319 If  $V^{\pi_m^*,\pi_{-m}} - V^{\pi_m,\pi_{-m}} < \epsilon$ , we terminate the entire algorithm and output the policy  $\otimes_{m=1}^n \pi_m$ . 320 Otherwise, we replace the policy of agent m with  $\pi_m^*$  and continue the procedure. Consequently, we 321 are tasked with developing a sample-efficient algorithm to obtain  $\pi_m^* = \arg \max_{\mu_m} V^{\mu_m, \pi_{-m}}$ . We 322 present Algorithm D.1, which fulfills this task. The algorithm consists of two main steps, which we 323 now explain in detail. For clarity in the presentation, we denote  $\overline{ch}(i) = \{l_1, \ldots, l_r\}$  with  $l_r = i$ .

1:	Initialize: $\mathcal{B}_m^1 = \{\hat{\theta}_m \in \Theta_m : \min_h \sigma_S(\mathbb{O}_m(\hat{\theta}_m)) \geq \alpha\}, \mathcal{D}_m = \{\} \text{ for all } m \in \overline{ch}(i)\}$
	$\tilde{\mathcal{B}}^1 = \{\theta_{m \in nch(i)} : \min_h \sigma_S(\mathbb{O}_m(\theta_m)) \ge \alpha, \forall m \notin \overline{ch(i)}\}, \tilde{\mathcal{D}} = \{\}, \text{ central agent } i, \text{ policy of } i\}$
	other agent $\pi_{-i}$ .
2:	for $k = 1 \dots T$ do
3:	Follow $\pi_{nch(i)}$ to collect trajectories $\tau_{ch(i)}^k = \{o_{1,m}^k, \dots, a_{H,m}^k\}_{m \in nch(i)}$ .
4:	Add $\tau^k_{nch(i)}$ into $\tilde{\mathcal{D}}$ and update confidence interval with eq. (4).
5:	for $k = 1 \dots T$ do
6:	$\text{compute } (\boldsymbol{\theta}^k, \pi_i^k) = \arg \max_{\{\hat{\theta}_m \in \mathcal{B}_m^k\}_m \in Chl(i), \tilde{\boldsymbol{\theta}}_i \in \tilde{\boldsymbol{\Theta}}_i, \mu_i} V^{\mu_i, \pi_{-i}}(\hat{\boldsymbol{\theta}})$
7:	for $m=1,2,\ldots,r$ do
8:	for $h = 1, \ldots, H$ do
9:	Agent $l \in nch(i)$ take action $a_{h,l}^I$ .
10:	Select an action $a_{h,l_j}^k \sim \pi_{h,l_j}(\cdot \mid \tau_{h-1,l_j}, o_{h,l_j})$ for all $j \in [r-1]$ .
11:	Select an action $a_{h,i}^k \sim \pi_{h,i}^k (\cdot \mid \tau_{h-1,i}, o_{h,i})$ .
12:	For agent $l_j$ with $j \in [m]$ , collect observation $o_{h+1,l_j}^k$ from the environment.
13:	For $j \in [r] \setminus [m]$ , sample dummy state $s_{h+1,l_j} \sim \mathbb{T}_{h,l_j}^k (\cdot \mid s_{h,l_j}, a_{h,\overline{pa}(l_j)})$ .
14:	Collect observation $o_{h+1,l_j}^k \sim \mathbb{O}_{h+1,l_j}^k (\cdot \mid s_{h+1,l_j})$ for $j \in [r] \setminus [m]$ .
15:	If $m \neq r$ , add $(\pi_m, \tau_{\overline{pa}(m) \cap \overline{ch}(i)}^k, \tau_{\overline{pa}(m) \setminus \overline{ch}(i)}^T)$ to $\mathcal{D}_m$ .
16:	Otherwise, add $(\pi_i^k, \tau_{\overline{pa}(i)\cap \overline{ch}(i)}^k, \tau_{\overline{pa}(i)\setminus \overline{ch}(i)}^T)$ to $\mathcal{D}_i$ .
17:	Update confidence interval with eq. (5) for all $j \in \overline{ch}(i)$ .
18:	<b>Output</b> $\hat{\pi}$ as uniform mixture of the policies $\pi_i^1, \pi_i^2, \ldots, \pi_i^K$ .

Estimate model parameters θ<sub>nch(i)</sub> (Lines 2-4): According to the factored structure, agent m ∉ ch(i) is not influenced by the actions of agents m ∈ ch(i), and the policies of agents m ∉ ch(i) are predetermined. Based on this observation, the key idea of our algorithm is to estimate the model parameters θ<sub>ch(i)</sub> and θ<sub>nch(i)</sub> separately in two loops over T episodes. For the estimation of the model parameter θ<sub>nch(i)</sub>, we construct B̃ as the set of model parameters for agents m ∉ ch(i). We then iteratively follow this process for T episodes: In the k-th episode, by executing policy π<sub>nch(i)</sub>, we collect trajectories τ<sup>k</sup><sub>nch(i)</sub>. Subsequently, we update the confidence interval with:

$$\tilde{\mathcal{B}}^{k+1} = \left\{ \theta_{\mathsf{nch}(i)} \in \tilde{\mathcal{B}}^{1} : \sum_{\tau_{\mathsf{nch}(i)} \in \tilde{\mathcal{D}}} \log f(\hat{\theta}_{\mathsf{nch}(i)}, \tau_{\mathsf{nch}(i)}) \geq \max_{\hat{\theta}'_{\mathsf{nch}(i)}} \sum_{\tau_{\mathsf{nch}(i)} \in \tilde{\mathcal{D}}} f(\hat{\theta}'_{\mathsf{nch}(i)}, \tau_{\mathsf{nch}(i)}) - \tilde{\beta} \right\},\tag{4}$$

where  $\forall \theta_{\mathsf{nch}(i)}$ , we define  $f(\theta_{\mathsf{nch}(i)}, \tau_{\mathsf{nch}(i)}) = \prod_{m \notin \overline{\mathsf{ch}}(i)} \mathbb{P}_{\theta_m}^{\pi_m}(\tau_m \mid \tau_{\overline{\mathsf{pa}}(m)})$ .

Estimate model parameters θ<sub>ch(i)</sub> and update policy (Lines 9-15): We proceed with another T episodes to estimate the remaining parameters and find the optimal policy. Similar to Algorithm 4.2, we construct separate confidence intervals to estimate each model parameter for the individual agents. At the beginning of the k-th episode, we select the parameter θ<sub>m</sub><sup>k</sup> from the confidence interval B<sub>m</sub><sup>k</sup>, the parameter θ<sub>nch(i)</sub> from B̃<sup>k</sup>, and the policy for agent i that optimizes the total value function. Next, for each m ∈ [r], we collect a joint trajectory using a similar sampling procedure as in Algorithm 4.2. Specifically, at step h, agent l ∈ nch(i) takes action a<sub>h,l</sub><sup>T</sup>, while agent l<sub>j</sub> (with j ∈ [m]) samples actions and observations from the true environment, and the remaining agents sample from the model corresponding to θ<sup>k</sup>. Eventually, we add a sample to each individual sample set and update the confidence intervals with:

$$\mathcal{B}_{j}^{k+1} = \left\{ \hat{\theta}_{j} \in \mathcal{B}_{j}^{1} : \sum_{(\pi_{j}, \tau_{\overline{\mathsf{pa}}(j)}) \in \mathcal{D}_{j}} \log \mathbb{P}_{\hat{\theta}_{j}}^{\pi_{j}}(\tau_{j} \mid \tau_{\mathsf{pa}(j)}) \geq \max_{\theta_{j}' \in \Theta_{j}} \sum_{(\pi_{j}, \tau_{\overline{\mathsf{pa}}(j)}) \in \mathcal{D}_{j}} \log \mathbb{P}_{\theta_{j}'}^{\pi_{j}}(\tau_{j} \mid \tau_{\mathsf{pa}(j)}) - \beta_{j} \right\}$$

$$(5)$$

**Theorem 4.3.** We define bonus parameter in eq. (15). Then, with probability at least  $1-\delta$ , Algorithm D.1 terminates within  $4H/\epsilon$  steps of the while loop, and outputs an  $\epsilon$ -approximate local optimal



Figure 1: An example illustrating the sample complexity of Algorithm D.1. For each  $m \in \mathcal{V}$ , we highlight the subgraph induced by  $\overline{ch}(m)$  in green. It can be observed that  $\max_{m \in [n]} d_m = 1$ .

policy. The total number of episodes is at most

$$K = \tilde{\mathcal{O}}\left(\sum_{m=1}^{n} S^4 O^2 A^{2d_m+2} (S^2 A^{d+1} + SO) \cdot \operatorname{poly}(H) / (\alpha^4 \epsilon^3)\right).$$

where we use  $d_m$  to denote the maximum in-degree of the sub-graph induced by  $\overline{ch}(m)$ , and we naturally have  $d_m \leq d$ .

Technical Insights: The novelty of Algorithm D.1 compared to Algorithm 4.2 lies in its care-398 ful utilization of the structural properties. Specifically, since only the policy of the central agent 399 (denoted as i) varies across episodes, and the trajectory probability of agents  $m \notin ch(i)$  remains un-400 affected by this variation, we can pre-estimate the model parameters  $\theta_{nch(i)}$ . Finally, we proceed 401 to estimate  $\theta_{\overline{ch}(i)}$ , where we adopt a similar sampling method as in Algorithm 4.2 and achieve 402 a sample complexity that is exponential only in  $d_m$ , rather than in d. If we were to directly 403 apply Algorithm 4.2 by adjusting the parameter and policy selection to  $(\theta_1^k, \theta_2^k, \dots, \theta_n^k, \pi_m^k) =$ 404  $\arg \max_{\hat{\theta}_1 \in \mathcal{B}_1^k, \hat{\theta}_2 \in \mathcal{B}_2^k, \dots, \hat{\theta}_n \in \mathcal{B}_n^k, \pi_m} V^{\pi_m, \pi_{-m}}(\hat{\theta})$ , while keeping the other proceduresd unchanged, the 405 total number of required episodes would be  $K = \tilde{\mathcal{O}}(S^4 A^{2d+2}(S^2 A^{d+1} + SO) \cdot \text{poly}(H)/(\alpha^4 \epsilon^3)).$ 406 Thus, Algorithm D.1 significantly reduces the sample complexity needed to achieve an  $\epsilon$ -407 approximate locally optimal policy. To illustrate this improvement, we provide an example in Figure 408 1. In this example, Algorithm D.1 requires a sample complexity of  $\mathcal{O}(A^8)$ , whereas applying Algo-409 rithm 4.2 directly would result in a sample complexity of  $\mathcal{O}(A^{12})$ . 410

# 411 4.4 APPLICATION: POMDP WITH KNAPSACK CONSTRAINTS

As a minor extension, we demonstrate the applicability of our approach to a specific problem domain: POMDP with knapsack constraints, akin to the example in (Chen et al., 2020). We consider a POMDP with a budget  $\mathbf{M} \in \mathbb{R}^d$ . At each time step h, the agent incurs a cost vector  $\mathbf{C}_h$ , and the total budget updates to  $\mathbf{M}_{h+1} = \mathbf{M}_h - \mathbf{C}_h$ . We model the transition of each budget component i as  $\mathbf{M}_{h+1,i} \sim \mathbb{T}_{h,\mathbf{M}}(\cdot \mid \mathbf{M}_{h,i}, o_h, a_h)$ . The episode terminates when any budget component reaches 0.

418 We formulate this problem as a factored DEC-POMDP with d+1 agents, 419 treating the budgets as observations of d dummy agents. Consequently, 420 Algorithm 4.2 can be directly applied. However, since the budgets are 421 directly observed, there are still opportunities for improvement. We es-422 timate the transition  $\mathbb{T}_{h,\mathbf{B}}$  using a confidence interval approach akin to 423 UCB-VI (Azar et al., 2017), and we can achieve a sharper bound. Due 424 to the space limit, we defer the complete discussion to Appendix F.

424 425 426

427

412

378 379 380

382

384

386 387

388

389

390 391

392 393 394

396 397

### 5 LEARNING DEC-POMDP WITH MEMORYLESS POLICY

In this section, we focus on the setting where agents adopt memoryless policies. Namely, each agent select her action base solely on her current observation. We define the policy class as  $\{\{\otimes_{m=1}^{n}\pi_{h,m}\}_{h\in[H]} \mid \pi_{h,m}: \mathcal{O}_{m} \to \Delta(\mathcal{A}_{m})\}$ . Notably, our results can be readily extended to settings where agents consider observations and actions from the preceding *L* steps. In such cases, the policy class broadens to  $\{\{\otimes_{m=1}^{n}\pi_{h,m}\}_{h\in[H]} \mid \pi_{h,m}: (\mathcal{O}_{m} \times \mathcal{A}_{m})^{\min(h,L)-1} \times \mathcal{O}_{m} \to \Delta(\mathcal{A}_{m})\}$ .

influential graph

452 453

454 455

456 457

458

459

460

461

462

463

468 469

470

471

472

# 4324335.1 Achieving Local Optimality with Memoryless Policy

We utilize a framework similar to that described in Section 4.3. Specifically, we iteratively update the policy of the  $m^{th}$  agent using  $\pi_m^* = \arg \max_{\mu_m} V^{\mu_m, \pi_{-m}}$  and terminate the procedure when further updates no longer produce a significant increase in the total value function. Our remaining task is to develop a sample-efficient method to obtain  $\pi_m^* = \arg \max_{\mu_m} V^{\mu_m, \pi_{-m}}$  given  $\pi_{-m}$ , for all  $m \in [n]$ . We present Algorithm 5.1, which addresses this task. Due to space constraints, a detailed description of the complete algorithm is provided in Appendix E.1 (Algorithm E.1).

**Algorithm Description:** For each agent  $m \in [n]$ , the parameter set  $\theta_m$  represents the model parameters of the joint probability distribution of trajectories for the *m*-th agent and the *i*-th agent. We denote  $\mathbb{P}_{\theta_i}^{\pi_i}(\tau_i)$  as the probability of the trajectory for the *i*-th agent and  $\mathbb{P}_{\theta_m}^{\pi_i,\pi_m}(\tau_i,\tau_m)$  as the joint probability of trajectories for the *i*-th agents, given the underlying DEC-POMDP with parameters  $\theta = (\theta_1, \ldots, \theta_n)$ . The formal definition is provided in Appendix E.1 (Equation equation 18). With these definitions, we now proceed to explain Algorithm 5.1 in detail.

• Update policy and parameters (Lines 3-4): We denote the value functions  $V_i^{\pi_i,\pi_{-i}}(\theta_i)$  and  $V_m^{\pi_i,\pi_{-i}}(\theta_m)$  for all  $m \in [n] \setminus \{i\}$  as:

$$V_{i}^{\pi_{i},\pi_{-i}}(\theta_{i}) = \sum_{\tau_{i}} \mathbb{P}_{\theta_{i}}^{\pi_{i}}(\tau_{i}) \Big( \sum_{h=1}^{H} r_{h,i}(o_{h,i}) \Big),$$
  
$$V_{m}^{\pi_{i},\pi_{-i}}(\theta_{m}) = \sum_{\tau_{m},\tau_{i}} \mathbb{P}_{\theta_{m}}^{\pi_{i},\pi_{m}}(\tau_{m},\tau_{i}) \Big( \sum_{h=1}^{H} r_{h,m}(o_{h,m}) \Big).$$

The observation probabilities  $\{\mathbb{O}_{h,m}\}_{h=1}^{H}$  and the transition probabilities  $\{\mathbb{T}_{h,m}\}_{h=1}^{H}$  corresponding to the real model  $\theta_m = \theta_m^*$  are defined as per Equations equation 19 and equation 20 in Appendix E.1. In this context, the value function can be decomposed as  $V^{\pi} = \sum_{m=1}^{n} V_m^{\pi_i,\pi_{-i}}(\theta_m^*)$ . Thus, we decompose the value function into *n* distinct terms, each depending solely on the parameter  $\theta_m$ . For each  $m \in [n]$ , we select  $\theta_m^k \in \mathcal{B}_m^k$  as the optimal parameter that maximizes  $V_m^{\pi_i,\pi_{-i}}(\theta_m)$ . We then determine the policy  $\pi_i^k$  as the optimal policy that maximizes the total value function. Subsequently, we use the policy  $\pi^k = (\pi_i^k, \pi_{-i})$  to collect a trajectory  $\tau^k = (\mathbf{o}_1^k, \mathbf{a}_1^k, \dots, \mathbf{o}_H^k, \mathbf{a}_H^k)$ .

• Construct Confidence Intervals (Lines 5-6): We construct n different confidence intervals to estimate the n model parameters separately. For each agent  $m \in [n] \setminus \{i\}$ , we add the newly collected policy-trajectory pair  $(\pi_i^k, \tau_i^k, \tau_m^k)$  to the dataset  $\mathcal{D}_m$ . Similarly, for agent i, we add the policy-trajectory pair  $(\pi_i^k, \tau_i^k)$  to the dataset  $\mathcal{D}_i$ . Subsequently, we update each of the n confidence intervals separately according to the following equations for agent i and  $m \in [n] \setminus \{i\}$ :

$$\mathcal{B}_{i}^{k+1} = \left\{ \hat{\theta}_{i} \in \mathcal{B}_{i}^{1} : \sum_{(\pi_{i},\tau_{i})\in\mathcal{D}_{i}} \log \mathbb{P}_{\hat{\theta}_{i}}^{\pi_{i}}(\tau_{i}) \geq \max_{\theta_{i}'\in\Theta_{i}} \sum_{(\pi_{i},\tau_{i})\in\mathcal{D}_{i}} \log \mathbb{P}_{\theta_{i}'}^{\pi_{i}}(\tau_{i}) - \beta_{i} \right\}$$
$$\mathcal{B}_{m}^{k+1} = \left\{ \hat{\theta}_{m} \in \mathcal{B}_{m}^{1} : \sum_{(\pi_{i},\tau_{i},\tau_{m})\in\mathcal{D}_{m}} \log \mathbb{P}_{\hat{\theta}_{m}}^{\pi_{i},\pi_{m}}(\tau_{i},\tau_{m}) \geq \max_{\theta_{m}'\in\Theta_{m}} \sum_{(\pi_{i},\tau_{i},\tau_{m})\in\mathcal{D}_{m}} \log \mathbb{P}_{\theta_{m}'}^{\pi_{i},\pi_{m}}(\tau_{i},\tau_{m}) - \beta_{m} \right\}$$

**Technical Challenge and Novelty:** To showcase our novel approach, we highlight the challenges in developing a sample-efficient algorithm to find  $\pi_m^* = \arg \max_{\mu_m} V^{\mu_m, \pi_{-m}}$  with given  $\pi_{-m}$ , along with our solutions to these challenges:

1. In DEC-MDPs, given  $\pi_{-m}$ , the model reduces to a single-agent problem with action space  $\mathcal{A}_m$ . However, this reduction does not apply to DEC-POMDPs, precluding the use of single-agent algorithms for sample-efficient guarantees as in MDP setting.

476 2. Another challenge arises from the exponential growth of the joint action and observation space 477 with the number of agents, resulting in a model dimension that scales as  $\mathcal{O}(A^n O^n)$ . Conse-478 quently, constructing a single confidence interval to estimate the model parameters leads to a 479 sample complexity of  $\mathcal{O}(A^n O^n)$ .

480 We overcome the technical challenges by assigning a parameter for the trajectory probability of each 481 agent, subsequently estimating and updating these parameters with separate confidence intervals. 482 Since the dimension of parameter  $\theta_m$  ( $m \in [n]$ ) is at most  $H(S^2A^2O^2 + SO^2) + S$ , we achieve an 483  $\epsilon$ -approximate local optimum with a sample complexity that avoids exponential scaling in n.

**Theorem 5.1.** Let the central agent for Algorithm 5.1 be agent *i*. We define the bonus parameter as eq. (17). Then, with a probability of at least  $1 - \delta$ , Algorithm E.1 terminates within  $4H/\epsilon$  steps of the while loop and outputs an  $\epsilon$ -approximate local optimal policy. The total number of episodes

Algorit	thm 3 OMLE for memoryless policy
1: Inp	<b>put</b> : Central agent <i>i</i> , and the policy for agent $[n] \setminus \{i\}, \pi_1, \pi_2, \ldots, \pi_{i-1}, \pi_{i+1}, \ldots, \pi_n$ .
2: Ini	$\textbf{tialize: } \mathcal{B}_i^1 = \{ \hat{\theta}_i \in \Theta_i : \min_h \sigma_S(\mathbb{O}_i(\hat{\theta}_i) \ge \alpha) \}, \mathcal{B}_m^1 = \{ \hat{\theta}_m \in \Theta_m : \min_h \sigma_S(\mathbb{O}_m(\hat{\theta}_m) \ge \alpha) \}$
$\alpha \setminus$	$(\sqrt{O})$ for all $m \in [n] \setminus \{i\}$ . Set $\mathcal{D}_m = \{\}$ , for all agents $m \in [n]$ .
3: <b>for</b>	$k = 1 \dots T$ do
4:	Compute $(\theta_1^k, \theta_2^k, \dots, \theta_n^k, \pi_i^k) = \arg \max_{\hat{\theta}_1 \in \mathcal{B}_1^k, \hat{\theta}_2 \in \mathcal{B}_2^k, \dots, \hat{\theta}_n \in \mathcal{B}_n^k, \pi_i} \sum_{m=1}^n V_m^{\pi_i, \pi_{-i}}(\hat{\theta}_m).$
5:	Follow $\pi^k$ to collect a trajectory $\boldsymbol{\tau}^k = (\mathbf{o}_1^k, \mathbf{a}_1^k, \dots, \mathbf{o}_H^k, \mathbf{a}_H^k)$ .
6:	Add $(\pi_i^k, \tau_i^k, \tau_m^k)$ into $\mathcal{D}_m$ for $m \in [n] \setminus \{i\}$ and add $(\pi_i^k, \tau_i^k)$ into $\mathcal{D}_i$ .
7:	Update $\mathcal{B}_i^{k+1}$ and $\mathcal{B}_m^{k+1}$ for all $m \in [n] \setminus \{i\}$ with eq. (6).
8: <b>Ou</b>	<b>tput</b> $\hat{\pi}$ , which is selected uniformly from the policies $\pi_i^1, \pi_i^2, \ldots, \pi_i^T$ .

played by Algorithm E.1 is at most

 $K = \tilde{\mathcal{O}} \left( S^4 O^4 A^4 (S^2 A^2 O^2 + SO^2) \cdot \operatorname{poly}(H) / (\alpha^4 \epsilon^3) \right).$ 

**Remark 5.1.** In a commonly studied model (where agents adopt memoryless policies), we derive an algorithm capable of achieving an  $\epsilon$ -approximate local optimal policy for DEC-POMDPs. Importantly, the sample complexity of this algorithm does not scale exponentially with n.

Moreover, since DEC-POMDPs can be seen as a special case of a partially observable version of a Markov potential game, our framework extends to the analysis of achieving Nash equilibrium within this partially observable version.

### 5.2 HARDNESS RESULT FOR ACHIEVING GLOBAL OPTIMALITY

In addition to local optimality, we now explore the attainment of global optimality with memoryless policies. However, the following theorem reveals that, without additional assumptions on the model, deriving an algorithm to achieve global optimality with regret not exponential in n is unattainable.

**Theorem 5.2.** For any randomized or deterministic algorithms, there exists an instance of DEC-MDP with horizon H = 2 wherein the regret scales at least as  $O(\sqrt{A^nT})$ . This result underscores the limitation of achieving sample efficiency in algorithms for DEC-POMDP without imposing assumptions on the transition model, either when agents adopt memoryless policies.

519 520

521

500

501 502

503

504

505

506

507

508 509 510

511

6 CONCLUSION AND DISCUSSION

522 **Conclusion and Summary:** This work introduces a sample-efficient algorithm and provides the-523 oretical guarantees for DEC-POMDPs. Theorem 5.2 highlights the challenges associated with de-524 veloping sample-efficient algorithms for DEC-POMDPs without making any assumptions about the model. Consequently, our focus shifts towards identifying such algorithms under specific condi-526 tions rather than for the general model. Initially, we present a sample-efficient algorithm for a 527 commonly studied scenario where agents utilize memoryless policies. Furthermore, inspired by em-528 pirical methods that leverage value decomposition to address exponential complexity, we propose 529 a factored structural model as a sufficient condition for value decomposition and derive a sampleefficient algorithm based on this assumption. This analysis provides a theoretical foundation for the 530 empirical strategies currently employed. 531

532 **Open Directions**: One open question is whether it is possible to derive a sample-efficient algorithm 533 that achieves local optimality without imposing any assumptions on the model. When Algorithm 5.1 534 is applied to a full-memory setting, the sample complexity upper bound scales as  $\mathcal{O}(A^{\overline{H}})$ . In contrast, applying the vanilla OMLE algorithm from Liu et al. (2022a) results in a sample complexity 535 536 of  $\mathcal{O}(A^n)$ . Therefore, it remains an open problem to determine whether a lower bound on sample complexity can scale as  $\mathcal{O}(A^{\min\{H,n\}})$ , or if it is possible to overcome the multi-agent curse without 537 imposing assumptions on the model. Another avenue for future research is to explore whether addi-538 tional reasonable assumptions about the model could facilitate the development of sample-efficient algorithms. We leave these directions for future investigation.

# 540 REFERENCES

547

553

559

565

571

577

578

579

580

584

585

586

- Altabaa, A. and Yang, Z. (2024). On the role of information structure in reinforcement learning for
   partially-observable sequential teams and games.
- 544 https://arxiv.org/abs/2403.00993
- Azar, M. G., Osband, I. and Munos, R. (2017). Minimax regret bounds for reinforcement learning.
   In *International conference on machine learning*. PMLR.
- Burago, D., De Rougemont, M. and Slissenko, A. (1996). On the complexity of partially observed markov decision processes. *Theoretical Computer Science*, **157** 161–183.
- Cai, Q., Yang, Z. and Wang, Z. (2022). Reinforcement learning from partial observation: Linear function approximation with provable sample efficiency. In *International Conference on Machine Learning*. PMLR.
- Chakraborty, D. and Stone, P. (2011). Structure learning in ergodic factored mdps without knowl edge of the transition function's in-degree. In *Proceedings of the 28th International Conference on Machine Learning (ICML-11).*
- <sup>557</sup> Chen, F., Bai, Y. and Mei, S. (2022). Partially observable rl with b-stability: Unified structural
   <sup>558</sup> condition and sharp sample-efficient algorithms. *arXiv preprint arXiv:2209.14990*.
- Chen, X., Hu, J., Li, L. and Wang, L. (2020). Efficient reinforcement learning in factored mdps with application to constrained rl. *arXiv preprint arXiv:2008.13319*.
- 562 Cui, Q., Zhang, K. and Du, S. (2023). Breaking the curse of multiagents in a large state space: Rl
   563 in markov games with independent linear function approximation. In *The Thirty Sixth Annual* 564 *Conference on Learning Theory*. PMLR.
- Daskalakis, C., Golowich, N. and Zhang, K. (2023). The complexity of markov equilibrium in stochastic games. In *The Thirty Sixth Annual Conference on Learning Theory*. PMLR.
- Diuk, C., Li, L. and Leffler, B. R. (2009). The adaptive k-meteorologists problem and its application
   to structure learning and feature selection in reinforcement learning. In *Proceedings of the 26th Annual International Conference on Machine Learning*.
- Efroni, Y., Jin, C., Krishnamurthy, A. and Miryoosefi, S. (2022). Provable reinforcement learning with a short-term memory. In *International Conference on Machine Learning*. PMLR.
- Foerster, J., Song, F., Hughes, E., Burch, N., Dunning, I., Whiteson, S., Botvinick, M. and
  Bowling, M. (2019). Bayesian action decoder for deep multi-agent reinforcement learning. In *International Conference on Machine Learning*. PMLR.
  - Gu, H., Guo, X., Wei, X. and Xu, R. (2021). Mean-field controls with q-learning for cooperative marl: convergence and complexity analysis. *SIAM Journal on Mathematics of Data Science*, **3** 1168–1196.
- Guestrin, C., Koller, D. and Parr, R. (2001). Solving factored pomdps with linear value functions.
   In Seventeenth International Joint Conference on Artificial Intelligence (IJCAI-01) workshop on Planning under Uncertainty and Incomplete Information. Citeseer.
  - Hu, H. and Foerster, J. N. (2019). Simplified action decoder for deep multi-agent reinforcement learning. *arXiv preprint arXiv:1912.02288*.
- Jin, C., Kakade, S., Krishnamurthy, A. and Liu, Q. (2020). Sample-efficient reinforcement learning of undercomplete pomdps. *Advances in Neural Information Processing Systems*, 33 18530–18539.
- Jin, C., Liu, Q., Wang, Y. and Yu, T. (2021). V-learning–a simple, efficient, decentralized algorithm for multiagent rl. *arXiv preprint arXiv:2110.14555*.
- 593 Kara, A. and Yuksel, S. (2022). Near optimality of memory-finite feedback policies in partially observed markov decision processes. *Journal of Machine Learning Research*, **23** 1–46.

594 595 596	Kara, A. D. and Yüksel, S. (2023). Convergence of memory-finite q learning for pomdps and near optimality of learned policies under filter stability. <i>Mathematics of Operations Research</i> , 48 2066–2093.
597 598 599	Katt, S., Oliehoek, F. and Amato, C. (2018). Bayesian reinforcement learning in factored pomdps. <i>arXiv preprint arXiv:1811.05612</i> .
600 601	Kraemer, L. and Banerjee, B. (2016). Multi-agent reinforcement learning as a rehearsal for decentralized planning. <i>Neurocomputing</i> , <b>190</b> 82–94.
603	Kreps, D. M. (1989). Nash equilibrium. In Game theory. Springer, 167–177.
604 605 606	Krishnamurthy, A., Agarwal, A. and Langford, J. (2016). Pac reinforcement learning with rich observations. <i>Advances in Neural Information Processing Systems</i> , <b>29</b> .
607 608 609	Lerer, A., Hu, H., Foerster, J. and Brown, N. (2020). Improving policies via search in cooperative partially observable games. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , vol. 34.
610 611	Littman, M. L. (1994). Memoryless policies: Theoretical limitations and practical results.
612 613	Liu, Q., Chung, A., Szepesvári, C. and Jin, C. (2022a). When is partially observable reinforcement learning not scary? In <i>Conference on Learning Theory</i> . PMLR.
614 615 616 617	Liu, Q., Netrapalli, P., Szepesvari, C. and Jin, C. (2023). Optimistic mle: A generic model-based algorithm for partially observable sequential decision making. In <i>Proceedings of the 55th Annual ACM Symposium on Theory of Computing</i> .
618 619	Liu, Q., Szepesvári, C. and Jin, C. (2022b). Sample-efficient reinforcement learning of partially observable markov games. <i>Advances in Neural Information Processing Systems</i> , <b>35</b> 18296–18308.
620 621 622	Liu, X. and Zhang, K. (2023). Partially observable multi-agent rl with (quasi-) efficiency: the bless- ing of information sharing. In <i>International Conference on Machine Learning</i> . PMLR.
623 624	Lusena, C., Goldsmith, J. and Mundhenk, M. (2001). Nonapproximability results for partially observable markov decision processes. <i>Journal of artificial intelligence research</i> , <b>14</b> 83–103.
625 626 627	Mannor, S. and Tsitsiklis, J. N. (2004). The sample complexity of exploration in the multi-armed bandit problem. <i>Journal of Machine Learning Research</i> , <b>5</b> 623–648.
628 629	Mao, W., Zhang, K., Yang, Z. and Başar, T. (2023). Decentralized learning of finite-memory policies in dec-pomdps. <i>IFAC-PapersOnLine</i> , <b>56</b> 2601–2607.
630 631 632 633	Mondal, W. U., Agarwal, M., Aggarwal, V. and Ukkusuri, S. V. (2022). On the approximation of cooperative heterogeneous multi-agent reinforcement learning (marl) using mean field control (mfc). <i>Journal of Machine Learning Research</i> , <b>23</b> 1–46.
634 635	Mossel, E. and Roch, S. (2005). Learning nonsingular phylogenies and hidden markov models. In <i>Proceedings of the thirty-seventh annual ACM symposium on Theory of computing.</i>
636 637 638 639	Oliehoek, F. A., Spaan, M. T. J. and Vlassis, N. (2008a). Optimal and approximate q-value functions for decentralized pomdps. <i>Journal of Artificial Intelligence Research</i> , <b>32</b> 289–353. http://dx.doi.org/10.1613/jair.2447
640 641 642 643	Oliehoek, F. A., Spaan, M. T. J., Whiteson, S. and Vlassis, N. (2008b). Exploiting locality of in- teraction in factored dec-pomdps. In <i>Proceedings of the 7th International Joint Conference on</i> <i>Autonomous Agents and Multiagent Systems - Volume 1</i> . AAMAS '08, International Foundation for Autonomous Agents and Multiagent Systems, Richland, SC.
644 645 646	Osband, I. and Van Roy, B. (2014). Near-optimal reinforcement learning in factored mdps. <i>Advances in Neural Information Processing Systems</i> , <b>27</b> .
	<b>D</b> $(1)$

647 Papadimitriou, C. H. and Tsitsiklis, J. N. (1987). The complexity of markov decision processes. *Mathematics of operations research*, **12** 441–450. 648 Pasztor, B., Bogunovic, I. and Krause, A. (2021). Efficient model-based multi-agent mean-field 649 reinforcement learning. arXiv preprint arXiv:2107.04050. 650 Oiu, D., Wang, J., Dong, Z., Wang, Y. and Strbac, G. (2022). Mean-field multi-agent reinforcement 651 learning for peer-to-peer multi-energy trading. IEEE Transactions on Power Systems, 38 4853-652 4866. 653 654 Rashid, T., Samvelyan, M., de Witt, C. S., Farquhar, G., Foerster, J. and Whiteson, S. (2018). Qmix: 655 Monotonic value function factorisation for deep multi-agent reinforcement learning. 656 Strehl, A. L., Diuk, C. and Littman, M. L. (2007). Efficient structure learning in factored-state mdps. 657 In AAAI, vol. 7. 658 659 Sunehag, P., Lever, G., Gruslys, A., Czarnecki, W. M., Zambaldi, V., Jaderberg, M., Lanctot, M., 660 Sonnerat, N., Leibo, J. Z., Tuyls, K. and Graepel, T. (2017). Value-decomposition networks for 661 cooperative multi-agent learning. 662 Tian, Y., Qian, J. and Sra, S. (2020). Towards minimax optimal reinforcement learning in factored 663 markov decision processes. Advances in Neural Information Processing Systems, 33 19896– 664 19907. 665 666 Uehara, M., Sekhari, A., Lee, J. D., Kallus, N. and Sun, W. (2022). Provably efficient reinforcement 667 learning in partially observable dynamical systems. Advances in Neural Information Processing Systems, 35 578–592. 668 669 Uehara, M., Sekhari, A., Lee, J. D., Kallus, N. and Sun, W. (2023). Computationally efficient pac rl 670 in pomdps with latent determinism and conditional embeddings. In International Conference on 671 Machine Learning. PMLR. 672 Wang, J., Ren, Z., Liu, T., Yu, Y. and Zhang, C. (2021). Qplex: Duplex dueling multi-agent q-673 learning. 674 675 Wang, L., Cai, Q., Yang, Z. and Wang, Z. (2022). Embed to control partially observed systems: 676 Representation learning with provable sample efficiency. arXiv preprint arXiv:2205.13476. 677 Wang, Y., Liu, Q., Bai, Y. and Jin, C. (2023). Breaking the curse of multiagency: Provably efficient 678 decentralized multi-agent rl with function approximation. In The Thirty Sixth Annual Conference 679 on Learning Theory. PMLR. 680 681 Xu, Z. and Tewari, A. (2020). Reinforcement learning in factored mdps: Oracle-efficient algorithms 682 and tighter regret bounds for the non-episodic setting. Advances in Neural Information Processing 683 Systems, 33 18226–18236. 684 Yang, Y., Luo, R., Li, M., Zhou, M., Zhang, W. and Wang, J. (2018). Mean field multi-agent rein-685 forcement learning. In International conference on machine learning. PMLR. 686 687 Zhang, H., Ren, T., Xiao, C., Schuurmans, D. and Dai, B. (2023). Provable representation with effi-688 cient planning for partially observable reinforcement learning. arXiv preprint arXiv:2311.12244. 689 Zhang, K., Miehling, E. and Başar, T. (2019). Online planning for decentralized stochastic control 690 with partial history sharing. In 2019 American Control Conference (ACC). IEEE. 691 692 693 694 696 697 699 700

# Appendix

# Table of Contents

A	Addi	tional Related Work	15
B	Proo	f of Theorem 3.1	16
С	Supp	elementary Details for Section 4.2	16
	C.1	Proof for Theorem 4.1	16
	C.2	Proof for Theorem 4.2	24
D	Supp	elementary Details for Section 4.3	25
	D.1	Complete Algorithm for Achieving Local Optimal Under Factored Structure Model	25
	D.2	Proof for Theorem 4.3	26
Е	Sup	plementary Details for Section 5	34
	E.1	Complete Algorithm for Achieving Local Optimal with Memoryless Policies	34
	E.2	Proof for Theorem 5.1	36
	E.3	Proof for Theorem 5.2	42
F	POM	IDP with Knapsack Constraints	42
	F.1	Model	42
	F.2	Algorithm	43
	F.3	Theoretical Guarantee	44
			10
	F.4	Improvement to Achieve Sharper Bound	46

#### 756 ADDITIONAL RELATED WORK А

757 758

Learning POMDPs Learning partially observable Markov decision processes (POMDPs) 759 presents significant challenges due to the lack of the Markov property in observations and the depen-760 dence of policies on the full observation history. This complexity is underscored by lower bounds, 761 such as those established by Mossel and Roch (2005) and Krishnamurthy et al. (2016), which show 762 exponential complexity in the horizon for learning near-optimal policies in POMDPs. Given the dif-763 ficulty of learning POMDPs in the general case, recent research has explored learning under various 764 structural conditions. Some works, like Jin et al. (2020) and Liu et al. (2022a), have investigated weakly revealing conditions, while others, such as Cai et al. (2022) and Wang et al. (2022), have 765 focused on low-rank POMDPs. Efroni et al. (2022) and Zhang et al. (2023) have delved into learn-766 ing under decodable conditions, while Uehara et al. (2022) and Uehara et al. (2023) have proposed 767 algorithms for learning with memoryless policies and deterministic transition models, respectively. 768 Chen et al. (2022) have introduced the B-stability condition as a comprehensive framework that en-769 compasses previous structural conditions. In our work, we demonstrate our results under a weakly 770 revealing condition akin to that of Jin et al. (2020) and Liu et al. (2022a). However, it's important to 771 note that our framework can be extended to incorporate other conditions proposed in previous works, 772 such as the B-stability condition introduced by Chen et al. (2022). This flexibility underscores the 773 applicability and generality of our approach within the broader landscape of learning POMDPs.

774

775 Learning POMGs Liu et al. (2022b) present the OMLE algorithm for finding approximate 776 Nash equilibria, correlated equilibria, as well as coarse correlated equilibrium of weakly revealing 777 POMGs in a polynomial number of samples, particularly when the number of agents is small. On a 778 related note, Liu and Zhang (2023) develop a partially observable multi-agent reinforcement learn-779 ing (MARL) algorithm that is both statistically and computationally quasi-efficient, incorporating information sharing under the general framework of partially observable stochastic games. 780

781

782 Learning MDPs and POMDPs with Specific Structures We consider a factorized structure 783 model in this work. In the context of factored MDPs, Osband and Van Roy (2014) first proposed the 784 factored MDP model and introduced PSRL and UCRL-style algorithms with near-optimal Bayesian and frequentist regret bounds. Xu and Tewari (2020) extended the results of Osband and Van Roy 785 (2014) to the infinite horizon setting. Tian et al. (2020) applied the UCBVI algorithm (Azar et al., 786 2017) to the factored MDP framework, while Chen et al. (2020) further refined the approach by 787 applying the UCB-VI algorithm and developing the FMDP-BF algorithm, which achieves a sharper 788 bound compared to Tian et al. (2020). Additionally, Chen et al. (2020) introduced reinforcement 789 learning with knapsack constraints as an example of factored MDPs. Diuk et al. (2009) proposes an 790 algorithmic framework based on the KWIK principle for learning probabilistic concepts, and applies 791 this framework to reinforcement learning in factored models. The authors provide empirical insights 792 that suggest more efficient algorithms can be derived when restricted to factored structure models. 793 Strehl et al. (2007) addresses the reinforcement learning problem in factored MDPs and proposes an 794 efficient algorithm leveraging dynamic Bayesian networks (DBNs). Chakraborty and Stone (2011) 795 studies factored-state MDPs and aims to develop an algorithm that guarantees a return close to the optimal in factored MDPs. Since factored-state MDPs are a special case of Markov chains, they 796 utilize the properties of ergodic stochastic processes to analyze factored MDPs. 797

798 In terms of POMDPs, several prior studies have explored POMDPs with specific factored structures. 799 For example, Katt et al. (2018) introduced the Factored Bayes-Adaptive POMDP model, along with 800 a method to learn both the factorization and the model parameters simultaneously. Guestrin et al. (2001) demonstrated that, for factored POMDPs, the value function can be represented as a linear 801 combination of basis functions, enabling the derivation of an efficient algorithm by leveraging the 802 decomposition of the value function. Similar to our setting, several previous studies have examined 803 factor structures in DEC-POMDPs. For instance, Oliehoek et al. (2008b) analyzed general factored 804 DEC-POMDPs, focusing on the model's dependencies over space and time, and formulated decom-805 posable value functions. 806

807 The aforementioned works primarily focus on analyzing the state structure of POMDPs. Beyond considering POMDPs with specific state structures, Altabaa and Yang (2024) investigated the role 808 of information structure, which describes how events in the system occurring at different points in time influence each other. Altabaa and Yang (2024) also provided an upper bound on the sample 810 complexity for learning a general sequential decision-making problem with a directed acyclic graph 811 (DAG) information structure. 812

813 Learning with Memoryless Policies Since learning POMDPs is known to be PSPACE-complete, 814 many works focus on developing algorithms to learn optimal memoryless policies, which can be 815 viewed as a special case of general POMDPs. Kara and Yuksel (2022) studied learning optimal 816 memoryless policies for POMDPs by approximating the belief model through discretizing the belief 817 space. Kara and Yüksel (2023) provided convergence analysis for a Q-learning algorithm tailored 818 to POMDPs with memoryless policies.

819 820

Learning Multi-agent System: In multi-agent reinforcement learning (MARL), the action space 821 grows exponentially with the number of agents, making it crucial to derive algorithms whose sam-822 ple complexity is not exponential in the number of agents-a challenge commonly referred to as 823 breaking the curse of multi-agency. Daskalakis et al. (2023); Jin et al. (2021) derive sample-efficient 824 algorithms with non-exponential sample complexity, while Wang et al. (2023); Cui et al. (2023) 825 further generalize this approach to settings with linear function approximation. Another method to address the exponential growth of the action space is the mean-field approach, which assumes that 826 each agent's decision is influenced by the mean field (i.e., the average behavior of other agents) 827 rather than by the individual actions of each agent. Previous work utilizing the mean-field method to 828 tackle exponential growth in multi-agent RL includes Yang et al. (2018); Pasztor et al. (2021); Qiu 829 et al. (2022). Gu et al. (2021) demonstrates that if all agents are homogeneous and exchangeable, 830 mean-field control can provide a good approximation to an N-agent problem. Similarly, Mondal 831 et al. (2022) provides a comparable approximation for a K-class of heterogeneous agents. In con-832 trast to the mean-field method, we adopt a similar approach of optimizing the performance of a 833 single agent while considering the overall effect of others. This allows the optimization problem to be approximated as a single-agent problem, thereby mitigating the exponential growth of the action 835 space in the environment.

836 837

838 839 840

841

842 843

847

848

849

851

#### В **PROOF OF THEOREM 3.1**

**Theorem B.1.** For any randomized or deterministic algorithms, there exists an instance of DEC-MDP wherein the regret scales at least as  $\Omega(\sqrt{A^nT})$ .

- 844 *Proof.* The proof for Theorem 3.1 proceeds straightforwardly. We consider a two-step DEC-MDP, commencing from an initial state  $s_1$ . For all  $s \in S$ , we assume the reward function satisfies 845  $r_{h,1}(s) = r_{h,2}(s) = \cdots = r_{h,n}(s)$  for all  $h \in [2]$ . Consequently, the entire DEC-MDP reduces 846 to a multi-armed bandit problem. By leveraging a classic result on the lower bound of regret for the multi-armed bandit problem (Mannor and Tsitsiklis, 2004), it follows that for any randomized or deterministic algorithm, there exists an instance of the multi-arm bandit problem such that the regret is at least  $\mathcal{O}(\sqrt{A}T)$ , where  $\tilde{A}$  denotes the number of arms. Consequently, for any random-850 ized or deterministic algorithm, there exists an instance of DEC-MDP such that the regret is at least  $\mathcal{O}(\sqrt{A^nT}).$  $\square$
- 852 853 854

#### С SUPPLEMENTARY DETAILS FOR SECTION 4.2

855 856 857

C.1 PROOF FOR THEOREM 4.1

858 In this section, we present the proof of Theorem 4.1. To ensure clarity, we begin by defining several 859 notations that will be useful throughout the proof. 860

**Definition C.1.** For all  $(m, i) \in [n] \times [T]$ ,  $\theta_m \in \Theta_m$ , and any policy  $\pi_m$  of agent m, we denote 861  $f_m(\theta_m, \pi_m)$  as  $f_m(\theta_m, \pi_m) = \mathbb{P}_{\theta_m}^{\pi_m}(\tau_m \mid \{\tau_r\}_{r \in \mathsf{pa}(m)})$ . Additionally, we use  $\tilde{f}_m^i(\theta_m, \pi_m)$  to 862 denote  $\tilde{f}_m^i(\theta_m, \pi_m) = \mathbb{P}_{\theta_m, m}^{\pi_m} \left( \tau_m^i \mid \{\tau_r^i\}_{r \in \mathsf{pa}(m)} \right).$ 863

**Lemma C.1.** For all  $(\theta_m, t) \in \Theta_m \times [T]$  and agent  $m \in [n]$ , the following inequality holds with probability at least  $1 - \delta$ : 

$$\sum_{i=1}^{t} \log \left( \tilde{f}_m^i(\theta_m, \pi_m) / \tilde{f}_m^i(\theta_m^{\star}, \pi_m) \right) \le \beta_m,$$

where we define bonus term  $\beta_m = c(H^2(S^2A^{|\mathsf{pa}(m)|+1} + SO)\log(TSAOH) + \log(Tn/\delta))$  for some absolute constant c. 

*Proof.* Initially, we can view  $\Theta_m$  as a subset of  $\mathbb{R}^{d_m}$  with  $d_m = H(S^2 A^{|\mathsf{pa}(m)|} + SO) + S$ . We denote  $\theta_m$  as the optimistic  $\epsilon$ -discretelization of  $\theta_m$ , so that  $\overline{\theta}_{m,i} = \lceil \theta_{m,i}/\epsilon \rceil \times \epsilon$  for all coordinates i. Selecting  $\epsilon \leq 1/(c(S+O+A)HTO^{H}A^{H})$ , we obtain the following relationship:

$$f_m(\bar{\theta}_m, \pi_m) \ge f_m(\theta_m, \pi_m), \quad \left| f_m(\bar{\theta}_m, \pi_m) - f_m(\theta_m, \pi_m) \right| \le 1/(TO^H A^H).$$

The inequalities holds for all trajectories  $\tau_{H,m} \in (\mathcal{O}_m \times \mathcal{A}_m)^H$ . We use  $\overline{\Theta}_m$  to represent the collections of all such  $\bar{\theta}_m$ , then, the log-cardinality of  $\bar{\Theta}_m$  is bounded by

$$\log \left|\bar{\Theta}_{m}\right| \leq \mathcal{O}\left(H^{2}\left(S^{2}A^{|\mathsf{pa}(m)|+1}+SO\right)\log(TSAOH)\right).$$

In the following step, we aim to apply Markov inequality to bound the following expectation:  $\mathbb{E}[\exp(\sum_{i=1}^{t}\log(\tilde{f}_m^i(\bar{\theta}_m,\pi_m)/\tilde{f}_m^i(\theta_m^{\star},\pi_m)))]. \text{ We denote } \mathbb{E}_t[\cdot] = \mathbb{E}\left[\cdot \mid \{\boldsymbol{\pi}^i,\boldsymbol{\tau}^i\}_{i=1}^{t-1} \cup \{\boldsymbol{\pi}^t\}\right]. \text{ We denote } \mathbb{E}_t[\cdot] = \mathbb{E}\left[\cdot \mid \{\boldsymbol{\pi}^i,\boldsymbol{\tau}^i\}_{i=1}^{t-1} \cup \{\boldsymbol{\pi}^t\}\right].$ then have

$$\mathbb{E}\left[\exp\left(\sum_{i=1}^{t}\log\left(\tilde{f}_{m}^{i}\left(\bar{\theta}_{m},\pi_{m}^{i}\right)/\tilde{f}_{m}^{i}\left(\theta_{m}^{\star},\pi_{m}^{i}\right)\right)\right)\right]$$
$$=\mathbb{E}\left[\exp\left(\sum_{i=1}^{t-1}\log\left(\tilde{f}_{m}^{i}\left(\bar{\theta}_{m},\pi_{m}^{i}\right)/\tilde{f}_{m}^{i}\left(\theta_{m}^{\star},\pi_{m}^{i}\right)\right)\right) \cdot \mathbb{E}_{t}\left[\exp\left(\log\left(\tilde{f}_{m}^{t}(\bar{\theta}_{m},\pi_{m}^{t})/\tilde{f}_{m}^{t}\left(\theta_{m}^{\star},\pi_{m}^{t}\right)\right)\right)\right]\right]$$
$$=\mathbb{E}\left[\exp\left(\sum_{i=1}^{t-1}\log\left(\tilde{f}_{m}^{i}\left(\bar{\theta}_{m},\pi_{m}^{i}\right)/\tilde{f}_{m}^{i}\left(\theta_{m}^{\star},\pi_{m}^{i}\right)\right)\right) \cdot \mathbb{E}_{t}\left(\tilde{f}_{m}^{t}(\bar{\theta}_{m},\pi_{m}^{t})/\tilde{f}_{m}^{t}\left(\theta_{m}^{\star},\pi_{m}^{t}\right)\right)\right)\right]$$
(7)

According to the Definition C.1, we further have

$$\mathbb{E}_{t}\left(\frac{\tilde{f}_{m}^{t}(\bar{\theta}_{m},\pi_{m}^{t})}{\tilde{f}_{m}^{t}(\theta_{m}^{\star},\pi_{m}^{t})}\right) = \mathbb{E}_{t}\left[\sum_{\boldsymbol{\tau}_{H}}f_{m}\left(\bar{\theta}_{m},\pi_{m}^{t}\right)\left(\prod_{j=1}^{m-1}f_{j}\left(\theta_{j}^{\star},\pi_{j}^{t}\right)\right)\left(\prod_{j=m+1}^{n}f_{j}\left(\theta_{j}^{t},\pi_{j}^{t}\right)\right)\right] \leq \left(1+\frac{1}{T}\right)$$

$$(8)$$

We insert eq. (8) back into eq. (7), and we can obtain that

$$\mathbb{E}\left[\exp\left(\sum_{i=1}^{t}\log\left(\tilde{f}_{m}^{i}\left(\bar{\theta}_{m},\pi_{m}\right)/\tilde{f}_{m}^{i}\left(\theta_{m}^{\star},\pi_{m}\right)\right)\right)\right] \leq e$$

We then use Markov inequality and take a union bound for all  $(\theta_m, t) \in \Theta_m \times [T]$  and  $m \in [n]$ , and we can conclude that the following event holds with probability at least  $1 - \delta$  for all  $m \in [n]$ :

$$\max_{(\bar{\theta}_m,t)\in\bar{\Theta}_m\times[T]}\sum_{i=1}^t \log\left(\tilde{f}^i_m(\bar{\theta}_m,\pi_m)/\tilde{f}^i_m(\theta^\star_m,\pi_m)\right) \leq \beta_m,$$

According to the definition of optimistic discretization, we obtain that the following inequality holds with probability at least  $1 - \delta$  for all  $m \in [n]$ :

914  
915  
916  
917  

$$\max_{(\theta_m,t)\in\bar{\Theta}_m\times[T]}\sum_{i=1}^t \log\left(\tilde{f}^i_m(\theta_m,\pi_m)/\tilde{f}^i_m(\theta^\star_m,\pi_m)\right) \leq \beta_m,$$

**18 Lemma C.2.** There exists a universal constant c such that for any  $\delta \in (0, 1]$ , with probability at least for all  $t \in [T]$  and all  $\theta_m \in \Theta_m$ ,  $m \in [n]$ , it holds that

$$\sum_{i=1}^{t} \left( \sum_{\boldsymbol{\tau}} \left| f_m\left(\theta_m, \pi_m^i\right) - f_m\left(\theta_m^{\star}, \pi_m^i\right) \right| \left[ \prod_{l=1}^{m-1} f_l\left(\theta_l^{\star}, \pi_l^i\right) \right] \left[ \prod_{j=m+1}^{n} f_l\left(\theta_j^i, \pi_j^i\right) \right] \right)^2$$
$$\lesssim \left( \sum_{i=1}^{t} \log\left( \tilde{f}_m^i\left(\theta_m^{\star}, \pi_m^i\right) / \tilde{f}_m^i\left(\theta_m, \pi_m^i\right) \right) + \beta_m \right)$$

*Proof.* We define tangent trajectory sample  $\hat{\tau}^i$  that satisfies  $\hat{\tau}^i \sim \prod_{i=1}^m \mathbb{P}_{\theta*,i}^{\pi_i^k} \prod_{j=m+1}^n \mathbb{P}_{\theta_j^k,j}^{\pi_j^k}$  but are independent with  $\tau^i$ . With similar analysis as Lemma 15 of Liu et al. (2022a), we obtain that

$$\mathbb{E}\left[\exp\left(\sum_{i=1}^{t} \frac{1}{2} \log\left(\frac{\tilde{f}_{m}^{i}\left(\bar{\theta}_{m}, \pi_{m}^{i}\right)}{\tilde{f}_{m}^{i}\left(\theta_{m}^{\star}, \pi_{m}^{i}\right)}\right) - \log\mathbb{E}\left[\exp\left(\frac{1}{2} \log\left(\frac{\hat{f}_{m}^{i}\left(\bar{\theta}_{m}, \pi_{m}^{i}\right)}{\hat{f}_{m}^{i}\left(\theta_{m}^{\star}, \pi_{m}^{i}\right)}\right)\right) \middle| \mathcal{E}_{m}\right]\right)\right] = 1,$$

where for all  $(\theta_m, m) \in \Theta_m \times [n]$ , we denote  $\widehat{f}_m^i(\theta_m, \pi_m^i)$  as  $\widehat{f}_m^i(\theta_m, \pi_m^i) = \mathbb{P}_{\theta_m, m}^{\pi_m}(\tau_m^i | \{\tau_r^i\}_{r \in \mathsf{pa}(m)})$ , and we denote  $\mathcal{E}_m$ ,  $\widehat{\mathcal{E}}_m$  as  $\mathcal{E}_m = \{(\pi^i, \tau^i)\}_{i=1}^t$ ,  $\widehat{\mathcal{E}}_m = \{(\pi^i, \widehat{\tau}^i)\}_{i=1}^t$ . With Chernoff's method, we can obtain that with probability at least  $1 - \delta$ , for all  $\overline{\theta}_m \in \overline{\Theta}_m$  we have

$$-\log \mathbb{E}_{\widehat{\mathcal{E}}_m}\left[\exp\left(\sum_{i=1}^t \frac{1}{2} \log\left(\frac{\widehat{f}_m^i\left(\bar{\theta}_m, \pi_m^i\right)}{\widehat{f}_m^i\left(\theta_m^\star, \pi_m^i\right)}\right)\right) \middle| \mathcal{E}_m\right] \le -\sum_{i=1}^t \frac{1}{2} \log\left(\frac{\widehat{f}_m^i\left(\bar{\theta}_m, \pi_m^i\right)}{\widehat{f}_m^i\left(\theta_m^\star, \pi_m^i\right)}\right) + \beta_m.$$
(9)

Then, we apply elementary inequality  $-\log x \ge 1 - x$ , and we can obtain that

$$-\log \mathbb{E}_{\widehat{\mathcal{E}}_{m}}\left[\exp\left(\sum_{i=1}^{t}\frac{1}{2}\log\left(\frac{\widehat{f}_{m}^{i}\left(\bar{\theta}_{m},\pi_{m}^{i}\right)}{\widehat{f}_{m}^{i}\left(\theta_{m}^{\star},\pi_{m}^{i}\right)}\right)\right) \middle| \mathcal{E}_{m}\right]$$

$$=-\sum_{i=1}^{t}\log \mathbb{E}_{\boldsymbol{\tau}\sim\left(\prod_{l=1}^{m}\mathbb{P}_{\theta_{l}^{\star},l}^{\pi_{l}^{i}}\right)\left(\prod_{j=m+1}^{n}\mathbb{P}_{\theta_{j}^{j},j}^{\pi_{j}^{i}}\right)\left[\sqrt{\frac{f_{m}\left(\bar{\theta}_{m},\pi_{m}^{i}\right)}{f_{m}\left(\theta_{m}^{\star},\pi_{m}^{i}\right)}}\right]$$

$$\geq \sum_{i=1}^{t}\left(1-\mathbb{E}_{\boldsymbol{\tau}\sim\left(\prod_{l=1}^{m}\mathbb{P}_{\theta_{l}^{\star},l}^{\pi_{l}^{i}}\right)\left(\prod_{j=m+1}^{n}\mathbb{P}_{\theta_{j}^{j},j}^{\pi_{j}^{i}}\right)\left[\sqrt{\frac{f_{m}\left(\bar{\theta}_{m},\pi_{m}^{i}\right)}{f_{m}\left(\theta_{m}^{\star},\pi_{m}^{i}\right)}}\right]\right)$$

$$=\sum_{i=1}^{t}\left(1-\sum_{\boldsymbol{\tau}}\sqrt{f_{m}\left(\bar{\theta}_{m},\pi_{m}^{i}\right)\cdot f_{m}\left(\theta_{m}^{\star},\pi_{m}^{i}\right)}\left[\prod_{l=1}^{m-1}f_{l}\left(\theta_{l}^{\star},\pi_{l}^{i}\right)\right]\left[\prod_{j=m+1}^{n}f_{j}\left(\theta_{j}^{i},\pi_{j}^{i}\right)\right]\right).$$

To continue, we aim to achieve the lower bound for the following term of interest:

$$-\log \mathbb{E}_{\widehat{\mathcal{E}}_m}\left[\exp\left(\sum_{i=1}^t \frac{1}{2} \log\left(\frac{\widehat{f}_m^i(\bar{\theta}_m, \pi_m^i)}{\widehat{f}_m^i(\theta_m^\star, \pi_m^i)}\right)\right) \middle| \mathcal{E}_m\right] + \frac{1}{2}.$$

We have the following inequalities:

$$-\log \mathbb{E}_{\widehat{\mathcal{E}}_{m}}\left[\exp\left(\sum_{i=1}^{t}\frac{1}{2}\log\left(\frac{\widehat{f}_{m}^{i}\left(\bar{\theta}_{m},\pi_{m}^{i}\right)}{\widehat{f}_{m}^{i}\left(\theta_{m}^{\star},\pi_{m}^{i}\right)}\right)\right) \middle| \mathcal{E}_{m}\right] + \frac{1}{2}$$

$$\geq \sum_{i=1}^{t}\left(1 - \sum_{\tau}\sqrt{f_{m}\left(\bar{\theta}_{m}^{-},\pi_{m}^{i}\right) \cdot f_{m}\left(\theta_{m}^{\star},\pi_{m}^{i}\right)}\left[\prod_{l=1}^{m-1}f_{l}\left(\theta_{l}^{\star},\pi_{l}^{i}\right)\right]\left[\prod_{j=m+1}^{n}f_{j}\left(\theta_{j}^{i},\pi_{j}^{i}\right)\right]\right) + \frac{1}{2}$$

$$\geq \frac{1}{2}\sum_{i=1}^{t}\sum_{\tau}\left(\sqrt{\left[\prod_{l=1}^{m-1}f_{l}\left(\theta_{l}^{\star},\pi_{l}^{i}\right)\right]f_{m}\left(\bar{\theta}_{m},\pi_{m}^{i}\right)}\left[\prod_{j=m+1}^{n}f_{j}\left(\theta_{j}^{i},\pi_{j}^{i}\right)\right] - \sqrt{\left[\prod_{l=1}^{m}f_{l}\left(\theta_{l}^{\star},\pi_{l}^{i}\right)\right]\left[\prod_{j=m+1}^{n}f_{j}\left(\theta_{j}^{i},\pi_{j}^{i}\right)\right]}\right)^{2}$$

With elementary calculation, we have 

$$-\log \mathbb{E}_{\widehat{\mathcal{E}}_{m}}\left[\exp\left(\sum_{i=1}^{t}\frac{1}{2}\log\left(\frac{\widehat{f}_{m}^{i}\left(\bar{\theta}_{m},\pi_{m}^{i}\right)}{\widehat{f}_{m}^{i}\left(\theta_{m}^{\star},\pi_{m}^{i}\right)}\right)\right) \middle| \mathcal{E}_{m}\right] + \frac{1}{2}$$

$$\geq \frac{1}{12}\sum_{i=1}^{t}\left[\sum_{\tau}\left(\sqrt{\left[\prod_{l=1}^{m-1}f_{l}\left(\theta_{l}^{\star},\pi_{l}^{i}\right)\right]f_{m}\left(\bar{\theta}_{m},\pi_{m}^{i}\right)\left[\prod_{j=m+1}^{n}f_{j}\left(\theta_{j}^{i},\pi_{j}^{i}\right)\right]} - \sqrt{\left[\prod_{l=1}^{m}f_{l}\left(\theta_{l}^{\star},\pi_{l}^{i}\right)\right]\left[\prod_{j=m+1}^{n}f_{j}\left(\theta_{j}^{i},\pi_{j}^{i}\right)\right]}\right)^{2}\right]$$

$$\cdot\left[\sum_{\tau}\left(\sqrt{\left[\prod_{l=1}^{m-1}f_{l}\left(\theta_{l}^{\star},\pi_{l}^{i}\right)\right]f_{m}\left(\bar{\theta}_{m},\pi_{m}^{i}\right)\left[\prod_{j=m+1}^{n}f_{j}\left(\theta_{j}^{i},\pi_{j}^{i}\right)\right]} + \sqrt{\left[\prod_{l=1}^{m}f_{l}\left(\theta_{l}^{\star},\pi_{l}^{i}\right)\right]\left[\prod_{j=m+1}^{n}f_{j}\left(\theta_{j}^{i},\pi_{j}^{i}\right)\right]}\right)^{2}\right]$$

We then apply Cauchy-Schawarz inequality, and we arrive at

$$-\log \mathbb{E}_{\widehat{\mathcal{E}}_{m}}\left[\exp\left(\sum_{i=1}^{t}\frac{1}{2}\log\left(\frac{\widehat{f}_{m}^{i}\left(\bar{\theta}_{m},\pi_{m}^{i}\right)}{\widehat{f}_{m}^{i}\left(\theta_{m}^{\star},\pi_{m}^{i}\right)}\right)\right) \middle| \mathcal{E}_{m}\right] + \frac{1}{2}$$

$$\geq \frac{1}{12}\sum_{i=1}^{t}\left(\sum_{\boldsymbol{\tau}}\left|f_{m}\left(\bar{\theta}_{m},\pi_{m}^{i}\right) - f_{m}\left(\theta_{m}^{\star},\pi_{m}^{i}\right)\right| \left[\prod_{l=1}^{m-1}f_{l}\left(\theta_{l}^{\star},\pi_{l}^{i}\right)\right] \left[\prod_{j=m+1}^{n}f_{j}\left(\theta_{j}^{i},\pi_{j}^{i}\right)\right]\right)^{2} - \frac{1}{2}$$

$$(10)$$

We insert eq. (10) back into eq. (9), and we obtain that there exist a universal constant c such that for any  $\delta \in (0, 1]$ , with probability at least  $1 - \delta$  for all  $t \in [T]$ ,  $m \in [n]$ , and all  $\theta_m \in \Theta_m$ , it holds that

$$\sum_{i=1}^{t} \left( \sum_{\boldsymbol{\tau}} \left| f_m\left(\theta_m, \pi_m^i\right) - f_m\left(\theta_m^{\star}, \pi_m^i\right) \right| \left[ \prod_{l=1}^{m-1} f_l\left(\theta_l^{\star}, \pi_l^i\right) \right] \left[ \prod_{j=m+1}^{n} f_l\left(\theta_j^i, \pi_j^i\right) \right] \right)^2$$
$$\lesssim \left( \sum_{i=1}^{t} \log\left( \tilde{f}_m^i\left(\theta_m^{\star}, \pi_m^i\right) / \tilde{f}_m^i\left(\theta_m, \pi_m^i\right) \right) + \beta_m \right)$$

We combine this result with the update rule of Algorithm 4.2.

**Corollary C.1.** With probability at least  $1 - \delta$ , for all  $k \in [K]$ ,  $m \in [n]$ , the following inequality holds. 

$$\begin{array}{l} \begin{array}{l} \text{1011} \\ \text{1012} \\ \text{1013} \\ \text{1014} \end{array} & \sum_{t=1}^{k-1} \sum_{\{\tau_{H,r}\}_{r \in [n]}} \left| f_m\left(\theta_m^k, \pi_m^t\right) - f_m\left(\theta_m^\star, \pi_m^t\right) \right| \left[ \prod_{l=1}^{m-1} f_l\left(\theta_l^\star, \pi_l^t\right) \right] \left[ \prod_{j=m+1}^n f_j\left(\theta_j^t, \pi_j^t\right) \right] \lesssim \sqrt{\beta_m k}. \end{array}$$

According to Lemma C.2 and Lemma D.1, with probability at least  $1 - \delta$ , for all  $k \in [K]$ ,  $m \in [n]$ , the following inequality holds.

$$\begin{array}{l} \begin{array}{l} \begin{array}{l} 1018\\ 1019\\ 1020\\ 1021 \end{array} & \sum_{t=1}^{k-1} \left( \sum_{\{\tau_{H,r}\}_{r\in[n]}} \left| f_m\left(\theta_m^k, \pi_m^t\right) - f_m\left(\theta_m^\star, \pi_m^t\right) \right| \left[ \prod_{l=1}^{m-1} f_l\left(\theta_l^\star, \pi_l^t\right) \right] \left[ \prod_{j=m+1}^n f_j\left(\theta_j^t, \pi_j^t\right) \right] \right)^2 \lesssim \beta_m. \end{array}$$

We apply Cauchy-Schawarz inequality, and we can obtain that

1023  
1024 
$$\sum_{t=1}^{k-1} \sum_{\{\tau_{H,r}\}_{r\in[n]}} \left| f_m\left(\theta_m^k, \pi_m^t\right) - f_m\left(\theta_m^\star, \pi_m^t\right) \right| \left[ \prod_{l=1}^{m-1} f_l\left(\theta_l^\star, \pi_l^t\right) \right] \left[ \prod_{j=m+1}^n f_j\left(\theta_j^t, \pi_j^t\right) \right] \lesssim \sqrt{\beta_m k}.$$

**Definition C.2.** For all  $(m, h, k) \in [n] \times [H] \times [T]$ , we define the matrix notations  $\mathbb{M}_{h,m}^{\star} \in \mathbb{R}^{O \times O}$ and  $\mathbb{M}_{h.m}^k \in \mathbb{R}^{O \times O}$  as follows:  $\mathbb{M}_{0,m}^{\star} = \mathbb{O}_{1,m}^{\star} \mu_m^{\star} \in \mathbb{R}^O, \quad \mathbb{M}_{0,m}^k = \mathbb{O}_{1,m}^k \mu_m^k \in \mathbb{R}^O,$  $\mathbb{M}_{h,m}^{\star}(o_m, a_m, \{a_r\}_{r \in \mathsf{pa}(m)}) = \mathbb{O}_{h+1,m}^{\star} \mathbb{T}_{h,m,a_m,\{a_r\}_{r \in \mathsf{pa}(m)}}^{\star} \cdot \operatorname{diag}\left(\mathbb{O}_{h,m}^{\star}(o_m \mid \cdot)\right) \left(\mathbb{O}_{h,m}^{\star}\right)^{\dagger} \in \mathbb{R}^{O \times O},$  $\mathbb{M}_{h,m}^k(o_m, a_m, \{a_r\}_{r \in \mathsf{pa}(m)}) = \mathbb{O}_{h+1,m}^k \mathbb{T}_{h,m,a_m,\{a_r\}_{r \in \mathsf{pa}(m)}}^k \cdot diag\left(\mathbb{O}_{h,m}^k(o_m \mid \cdot)\right) \left(\mathbb{O}_{h,m}^k\right)^{\dagger} \in \mathbb{R}^{O \times O},$ where  $\{\mathbb{O}_{h,m}^{\star}\}_{(h,m)\in[H]\times[n]}$  and  $\{\mathbb{T}_{h,m,a_m,\{a_r\}_{r\in pa(m)}}^{\star}\}_{(h,m)\in[H]\times[n]}$  denote the observation and transition matrices corresponding to the true transition model, and  $\{\mathbb{O}_{h,m}^k\}_{(h,m)\in[H]\times[n]}$  and  $\{\mathbb{T}_{h,m,a_m,\{a_r\}_{r\in \mathsf{pa}(m)}}^k\}_{(h,m)\in[H]\times[n]}$  denote the observation and transition matrices correspond-ing to model parameter  $\theta_m^k$  for all  $k \in [T]$ . When no confusion arises, we simplify the notation by using  $\mathbb{M}_{h,m}^{\star}$  to represent  $\mathbb{M}_{h,m}^{\star}(o_m, a_m, \{a_r\}_{r \in \mathsf{pa}(m)})$  and  $\mathbb{M}_{h,m}^k$  to represent  $\mathbb{M}_{h}^{k}(o_{m}, a_{m}, \{a_{r}\}_{r \in \mathsf{pa}(m)}).$ 

Since marginalizing two distribution will not increase their TV distance, so for all  $(k,h) \in [K] \times$ [H-1], we have the following corollary. 

**Corollary C.2.** With probability at least  $1-\delta$ , for all  $(k,h) \in [T] \times [H-1]$ , the following inequality holds true. 

$$\sum_{t=1}^{k-1} \sum_{\boldsymbol{\tau}_h} \pi_m^t(\boldsymbol{\tau}_{h,m}) \left\| \left[ \prod_{h'=0}^h \mathbb{M}_{h',m}^k \right] - \left[ \prod_{h'=0}^h \mathbb{M}_{h',m}^\star \right] \right\|_1 \left[ \prod_{l=1}^{m-1} \left\| \prod_{h'=0}^h \mathbb{M}_{h',l}^\star \right\|_1 \pi_l^t(\boldsymbol{\tau}_{h,l}) \right] \\ \cdot \left[ \prod_{j=m+1}^n \left\| \prod_{h'=0}^h \mathbb{M}_{h',j}^t \right\|_1 \pi_j^t(\boldsymbol{\tau}_{h,j}) \right] \lesssim \sqrt{k\beta_m}.$$

**Lemma C.3.** With probability at least  $1 - \delta$ , for all  $(k, h, m) \in [T] \times [H - 1] \times [n]$ ,

$$\sum_{t=1}^{k-1} \sum_{\tau_h} \pi_m^t(\tau_{h,m}) \cdot \left\| \left( \mathbb{M}_{h,m}^k - \mathbb{M}_{h,m}^\star \right) \left[ \prod_{h'=0}^{h-1} \mathbb{M}_{h',m}^\star \right] \right\|_1 \cdot \left[ \prod_{l=1}^{m-1} \left\| \prod_{h'=0}^h \mathbb{M}_{h',l}^\star \right\|_1 \cdot \pi_l^t(\tau_{h,l}) \right] \\ \cdot \left[ \prod_{j=m+1}^n \left\| \prod_{h'=0}^h \mathbb{M}_{h',j}^\star \right\|_1 \cdot \pi_j^t(\tau_{h,j}) \right] \lesssim \sqrt{Sk\beta_m} / \alpha.$$

$$\begin{array}{l} 1062\\ 1063\\ 1064\\ 1065 \end{array} = \sum_{t=1}^{k-1} \sum_{\tau_h} \pi_m^t(\tau_{h,m}) \left\| \left( \mathbb{M}_{h,m}^k - \mathbb{M}_{h,m}^\star \right) \left[ \prod_{h'=0}^{h-1} \mathbb{M}_{h',m}^\star \right] \right\|_1 \left[ \prod_{l=1}^{m-1} \left\| \prod_{h'=0}^h \mathbb{M}_{h',l}^\star \right\|_1 \pi_l^t(\tau_{h,l}) \right] \left[ \prod_{j=m+1}^n \left\| \prod_{h'=0}^h \mathbb{M}_{h',j}^\star \right\|_1 \pi_j^t(\tau_{h,m}) \right\|_1 \pi_j^t(\tau_{h,m}) \left\| \prod_{j=m+1}^n \mathbb{M}_{h',j}^\star \right\|_1 \pi_j^t(\tau_{h',j}) \right\|_1 \pi_j^t(\tau_{h',j}) \left\| \prod_{j=m+1}^n \mathbb{M}_{h',j}^\star \right\|_1 \pi_j^t(\tau_{h',j}) \right$$

We initially have the following decomposition:

*Proof.* We intend to bound the following term:

According to the result in Corollary C.2, we obtain that

$$\sum_{t=1}^{k-1} \sum_{\tau_h} \pi_m^t(\tau_{h,m}) \left\| \left[ \prod_{h'=0}^h \mathbb{M}_{h',m}^k \right] - \left[ \prod_{h'=0}^h \mathbb{M}_{h',m}^\star \right] \right\|_1 \\ \cdot \left[ \prod_{l=1}^{m-1} \left\| \prod_{h'=0}^h \mathbb{M}_{h',l}^\star \right\|_1 \pi_l^t(\tau_{h,l}) \right] \left[ \prod_{j=m+1}^n \left\| \prod_{h'=0}^h \mathbb{M}_{h',j}^\star \right\|_1 \pi_j^t(\tau_{h,j}) \right] \lesssim \sqrt{k\beta_m}.$$

$$(11)$$

According to the definition of matrix operator, we have 

$$\sum_{t=1}^{k-1} \sum_{\tau_h} \pi_m^t(\tau_{h,m}) \left\| \mathbb{M}_{h,m}^k \left( \left[ \prod_{h'=0}^{h-1} \mathbb{M}_{h',m}^k \right] - \left[ \prod_{h'=0}^{h-1} \mathbb{M}_{h',m}^\star \right] \right) \right\|_1 \\ \cdot \left[ \prod_{l=1}^{m-1} \left\| \prod_{h'=0}^h \mathbb{M}_{h',l}^\star \right\|_1 \pi_l^t(\tau_{h,l}) \right] \left[ \prod_{j=m+1}^n \left\| \prod_{h'=0}^h \mathbb{M}_{h',j}^\star \right\|_1 \pi_j^t(\tau_{h,j}) \right] \lesssim \sqrt{Sk\beta_m} / \alpha.$$
(12)  
We combine eq. (11) and eq. (12), and we enventually arrive at:

$$\sum_{t=1}^{k-1} \sum_{\tau_h} \pi_m^t(\tau_{h,m}) \cdot \left\| \left( \mathbb{M}_{h,m}^k - \mathbb{M}_{h,m}^\star \right) \left[ \prod_{h'=0}^{h-1} \mathbb{M}_{h',m}^\star \right] \right\|_1 \cdot \left[ \prod_{l=1}^{m-1} \left\| \prod_{h'=0}^h \mathbb{M}_{h',l}^\star \right\|_1 \cdot \pi_l^t(\tau_{h,l}) \right] \\ \cdot \left[ \prod_{j=m+1}^n \left\| \prod_{h'=0}^h \mathbb{M}_{h',j}^\star \right\|_1 \cdot \pi_j^t(\tau_{h,j}) \right] \lesssim \sqrt{Sk\beta_m} / \alpha,$$

holds with probability at least  $1 - \delta$  for all  $(k, h) \in [T] \times [H - 1]$ .

**Lemma C.4.** The regret is bounded by the following inequality:

$$Regret(k) = \sum_{t=1}^{k} V^{\boldsymbol{\pi}^{*}} - V^{\boldsymbol{\pi}^{t}} \leq nH \sum_{t=1}^{k} \sum_{\boldsymbol{\tau}_{H}} \left| \mathbb{P}_{\boldsymbol{\theta}^{t}}^{\boldsymbol{\pi}^{t}}(\boldsymbol{\tau}_{H}) - \mathbb{P}_{\boldsymbol{\theta}^{*}}^{\boldsymbol{\pi}^{t}}(\boldsymbol{\tau}_{H}) \right|,$$

where we define 

$$\mathbb{P}_{\theta^{t}}^{\pi^{t}}(\tau_{H}) = \prod_{m=1}^{n} \mathbb{P}_{\theta_{m}^{t}}^{\pi^{t}_{m}}(\tau_{H,m} \mid \{\tau_{H,r}\}_{r \in \mathsf{pa}(m)}), \quad \mathbb{P}_{\theta^{*}}^{\pi^{t}}(\tau_{H}) = \prod_{m=1}^{n} \mathbb{P}_{\theta_{m}^{*}}^{\pi^{t}_{m}}(\tau_{H,m} \mid \{\tau_{H,r}\}_{r \in \mathsf{pa}(m)}).$$

 $\leq nH\sum_{t=1}^{k}\sum_{m=1}^{n}\sum_{h=1}^{H-1}\sum_{\{\tau_{h,r}\}_{r\in[n]}}\frac{\sqrt{S}}{\alpha}\cdot\pi_{m}^{t}(\tau_{h,m})\cdot\left\|\left(\mathbb{M}_{h,m}^{t}-\mathbb{M}_{h,m}^{\star}\right)\left[\prod_{h'=0}^{h}\mathbb{M}_{h',m}^{\star}\right]\right\|_{1}$ 

 $\cdot \left[\prod_{l=1}^{m-1} \left\|\prod_{h'=0}^{h} \mathbb{M}_{h',l}^{\star}\right\|_{1} \cdot \pi_{l}^{t}(\tau_{h,l})\right] \cdot \left[\prod_{j=m+1}^{n} \left\|\prod_{h'=0}^{h} \mathbb{M}_{h',j}^{t}\right\|_{1} \cdot \pi_{j}^{t}(\tau_{h,j})\right]$ 

Proof. We can strightforwardly achieve this result according to the definition of value function and regret. 

**Lemma C.5.** The regret is bounded by the following inequality: 

*Proof.* According to the definition of transition model factorization, we have

$$\sum_{t=1}^k \sum_{oldsymbol{ au}_H} \left| \mathbb{P}_{oldsymbol{ heta}^t}^{oldsymbol{\pi}^t}(oldsymbol{ au}_H) - \mathbb{P}_{oldsymbol{ heta}^*}^{oldsymbol{\pi}^t}(oldsymbol{ au}_H) 
ight|$$

 $Regret(k) = \sum_{i=1}^{n} V^{\pi^*} - V^{\pi^t}$ 

$$= \sum_{t=1}^{k} \sum_{\tau_{H}} \left| \prod_{m=1}^{n} \mathbb{P}_{\theta_{m}^{t}}^{\pi_{m}^{t}} \left( \tau_{H,m} \mid \{\tau_{H,r}\}_{r \in \mathsf{pa}(m)} \right) - \prod_{m=1}^{n} \mathbb{P}_{\theta_{m}^{*}}^{\pi_{m}^{t}} (\tau_{H,m} \left( \tau_{H,m} \mid \{\tau_{H,r}\}_{r \in \mathsf{pa}(m)} \right) \right) \right|$$

Moreover, we have the following inequality of difference between transition probability measure.

1128 Moreover, we have the following metalling 
$$\sum_{t=1}^{k} \sum_{\tau_{H}} \left| \mathbb{P}_{\theta^{t}}^{\pi^{t}}(\tau_{H}) - \mathbb{P}_{\theta^{*}}^{\pi^{t}}(\tau_{H}) \right|$$
1130

$$\frac{1}{t=1} \frac{1}{\tau_H}$$

$$= \sum_{t=1}^{k} \sum_{\tau_{H}} \sum_{m=1}^{n} \left[ \prod_{j=1}^{m-1} f_{j} \left( \theta_{j}^{\star}, \pi_{j}^{t} \right) \right] \left| f_{m} \left( \theta_{m}^{t}, \pi_{m}^{t} \right) - f_{m} \left( \theta_{m}^{\star}, \pi_{m}^{t} \right) \right| \left[ \prod_{l=m+1}^{n} f_{l} \left( \theta_{l}^{t}, \pi_{l}^{t} \right) \right]$$
(13)

Moreover, according to the definition of matrix notation, we can rewrite the term  $|f_m(\theta_m^t, \pi_m^t)|$  $f_m(\theta_m^\star, \pi_m^t)$  as 

$$\left|f_m\left(\theta_m^t, \pi_m^t\right) - f_m\left(\theta_m^\star, \pi_m^t\right)\right| = \left\|\left[\prod_{h'=0}^{H-1} \mathbb{M}_{h,m}^t\right] - \left[\prod_{h'=0}^{H-1} \mathbb{M}_{h,m}^\star\right]\right\|_1 \pi_m^t(\tau_{H,m})$$

Then we have the following inequality: 

We insert eq. (14) back into eq. (13), and we can obtain that

$$\begin{split} & \stackrel{1147}{1148} & \sum_{\tau_{H}} \left[ \prod_{j=1}^{m-1} f_{j}\left(\theta_{j}^{\star}, \pi_{j}^{t}\right) \right] \left| f_{m}\left(\theta_{m}^{t}, \pi_{m}^{t}\right) - f_{m}\left(\theta_{m}^{\star}, \pi_{m}^{t}\right) \right| \left[ \prod_{l=m+1}^{n} f_{l}\left(\theta_{l}^{t}, \pi_{l}^{t}\right) \right] \\ & \stackrel{1151}{1152} & \leq \sum_{\tau_{H}} \sum_{j=1}^{H-1} \left\| \left[ \prod_{h=j+1}^{H-1} \mathbb{M}_{h,m}^{t} \right] \left( \mathbb{M}_{j,m}^{t} - \mathbb{M}_{j,m}^{\star} \right) \left[ \prod_{h=0}^{j-1} \mathbb{M}_{h,m}^{\star} \right] \right\|_{1} \pi_{m}^{t}(\tau_{H,m}) \left[ \prod_{j=1}^{m-1} f_{j}\left(\theta_{j}^{\star}, \pi_{j}^{t}\right) \right] \left[ \prod_{l=m+1}^{n} f_{l}\left(\theta_{l}^{t}, \pi_{l}^{t}\right) \right] \\ & \stackrel{1154}{1155} & \leq \sum_{h=1}^{H-1} \sum_{\tau_{h}} \frac{\sqrt{S}}{\alpha} \left\| \left( \mathbb{M}_{h,m}^{t} - \mathbb{M}_{h,m}^{\star} \right) \left[ \prod_{h'}^{h-1} \mathbb{M}_{h',m}^{\star} \right] \right\|_{1} \cdot \pi_{m}^{t}(\tau_{H,m}) \\ & \quad \cdot \left[ \prod_{j=1}^{m-1} \left[ \prod_{h'=0}^{h} \mathbb{M}_{j,m}^{\star} \right] \pi_{j}^{t}(\tau_{h,j}) \right] \cdot \left[ \prod_{l=m+1}^{n} \left[ \prod_{h'=0}^{h} \mathbb{M}_{l,m}^{t} \right] \pi_{l}^{t}(\tau_{h,l}) \right] \\ & \quad 1160 \end{split}$$

Eventually, the target regret can be bounded with 

$$\begin{aligned} & \text{Regret}(k) \leq \sum_{t=1}^{k} \sum_{\tau_{H}} \left| \mathbb{P}_{\theta^{t}}^{\pi^{t}}(\tau_{H}) - \mathbb{P}_{\theta^{*}}^{\pi^{t}}(\tau_{H}) \right| \\ & = \sum_{t=1}^{k} \sum_{\tau_{H}} \left| \prod_{m=1}^{n} \mathbb{P}_{\theta^{t}_{m}}^{\pi^{t}}(\tau_{H,m} | \tau_{H,m-1}) - \prod_{m=1}^{n} \mathbb{P}_{\theta^{*}_{m}}^{\pi^{t}_{m}}(\tau_{H,m} | \tau_{H,m-1}) \right| \\ & = \sum_{t=1}^{k} \sum_{\tau_{H}} \left| \prod_{m=1}^{n} \mathbb{P}_{\theta^{t}_{m}}^{\pi^{t}_{m}}(\tau_{H,m} | \tau_{H,m-1}) - \prod_{m=1}^{n} \mathbb{P}_{\theta^{*}_{m}}^{\pi^{t}_{m}}(\tau_{H,m} | \tau_{H,m-1}) \right| \\ & \leq \sum_{t=1}^{k} \sum_{\tau_{H}} \sum_{m=1}^{n} \left[ \prod_{j=1}^{m-1} f_{j}\left(\theta^{*}_{j}, \pi^{t}_{j}\right) \right] \left| f_{m}\left(\theta^{t}_{m}, \pi^{t}_{m}\right) - f_{m}\left(\theta^{*}_{m}, \pi^{t}_{m}\right) \right| \left[ \prod_{l=m+1}^{n} f_{l}\left(\theta^{t}_{l}, \pi^{t}_{l}\right) \right] \\ & \leq nH \sum_{t=1}^{k} \sum_{m=1}^{n} \sum_{h=1}^{H-1} \sum_{\{\tau_{h,r}\}_{r\in[n]}} \frac{\sqrt{S}}{\alpha} \cdot \pi^{t}_{m}(\tau_{h,m}) \cdot \left\| \left(\mathbb{M}_{h,m}^{t} - \mathbb{M}_{h,m}^{*}\right) \left[ \prod_{h'=0}^{h} \mathbb{M}_{h',m}^{*} \right] \right\|_{1} \\ & \cdot \left[ \prod_{l=1}^{m-1} \left\| \prod_{h'=0}^{h} \mathbb{M}_{h',l}^{*} \right\|_{1} \cdot \pi^{t}_{l}(\tau_{h,l}) \right] \cdot \left[ \prod_{j=m+1}^{n} \left\| \prod_{h'=0}^{h} \mathbb{M}_{h',j}^{t} \right\|_{1} \cdot \pi^{t}_{j}(\tau_{h,j}) \right] \\ & \prod_{l=1}^{H} \left\| \prod_{l=1}^{h} \mathbb{M}_{h',l}^{*} \right\|_{1} \cdot \pi^{t}_{l}(\tau_{h,l}) \right] \cdot \left[ \prod_{l=1}^{n} \left\| \prod_{l=1}^{h} \mathbb{M}_{h',l}^{t} \right\|_{1} \cdot \pi^{t}_{j}(\tau_{h,l}) \right] \\ & \prod_{l=1}^{H} \left\| \prod_{l=1}^{h} \mathbb{M}_{h',l}^{*} \right\|_{1} \cdot \pi^{t}_{l}(\tau_{h,l}) \right] \cdot \left[ \prod_{l=1}^{n} \mathbb{M}_{h',l}^{*} \right\|_{1} \cdot \pi^{t}_{j}(\tau_{h,l}) \right] \\ & \prod_{l=1}^{H} \left\| \prod_{l=1}^{h} \mathbb{M}_{h',l}^{*} \right\|_{1} \cdot \pi^{t}_{l}(\tau_{h,l}) \right] \cdot \left[ \prod_{l=1}^{n} \mathbb{M}_{h',l}^{*} \right\|_{1} \cdot \pi^{t}_{j}(\tau_{h,l}) \right] \\ & \prod_{l=1}^{H} \left\| \prod_{l=1}^{h} \mathbb{M}_{h',l}^{*} \right\|_{1} \cdot \pi^{t}_{l}(\tau_{h,l}) \right] \cdot \left[ \prod_{l=1}^{n} \mathbb{M}_{h',l}^{*} \right\|_{1} \cdot \pi^{t}_{j}(\tau_{h,l}) \right] \\ & \prod_{l=1}^{H} \left\| \prod_{l=1}^{h} \mathbb{M}_{h',l}^{*} \right\|_{1} \cdot \pi^{t}_{l}(\tau_{h,l}) \right] \cdot \left[ \prod_{l=1}^{n} \mathbb{M}_{h',l}^{*} \right\|_{1} \cdot \pi^{t}_{j}(\tau_{h,l}) \right] \\ & \prod_{l=1}^{H} \left\| \prod_{l=1}^{H} \mathbb{M}_{h',l}^{*} \right\|_{1} \cdot \pi^{t}_{l}(\tau_{h,l}) \right] \\ & \prod_{l=1}^{H} \left\| \prod_{l=1}^{H} \mathbb{M}_{l}^{*} \right\|_{1} \cdot \pi^{t}_{l}(\tau_{h,l}) \right\|_{1} \cdot \pi^{t}_{l}(\tau_{h,l}) \right\|_{1} \cdot \pi^{t}_{l}(\tau_{h,l}) \\ & \prod_{l=1}^{H} \left\| \prod_{l=1}^{H} \mathbb{M}_{l}^{*} \right\|_{1} \cdot \pi^{t}_{l}(\tau_{h,l}) \right\|_{1} \cdot \pi^{t}_{l}(\tau_{h,l}) \right\|_{1} \cdot \pi^{t}_{l}(\tau_{h,l}) \right\|_$$

**Proof for Theorem 4.1** With the Lemmas provided above, we now present the proof for Theorem 4.1, We first restate the theorem as follows:

**Theorem C.1.** We select bonus parameter as  $\beta_m = H^2(S^2A^{|\mathsf{pa}(m)|+1} + SO)\log(TSAOH) +$  $\log(Tn/\delta)$  for some constant c. Then with probability at least  $1 - \delta$ , Algorithm 4.2 guarantees that the following inequality holds true. 

$$Regret(k) = \sum_{t=1}^{k} V^{\pi^{*}} - V^{\pi^{t}} \le \tilde{\mathcal{O}} \Big( \sum_{m=1}^{n} \frac{S^{2}OA^{|\mathsf{pa}(m)|+1}}{\alpha^{2}} \sqrt{k(S^{2}A^{|\mathsf{pa}(m)|+1} + SO)} \Big).$$

where we define  $\pi^*$  as  $\pi^* = \arg \max_{\pi} V^{\pi}$ .

*Proof.* According to Lemma C.5, it is sufficient to obtain an upper bound for the following term: 

$$\sum_{t=1}^{k} \sum_{h=1}^{H-1} \sum_{\tau_h} \pi_m^t(\tau_{h,m}) \left\| \left( \mathbb{M}_{h,m}^t - \mathbb{M}_{h,m}^\star \right) \left[ \prod_{h'=0}^{h} \mathbb{M}_{h',m}^\star \right] \right\|_{1}$$

 $\left\|\prod_{l=1}^{m-1} \left\|\prod_{h'=0}^{n} \mathbb{M}_{h',l}^{\star}\right\|_{1} \cdot \pi_{l}^{t}(\tau_{h,l})\right\| \cdot \left\|\prod_{j=m+1}^{n} \left\|\prod_{h'=0}^{n} \mathbb{M}_{h',j}^{t}\right\|_{1} \cdot \pi_{j}^{t}(\tau_{h,j})\right\|$ 

For probability at least  $1 - \delta$ , we have

$$\sum_{t=1}^{k} \sum_{h=1}^{H-1} \sum_{\tau_{h}} \pi_{m}^{t}(\tau_{h,m}) \left\| \left( \mathbb{M}_{h,m}^{k} - \mathbb{M}_{h,m}^{\star} \right) \left[ \prod_{h'=0}^{h} \mathbb{M}_{h',m}^{\star} \right] \right\|_{1} \cdot \left[ \prod_{l=1}^{m-1} \left\| \prod_{h'=0}^{h} \mathbb{M}_{h',l}^{\star} \right\|_{1} \cdot \pi_{l}^{t}(\tau_{h,l}) \right] \\ \cdot \left[ \prod_{j=m+1}^{n} \left\| \prod_{h'=0}^{h} \mathbb{M}_{h',j}^{t} \right\|_{1} \cdot \pi_{j}^{t}(\tau_{h,j}) \right] \lesssim \sqrt{Sk\beta_{m}}/\alpha$$

For  $m \in [n]$ , we fix  $(o, a, \{a_r\}_{r \in pa(m)}) \in \mathcal{O}_m \times \mathcal{A}_m \times (\times_{r \in pa(m)} \mathcal{A}_r)$ . We define the set of trajectories  $\{\tau_{h,r}\}_{r\in[n]}$ , denoted by  $\mathcal{C}_m$ , as: 

$$\mathcal{C}_m = \left\{ \tau_h \mid \{\tau_{h,r}\}_{r \in [n]} : (o_{h,m}, a_{h,m}, \{a_{h,r}\}_{r \in \mathsf{pa}(m)}) = (o, a, \{a_r\}_{r \in \mathsf{pa}(m)}) \right\}$$

Then, we have the following condition:

$$\sum_{t=1}^{k} \sum_{h=1}^{H-1} \sum_{\tau_{h}} \pi_{m}^{t}(\tau_{h,m}) \left\| \left[ \left( \mathbb{M}_{h,m}^{k} - \mathbb{M}_{h,m}^{\star} \right) \mathbb{O}_{h,m}^{\star} \right] \left( \mathbb{O}_{h,m}^{\star} \right)^{\dagger} \left[ \prod_{h'=0}^{h} \mathbb{M}_{h',m}^{\star} \right] \right\|_{1} \\ \cdot \left[ \prod_{l=1}^{m-1} \left\| \prod_{h'=0}^{h} \mathbb{M}_{h',l}^{\star} \right\|_{1} \cdot \pi_{l}^{t}(\tau_{h,l}) \right] \cdot \left[ \prod_{j=m+1}^{n} \left\| \prod_{h'=0}^{h} \mathbb{M}_{h',j}^{t} \right\|_{1} \cdot \pi_{j}^{t}(\tau_{h,j}) \right] \lesssim \sqrt{Sk\beta_{m}} / \alpha$$

We define  $\{w_{t,l}\}_{(t,l)\in[T]\times[O]}$  that satisfies:

$$w_{t,l} = \left[ (\mathbb{M}_{h,m}^t(o, a, \{a_r\}_{r \in \mathsf{pa}(m)}) - \mathbb{M}_{h,m}(o, a, \{a_r\}_{r \in \mathsf{pa}(m)})) \mathbb{O}_{h,m} \right]_l.$$

We denote the sequence

$$\pi_{m}^{t}(\tau_{h,m})\left(\mathbb{O}_{h,m}^{\star}\right)^{\dagger}\left[\prod_{h'=0}^{h}\mathbb{M}_{h',m}^{\star}\right]\cdot\left[\prod_{l=1}^{m-1}\left\|\prod_{h'=0}^{h}\mathbb{M}_{h',l}^{\star}\right\|_{1}\cdot\pi_{l}^{t}(\tau_{h,l})\right]\cdot\left[\prod_{j=m+1}^{n}\left\|\prod_{h'=0}^{h}\mathbb{M}_{h',j}^{t}\right\|_{1}\cdot\pi_{j}^{t}(\tau_{h,j})\right]$$

for all  $\tau_h$ :  $(o_{h,m}, a_{h,m}, \{a_{h,r}\}_{r \in pa(m)}) = (o, a, \{a_r\}_{r \in pa(m)})$  by  $x_{t,1}, x_{t,2}, \ldots, x_{t,N}$ , where 

$$N = \left| \left\{ \tau_h \mid \tau_h : (o_{h,m}, a_{h,m}, \{a_{h,r}\}_{r \in \mathsf{pa}(m)}) = (o, a, \{a_r\}_{r \in \mathsf{pa}(m)}) \right\} \right|$$

Then we have two observations about the x, w sequence:

The vector sequence  $\{x_{t,i}\}_{i=1}^N$  satisfies  $\sum_{i=1}^N ||x_{t,i}||_1 \le 1$  for all t because

$$\begin{aligned} \pi_{m}^{t}(\tau_{h,m}) \left\| \left( \mathbb{O}_{h,m}^{\star} \right)^{\dagger} \left[ \prod_{h'=0}^{h} \mathbb{M}_{h',m}^{\star} \right] \right\|_{1} \cdot \left[ \prod_{l=1}^{m-1} \left\| \prod_{h'=0}^{h} \mathbb{M}_{h',l}^{\star} \right\|_{1} \cdot \pi_{l}^{t}(\tau_{h,l}) \right] \cdot \left[ \prod_{j=m+1}^{n} \left\| \prod_{h'=0}^{h} \mathbb{M}_{h',j}^{t} \right\|_{1} \cdot \pi_{j}^{t}(\tau_{h,j}) \right] \\ &\leq \sum_{\tau_{h}:(o_{h,m},a_{h,m})=(o,a)} \pi_{m}^{t}(\tau_{h,m}) \left\| \left( \mathbb{O}_{h,m}^{\star} \right)^{\dagger} \left[ \prod_{h'=0}^{h} \mathbb{M}_{h',m}^{\star} \right] \right\|_{1} \left[ \prod_{l=1}^{m-1} \left\| \prod_{h'=0}^{h} \mathbb{M}_{h',l}^{\star} \right\|_{1} \pi_{l}^{t}(\tau_{h,l}) \right] \left[ \prod_{j=m+1}^{n} \left\| \prod_{h'=0}^{h} \mathbb{M}_{h',j}^{t} \right\|_{1} \pi_{j}^{t}(\tau_{h,j}) \right] \\ &\leq \sum_{\{\tau_{h,r}\}_{r\neq m}, \tau_{h-1,m}} \pi_{m}^{t}(\tau_{h,m}) \left\| \left( \mathbb{O}_{h,m}^{\star} \right)^{\dagger} \left[ \prod_{h'=0}^{h} \mathbb{M}_{h',m}^{\star} \right] \right\|_{1} \cdot \left[ \prod_{l=1}^{m-1} \left\| \prod_{h'=0}^{h} \mathbb{M}_{h',l}^{\star} \right\|_{1} \cdot \pi_{l}^{t}(\tau_{h,l}) \right] \\ &\cdot \left[ \prod_{j=m+1}^{n} \left\| \prod_{h'=0}^{h} \mathbb{M}_{h',j}^{t} \right\|_{1} \cdot \pi_{j}^{t}(\tau_{h,j}) \right] = 1. \end{aligned}$$

The vectors 
$$\{w_{t,l}\}_{l=1}^O$$
 satisfy  $\sum_{l=1}^O ||w_{t,l}||_1 \le 2S^{1.5}/\alpha$  for all t, since we have

$$\sum_{l=1}^{O} \|w_{t,l}\|_{1} = \left\| \left( \mathbb{M}_{h,m}^{t} - \mathbb{M}_{h,m}^{\star} \right) \mathbb{O}_{h,m}^{\star} \right\|_{1} \le S \left( \left\| \mathbb{M}_{h,m}^{t} \right\|_{1} + \left\| \mathbb{M}_{h,m}^{\star} \right\|_{1} \right) \le 2S^{1.5} / \alpha.$$

Using the notation of  $\{x_{t,i}\}_{i=1}^n$  and  $\{w_{t,l}\}_{l=1}^O$ , we have

$$\sum_{t=1}^{k-1} \sum_{l=1}^{O} \sum_{i=1}^{n} |w_{k,l}^T x_{t,i}| = \mathcal{O}\left(\frac{\sqrt{S}}{\alpha} \sqrt{k\beta_m}\right).$$

Therefore, we can bind the target term with Eluder-Dimension lemma (Proposition 22 of Liu et al. (2022a)). We have the following result.

$$\sum_{t=1}^{k} \sum_{l=1}^{O} \sum_{i=1}^{n} |w_{t,l}^{T} x_{t,i}| = \mathcal{O}\left(\frac{S^{1.5} H^{2}}{\alpha} \sqrt{k\beta_{m}}\right).$$

1258 The equation holds for all  $k \in [T]$ . We represent the result with matrix operator, and we arrive at

$$\begin{array}{ll} 1259 \\ 1260 \\ 1261 \\ 1262 \\ 1262 \\ 1263 \\ 1264 \\ 1265 \end{array} & \sum_{t=1}^{k} \sum_{h=1}^{H-1} \sum_{\tau_h \in \mathcal{C}_m} \pi_m^t(\tau_{h,m}) \left\| \left( \mathbb{M}_{h,m}^k - \mathbb{M}_{h,m}^\star \right) \left[ \prod_{h'=0}^{h} \mathbb{M}_{h',m}^\star \right] \right\|_1 \cdot \left[ \prod_{l=1}^{m-1} \left\| \prod_{h'=0}^{h} \mathbb{M}_{h',l}^\star \right\|_1 \cdot \pi_l^t(\tau_{h,l}) \right] \\ \cdot \left[ \prod_{j=m+1}^{n} \left\| \prod_{h'=0}^{h} \mathbb{M}_{h',j}^t \right\|_1 \cdot \pi_j^t(\tau_{h,j}) \right] \lesssim \sqrt{S^3 H^4 k \beta_m} / \alpha \end{array}$$

We sum up both the left-hand side and the right-hand side for all  $(o, a, \{a_r\}_{r \in pa(m)})$ , and we can obtain that

$$\sum_{t=1}^{k} \sum_{h=1}^{H-1} \sum_{\boldsymbol{\tau}_{h}} \pi_{m}^{t}(\boldsymbol{\tau}_{h,m}) \left\| \left( \mathbb{M}_{h,m}^{k} - \mathbb{M}_{h,m}^{\star} \right) \left[ \prod_{h'=0}^{h} \mathbb{M}_{h',m}^{\star} \right] \right\|_{1} \cdot \left[ \prod_{l=1}^{m-1} \left\| \prod_{h'=0}^{h} \mathbb{M}_{h',l}^{\star} \right\|_{1} \cdot \pi_{l}^{t}(\boldsymbol{\tau}_{h,l}) \right] \\ \cdot \left[ \prod_{j=m+1}^{n} \left\| \prod_{h'=0}^{h} \mathbb{M}_{h',j}^{t} \right\|_{1} \cdot \pi_{j}^{t}(\boldsymbol{\tau}_{h,j}) \right] \lesssim \frac{S^{1.5} H^{2} O A^{\mathsf{pa}(m)+1}}{\alpha} \sqrt{k\beta_{m}}.$$

1276 1277

1278

1281

1285 1286 1287

1266

1267

1244 1245 1246

1249 1250 1251

1255 1256 1257

### C.2 PROOF FOR THEOREM 4.2

**Theorem C.2.** For both randomized and deterministic algorithms, there exists an instance of DEC-POMDP with factorization such that the regret is at least  $\mathcal{O}(\sqrt{A^{ind(G)+1}T})$ .

1282 *Proof.* We consider the scenario where the state is directly observable to the agents, namely, a DEC-1283 MDP under a factored structure model. We further assume that the transition probability satisfies 1284 the following structure: for all  $\mathbf{s}', \mathbf{s} \in S$ ,  $\mathbf{a} \in A$ , and  $h \in [H]$ ,

$$\mathbb{T}_{h,\mathbf{a}}(\mathbf{s}' \mid \mathbf{s}) = \left[\prod_{m=1}^{n-1} \mathbb{T}_{h,m}(s'_m \mid s_m, a_m)\right] \mathbb{T}_{h,n}(s'_n \mid s_n, \{a_r\}_{r \in [n]}).$$

We further assume that the episode length H = 2, and the reward function satisfies  $r_{h,m}(s_m) = 0$  for all  $h \in [H]$  and  $m \in [n-1]$ . Thus, the entire model is equivalent to an  $A^n$ -armed bandit problem. By leveraging a classic result on the lower bound of regret for the multi-armed bandit problem (Mannor and Tsitsiklis, 2004), it follows that for any randomized or deterministic algorithm, there exists an instance of the multi-armed bandit problem such that the regret is at least  $\mathcal{O}(\sqrt{AT})$ , where  $\tilde{A}$  denotes the number of arms. Consequently, for any randomized or deterministic algorithm, there exists an instance of a factored DEC-POMDP such that the regret for achieving the global optimum is at least  $\mathcal{O}(\sqrt{A^nT}) = \mathcal{O}(\sqrt{A^{ind(G)+1}T})$ .

# <sup>1296</sup> D SUPPLEMENTARY DETAILS FOR SECTION 4.3

# D.1 COMPLETE ALGORITHM FOR ACHIEVING LOCAL OPTIMAL UNDER FACTORED STRUCTURE MODEL

We present Algorithm D.1 as the complete algorithm for achieving local optimal under factored structure model.

1304

1349

1298

1299

1300

1305 Algorithm 4 NASH-CA for Achieving Local Optimal Under Factored Structure Model 1: Initialize  $\pi = {\pi_i}_{i \in [n]}$ , where  $\pi_i = {\pi_{h,i}}_{(h,i) \in [H] \times [m]}$ . 1307 2: while true do Execute policy  $\pi$  for  $N = \frac{CH^2}{\epsilon^2} \log \left( (nHSK \max_{i \in [n]} A_i) / (\epsilon \delta) \right)$  episodes and obtain 3: 1309  $\hat{V}_{1,i}(\pi)$  which is the empirical average of the total return under policy  $\pi$ . 1310 4: for agent  $i = 1, \ldots, m$  do 1311 Appoint agent i as the central agent and fix  $\pi_{-i}$  to run Algorithm ?? for  $K_i =$ 5: 1312  $\tilde{\mathcal{O}}(S^4 A^{2 \cdot \mathsf{ind}(G[\mathsf{ch}(i) \cup \{i\}])}(S^2 A^{\mathsf{ind}(G)} + SO) \cdot \mathsf{poly}(H)/(\alpha^4 \epsilon^2))$  episodes and get a new policy 1313  $\hat{\pi}_i$ . Execute policy  $(\hat{\pi}_i, \pi_{-i})$  for  $N = \frac{CH^2}{\epsilon^2} \log \left( (nHSK \max_{i \in [n]} A_i) / (\epsilon \delta) \right)$  episodes and 6: 1315 obtain  $\hat{V}(\hat{\pi}_i, \pi_{-i})$  which is the empirical average of the total return under policy  $(\hat{\pi}_i, \pi_{-i})$ . 1316 Set  $\Delta_i \leftarrow \hat{V}(\hat{\pi}_i, \pi_{-i}) - \hat{V}(\pi)$ . 7: 1317 if  $\max_{i \in [n]} \Delta_i > \epsilon/2$  then 8: 1318 Update  $\pi_j \leftarrow \hat{\pi}_j$  where  $j = \arg \max_{i \in [n]} \Delta_i$ . 9: 1319 10: else 1320 11: return  $\pi$ 1321 1322 1323 1324 Algorithm 5 OMLE for Achieving Local Optimal Under Factored Model 1326 1: Initialize:  $\mathcal{B}_m^1 = \{\hat{\theta}_m \in \Theta_m : \min_h \sigma_S(\mathbb{O}_m(\hat{\theta}_m)) \ge \alpha\}, \mathcal{D}_m = \{\}$  for all  $m \in \overline{\mathsf{ch}}(i),$ 1327  $\tilde{\mathcal{B}}^1 = \{\theta_{m \in \mathsf{nch}(i)} : \min_h \sigma_S(\mathbb{O}_m(\theta_m)) \ge \alpha, \forall m \notin \overline{\mathsf{ch}}(i)\}, \tilde{\mathcal{D}} = \{\}, \text{ central agent } i, \text{ policy of } i\}$ 1328 other agent  $\pi_{-i}$ . 2: for k = 1 ... T do 1330 Follow  $\pi_{\mathsf{nch}(i)}$  to collect trajectories  $\tau_{\mathsf{ch}(i)}^k = \{o_{1,m}^k, \dots, a_{H,m}^k\}_{m \in \mathsf{nch}(i)}$ . 3: 1331 Add  $\tau^k_{\mathsf{nch}(i)}$  into  $\tilde{\mathcal{D}}$  and update confidence interval with eq. (4). 4: 1332 5: for k = 1 ... T do 1333 compute  $(\boldsymbol{\theta}^k, \pi_i^k) = \arg \max_{\{\hat{\boldsymbol{\theta}}_m \in \mathcal{B}_m^k\}_m \in \mathsf{Chl}(i), \tilde{\boldsymbol{\theta}}_i \in \tilde{\boldsymbol{\Theta}}_i, \mu_i} V^{\mu_i, \pi_{-i}}(\hat{\boldsymbol{\theta}})$ 6: 1334 7: for m = 1, 2, ..., r do 1335 for h = 1, ..., H do 8: 1336 Agent  $l \in \operatorname{nch}(i)$  take action  $a_{h,l}^T$ . 9: 1337 Select an action  $a_{h,l_j}^k \sim \pi_{h,l_j}(\cdot \mid \tau_{h-1,l_j}, o_{h,l_j})$  for all  $j \in [r-1]$ . 1338 10: 1339 Select an action  $a_{h,i}^k \sim \pi_{h,i}^k (\cdot \mid \tau_{h-1,i}, o_{h,i}).$ 11: 1340 For agent  $l_j$  with  $j \in [m]$ , collect observation  $o_{h+1,l_j}^k$  from the environment. 12: 1341 For  $j \in [r] \setminus [m]$ , sample dummy state  $s_{h+1,l_j} \sim \mathbb{T}_{h,l_j}^{k} (\cdot \mid s_{h,l_j}, a_{h,\overline{pa}(l_j)})$ . 13: 1342 Collect observation  $o_{h+1,l_j}^k \sim \mathbb{O}_{h+1,l_j}^k (\cdot \mid s_{h+1,l_j})$  for  $j \in [r] \setminus [m]$ . 14: 1343 If  $m \neq r$ , add  $(\pi_m, \tau^k_{\overline{\mathsf{pa}}(m) \cap \overline{\mathsf{ch}}(i)}, \tau^T_{\overline{\mathsf{pa}}(m) \setminus \mathsf{ch}(i)})$  to  $\mathcal{D}_m$  for  $m \neq r$ . 15: 1345 Otherwise, add  $(\pi_i^k, \tau_{\overline{pa}(i)\cap \overline{ch}(i)}^k, \tau_{\overline{pa}(i)\setminus ch(i)}^T)$  to  $\mathcal{D}_i$ . 16: Update confidence interval with eq. (5) for all  $m \in \overline{pa}(i)$ . 17: 1347 18: **Output**  $\hat{\pi}$  as uniform mixture of the policies  $\pi_i^1, \pi_i^2, \ldots, \pi_i^K$ . 1348

25

**Bonums Term** : For  $m \in [n]$ , the bonus term  $\beta_m$  and  $\tilde{\beta}$  is defined as 

$$\beta_m = c(H^2(S^2A^{|\mathsf{pa}(m)|+1} + SO)\log(TSAOH) + \log(Tn/\delta)),$$

$$\tilde{\beta} = c\left(\left(\sum_{\substack{r\notin c\bar{h}(n)}} H(S^2A^{|\mathsf{pa}(r)|+1} + SO)\log(TSAOH)\right) + \log(Tn/\delta)\right).$$
(15)

#### D.2 PROOF FOR THEOREM 4.3

We first proof the following Theorem D.1, and then we will prove Theorem 4.3. **Theorem D.1.** If Algorithm 5.1 take agent i as central agent and take policies  $\{\pi_m\}_{m \in [n]/\{i\}}$  as input, then Algorithm 5.1 guarantees that with probability at least  $1 - \delta$  for all  $k \in [T]$ 

$$\left[\sum_{t=1}^{k} V^{\pi_i^*, \boldsymbol{\pi}_{-i}} - V^{\pi_i^k, \boldsymbol{\pi}_{-i}}\right] \leq \tilde{\mathcal{O}}\left(\frac{S^2 O A^{d_i+1}}{\alpha^2} \sqrt{K\left(S^2 A^{d+1} + SO\right)} \cdot \operatorname{poly}(H)\right),$$

where  $\pi_i^*$  is defined as  $\pi_i^* = \arg \max_{\pi_i} V^{\pi_i, \pi_{-i}}$ , and recall that we use  $d_i$  to denote the maximum indegree of the subgraph induced by ch(i)

#### D.2.1 PROOF FOR THEOREM D.1

Without loss of generosity, we consider the case where the central agent is agent n. 

**Lemma D.1.** There exists an absolute constant c such that for any  $\delta \in (0, 1]$ , with probability at least  $1 - \delta$ : the following inequality holds true. 

$$\max_{(\theta_m,t)\in\Theta_m\times[T]}\sum_{i=1}^t \log\left(\frac{\mathbb{P}_{\theta_m}^{\pi_m}\left(\tau_m^i \mid \{\tau_r^i\}_{r\in\mathsf{pa}(m)\cap\overline{\mathsf{ch}}(n)}, \{\tau_r^T\}_{r\in\mathsf{pa}(m)}^{r\notin\overline{\mathsf{ch}}(n)}\right)}{\mathbb{P}_{\theta_m^*}^{\pi_m}\left(\tau_m^i \mid \{\tau_r^i\}_{r\in\mathsf{pa}(m)\cap\overline{\mathsf{ch}}(n)}, \{\tau_r^T\}_{r\in\mathsf{pa}(m)}^{r\notin\overline{\mathsf{ch}}(n)}\right)}\right) < \beta_m, m \in \overline{\mathsf{ch}}(n),$$

$$\max_{(\theta_n,t)\in\Theta_n\times[T]}\sum_{i=1}^t \log\left(\frac{\mathbb{P}_{\theta_n}^{\pi_n^i}\left(\tau_n^i \mid \{\tau_r^i\}_{r\in\mathsf{pa}(n)\cap\overline{\mathsf{ch}}(n)}, \{\tau_r^T\}_{r\in\mathsf{pa}(n)}^{r\notin\overline{\mathsf{ch}}(n)}\right)}{\mathbb{P}_{\theta_n^*}^{\pi_n^i}\left(\tau_n^i \mid \{\tau_r^i\}_{r\in\mathsf{pa}(n)\cap\overline{\mathsf{ch}}(n)}, \{\tau_r^T\}_{r\in\mathsf{pa}(n)}^{r\notin\overline{\mathsf{ch}}(n)}\right)}\right) < \beta_n,$$

$$\max_{\substack{(\times_{j\notin\overline{\mathsf{ch}}(n)}\theta_j,t)\in\times_{j\notin\overline{\mathsf{ch}}(n)}\Theta_j\times[T]}}\sum_{i=1}^t \log\left(\frac{\prod_{j\notin\overline{\mathsf{ch}}(n)}\mathbb{P}^{\pi_j}_{\theta_j}(\tau_j^i\mid\{\tau_r^i\}_{r\in\mathsf{pa}(n)})}{\prod_{j\notin\overline{\mathsf{ch}}(n)}\mathbb{P}^{\pi_j}_{\theta_j^*}(\tau_j^i\mid\{\tau_r^i\}_{r\in\mathsf{pa}(n)})}\right)<\tilde{\beta}.$$

*Proof.* The proof of the lemma is similar to the proof for Lemma D.1, so we omit it here for clarity. 

**Lemma D.2.** There exists a universal constant c such that for any  $\delta \in (0, 1]$ , with probability at least for all  $t \in [T]$  and all  $\theta_m \in \Theta_m$ ,  $m \in |ch(n)| - 1$ , the following inequalities hold true. Initially, for agent  $m \in [|ch(n)|]$ , we have 

$$\sum_{i=1}^{t} \left( \sum_{\{\tau_r\}_{r\in\overline{\mathbf{h}}(n)}} \left| \mathbb{P}_{\theta_{c_n,m}}^{\pi_{c_n,m}} - \mathbb{P}_{\theta_{c_n,m}^*}^{\pi_{c_n,m}} \right| \left[ \prod_{l=1}^{m-1} \mathbb{P}_{\theta_{c_n,l}^*}^{\pi_{c_n,l}} \right] \left[ \prod_{j=m+1}^{|\overline{\mathbf{h}}(n)|-1} \mathbb{P}_{\theta_{c_n,j}}^{\pi_{c_n,j}} \right] \mathbb{P}_{\theta_n^*}^{\pi_n^*} \right)^2$$

$$\left( t \left( \mathbb{P}_{\theta^*}^{\pi_{c_n,m}} \left( \tau_c^i - |\{\tau_r^i\}_{r\in\mathbf{p}(c_n)} \cup \overline{\mathbf{o}}_{\overline{\mathbf{h}}(n)} \right) \{\tau_r^T\}_{r\in\mathbf{p}(c_n)}^{r\notin\overline{\mathbf{c}}_{\overline{\mathbf{h}}(n)}} \right) \right) \right)$$

$$\leq c \left( \sum_{i=1}^{t} \log \left( \frac{\mathbb{P}_{\theta_{c_{n,m}}^{\tau_{c_{n,m}}}}\left(\tau_{c_{n,m}}^{i} \mid \{\tau_{r}^{i}\}_{r \in \mathsf{pa}(c_{n,m}) \cap \overline{\mathsf{ch}}(n)}, \{\tau_{r}^{T}\}_{r \in \mathsf{pa}(c_{n,m})}^{r \notin \overline{\mathsf{ch}}(n)}\right)}{\mathbb{P}_{\theta_{c_{n,m}}}^{\pi_{c_{n,m}}}\left(\tau_{c_{n,m}}^{i} \mid \{\tau_{r}^{i}\}_{r \in \mathsf{pa}(c_{n,m}) \cap \overline{\mathsf{ch}}(n)}, \{\tau_{r}^{T}\}_{r \in \mathsf{pa}(c_{n,m})}^{r \notin \overline{\mathsf{ch}}(n)}\right)} \right) + \beta_{c_{n,m}}} \right).$$

For central agent n, we have

1397  
1398  
1399  
1400  

$$\sum_{i=1}^{t} \left( \sum_{\{\tau_r\}_{r\in\overline{ch}(n)}} \left| \mathbb{P}_{\theta_n}^{\pi_n^i} - \mathbb{P}_{\theta_n^*}^{\pi_n^i} \right| \left[ \prod_{l\in\mathsf{ch}(n)} \mathbb{P}_{\theta_l^*}^{\pi_l} \right] \right)^2$$

$$\leq c \left( \sum_{i=1}^{t} \log \left( \frac{\mathbb{P}_{\theta_{n}^{*}}^{\pi_{n}^{i}} \left(\tau_{n}^{i} \mid \{\tau_{r}^{i}\}_{r \in \mathsf{pa}(n) \cap \overline{\mathsf{ch}}(n)}, \{\tau_{r}^{T}\}_{r \in \mathsf{pa}(n)}^{r \notin \overline{\mathsf{ch}}(n)} \right)}{\mathbb{P}_{\theta_{n}^{*}}^{\pi_{n}^{i}} \left(\tau_{n}^{i} \mid \{\tau_{r}^{i}\}_{r \in \mathsf{pa}(n) \cap \overline{\mathsf{ch}}(n)}, \{\tau_{r}^{T}\}_{r \in \mathsf{pa}(n)}^{r \notin \overline{\mathsf{ch}}(n)} \right)} \right) + \beta_{n} \right)$$

For estimation error of  $\{\theta_l\}_{l\notin \overline{ch}(n)}$ , we have 

$$\leq c \left( \sum_{i=1}^{t} \log \left( \frac{\prod_{l \notin \overline{\mathsf{ch}}(n)} \mathbb{P}_{\theta_l^i}^{\pi_l} \left( \tau_l^i \mid \{\tau_r^i\}_{r \in \mathsf{pa}(l)} \right)}{\prod_{l \notin \overline{\mathsf{ch}}(n)} \mathbb{P}_{\theta_l}^{\pi_l} \left( \tau_l^i \mid \{\tau_r^i\}_{r \in \mathsf{pa}(l)} \right)} \right) + \tilde{\beta} \right)$$

where we for all  $m \in |\overline{ch}(n)| - 1$ ,  $\theta_m \in \Theta_m$ , and policy  $\pi_m$ , we denote  $\mathbb{P}_{\theta_m}^{\pi_m}$  as  $\mathbb{P}_{\theta_m}^{\pi_m} =$  $\mathbb{P}_{\theta_m}^{\pi_m}\left(\tau_m \mid \{\tau_r\}_{r \in \mathsf{pa}(m) \cap \overline{\mathsf{ch}}(n)}, \{\tau_r^T\}_{r \in \mathsf{pa}(m)}^{r \notin \overline{\mathsf{ch}}(n)}\right).$ 

 $\sum_{i=1}^{t} \left( \sum_{\{\tau_r\}_{r \in \overline{\mathsf{ch}}(r)}} \left| \prod_{l \notin \overline{\mathsf{ch}}(r)} \mathbb{P}_{\theta_l}^{\pi_l} \left( \tau_l \mid \{\tau_r\}_{r \in \mathsf{pa}(l)} \right) - \prod_{l \notin \overline{\mathsf{ch}}(r)} \mathbb{P}_{\theta_l}^{\pi_l} \left( \tau_l \mid \{\tau_r\}_{r \in \mathsf{pa}(l)} \right) \right| \right)$ 

*Proof.* The proof of the lemma is similar to the proof for Lemma C.2, so we omit it here for clarity. 

**Definition D.1.** For all  $(m, h, k) \in [n] \times [H] \times [T]$ , we define the matrix notations  $\mathbb{M}_{h,m}^{\star} \in \mathbb{R}^{O \times O}$ and  $\mathbb{M}_{h,m}^k \in \mathbb{R}^{O \times O}$  as follows: 

 $\mathbb{M}_{0,m}^{\star} = \mathbb{O}_{1,m}^{\star} \mu_m^{\star} \in \mathbb{R}^O, \quad \mathbb{M}_{0,m}^k = \mathbb{O}_{1,m}^k \mu_m^k \in \mathbb{R}^O,$  $\mathbb{M}_{h,m}^{\star}(o_m, a_m, \{a_r\}_{r \in \mathsf{pa}(m)}) = \mathbb{O}_{h+1,m}^{\star} \mathbb{T}_{h,m,a_m,\{a_r\}_{r \in \mathsf{pa}(m)}}^{\star} \cdot \operatorname{diag}\left(\mathbb{O}_{h,m}^{\star}(o_m \mid \cdot)\right) \left(\mathbb{O}_{h,m}^{\star}\right)^{\dagger} \in \mathbb{R}^{O \times O},$  $\mathbb{M}_{h,m}^{k}(o_{m}, a_{m}, \{a_{r}\}_{r \in \mathsf{pa}(m)}) = \mathbb{O}_{h+1,m}^{k} \mathbb{T}_{h,m,a_{m},\{a_{r}\}_{r \in \mathsf{pa}(m)}}^{k} \cdot diag\left(\mathbb{O}_{h,m}^{k}(o_{m} \mid \cdot)\right)\left(\mathbb{O}_{h,m}^{k}\right)^{\dagger} \in \mathbb{R}^{O \times O},$ where  $\{\mathbb{O}_{h,m}^{\star}\}_{(h,m)\in[H]\times[n]}$  and  $\{\mathbb{T}_{h,m,a_m,\{a_r\}_{r\in pa(m)}}^{\star}\}_{(h,m)\in[H]\times[n]}$  denote the observation and transition matrices corresponding to the true transition model, and  $\{\mathbb{O}_{h,m}^k\}_{(h,m)\in[H]\times[n]}$  and  $\{\mathbb{T}_{h,m,a_m,\{a_r\}_{r\in pa(m)}}^k\}_{(h,m)\in[H]\times[n]}$  denote the observation and transition matrices corresponding to model parameter  $\theta_k$  for all  $k \in [T]$ . When no confusion arises, we simplify the notation by us-ing  $\mathbb{M}_{h,m}^{\star}(\{a_{h,r}\}_{r\in\mathsf{pa}(m)}^{r\notin\overline{\mathsf{ch}}(n)})$  to represent  $\mathbb{M}_{h,m}^{\star}(o_m, a_m, \{a_r\}_{r\in\mathsf{pa}(m)})$  and  $\mathbb{M}_{h,m}^k(\{a_{h,r}\}_{r\in\mathsf{pa}(m)}^{r\notin\overline{\mathsf{ch}}(n)})$  to represent  $\mathbb{M}_{h,m}^k(o_m, a_m, \{a_r\}_{r \in \mathsf{pa}(m)})$ . 

According to the Definition D.1, we can directly achieve the following result: 

**Lemma D.3.** With probability at least  $1 - \delta$ , for all  $(k, h, m) \in [K] \times [H - 1] \times |ch(m)| - 1$ , the following probability holds true. 

$$\begin{aligned} & \underset{1439}{1438} & \sum_{t=1}^{k-1} \sum_{\{\tau_r\}_{r \in \overline{ch}(n)}} \left\| \left( \mathbb{M}_{h,c_{n,m}}^k \left( \{a_{h,r}^T\}_{r \in \mathsf{pa}(c_{n,m})}^{r \notin \overline{ch}(n)} \right) - \mathbb{M}_{h,c_{n,m}}^* \left( \{a_{h,r}^T\}_{r \in \mathsf{pa}(c_{n,m})}^{r \notin \overline{ch}(n)} \right) \right) \left[ \prod_{h'=0}^{h-1} \mathbb{M}_{h',c_{n,m}}^* \left( \{a_{h',r}^T\}_{r \in \mathsf{pa}(c_{n,m})}^{r \notin \overline{ch}(n)} \right) \right) \right\|_{1} \\ & \underset{1441}{1442} & \cdot \pi_{c_{n,m}}(\tau_{h,c_{n,m}}) \left[ \prod_{j=m+1}^{n-1} \left\| \mathbf{m}_{h,c_{n,j}}^t \|_{1} \pi_{c_{n,j}}(\tau_{h,c_{n,j}}) \right] \left[ \prod_{j=1}^{m-1} \left\| \mathbf{m}_{h,c_{n,j}}^t \|_{1} \pi_{c_{n,j}}(\tau_{h,c_{n,j}}) \right] \left\| \mathbf{m}_{h,n} \|_{1} \pi_{n}^t(\tau_{h,n}) \lesssim \frac{\sqrt{Sk\beta_{c_{n,m}}}}{\alpha} \\ & \frac{1444}{1444} & \sum_{t=1}^{k-1} \sum_{\{\tau_r\}_{r \in \overline{ch}(n)}} \pi_n^t(\tau_{h,n}) \left\| \left( \mathbb{M}_{h,n}^k \left( \{a_{h,r}^T\}_{r \in \overline{cp}(n)}^{r \notin \overline{ch}(n)} \right) - \mathbb{M}_{h,n}^* \left( \{a_{h,r}^T\}_{r \in \overline{cp}(n)}^{r \notin \overline{ch}(n)} \right) \right) \mathbf{m}_{h-1,n} \right\|_{1} \\ & \frac{1444}{1447} & \sum_{t=1}^{k-1} \sum_{\{\tau_r\}_{r \in \overline{ch}(n)}} \pi_n^t(\tau_{h,n}) \left\| \left( \mathbb{M}_{h,n}^k \left( \{a_{h,r}^T\}_{r \in \overline{cp}(n)}^{r \notin \overline{ch}(n)} \right) - \mathbb{M}_{h,n}^* \left( \{a_{h,r}^T\}_{r \in \overline{cp}(n)}^{r \notin \overline{ch}(n)} \right) \right) \mathbf{m}_{h-1,n} \right\|_{1} \\ & \frac{1448}{1449} & \cdot \left[ \prod_{l=1}^{n-1} \left\| \mathbf{m}_{h,c_{n,l}} \right\|_{1} \pi_{c_{n,l}}(\tau_{h,c_{n,l}}) \right] \lesssim \sqrt{Sk\beta_n} / \alpha \\ & \frac{1451}{1452} & \sum_{t=1}^{k-1} \sum_{\{\tau_r\}_{r \notin \overline{ch}(n)}} \left\| \prod_{l \notin \overline{ch}(n)} \mathbb{P}_{\theta_l^k}^{\pi_l} \left( \tau_l \mid \{\tau_r\}_{r \in pa(l)} \right) - \prod_{l \notin \overline{ch}(n)} \mathbb{P}_{\theta_l^k}^{\pi_l} \left( \tau_l \mid \{\tau_r\}_{r \in pa(l)} \right) \right\| = \mathcal{O} \left( \sqrt{k\tilde{\beta}} \right). \\ & \text{where for all } m \in \overline{ch}(n), h \in [H], t \in [T], we define \mathbf{m}_{h,m} \text{ and } \mathbf{m}_{h,m}^t \text{ as} \\ & \frac{1}{1456} & \frac{1}{2} \sum_{t=1}^{k-1} \sum_{t=1}^{k-1}$$

$$\mathbf{m}_{h,m} = \prod_{h'=0}^{h} \mathbb{M}_{h,m}^{\star} \left( \{a_{h,r}^{T}\}_{r \in \mathsf{pa}(m)}^{r \notin \overline{\mathsf{ch}}(n)} \right), \quad \mathbf{m}_{h,m}^{t} = \prod_{h'=0}^{h} \mathbb{M}_{h,m}^{t} \left( \{a_{h,r}^{T}\}_{r \in \mathsf{pa}(m)}^{r \notin \overline{\mathsf{ch}}(n)} \right).$$

According to the definition of regret, we can bound the regret with trajectory probability in the following way: 

**Lemma D.4.** *The regret is bounded by the following inequality:* 

$$Regret(k) \le nH \sum_{t=1}^{k} \sum_{\boldsymbol{\tau}_{H}} \left| \mathbb{P}_{\boldsymbol{\theta}^{t}}^{\boldsymbol{\pi}_{-n}, \boldsymbol{\pi}_{n}^{t}}(\boldsymbol{\tau}_{H}) - \mathbb{P}_{\boldsymbol{\theta}^{*}}^{\boldsymbol{\pi}_{-n}, \boldsymbol{\pi}_{n}^{t}}(\boldsymbol{\tau}_{H}) \right|,$$

where we define

$$\mathbb{P}_{\theta^{t}}^{\pi_{-n},\pi_{n}^{t}}(\boldsymbol{\tau}_{H}) = \left[\prod_{m=1}^{n-1} \mathbb{P}_{\theta_{m}^{t}}^{\pi_{m}}(\tau_{H,m} \mid \{\tau_{H,r}\}_{r \in \mathsf{pa}(m)})\right] \mathbb{P}_{\theta_{n}^{t}}^{\pi_{n}^{t}}(\tau_{H,n} \mid \{\tau_{H,r}\}_{r \in \mathsf{pa}(n)}),$$

$$\mathbb{P}_{\boldsymbol{\theta}^*}^{\boldsymbol{\pi}_{-n}, \boldsymbol{\pi}_n^t}(\boldsymbol{\tau}_H) = \left[\prod_{m=1}^{n-1} \mathbb{P}_{\boldsymbol{\theta}_m^*}^{\boldsymbol{\pi}_m}(\boldsymbol{\tau}_{H,m} \mid \{\boldsymbol{\tau}_{H,r}\}_{r \in \mathsf{pa}(m)})\right] \mathbb{P}_{\boldsymbol{\theta}_n^*}^{\boldsymbol{\pi}_n^t}\left(\boldsymbol{\tau}_{H,n} \mid \{\boldsymbol{\tau}_{H,r}\}_{r \in \mathsf{pa}(n)}\right).$$

Lemma D.5. The regret is bounded by the following inequality: 

 $\sum_{t=1}^{2} \sum_{\{\tau_{H,r}\}_{r\notin \overline{\mathsf{ch}}(n)}} \left| \underset{l\notin \overline{\mathsf{ch}}(n)}{\overset{\mathbf{II}}{\mathsf{ch}}} \right|_{0}$  $l\notin \overline{ch}(n)$ where for any  $m \in \overline{ch}(n)$ ,  $\theta_m \in \Theta_m$ , any policy  $\pi_m$ , we denote  $\mathbb{P}_{\theta_m}^{\pi_m}$  as  $\mathbb{P}_{\theta_m}^{\pi_m}$ =

1481 
$$\mathbb{P}_{\theta_m}^{\pi_m} \left( \tau_{H,m} \mid \{\tau_{H,r}\}_{r \in \mathsf{pa}(m)} \right).$$

Proof. According to the factorization of trajectory probability, we bound can  $\sum_{\boldsymbol{\tau}_{H}} \left| \mathbb{P}_{\boldsymbol{\theta}^{t}}^{\boldsymbol{\pi}_{-n}, \boldsymbol{\pi}_{n}^{t}}(\boldsymbol{\tau}_{H}) - \mathbb{P}_{\boldsymbol{\theta}^{*}}^{\boldsymbol{\pi}_{-n}, \boldsymbol{\pi}_{n}^{t}}(\boldsymbol{\tau}_{H}) \right| \text{ with the following inequalities.}$ 

$$\begin{aligned}
& \sum_{\tau_{H}} \left| \mathbb{P}_{\theta^{t}}^{\pi_{-n},\pi_{n}^{t}}(\tau_{H}) - \mathbb{P}_{\theta^{*}}^{\pi_{-n},\pi_{n}^{t}}(\tau_{H}) \right| \\
& = \sum_{\tau_{H}} \left| \left[ \prod_{l \in ch(n)} \mathbb{P}_{\theta^{l}_{l},l}^{\pi_{l}} \right] \mathbb{P}_{\theta^{t}_{n},n}^{\pi_{n}} - \left[ \prod_{l \in ch(n)} \mathbb{P}_{\theta^{l}_{n},l}^{\pi_{l}} \right] \mathbb{P}_{\theta^{t}_{n},n}^{\pi_{n}} \right| \left[ \prod_{l \notin ch(n)} \mathbb{P}_{\theta^{l}_{l},l}^{\pi_{l}} \right] \\
& + \sum_{t=1}^{k} \sum_{\tau_{H}} \left| \prod_{l \notin ch(n)} \mathbb{P}_{\theta^{l}_{l},l}^{\pi_{l}} \left( \tau_{H,l} \mid \{\tau_{H,r}\}_{r \in pa(l)} \right) - \prod_{l \notin ch(n)} \mathbb{P}_{\theta^{l}_{l},l}^{\pi_{l}} \left( \tau_{H,l} \mid \{\tau_{H,r}\}_{r \in pa(l)} \right) \right| \left[ \prod_{l \in ch(n)} \mathbb{P}_{\theta^{l}_{l},l}^{\pi_{n}} \right] \\
& = \sum_{t=1}^{k} \sum_{\tau_{H}} \left| \prod_{l \notin ch(n)} \mathbb{P}_{\theta^{l}_{l},l}^{\pi_{l}} \right| \mathbb{P}_{\theta^{t}_{n},n}^{\pi_{n}} - \left[ \prod_{l \in ch(n)} \mathbb{P}_{\theta^{l}_{l},l}^{\pi_{l}} \right] \mathbb{P}_{\theta^{t}_{n},n}^{\pi_{n}} \right| \left[ \prod_{l \notin ch(n)} \mathbb{P}_{\theta^{l}_{l},l}^{\pi_{l}} \right] \\
& = \sum_{t=1}^{k} \sum_{\tau_{H}} \left| \left[ \prod_{l \in ch(n)} \mathbb{P}_{\theta^{l}_{l},l}^{\pi_{l}} \right] \mathbb{P}_{\theta^{t}_{n},n}^{\pi_{l}} - \left[ \prod_{l \in ch(n)} \mathbb{P}_{\theta^{l}_{l},l}^{\pi_{l}} \right] \mathbb{P}_{\theta^{t}_{n},n}^{\pi_{n}} \right| \left[ \prod_{l \notin ch(n)} \mathbb{P}_{\theta^{l}_{l},l}^{\pi_{l}} \right] \\
& = \sum_{t=1}^{k} \sum_{\tau_{H}} \left| \prod_{l \in ch(n)} \mathbb{P}_{\theta^{l}_{l},l}^{\pi_{l}} \right| \mathbb{P}_{\theta^{t}_{l},l}^{\pi_{l}} \left( \tau_{H,l} \mid \{\tau_{H,r}\}_{r \in pa(l)} \right) - \prod_{l \notin ch(n)} \mathbb{P}_{\theta^{l}_{l},l}^{\pi_{l}} \left( \tau_{H,l} \mid \{\tau_{H,r}\}_{r \in pa(l)} \right) \right|. \\
& = \sum_{t=1}^{k} \sum_{\tau_{H,r},\tau_{r \notin ch(n)}} \left| \prod_{l \notin ch(n)} \mathbb{P}_{\theta^{l}_{l},l}^{\pi_{l}} \left( \tau_{H,l} \mid \{\tau_{H,r}\}_{r \in pa(l)} \right) - \prod_{l \notin ch(n)} \mathbb{P}_{\theta^{l}_{l},l}^{\pi_{l}} \left( \tau_{H,l} \mid \{\tau_{H,r}\}_{r \in pa(l)} \right) \right|. \\
& = \sum_{t=1}^{k} \sum_{\tau_{H,r},\tau_{r \notin ch(n)}} \left| \prod_{l \notin ch(n)} \mathbb{P}_{\theta^{l}_{l},l}^{\pi_{l}} \left( \tau_{H,l} \mid \{\tau_{H,r}\}_{r \in pa(l)} \right) - \prod_{l \notin ch(n)} \mathbb{P}_{\theta^{l}_{l},l}^{\pi_{l}} \left( \tau_{H,l} \mid \{\tau_{H,r}\}_{r \in pa(l)} \right) \right|. \\
& = \sum_{t=1}^{k} \sum_{\tau_{H,r},\tau_{T \notin ch(n)}} \mathbb{P}_{\theta^{l}_{L},\ell_{T}}^{\pi_{L}} \left\{ \tau_{H,r}\}_{r \in pa(l)} \mathbb{P}_{\theta^{l}_{L},\ell_{T}}^{\pi_{L}} \left\{ \tau_{H,r}\}_{r \in pa(l)} \right\} \right|. \\
& = \sum_{t=1}^{k} \sum_{\tau_{H,r},\tau_{T \notin ch(n)}} \mathbb{P}_{\theta^{l}_{L},\ell_{T}}^{\pi_{L}} \left\{ \tau_{H,r}\}_{r \in pa(l)} \mathbb{P}_{\theta^{l}_{L},\ell_{T}}^{\pi_{L}} \left\{ \tau_{H,r}\}_{r \in pa(l)} \mathbb{P}_{\theta^{l}_{L},\ell_{T}}^{\pi_{L}} \right\} \right|. \\
& = \sum_{t=1}^{k} \sum_{\tau_{$$

For the clarity of presentation, we define the following notations: **Definition D.2.** We define  $R_1(k)$  and  $R_2(k)$  as

$$\begin{array}{c} \mathbf{1506} \\ \mathbf{1507} \\ \mathbf{1508} \\ \mathbf{1509} \\ \mathbf{1509} \\ \mathbf{1510} \\ \mathbf{1510} \\ \mathbf{1510} \\ \mathbf{1510} \\ \mathbf{1510} \end{array} \xrightarrow{k} \left\| \left[ \prod_{l \in \mathsf{ch}(n)} \mathbb{P}_{\theta_{l}^{t}, l}^{\pi_{l}} \right] \mathbb{P}_{\theta_{n}^{t}}^{\pi_{n}^{t}} - \left[ \prod_{l \in \mathsf{ch}(n)} \mathbb{P}_{\theta_{n}^{t}}^{\pi_{l}} \right] \mathbb{P}_{\theta_{n}^{t}}^{\pi_{n}^{t}} \left\| \left[ \prod_{l \notin \overline{\mathsf{ch}}(n)} \mathbb{P}_{\theta_{l}^{t}}^{\pi_{l}} \right] \right\|,$$

1510  
1511 
$$R_{2}(k) = \sum_{t=1}^{k} \sum_{\{\tau_{H,r}\}_{r\notin \overline{ch}(n)}} \left| \prod_{l\notin \overline{ch}(n)} \mathbb{P}_{\theta_{l}^{t}}^{\pi_{l}} \left( \tau_{H,l} \mid \{\tau_{H,r}\}_{r\in \mathsf{pa}(l)} \right) - \prod_{l\notin \overline{ch}(n)} \mathbb{P}_{\theta_{l}^{t}}^{\pi_{l}} \left( \tau_{H,l} \mid \{\tau_{H,r}\}_{r\in \mathsf{pa}(l)} \right) \right|$$

We then have 

$$\operatorname{Regret}(k) \leq H \cdot (R_1(k) + R_2(k))$$

**Lemma D.6.** With probability at least  $1 - \delta$ , we have the following bound on  $R_2(k)$ , for all  $k \in [T]$ .

$$R_2(k) \leq \mathcal{O}\left(\sqrt{k\tilde{\beta}}\right).$$

*Proof.* According to Lemma D.3, we have with probability at least  $1 - \delta$  for all  $k \in [T]$ ,

$$\sum_{t=1}^{k-1} \sum_{\{\tau_r\}_{r\notin\overline{ch}(n)}} \left| \prod_{l\notin\overline{ch}(n)} \mathbb{P}_{\theta_l^k}^{\pi_l} \left( \tau_l \mid \{\tau_r\}_{r\in\mathsf{pa}(l)} \right) - \prod_{l\notin\overline{ch}(n)} \mathbb{P}_{\theta_l^*}^{\pi_l} \left( \tau_l \mid \{\tau_r\}_{r\in\mathsf{pa}(l)} \right) \right| = \mathcal{O}\left(\sqrt{k\tilde{\beta}}\right).$$

I

Then we can straightforwardly obtain that  $R_2(k) \leq \mathcal{O}\left(\sqrt{k\tilde{\beta}}\right)$ . 

**Lemma D.7.**  $R_1(k)$  is bounded by the following inequality: 

$$\begin{aligned} & 1530 \\ & 1531 \\ & 1532 \\ & 1532 \\ & 1532 \\ & 1532 \\ & 1532 \\ & 1533 \\ & 1534 \\ & 1534 \\ & 1535 \\ & 1535 \\ & 1535 \\ & 1536 \\ & 1536 \\ & 1537 \\ & 1537 \\ & 1538 \end{aligned} \\ & \left\{ \frac{\sqrt{S}}{\alpha} \left\| \left( \mathbb{M}_{h,n}^{t} \left( \{a_{h,r}^{T}\}_{r \in \mathsf{pa}(n)}^{r \notin \overline{\mathsf{ch}}(n)} \right) - \mathbb{M}_{h,n}^{\star} \left( \{a_{h,r}^{T}\}_{r \in \mathsf{pa}(n)}^{r \notin \overline{\mathsf{ch}}(n)} \right) \right) - \mathbb{M}_{h,c_{n,m}}^{\star} \left( \{a_{h,r}^{T}\}_{r \in \mathsf{pa}(c_{n,m})}^{r \notin \overline{\mathsf{ch}}(n)} \right) \right) \right\| \\ & \left[ \prod_{j=m+1}^{n-1} \left\| \mathbf{m}_{h,c_{n,j}}^{t} \right\|_{1}^{j} \pi_{c_{n,j}}^{t} (\tau_{h,c_{n,j}}) \right] \left[ \prod_{j=1}^{m-1} \left\| \mathbf{m}_{h,c_{n,j}}^{t} \right\|_{1}^{j} \pi_{c_{n,j}}^{t} (\tau_{h,c_{n,j}}) \right] \left\| \mathbf{m}_{h,n} \right\|_{1}^{j} \pi_{n}^{t} (\tau_{h,n}) \\ & \left[ \prod_{j=1}^{n-1} \left\| \mathbf{m}_{h,c_{n,j}} \right\|_{1}^{j} \pi_{c_{n,j}}^{t} (\tau_{h,c_{n,j}}) \right] \right\| \\ & \left[ \prod_{j=1}^{n-1} \left\| \mathbf{m}_{h,c_{n,j}} \right\|_{1}^{j} \pi_{c_{n,j}}^{t} (\tau_{h,c_{n,j}}) \right] \right] \\ & \left[ \prod_{j=1}^{n-1} \left\| \mathbf{m}_{h,c_{n,j}} \right\|_{1}^{j} \pi_{c_{n,j}}^{t} (\tau_{h,c_{n,j}}) \right] \right] \\ & \left[ \prod_{j=1}^{n-1} \left\| \mathbf{m}_{h,c_{n,j}} \right\|_{1}^{j} \pi_{c_{n,j}}^{t} (\tau_{h,c_{n,j}}) \right] \right] \\ & \left[ \prod_{j=1}^{n-1} \left\| \mathbf{m}_{h,c_{n,j}} \right\|_{1}^{j} \pi_{c_{n,j}}^{t} (\tau_{h,c_{n,j}}) \right] \right] \\ & \left[ \prod_{j=1}^{n-1} \left\| \mathbf{m}_{h,c_{n,j}} \right\|_{1}^{j} \pi_{c_{n,j}}^{t} (\tau_{h,c_{n,j}}) \right] \right] \\ & \left[ \prod_{j=1}^{n-1} \left\| \mathbf{m}_{h,c_{n,j}} \right\|_{1}^{j} \pi_{c_{n,j}}^{t} (\tau_{h,c_{n,j}}) \right] \right\|_{1}^{j} \pi_{c_{n,j}}^{t} (\tau_{h,c_{n,j}}) \right\|_{1}^{j} \pi_{c_{n,j}}^{t} (\tau_{h,c_{n,j}}) \left\| \mathbf{m}_{h,c_{n,j}} \right\|_{1}^{j} \pi_{c_{n,j}}^{t} (\tau_{h,c_{n,j}}) \right\|_{1}^{j} \pi_{c_{n,j}}^{t} (\tau_{h,c_{n,j}}) \|_{1}^{j} \pi_{c_{n,j}}^{t} (\tau_{h,c_{n$$

where for all  $m \in \overline{ch}(n)$ ,  $h \in [H]$ ,  $t \in [T]$ , we define  $\mathbf{m}_{h,m}$  and  $\mathbf{m}_{h,m}^t$  as

$$\mathbf{m}_{h,m} = \prod_{h'=0}^{h} \mathbb{M}_{h,m}^{\star} \left( \left\{ a_{h,r}^{T} \right\}_{r \in \mathsf{pa}(m)}^{r \notin \overline{\mathsf{ch}}(n)} \right), \quad \mathbf{m}_{h,m}^{t} = \prod_{h'=0}^{h} \mathbb{M}_{h,m}^{t} \left( \left\{ a_{h,r}^{T} \right\}_{r \in \mathsf{pa}(m)}^{r \notin \overline{\mathsf{ch}}(n)} \right).$$

*Proof.* According to the definition of  $R_1(k)$ , we can obtain the following bound on  $R_1(k)$ :

$$R_{1}(k) \leq \sum_{t=1}^{k} \sum_{\tau_{H}} \sum_{m=1}^{\mathsf{chl}(n)-1} \left[ \prod_{l=1}^{m-1} \mathbb{P}_{\theta_{l}^{t}}^{\pi_{l}} \right] \left| \mathbb{P}_{\theta_{m}^{t}}^{\pi_{m}} - \mathbb{P}_{\theta_{m}^{t}}^{\pi_{m}} \right| \left[ \prod_{j=m+1}^{n-1} \mathbb{P}_{\theta_{j}^{t}}^{\pi_{j}} \right] \mathbb{P}_{\theta_{n}^{t}}^{\pi_{n}^{t}} \left[ \prod_{l \notin \overline{\mathsf{ch}}(n)} \mathbb{P}_{\theta_{l}^{t}}^{\pi_{l}} \right] + \sum_{t=1}^{k} \sum_{\tau_{H}} \left[ \prod_{l=1}^{n-1} \mathbb{P}_{\theta_{l}^{t}}^{\pi_{l}} \right] \left| \mathbb{P}_{\theta_{n}^{t}}^{\pi_{n}^{t}} - \mathbb{P}_{\theta_{n}^{t}}^{\pi_{n}^{t}} \right| \left[ \prod_{l \notin \overline{\mathsf{ch}}(n)} \mathbb{P}_{\theta_{l}^{t}}^{\pi_{l}} \right].$$

We can deduce the result of lemma by representing this inequality with matrix notations. 

**Lemma D.8.** With probability at least 
$$1 - \delta$$
, for all  $(k, h, m) \in [K] \times [H - 1] \times |ch(n)| - 1$ , the following inequality holds true.

$$\begin{split} & \sum_{t=1}^{1558} \sum_{t=1}^{k-1} \sum_{\{\tau_r\}_{r \in \overline{ch}(n)}} \pi_{c_{n,m}}^t(\tau_{h,c_{n,m}}) \left\| \left( \mathbb{M}_{h,c_{n,m}}^t \left( \{a_{h,r}^T\}_{r \in \mathsf{pa}(c_{n,m})}^{r \notin \overline{ch}(n)} \right) - \mathbb{M}_{h,c_{n,m}}^{\star} \left( \{a_{h,r}^T\}_{r \in \mathsf{pa}(c_{n,m})}^{r \notin \overline{ch}(n)} \right) \right) \mathbf{m}_{h-1,c_{n,m}} \right\|_{1} \\ & \sum_{t=1}^{1561} \sum_{i=m+1}^{n-1} \left\| \mathbf{m}_{h,c_{n,i}}^t \right\|_{1} \pi_{c_{n,i}}(\tau_{h,c_{n,i}}^t) \right\| \left[ \prod_{j=1}^{m-1} \left\| \mathbf{m}_{h,c_{n,j}}^t \right\|_{1} \pi_{c_{n,j}}^t(\tau_{h,c_{n,j}}) \right] \left\| \mathbf{m}_{h,n} \right\|_{1} \pi_{n}^t(\tau_{h,n}) \\ & \sum_{t=1}^{1563} \frac{S^{1.5}OA^{|\mathsf{pa}(c_{n,m})\cap\overline{ch}(n)|+1}H^2}{\alpha} \sqrt{k\beta_{c_{n,m}}}, \end{split}$$

and for agent n, the following inequality holds true:  $\sum_{t=1}^{n} \sum_{\{\tau_{n}\}_{n=1}} \pi_{n}^{t}(\tau_{h,n}) \left\| \left( \mathbb{M}_{h,n}^{t} \left( \{a_{h,r}^{T}\}_{r\in\mathsf{pa}(n)}^{r\notin\overline{\mathsf{ch}}(n)} \right) - \mathbb{M}_{h,n}^{\star} \left( \{a_{h,r}^{T}\}_{r\in\mathsf{pa}(n)}^{r\notin\overline{\mathsf{ch}}(n)} \right) \right) \mathbf{m}_{h-1,n} \right\|_{1} \left\| \prod_{i=1}^{n-1} \| \mathbf{m}_{h,c_{n,i}} \|_{1} \pi_{c_{n,i}}^{t}(\tau_{h,c_{n,i}}) \right\|$  $\lesssim \frac{S^{1.5}OA^{|\mathsf{pa}(n)\cap\overline{\mathsf{ch}}(n)|+1}H^2}{\alpha}\sqrt{k\beta_n},$ where for all  $m \in \overline{ch}(n)$ ,  $h \in [H]$ ,  $t \in [T]$ , we define  $\mathbf{m}_{h,m}$  and  $\mathbf{m}_{h,m}^t$  as  $\mathbf{m}_{h,m} = \prod_{l \neq 0}^{h} \mathbb{M}_{h,m}^{\star} \left( \{a_{h,r}^{T}\}_{r \in \mathsf{pa}(m)}^{r \notin \overline{\mathsf{ch}}(n)} \right), \quad \mathbf{m}_{h,m}^{t} = \prod_{l \neq 0}^{h} \mathbb{M}_{h,m}^{t} \left( \{a_{h,r}^{T}\}_{r \in \mathsf{pa}(m)}^{r \notin \overline{\mathsf{ch}}(n)} \right).$ *Proof.* Initially, for  $m \in |\overline{ch}(n)| - 1$ , we fix  $(o, a, \{a_r\}_{r \in pa(c_{n,m}) \cap \overline{ch}(n)}) \in \mathcal{O}_{c_{n,m}} \times \mathcal{A}_{c_{n,m}} \times \mathcal{A}_{c_$  $\left(\times_{j\in \mathsf{pa}(c_{n,m})\cap \overline{\mathsf{ch}}(n)}\mathcal{A}_{j}\right)$ , we first define set  $S_{c_{n,m}}$  as  $\mathcal{C}_{c_{n,m}} = \left\{ \{\tau_{h,r}\}_{r \in \overline{\mathsf{ch}}(n)} \mid \{\tau_{h,r}\}_{r \in \overline{\mathsf{ch}}(n)} : \left(o_{h,c_{n,m}}, a_{h,c_{n,m}}, \{a_{h,r}\}_{r \in \mathsf{pa}(c_{n,m}) \cap \overline{\mathsf{ch}}(n)}\right) = \left(o,a, \{a_r\}_{r \in \mathsf{pa}(c_{n,m}) \cap \overline{\mathsf{ch}}(n)}\right) \right\}.$ We assume that  $w_{t,l} = \left[ \left( \tilde{\mathbb{M}}_{h,c_{n,m}}^{t} - \tilde{\mathbb{M}}_{h,c_{n,m}}^{\star} \right) \mathbb{O}_{h,c_{n,m}}^{\star} \right],$ where we denote  $\tilde{\mathbb{M}}_{h,c_{n,m}}^t$  and  $\tilde{\mathbb{M}}_{h,c_{n,m}}^{\star}$  as  $\tilde{\mathbb{M}}_{h,c_{n,m}}^t = \mathbb{M}_{h,c_{n,m}}^t \left( o, a, \{a_r\}_{r \in \mathsf{pa}(c_{n,m}) \cap \overline{\mathsf{ch}}(n)}, \{a_{h,r}^T\}_{r \in \mathsf{pa}(c_{n,m})}^{r \notin \overline{\mathsf{ch}}(n)} \right),$  $\tilde{\mathbb{M}}_{h,c_{n,m}}^{\star} = \mathbb{M}_{h,c_{n,m}}^{\star} \left( o, a, \{a_r\}_{r \in \mathsf{pa}(c_{n,m}) \cap \overline{\mathsf{ch}}(n)}, \{a_{h,r}^T\}_{r \in \mathsf{pa}(c_{n,m})}^{r \notin \overline{\mathsf{ch}}(n)} \right).$ We denote the sequence  $\pi_{c_{n,m}}^t(\tau_{h,c_{n,m}})\mathbb{O}_{h,c_{n,m}}^{\dagger}\mathbf{m}_{h-1,c_{n,m}} \cdot [\prod_{l=1}^{m-1} \|\mathbf{m}_{h,c_{n,l}}\|_1 \pi_{c_{n,l}}(\tau_{h,c_{n,l}})]$ .  $[\prod_{j=m+1}^{n-1} \|\mathbf{m}_{h,c_{n,j}}^t\|_1 \pi_{c_{n,j}}^t(\tau_{h,c_{n,j}})] \cdot \|\mathbf{m}_{h,n}^t\|_1 \pi_n^t(\tau_{h,n}) \text{ for all } \{\tau_{h,r}\}_{r\in\overline{ch}(n)} \in \mathcal{C}_{c_n,m} \text{ by } x_{t,1}, x_{t,2}, \ldots, x_{t,N}, \text{ where } N = |\mathcal{C}_{c_{n,m}}|. \text{ Then we have two observations about the } x, w \text{ sequence:}$ The vector sequence  $\{x_{t,i}\}_{i=1}^N$  satisfies  $\sum_{i=1}^N ||x_{t,i}||_1 \le 1$  for all t because  $\sum_{\{\tau_{h,r}\}_{r\in\overline{ch}(n)}\in\mathcal{C}_{c_{n,m}}}\pi_{c_{n,m}}(\tau_{h,c_{n,m}})\left\|\mathbb{O}_{h,c_{n,m}}^{\dagger}\mathbf{m}_{h-1,c_{n,m}}\right\|_{1}\left[\prod_{l=1}^{m-1}\left\|\mathbf{m}_{h,c_{n,l}}\right\|_{1}\pi_{l}(\tau_{h,c_{n,l}})\right]$  $\cdot \left[ \prod_{i=m+1}^{n-1} \left\| \mathbf{m}_{h,c_{n,j}}^t \right\|_1 \pi_j(\tau_{h,c_{n,j}}) \right] \left\| \mathbf{m}_{h,n}^t \right\|_1 \pi_n^t(\tau_{h,n})$  $\leq \sum_{\{\tau_{h,r}\}_{r\in\overline{ch}(n)}:(o_{h,c_{n,m}},a_{h,c_{n,m}})=(o,a)} \pi_{c_{n,m}}(\tau_{h,c_{n,m}}) \left\| \mathbb{O}_{h,c_{n,m}}^{\dagger} \mathbf{m}_{h-1,c_{n,m}} \right\|_{1} \left[ \prod_{l=1}^{m-1} \left\| \mathbf{m}_{h,c_{n,l}} \right\|_{1} \pi_{l}(\tau_{h,c_{n,l}}) \right]$  $\cdot \left| \prod_{i=m+1}^{n-1} \left\| \mathbf{m}_{h,c_{n,j}}^{t} \right\|_{1} \pi_{j}(\tau_{h,c_{n,j}}) \right| \cdot \left\| \mathbf{m}_{h,n}^{t} \right\|_{1} \pi_{n}^{t}(\tau_{h,n})$  $\leq \sum_{\{\tau_{h,r}\}_{r\in\overline{ch}(n)/\{c_{n-m}\}},\tau_{h-1,c_{n-m}}} \pi_{c_{n,m}}(\tau_{h-1,c_{n,m}}) \left\| \mathbb{O}_{h,c_{n,m}}^{\dagger} \mathbf{m}_{h-1,c_{n,m}} \right\|_{1} \left[ \prod_{r=1}^{m-1} \left\| \mathbf{m}_{h,c_{n,l}} \right\|_{1} \pi_{l}(\tau_{h,c_{n,l}}) \right]$  $\cdot \left| \prod_{i=m+1}^{n-1} \left\| \mathbf{m}_{h,c_{n,j}}^t \right\|_1 \pi_j(\tau_{h,c_{n,j}}) \right| = 1.$ 

1620 The vectors 
$$\{w_{t,l}\}_{l=1}^{O}$$
 satisfy  $\sum_{l=1}^{O} \|w_{t,l}\|_1 \le 2S^{1.5}/\alpha$  for all  $t$ , since we have  
1621  $\sum_{l=1}^{O} \|w_{t,l}\|_1 = \left\| \left( \tilde{\mathbb{M}}_{h,c_{n,m}}^t - \tilde{\mathbb{M}}_{h,c_{n,m}}^\star \right) \mathbb{O}_{h,c_{n,m}^\star} \right\|_1 \le S \left( \left\| \tilde{\mathbb{M}}_{h,c_{n,m}}^t \right\|_1 + \left\| \tilde{\mathbb{M}}_{h,c_{n,m}}^\star \right\|_1 \right) \le \frac{2S^{1.5}}{\alpha}.$ 
1624 1621

Using the notation of  $\{x_{t,i}\}_{i=1}^n$  and  $\{w_{t,l}\}_{l=1}^O$ , we have 

$$\sum_{t=1}^{k-1} \sum_{l=1}^{O} \sum_{i=1}^{n} |w_{k,l}^T x_{t,i}| = \mathcal{O}\left(\frac{\sqrt{S}}{\alpha} \sqrt{k\beta_{c_{n,m}}}\right).$$

Therefore, we can bind the target term with Eluder-Dimension lemma (Proposition 22 of Liu et al. (2022a)). We have the following result.

$$\sum_{t=1}^{k} \sum_{l=1}^{O} \sum_{i=1}^{n} |w_{t,l}^{T} x_{t,i}| = \mathcal{O}\left(\frac{S^{1.5} H^{2}}{\alpha} \sqrt{k\beta_{c_{n,m}}}\right).$$

The equation holds for all  $k \in [K]$ . We represent it with matrix operator, and we have 

$$\begin{aligned}
& \text{1637} \\
& \text{1638} \\
& \text{1639} \\
& \text{1639} \\
& \text{1640} \\
& \text{1640} \\
& \text{1641} \\
& \text{1641} \\
& \text{1642} \\
& \text{1642} \\
& \text{1643} \\
& \text{1643} \\
& \text{1643} \\
& \text{1644} \\
& \text{We sum up both the left head side and the right head side for all (a, a, (a, b)) = \mathcal{O}\left(\frac{S^{1.5}H^2}{\alpha}\sqrt{k\beta_{c_{n,m}}}\right). \\
& \text{1644} \\
& \text{1644} \\
& \text{We sum up both the left head side and the right head side for all (a, a, (a, b)) = \mathcal{O}\left(\frac{S^{1.5}H^2}{\alpha}\sqrt{k\beta_{c_{n,m}}}\right).
\end{aligned}$$

We sum up both the left-hand side and the right-hand side for all  $(o, a, \{a_r\}_{r \in \mathsf{pa}(c_{n,m}) \cap \overline{\mathsf{ch}}(n)}) \in \mathbb{R}$ × A × ( ×  $(A_i)$  and we can obtain that  $\mathcal{O}$ 

Then, with similar techniques, we consider the term

$$\sum_{t=1}^{k-1} \sum_{\{\tau_r\}_{r\in\overline{ch}(n)}} \pi_n^t(\tau_{h,n}) \left\| \left( \mathbb{M}_{h,n}^t \left( \{a_{h,r}^T\}_{r\in\mathsf{pa}(n)}^{r\notin\overline{ch}(n)} \right) - \mathbb{M}_{h,n}^\star \left( \{a_{h,r}^T\}_{r\in\mathsf{pa}(n)}^{r\notin\overline{ch}(n)} \right) \right) \mathbf{m}_{h-1,n} \right\|_1 \left[ \prod_{l=1}^{n-1} \left\| \mathbf{m}_{h,c_{n,l}} \right\|_1 \pi_{c_{n,l}}^t(\tau_{h,c_{n,l}}) \right].$$

We define set  $S_n$  as

$$\mathcal{C}_n = \left\{ \{\tau_{h,r}\}_{r \in \overline{\mathsf{ch}}(n)} \mid \{\tau_{h,r}\}_{r \in \overline{\mathsf{ch}}(n)} : \left(o_{h,n}, a_{h,n}, \{a_{h,r}\}_{r \in \mathsf{pa}(n) \cap \overline{\mathsf{ch}}(n)}\right) = \left(o, a, \{a_r\}_{r \in \mathsf{pa}(n) \cap \overline{\mathsf{ch}}(n)}\right) \right\}$$

We assume that

$$w_{t,l} = \left[ \left( \tilde{\mathbb{M}}_{h,n}^t - \tilde{\mathbb{M}}_{h,n}^\star \right) \mathbb{O}_{h,n}^\star \right]_l$$

where we denote  $\tilde{\mathbb{M}}_{h,n}^t$  and  $\tilde{\mathbb{M}}_{h,n}^{\star}$  as 

1668  
1669 
$$\tilde{\mathbb{M}}_{h,n}^{t} = \mathbb{M}_{h,n}^{t} \left( o, a, \{a_r\}_{r \in \mathsf{pa}(n) \cap \overline{\mathsf{ch}}(n)}, \{a_{h,r}^T\}_{r \in \mathsf{pa}(n)}^{r \notin \overline{\mathsf{ch}}(n)} \right),$$

1670  
1671 
$$\tilde{\mathbb{M}}_{h,n}^{\star} = \mathbb{M}_{h,n}^{\star} \left( o, a, \{a_r\}_{r \in \mathsf{pa}(n) \cap \overline{\mathsf{ch}}(n)}, \{a_{h,r}^T\}_{r \in \mathsf{pa}(n)}^{r \notin \overline{\mathsf{ch}}(n)} \right).$$

We denote the target sequence  $\pi_n^t(\tau_{h,n}) \mathbb{O}_{h,n}^{\dagger} \mathbf{m}_{h-1,n} \cdot [\prod_{l=1}^{n-1} ||\mathbf{m}_{h,c_{n,l}}||_1 \pi_{c_{n,l}}(\tau_{h,c_{n,l}})]$  for all  $\{\tau_{h,r}\}_{r\in\overline{ch}(n)} \in S_n$  by  $x_{t,1}, x_{t,2}, \ldots, x_{t,N}$ , where  $N = |\mathcal{C}_n|$ . Then we have two observations 

about the x, w sequence: The vector sequence  $\{x_{t,i}\}_{i=1}^N$  satisfies  $\sum_{i=1}^N ||x_{t,i}||_1 \le 1$  for all t because  $\sum_{\{\tau_{h,r}\}_{r\in\overline{ch}(n)}\in\mathcal{C}_{n}}\pi_{n}^{t}(\tau_{h,n})\left\|\mathbb{O}_{h,n}^{\dagger}\mathbf{m}_{h,n}\right\|_{1}\left[\prod_{l=1}^{n-1}\left\|\mathbf{m}_{h,c_{n,l}}\right\|_{1}\pi_{l}^{t}(\tau_{h,c_{n,l}})\right]$  $\leq \sum_{\{\tau_{h,r}\}_{r\in\overline{ch}(n)}:(o_{h,n},a_{h,n})=(o,a)} \pi_{n}^{t}(\tau_{h,n}) \left\| \mathbb{O}_{h,n}^{\dagger}\mathbf{m}_{h-1,n} \right\|_{1} \left[ \prod_{l=1}^{n-1} \left\| \mathbf{m}_{h,c_{n,l}} \right\|_{1} \pi_{l}^{t}(\tau_{h,c_{n,l}}) \right]$  $\leq \sum_{\{\tau_{h,r}\}_{r \in \mathsf{ch}(n)}, \tau_{h-1,n}} \pi_n^t(\tau_{h-1,n}) \left\| \mathbb{O}_{h,n}^{\dagger} \mathbf{m}_{h-1,n} \right\|_1 \left[ \prod_{l=1}^{n-1} \left\| \mathbf{m}_{h,c_{n,l}} \right\|_1 \pi_l^t(\tau_{h,c_{n,l}}) \right] = 1,$ 

The vectors  $\{w_{t,l}\}_{l=1}^{O}$  satisfy  $\sum_{l=1}^{O} ||w_{t,l}||_1 \leq 2S^{1.5}/\alpha$  for all t, since we have

$$\sum_{l=1}^{O} \|w_{t,l}\|_{1} = \left\| \left( \tilde{\mathbb{M}}_{h,n}^{t} - \tilde{\mathbb{M}}_{h,n}^{\star} \right) \mathbb{O}_{h,n}^{\star} \right\|_{1} \le S \left( \left\| \tilde{\mathbb{M}}_{h,n}^{t} \right\|_{1} + \left\| \tilde{\mathbb{M}}_{h,n}^{\star} \right\|_{1} \right) \le 2S^{1.5} / \alpha.$$

Using the notation of  $\{x_{t,i}\}_{i=1}^n$  and  $\{w_{t,l}\}_{l=1}^O$ , we have 

$$\sum_{k=1}^{n} \sum_{l=1}^{O} \sum_{i=1}^{n} |w_{k,l}^T x_{t,i}| = \mathcal{O}\left(\frac{\sqrt{S}}{\alpha} \sqrt{k\beta_n}\right).$$

Therefore, we can bind the target term with Eluder-Dimension lemma (Proposition 22 of Liu et al. (2022a)). We have the following result. 

$$\sum_{i=1}^{k} \sum_{l=1}^{O} \sum_{i=1}^{n} |w_{t,l}^T x_{t,i}| = \mathcal{O}\left(\frac{S^{1.5} H^2}{\alpha} \sqrt{k\beta_n}\right)$$

The equation holds for all  $k \in [K]$ . We represent it with *B*-operator, and we have 

$$\sum_{t=1}^{1700} \sum_{t=1}^{k} \sum_{\{\tau_{h,r}\}_{r\in\overline{ch}(n)}\in\mathcal{C}_{n}} \left\| \left( \tilde{\mathbb{M}}_{h,n}^{t} - \tilde{\mathbb{M}}_{h,n}^{\star} \right) \mathbf{m}_{h-1,n} \right\|_{1} \cdot \pi_{n}^{t}(\tau_{h,n}) \left[ \prod_{l=1}^{n-1} \| \mathbf{m}_{h,c_{n,l}} \|_{1} \pi_{c_{n,l}}^{t}(\tau_{h,c_{n,l}}) \right] = \mathcal{O}\left( \frac{S^{1.5}H^{2}}{\alpha} \sqrt{k\beta_{c_{n,m}}} \right)$$

$$\text{We sum up both the left-hand side and the right-hand side for all } \left( o, a, \{a_{r}\}_{c_{n,l}} \circ (\cdots) \circ \overline{(\cdot)} \right) \in \mathbb{C}$$

We sum up both the left-hand side and the right-hand side for all  $(o, a, \{a_r\}_{r \in pa(c_{n,m}) \cap \overline{ch}(n)}) \in C_{r,m}$ 

 $\mathcal{O}_{c_{n,m}} \times \mathcal{A}_{c_{n,m}} \times \left( \times_{j \in \mathsf{pa}(c_{n,m}) \cap \overline{\mathsf{ch}}(n)} \mathcal{A}_{j} \right)$ , and we can obtain that 

$$\sum_{t=1}^{\kappa} \sum_{\boldsymbol{\tau}_{h}} \pi_{n}^{t}(\boldsymbol{\tau}_{h,n}) \left\| \left( \mathbb{M}_{h,n}^{t} \left( \{a_{h,r}^{T}\}_{r \in \mathsf{pa}(n)}^{r \notin \overline{\mathsf{ch}}(n)} \right) - \mathbb{M}_{h,n}^{\star} \left( \{a_{h,r}^{T}\}_{r \in \mathsf{pa}(n)}^{r \notin \overline{\mathsf{ch}}(n)} \right) \right) \mathbf{m}_{h-1,n} \right\|_{1}$$
$$\cdot \left[ \prod_{l=1}^{n-1} \| \mathbf{m}_{h,l} \|_{1} \pi_{l}^{t}(\boldsymbol{\tau}_{h,l}) \right] = \mathcal{O} \left( \frac{S^{1.5} H^{2} O A^{1+|\mathsf{pa}(c_{n,m}) \cap \overline{\mathsf{ch}}(n)|}}{\alpha} \sqrt{k\beta_{c_{n,m}}} \right).$$

**Corollary D.1.** With probability at least  $1 - \delta$ , for all  $k \in [T]$ ,  $R_1(k)$  is bounded by the following inequality: 

$$R_1(k) \le \mathcal{O}\left(\frac{S^2 H^3 O A}{\alpha^2} \left[\sum_{m \in \overline{\mathsf{ch}}(n)} A^{|\mathsf{pa}(m) \cap \overline{\mathsf{ch}}(n)| + 1} \sqrt{k\beta_m}\right]\right).$$

Proof. We only need to combine the result in Lemma D.3 and Lemma D.7, and we can achieve the result. 

**Corollary D.2.** With probability at least  $1 - \delta$ , for all  $k \in [T]$ ,  $R_1(k)$  is bounded by the following inequality: 

$$R_1(k) \le \mathcal{O}\left(\frac{S^2 H^3 O A}{\alpha^2} \left[\sum_{m \in \overline{\mathsf{ch}}(n)} A^{|\mathsf{pa}(m) \cap \overline{\mathsf{ch}}(n)| + 1} \sqrt{k\beta_m}\right]\right).$$

Proof. We only need to combine the result in Lemma D.3 and Lemma D.7, and we can achieve the result. 

Proof of Theorem D.1 With the lemmas presented above, we are now ready to prove Theorem D.1.
D.1.

**Theorem D.2.** If Algorithm 5.1 take agent *i* as central agent and take policies  $\{\pi_m\}_{m \in [n]/\{i\}}$  as input, then Algorithm 5.1 guarantees that with probability at least  $1 - \delta$  for all  $k \in [T]$ 

$$\left[\sum_{t=1}^{k} V^{\pi_i^*, \boldsymbol{\pi}_{-i}} - V^{\pi_i^k, \boldsymbol{\pi}_{-i}}\right] \leq \tilde{\mathcal{O}}\left(\frac{S^2 O A^{d_i+1}}{\alpha^2} \sqrt{K\left(S^2 A^{\mathsf{ind}(G)+1} + SO\right)} \cdot \operatorname{poly}(H)\right),$$

1736 where  $\pi_i^*$  is defined as  $\pi_i^* = \arg \max_{\pi_i} V^{\pi_i, \pi_{-i}}$ .

*Proof.* According to Corollary D.2 and Lemma D.6, we can obtain that with probability at least  $1 - \delta$ , for all  $k \in [T]$ ,

$$R^k \le H \cdot (R_1(k) + R_2(k))$$

$$\begin{split} &\leq \mathcal{O}\left(\frac{S^2 H^4 O A}{\alpha^2} \left[\sum_{m \in \overline{\mathsf{ch}}(n)} A^{|\mathsf{pa}(m) \cap \overline{\mathsf{ch}}(n)| + 1} \sqrt{k\beta_m}\right] + H \sqrt{k\tilde{\beta}}\right) \\ &\leq \tilde{\mathcal{O}}\left(\frac{S^2 O A^{\mathrm{indeg}(F_n) + 1}}{\alpha^2} \sqrt{K \left(S^2 A^m + SO\right)} \cdot \mathrm{poly}(H)\right). \end{split}$$

		. 1
		. 1

1750 D.2.2 PROOF FOR THEOREM 4.3

After proving Theorem D.1, we are ready to prove Theorem 4.3. For the readers' convenience, we first restate the theorem here.

**Theorem D.3.** If central agent for Algorithm D.1 is i, we define bonus parameter as  $\beta_m = c(H^2(S^2A^{|pa(m)|+1} + SO)\log(TSAOH) + \log(Tn/\delta)), \forall m \in [n], \tilde{\beta} = c((\sum_{r\notin ch(i)} H(S^2A^{ind(r)+1} + SO)\log(TSAOH) + \log(Tn/\delta))))$ . Then, with probability at least  $1-\delta$ , Algorithm D.1 terminates within  $4H/\epsilon$  steps of the while loop, and outputs an  $\epsilon$ -approximate local optimal policy. The total episodes of play in Algorithm D.1 is at most

$$K = \tilde{\mathcal{O}}\left(\sum_{m=1}^{n} S^4 O^2 A^{2 \cdot \operatorname{ind}(G[\operatorname{ch}(m) \cup \{m\}]) + 2} (S^2 A^{\operatorname{ind}(G) + 1} + SO) \cdot \operatorname{poly}(H) / (\alpha^4 \epsilon^3)\right).$$

Proof. We use superscript t to represent variables at the  $t^{th}$  step (before  $\pi$  is updated) of the while 

$$\max_{\mu_i} V(\mu_i, \pi_{-i}^t) - V(\hat{\pi}_i^t, \pi_{-i}^t) \le \frac{\epsilon}{4}.$$

1768 We then take a union bound over all  $t \le 4H/\epsilon$  and for all  $i \in [n]$ , and we have with probability at 1769 least  $1 - \delta$  the following inequality holds for all  $i \in [n]$  and  $t \le 4H/\epsilon$ 

$$\max_{\mu_i} V(\mu_i, \pi_{-i}^t) - V(\hat{\pi}_i^t, \pi_{-i}^t) \le \frac{\epsilon}{4}.$$
(16)

For the empirical estimator  $\hat{V}^t$ , it's bounded in [0, H]. Thus, by Hoeffding's inequality, for fixed  $i \in [n]$  and t

1775  
1776  
1777  
1777  

$$\mathbb{P}\left(\left|\hat{V}^t - V^t\right| \ge \frac{\epsilon}{8}\right) \le 2\exp\left(-\frac{N\epsilon^2}{32H^2}\right)$$

1779 Choosing  $N = \frac{CH^2}{\epsilon^2} \log\left(\frac{nHSK \max_{i \in [n]} A_i}{\epsilon \delta}\right)$  for some large constant C, we have 

1781 
$$\mathbb{P}\left(\left|\hat{V}^t - V^t\right| \ge \frac{\epsilon}{8}\right) \le \frac{\epsilon\delta}{16nH}.$$

1782 Applying this inequality for  $\hat{V}^t(\hat{\pi}_i^t, \pi_{-i}^t)$  and  $\hat{V}^t(\pi^t)$  and taking a union bound over  $i \in [n]$  and  $t \leq 4H/\epsilon$ , we can achieve that

$$|\hat{V}(\hat{\pi}_{i}^{t}, \pi_{-i}^{t}) - V(\hat{\pi}_{i}^{t}, \pi_{-i}^{t})| \le \epsilon/8, \quad |\hat{V}^{t}(\pi^{t}) - V(\pi^{t})| \le \epsilon/8.$$

We combine this result with equation 22, and we can obtain that with probability at least  $1 - \delta$ 

$$\max_{\mu_i} V(\mu_i, \pi_{-i}^t) - V(\hat{\pi}_i^t, \pi_{-i}^t) \le \frac{\epsilon}{4}, \quad |\hat{V}(\hat{\pi}_i^t, \pi_{-i}^t) - V(\hat{\pi}_i^t, \pi_{-i}^t)| \le \epsilon/8, \quad |\hat{V}^t(\pi^t) - V(\pi^t)| \le \epsilon/8$$

holds for all  $i \in [n]$  and  $t \leq \frac{4H}{\epsilon}$ . On this event,

$$\delta_i^t = \hat{V}^t(\hat{\pi}_i, \pi_{-i}) - \hat{V}^t(\pi^t) \le V(\hat{\pi}_i^t, \pi_{-i}^t) - V(\pi^t) + \epsilon/4.$$

1793 If the while loop doesn't end after the t-th iteration and  $t \le 4H/\epsilon$ , there exists  $j^t$  s.t.  $\Delta_{j^t}^t \ge \epsilon/2$ , 1794 so we have

$$V(\hat{\pi}_{j^t}^t, \pi_{-j^t}^t) - V(\pi^t) \ge \Delta_{j^t}^t - \epsilon/4 \ge \epsilon.$$

Since the value function is bound by H, so the while loop ends within 4H/epsilon steps. Therefore, the inequality above that holds for all  $i \in [n]$  and  $t \leq 4H/epsilon$  holds for simultaneously before the end of the while loop with probability at least  $1 - \delta$ . Again, on this event, if the while loop stops at the end of  $t^{th}$  step, we have  $\max_{i \in [n]} \delta_i^t \leq \epsilon/2$ , then

$$\begin{aligned} \max_{\mu_{i}} V(\mu_{i}, \pi_{-i}^{t}) - V(\pi^{t}) &= \max_{\mu_{i}} V(\mu_{i}, \pi_{-i}^{t}) - V(\hat{\pi}_{i}^{t}, \pi_{-i}^{t}) + V(\hat{\pi}_{i}^{t}, \pi_{-i}^{t}) - V(\pi^{t}) \\ &\leq \epsilon/4 + \hat{V}^{t}(\hat{\pi}_{i}^{t}, \pi_{-i}^{t}) - \hat{V}^{t}(\pi^{t}) + 2\epsilon/8 \\ &\leq \epsilon/2 + \Delta_{i}^{t} \\ &\leq \epsilon. \end{aligned}$$

So the returned policy  $\pi^t$  is an  $\epsilon$ -approximate local optimal. Therefore, we can conclude that probability at least  $1 - \delta$ , within  $4H/\epsilon$  steps of the while loop, Algorithm E.1 outputs an  $\epsilon$ -approximate local optimal policy.

Eventually, we compute the total number of episodes as the total sample complexity. According to the definition of N and  $K_i$  for all  $i \in [n]$ , we can obtain that

$$K = \frac{4H}{\epsilon} \left( N + \sum_{i=1}^{n} (K_i + N) \right) = \tilde{\mathcal{O}} \left( \sum_{i=1}^{n} \frac{S^4 O^2 A^{2 \cdot \mathsf{ind}(G[\mathsf{ch}(i) \cup \{i\}]) + 2} (S^2 A^{\mathsf{ind}(G) + 1} + SO) \cdot \mathsf{poly}(H)}{\alpha^4 \epsilon^3} \right).$$

1815 1816 1817

1818

1785

1787 1788

1791 1792

1795

1801 1802 1803

1805

### E SUPPLEMENTARY DETAILS FOR SECTION 5

# 1819 E.1 COMPLETE ALGORITHM FOR ACHIEVING LOCAL OPTIMAL WITH MEMORYLESS 1820 POLICIES

The complete Algorithm for achieving local optimal with memoryless policies is presented in Algorithm E.1. If the central agent of Algorithm 5.1 is n, then the bonus parameter  $\beta_m$  for  $m \in [n]$  is defined as

$$\beta_m = c \big( H(S^2 A^2 O^2 + O^2 S) \log(TSAOH) + \log(Tn/\delta) \big), m \in [n-1]$$
  
$$\beta_n = c \big( H(S^2 AO + OS) \log(TSAOH) + \log(Tn/\delta) \big).$$
(17)

1826 1827 1828

1825

**Trajectory Probability**  $\forall m \in [n]$ , we use  $\tau_m = (o_{1,m}, a_{1,m}, \dots, o_{H,m}, a_{H,m})$  to denote the trajectory of the  $m^{th}$  agent, and we denote the parameter  $\theta_m = (\mathbb{T}_m, \mathbb{O}_m, \mu_m)$  as the collection of parameters representing the joint probability of the  $m^{th}$  and  $i^{th}$  agent's trajectory. We define trajectory probability  $\mathbb{P}_{\theta_i,i}^{\pi_i}(\tau_i)$  and  $\mathbb{P}_{\theta_m,m}^{\pi_i,\pi_m}(\tau_i, \tau_m)$  for all  $m \in [n] \setminus \{i\}$  as follows:

$$\mathbb{P}_{\theta_{i},i}^{\pi_{i}}(\tau_{i}) = \sum_{s_{1},...,s_{H}} \mu(s_{1}) \mathbb{O}_{1,i}\pi_{1,i} \big[ \prod_{h=1}^{H-1} \mathbb{T}_{h,i,a_{h,i}} \mathbb{O}_{h+1,i}\pi_{h+1,i} \big], \\ \mathbb{P}_{\theta_{m},m}^{\pi_{i},\pi_{m}}(\tau_{i},\tau_{m}) = \sum_{s_{1},...,s_{H}} \mu(s_{1}) \mathbb{O}_{1,i}\pi_{1,i}\pi_{1,m} \big[ \prod_{h=1}^{H-1} \mathbb{T}_{h,i,a_{h,i},a_{h,m}} \mathbb{O}_{h+1,m}\pi_{h+1,i}\pi_{h+1,m} \big].$$
(18)

Alg	orithm 6 NASH-CA for Achieving Local Optimal with Memoryless Policies
1:	Initialize $\pi = {\pi_i}_{i \in [n]}$ , where $\pi_i = {\pi_{h,i}}_{(h,i) \in [H] \times [m]}$ .
2:	while true do $CH^2$ 1 ( $HCH^2$ 1 ( $HCH^2$ 1 ( $HCH^2$
3:	Execute policy $\pi$ for $N = \frac{GH}{\epsilon^2} \log \left( nHSK \max_{i \in [n]} A_i/(\epsilon \delta) \right)$ episodes and obtain $V_{1,i}(\pi)$ which is the empirical average of the total return under policy $\pi$
4:	for agent $i = 1,, m$ do
5:	Fix $\pi_{-i}$ let the $i^{th}$ agent be the central agent to run Algorithm E.1 for $K_i$ =
	$\mathcal{O}(S^4O^4A^4(S^2A^2O^2+SO^2)\times \operatorname{poly}(H)/(\alpha^4\epsilon^2))$ episodes and get a new policy $\hat{\pi}_i$ .
6:	Execute policy $(\hat{\pi}_i, \pi_{-i})$ for $N = \frac{CH^2}{\epsilon^2} \log(nHSK \max_{i \in [n]} A_i / (\epsilon \delta))$ episodes and
	obtain $V(\hat{\pi}_i, \pi_{-i})$ which is the empirical average of the total return under policy $(\hat{\pi}_i, \pi_{-i})$ .
7:	Set $\Delta_i \leftarrow V(\hat{\pi}_i, \pi_{-i}) - V(\pi)$ .
8:	if $\max_{i \in [n]} \Delta_i > \epsilon/2$ then Update $\pi_i$ ( $\hat{\pi}_i$ where $i = \arg \max_{i \in [n]} \Delta_i$
9: 10:	else opuate $\pi_j \leftarrow \pi_j$ where $j = \arg \max_{i \in [n]} \Delta_i$ .
1:	return $\pi$
١Į	orithm 7 OMLE for memoryless policy
	<b>Input:</b> central agent i and the policy for agent $[n]/(i)$ $\pi_i$ $\pi_2$ $\pi_2$ $\pi_3$ $\pi_4$ $\pi_5$
1:	<b>input</b> : central agent <i>i</i> , and the poincy for agent $[n]/\{i\}, n_1, n_2, \ldots, n_{i-1}, n_{i+1}, \ldots, n_n$ .
1: 2:	<b>Initialize:</b> $\mathcal{B}_i^1 = \{\hat{\theta}_i \in \Theta_i : \min_h \sigma_S(\mathbb{O}_i(\hat{\theta}_i) \ge \alpha)\}, \mathcal{B}_m^1 = \{\hat{\theta}_m \in \Theta_m : \min_h \sigma_S(\mathbb{O}_m(\hat{\theta}_m) \ge \alpha)\}$
1: 2:	<b>Initialize:</b> $\mathcal{B}_{i}^{1} = \{\hat{\theta}_{i} \in \Theta_{i} : \min_{h} \sigma_{S}(\mathbb{O}_{i}(\hat{\theta}_{i}) \geq \alpha)\}, \mathcal{B}_{m}^{1} = \{\hat{\theta}_{m} \in \Theta_{m} : \min_{h} \sigma_{S}(\mathbb{O}_{m}(\hat{\theta}_{m}) \geq \alpha/\sqrt{O})\}$ for all $m \in [n]/\{i\}$ . Set $\mathcal{D}_{m} = \{\}$ , for all agents $m \in [n]$ .
1: 2: 3:	<b>Initialize:</b> $\mathcal{B}_{i}^{1} = \{\hat{\theta}_{i} \in \Theta_{i} : \min_{h} \sigma_{S}(\mathbb{O}_{i}(\hat{\theta}_{i}) \geq \alpha)\}, \mathcal{B}_{m}^{1} = \{\hat{\theta}_{m} \in \Theta_{m} : \min_{h} \sigma_{S}(\mathbb{O}_{m}(\hat{\theta}_{m}) \geq \alpha/\sqrt{O})\}$ for all $m \in [n]/\{i\}$ . Set $\mathcal{D}_{m} = \{\}$ , for all agents $m \in [n]$ . <b>for</b> $k = 1 \dots T$ <b>do</b> <b>compute</b> $(\theta^{k}, \theta^{k}, \dots, \theta^{k}, \pi^{k}) = \arg\max_{i} \dots i = 1 \dots \sum_{i} n \dots V^{\pi_{i}, \pi_{-i}}(\hat{\theta}_{i})$
1: 2: 3: 4:	Initialize: $\mathcal{B}_{i}^{1} = \{\hat{\theta}_{i} \in \Theta_{i} : \min_{h} \sigma_{S}(\mathbb{O}_{i}(\hat{\theta}_{i}) \geq \alpha)\}, \mathcal{B}_{m}^{1} = \{\hat{\theta}_{m} \in \Theta_{m} : \min_{h} \sigma_{S}(\mathbb{O}_{m}(\hat{\theta}_{m}) \geq \alpha/\sqrt{O})\}$ for all $m \in [n]/\{i\}$ . Set $\mathcal{D}_{m} = \{\}$ , for all agents $m \in [n]$ . for $k = 1 \dots T$ do compute $(\theta_{1}^{k}, \theta_{2}^{k}, \dots, \theta_{n}^{k}, \pi_{i}^{k}) = \arg\max_{\hat{\theta}_{1} \in \mathcal{B}_{1}^{k}, \hat{\theta}_{2} \in \mathcal{B}_{2}^{k}, \dots, \hat{\theta}_{n} \in \mathcal{B}_{n}^{k}, \pi_{i}} \sum_{m=1}^{n} V_{m}^{\pi_{i}, \pi_{-i}}(\hat{\theta}_{m})$
1: 2: 3: 4: 5:	Initialize: $\mathcal{B}_{i}^{1} = \{\hat{\theta}_{i} \in \Theta_{i} : \min_{h} \sigma_{S}(\mathbb{O}_{i}(\hat{\theta}_{i}) \geq \alpha)\}, \mathcal{B}_{m}^{1} = \{\hat{\theta}_{m} \in \Theta_{m} : \min_{h} \sigma_{S}(\mathbb{O}_{m}(\hat{\theta}_{m}) \geq \alpha/\sqrt{O})\}$ for all $m \in [n]/\{i\}$ . Set $\mathcal{D}_{m} = \{\}$ , for all agents $m \in [n]$ . for $k = 1 \dots T$ do compute $(\theta_{1}^{k}, \theta_{2}^{k}, \dots, \theta_{n}^{k}, \pi_{i}^{k}) = \arg\max_{\hat{\theta}_{1} \in \mathcal{B}_{1}^{k}, \hat{\theta}_{2} \in \mathcal{B}_{2}^{k}, \dots, \hat{\theta}_{n} \in \mathcal{B}_{n}^{k}, \pi_{i}} \sum_{m=1}^{n} V_{m}^{\pi_{i}, \pi_{-i}}(\hat{\theta}_{m})$ follow $\pi^{k}$ to collect a trajectory $\tau^{k} = (\mathbf{o}_{1}^{k}, \mathbf{a}_{1}^{k}, \dots, \mathbf{o}_{k}^{k}, \mathbf{a}_{k}^{k})$ add $(\pi^{k}, \pi^{k}, \pi^{k})$ into $\mathcal{D}_{i}$ for $m \in [n]/\{i\}$ and add $(\pi^{k}, \pi^{k}, \pi^{k})$ into $\mathcal{D}_{i}$ for $m \in [n]/\{i\}$ and $d(\pi^{k}, \pi^{k})$ into $\mathcal{D}_{i}$ and under $\mathcal{B}^{k+1}$ and
1: 2: 3: 4: 5: 6:	Initialize: $\mathcal{B}_{i}^{1} = \{\hat{\theta}_{i} \in \Theta_{i} : \min_{h} \sigma_{S}(\mathbb{O}_{i}(\hat{\theta}_{i}) \geq \alpha)\}, \mathcal{B}_{m}^{1} = \{\hat{\theta}_{m} \in \Theta_{m} : \min_{h} \sigma_{S}(\mathbb{O}_{m}(\hat{\theta}_{m}) \geq \alpha/\sqrt{O})\}$ for all $m \in [n]/\{i\}$ . Set $\mathcal{D}_{m} = \{\}$ , for all agents $m \in [n]$ . for $k = 1 \dots T$ do compute $(\theta_{1}^{k}, \theta_{2}^{k}, \dots, \theta_{n}^{k}, \pi_{i}^{k}) = \arg \max_{\hat{\theta}_{1} \in \mathcal{B}_{1}^{k}, \hat{\theta}_{2} \in \mathcal{B}_{2}^{k}, \dots, \hat{\theta}_{n} \in \mathcal{B}_{n}^{k}, \pi_{i}} \sum_{m=1}^{n} V_{m}^{\pi_{i}, \pi_{-i}}(\hat{\theta}_{m})$ follow $\pi^{k}$ to collect a trajectory $\boldsymbol{\tau}^{k} = (\mathbf{o}_{1}^{k}, \mathbf{a}_{1}^{k}, \dots, \mathbf{o}_{H}^{k}, \mathbf{a}_{H}^{k})$ add $(\pi_{i}^{k}, \tau_{i}^{k}, \tau_{m}^{k})$ into $\mathcal{D}_{m}$ for $m \in [n]/\{i\}$ and add $(\pi_{i}^{k}, \tau_{i}^{k})$ into $\mathcal{D}_{i}$ , and update $\mathcal{B}_{i}^{k+1}$ and $\mathcal{B}^{k+1}$ for $m \in [n]/\{i\}$ as follows:
1: 2: 3: 4: 5: 6:	Initialize: $\mathcal{B}_{i}^{1} = \{\hat{\theta}_{i} \in \Theta_{i} : \min_{h} \sigma_{S}(\mathbb{O}_{i}(\hat{\theta}_{i}) \geq \alpha)\}, \mathcal{B}_{m}^{1} = \{\hat{\theta}_{m} \in \Theta_{m} : \min_{h} \sigma_{S}(\mathbb{O}_{m}(\hat{\theta}_{m}) \geq \alpha/\sqrt{O})\}$ for all $m \in [n]/\{i\}$ . Set $\mathcal{D}_{m} = \{\}$ , for all agents $m \in [n]$ . for $k = 1 \dots T$ do compute $(\theta_{1}^{k}, \theta_{2}^{k}, \dots, \theta_{n}^{k}, \pi_{i}^{k}) = \arg \max_{\hat{\theta}_{1} \in \mathcal{B}_{1}^{k}, \hat{\theta}_{2} \in \mathcal{B}_{2}^{k}, \dots, \hat{\theta}_{n} \in \mathcal{B}_{n}^{k}, \pi_{i}} \sum_{m=1}^{n} V_{m}^{\pi_{i}, \pi_{-i}}(\hat{\theta}_{m})$ follow $\pi^{k}$ to collect a trajectory $\boldsymbol{\tau}^{k} = (\mathbf{o}_{1}^{k}, \mathbf{a}_{1}^{k}, \dots, \mathbf{o}_{H}^{k}, \mathbf{a}_{H}^{k})$ add $(\pi_{i}^{k}, \tau_{i}^{k}, \tau_{m}^{k})$ into $\mathcal{D}_{m}$ for $m \in [n]/\{i\}$ and add $(\pi_{i}^{k}, \tau_{i}^{k})$ into $\mathcal{D}_{i}$ , and update $\mathcal{B}_{i}^{k+1}$ and $\mathcal{B}_{m}^{k+1}$ for $m \in [n]/\{i\}$ as follows:
1: 2: 3: 4: 5: 6:	Initialize: $\mathcal{B}_{i}^{1} = \{\hat{\theta}_{i} \in \Theta_{i} : \min_{h} \sigma_{S}(\mathbb{O}_{i}(\hat{\theta}_{i}) \geq \alpha)\}, \mathcal{B}_{m}^{1} = \{\hat{\theta}_{m} \in \Theta_{m} : \min_{h} \sigma_{S}(\mathbb{O}_{m}(\hat{\theta}_{m}) \geq \alpha/\sqrt{O})\}$ for all $m \in [n]/\{i\}$ . Set $\mathcal{D}_{m} = \{\}$ , for all agents $m \in [n]$ . for $k = 1 \dots T$ do compute $(\theta_{1}^{k}, \theta_{2}^{k}, \dots, \theta_{n}^{k}, \pi_{i}^{k}) = \arg \max_{\hat{\theta}_{1} \in \mathcal{B}_{1}^{k}, \hat{\theta}_{2} \in \mathcal{B}_{2}^{k}, \dots, \hat{\theta}_{n} \in \mathcal{B}_{n}^{k}, \pi_{i}} \sum_{m=1}^{n} V_{m}^{\pi_{i}, \pi_{-i}}(\hat{\theta}_{m})$ follow $\pi^{k}$ to collect a trajectory $\tau^{k} = (\mathbf{o}_{1}^{k}, \mathbf{a}_{1}^{k}, \dots, \mathbf{o}_{H}^{k}, \mathbf{a}_{H}^{k})$ add $(\pi_{i}^{k}, \tau_{i}^{k}, \tau_{m}^{k})$ into $\mathcal{D}_{m}$ for $m \in [n]/\{i\}$ and add $(\pi_{i}^{k}, \tau_{i}^{k})$ into $\mathcal{D}_{i}$ , and update $\mathcal{B}_{i}^{k+1}$ and $\mathcal{B}_{m}^{k+1}$ for $m \in [n]/\{i\}$ as follows: $\mathcal{B}_{i}^{k+1} = \{\hat{\theta}_{i} \in \mathcal{B}_{i}^{1} : \sum_{(\pi_{i}, \tau_{i}) \in \mathcal{D}_{i}} \log \mathbb{P}_{\hat{\theta}_{i}, i}^{\pi_{i}}(\tau_{i}) \geq \max_{\theta_{i}' \in \Theta_{i}} \sum_{(\pi_{i}, \tau_{i}) \in \mathcal{D}_{i}} \log \mathbb{P}_{\theta_{i}', i}^{\pi_{i}}(\tau_{i}) - \beta_{i}\}$
1: 2: 3: 4: 5: 6:	$\begin{aligned} \text{Initialize: } &\mathcal{B}_{i}^{1} = \{\hat{\theta}_{i} \in \Theta_{i} : \min_{h} \sigma_{S}(\mathbb{O}_{i}(\hat{\theta}_{i}) \geq \alpha)\}, \mathcal{B}_{m}^{1} = \{\hat{\theta}_{m} \in \Theta_{m} : \min_{h} \sigma_{S}(\mathbb{O}_{m}(\hat{\theta}_{m}) \geq \alpha/\sqrt{O})\} \text{ for all } m \in [n]/\{i\}. \text{ Set } \mathcal{D}_{m} = \{\}, \text{ for all agents } m \in [n]. \end{aligned}$ $\begin{aligned} &for \ k = 1 \dots T \ \mathbf{do} \\ &compute \ (\theta_{1}^{k}, \theta_{2}^{k}, \dots, \theta_{n}^{k}, \pi_{i}^{k}) = \arg\max_{\hat{\theta}_{1} \in \mathcal{B}_{1}^{k}, \hat{\theta}_{2} \in \mathcal{B}_{2}^{k}, \dots, \hat{\theta}_{n} \in \mathcal{B}_{n}^{k}, \pi_{i} \sum_{m=1}^{n} V_{m}^{\pi_{i}, \pi_{-i}}(\hat{\theta}_{m}) \\ &follow \ \pi^{k} \ \text{ to collect a trajectory } \mathbf{\tau}^{k} = (\mathbf{o}_{1}^{k}, \mathbf{a}_{1}^{k}, \dots, \mathbf{o}_{H}^{k}, \mathbf{a}_{H}^{k}) \\ &add \ (\pi_{i}^{k}, \tau_{i}^{k}, \tau_{m}^{k}) \ \text{ into } \mathcal{D}_{m} \ \text{ for } m \in [n]/\{i\} \ \text{ and } add \ (\pi_{i}^{k}, \tau_{i}^{k}) \ \text{ into } \mathcal{D}_{i}, \text{ and update } \mathcal{B}_{i}^{k+1} \ \text{ and } \mathcal{B}_{m}^{k+1} \ \text{ for } m \in [n]/\{i\} \ \text{ as follows:} \end{aligned}$ $\begin{aligned} &\mathcal{B}_{i}^{k+1} = \left\{\hat{\theta}_{i} \in \mathcal{B}_{1}^{1} : \sum_{(\pi_{i}, \tau_{i}) \in \mathcal{D}_{i}} \log \mathbb{P}_{\hat{\theta}_{i,i}}^{\pi_{i}}(\tau_{i}) \geq \max_{\theta_{i}' \in \Theta_{i}} \sum_{(\pi_{i}, \tau_{i}) \in \mathcal{D}_{i}} \log \mathbb{P}_{\hat{\theta}_{i,i}'}^{\pi_{i}, \pi_{m}}(\tau_{i}, \tau_{m}) \geq \max_{\theta_{m}' \in \Theta_{m}} \sum_{(\pi_{i}, \tau_{i}, \tau_{m}) \in \mathcal{D}_{m}} \log \mathbb{P}_{\theta_{m}', m}^{\pi_{i}, \pi_{m}}(\tau_{i}, \tau_{m}) \end{aligned}$

1890 1891 1892 1892 1893 1894  $\forall h \in [H]$ , the notation in the equations are defined as: for agent  $i, \pi_{h,i} = \pi_{h,i}(a_{h,i} \mid o_{h,i}),$  $\mathbb{T}_{h,i,a_{h,i}} = \mathbb{T}_{h,i,a_{h,i}}(s_{h+1} \mid s_h, o_{h,i}), \mathbb{O}_{h,i} = \mathbb{O}_{h,i}(o_{h,i} \mid s_h), \text{ and } \forall m \in [n] \setminus \{i\}, \pi_{h,m} = \pi_{h,m}(a_{h,m} \mid o_{h,m}), \mathbb{T}_{h,m,a_{h,m}} = \mathbb{T}_{h,m,a_{h,m},a_{h,i}}(s_{h+1} \mid s_h, o_{h,m}, o_{h,i}), \mathbb{O}_{h,m} = \mathbb{O}_{h,m}(o_{h,m}, o_{h,i} \mid s_h).$ 

For the real transition model  $\theta_i = \theta_i^*$ , the observation probability  $\{\mathbb{O}_{h,i}\}_{h \in [H]}$  and transition probability  $\{\mathbb{T}_{h,i}\}_{h \in [H]}$  are defined as  $(\forall s_h \in \mathcal{S}, s_{h+1} \in \mathcal{S}, o_{h,i} \in \mathcal{O}_i, a_{h,i} \in \mathcal{A}_i.)$ 

1897 
$$\mathbb{O}_{h,i}(o_{h,i} \mid s_h) = \sum_{\{o_{h,m}\}_{m \in [n]/\{i\}}} \mathbb{O}_h(\mathbf{o}_h \mid s_h)$$
1899

1900 1901 1902

1920 1921 1922

1925 1926

1931 1932

1934 1935

$$\mathbb{T}_{h,i,a_{h,i}}(s_{h+1} \mid s_h, o_{h,i}) = \sum_{\{(o_{h,m},a_{h,m})\}_{m \in [n]/\{i\}}} \left[ \prod_{j \in [n]/\{i\}} \pi_{h,j}(a_{h,j} \mid o_{h,j}) \right] \frac{\mathbb{O}_h(\mathbf{o}_h \mid s_h)}{\mathbb{O}_{h,i}(o_{h,i} \mid s_h)} \mathbb{T}_{h,\mathbf{a}_h}(s_{h+1} \mid s_h)$$
(19)

Similarly, for all  $m \in [n]/\{i\}$ , the real transition model  $\theta_m = \theta_m^*$ , the observation probability  $\{\mathbb{O}_{h,m}\}_{h\in[H]}$  and transition probability  $\{\mathbb{T}_{h,m}\}_{h\in[H]}$  are defined as  $(\forall s_h \in \mathcal{S}, s_{h+1} \in \mathcal{S}, o_{h,m} \in \mathcal{O}_m, a_{h,m} \in \mathcal{A}_m.)$ 

1907 
$$\mathbb{O}_{h,m}(o_{h,i}, o_{h,m} \mid s_h) = \sum_{\{o_{h,r}\}_{r \in [n]/\{i,m\}}} \mathbb{O}_h(\mathbf{o}_h \mid s_h)$$
  
1909

$$\mathbb{T}_{h,m,a_{h,i},a_{h,m}}(s_{h+1} \mid s_h, o_{h,i}, o_{h,m}) = \sum_{\{(o_{h,r},a_{h,r})\}_{m \in [n]/\{i,m\}}} \left[ \prod_{j \in [n]/\{i,m\}} \pi_{h,j}(a_{h,j} \mid o_{h,j}) \right] \frac{\mathbb{O}_h(\mathbf{o}_h \mid s_h)}{\mathbb{O}_{h,m}(o_{h,i}, o_{h,m} \mid s_h)} \mathbb{T}_{(20)}$$

г

٦

1913 where we denote  $\mathbb{T}_{h,\mathbf{a}_h}$  as  $\mathbb{T}_{h,\mathbf{a}_h}(s_{h+1} \mid s_h)$ . 1914

# 1915 E.2 PROOF FOR THEOREM 5.11916

1917 We first proof the following theorem, which can be seen as the bound of regret for Algorithm E.1. 1918 **Theorem E.1.** If Algorithm 5.1 take agent *i* as central agent and take policies  $\{\pi_m\}_{m \in [n]/\{i\}}$  as 1919 input, then Algorithm 5.1 guarantees that with probability at least  $1 - \delta$  for all  $k \in [T]$ 

$$\left[\sum_{t=1}^{k} V^{\pi_i^*, \boldsymbol{\pi}_{-i}} - V^{\pi_i^k, \boldsymbol{\pi}_{-i}}\right] \leq \tilde{\mathcal{O}}\left(\frac{S^2 O^2 A^2}{\alpha^2} \sqrt{k(S^2 A^2 O^2 + SO^2)} \times \operatorname{poly}(H)\right),$$

1923 where  $\pi_i^*$  is defined as

$$\pi_i^* = \arg\max_{\pi_i} V^{\pi_i, \boldsymbol{\pi}_{-i}}.$$

1927 According to the symmetric principle, we assume the central agent in Algorithm E.1 is agent n1928 without loss of generosity.

**1929 Definition E.1.** For all  $(m, i) \in [n-1] \times [T]$ ,  $\theta_m \in \Theta_m$ , and any policy  $\pi_m$  of agent m, we define **1930**  $f_m(\theta_m, \pi_n, \pi_m)$  and  $\tilde{f}_m^i(\theta_m, \pi_n, \pi_m)$  as:

$$f_m(\theta_m, \pi_n, \pi_m) = \mathbb{P}_{\theta_m, m}^{\pi_n, \pi_m}(\tau_n, \tau_m), \quad \tilde{f}_m^i(\theta_m, \pi_n, \pi_m) = \mathbb{P}_{\theta_m, m}^{\pi_n, \pi_m}\left(\tau_n^i, \tau_m^i\right).$$

1933 For agent n, we define  $f_n(\theta_n, \pi_n)$  and  $\tilde{f}_n^i(\theta_n, \pi_n)$  as:

$$f_n(\theta_n, \pi_n) = \mathbb{P}_{\theta_n, n}^{\pi_n}(\tau_n), \quad \tilde{f}_n^i(\theta_n, \pi_n) = \mathbb{P}_{\theta_n, n}^{\pi_n}\left(\tau_n^i\right).$$

**Lemma E.1.** There exist an absolute constant c such that for any  $\delta \in (0, 1]$ , with probability at least  $1 - \delta$ , the following inequality holds for all  $t \in [T]$  and all  $\theta_m \in \Theta_m$ ,  $m \in [n]$ .

$$\sum_{i=1}^{1938} \log\left(\tilde{f}_{m}^{i}\left(\theta_{m}, \pi_{n}^{i}, \pi_{m}\right) / \tilde{f}_{m}^{i}\left(\theta_{m}^{\star}, \pi_{n}^{i}, \pi_{m}\right)\right) \leq \beta_{m}, \quad \sum_{i=1}^{t} \log\left(\tilde{f}_{n}^{i}\left(\theta_{n}, \pi_{n}^{i}\right) / \tilde{f}_{n}^{i}\left(\theta_{n}^{\star}, \pi_{n}^{i}\right)\right) \leq \beta_{n},$$

$$1940$$

1941 where we define bonus term  $\beta_n$  and  $\beta_m$  for all  $m \in [n-1]$  as:

1942  
1942  
1943
$$\beta_m = c \left( H(S^2 A^2 O^2 + O^2 S) \log(TSAOH) + \log(Tn/\delta) \right),$$

$$\beta_n = c \left( H(S^2 A O + OS) \log(TSAOH) + \log(Tn/\delta) \right).$$

$$\log \left| \bar{\theta}_n \right| \le \mathcal{O} \left( H(S^2 A O + S O) \log(T S A O H) \right)$$

We denote 
$$\mathbb{E}_{t}[\cdot] = \mathbb{E} = [\cdot|\{(\pi_{n}^{i}, \boldsymbol{\tau}^{i})\}_{i=1}^{t-1} \cup \{\pi_{n}^{i}\}].$$
  

$$\mathbb{E}\left[\exp\left(\sum_{i=1}^{t}\log\left(\tilde{f}_{n}^{i}\left(\bar{\theta}_{n}, \pi_{n}^{i}\right)/\tilde{f}_{n}^{i}\left(\theta_{n}^{\star}, \pi_{n}^{i}\right)\right)\right)\right]$$

$$=\mathbb{E}\left[\exp\left(\sum_{i=1}^{t-1}\log\left(\tilde{f}_{n}^{i}\left(\bar{\theta}_{n}, \pi_{n}^{i}\right)/\tilde{f}_{n}^{i}\left(\theta_{n}^{\star}, \pi_{n}^{i}\right)\right)\right) \cdot \mathbb{E}_{t}\left[\exp\left(\log\left(\tilde{f}_{n}^{t}(\bar{\theta}_{n}, \pi_{n}^{t})/\tilde{f}_{n}^{t}(\theta_{n}^{\star}, \pi_{n}^{t})\right)\right)\right]\right]$$

$$=\mathbb{E}\left[\exp\left(\sum_{i=1}^{t-1}\log\left(\tilde{f}_{n}^{i}\left(\bar{\theta}_{n}, \pi_{n}^{i}\right)/\tilde{f}_{n}^{i}\left(\theta_{n}^{\star}, \pi_{n}^{i}\right)\right)\right) \cdot \mathbb{E}_{t}\left(\tilde{f}_{n}^{t}(\bar{\theta}_{n}, \pi_{n}^{t})/\tilde{f}_{n}^{t}(\theta_{n}^{\star}, \pi_{n}^{t})\right)\right]$$

$$(21)$$

1963 Since we have

$$\mathbb{E}_t\left(\tilde{f}_n^t(\bar{\theta}_n, \pi_n^t) / \tilde{f}_n^t(\theta_n^\star, \pi_n^t)\right) = \sum_{\tau_n} \tilde{f}_n^t\left(\bar{\theta}_n, \pi_n^t\right) \le (1 + 1/T).$$

we obtain that

$$\mathbb{E}\left[\exp\left(\sum_{i=1}^{t}\log\left(\tilde{f}_{n}^{i}\left(\bar{\theta}_{n},\pi_{n}^{i}\right)/\tilde{f}_{n}^{i}\left(\theta_{n}^{\star},\pi_{n}^{i}\right)\right)\right)\right] \leq e.$$

Therefore, by Markov's inequality, we have

$$\mathbb{P}\left(\sum_{i=1}^{t} \log\left(\frac{\tilde{f}_{n}^{i}\left(\bar{\theta}_{n}, \pi_{n}^{i}\right)}{\tilde{f}_{n}^{i}\left(\theta_{n}^{*}, \pi_{n}^{i}\right)}\right) > \log(1/\delta)\right) \leq \mathbb{E}\left[\exp\sum_{i=1}^{t} \log\left(\frac{\tilde{f}_{n}^{i}\left(\bar{\theta}_{n}, \pi_{n}^{i}\right)}{\tilde{f}_{n}^{i}\left(\theta_{n}^{*}, \pi_{n}^{i}\right)}\right)\right] \cdot \exp(-\log(1/\delta)) \leq e\delta$$
We take a union bound over all  $(\bar{\theta}_{n}, t) \in \bar{\Theta}$ ,  $\forall [T]$  and measuring  $\delta$ , we obtain

We take a union bound over all  $(\theta_n, t) \in \Theta_n \times [T]$  and rescaling  $\delta$ , we obtain

$$\begin{aligned} & \frac{1976}{1977} & \mathbb{P}\left(\max_{(\bar{\theta}_n, t)\in\bar{\Theta}_n\times[T]}\sum_{i=1}^t \log\left(\frac{\tilde{f}_n^i\left(\bar{\theta}_n, \pi_n^i\right)}{\tilde{f}_n^i\left(\theta_n^\star, \pi_n^i\right)}\right) > c(H(S^2AO + SO)\log(TSAOH) + \log(T/\delta))\right) \leq \delta. \end{aligned}$$

1979 Since  $\bar{\theta}_n$  is an optimistic discretization of  $\theta_n$ , which implies that  $\mathbb{P}_{\theta_n,n}^{\pi_n}(\tau_n) \leq \mathbb{P}_{\bar{\theta}_n,n}^{\pi_n}(\tau_n)$  for all  $\theta_n, \pi_n, \tau_n$ . As a result, we obtain that

$$\mathbb{P}\left(\max_{(\theta_n,t)\in\Theta_n\times[T]}\sum_{i=1}^t \log\left(\frac{\tilde{f}_n^i\left(\theta_n,\pi_n^i\right)}{\tilde{f}_n^i\left(\theta_n^\star,\pi_n^i\right)}\right) > c(H(S^2AO + SO)\log(TSAOH) + \log(T/\delta))\right) \le \delta.$$

We similarly consider  $f_m(\theta_m, \pi_n, \pi_m)$  for all  $m \in [n-1]$ . We use a  $d_m = H(S^2A^2O^2 + SO^2) + S$ dimension parameter  $\theta_m = (\mathbb{T}_m, \mathbb{O}_m, \mu)$  to denote the ensemble of all the parameters of the probability of trajectories  $(\tau_n^i, \tau_m^i)$ . We denote  $\Theta_m$  as the collection of all the  $\epsilon$ -optimistic discretelization,  $\bar{\theta}_m$ . We still choose  $\epsilon \leq 1/(c(S + O + A)HT)$  for large constant c so that  $\sum_{\tau_n, \tau_m} |f_m(\bar{\theta}_m, \pi_n, \pi_m) - f_m(\theta_m, \pi_n, \pi_m)| \leq 1/T$ . With similar analysis as above, we can derive the following inequality:

$$\mathbb{P}\left(\max_{(\theta_m,t)\in\Theta_m\times[T]}\sum_{i=1}^t \log\left(\tilde{f}^i_m\left(\theta_m,\pi^i_n,\pi_m\right)/\tilde{f}^i_m\left(\theta^\star_m,\pi^i_n,\pi_m\right)\right) > \beta_m\right) \le \delta$$

Eventually, we can obtain that with probability at least  $1 - \delta$ , the following events hold true:

$$\sum_{i=1}^{1995} \log \left( \tilde{f}_m^i \left( \theta_m, \pi_n^i, \pi_m \right) / \tilde{f}_m^i \left( \theta_m^\star, \pi_n^i, \pi_m \right) \right) \le \beta_m, \quad \sum_{i=1}^t \log \left( \tilde{f}_n^i \left( \theta_n, \pi_n^i \right) / \tilde{f}_n^i \left( \theta_n^\star, \pi_n^i \right) \right) \le \beta_n,$$

$$\square$$

**Lemma E.2.** There exists a universal constant c such that for any  $\delta \in (0, 1]$ , with probability at least  $1 - \delta$ , for all  $t \in [T]$  and all  $\theta_m \in \Theta_m$ ,  $m \in [n - 1]$ , it holds that 

$$\sum_{t=1}^{k} \left( \sum_{\tau_n \in (\mathcal{O}_n \times \mathcal{A}_n)^H, \tau_m \in (\mathcal{O}_m \times \mathcal{A}_m)^H} \left| f_m \left( \theta_m, \pi_n^i, \pi_m \right) - f_m \left( \theta_m^\star, \pi_n^i, \pi_m \right) \right| \right)^2$$
$$\leq c \left( \sum_{t=1}^{k} \log \left( \tilde{f}_m^i \left( \theta_m^\star, \pi_n^i, \pi_m \right) / \tilde{f}_m^i \left( \theta_m, \pi_n^i, \pi_m \right) \right) + \beta_m \right),$$

and for  $\theta_n \in \Theta_n$ , it holds that

 $\sum_{k=1}^{k}$ 

$$\sum_{t=1}^{k} \left( \sum_{\tau_n \in (\mathcal{O}_n \times \mathcal{A}_n)^H} \left| f_n\left(\theta_n, \pi_n^i\right) - f_n\left(\theta_n^\star, \pi_n^i\right) \right| \right)^2 \le c \left( \sum_{i=1}^{t} \log \left( \frac{\tilde{f}_n^i\left(\theta_n^\star, \pi_n^i\right)}{\tilde{f}_n^i\left(\theta_n, \pi_n^i\right)} \right) + \beta_n \right)$$

*Proof.* The proof of this lemma is very similar to the proof for Lemma C.2, so we omit it here for clarity.

**Lemma E.3.** We have the following bound on the regret of Algorithm 5.1. 

where we define  $\pi_i^*$  as  $\pi_i^* = \arg \max_{\pi_i} V^{\pi_i, \pi_{-i}}$ .

*Proof.* For any policy  $\pi = (\pi_1, \pi_2, \dots, \pi_n)$ , we can decompose the value function  $V^{\pi}$  as follows:

For  $m \in [n-1]$ , we further define

2038  
2039  
2040
$$V_{n}^{\pi_{n},\pi_{-n}}(\theta_{n}) = \sum_{\tau_{H,n}} \mathbb{P}_{\theta_{n},n}^{\pi_{n}}(\tau_{H,n}) \cdot \left(\sum_{h=1}^{H} r_{h}(o_{h,n})\right), \quad V_{m}^{\pi_{n},\pi_{-n}}(\theta_{m}) \sum_{\tau_{H,m},\tau_{H,n}} \mathbb{P}_{\theta_{m}^{*},m}^{\pi_{n},\pi_{m}}(\tau_{H,m},\tau_{H,n}) \cdot \left(\sum_{h=1}^{H} r_{h}(o_{h,m})\right).$$
2041

Then we can decompose the value function as 

$$V^{\pi_n^*, \pi_{-n}} = \sum_{m=1}^n V_m^{\pi_n^t, \pi_{-n}}(\theta_m^*), \quad V^{\pi_n^k, \pi_{-n}} = \sum_{m=1}^n V_m^{\pi_n^*, \pi_{-n}}(\theta_m^*).$$

According to Lemma E.1 and Lemma E.2, we can deduce that  $\theta_m^* \in \bigcap_{t \in [K]} \mathcal{B}_m^t$ , holds for all  $m \in [n]$ with probability at least  $1 - \delta$ . In the following analysis, we assume that  $\theta_m^* \in \bigcap_{t \in [K]} \mathcal{B}_m^t$ . On this event, according to the optimism of  $\{\theta_m\}_{m \in [n]}$  and  $\pi_n^t$  for  $t \in [K]$ , we can obtain that 

2050  
2051 
$$V^{\pi_n^*, \pi_{-n}} - V^{\pi_n^k, \pi_{-n}} = \sum_{m=1}^n V_m^{\pi_n^*, \pi_{-n}}(\theta_m^*) - \sum_{m=1}^n V_m^{\pi_n^t, \pi_{-n}}(\theta_m^*) \le \sum_{m=1}^n V_m^{\pi_n^t, \pi_{-n}}(\theta_m^t) - \sum_{m=1}^n V_m^{\pi_n^t, \pi_{-n}}(\theta_m^*).$$

L

According to the definition of 
$$\left\{V_m^{\pi_n^t, \pi_{-n}}(\theta_m^*)\right\}_{m \in [n]}$$
 and  $\left\{V_m^{\pi_n^t, \pi_{-n}}(\theta_m^t)\right\}_{m \in [n]}$ , we further have the following bound on the regret.

$$\sum_{t=1}^{K} \left[ V^{\pi_{n}^{*}, \pi_{-n}} - V^{\pi_{n}^{k}, \pi_{-n}} \right]$$

$$\leq \sum_{t=1}^{K} \left[ \sum_{m=1}^{n} V_{m}^{\pi_{n}^{t}, \pi_{-n}}(\theta_{m}^{t}) - \sum_{m=1}^{n} V_{m}^{\pi_{n}^{t}, \pi_{-n}}(\theta_{m}^{*}) \right]$$

$$\leq H \cdot \sum_{t=1}^{K} \sum_{\tau_{H,n}} \left| \mathbb{P}_{\theta_{n}^{t}}^{\pi_{n}^{t}}(\tau_{H,n}) - \mathbb{P}_{\theta_{n}^{*}}^{\pi_{n}^{t}}(\tau_{H,n}) \right| + H \cdot \sum_{m=1}^{n-1} \sum_{t=1}^{K} \sum_{\tau_{H,n}, \tau_{H,m}} \left| \mathbb{P}_{\theta_{m}^{t}}^{\pi_{n}^{t}, \pi_{m}}(\tau_{H,n}, \tau_{H,m}) - \mathbb{P}_{\theta_{m}^{*}}^{\pi_{n}^{t}, \pi_{m}}(\tau_{H,n}, \tau_{H,m}) \right|.$$

**Definition E.2.** For all  $(m, h, k) \in [n-1] \times [H] \times [T]$ , we define the matrix notations  $\mathbb{M}_{h,m}^{\star} \in$  $\mathbb{R}^{O^2 \times O^2}$  and  $\mathbb{M}_{h_m}^k \in \mathbb{R}^{O^2 \times O^2}$  as follows:

$$\begin{split} \mathbb{M}_{0,m}^{\star} &= \mathbb{O}_{1,m}^{\star} \mu_{m}^{\star} \in \mathbb{R}^{O}, \quad \mathbb{M}_{0,m}^{k} = \mathbb{O}_{1,m}^{k} \mu_{m}^{k} \in \mathbb{R}^{O}, \\ \mathbb{M}_{h,m}^{\star}(o_{m}, a_{m}, a_{n}) &= \mathbb{O}_{h+1,m}^{\star} \mathbb{T}_{h,m,a_{m},a_{n}}^{\star} \cdot diag \left( \mathbb{O}_{h,m}^{\star}(o_{m}, o_{n} \mid \cdot) \right) \left( \mathbb{O}_{h,m}^{\star} \right)^{\dagger} \in \mathbb{R}^{O^{2} \times O^{2}}, \\ \mathbb{M}_{h,m}^{k}(o_{m}, a_{m}, a_{n}) &= \mathbb{O}_{h+1,m}^{k} \mathbb{T}_{h,m,a_{m},a_{n}}^{k} \cdot diag \left( \mathbb{O}_{h,m}^{k}(o_{m}, o_{n} \mid \cdot) \right) \left( \mathbb{O}_{h,m}^{k} \right)^{\dagger} \in \mathbb{R}^{O^{2} \times O^{2}}, \end{split}$$

and for agent n, we define matrix notations  $\mathbb{M}_{h,n}^{\star} \in \mathbb{R}^{O \times O}$  and  $\mathbb{M}_{h,n}^{k} \in \mathbb{R}^{O \times O}$  as follows for all  $(h,k) \in [H] \times [T]$ 

$$\begin{split} \mathbb{M}_{0,n}^{\star} &= \mathbb{O}_{1,n}^{\star} \mu_{n}^{\star} \in \mathbb{R}^{O}, \quad \mathbb{M}_{0,m}^{k} = \mathbb{O}_{1,n}^{k} \mu_{n}^{k} \in \mathbb{R}^{O}, \\ \mathbb{M}_{h,n}^{\star}(o_{n}, a_{n}) &= \mathbb{O}_{h+1,n}^{\star} \mathbb{T}_{h,n,a_{n}}^{\star} \cdot diag \left(\mathbb{O}_{h,n}^{\star}(o_{n} \mid \cdot)\right) \left(\mathbb{O}_{h,n}^{\star}\right)^{\dagger} \in \mathbb{R}^{O \times O}, \\ \mathbb{M}_{h,n}^{k}(o_{n}, a_{n}) &= \mathbb{O}_{h+1,n}^{k} \mathbb{T}_{h,n,a_{n}}^{k} \cdot diag \left(\mathbb{O}_{h,n}^{k}(o_{n} \mid \cdot)\right) \left(\mathbb{O}_{h,n}^{k}\right)^{\dagger} \in \mathbb{R}^{O \times O}, \end{split}$$

where  $\{\mathbb{O}_{h,m}^{\star}\}_{(h,m)\in[H]\times[n]}$  and  $\{\mathbb{T}_{h,m}^{\star}\}_{(h,m)\in[H]\times[n]}$  denote the observation and transition matri-ces corresponding to the true transition model, and  $\{\mathbb{O}_{h,m}^k\}_{(h,m)\in[H]\times[n]}$  and  $\{\mathbb{T}_{h,m}^k\}_{(h,m)\in[H]\times[n]}$ denote the observation and transition matrices corresponding to model parameter  $\theta_m^k$  for all  $k \in [T]$ . When no confusion arises, we simplify the notation by using  $\mathbb{M}_{h,m}^{\star}$  to represent  $\mathbb{M}_{h,m}^{\star}(o_m, a_m, a_n)$  and  $\mathbb{M}_{h,m}^k$  to represent  $\mathbb{M}_{h,m}^k(o_m, a_m, a_n)$  for  $m \in [n-1]$ . We also simplify by using  $\mathbb{M}_{h,n}^{\star}$  to represent  $\mathbb{M}_{h,n}^{\star}(o_n, a_n)$  and  $\mathbb{M}_{h,n}^k$  to represent  $\mathbb{M}_{h,n}^k(o_n, a_n)$ . 

**Lemma E.4.** Given  $O \times S$  matrix  $\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_n$ . We further define matrix  $\mathbf{B}$  and matrix  $\mathbf{C}$ as  $\mathbf{B}^{\intercal} = (\mathbf{A}_1^{\intercal}, \mathbf{A}_2^{\intercal}, \dots, \mathbf{A}_n^{\intercal}), \mathbf{C} = \mathbf{A}_1 + \mathbf{A}_2 + \dots + \mathbf{A}_n$ . Then if  $\sigma_S(\mathbf{C}) \geq \alpha$ , then we have  $\sigma_S(\mathbf{C}) \geq \alpha/\sqrt{n}.$ 

*Proof.* We only need to prove that for any given unit vector  $\mathbf{x} \in \mathbb{R}^S$ ,  $\|\mathbf{B}\mathbf{x}\|_2 \geq \frac{\alpha}{\sqrt{O}}$ . Since we have

2093  
2094 
$$\|\mathbf{B}\mathbf{x}\|_{2} = \|(\mathbf{A}_{1}\mathbf{x}^{\mathsf{T}}, \dots, \mathbf{A}_{n}\mathbf{x}^{\mathsf{T}})^{\mathsf{T}}\|_{2} = \sqrt{\|\mathbf{A}_{1}\mathbf{x}\|_{2}^{2} + \dots + \|\mathbf{A}_{n}\mathbf{x}\|_{2}^{2}} \ge \frac{1}{\sqrt{n}}\|\mathbf{A}_{1}\mathbf{x} + \dots + \mathbf{A}_{n}\mathbf{x}\|_{2} = \frac{1}{\sqrt{n}}\|\mathbf{C}\mathbf{x}\|_{2} \ge \frac{\alpha}{\sqrt{n}}$$
2095  
2096 Thus, we finished the proof of the lemma.

Thus, we finished the proof of the lemma. 

**Corollary E.1.** According to the definition of observable condition, we have  $\sigma_S(\mathbb{O}_{h,n}) \geq \alpha$ , for all  $h \in [H]$ , and  $\sigma_S(\mathbb{O}_{h,m}) \ge \alpha/\sqrt{O}$  for all  $h \in [H]$  and all  $m \in [n-1]$ . 

According to the definition of matrix notation, we can directly obtain the following lemma:

**Lemma E.5.** (Bound the regret of Operator Estimates) The following two inequalities holds true for all  $h \in [H]$ : 

$$\sum_{\tau_{h,n}} \left\| \begin{bmatrix} h \\ j=0 \end{bmatrix} \mathbb{I}_{j=0}^{h} \mathbb{M}_{j,n}^{k} \right\| - \begin{bmatrix} h \\ j=0 \end{bmatrix} \|_{1}^{\star} \pi_{n}^{t}(\tau_{h,n}) \leq \frac{\sqrt{S}}{\alpha} \left( \sum_{j=0}^{h} \sum_{\tau_{j,n}} \left\| (\mathbb{M}_{j,n}^{k} - \mathbb{M}_{j,n}^{\star}) \begin{bmatrix} j-1 \\ \prod_{h'=0}^{\star} \mathbb{M}_{h',n}^{\star} \right\|_{1}^{\star} \pi_{n}^{t}(\tau_{j,n}) \right),$$

and for all agent  $m \in [n-1]$  we further have 

$$\sum_{\tau_{h,n},\tau_{h,m}} \left\| \left[ \prod_{j=0}^{h} \mathbb{M}_{j,m}^{k} \right] \mathbf{b}_{0,m}^{k} - \left[ \prod_{j=0}^{h} \mathbb{M}_{j,m}^{\star} \right] \right\|_{1} \pi_{n}^{t}(\tau_{h,n}) \pi_{m}(\tau_{h,m})$$

2111  
2112  
2113  
2114
$$\leq \frac{\sqrt{S}}{\alpha} \sum_{j=0}^{h} \sum_{\tau_{j,n},\tau_{j,m}} \left\| \left( \mathbb{M}_{j,m}^{k} - \mathbb{M}_{j,m}^{\star} \right) \left[ \prod_{h'=0}^{j-1} \mathbb{M}_{h',m}^{\star} \right] \right\|_{1} \pi_{n}^{t}(\tau_{h,n}) \pi_{m}(\tau_{h,m})$$
2114

**Lemma E.6.** (Constraints for the Operator Estimates from OMLE) With probability at least  $1 - \delta$ , for all  $m \in [n-1]$ , the following events hold true. 

1 -

\

$$= \sum_{t=1}^{k-1} \sum_{\tau_{h,n}} \pi_n^t(\tau_{h,n}) \cdot \left\| \left( \mathbb{M}_{h,n}^k - \mathbb{M}_{h,n}^\star \right) \left[ \prod_{h'=0}^{h-1} \mathbb{M}_{h',n}^\star \right] \right\|_1 = \mathcal{O}\left( \frac{\sqrt{S}}{\alpha} \sqrt{k\beta_n} \right),$$
$$= \sum_{t=1}^{k-1} \sum_{\tau_{h,n},\tau_{h,m}} \pi_n^t(\tau_{h,n}) \pi_m(\tau_{h,m}) \cdot \left\| \left( \mathbb{M}_{h,m}^k - \mathbb{M}_{h,m}^\star \right) \left[ \prod_{h'=0}^{h-1} \mathbb{M}_{h',m}^\star \right] \right\|_1 = \mathcal{O}\left( \frac{\sqrt{S}}{\alpha} \sqrt{k\beta_m} \right),$$

where  $\beta_m$  and  $\beta_n$  is defined as

$$\beta_n = c(H(S^2A^2O^2 + SO^2)\log(TSAOH) + \log(nT/\delta)),$$
  
$$\beta_n = c(H(S^2AO + SO)\log(TSAOH) + \log(Tn/\delta)).$$

Proof. The proof of this lemma is very similar to the proof for Lemma C.3, so we omit it here for clarity. 

**Proof of Theorem E.1** : We only need to consider the following problem: We are required to bound the following target term for  $m \in [n-1]$ .

$$\sum_{t=1}^{k} \left( \sum_{h=0}^{H-1} \sum_{\tau_{h,n}} \left\| \left( \mathbb{M}_{h,n}^{t} - \mathbb{M}_{h,n}^{\star} \right) \left[ \prod_{h'=0}^{h-1} \mathbb{M}_{h',n}^{\star} \right] \right\|_{1} \cdot \pi_{n}^{t}(\tau_{h,n}) \right)$$
$$\sum_{t=1}^{k} \left( \sum_{h=0}^{H-1} \sum_{\tau_{h,n},\tau_{h,m}} \left\| \left( \mathbb{M}_{h,m}^{t} - \mathbb{M}_{h,m}^{\star} \right) \left[ \prod_{h'=0}^{h-1} \mathbb{M}_{h',m}^{\star} \right] \right\|_{1} \cdot \pi_{n}^{t}(\tau_{h,n}) \pi_{m}(\tau_{h,m}) \right),$$

For agent n we have the following condition:

$$\sum_{t=1}^{k-1} \sum_{\tau_{h,n}} \left\| \left( \mathbb{M}_{h,n}^k - \mathbf{M}_{h,n}^\star \right) \left[ \prod_{h'=0}^{h-1} \mathbb{M}_{h',n}^\star \right] \right\|_1 \times \pi_n^t(\tau_{h,n}) = \mathcal{O}\left( \frac{\sqrt{S}}{\alpha} \sqrt{k\beta_n} \right), \forall k \in [T].$$

Corresponding to agent  $m \in [n-1]$  we have the following condition:

$$\sum_{t=1}^{k-1} \sum_{\tau_{h,n},\tau_{h,m}} \left\| \left( \mathbb{M}_{h,m}^{k} - \mathbb{M}_{h,m}^{\star} \right) \left[ \prod_{h'=0}^{h-1} \mathbb{M}_{h',m}^{\star} \right] \right\|_{1} \cdot \pi_{n}^{t}(\tau_{h,n}) \pi_{m}(\tau_{h,m}) = \mathcal{O}\left( \frac{\sqrt{S}}{\alpha} \sqrt{k\beta_{m}} \right).$$

We apply Eluder-Dimension Lemma (Proposition 22 of Liu et al. (2022a)), and we can obtain the following bound on the regret with probability at least  $1 - \delta$  for all  $m \in [n - 1]$ .

2153  
2154  
2155  
2155  
2156  
2156  
2157  
2156  
2157  
2158  
2158  
2158  
2159  

$$k = \left( \sum_{h=0}^{H-1} \sum_{\tau_{h,n}, \tau_{h,m}} \left\| \left( \mathbb{M}_{h,m}^{t} - \mathbb{M}_{h,n}^{\star} \right) \left[ \prod_{h'=0}^{h-1} \mathbb{M}_{h',m}^{\star} \right] \right\|_{1} \pi_{n}^{t}(\tau_{h,n}) \pi_{m}(\tau_{h,m}) \right) = \tilde{\mathcal{O}} \left( \frac{S^{1.5}O^{2}A^{2}H^{3}}{\alpha} \sqrt{k\beta_{m}} \right)$$
2159  
2159

Therefore, we achieve the bound of the regret.

**Proof for Theorem 5.1** We are now ready to prove Theorem 5.1. We begin by restating the theorem here for the reader's convenience.

**Theorem E.2.** If the central agent for Algorithm 5.1 is i, we define bonus parameter as  $\beta_i =$  $c(H(S^2AO + SO)\log(TSAOH) + \log(Tn/\delta)), \beta_m = c(H(S^2A^2O^2 + SO^2)\log(TSAOH) + \log(TSAOH)), \beta_m = c(H(S^2A^2O^2 + SO^2)\log(TSAOH) + \log(TSAOH)))$  $\log(Tn/\delta)$  ( $\forall m \in [n]/\{i\}$ ) for some constant c. Then, with probability at least  $1 - \delta$ , Algorithm terminates within  $4H/\epsilon$  steps of the while loop, and outputs an  $\epsilon$ -approximate local optimal policy. The total episodes of play is at most

$$K = \tilde{\mathcal{O}} \left( S^4 O^4 A^4 (S^2 A^2 O^2 + S O^2) \times poly(H) / (\alpha^4 \epsilon^3) \right)$$

*Proof.* We use superscript t to represent variables at the  $t^{th}$  step (before  $\pi$  is updated) of the while loop. We set  $K_i = \tilde{\mathcal{O}}(S^4 O^4 A^4 (S^2 A^2 O^2 + SO^2) \times \text{poly}(H) / (\alpha^4 \epsilon^2))$ . According to Theorem E.1, for fixed i and t, we have with probability at least  $1 - \frac{\delta \epsilon}{8\pi H}$ 

$$\max_{\mu_{i}} V(\mu_{i}, \pi_{-i}^{t}) - V(\hat{\pi}_{i}^{t}, \pi_{-i}^{t}) \le \frac{\epsilon}{4}$$

We then take a union bound over all  $t \le 4H/\epsilon$  and for all  $i \in [n]$ , and we have with probability at least  $1 - \delta$  the following inequality holds for all  $i \in [n]$  and  $t \leq 4H/\epsilon$ 

$$\max_{\mu_i} V(\mu_i, \pi_{-i}^t) - V(\hat{\pi}_i^t, \pi_{-i}^t) \le \frac{\epsilon}{4}.$$
(22)

For the empirical estimator  $\hat{V}^t$ , it's bounded in [0, H]. Thus, by Hoeffding's inequality, for fixed  $i \in [n]$  and t 

$$\mathbb{P}\left(\left|\hat{V}^t - V^t\right| \ge \frac{\epsilon}{8}\right) \le 2\exp\left(-\frac{N\epsilon^2}{32H^2}\right).$$

Choosing  $N = \frac{CH^2}{\epsilon^2} \log \left( \frac{nHSK \max_{i \in [n]} A_i}{\epsilon \delta} \right)$  for some large constant C, we have 

$$\mathbb{P}\left(\left|\hat{V}^t - V^t\right| \ge \frac{\epsilon}{8}\right) \le \frac{\epsilon\delta}{16nH}.$$

Applying this inequality for  $\hat{V}^t(\hat{\pi}_i^t, \pi_{-i}^t)$  and  $\hat{V}^t(\pi^t)$  and taking a union bound over  $i \in [n]$  and  $t \leq 4H/\epsilon$ , we can achieve that 

$$|\hat{V}(\hat{\pi}_{i}^{t}, \pi_{-i}^{t}) - V(\hat{\pi}_{i}^{t}, \pi_{-i}^{t})| \le \epsilon/8, \quad |\hat{V}^{t}(\pi^{t}) - V(\pi^{t})| \le \epsilon/8$$

We combine this result with equation 22, and we can obtain that with probability at least  $1 - \delta$ 

$$\max_{\mu_i} V(\mu_i, \pi_{-i}^t) - V(\hat{\pi}_i^t, \pi_{-i}^t) \le \frac{\epsilon}{4}, \qquad |\hat{V}(\hat{\pi}_i^t, \pi_{-i}^t) - V(\hat{\pi}_i^t, \pi_{-i}^t)| \le \epsilon/8, \quad |\hat{V}^t(\pi^t) - V(\pi^t)| \le \epsilon/8$$

holds for all  $i \in [n]$  and  $t \leq \frac{4H}{\epsilon}$ . On this event, 

$$\delta_i^t = \hat{V}^t(\hat{\pi}_i, \pi_{-i}) - \hat{V}^t(\pi^t) \le V(\hat{\pi}_i^t, \pi_{-i}^t) - V(\pi^t) + \epsilon/4.$$

If the while loop doesn't end after the t-th iteration and  $t \leq 4H/\epsilon$ , there exists  $j^t$  s.t.  $\Delta_{j^t}^t \geq \epsilon/2$ , so we have 

 $V(\hat{\pi}_{jt}^t, \pi_{-jt}^t) - V(\pi^t) \ge \Delta_{jt}^t - \epsilon/4 \ge \epsilon.$ 

Since the value function is bound by H, so the while loop ends within 4H/epsilon steps. Therefore, the inequality above that holds for all  $i \in [n]$  and  $t \leq 4H/epsilon$  holds for simultaneously before the end of the while loop with probability at least  $1 - \delta$ . Again, on this event, if the while loop stops at the end of  $t^{th}$  step, we have  $\max_{i \in [n]} \delta_i^t \leq \epsilon/2$ , then 

$$\max_{\mu_i} V(\mu_i, \pi_{-i}^t) - V(\pi^t) = \max_{\mu_i} V(\mu_i, \pi_{-i}^t) - V(\hat{\pi}_i^t, \pi_{-i}^t) + V(\hat{\pi}_i^t, \pi_{-i}^t) - V(\pi^t)$$

2211 
$$\leq \epsilon/4 + \hat{V}^t(\hat{\pi}_i^t, \pi_{-i}^t) - \hat{V}^t(\pi^t) + 2\epsilon/8$$

$$\frac{2212}{\leq \epsilon/2 + \Delta_i^t}$$

$$\frac{2213}{\leq \epsilon}.$$

So the returned policy  $\pi^t$  is an  $\epsilon$ -approximate local optimal. Therefore, we can conclude that probability at least  $1 - \delta$ , within  $4H/\epsilon$  steps of the while loop, Algorithm E.1 outputs an  $\epsilon$ -approximate local optimal policy.

Eventually, we compute the total number of episodes as the total sample complexity. We first compute the number of episodes within each step of the while loop.

$$N + \sum_{i=1}^{n} (K_i + N) = \mathcal{O}\left(\frac{nH^2}{\epsilon^2} \log\left(\frac{nHSK\max_{i\in[n]}A_i}{\epsilon\delta}\right)\right) + \tilde{\mathcal{O}}\left(\frac{S^4O^4A^4(S^2A^2O^2 + SO^2) \times \operatorname{poly}(H)}{\alpha^4\epsilon^2}\right).$$

Since the algorithm ends within at most  $4H/\epsilon$ , we can compute the total sample complexity.

$$K = \frac{4H}{\epsilon} \left( N + \sum_{i=1}^{n} (K_i + N) \right) = \tilde{\mathcal{O}} \left( \frac{S^4 O^4 A^4 (S^2 A^2 O^2 + S O^2) \times \operatorname{poly}(H)}{\alpha^4 \epsilon^3} \right).$$

2227 2228 2229

2230

2235

### E.3 PROOF FOR THEOREM 5.2

**Theorem E.3.** For both randomized and deterministic algorithms, there exists an instance of DEC-MDP wherein the regret scales at least as  $\mathcal{O}(\sqrt{A^nT})$ . This result underscores the limitation of achieving sample efficiency in algorithms for DEC-POMDP without imposing assumptions on the transition model, even with a memoryless policy.

*Proof.* The proof for Theorem 5.2 proceeds straightforwardly. We consider a two-step DEC-MDP, 2236 commencing from an initial state  $s_1$ . For all  $s \in S$ , we assume the reward function satisfies 2237  $r_{h,1}(s) = r_{h,2}(s) = \cdots = r_{h,n}(s)$  for all  $h \in [2]$ . Consequently, the entire DEC-MDP reduces to a multi-armed bandit problem. By leveraging a classic result on the lower bound of regret for 2239 the multi-armed bandit problem (Mannor and Tsitsiklis, 2004), it follows that for any randomized 2240 or deterministic algorithm, there exists an instance of the multi-arm bandit problem such that the 2241 regret is at least  $\mathcal{O}(\sqrt{AT})$ , where A denotes the number of arms. Consequently, for any random-2242 ized or deterministic algorithm, there exists an instance of DEC-MDP such that the regret is at least 2243  $\mathcal{O}(\sqrt{A^nT}).$ 

2244 2245 2246

2247

# F POMDP WITH KNAPSACK CONSTRAINTS

# 2248 F.1 MODEL

2249 We commence by formally defining the model. We consider the framework of tabular Partially Observable Markov Decision Processes (POMDPs), denoted as  $(\mathcal{S}, \mathcal{A}, \mathcal{O}, H, \mathbb{T}, r, \mathcal{M})$ , which extends 2251 to an episodic POMDP with a d-dimensional budget. Each component  $\mathbf{M}_i$  of the budget vector  $\mathbf{M}$ 2252 represents the total budget of the  $i^{th}$  cost. At the onset of each episode, the agent is endowed with 2253 a budget  $\mathbf{M}_1 = \mathbf{M} = (M, M, \dots, M)$ . During the  $h^{th}$  step, the agent incurs a cost vector, thereby decrementing the total budget to  $\mathbf{M}_{h+1} = \mathbf{M}_h - \mathbf{C}_h$ . Subsequently, the budget for the  $(h+1)^{th}$ 2254 step follows a transition probability  $\mathbf{M}_{h+1,i} \sim \mathbb{T}_h(\cdot | \mathbf{M}_{h,i}, o_h, a_h)$ . An episode concludes after H2255 steps or when the budget of any dimension i reaches 0. The primary objective of the agent is to 2256 maximize its cumulative reward  $\sum_{k=1}^{K} \sum_{h=1}^{H} r_{k,h}$  over K episodes. Furthermore, we impose the 2257 knapsack assumption on the cost: 2258

**Assumption F.1.** Both the budget  $\mathbf{M}_i$  and the possible values of costs  $\mathbf{C}_i$  are integral multiples of the unit cost  $\frac{1}{m}$ .

2262 We conceptualize the POMDP model as a factored Decentralized POMDP (DEC-POMDP). Initially, 2263 the policy class is defined as follows:

2261

$$\Pi = \left\{ \{\pi_h\}_{h=1}^H \mid \pi_h : (\mathcal{A} \times \mathcal{O} \times \mathcal{M}^d)^{H-1} \times \mathcal{O} \times \mathcal{M}^d \to \mathcal{A} \right\}$$

2266 We define the joint state space as  $\tilde{S} = S \times M$ , where the tuple consists of the true state and the 2267 budget. We introduce d dummy agents, where the local state of the *i*<sup>th</sup> dummy agent corresponds to the *i*<sup>th</sup> entry in the budget vector. The true agent is denoted as the  $(d+1)^{\text{th}}$  agent. The state transition 2268 of the  $(d+1)^{\text{th}}$  agent follows  $s_{h+1} \sim \mathbb{T}_{h,i}(\cdot \mid s_h, a_h)$ , and its observation follows  $o_h \sim \mathbb{O}_h(\cdot \mid s_h)$ . 2269 The transition of the *i*<sup>th</sup> agent is given by:  $\mathbf{M}_{h+1,i} \sim \mathbb{T}_h(\cdot \mid \mathbf{M}_{h,i}, o_h, a_h)$ . Therefore, the transition 2270 of the joint state is: 2271

$$(s_{h+1}, \mathbb{M}_{h+1}) \sim \left[\prod_{m=1}^{d} \mathbb{T}_{h,m}(\cdot \mid M_{h,m}, o_h, a_h)\right] \mathbb{T}_h(\cdot \mid s_h, a_h).$$

2274 The probability of trajectory  $(o_1, a_1, \mathbf{M}_1, o_2, a_2, \mathbf{M}_2, \dots, o_H, a_H, \mathbf{M}_H)$  for policy  $\pi$  is defined as: 2275

$$\mathbb{P}^{2276} \qquad \mathbb{P}^{\pi}(\tau_{H}) = \sum_{s_{1}, s_{2}, \dots, s_{H}} \mu_{1}(s_{1}) \mathbb{O}_{1}(o_{1} \mid s_{1}) \pi_{1}(a_{1} \mid o_{1}, \mathbf{M}_{1}) \mathbb{T}_{1, a_{1}}(s_{2} \mid s_{1}) \prod_{i=1}^{d} \mathbb{T}_{1, i}(\mathbf{M}_{2, i} \mid \mathbf{M}_{1, i}, o_{1}, a_{1}) \times \cdots$$

$$\times \mathbb{O}(o_{H-1} \mid s_{H-1}) \pi_{H-1}(a_{H-1} \mid o_{1}, a_{1}, \mathbf{M}_{1}, \dots, o_{H-2}, a_{H-2}, \mathbf{M}_{H-2}, o_{H-1}) \mathbb{T}_{H-1, a_{H-1}}(s_{H} \mid s_{H-1})$$

$$\times \prod_{i=1}^{d} \mathbb{T}_{H-1, i}(\mathbf{M}_{H, i} \mid \mathbf{M}_{H-1, i}, o_{H-1}, a_{H-1}) \times \mathbb{O}(o_{H} \mid s_{H}) \times \pi_{H-1}(a_{H} \mid o_{1}, a_{1}, \mathbf{M}_{1}, \dots, a_{H-1}, \mathbf{M}_{H-1}, o_{H}).$$

The reward function is defined as (for all  $h \in [H]$ ):

$$\tilde{r}_h(o_h, \mathbf{M}_h) = r_h(o_h)$$
 if  $M_{h,i} > 0$  for all  $i \in [d]$ .

2286 Hence, we can model the POMDP setting with knapsack constraints as a DEC-POMDP with d +2287 1 agents. For the  $(d+1)^{th}$  agent, o, a, and s are defined as its observation, action, and state, 2288 respectively. We interpret  $\mathbf{B}_i$  as both the individual state and the observation of the  $i^{th}$  agent. The 2289 transition of the *i*<sup>th</sup> agent is defined as:  $\mathbb{T}_{h,i}(\mathbf{M}_{h+1,i}|\mathbf{M}_{h,i}, o_h, a_h)$ , 2290

which is influenced by the observation and action of the  $(d+1)^{th}$  agent. 2291 The actions of agents  $1, 2, \ldots, d$  do not affect the model, hence we do 2292 not need to consider their actions. The reward function  $\tilde{r}_h(o_h, \mathbf{M}_h)$  is 2293 a function of the observation of all individuals. Thus, the POMDPwk 2294 model can be formulated as a factored DEC-POMDP. The influence 2295 graph of POMDPwk is depicted in Figure 2. The maximum indegree 2296 is 1. 2297



Figure 2: Influential graph for POMDP with constraints

#### F.2 ALGORITHM

2300 We can apply Algorithm 4.2 to the the setting of POMDP with con-2301 straints. We introduce Algorithm F.2. We define the complete trajectory  $\tau_H$  and trajectory  $\tau_{H,i}$ 2302

Algorithm 8 OMLE for POMDPwk 1: Initialize:  $\mathcal{B}_{d+1}^1 = \{\hat{\theta}_{d+1} \in \Theta_{d+1} : \min_h \sigma_S(\mathbb{O}_i(\hat{\theta}_{d+1}) \ge \alpha)\}, \mathcal{B}_i = \{\Theta_i\}, \mathcal{D}_i = \{\}, \forall i \in [d].$ 2: for  $k = 1, \dots, K$  do compute  $(\theta_1^k, \theta_2^k, \dots, \theta_n^k, \pi^k) = \arg \max_{\hat{\theta}_1 \in \mathcal{B}_1^k, \hat{\theta}_2 \in \mathcal{B}_2^k, \dots, \hat{\theta}_n \in \mathcal{B}_n^k, \pi} V^{\pi}(\hat{\theta})$ 3: follow  $\pi^k$  to collect a trajectory  $\boldsymbol{\tau}^k = (o_1^k, a_1^k, \mathbf{M}_1^k, \dots, o_H^k, a_H^k, \mathbf{M}_H^k)$ 4: add  $(\pi^k, \boldsymbol{\tau}^k)$  into  $\mathcal{D}_{d+1}$ , and then update 5:  $\mathcal{B}_{d+1}^{k+1} = \left\{ \hat{\theta}_{d+1} \in \mathcal{B}_{d+1}^1 : \sum_{(\pi, \tau) \in \mathcal{D}_{d+1}} \log \mathbb{P}_{\hat{\theta}_{d+1}, d+1}^{\pi} \ge \max_{\theta'_{d+1}} \sum_{(\pi, \tau) \in \mathcal{D}_{d+1}} \log \mathbb{P}_{\theta'_{d+1}, d+1}^{\pi} - \beta_{d+1} \right\}$ for i = 1, ..., d do 6: follow  $\pi^k$  and model  $\left[\prod_{j=1}^i \mathbb{P}_{\theta_j^*, j}\right] \left[\prod_{j=i+1}^d \mathbb{P}_{\theta_j^k}\right] \mathbb{P}_{\theta_{d+1}^k}^{\pi^k}$  to collect a trajectory  $\boldsymbol{\tau}^k$ . add  $(\tau_i^k, \tau_{d+1}^k)$  into  $\mathcal{D}_i$  and then update 7: 8:  $\mathcal{B}_{i}^{k+1} = \left\{ \hat{\theta} \in \mathcal{B}_{i}^{1} : \sum_{(\tau_{i}, \tau_{d+1}) \in \mathcal{D}_{i}} \log \mathbb{P}_{\hat{\theta}_{i}, i}\left(\tau_{i} | \tau_{d+1}\right) \geq \max_{\theta_{i}' \in \Theta_{i}} \sum_{(\tau_{i}, \tau_{d+1}) \in \mathcal{D}_{i}} \log \mathbb{P}_{\theta_{i}', i}^{\pi}\left(\tau_{i} | \tau_{d+1}\right) - \beta_{i} \right\}$ 

2317 2318 2319

2272 2273

2281 2282

2283 2284 2285

2298

2299

2303

2305 2306

2307

2308

2309

2310 2311 2312

2313

2314 2315

2316

2320 corresponding to agent  $i \in [d+1]$  as 2321

 $\boldsymbol{\tau}_{H} = (o_{1}, a_{1}, \mathbf{M}_{1}, \dots, o_{H}, a_{H}, \mathbf{M}_{H}), \quad \boldsymbol{\tau}_{H,d+1} = (o_{1}, a_{1}, \dots, o_{H}, a_{H}), \quad \boldsymbol{\tau}_{H,i} = (M_{1,i}, M_{2,i}, \dots, M_{H,i}), \quad i \in [d].$ 

The probability of individual  $\mathbb{P}_{\theta_{d+1},d+1}^{\pi}(\tau_{d+1}|\{\tau_i\}_{i=1}^d)$  and  $\mathbb{P}_{\theta_i,i}(\tau_i|\tau_{d+1}), i \in [d]$  is defined as:  $\mathbb{P}_{\theta_{d+1},d+1}^{\pi}\left(\tau_{d+1} \mid \{\tau_i\}_{i=1}^d\right) = \sum_{s_1,s_2,\dots,s_H} \mu_1(s_1) \mathbb{O}_1(o_1 \mid s_1) \pi_1(a_1 \mid o_1, \mathbf{M}_1) \mathbb{T}_{1,a_1}(s_2 \mid s_1) \times \cdots$  $\times \mathbb{O}(o_{H-1} \mid s_{H-1})\pi_{H-1}(a_{H-1} \mid o_1, a_1, \mathbf{M}_1, \dots, o_{H-2}, a_{H-2}, \mathbf{M}_{H-2}, o_{H-1})\mathbb{T}_{H-1, a_{H-1}}(s_H \mid s_{H-1})$  $\times \mathbb{O}(o_H \mid s_H) \times \pi_{H-1}(a_H \mid o_1, a_1, \mathbf{M}_1, \dots, o_{H-1}, a_{H-1}, \mathbf{M}_{H-1}, o_H),$  $\mathbb{P}_{\theta_{i},i}(\tau_{i} \mid \tau_{d+1}) = \mathbb{T}_{1,i}(\mathbf{M}_{2,i} \mid \mathbf{M}_{1,i}, o_{1}, a_{1}) \times \mathbb{T}_{2,i}(\mathbf{M}_{3,i} \mid \mathbf{M}_{2,i}, o_{2}, a_{2}) \times \dots \times \mathbb{T}_{H-1}(\mathbf{M}_{H,i} \mid \mathbf{M}_{H-1,i}, o_{H-1}, a_{H-1}).$ 

#### F.3 THEORETICAL GUARANTEE

We establish the following theoretical guarantee concerning the regret of Algorithm F.2. 

**Theorem F.1.** Let  $\beta_{d+1} = c(H(SA + SO)\log(TSAOH) + \log(Td/\delta))$  and  $\beta_i =$  $c(HSM^2m^2OA\log(TSMmAOH) + \log(dT/\delta))$  for all  $i \in [d]$ , where c is a constant. Then, with probability at least  $1 - \delta$ , Algorithm F.2 ensures the following inequality:

$$\textit{Regret}(k) \leq \tilde{\mathcal{O}}\left(\frac{S^2AO}{\alpha^2}\sqrt{k(S^2A + SO)} \times \textit{poly}(H) + d(Mm)^2OA\sqrt{k(Mm)^2OA} \times \textit{poly}(H)\right).$$

The proof of this theorem closely follows the derivation outlined in previous sections. Here, we present key steps while omitting detailed proofs for some of the lemmas for clarity. 

**Lemma F.1.** With probability at least  $1 - \delta$ , the following events hold:

$$\max_{(\theta_i,k)\in\Theta_i\times[T]}\sum_{t=1}^k \log\left(\frac{\mathbb{P}_{\theta_i,i}(\tau_i^t\mid\tau_{d+1}^t)}{\mathbb{P}_{\theta_i^*,i}(\tau_i^t\mid\tau_{d+1}^t)}\right) > c\left(H(SM^2m^2OA)\log(TSMmAOH) + \log(dT/\delta)\right), i \in [d].$$

**Lemma F.2.** The following event holds with probability at least  $1 - \delta$  for all  $\theta_i \in \Theta_i$ ,  $i \in [d]$ .

$$\begin{aligned}
&\sum_{t=1}^{2352} \sum_{t=1}^{k} \left( \sum_{\tau} \left| \mathbb{P}_{\theta_{d+1},d+1}^{\pi^{t}}(\tau_{d+1} \mid \{\tau_{i}\}_{i=1}^{d}) - \mathbb{P}_{\theta_{d+1},d+1}^{\pi^{t}}(\tau_{d+1} \mid \{\tau_{i}\}_{i=1}^{d}) \right| \prod_{i=1}^{d} \mathbb{P}_{\theta_{i}^{*},i}(\tau_{i} \mid \tau_{d+1}) \right)^{2} \\
&\sum_{t=1}^{2355} \leq c \left( \sum_{t=1}^{k} \log \left( \frac{\mathbb{P}_{\theta_{d+1},d+1}^{\pi^{t}}(\tau_{d+1}^{t} \mid \{\tau_{i}^{t}\}_{i=1}^{d})}{\mathbb{P}_{\theta_{d+1},d+1}^{\pi^{t}}(\tau_{d+1}^{t} \mid \{\tau_{i}^{t}\}_{i=1}^{d})} \right) + H(SA + SO) \log(TSAOH) + \log(Td/\delta) \right), \\
&\sum_{t=1}^{2358} \sum_{t=1}^{k} \left( \sum_{\tau} \left| \mathbb{P}_{\theta_{i},i}(\tau_{i} \mid \tau_{d+1}) - \mathbb{P}_{\theta_{i}^{*},i}(\tau_{i} \mid \tau_{d+1}) \right| \left[ \prod_{j=1}^{i-1} \mathbb{P}_{\theta_{j}^{*},j} \right] \left[ \prod_{j=i+1}^{d} \mathbb{P}_{\theta_{j}^{t},j} \right] \mathbb{P}_{\theta_{d+1},d+1}^{\pi^{t}}(\tau_{d+1} \mid \{\tau_{i}\}_{i=1}^{d}) \right)^{2} \\
&\sum_{t=1}^{2362} \sum_{t=1}^{k} \log \left( \frac{\mathbb{P}_{\theta_{i},i}(\tau_{i}^{t} \mid \tau_{d+1}^{t})}{\mathbb{P}_{\theta_{i}^{*},i}(\tau_{i}^{t} \mid \tau_{d+1}^{t})} \right) + H(SM^{2}m^{2}OA) \log(TSMmOAH) + \log(dT/\delta) \right), \\
&\sum_{t=1}^{2365} \sum_{t=1}^{k} \log \left( \frac{\mathbb{P}_{\theta_{i},i}(\tau_{i}^{t} \mid \tau_{d+1}^{t})}{\mathbb{P}_{\theta_{i}^{*},i}(\tau_{i}^{t} \mid \tau_{d+1}^{t})} \right) + H(SM^{2}m^{2}OA) \log(TSMmOAH) + \log(dT/\delta) \right), \\
&\sum_{t=1}^{2365} \sum_{t=1}^{k} \sum_{t=1}^{k} \log \left( \frac{\mathbb{P}_{\theta_{i},i}(\tau_{i}^{t} \mid \tau_{d+1}^{t})}{\mathbb{P}_{\theta_{i}^{*},i}(\tau_{i}^{t} \mid \tau_{d+1}^{t})} \right) + H(SM^{2}m^{2}OA) \log(TSMmOAH) + \log(dT/\delta) \right), \\
&\sum_{t=1}^{2365} \sum_{t=1}^{k} \sum_{t=1}^{k} \mathbb{P}_{\theta_{i},i}^{t}(\tau_{i}^{t} \mid \tau_{d+1}^{t}) + \mathbb{P}_{\theta_{i}^{*},i}^{t}(\tau_{i}^{t} \mid \tau_{d+1}^{t})} \right) \\
&\sum_{t=1}^{k} \sum_{t=1}^{k} \sum_{t=1}^{k} \mathbb{P}_{\theta_{i},i}^{t}(\tau_{i}^{t} \mid \tau_{d+1}^{t}) + \mathbb{P}_{\theta_{i}^{*},i}^{t}(\tau_{i}^{t} \mid \tau_{d+1}^{t}) + \mathbb{P}_{\theta_{i}^{*},i}^{t}(\tau_{d+1}^{t} \mid \tau_{d+1}^{t}) + \mathbb{P}_{\theta_{i}^{*},i}^{t}(\tau_{d+1}^{t}) + \mathbb{P}_{\theta_{i}^{*},i}^{t}(\tau_{d+1}^{t}) + \mathbb{P}_{\theta_{i}^{*},i}^{t}(\tau_{d+1}^{t} \mid \tau_{d+1}^{t}) + \mathbb{P}_{\theta_{i}^{*},i}^{t}(\tau_{d+1}^{t} \mid \tau_{d+1}^{t}) + \mathbb{P}_{\theta_{i}^{*},i}^{t}(\tau_{d+1}^{t}) + \mathbb{P}_{\theta_{i}^{*},i}^{t}(\tau_{d+1}^{t}) + \mathbb{P}_{\theta_{i}^{*},i}^{t}(\tau_{d+1}^{t} \mid \tau_{d+1}^{t}) + \mathbb{P$$

where for all  $m \in [d]$ ,  $\theta_m \in \Theta_m$ , we denote  $\mathbb{P}_{\theta_m,m}$  as  $\mathbb{P}_{\theta_m,m}(\tau_m \mid \tau_{d+1})$ . **Definition F.1.** For  $m \in [n-1]$ , we define the *B*-operator as follows:

h'=1

$$\begin{aligned} \mathbf{b}_{0,d+1} &= \mathbb{O}_{1}\mu_{1} \in \mathbb{R}^{O}, \\ \mathbf{B}_{h,d+1}(o,a) &= \mathbb{O}_{h+1}\mathbb{T}_{h,a} \cdot diag(\mathbb{O}_{h}(o \mid \cdot))\mathbb{O}_{h}^{\dagger} \in \mathbb{R}^{O \times O}, \forall h \in [H], \\ \mathbf{b}_{h,i}(\tau_{h+1,i}) &= \left[\prod_{h'=1}^{h} \mathbb{T}_{h',i}(M_{h'+1,i} \mid M_{h',i}, o_{h'}, a_{h'})\right], \\ \mathbf{b}_{h,d+1}(\tau_{h,d+1}) &= \left[\prod_{h'=1}^{h} \mathbf{B}_{h',d+1}(o_{h'}, a_{h'})\right] \mathbf{b}_{0,d+1}. \end{aligned}$$

]

**Lemma F.3.** With probability at least  $1 - \delta$ , the following inequality holds true for all  $h \in [H], k \in$  $[T], i \in [d]$ :  $\sum_{t=1}^{k-1} \sum_{\tau} \pi^{t}(\tau_{h}) \| \left( \mathbf{B}_{h,d+1}^{k}(o_{h},a_{h}) - \mathbf{B}_{h,d+1}(o_{h},a_{h}) \right) \mathbf{b}_{h-1,d+1}(\tau_{h-1,d+1}) \|_{1} \left[ \prod_{i=1}^{d} \mathbf{b}_{h-1,i}(\tau_{h,i}) \right] = \mathcal{O}\left( \frac{\sqrt{S}}{\alpha} \sqrt{k\beta_{d+1}} \right),$  $\sum_{i=1}^{k-1} \sum_{-} \| \left( \mathbb{T}_{h-1,i}^{k}(M_{h,i} \mid M_{h-1,i}, o_{h-1}, a_{h-1}) - \mathbb{T}_{h-1,i}(M_{h,i} \mid M_{h-1,i}, o_{h-1}, a_{h-1}) \mathbf{b}_{h-2,i}(\tau_{h-1,i}) \right) \|_{1} \times \mathbf{T}_{h-1,i}^{k}(M_{h,i} \mid M_{h-1,i}, o_{h-1}, a_{h-1}) - \mathbf{T}_{h-1,i}(M_{h,i} \mid M_{h-1,i}, o_{h-1}, a_{h-1}) \mathbf{b}_{h-2,i}(\tau_{h-1,i}) \right) \|_{1} \times \mathbf{T}_{h-1,i}^{k}(M_{h,i} \mid M_{h-1,i}, o_{h-1}, a_{h-1}) - \mathbf{T}_{h-1,i}(M_{h,i} \mid M_{h-1,i}, o_{h-1}, a_{h-1}) \mathbf{b}_{h-2,i}(\tau_{h-1,i}) \|_{1} \times \mathbf{T}_{h-1,i}^{k}(M_{h,i} \mid M_{h-1,i}, o_{h-1}, a_{h-1}) - \mathbf{T}_{h-1,i}(M_{h,i} \mid M_{h-1,i}, o_{h-1}, a_{h-1}) \mathbf{b}_{h-2,i}(\tau_{h-1,i}) \|_{1} \times \mathbf{T}_{h-1,i}^{k}(M_{h,i} \mid M_{h-1,i}, o_{h-1}, a_{h-1}) - \mathbf{T}_{h-1,i}(M_{h,i} \mid M_{h-1,i}, o_{h-1}, a_{h-1}) \mathbf{b}_{h-2,i}(\tau_{h-1,i}) \|_{1} \times \mathbf{T}_{h-1,i}^{k}(M_{h,i} \mid M_{h-1,i}, o_{h-1}, a_{h-1}) \|_{1} \times \mathbf{T}_{h-1,i}^{k}(M_{h,i} \mid M_{h-1,i}, a_{h-1}) \|_{1} \times \mathbf{T}_{h-1,i}^{k}(M_{h,i} \mid M_{h-1,i}, a_{h-1}) \|_{1} \times \mathbf{T}_{h-1,i}^{k}(M_{h-1,i}, a_{h-1}) \|_{1} \times \mathbf{T}_{h-1,i}^{k}(M_{h-1,i}, a_{h-1}) \|_{1} \times \mathbf{T}_{h-1,i}^{k}(M_{h-1,i}, a_{h-1,i}) \|_{1} \times \mathbf{T}_{h-1,i}^{k}(M_{h-1,i}, a_{h-1}) \|_{1} \times \mathbf{T}_{h-1,i}^{k}(M_{h-1,i}, a_{h-1}) \|_{1} \times \mathbf{T}_{h-1,i}^{k}(M_{h-1,i}, a_{h-1,i}) \|_{1} \times \mathbf{T}_{h-1,i}^{k}(M_{h-1$  $\cdot \left[\prod_{i=1}^{i-1} \mathbf{b}_{h-1,j}(\tau_{h,j})\right] \left[\prod_{i=i+1}^{d} \mathbf{b}_{h-1,j}^{t}(\tau_{h,j})\right] \|\mathbf{b}_{h,d+1}^{t}(\tau_{h,d+1})\|_{1} \pi_{d+1}^{t}(\boldsymbol{\tau}_{h}) = \mathcal{O}\left(\sqrt{k\beta_{i}}\right).$ 

Lemma F.4. the regret by the error of operator estimates is bounded by the following term:

$$Regret(k) \leq \sum_{t=1}^{k} \sum_{h=1}^{H-1} \sum_{\tau_h} \frac{S^{1.5}}{\alpha} \| (\mathbf{B}_{h,d+1}^t - \mathbf{B}_{h,d+1}) \mathbf{b}_{h-1,d+1}(\tau_{h-1,d+1}) \|_1 \cdot \pi^t(\tau_h) \cdot \left[ \prod_{i=1}^{d} \mathbf{b}_{h-1,i}(\tau_{h,i}) \right] \\ + \sum_{t=1}^{k} \sum_{h=1}^{H-1} \sum_{i=1}^{d} \sum_{\tau_h} \| (\mathbb{T}_{h-1,i}^t - \mathbb{T}_{h-1,i}) \mathbf{b}_{h-2,i} \|_1 \left[ \prod_{j=1}^{i-1} \mathbf{b}_{h-1,j}(\tau_{h,j}) \right] \left[ \prod_{j=i+1}^{d} \mathbf{b}_{h-1,j}^t(\tau_{h,j}) \right] \| \mathbf{b}_{h,d+1}^t(\tau_{h,d+1}) \|_1 \pi^t(\tau_h),$$

$$Regret(k) \leq \sum_{t=1}^{k} \sum_{h=1}^{H-1} \sum_{i=1}^{d} \sum_{\tau_h} \| (\mathbb{T}_{h-1,i}^t - \mathbb{T}_{h-1,i}) \mathbf{b}_{h-2,i} \|_1 \left[ \prod_{j=1}^{i-1} \mathbf{b}_{h-1,j}(\tau_{h,j}) \right] \left[ \prod_{j=i+1}^{d} \mathbf{b}_{h-1,j}^t(\tau_{h,j}) \right] \| \mathbf{b}_{h,d+1}^t(\tau_{h,d+1}) \|_1 \pi^t(\tau_h),$$

where for all  $h \in [H]$ , we denote  $\mathbf{B}_{h,d+1}^t$  as  $\mathbf{B}_{h,d+1}^t(o_h, a_h)$  and denote  $\mathbf{B}_{h,d+1}$  as  $\mathbf{B}_{h,d+1}(o_h, a_h)$ . For all  $i \in [d]$  and  $h \in [H]$ , we denote  $\mathbb{T}_{h,i}^t$  as  $\mathbb{T}_{h,i}^t$   $(M_{h,i} \mid M_{h-1,i}, o_{h-1,a_{h-1}})$ , and denote  $\mathbb{T}_{h,i}$  as  $\mathbb{T}_{h,i}(M_{h,i} \mid M_{h-1,i}, o_{h-1,a_{h-1}})$ . Moreover, for all  $h \in [H], i \in [d]$ , we denote  $\mathbf{b}_{h,i}$  as  $\mathbf{b}_{h,i}(\tau_{h+1,i})$ . **Theorem F.2.** Let  $\beta_{d+1} = c(H(SA + SO)\log(TSAOH) + \log(Td/\delta))$  and  $\beta_i =$  $c(HSM^2m^2OA\log(TSMmAOH) + \log(dT/\delta))$  for all  $i \in [d]$ , where c is a constant. Then, with probability at least  $1 - \delta$ , Algorithm F.2 ensures the following inequality: 

*Proof.* The target is to bound the following d + 1 terms: 

2407 2408 
$$\sum_{t=1}^{k} \sum_{h=1}^{H-1} \sum_{\tau_{h}} \| (\mathbf{B}_{h,d+1}^{t} - \mathbf{B}_{h,d+1}) \mathbf{b}_{h-1,d+1}(\tau_{h-1,d+1}) \|_{1} \cdot \pi^{t}(\boldsymbol{\tau}_{h}) \cdot \left[ \prod_{i=1}^{d} \mathbf{b}_{h-1,i}(\tau_{h,i}) \right],$$
2409 2409 2410 
$$\sum_{t=1}^{k} \sum_{h=1}^{H-1} \sum_{\tau_{h}} \| (\mathbb{T}_{h-1,i}^{t} - \mathbb{T}_{h-1,i}) \mathbf{b}_{h-2,i} \|_{1} \left[ \prod_{j=1}^{i-1} \mathbf{b}_{h-1,j}(\tau_{h,j}) \right] \left[ \prod_{j=i+1}^{d} \mathbf{b}_{h-1,j}^{t}(\tau_{h,j}) \right] \| \mathbf{b}_{h,d+1}^{t}(\tau_{h,d+1}) \|_{1} \pi^{t}(\boldsymbol{\tau}_{h}), i \in [d].$$
2413 The provide the state of the s

The condition we have is that with probability at least  $1 - \delta$ , 

2415 
$$\sum_{t=1}^{k-1} \sum_{h=1}^{H-1} \sum_{\tau_h} \| (\mathbf{B}_{h,d+1}^k - \mathbf{B}_{h,d+1}) \mathbf{b}_{h-1,d+1}(\tau_{h-1,d+1}) \|_1 \cdot \pi^t(\tau_h) \cdot \left[ \prod_{i=1}^d \mathbf{b}_{h-1,i}(\tau_{h,i}) \right] = \mathcal{O}\left( \frac{\sqrt{S}}{\alpha} \sqrt{k\beta_{d+1}} \right),$$
2417 
$$\sum_{t=1}^{k-1} \sum_{h=1}^{H-1} \sum_{\tau_h} \| (\mathbb{T}_{h-1,i}^k - \mathbb{T}_{h-1,i}) \mathbf{b}_{h-2,i} \|_1 \left[ \prod_{j=1}^{i-1} \mathbf{b}_{h-1,j} \right] \left[ \prod_{j=i+1}^d \mathbf{b}_{h-1,j}^t \right] \| \mathbf{b}_{h,d+1}^t(\tau_{h,d+1}) \|_1 \pi^t(\tau_h) = \mathcal{O}\left( \sqrt{k\beta_i} \right), i \in [d]$$
2421 Then, we can easily Elyder dimension lemme (Lemma 22) to obtain that the tenset is bounded by:

Then, we can apply Eluder dimension lemma (Lemma ??) to obtain that the target is bounded by 

$$\begin{aligned}
& 2423 \\
& 2424 \\
& 2424 \\
& 2425 \\
& 2425 \\
& 2425 \\
& 2425 \\
& 2425 \\
& 2425 \\
& 2426 \\
& 2426 \\
& 2426 \\
& 2427 \\
& 2428 \\
& 2428 \\
& 2428 \\
& 2428 \\
& 2428 \\
& 2429 \\
& = \mathcal{O}\left((Mm)^2 OA\sqrt{k\beta_i}\right), i \in [d].
\end{aligned}$$

2430 Therefore, we can achieve that the regret is bounded by 2431

$$R^k \leq \tilde{\mathcal{O}}\left(\frac{S^2AO}{\alpha^2}\sqrt{k(S^2A+SO)}\times \operatorname{poly}(H) + d(Mm)^2OA\sqrt{k(Mm)^2OA}\times \operatorname{poly}(H)\right).$$

2433 2434 2435

2432

#### F.4 IMPROVEMENT TO ACHIEVE SHARPER BOUND 2436

#### 2437 F.4.1 MOTIVATION 2438

2439 Algorithm F.2 does not appear to achieve the optimal sample complexity. The intuition is that 2440 dummy agents  $1, 2, \ldots, d$  can observe their exact state. Therefore, they experience an MDP process. 2441 If we directly apply OMLE to the single-agent MDP setting, the algorithm yields the regret forms 2442 like  $R^k \leq \tilde{\mathcal{O}}(S^2A\sqrt{k(S^2A)} \times \text{poly}(H))$ . The regret is scaled in  $A^{1.5}$ . However, if we apply the 2443 UCB-VI algorithm to the single-agent MDP, we can obtain:  $R^k \leq \tilde{\mathcal{O}}(H^2 \sqrt{S^2 A k})$ . The regret is 2444 scaled in A<sup>0.5</sup>. Therefore, by combining UCB-VI with OMLE, we might achieve a sharper bound 2445 on the regret. 2446

We still use OMLE to estimate the model parameter  $\theta_{d+1}$  for the  $d + 1^{th}$  agent. For dummy agent 2447  $1, 2, \ldots, d$ , we first use number of times each state tuple  $(M_i, o, a, M'_i)$  and  $(M_i, o, a)$  is visited. 2448 Namely we have 2449

2450  
2451 
$$N_{h}^{k}(M_{i}, o, a, M_{i}') = \sum_{t=1}^{k} 1_{(M_{h,i}^{t}, o_{h}^{t}, a_{h}^{t}, M_{h+1,i}^{t}) = (M_{i}, o, a, M_{i}')}, N_{h}^{k}(M_{i}, o, a) = \sum_{t=1}^{k} 1_{(M_{h,i}^{t}, o_{h}^{t}, a_{h}^{t}) = (M_{i}, o, a)},$$
2452  
2453 
$$\hat{\mathbb{T}}_{h}^{k}(M_{i}' \mid M_{i}, o, a) = \frac{N_{h}^{k}(M_{i}, o, a, M_{i}')}{N_{h}^{k}(M_{i}, o, a)}.$$
2454 We define the bonus as follows:  
2456 
$$b_{h}^{k}(M_{i}, o, a) = \sqrt{\frac{2Mm\ln(MmOAKHd/\delta)}{N_{h}^{k}(M_{i}, o, a)}} + \frac{2\ln(MmOAKHd/\delta)}{N_{h}^{k}((M_{i}, o, a))},$$
2458 
$$\beta_{d+1} = c(H(SA + SO)\log(TSAOH) + \log(Td/\delta)).$$

245 2460 2461

2462

# Algorithm 9 OMLE for POMDPwk

1: Initialize:  $\mathcal{B}_{d+1}^1 = \{\hat{\theta}_{d+1} \in \Theta_{d+1} : \min_h \sigma_S(\mathbb{O}_i(\hat{\theta}_{d+1}) \ge \alpha)\}, \mathcal{B}_i^1 = \{\Theta_i\}, \mathcal{D}_i = \{\}, \text{ for all } i \le n \}$ 2463 players  $i \in [d]$ . 2464 2: for k = 1, ..., K do 2465 compute  $(\theta_1^k, \theta_2^k, \dots, \theta_n^k, \pi^k) = \arg \max_{\hat{\theta}_1 \in \mathcal{B}_1^k, \hat{\theta}_2 \in \mathcal{B}_2^k, \dots, \hat{\theta}_n \in \mathcal{B}_n^k, \pi} V^{\pi}(\hat{\theta})$ follow  $\pi^k$  to collect a trajectory  $\boldsymbol{\tau}^k = (o_1^k, a_1^k, \mathbf{M}_1^k, \dots, o_H^k, a_H^k, \mathbf{M}_H^k)$ 3: 2466 2467 4: 2468 add  $(\pi^k, \boldsymbol{\tau}^k)$  into  $\mathcal{D}_{d+1}$ , and then update 5: 2469  $\mathcal{B}_{d+1}^{k+1} = \left\{ \hat{\theta}_{d+1} \in \mathcal{B}_{d+1}^1 : \sum_{(\pi, \tau) \in \mathcal{D}_{d+1}} \log \mathbb{P}_{\hat{\theta}_{d+1}, d+1}^{\pi} \ge \max_{\theta'_{d+1}} \sum_{(\pi, \tau) \in \mathcal{D}_{d+1}} \log \mathbb{P}_{\theta'_{d+1}, d+1}^{\pi} - \beta_{d+1} \right\}$ 2470 2471 2472 for i = 1, 2, ..., d do 6: 2473  $\mathcal{B}_{i}^{k+1} = \left\{ \hat{\theta}_{i} \in \mathcal{B}_{i}^{1} : \sum_{M_{i}, i \neq i} \left| \mathbb{T}_{\hat{\theta}_{d+1}, h}(M_{h+1, i} \mid M_{h, i}, o, a) - \hat{\mathbb{T}}_{h}^{k}(M_{h+1, i} \mid M_{h, i}, o, a) \right| \le b_{h}^{k}(M_{h, i}, o, a), \forall (h, o, a, M_{h, i}) \right\}$ 

2474 2475

2483

### F.5 THEORETICAL GUARANTEE

2480 **Theorem F.3.** With probability at least  $1 - \delta$ , Algorithm F.4.1 guarantees the following bound on 2481 the regret: 2482

$$\operatorname{Regret}(k) \leq \tilde{\mathcal{O}}\left(\frac{S^2AO}{\alpha^2}\sqrt{k(S^A+SO)} \times \operatorname{poly}(H) + dHBm\sqrt{kOA}\right).$$

The proof of this theorem closely follows the derivation outlined in previous sections. Here, we present key steps while omitting detailed proofs for some of the lemmas for clarity.
 Window Line and the lemma for clarity.

**Lemma F.5.** With probability at least  $1 - \delta$ , for all  $o, a, M_{h,i}, k, h$ , we have

$$\sum_{M_{h+1,i}} \left| \mathbb{T}_{\theta_{d+1}^*,h}(M_{h+1,i} \mid M_{h,i}, o, a) - \hat{\mathbb{T}}_h^k(M_{h+1,i} \mid M_{h,i}, o, a) \right| \le b_h^k(M_{h,i}, o, a).$$

*Proof.* Consider a fixed tuple  $(o, a, M_{h,i}, h, k)$ . Define  $\mathcal{H}_{h,t}$  as the history starting from the begin-2492 ning of episode 1 to step h at episode t (including step h, i.e. up to  $s_h^t, a_h^t$ ). Define random variables 2493  $\{X_t\}_{t\geq 0}$  as

2494  
2495 
$$X_t = \mathbb{1}_{(o_h^t, a_h^t, M_{h,i}^t, M_{h+1,i}^t) = (o, a, M_{h,i}, M_{h+1,i})} - \mathbb{1}_{(o_h^t, a_h^t, M_{h,i}^t) = (o, a, M_{h,i})} \mathbb{T}_{\theta_{d+1}^*, h}(M_{h+1,i} \mid M_{h,i}, o, a).$$

2496 We now show that  $\{X_t\}_{t\geq 0}$  is a martingale sequence adapted to filtration  $\{\mathcal{H}_{h,t}\}_{t\geq 0}$ . Note that 2497  $\mathbb{E}[X_t \mid \mathcal{H}_{h,t}] = 0$ . We have  $|X_t| \leq 1$ . To use Azuma-Bernstein's inequality, we note that 2498  $\mathbb{E}[X_t^2 \mid \mathcal{H}_{h,t}]$  is bounded as:

$$\mathbb{E}\left[X_{t}^{2} \mid \mathcal{H}_{h,t}\right] = \mathbb{1}_{(o_{h}^{t}, a_{h}^{t}, M_{h,i}^{t}) = (o, a, M_{h,i})} \mathbb{T}_{\theta_{d+1}^{*}, h}(M_{h+1,i} \mid M_{h,i}, o, a) \left(1 - \mathbb{T}_{\theta_{d+1}^{*}, h}(M_{h+1,i} \mid M_{h,i}, o, a)\right),$$

where we use the fact that the variance of a Bernoulli with parameter p is p(1 - p). This means that

2503  
2504 
$$\sum_{t=1}^{k} \mathbb{E} \left[ X_t^2 \mid \mathcal{H}_{h,t} \right] = N_h^k(M_i, o, a) \mathbb{T}_{\theta_{d+1}^*, h}(M_{h+1,i} \mid M_{h,i}, o, a) \left( 1 - \mathbb{T}_{\theta_{d+1}^*, h}(M_{h+1,i} \mid M_{h,i}, o, a) \right).$$
2505

Now we apply Bernstein's inequality on the martingale difference sequence  $\{X_t\}_{t\geq 0}$ , we have

$$\begin{split} \left| \sum_{t=1}^{k} X_{t} \right| &\leq \sqrt{\frac{2\mathbb{T}_{\theta_{d+1}^{*},h}(M_{h+1,i} \mid M_{h,i}, o, a) \left(1 - \mathbb{T}_{\theta_{d+1}^{*},h}(M_{h+1,i} \mid M_{h,i}, o, a)\right) \ln(1/\delta)}{N_{h}^{k}(M_{h,i}, o, a)}} + \frac{2\ln(1/\delta)}{N_{h}^{k}(M_{h,i}, o, a)} \\ &\leq \sqrt{\frac{2\mathbb{T}_{\theta_{d+1}^{*},h}(M_{h+1,i} \mid M_{h,i}, o, a) \ln(1/\delta)}{N_{h}^{k}(M_{h,i}, o, a)}} + \frac{2\ln(1/\delta)}{N_{h}^{k}(M_{h,i}, o, a)}. \end{split}$$

2514 We apply a union bound over all  $B_{h,i} \in \mathcal{B}_i, o \in \mathcal{O}, a \in \mathcal{A}, h \in [H], k \in [K], i \in [d]$ , and we can achieve that

$$\begin{aligned} \left| \mathbb{T}_{\theta_{d+1}^*,h}(M_{h+1,i} \mid M_{h,i}, o, a) - \hat{\mathbb{T}}_h^k(M_{h+1,i} \mid M_{h,i}, o, a) \right| \\ = \left| \sum_{t=1}^k X_t \right| &\leq \sqrt{\frac{2\mathbb{T}_{\theta_{d+1}^*,h}(M_{h+1,i} \mid M_{h,i}, o, a)L}{N_h^k(M_{h,i}, o, a)}} + \frac{2L}{N_h^k(M_{h,i}, o, a)}, \end{aligned}$$

where we define  $\ln(MmOAKHd/\delta)$ .

**2522 Corollary F.1.** With probability at least  $1 - \delta$ ,

$$(\theta_1^*, \theta_2^*, \dots, \theta_d^*) \in \bigcap_{k \in [K]} (\mathcal{B}_1^k \times \mathcal{B}_1^k \times \dots \times \mathcal{B}_{d+1}^k).$$

**Lemma F.6.** With probability at least  $1 - \delta$ , the following event holds.

$$\max_{(\theta_{d+1},k)\in\Theta_{d+1}\times[T]}\sum_{t=1}^k \log\left(\frac{\mathbb{P}_{\theta_{d+1},d+1}^{\pi^t}(\tau_{d+1}^t \mid \{\tau_i^t\}_{i=1}^d)}{\mathbb{P}_{\theta_{d+1},d+1}^{\pi^t}(\tau_{d+1}^t \mid \{\tau_i^t\}_{i=1}^d)}\right) > c(H(SA+SO)\log(TSAOH) + \log(Td/\delta)).$$

**Lemma F.7.** We can obtain that the following event holds with probability at least  $1 - \delta$  for all  $\theta_{d+1} \in \Theta_{d+1}$  and  $k \in [T]$ .

$$\sum_{t=1}^{k} \left( \sum_{\tau} \left| \mathbb{P}_{\theta_{d+1},d+1}^{\pi^{t}}(\tau_{d+1} \mid \{\tau_{i}\}_{i=1}^{d}) - \mathbb{P}_{\theta_{d+1},d+1}^{\pi^{t}}(\tau_{d+1} \mid \{\tau_{i}\}_{i=1}^{d}) \right| \left[ \prod_{i=1}^{d} \mathbb{P}_{\theta_{i}^{*},i}(\tau_{i} \mid \tau_{d+1}) \right] \right)^{2} \\ \leq c \left( \sum_{t=1}^{k} \log \left( \frac{\mathbb{P}_{\theta_{d+1},d+1}^{\pi^{t}}(\tau_{d+1}^{t} \mid \{\tau_{i}^{t}\}_{i=1}^{d})}{\mathbb{P}_{\theta_{d+1}^{*},d+1}^{\pi^{t}}(\tau_{d+1}^{t} \mid \{\tau_{i}^{t}\}_{i=1}^{d})} \right) + H(SA + SO) \log(TSAOH) + \log(Td/\delta) \right).$$

2538 Lemma F.8. With probability at least  $1 - \delta$ , the following event holds: 2540  $\sum_{t=1}^{k-1} \sum_{\tau_h} \pi^t(\tau_h) \| (\mathbf{B}_{h,d+1}^k(o_h, a_h) - \mathbf{B}_{h,d+1}(o_h, a_h)) \mathbf{b}_{h-1,d+1}(\tau_{h-1,d+1}) \|_1 \left[ \prod_{i=1}^d \mathbf{b}_{h-1,i}(\tau_{h,i}) \right] = \mathcal{O}\left( \frac{\sqrt{S}}{\alpha} \sqrt{k\beta_{d+1}} \right).$ 2542 Theorem F.4. With probability at least  $1 - \delta$ . Algorithm F.4. Lease the following bound on

**Theorem F.4.** With probability at least  $1 - \delta$ , Algorithm F.4.1 guarantees the following bound on the regret:

$$\operatorname{Regret}(k) \leq \tilde{\mathcal{O}}\left(\frac{S^2AO}{\alpha^2}\sqrt{k(S^A+SO)} \times \operatorname{poly}(H) + dHBm\sqrt{kOA}\right).$$

*Proof.* Initially, using similar techniques as in the section of finding global optimal for factored DEC-POMDP, we have with probability at least  $1 - \delta$ ,  $\theta_{d+1} \in \bigcap_{k \in [K]} B_{d+1}^k$ . We combine this with result in F.1, and we can obtain that with probability at least  $1 - \delta$ ,

$$(\theta_1^*, \theta_2^*, \dots, \theta_{d+1}^*) \in \bigcap_{k \in [K]} \left( \mathcal{B}_1^k \times \mathcal{B}_2^k \times \dots \times \mathcal{B}_n^k \right).$$

2553 Therefore, we can bind the regret as:

$$R^{k} = \sum_{t=1}^{k} V_{\theta^{*}}^{*} - V_{\theta^{*}}^{\pi^{t}} \leq \sum_{t=1}^{k} V_{\theta^{t}}^{\pi^{t}} - V_{\theta^{*}}^{\pi^{t}} \leq \sum_{t=1}^{k} \sum_{\tau_{H}} H \left| \mathbb{P}_{\theta^{t}}^{\pi^{t}}(\tau_{H}) - \mathbb{P}_{\theta^{*}}^{\pi^{t}}(\tau_{H}) \right|.$$

2557 According to the factorization condition, we have

$$\begin{split} \sum_{i=1}^{k} \sum_{\tau_{il}} \left| \mathbb{P}_{\theta_{i}^{\star}}^{\pi_{i}}(\tau_{il}) - \mathbb{P}_{\theta_{i}^{\star}}^{\pi_{i}}(\tau_{il}) \right| \\ \sum_{i=1}^{k} \sum_{\tau_{il}} \left| \left[ \prod_{i=1}^{d} \mathbb{P}_{\theta_{i}^{\star},i}^{*}(\tau_{i} + \tau_{d+1}) \right] \mathbb{P}_{\theta_{d+1},d+1}^{\pi_{i}^{\star}}(\tau_{d+1} + |\{\tau_{i}\}_{i=1}^{d}) - \left[ \prod_{i=1}^{d} \mathbb{P}_{\theta_{i}^{\star},i}^{*}(\tau_{i} + \tau_{d+1}) \right] \mathbb{P}_{\theta_{d+1},d+1}^{\pi_{i}^{\star}}(\tau_{i} + \tau_{d+1}) \right] \\ \sum_{i=1}^{k} \sum_{\tau_{il}} \left| \left[ \prod_{i=1}^{d} \mathbb{P}_{\theta_{i}^{\star},i}^{*}(\tau_{i} + \tau_{d+1}) \right] \mathbb{P}_{\theta_{d+1},d+1}^{\pi_{i}^{\star}}(\tau_{d+1}) - \mathbb{P}_{\theta_{i}^{\star}}^{\pi_{i}^{\star}}(\tau_{H}) \right| \\ \sum_{i=1}^{k} \sum_{\tau_{il}} \left| \prod_{i=1}^{d} \mathbb{P}_{\theta_{i}^{\star},i}^{*}(\tau_{i} + \tau_{d+1}) \right| \mathbb{P}_{\theta_{d+1},d+1}^{\pi_{i}^{\star}}(\tau_{d+1} + |\{\tau_{i}\}_{i=1}^{d}) - \left[ \prod_{i=1}^{d} \mathbb{P}_{\theta_{i}^{\star},i}^{*}(\tau_{i} + \tau_{d+1}) \right] \mathbb{P}_{\theta_{d+1},d+1}^{\pi_{i}^{\star}}(\tau_{d+1} + |\{\tau_{i}\}_{i=1}^{d}) \right| \\ \sum_{i=1}^{k} \sum_{\tau_{il}} \prod_{i=1}^{k} \mathbb{P}_{\theta_{i}^{\star},i}^{*}(\tau_{i} + \tau_{d+1}) - \prod_{i=1}^{d} \mathbb{P}_{\theta_{i}^{\star},i}^{*}(\tau_{i} + \tau_{d+1}) \left| \mathbb{P}_{\theta_{d+1},d+1}^{\pi_{i}^{\star}}(\tau_{d+1} + |\{\tau_{i}\}_{i=1}^{d}) \right| \\ + \sum_{i=1}^{k} \sum_{\tau_{il}} \prod_{i=1}^{k} \mathbb{P}_{\theta_{i}^{\star},i}^{*}(\tau_{i} + \tau_{d+1}) - \prod_{i=1}^{d} \mathbb{P}_{\theta_{i}^{\star},i}^{\pi_{i}^{\star}}(\tau_{i} + \tau_{d+1}) \left| \prod_{i=1}^{d} \mathbb{P}_{\theta_{i}^{\star},i}^{*}(\tau_{i} + \tau_{d+1}) \right| \\ + \sum_{i=1}^{k} \sum_{\tau_{il}} \prod_{i=1}^{k} \mathbb{P}_{\theta_{i}^{\star},i}^{\pi_{i}^{\star}}(\tau_{i} + \tau_{d+1}) - \mathbb{P}_{\theta_{i}^{\star},i}^{\pi_{i}^{\star}}(\tau_{i} + \tau_{d+1}) \left| \prod_{i=1}^{d} \mathbb{P}_{\theta_{i}^{\star},i}^{*}(\tau_{i} + \tau_{d+1}) \right| \\ + \sum_{i=1}^{k} \sum_{\tau_{il}} \prod_{i=1}^{k} \mathbb{P}_{\theta_{i}^{\star},i}^{\pi_{i}^{\star}}(\tau_{d+1} + |\{\tau_{i}\}_{i=1}^{d}) - \mathbb{P}_{\theta_{i+1},d+1}^{\pi_{i}^{\star}}(\tau_{d+1} + |\{\tau_{i}\}_{i=1}^{d}) \right| \\ + \sum_{i=1}^{k} \sum_{\tau_{il}} \mathbb{P}_{\theta_{i}^{\star},i}^{\pi_{i}^{\star}}(\tau_{d+1} + |\{\tau_{i}\}_{i=1}^{d}) - \mathbb{P}_{\theta_{i+1},d+1}^{\pi_{i}^{\star}}(\tau_{d+1} + |\{\tau_{i}\}_{i=1}^{d}) \right| \\ + \sum_{i=1}^{k} \sum_{\tau_{il}} \mathbb{P}_{\theta_{i}^{\star},i}^{\pi_{i}^{\star}}(\tau_{d+1} + |\{\tau_{i}\}_{i=1}^{d}) - \mathbb{P}_{\theta_{i+1},d+1}^{\pi_{i}^{\star}}(\tau_{d+1} + |\{\tau_{i}\}_{i=1}^{d}) \right| \\ + \sum_{i=1}^{k} \sum_{\tau_{il}} \mathbb{P}_{\theta_{i}^{\star},i}^{\pi_{i}^{\star}}(\tau_{d+1} + |\{\tau_{i}\}_{i=1}^{d}) - \mathbb{P}_{\theta_{i+1}^{\star},i}^{\pi_{i}^{\star}}(\tau_{d+1} + |\{\tau_{i}\}_{i=1}^{d}) \right| \\ + \sum_{i=1}^{k} \sum_{\tau_{il}} \mathbb{$$

Now we consider the term  $\left|\mathbb{P}_{\theta_{d+1}^{t},d+1}^{\pi^{t}}(\tau_{d+1} \mid \{\tau_{i}\}_{i=1}^{d}) - \mathbb{P}_{\theta_{d+1}^{*},d+1}^{\pi^{t}}(\tau_{d+1} \mid \{\tau_{i}\}_{i=1}^{d})\right|$ . According to the selection of parameter in the implementation of the algorithm, we can derive that for all  $o \in \mathcal{O}, a \in \mathcal{A}, M_{h,i} \in \mathcal{M}_i, k \in [K], h \in [H], i \in [d],$ 

$$\begin{split} & \sum_{M_{h+1,i}} \left| \mathbb{T}_{\theta_{d+1}^*,h}(M_{h+1,i} \mid M_{h,i},o,a) - \mathbb{T}_{\theta_{d+1}^k,h}(M_{h+1,i} \mid M_{h,i},o,a) \right| \\ & \leq \sum_{M_{h+1,i}} \left| \mathbb{T}_{\theta_{d+1}^*,h}(M_{h+1,i} \mid M_{h,i},o,a) - \hat{\mathbb{T}}_h^k(M_{h+1,i} \mid M_{h,i},o,a) \right| \\ & + \sum_{M_{h+1,i}} \left| \mathbb{T}_{\theta_{d+1}^k,h}(M_{h+1,i} \mid M_{h,i},o,a) - \hat{\mathbb{T}}_h^k(M_{h+1,i} \mid M_{h,i},o,a) \right| \\ & \leq 2b_h^k(M_{h,i},o,a). \end{split}$$

Moreover, we have

$$\begin{aligned} & \text{Moteovel, we have} \end{aligned}$$

$$\begin{aligned} & \sum_{t=1}^{k} \sum_{\tau_{H}} \sum_{i=1}^{d} \left[ \prod_{j=1}^{i-1} \mathbb{P}_{\theta_{j}^{*}, j} \right] \left| \mathbb{P}_{\theta_{i}^{*}, i}(\tau_{i} \mid \tau_{d+1}) - \mathbb{P}_{\theta_{i}^{t}, i}(\tau_{i} \mid \tau_{d+1}) \right| \left[ \prod_{j=i+1}^{d} \mathbb{P}_{\theta_{j}^{*}, j} \right] \mathbb{P}_{\theta_{d+1}^{t}, d+1}^{\pi^{t}}(\tau_{d+1} \mid \{\tau_{i}\}_{i=1}^{d}) \\ & \sum_{t=1}^{k} \sum_{h=1}^{H-1} \sum_{i=1}^{d} \sum_{\tau_{h}} \left\| (\mathbb{T}_{h-1, i}^{t} - \mathbb{T}_{h-1, i}) \mathbf{b}_{h-1, i}(\tau_{h-1, i}) \right\|_{1} \left[ \prod_{j=1}^{i-1} \mathbf{b}_{h-1, j} \right] \left[ \prod_{j=i+1}^{d} \mathbf{b}_{h-1, j}^{t} \right] \left\| \mathbf{b}_{h, d+1}^{t}(\tau_{h, d+1}) \right\|_{1} \pi^{t}(\tau_{h}) \\ & \sum_{t=1}^{k} \sum_{h=1}^{H-1} \sum_{i=1}^{d} \sum_{M_{h, i}} \mathbb{E}_{M_{h-1, i}, o_{h-1, a_{h-1}}} \left[ \left| \mathbb{T}_{\theta_{d+1}^{*}, h-1}(M_{h, i} \mid M_{h-1, i}, o_{h-1}, a_{h-1}) - \mathbb{T}_{\theta_{d+1}^{t}, h-1}(M_{h, i} \mid M_{h-1, i}, o_{h-1}, a_{h-1}) \right| \right] \\ & \sum_{t=1}^{k} \sum_{h=1}^{H-1} \sum_{i=1}^{d} \mathbb{E}_{M_{h-1, i}, o_{h-1, a_{h-1}}} \left[ 2b_{h-1}^{t}(M_{h-1, i}, o_{h-1}, a_{h-1}) \right], \end{aligned}$$

where for all  $i \in [d]$ ,  $\theta_i \in \Theta_i$ , we denote  $\mathbb{P}_{\theta_i,i}$  as  $\mathbb{P}_{\theta_i,i}(\tau_i \mid \tau_{d+1})$ . For all  $h \in [H], i \in [d]$ , we denote  $\mathbf{b}_{h,i}$  as  $\mathbf{b}_{h,i}(\tau_{h+1,i})$ , and we denote  $\mathbf{b}_{h,i}^t$  as  $\mathbf{b}_{h,i}^t(\tau_{h+1,i})$ . 

We can then bound the summation of the bonus term with the following inequality:

$$\sum_{t=1}^{k} \sum_{M_{i},o,a} \frac{1}{\sqrt{N_{h}^{t}(M_{i},o,a)}} = \sum_{M_{i},o,a} \sum_{r=1}^{N_{h}^{k}(M_{i},o,a)} \frac{1}{\sqrt{r}} \le \sum_{M_{i},o,a} 2\sqrt{\sum_{M_{i},o,a}} \le \sqrt{MmOA\sum_{M_{i},o,a} N_{h}^{k}(M_{i},o,a)} = \sqrt{kMmOA}.$$

Therefore, we can deduce that

$$\sum_{t=1}^{k} \sum_{\tau_{H}} \sum_{i=1}^{d} \left[ \prod_{j=1}^{i-1} \mathbb{P}_{\theta_{j}^{*}, j} \right] \left| \mathbb{P}_{\theta_{i}^{*}, i}(\tau_{i} \mid \tau_{d+1}) - \mathbb{P}_{\theta_{i}^{t}, i}(\tau_{i} \mid \tau_{d+1}) \right| \left[ \prod_{j=i+1}^{d} \mathbb{P}_{\theta_{j}^{*}, j} \right] \mathbb{P}_{\theta_{d+1}^{t}, d+1}^{\pi^{t}}(\tau_{d+1} \mid \{\tau_{i}\}_{i=1}^{d})$$

$$\leq \sum_{t=1}^{k} \sum_{h=1}^{H-1} \sum_{i=1}^{d} \mathbb{E}_{M_{h-1, i}, o_{h-1}, a_{h-1}} \left[ 2b_{h-1}^{t}(M_{h-1, i}, o_{h-1}, a_{h-1}) \right]$$

$$\leq \mathcal{O} \left( dH \sqrt{kOA} Mm \ln(MmOAKHd/\delta) \right).$$

We are only left to bound the term

$$\sum_{t=1}^{k} \sum_{h=1}^{H-1} \sum_{\boldsymbol{\tau}_{h}} \frac{S^{1.5}}{\alpha} \| (\mathbf{B}_{h,d+1}^{t} - \mathbf{B}_{h,d+1}) \mathbf{b}_{h-1,d+1}(\boldsymbol{\tau}_{h-1,d+1}) \|_{1} \cdot \pi^{t}(\boldsymbol{\tau}_{h}) \cdot \left[ \prod_{i=1}^{d} \mathbf{b}_{h-1,i}(\boldsymbol{\tau}_{h,i}) \right].$$

The condition we have is that with probability at least  $1 - \delta$ , 

2644 
$$\sum_{t=1}^{k-1} \sum_{h=1}^{H-1} \sum_{\tau_h} \| (\mathbf{B}_{h,d+1}^k - \mathbf{B}_{h,d+1}) \mathbf{b}_{h-1,d+1}(\tau_{h-1,d+1}) \|_1 \cdot \pi^t(\tau_h) \cdot \left[ \prod_{i=1}^d \mathbf{b}_{h-1,i}(\tau_{h,i}) \right] = \mathcal{O}\left( \frac{\sqrt{S}}{\alpha} \sqrt{k\beta_{d+1}} \right).$$

We similarly apply Eluder-Dimension lemma (Lemma ??), and we can achieve that

$$\sum_{t=1}^{k} \sum_{h=1}^{H-1} \sum_{\tau_h} \| (\mathbf{B}_{h,d+1}^t - \mathbf{B}_{h,d+1}) \mathbf{b}_{h-1,d+1}(\tau_{h-1,d+1}) \|_1 \cdot \pi^t(\boldsymbol{\tau}_h) \cdot \left[ \prod_{i=1}^d \mathbf{b}_{h-1,i}(\tau_{h,i}) \right] = \mathcal{O}\left( \frac{S^{1.5} H^2 OA}{\alpha} \sqrt{k\beta_{d+1}} \right).$$

2651 Therefore, we can achieve that the regret is bounded by

$$R^k \leq \tilde{\mathcal{O}}\left(\frac{S^2AO}{\alpha^2}\sqrt{k(S^2A+SO)} \times \operatorname{poly}(H) + dHMm\sqrt{kOA}\right).$$

2699

2652